

Springer Series in Statistics

Frank E. Harrell, Jr.

Regression Modeling Strategies

With Applications to Linear Models,
Logistic and Ordinal Regression,
and Survival Analysis

Second Edition

 Springer

Springer Series in Statistics

Advisors:

P. Bickel, P. Diggle, S.E. Feinberg, U. Gather,
I. Olkin, S. Zeger

More information about this series at <http://www.springer.com/series/692>

Frank E. Harrell, Jr.

Regression Modeling Strategies

With Applications to Linear Models,
Logistic and Ordinal Regression,
and Survival Analysis

Second Edition

 Springer

Frank E. Harrell, Jr.
Department of Biostatistics
School of Medicine
Vanderbilt University
Nashville, TN, USA

ISSN 0172-7397 ISSN 2197-568X (electronic)
Springer Series in Statistics
ISBN 978-3-319-19424-0 ISBN 978-3-319-19425-7 (eBook)
DOI 10.1007/978-3-319-19425-7

Library of Congress Control Number: 2015942921

Springer Cham Heidelberg New York Dordrecht London

© Springer Science+Business Media New York 2001

© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media (www.springer.com)

*To the memories of Frank E. Harrell, Sr.,
Richard Jackson, L. Richard Smith, John
Burdeshaw, and Todd Nick, and with
appreciation to Liana and Charlotte
Harrell, two high school math teachers:
Carolyn Wailes (née Gaston) and Floyd
Christian, two college professors: David
Hurst (who advised me to choose the field
of biostatistics) and Doug Stocks, and my
graduate advisor P. K. Sen.*

Preface

There are many books that are excellent sources of knowledge about individual statistical tools (survival models, general linear models, etc.), but the art of data analysis is about choosing and using multiple tools. In the words of Chatfield [100, p. 420] “. . . students typically know the technical details of regression for example, but not necessarily when and how to apply it. This argues the need for a better balance in the literature and in statistical teaching between *techniques* and problem solving *strategies*.” Whether analyzing risk factors, adjusting for biases in observational studies, or developing predictive models, there are common problems that few regression texts address. For example, there are missing data in the majority of datasets one is likely to encounter (other than those used in textbooks!) but most regression texts do not include methods for dealing with such data effectively, and most texts on missing data do not cover regression modeling.

This book links standard regression modeling approaches with

- methods for relaxing linearity assumptions that still allow one to easily obtain predictions and confidence limits for future observations, and to do formal hypothesis tests,
- non-additive modeling approaches not requiring the assumption that interactions are always linear \times linear,
- methods for imputing missing data and for penalizing variances for incomplete data,
- methods for handling large numbers of predictors without resorting to problematic stepwise variable selection techniques,
- data reduction methods (unsupervised learning methods, some of which are based on multivariate psychometric techniques too seldom used in statistics) that help with the problem of “too many variables to analyze and not enough observations” as well as making the model more interpretable when there are predictor variables containing overlapping information,
- methods for quantifying predictive accuracy of a fitted model,

- powerful model validation techniques based on the bootstrap that allow the analyst to estimate predictive accuracy nearly unbiasedly without holding back data from the model development process, and
- graphical methods for understanding complex models.

On the last point, this text has special emphasis on what could be called “presentation graphics for fitted models” to help make regression analyses more palatable to non-statisticians. For example, nomograms have long been used to make equations portable, but they are not drawn routinely because doing so is very labor-intensive. An R function called `nomogram` in the package described below draws nomograms from a regression fit, and these diagrams can be used to communicate modeling results as well as to obtain predicted values manually even in the presence of complex variable transformations.

Most of the methods in this text apply to all regression models, but special emphasis is given to some of the most popular ones: multiple regression using least squares and its generalized least squares extension for serial (repeated measurement) data, the binary logistic model, models for ordinal responses, parametric survival regression models, and the Cox semiparametric survival model. There is also a chapter on nonparametric transform-both-sides regression. Emphasis is given to detailed case studies for these methods as well as for data reduction, imputation, model simplification, and other tasks. Except for the case study on survival of Titanic passengers, all examples are from biomedical research. However, the methods presented here have broad application to other areas including economics, epidemiology, sociology, psychology, engineering, and predicting consumer behavior and other business outcomes.

This text is intended for Masters or PhD level graduate students who have had a general introductory probability and statistics course and who are well versed in ordinary multiple regression and intermediate algebra. The book is also intended to serve as a reference for data analysts and statistical methodologists. Readers without a strong background in applied statistics may wish to first study one of the many introductory applied statistics and regression texts that are available. The author’s course notes *Biostatistics for Biomedical Research* on the text’s web site covers basic regression and many other topics. The paper by Nick and Hardin [476] also provides a good introduction to multivariable modeling and interpretation. There are many excellent intermediate level texts on regression analysis. One of them is by Fox, which also has a companion software-based text [200, 201]. For readers interested in medical or epidemiologic research, Steyerberg’s excellent text *Clinical Prediction Models* [586] is an ideal companion for *Regression Modeling Strategies*. Steyerberg’s book provides further explanations, examples, and simulations of many of the methods presented here. And no text on regression modeling should fail to mention the seminal work of John Nelder [450].

The overall philosophy of this book is summarized by the following statements.

- Satisfaction of model assumptions improves precision and increases statistical power.
- It is more productive to make a model fit step by step (e.g., transformation estimation) than to postulate a simple model and find out what went wrong.
- Graphical methods should be married to formal inference.
- Overfitting occurs frequently, so data reduction and model validation are important.
- In most research projects, the cost of data collection far outweighs the cost of data analysis, so it is important to use the most efficient and accurate modeling techniques, to avoid categorizing continuous variables, and to not remove data from the estimation sample just to be able to validate the model.
- The bootstrap is a breakthrough for statistical modeling, and the analyst should use it for many steps of the modeling strategy, including derivation of distribution-free confidence intervals and estimation of optimism in model fit that takes into account variations caused by the modeling strategy.
- Imputation of missing data is better than discarding incomplete observations.
- Variance often dominates bias, so biased methods such as penalized maximum likelihood estimation yield models that have a greater chance of accurately predicting future observations.
- Software without multiple facilities for assessing and fixing model fit may only seem to be user-friendly.
- Carefully fitting an improper model is better than badly fitting (and overfitting) a well-chosen one.
- Methods that work for all types of regression models are the most valuable.
- Using the data to guide the data analysis is almost as dangerous as not doing so.
- There are benefits to modeling by deciding how many degrees of freedom (i.e., number of regression parameters) can be “spent,” deciding where they should be spent, and then spending them.

On the last point, the author believes that significance tests and P -values are problematic, especially when making modeling decisions. Judging by the increased emphasis on confidence intervals in scientific journals there is reason to believe that hypothesis testing is gradually being de-emphasized. Yet the reader will notice that this text contains many P -values. How does that make sense when, for example, the text recommends against simplifying a model when a test of linearity is not significant? First, some readers may wish to emphasize hypothesis testing in general, and some hypotheses have special interest, such as in pharmacology where one may be interested in whether the effect of a drug is linear in log dose. Second, many of the more interesting hypothesis tests in the text are tests of complexity (nonlinearity, interaction) of the overall model. Null hypotheses of linearity of effects in particular are

frequently rejected, providing formal evidence that the analyst's investment of time to use more than simple statistical models was warranted.

The rapid development of Bayesian modeling methods and rise in their use is exciting. Full Bayesian modeling greatly reduces the need for the approximations made for confidence intervals and distributions of test statistics, and Bayesian methods formalize the still rather ad hoc frequentist approach to penalized maximum likelihood estimation by using skeptical prior distributions to obtain well-defined posterior distributions that automatically deal with shrinkage. The Bayesian approach also provides a formal mechanism for incorporating information external to the data. Although Bayesian methods are beyond the scope of this text, the text is Bayesian in spirit by emphasizing the careful use of subject matter expertise while building statistical models.

The text emphasizes predictive modeling, but as discussed in Chapter 1, developing good predictions goes hand in hand with accurate estimation of effects and with hypothesis testing (when appropriate). Besides emphasis on multivariable modeling, the text includes a Chapter 17 introducing survival analysis and methods for analyzing various types of single and multiple events. This book does not provide examples of analyses of one common type of response variable, namely, cost and related measures of resource consumption. However, least squares modeling presented in Chapter 15.1, the robust rank-based methods presented in Chapters 13, 15, and 20, and the transform-both-sides regression models discussed in Chapter 16 are very applicable and robust for modeling economic outcomes. See [167] and [260] for example analyses of such dependent variables using, respectively, the Cox model and nonparametric additive regression. The central Web site for this book (see the Appendix) has much more material on the use of the Cox model for analyzing costs.

This text does not address some important study design issues that if not respected can doom a predictive modeling or estimation project to failure. See Laupacis, Sekar, and Stiell [378] for a list of some of these issues.

Heavy use is made of the S language used by R. R is the focus because it is an elegant object-oriented system in which it is easy to implement new statistical ideas. Many R users around the world have done so, and their work has benefited many of the procedures described here. R also has a uniform syntax for specifying statistical models (with respect to categorical predictors, interactions, etc.), no matter which type of model is being fitted [96].

The free, open-source statistical software system R has been adopted by analysts and research statisticians worldwide. Its capabilities are growing exponentially because of the involvement of an ever-growing community of statisticians who are adding new tools to the base R system through contributed packages. All of the functions used in this text are available in R. See the book's Web site for updated information about software availability.

Readers who don't use R or any other statistical software environment will still find the statistical methods and case studies in this text useful, and it is hoped that the code that is presented will make the statistical methods more

concrete. At the very least, the code demonstrates that all of the methods presented in the text are feasible.

This text does not teach analysts how to use R. For that, the reader may wish to see reading recommendations on www.r-project.org as well as Venables and Ripley [635] (which is also an excellent companion to this text) and the many other excellent texts on R. See the Appendix for more information.

In addition to powerful features that are built into R, this text uses a package of freely available R functions called `rms` written by the author. `rms` tracks modeling details related to the expanded X or design matrix. It is a series of over 200 functions for model fitting, testing, estimation, validation, graphics, prediction, and typesetting by storing enhanced model design attributes in the fit. `rms` includes functions for least squares and penalized least squares multiple regression modeling in addition to functions for binary and ordinal regression, generalized least squares for analyzing serial data, quantile regression, and survival analysis that are emphasized in this text. Other freely available miscellaneous R functions used in the text are found in the `Hmisc` package also written by the author. Functions in `Hmisc` include facilities for data reduction, imputation, power and sample size calculation, advanced table making, recoding variables, importing and inspecting data, and general graphics. Consult the Appendix for information on obtaining `Hmisc` and `rms`.

The author and his colleagues have written SAS macros for fitting restricted cubic splines and for other basic operations. See the Appendix for more information. It is unfair not to mention some excellent capabilities of other statistical packages such as Stata (which has also been extended to provide regression splines and other modeling tools), but the extendability and graphics of R makes it especially attractive for all aspects of the comprehensive modeling strategy presented in this book.

Portions of Chapters 4 and 20 were published as reference [269]. Some of Chapter 13 was published as reference [272].

The author may be contacted by electronic mail at f.harrell@vanderbilt.edu and would appreciate being informed of unclear points, errors, and omissions in this book. Suggestions for improvements and for future topics are also welcome. As described in the Web site, instructors may contact the author to obtain copies of quizzes and extra assignments (both with answers) related to much of the material in the earlier chapters, and to obtain full solutions (with graphical output) to the majority of assignments in the text.

Major changes since the first edition include the following:

1. Creation of a now mature R package, `rms`, that replaces and greatly extends the `Design` library used in the first edition
2. Conversion of all of the book's code to R
3. Conversion of the book source into `knitr` [677] reproducible documents
4. All code from the text is executable and is on the web site
5. Use of color graphics and use of the `ggplot2` graphics package [667]
6. Scanned images were re-drawn

7. New text about problems with dichotomization of continuous variables and with classification (as opposed to prediction)
8. Expanded material on multiple imputation and predictive mean matching and emphasis on multiple imputation (using the Hmisc `aregImpute` function) instead of single imputation
9. Addition of redundancy analysis
10. Added a new section in Chapter 5 on bootstrap confidence intervals for rankings of predictors
11. Replacement of the U.S. presidential election data with analyses of a new diabetes dataset from NHANES using ordinal and quantile regression
12. More emphasis on semiparametric ordinal regression models for continuous Y , as direct competitors of ordinary multiple regression, with a detailed case study
13. A new chapter on generalized least squares for analysis of serial response data
14. The case study in imputation and data reduction was completely reworked and now focuses only on data reduction, with the addition of sparse principal components
15. More information about indexes of predictive accuracy
16. Augmentation of the chapter on maximum likelihood to include more flexible ways of testing contrasts as well as new methods for obtaining simultaneous confidence intervals
17. Binary logistic regression case study 1 was completely re-worked, now providing examples of model selection and model approximation accuracy
18. Single imputation was dropped from binary logistic case study 2
19. The case study in transform-both-sides regression modeling has been reworked using simulated data where true transformations are known, and a new example of the smearing estimator was added
20. Addition of 225 references, most of them published 2001–2014
21. New guidance on minimum sample sizes needed by some of the models
22. De-emphasis of bootstrap bumping [610] for obtaining simultaneous confidence regions, in favor of a general multiplicity approach [307].

Acknowledgments

A good deal of the writing of the first edition of this book was done during my 17 years on the faculty of Duke University. I wish to thank my close colleague Kerry Lee for providing many valuable ideas, fruitful collaborations, and well-organized lecture notes from which I have greatly benefited over the past years. Terry Therneau of Mayo Clinic has given me many of his wonderful ideas for many years, and has written state-of-the-art R software for survival analysis that forms the core of survival analysis software in my `rms` package. Michael Symons of the Department of Biostatistics of the University of North

Carolina at Chapel Hill and Timothy Morgan of the Division of Public Health Sciences at Wake Forest University School of Medicine also provided course materials, some of which motivated portions of this text. My former clinical colleagues in the Cardiology Division at Duke University, Robert Califf, Phillip Harris, Mark Hlatky, Dan Mark, David Pryor, and Robert Rosati, for many years provided valuable motivation, feedback, and ideas through our interaction on clinical problems. Besides Kerry Lee, statistical colleagues L. Richard Smith, Lawrence Muhlbaier, and Elizabeth DeLong clarified my thinking and gave me new ideas on numerous occasions. Charlotte Nelson and Carlos Alzola frequently helped me debug S routines when they thought they were just analyzing data.

Former students Bercedis Peterson, James Herndon, Robert McMahon, and Yuan-Li Shen have provided many insights into logistic and survival modeling. Associations with Doug Wagner and William Knaus of the University of Virginia, Ken Offord of Mayo Clinic, David Naftel of the University of Alabama in Birmingham, Phil Miller of Washington University, and Phil Goodman of the University of Nevada Reno have provided many valuable ideas and motivations for this work, as have Michael Schemper of Vienna University, Janez Stare of Ljubljana University, Slovenia, Ewout Steyerberg of Erasmus University, Rotterdam, Karel Moons of Utrecht University, and Drew Levy of Genentech. Richard Goldstein, along with several anonymous reviewers, provided many helpful criticisms of a previous version of this manuscript that resulted in significant improvements, and critical reading by Bob Edson (VA Cooperative Studies Program, Palo Alto) resulted in many error corrections. Thanks to Brian Ripley of the University of Oxford for providing many helpful software tools and statistical insights that greatly aided in the production of this book, and to Bill Venables of CSIRO Australia for wisdom, both statistical and otherwise. This work would also not have been possible without the S environment developed by Rick Becker, John Chambers, Allan Wilks, and the R language developed by Ross Ihaka and Robert Gentleman.

Work for the second edition was done in the excellent academic environment of Vanderbilt University, where biostatistical and biomedical colleagues and graduate students provided new insights and stimulating discussions. Thanks to Nick Cox, Durham University, UK, who provided from his careful reading of the first edition a very large number of improvements and corrections that were incorporated into the second. Four anonymous reviewers of the second edition also made numerous suggestions that improved the text.

Nashville, TN, USA
July 2015

Frank E. Harrell, Jr.

Contents

Typographical Conventions	xxv
1 Introduction	1
1.1 Hypothesis Testing, Estimation, and Prediction	1
1.2 Examples of Uses of Predictive Multivariable Modeling	3
1.3 Prediction vs. Classification	4
1.4 Planning for Modeling	6
1.4.1 Emphasizing Continuous Variables	8
1.5 Choice of the Model	8
1.6 Further Reading	11
2 General Aspects of Fitting Regression Models	13
2.1 Notation for Multivariable Regression Models	13
2.2 Model Formulations	14
2.3 Interpreting Model Parameters	15
2.3.1 Nominal Predictors	16
2.3.2 Interactions	16
2.3.3 Example: Inference for a Simple Model	17
2.4 Relaxing Linearity Assumption for Continuous Predictors ..	18
2.4.1 Avoiding Categorization	18
2.4.2 Simple Nonlinear Terms	21
2.4.3 Splines for Estimating Shape of Regression Function and Determining Predictor Transformations	22
2.4.4 Cubic Spline Functions	23
2.4.5 Restricted Cubic Splines	24
2.4.6 Choosing Number and Position of Knots	26
2.4.7 Nonparametric Regression	28
2.4.8 Advantages of Regression Splines over Other Methods	30

2.5	Recursive Partitioning: Tree-Based Models	30
2.6	Multiple Degree of Freedom Tests of Association	31
2.7	Assessment of Model Fit	33
2.7.1	Regression Assumptions	33
2.7.2	Modeling and Testing Complex Interactions	36
2.7.3	Fitting Ordinal Predictors	38
2.7.4	Distributional Assumptions	39
2.8	Further Reading	40
2.9	Problems	42
3	Missing Data	45
3.1	Types of Missing Data	45
3.2	Prelude to Modeling	46
3.3	Missing Values for Different Types of Response Variables	47
3.4	Problems with Simple Alternatives to Imputation	47
3.5	Strategies for Developing an Imputation Model	49
3.6	Single Conditional Mean Imputation	52
3.7	Predictive Mean Matching	52
3.8	Multiple Imputation	53
3.8.1	The <code>aregImpute</code> and Other Chained Equations Approaches	55
3.9	Diagnostics	56
3.10	Summary and Rough Guidelines	56
3.11	Further Reading	58
3.12	Problems	59
4	Multivariable Modeling Strategies	63
4.1	Prespecification of Predictor Complexity Without Later Simplification	64
4.2	Checking Assumptions of Multiple Predictors Simultaneously	67
4.3	Variable Selection	67
4.4	Sample Size, Overfitting, and Limits on Number of Predictors	72
4.5	Shrinkage	75
4.6	Collinearity	78
4.7	Data Reduction	79
4.7.1	Redundancy Analysis	80
4.7.2	Variable Clustering	81
4.7.3	Transformation and Scaling Variables Without Using Y	81
4.7.4	Simultaneous Transformation and Imputation	83
4.7.5	Simple Scoring of Variable Clusters	85
4.7.6	Simplifying Cluster Scores	87
4.7.7	How Much Data Reduction Is Necessary?	87

4.8	Other Approaches to Predictive Modeling	89
4.9	Overly Influential Observations	90
4.10	Comparing Two Models	92
4.11	Improving the Practice of Multivariable Prediction	94
4.12	Summary: Possible Modeling Strategies	94
4.12.1	Developing Predictive Models	95
4.12.2	Developing Models for Effect Estimation	98
4.12.3	Developing Models for Hypothesis Testing	99
4.13	Further Reading	100
4.14	Problems	102
5	Describing, Resampling, Validating, and Simplifying the Model	103
5.1	Describing the Fitted Model	103
5.1.1	Interpreting Effects	103
5.1.2	Indexes of Model Performance	104
5.2	The Bootstrap	106
5.3	Model Validation	109
5.3.1	Introduction	109
5.3.2	Which Quantities Should Be Used in Validation?	110
5.3.3	Data-Splitting	111
5.3.4	Improvements on Data-Splitting: Resampling	112
5.3.5	Validation Using the Bootstrap	114
5.4	Bootstrapping Ranks of Predictors	117
5.5	Simplifying the Final Model by Approximating It	118
5.5.1	Difficulties Using Full Models	118
5.5.2	Approximating the Full Model	119
5.6	Further Reading	121
5.7	Problem	124
6	R Software	127
6.1	The R Modeling Language	128
6.2	User-Contributed Functions	129
6.3	The <code>rms</code> Package	130
6.4	Other Functions	141
6.5	Further Reading	142
7	Modeling Longitudinal Responses using Generalized Least Squares	143
7.1	Notation and Data Setup	143
7.2	Model Specification for Effects on $E(Y)$	144
7.3	Modeling Within-Subject Dependence	144
7.4	Parameter Estimation Procedure	147
7.5	Common Correlation Structures	147
7.6	Checking Model Fit	148

7.7	Sample Size Considerations	148
7.8	R Software	149
7.9	Case Study	149
7.9.1	Graphical Exploration of Data	150
7.9.2	Using Generalized Least Squares	151
7.10	Further Reading	158
8	Case Study in Data Reduction	161
8.1	Data	161
8.2	How Many Parameters Can Be Estimated?	164
8.3	Redundancy Analysis	164
8.4	Variable Clustering	166
8.5	Transformation and Single Imputation Using <code>transcan</code>	167
8.6	Data Reduction Using Principal Components	170
8.6.1	Sparse Principal Components	175
8.7	Transformation Using Nonparametric Smoothers	176
8.8	Further Reading	177
8.9	Problems	178
9	Overview of Maximum Likelihood Estimation	181
9.1	General Notions—Simple Cases	181
9.2	Hypothesis Tests	185
9.2.1	Likelihood Ratio Test	185
9.2.2	Wald Test	186
9.2.3	Score Test	186
9.2.4	Normal Distribution—One Sample	187
9.3	General Case	188
9.3.1	Global Test Statistics	189
9.3.2	Testing a Subset of the Parameters	190
9.3.3	Tests Based on Contrasts	192
9.3.4	Which Test Statistics to Use When	193
9.3.5	Example: Binomial—Comparing Two Proportions	194
9.4	Iterative ML Estimation	195
9.5	Robust Estimation of the Covariance Matrix	196
9.6	Wald, Score, and Likelihood-Based Confidence Intervals	198
9.6.1	Simultaneous Wald Confidence Regions	199
9.7	Bootstrap Confidence Regions	199
9.8	Further Use of the Log Likelihood	203
9.8.1	Rating Two Models, Penalizing for Complexity	203
9.8.2	Testing Whether One Model Is Better than Another	204
9.8.3	Unitless Index of Predictive Ability	205
9.8.4	Unitless Index of Adequacy of a Subset of Predictors	207
9.9	Weighted Maximum Likelihood Estimation	208
9.10	Penalized Maximum Likelihood Estimation	209

9.11	Further Reading	213
9.12	Problems	216
10	Binary Logistic Regression	219
10.1	Model	219
10.1.1	Model Assumptions and Interpretation of Parameters	221
10.1.2	Odds Ratio, Risk Ratio, and Risk Difference	224
10.1.3	Detailed Example	225
10.1.4	Design Formulations	230
10.2	Estimation	231
10.2.1	Maximum Likelihood Estimates	231
10.2.2	Estimation of Odds Ratios and Probabilities	232
10.2.3	Minimum Sample Size Requirement	233
10.3	Test Statistics	234
10.4	Residuals	235
10.5	Assessment of Model Fit	236
10.6	Collinearity	255
10.7	Overly Influential Observations	255
10.8	Quantifying Predictive Ability	256
10.9	Validating the Fitted Model	259
10.10	Describing the Fitted Model	264
10.11	R Functions	269
10.12	Further Reading	271
10.13	Problems	273
11	Binary Logistic Regression Case Study 1	275
11.1	Overview	275
11.2	Background	275
11.3	Data Transformations and Single Imputation	276
11.4	Regression on Original Variables, Principal Components and Pretransformations	277
11.5	Description of Fitted Model	278
11.6	Backwards Step-Down	280
11.7	Model Approximation	287
12	Logistic Model Case Study 2: Survival of Titanic Passengers	291
12.1	Descriptive Statistics	291
12.2	Exploring Trends with Nonparametric Regression	294
12.3	Binary Logistic Model With Casewise Deletion of Missing Values	296
12.4	Examining Missing Data Patterns	302
12.5	Multiple Imputation	304
12.6	Summarizing the Fitted Model	307

13 Ordinal Logistic Regression	311
13.1 Background	311
13.2 Ordinality Assumption	312
13.3 Proportional Odds Model	313
13.3.1 Model	313
13.3.2 Assumptions and Interpretation of Parameters	313
13.3.3 Estimation	314
13.3.4 Residuals	314
13.3.5 Assessment of Model Fit	315
13.3.6 Quantifying Predictive Ability	318
13.3.7 Describing the Fitted Model	318
13.3.8 Validating the Fitted Model	318
13.3.9 R Functions	319
13.4 Continuation Ratio Model	319
13.4.1 Model	319
13.4.2 Assumptions and Interpretation of Parameters	320
13.4.3 Estimation	320
13.4.4 Residuals	321
13.4.5 Assessment of Model Fit	321
13.4.6 Extended CR Model	321
13.4.7 Role of Penalization in Extended CR Model	322
13.4.8 Validating the Fitted Model	322
13.4.9 R Functions	323
13.5 Further Reading	324
13.6 Problems	324
14 Case Study in Ordinal Regression, Data Reduction, and Penalization	327
14.1 Response Variable	328
14.2 Variable Clustering	329
14.3 Developing Cluster Summary Scores	330
14.4 Assessing Ordinality of Y for each X , and Unadjusted Checking of PO and CR Assumptions	333
14.5 A Tentative Full Proportional Odds Model	333
14.6 Residual Plots	336
14.7 Graphical Assessment of Fit of CR Model	338
14.8 Extended Continuation Ratio Model	340
14.9 Penalized Estimation	342
14.10 Using Approximations to Simplify the Model	348
14.11 Validating the Model	353
14.12 Summary	355
14.13 Further Reading	356
14.14 Problems	357

15	Regression Models for Continuous Y and Case Study in Ordinal Regression	359
15.1	The Linear Model	359
15.2	Quantile Regression.....	360
15.3	Ordinal Regression Models for Continuous Y	361
15.3.1	Minimum Sample Size Requirement	363
15.4	Comparison of Assumptions of Various Models	364
15.5	Dataset and Descriptive Statistics	365
15.5.1	Checking Assumptions of OLS and Other Models... ..	368
15.6	Ordinal Regression Applied to HbA_{1c}	370
15.6.1	Checking Fit for Various Models Using Age.....	370
15.6.2	Examination of BMI.....	374
15.6.3	Consideration of All Body Size Measurements.....	375
16	Transform-Both-Sides Regression	389
16.1	Background.....	389
16.2	Generalized Additive Models.....	390
16.3	Nonparametric Estimation of Y -Transformation	390
16.4	Obtaining Estimates on the Original Scale.....	391
16.5	R Functions.....	392
16.6	Case Study	393
17	Introduction to Survival Analysis	399
17.1	Background.....	399
17.2	Censoring, Delayed Entry, and Truncation	401
17.3	Notation, Survival, and Hazard Functions	402
17.4	Homogeneous Failure Time Distributions	407
17.5	Nonparametric Estimation of S and Λ	409
17.5.1	Kaplan–Meier Estimator	409
17.5.2	Altschuler–Nelson Estimator	413
17.6	Analysis of Multiple Endpoints.....	413
17.6.1	Competing Risks	414
17.6.2	Competing Dependent Risks	414
17.6.3	State Transitions and Multiple Types of Nonfatal Events	416
17.6.4	Joint Analysis of Time and Severity of an Event....	417
17.6.5	Analysis of Multiple Events	417
17.7	R Functions.....	418
17.8	Further Reading.....	420
17.9	Problems	421
18	Parametric Survival Models	423
18.1	Homogeneous Models (No Predictors).....	423
18.1.1	Specific Models	423
18.1.2	Estimation	424
18.1.3	Assessment of Model Fit	426

18.2	Parametric Proportional Hazards Models	427
18.2.1	Model	427
18.2.2	Model Assumptions and Interpretation of Parameters	428
18.2.3	Hazard Ratio, Risk Ratio, and Risk Difference	430
18.2.4	Specific Models	431
18.2.5	Estimation	432
18.2.6	Assessment of Model Fit	434
18.3	Accelerated Failure Time Models	436
18.3.1	Model	436
18.3.2	Model Assumptions and Interpretation of Parameters	436
18.3.3	Specific Models	437
18.3.4	Estimation	438
18.3.5	Residuals	440
18.3.6	Assessment of Model Fit	440
18.3.7	Validating the Fitted Model	446
18.4	Buckley–James Regression Model	447
18.5	Design Formulations	447
18.6	Test Statistics	447
18.7	Quantifying Predictive Ability	447
18.8	Time-Dependent Covariates	447
18.9	R Functions	448
18.10	Further Reading	450
18.11	Problems	451
19	Case Study in Parametric Survival Modeling and Model Approximation	453
19.1	Descriptive Statistics	453
19.2	Checking Adequacy of Log-Normal Accelerated Failure Time Model	458
19.3	Summarizing the Fitted Model	466
19.4	Internal Validation of the Fitted Model Using the Bootstrap	466
19.5	Approximating the Full Model	469
19.6	Problems	473
20	Cox Proportional Hazards Regression Model	475
20.1	Model	475
20.1.1	Preliminaries	475
20.1.2	Model Definition	476
20.1.3	Estimation of β	476
20.1.4	Model Assumptions and Interpretation of Parameters	478
20.1.5	Example	478

20.1.6	Design Formulations	480
20.1.7	Extending the Model by Stratification	481
20.2	Estimation of Survival Probability and Secondary Parameters	483
20.3	Sample Size Considerations	486
20.4	Test Statistics	486
20.5	Residuals	487
20.6	Assessment of Model Fit	487
20.6.1	Regression Assumptions	487
20.6.2	Proportional Hazards Assumption	494
20.7	What to Do When PH Fails	501
20.8	Collinearity	503
20.9	Overly Influential Observations	504
20.10	Quantifying Predictive Ability	504
20.11	Validating the Fitted Model	506
20.11.1	Validation of Model Calibration	506
20.11.2	Validation of Discrimination and Other Statistical Indexes	507
20.12	Describing the Fitted Model	509
20.13	R Functions	513
20.14	Further Reading	517
21	Case Study in Cox Regression	521
21.1	Choosing the Number of Parameters and Fitting the Model	521
21.2	Checking Proportional Hazards	525
21.3	Testing Interactions	527
21.4	Describing Predictor Effects	527
21.5	Validating the Model	529
21.6	Presenting the Model	530
21.7	Problems	531
A	Datasets, R Packages, and Internet Resources	535
	References	539
	Index	571

Typographical Conventions

Boxed numbers in the margins such as 1 correspond to numbers at the end of chapters in sections named “Further Reading.” Bracketed numbers and numeric superscripts in the text refer to the bibliography, while alphabetic superscripts indicate footnotes.

R language commands and names of R functions and packages are set in `typewriter font`, as are most variable names.

R code blocks are set off with a shadowbox, and R output that is not directly using `LATEX` appears in a box that is framed on three sides.

In the S language upon which R is based, $x \leftarrow y$ is read “ x gets the value of y .” The assignment operator \leftarrow , used in the text for aesthetic reasons (as are \leq and \geq), is entered by the user as `<-`. Comments begin with `#`, subscripts use brackets (`[]`), and the missing value is denoted by `NA` (not available).

In ordinary text and mathematical expressions, `[logical variable]` and `[logical expression]` imply a value of 1 if the logical variable or expression is true, and 0 otherwise.

Chapter 1

Introduction

1.1 Hypothesis Testing, Estimation, and Prediction

Statistics comprises among other areas study design, hypothesis testing, estimation, and prediction. This text aims at the last area, by presenting methods that enable an analyst to develop models that will make accurate predictions of responses for *future* observations. Prediction could be considered a superset of hypothesis testing and estimation, so the methods presented here will also assist the analyst in those areas. It is worth pausing to explain how this is so.

In traditional hypothesis testing one often chooses a *null hypothesis* defined as the absence of some effect. For example, in testing whether a variable such as cholesterol is a risk factor for sudden death, one might test the null hypothesis that an increase in cholesterol does not increase the risk of death. Hypothesis testing can easily be done within the context of a statistical model, but a model is not required. When one only wishes to assess whether an effect is zero, P -values may be computed using permutation or rank (non-parametric) tests while making only minimal assumptions. But there are still reasons for preferring a model-based approach over techniques that only yield P -values.

1. Permutation and rank tests do not easily give rise to estimates of *magnitudes* of effects.
2. These tests cannot be readily extended to incorporate complexities such as cluster sampling or repeated measurements within subjects.
3. Once the analyst is familiar with a model, that model may be used to carry out many different statistical tests; there is no need to learn specific formulas to handle the special cases. The two-sample t -test is a special case of the ordinary multiple regression model having as its sole X variable a dummy variable indicating group membership. The Wilcoxon-Mann-Whitney test is a special case of the proportional odds ordinal logistic

model.⁶⁶⁴ The analysis of variance (multiple group) test and the Kruskal–Wallis test can easily be obtained from these two regression models by using more than one dummy predictor variable.

Even without complexities such as repeated measurements, problems can arise when many hypotheses are to be tested. Testing too many hypotheses is related to fitting too many predictors in a regression model. One commonly hears the statement that “the dataset was too small to allow modeling, so we just did hypothesis tests.” It is unlikely that the resulting inferences would be reliable. If the sample size is insufficient for modeling it is often insufficient for tests or estimation. This is especially true when one desires to publish an estimate of the effect corresponding to the hypothesis yielding the smallest P -value. Ordinary point estimates are known to be badly biased when the quantity to be estimated was determined by “data dredging.” This can be remedied by the same kind of shrinkage used in multivariable modeling (Section 9.10).

Statistical estimation is usually model-based. For example, one might use a survival regression model to estimate the relative effect of increasing cholesterol from 200 to 250 mg/dl on the hazard of death. Variables other than cholesterol may also be in the regression model, to allow estimation of the effect of increasing cholesterol, holding other risk factors constant. But accurate estimation of the cholesterol effect will depend on how cholesterol as well as each of the adjustment variables is assumed to relate to the hazard of death. If linear relationships are incorrectly assumed, estimates will be inaccurate. Accurate estimation also depends on avoiding overfitting the adjustment variables. If the dataset contains 200 subjects, 30 of whom died, and if one adjusted for 15 “confounding” variables, the estimates would be “over-adjusted” for the effects of the 15 variables, as some of their apparent effects would actually result from spurious associations with the response variable (time until death). The overadjustment would reduce the cholesterol effect. The resulting unreliability of estimates equals the degree to which the overall model fails to validate on an independent sample.

It is often useful to think of effect estimates as differences between two predicted values from a model. This way, one can account for nonlinearities and interactions. For example, if cholesterol is represented nonlinearly in a logistic regression model, predicted values on the “linear combination of X ’s scale” are predicted log odds of an event. The increase in log odds from raising cholesterol from 200 to 250 mg/dl is the difference in predicted values, where cholesterol is set to 250 and then to 200, and all other variables are held constant. The point estimate of the 250:200 mg/dl odds ratio is the anti-log of this difference. If cholesterol is represented nonlinearly in the model, it does not matter how many terms in the model involve cholesterol as long as the overall predicted values are obtained.

Thus when one develops a reasonable multivariable predictive model, hypothesis testing and estimation of effects are byproducts of the fitted model. So predictive modeling is often desirable even when prediction is not the main goal.

1.2 Examples of Uses of Predictive Multivariable Modeling

There is an endless variety of uses for multivariable models. Predictive models have long been used in business to forecast financial performance and to model consumer purchasing and loan pay-back behavior. In ecology, regression models are used to predict the probability that a fish species will disappear from a lake. Survival models have been used to predict product life (e.g., time to burn-out of an mechanical part, time until saturation of a disposable diaper). Models are commonly used in discrimination litigation in an attempt to determine whether race or sex is used as the basis for hiring or promotion, after taking other personnel characteristics into account.

Multivariable models are used extensively in medicine, epidemiology, biostatistics, health services research, pharmaceutical research, and related fields. The author has worked primarily in these fields, so most of the examples in this text come from those areas. In medicine, two of the major areas of application are diagnosis and prognosis. There models are used to predict the probability that a certain type of patient will be shown to have a specific disease, or to predict the time course of an already diagnosed disease. In observational studies in which one desires to compare patient outcomes between two or more treatments, multivariable modeling is very important because of the biases caused by nonrandom treatment assignment. Here the simultaneous effects of several uncontrolled variables must be controlled (held constant mathematically if using a regression model) so that the effect of the factor of interest can be more purely estimated. A newer technique for more aggressively adjusting for nonrandom treatment assignment, the *propensity score*,^{116,530} provides yet another opportunity for multivariable modeling (see Section 10.1.4). The propensity score is merely the predicted value from a multivariable model where the response variable is the exposure or the treatment actually used. The estimated propensity score is then used in a second step as an adjustment variable in the model for the response of interest.

It is not widely recognized that multivariable modeling is extremely valuable even in well-designed randomized experiments. Such studies are often designed to make *relative* comparisons of two or more treatments, using odds ratios, hazard ratios, and other measures of relative effects. But to be able to estimate *absolute* effects one must develop a multivariable model of the response variable. This model can predict, for example, the probability that a patient on treatment A with characteristics X will survive five years, or it can

predict the life expectancy for this patient. By making the same prediction for a patient on treatment B with the same characteristics, one can estimate the absolute difference in probabilities or life expectancies. This approach recognizes that low-risk patients must have less absolute benefit of treatment (lower change in outcome probability) than high-risk patients,³⁵¹ a fact that has been ignored in many clinical trials. Another reason for multivariable modeling in randomized clinical trials is that when the basic response model is nonlinear (e.g., logistic, Cox, parametric survival models), the unadjusted estimate of the treatment effect is not correct if there is moderate heterogeneity of subjects, even with perfect balance of baseline characteristics across the treatment groups.^{a9, 24, 198, 588} So even when investigators are interested in simple comparisons of two groups' responses, multivariable modeling can be advantageous and sometimes mandatory.

Cost-effectiveness analysis is becoming increasingly used in health care research, and the "effectiveness" (denominator of the cost-effectiveness ratio) is always a measure of absolute effectiveness. As absolute effectiveness varies dramatically with the risk profiles of subjects, it must be estimated for individual subjects using a multivariable model^{90, 344}.

1.3 Prediction vs. Classification

For problems ranging from bioinformatics to marketing, many analysts desire to develop "classifiers" instead of developing predictive models. Consider an optimum case for classifier development, in which the response variable is binary, the two levels represent a sharp dichotomy with no gray zone (e.g., complete success vs. total failure with no possibility of a partial success), the user of the classifier is forced to make one of the two choices, the cost of misclassification is the same for every future observation, and the ratio of the cost of a false positive to that of a false negative equals the (often hidden) ratio implied by the analyst's classification rule. Even if all of those conditions are met, classification is still inferior to probability modeling for driving the development of a predictive instrument or for estimation or hypothesis testing. It is far better to use the full information in the data to develop a probability model, then develop classification rules on the basis of estimated probabilities. At the least, this forces the analyst to use a proper accuracy score²¹⁹ in finding or weighting data features.

When the dependent variable is ordinal or continuous, classification through forced up-front dichotomization in an attempt to simplify the problem results in arbitrariness and major information loss even when the optimum cut point

^a For example, unadjusted odds ratios from 2×2 tables are different from adjusted odds ratios when there is variation in subjects' risk factors within each treatment group, even when the distribution of the risk factors is identical between the two groups.

(the median) is used. Dichotomizing the outcome at a different point may require a many-fold increase in sample size to make up for the lost information¹⁸⁷. In the area of medical diagnosis, it is often the case that the disease is really on a continuum, and predicting the severity of disease (rather than just its presence or absence) will greatly increase power and precision, not to mention making the result less arbitrary.

It is important to note that two-group classification represents an artificial forced choice. It is not often the case that the user of the classifier needs to be limited to two possible actions. The best option for many subjects may be to refuse to make a decision or to obtain more data (e.g., order another medical diagnostic test). A gray zone can be helpful, and predictions include gray zones automatically.

Unlike prediction (e.g., of absolute risk), classification implicitly uses utility functions (also called loss or cost functions, e.g., cost of a false positive classification). Implicit utility functions are highly problematic. First, it is well known that the utility function depends on variables that are not predictive of outcome and are not collected (e.g., subjects' preferences) that are available only at the decision point. Second, the approach assumes every subject has the same utility function^b. Third, the analyst presumptuously assumes that the subject's utility coincides with his own.

Formal decision analysis uses subject-specific utilities and optimum predictions based on all available data^{62, 74, 183, 210, 219, 642c}. It follows that receiver

^b Simple examples to the contrary are the less weight given to a false negative diagnosis of cancer in the elderly and the aversion of some subjects to surgery or chemotherapy.

^c To make an optimal decision you need to know all relevant data about an individual (used to estimate the probability of an outcome), and the utility (cost, loss function) of making each decision. Sensitivity and specificity do not provide this information. For example, if one estimated that the probability of a disease given age, sex, and symptoms is 0.1 and the "cost" of a false positive equaled the "cost" of a false negative, one would act as if the person does not have the disease. Given other utilities, one would make different decisions. If the utilities are unknown, one gives the best estimate of the probability of the outcome to the decision maker and let her incorporate her own unspoken utilities in making an optimum decision for her.

Besides the fact that cutoffs that are not individualized do not apply to individuals, only to groups, individual decision making does not utilize sensitivity and specificity. For an individual we can compute $\text{Prob}(Y = 1|X = x)$; we don't care about $\text{Prob}(Y = 1|X > c)$, and an individual having $X = x$ would be quite puzzled if she were given $\text{Prob}(X > c|\text{future unknown } Y)$ when she already knows $X = x$ so X is no longer a random variable.

Even when group decision making is needed, sensitivity and specificity can be bypassed. For mass marketing, for example, one can rank order individuals by the estimated probability of buying the product, to create a lift curve. This is then used to target the k most likely buyers where k is chosen to meet total program cost constraints.

operating characteristic curve (ROC^d) analysis is misleading except for the special case of mass one-time group decision making with unknown utilities (e.g., launching a flu vaccination program).

1

An analyst's goal should be the development of the most accurate and reliable predictive model or the best model on which to base estimation or hypothesis testing. In the vast majority of cases, classification is the task of the user of the predictive model, at the point in which utilities (costs) and preferences are known.

1.4 Planning for Modeling

When undertaking the development of a model to predict a response, one of the first questions the researcher must ask is “will this model actually be used?” Many models are never used, for several reasons⁵²² including: (1) it was not deemed relevant to make predictions in the setting envisioned by the authors; (2) potential users of the model did not trust the relationships, weights, or variables used to make the predictions; and (3) the variables necessary to make the predictions were not routinely available.

Once the researcher convinces herself that a predictive model is worth developing, there are many study design issues to be addressed.^{18, 378} Models are often developed using a “convenience sample,” that is, a dataset that was not collected with such predictions in mind. The resulting models are often fraught with difficulties such as the following.

1. The most important predictor or response variables may not have been collected, tempting the researchers to make do with variables that do not capture the real underlying processes.
2. The subjects appearing in the dataset are ill-defined, or they are not representative of the population for which inferences are to be drawn; similarly, the data collection sites may not represent the kind of variation in the population of sites.
3. Key variables are missing in large numbers of subjects.
4. Data are not missing at random; for example, data may not have been collected on subjects who dropped out of a study early, or on patients who were too sick to be interviewed.
5. Operational definitions of some of the key variables were never made.
6. Observer variability studies may not have been done, so that the reliability of measurements is unknown, or there are other kinds of important measurement errors.

A predictive model will be more accurate, as well as useful, when data collection is planned prospectively. That way one can design data collection

^d The ROC curve is a plot of sensitivity vs. one minus specificity as one varies a cutoff on a continuous predictor used to make a decision.

instruments containing the necessary variables, and all terms can be given standard definitions (for both descriptive and response variables) for use at all data collection sites. Also, steps can be taken to minimize the amount of missing data.

In the context of describing and modeling health outcomes, Iezzoni³¹⁷ has an excellent discussion of the dimensions of risk that should be captured by variables included in the model. She lists these general areas that should be quantified by predictor variables:

1. age,
2. sex,
3. acute clinical stability,
4. principal diagnosis,
5. severity of principal diagnosis,
6. extent and severity of comorbidities,
7. physical functional status,
8. psychological, cognitive, and psychosocial functioning,
9. cultural, ethnic, and socioeconomic attributes and behaviors,
10. health status and quality of life, and
11. patient attitudes and preferences for outcomes.

Some baseline covariates to be sure to capture in general include

1. a baseline measurement of the response variable,
2. the subject's most recent status,
3. the subject's trajectory as of time zero or past levels of a key variable,
4. variables explaining much of the variation in the response, and
5. more subtle predictors whose distributions strongly differ between the levels of a key variable of interest in an observational study.

Many things can go wrong in statistical modeling, including the following.

1. The process generating the data is not stable.
2. The model is misspecified with regard to nonlinearities or interactions, or there are predictors missing.
3. The model is misspecified in terms of the transformation of the response variable or the model's distributional assumptions.
4. The model contains discontinuities (e.g., by categorizing continuous predictors or fitting regression shapes with sudden changes) that can be gamed by users.
5. Correlations among subjects are not specified, or the correlation structure is misspecified, resulting in inefficient parameter estimates and overconfident inference.
6. The model is overfitted, resulting in predictions that are too extreme or positive associations that are false.

7. The user of the model relies on predictions obtained by extrapolating to combinations of predictor values well outside the range of the dataset used to develop the model.
8. Accurate and discriminating predictions can lead to behavior changes that make future predictions inaccurate.

1.4.1 *Emphasizing Continuous Variables*

When designing the data collection it is important to emphasize the use of continuous variables over categorical ones. Some categorical variables are subjective and hard to standardize, and on the average they do not contain the same amount of statistical information as continuous variables. Above all, it is unwise to categorize naturally continuous variables during data collection,^e as the original values can then not be recovered, and if another researcher feels that the (arbitrary) cutoff values were incorrect, other cutoffs cannot be substituted. Many researchers make the mistake of assuming that categorizing a continuous variable will result in less measurement error. This is a false assumption, for if a subject is placed in the wrong interval this will be as much as a 100% error. Thus the magnitude of the error multiplied by the probability of an error is no better with categorization.

2

1.5 Choice of the Model

The actual method by which an underlying statistical model should be chosen by the analyst is not well developed. A. P. Dawid is quoted in Lehmann³⁹⁷ as saying the following.

Where do probability models come from? To judge by the resounding silence over this question on the part of most statisticians, it seems highly embarrassing. In general, the theoretician is happy to accept that his abstract probability triple (Ω, A, P) was found under a gooseberry bush, while the applied statistician's model "just grew".

3

In biostatistics, epidemiology, economics, psychology, sociology, and many other fields it is seldom the case that subject matter knowledge exists that would allow the analyst to pre-specify a model (e.g., Weibull or log-normal survival model), a transformation for the response variable, and a structure

^e An exception may be sensitive variables such as income level. Subjects may be more willing to check a box corresponding to a wide interval containing their income. It is unlikely that a reduction in the probability that a subject will inflate her income will offset the loss of precision due to categorization of income, but there will be a decrease in the number of refusals. This reduction in missing data can more than offset the lack of precision.

for how predictors appear in the model (e.g., transformations, addition of nonlinear terms, interaction terms). Indeed, some authors question whether the notion of a true model even exists in many cases.¹⁰⁰ We are for better or worse forced to develop models empirically in the majority of cases. Fortunately, careful and objective validation of the accuracy of model predictions against observable responses can lend credence to a model, if a good validation is not merely the result of overfitting (see Section 5.3).

There are a few general guidelines that can help in choosing the basic form of the statistical model.

1. The model must use the data efficiently. If, for example, one were interested in predicting the probability that a patient with a specific set of characteristics would live five years from diagnosis, an inefficient model would be a binary logistic model. A more efficient method, and one that would also allow for losses to follow-up before five years, would be a semi-parametric (rank based) or parametric survival model. Such a model uses individual times of events in estimating coefficients, but it can easily be used to estimate the probability of surviving five years. As another example, if one were interested in predicting patients' quality of life on a scale of excellent, very good, good, fair, and poor, a polytomous (multinomial) categorical response model would not be efficient as it would not make use of the ordering of responses.
2. Choose a model that fits overall structures likely to be present in the data. In modeling survival time in chronic disease one might feel that the importance of most of the risk factors is constant over time. In that case, a proportional hazards model such as the Cox or Weibull model would be a good initial choice. If on the other hand one were studying acutely ill patients whose risk factors wane in importance as the patients survive longer, a model such as the log-normal or log-logistic regression model would be more appropriate.
3. Choose a model that is robust to problems in the data that are difficult to check. For example, the Cox proportional hazards model and ordinal logistic models are not affected by monotonic transformations of the response variable.
4. Choose a model whose mathematical form is appropriate for the response being modeled. This often has to do with minimizing the need for interaction terms that are included only to address a basic lack of fit. For example, many researchers have used ordinary linear regression models for binary responses, because of their simplicity. But such models allow predicted probabilities to be outside the interval $[0, 1]$, and strange interactions among the predictor variables are needed to make predictions remain in the legal range.
5. Choose a model that is readily extendible. The Cox model, by its use of stratification, easily allows a few of the predictors, especially if they are categorical, to violate the assumption of equal regression coefficients over

time (proportional hazards assumption). The continuation ratio ordinal logistic model can also be generalized easily to allow for varying coefficients of some of the predictors as one proceeds across categories of the response.

R. A. Fisher as quoted in Lehmann³⁹⁷ had these suggestions about model building: “(a) We must confine ourselves to those forms which we know how to handle,” and (b) “More or less elaborate forms will be suitable according to the volume of the data.” Ameen [100, p. 453] stated that a good model is “(a) satisfactory in performance relative to the stated objective, (b) logically sound, (c) representative, (d) questionable and subject to on-line interrogation, (e) able to accommodate external or expert information and (f) able to convey information.”

It is very typical to use the data to make decisions about the form of the model as well as about how predictors are represented in the model. Then, once a model is developed, the entire modeling process is routinely forgotten, and statistical quantities such as standard errors, confidence limits, P -values, and R^2 are computed as if the resulting model were entirely pre-specified. However, Faraway,¹⁸⁶ Draper,¹⁶³ Chatfield,¹⁰⁰ Buckland et al.⁸⁰ and others have written about the severe problems that result from treating an empirically derived model as if it were pre-specified and as if it were the correct model. As Chatfield states [100, p. 426]: “It is indeed strange that we often admit model uncertainty by searching for a best model but then ignore this uncertainty by making inferences and predictions as if certain that the best fitting model is actually true.”

Stepwise variable selection is one of the most widely used and abused of all data analysis techniques. Much is said about this technique later (see Section 4.3), but there are many other elements of model development that will need to be accounted for when making statistical inferences, and unfortunately it is difficult to derive quantities such as confidence limits that are properly adjusted for uncertainties such as the data-based choice between a Weibull and a log-normal regression model.

Ye⁶⁷⁸ developed a general method for estimating the “generalized degrees of freedom” (GDF) for any “data mining” or model selection procedure based on least squares. The GDF is an extremely useful index of the amount of “data dredging” or overfitting that has been done in a modeling process. It is also useful for estimating the residual variance with less bias. In one example, Ye developed a regression tree using recursive partitioning involving 10 candidate predictor variables on 100 observations. The resulting tree had 19 nodes and GDF of 76. The usual way of estimating the residual variance involves dividing the pooled within-node sum of squares by $100 - 19$, but Ye showed that dividing by $100 - 76$ instead yielded a much less biased (and much higher) estimate of σ^2 . In another example, Ye considered stepwise variable selection using 20 candidate predictors and 22 observations. When there is no true association between any of the predictors and the response, Ye found that $GDF = 14.1$ for a strategy that selected the best five-variable model.

4

5

Given that the choice of the model has been made (e.g., a log-normal model), penalized maximum likelihood estimation has major advantages in the battle between making the model fit adequately and avoiding overfitting (Sections 9.10 and 13.4.7). Penalization lessens the need for model selection.

1.6 Further Reading

- [1] Briggs and Zaretzki⁷⁴ eloquently state the problem with ROC curves and the areas under them (AUC):

Statistics such as the AUC are not especially relevant to someone who must make a decision about a particular x_c . . . ROC curves lack or obscure several quantities that are necessary for evaluating the operational effectiveness of diagnostic tests. . . ROC curves were first used to check how radio *receivers* (like radar receivers) operated over a range of frequencies. . . This is not how must ROC curves are used now, particularly in medicine. The receiver of a diagnostic measurement . . . wants to make a decision based on some x_c , and is not especially interested in how well he would have done had he used some different cutoff.

In the discussion to their paper, David Hand states

When integrating to yield the overall AUC measure, it is necessary to decide what weight to give each value in the integration. The AUC implicitly does this using a weighting derived empirically from the data. This is nonsensical. The relative importance of misclassifying a case as a noncase, compared to the reverse, cannot come from the data itself. It must come externally, from considerations of the severity one attaches to the different kinds of misclassifications.

AUC, only because it equals the concordance probability in the binary Y case, is still often useful as a predictive discrimination measure.

- [2] More severe problems caused by dichotomizing continuous variables are discussed in [13, 17, 45, 82, 185, 294, 379, 521, 597].
- [3] See the excellent editorial by Mallows⁴³⁴ for more about model choice. See Breiman and discussants⁶⁷ for an interesting debate about the use of data models vs. algorithms. This material also covers interpretability vs. predictive accuracy and several other topics.
- [4] See [15, 80, 100, 163, 186, 415] for information about accounting for model selection in making final inferences. Faraway¹⁸⁶ demonstrated that the bootstrap has good potential in related although somewhat simpler settings, and Buckland et al.⁸⁰ developed a promising bootstrap weighting method for accounting for model uncertainty.
- [5] Tibshirani and Knight⁶¹¹ developed another approach to estimating the generalized degrees of freedom. Luo et al.⁴³⁰ developed a way to add noise of known variance to the response variable to tune the stopping rule used for variable selection. Zou et al.⁶⁸⁹ showed that the lasso, an approach that simultaneously selects variables and shrinks coefficients, has a nice property. Since it uses penalization (shrinkage), an unbiased estimate of its effective number of degrees of freedom is the number of nonzero regression coefficients in the final model.

Chapter 2

General Aspects of Fitting Regression Models

2.1 Notation for Multivariable Regression Models

The ordinary multiple linear regression model is frequently used and has parameters that are easily interpreted. In this chapter we study a general class of regression models, those stated in terms of a weighted sum of a set of independent or predictor variables. It is shown that after linearizing the model with respect to the predictor variables, the parameters in such regression models are also readily interpreted. Also, all the designs used in ordinary linear regression can be used in this general setting. These designs include analysis of variance (ANOVA) setups, interaction effects, and nonlinear effects. Besides describing and interpreting general regression models, this chapter also describes, in general terms, how the three types of assumptions of regression models can be examined.

First we introduce notation for regression models. Let Y denote the response (dependent) variable, and let $X = X_1, X_2, \dots, X_p$ denote a list or vector of predictor variables (also called covariables or independent, descriptor, or concomitant variables). These predictor variables are assumed to be constants for a given individual or subject from the population of interest. Let $\beta = \beta_0, \beta_1, \dots, \beta_p$ denote the list of regression coefficients (parameters). β_0 is an optional intercept parameter, and β_1, \dots, β_p are weights or regression coefficients corresponding to X_1, \dots, X_p . We use matrix or vector notation to describe a weighted sum of the X s:

$$X\beta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p, \quad (2.1)$$

where there is an implied $X_0 = 1$.

A regression model is stated in terms of a connection between the predictors X and the response Y . Let $C(Y|X)$ denote a property of the distribution of Y given X (as a function of X). For example, $C(Y|X)$ could be $E(Y|X)$,

the expected value or average of Y given X , or $C(Y|X)$ could be the probability that $Y = 1$ given X (where $Y = 0$ or 1).

2.2 Model Formulations

We define a regression function as a function that describes interesting properties of Y that may vary across individuals in the population. X describes the list of factors determining these properties. Stated mathematically, a general regression model is given by

$$C(Y|X) = g(X). \quad (2.2)$$

We restrict our attention to models that, after a certain transformation, are linear in the unknown parameters, that is, models that involve X only through a weighted sum of all the X s. The *general linear regression model* is given by

$$C(Y|X) = g(X\beta). \quad (2.3)$$

For example, the ordinary linear regression model is

$$C(Y|X) = E(Y|X) = X\beta, \quad (2.4)$$

and given X , Y has a normal distribution with mean $X\beta$ and constant variance σ^2 . The binary logistic regression model^{129,647} is

$$C(Y|X) = \text{Prob}\{Y = 1|X\} = (1 + \exp(-X\beta))^{-1}, \quad (2.5)$$

where Y can take on the values 0 and 1. In general the model, when stated in terms of the property $C(Y|X)$, may not be linear in $X\beta$; that is $C(Y|X) = g(X\beta)$, where $g(u)$ is nonlinear in u . For example, a regression model could be $E(Y|X) = (X\beta)^5$. The model may be made linear in the unknown parameters by a transformation in the property $C(Y|X)$:

$$h(C(Y|X)) = X\beta, \quad (2.6)$$

where $h(u) = g^{-1}(u)$, the inverse function of g . As an example consider the binary logistic regression model given by

$$C(Y|X) = \text{Prob}\{Y = 1|X\} = (1 + \exp(-X\beta))^{-1}. \quad (2.7)$$

If $h(u) = \text{logit}(u) = \log(u/(1-u))$, the transformed model becomes

$$h(\text{Prob}(Y = 1|X)) = \log(\exp(X\beta)) = X\beta. \quad (2.8)$$

The transformation $h(C(Y|X))$ is sometimes called a *link function*. Let $h(C(Y|X))$ be denoted by $C'(Y|X)$. The general linear regression model then becomes

$$C'(Y|X) = X\beta. \quad (2.9)$$

In other words, the model states that some property C' of Y , given X , is a weighted sum of the X s ($X\beta$). In the ordinary linear regression model, $C'(Y|X) = E(Y|X)$. In the logistic regression case, $C'(Y|X)$ is the logit of the probability that $Y = 1$, $\log \text{Prob}\{Y = 1\} / [1 - \text{Prob}\{Y = 1\}]$. This is the log of the odds that $Y = 1$ versus $Y = 0$.

It is important to note that the general linear regression model has two major components: $C'(Y|X)$ and $X\beta$. The first part has to do with a property or transformation of Y . The second, $X\beta$, is the *linear regression* or *linear predictor* part. The method of least squares can sometimes be used to fit the model if $C'(Y|X) = E(Y|X)$. Other cases must be handled using other methods such as maximum likelihood estimation or nonlinear least squares.

2.3 Interpreting Model Parameters

In the original model, $C(Y|X)$ specifies the way in which X affects a property of Y . Except in the ordinary linear regression model, it is difficult to interpret the individual parameters if the model is stated in terms of $C(Y|X)$. In the model $C'(Y|X) = X\beta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$, the regression parameter β_j is interpreted as the change in the property C' of Y per unit change in the descriptor variable X_j , all other descriptors remaining constant^a:

$$\beta_j = C'(Y|X_1, X_2, \dots, X_j + 1, \dots, X_p) - C'(Y|X_1, X_2, \dots, X_j, \dots, X_p). \quad (2.10)$$

In the ordinary linear regression model, for example, β_j is the change in expected value of Y per unit change in X_j . In the logistic regression model β_j is the change in log odds that $Y = 1$ per unit change in X_j . When a non-interacting X_j is a dichotomous variable or a continuous one that is linearly related to C' , X_j is represented by a single term in the model and its contribution is described fully by β_j .

In all that follows, we drop the $'$ from C' and assume that $C(Y|X)$ is the property of Y that is linearly related to the weighted sum of the X s.

^a Note that it is not necessary to “hold constant” all other variables to be able to interpret the effect of one predictor. It is sufficient to hold constant the weighted sum of all the variables other than X_j . And in many cases it is not physically possible to hold other variables constant while varying one, e.g., when a model contains X and X^2 (David Hoaglin, personal communication).

2.3.1 Nominal Predictors

Suppose that we wish to model the effect of two or more treatments and be able to test for differences between the treatments in some property of Y . A nominal or polytomous factor such as treatment group having k levels, in which there is no definite ordering of categories, is fully described by a series of $k - 1$ binary indicator variables (sometimes called *dummy variables*). Suppose that there are four treatments, J, K, L , and M , and the treatment factor is denoted by T . The model can be written as

$$\begin{aligned} C(Y|T = J) &= \beta_0 \\ C(Y|T = K) &= \beta_0 + \beta_1 \\ C(Y|T = L) &= \beta_0 + \beta_2 \\ C(Y|T = M) &= \beta_0 + \beta_3. \end{aligned} \tag{2.11}$$

The four treatments are thus completely specified by three regression parameters and one intercept that we are using to denote treatment J , the reference treatment. This model can be written in the previous notation as

$$C(Y|T) = X\beta = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3, \tag{2.12}$$

where

$$\begin{aligned} X_1 &= 1 \text{ if } T = K, 0 \text{ otherwise} \\ X_2 &= 1 \text{ if } T = L, 0 \text{ otherwise} \\ X_3 &= 1 \text{ if } T = M, 0 \text{ otherwise.} \end{aligned} \tag{2.13}$$

For treatment J ($T = J$), all three X s are zero and $C(Y|T = J) = \beta_0$. The test for any differences in the property $C(Y)$ between treatments is $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$.

This model is an *analysis of variance* or *k-sample-type* model. If there are other descriptor covariables in the model, it becomes an *analysis of covariance-type* model.

2.3.2 Interactions

Suppose that a model has descriptor variables X_1 and X_2 and that the effect of the two X s cannot be separated; that is the effect of X_1 on Y depends on the level of X_2 and vice versa. One simple way to describe this *interaction* is to add the constructed variable $X_3 = X_1 X_2$ to the model:

$$C(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2. \tag{2.14}$$

It is now difficult to interpret β_1 and β_2 in isolation. However, we may quantify the effect of a one-unit increase in X_1 if X_2 is held constant as

Table 2.1 Parameters in a simple model with interaction

Parameter	Meaning
β_0	$C(Y age = 0, sex = m)$
β_1	$C(Y age = x + 1, sex = m) - C(Y age = x, sex = m)$
β_2	$C(Y age = 0, sex = f) - C(Y age = 0, sex = m)$
β_3	$C(Y age = x + 1, sex = f) - C(Y age = x, sex = f) - [C(Y age = x + 1, sex = m) - C(Y age = x, sex = m)]$

$$\begin{aligned}
 C(Y|X_1 + 1, X_2) - C(Y|X_1, X_2) & \\
 &= \beta_0 + \beta_1(X_1 + 1) + \beta_2X_2 \\
 &+ \beta_3(X_1 + 1)X_2 \qquad (2.15) \\
 &- [\beta_0 + \beta_1X_1 + \beta_2X_2 + \beta_3X_1X_2] \\
 &= \beta_1 + \beta_3X_2.
 \end{aligned}$$

Likewise, the effect of a one-unit increase in X_2 on C if X_1 is held constant is $\beta_2 + \beta_3X_1$. Interactions can be much more complex than can be modeled with a product of two terms. If X_1 is binary, the interaction may take the form of a difference in shape (and/or distribution) of X_2 versus $C(Y)$ depending on whether $X_1 = 0$ or $X_1 = 1$ (e.g., logarithm vs. square root). When both variables are continuous, the possibilities are much greater (this case is discussed later). Interactions among more than two variables can be exceedingly complex.

2.3.3 Example: Inference for a Simple Model

Suppose we postulated the model

$$C(Y|age, sex) = \beta_0 + \beta_1age + \beta_2[sex = f] + \beta_3age[sex = f],$$

where $[sex = f]$ is a 0–1 indicator variable for sex = female; the reference cell is sex = male corresponding to a zero value of the indicator variable. This is a model that assumes

1. age is linearly related to $C(Y)$ for males,
2. age is linearly related to $C(Y)$ for females, and
3. whatever distribution, variance, and independence assumptions are appropriate for the model being considered.

We are thus assuming that the interaction between age and sex is simple; that is it only alters the slope of the age effect. The parameters in the model have interpretations shown in Table 2.1. β_3 is the difference in slopes (female – male).

There are many useful hypotheses that can be tested for this model. First let's consider two hypotheses that are seldom appropriate although they are routinely tested.

1. $H_0 : \beta_1 = 0$: This tests whether age is associated with Y for males.
2. $H_0 : \beta_2 = 0$: This tests whether sex is associated with Y for zero-year olds.

Now consider more useful hypotheses. For each hypothesis we should write what is being tested, translate this to tests in terms of parameters, write the alternative hypothesis, and describe what the test has maximum power to detect. The latter component of a hypothesis test needs to be emphasized, as almost every statistical test is focused on one specific pattern to detect. For example, a test of association against an alternative hypothesis that a slope is nonzero will have maximum power when the true association is linear. If the true regression model is exponential in X , a linear regression test will have some power to detect “non-flatness” but it will not be as powerful as the test from a well-specified exponential regression effect. If the true effect is U-shaped, a test of association based on a linear model will have almost no power to detect association. If one tests for association against a quadratic (parabolic) alternative, the test will have some power to detect a logarithmic shape but it will have very little power to detect a cyclical trend having multiple “humps.” In a quadratic regression model, a test of linearity against a quadratic alternative hypothesis will have reasonable power to detect a quadratic nonlinear effect but very limited power to detect a multiphase cyclical trend. Therefore in the tests in Table 2.2 keep in mind that power is maximal when linearity of the age relationship holds for both sexes. In fact it may be useful to write alternative hypotheses as, for example, “ H_a : age is associated with $C(Y)$, powered to detect a *linear* relationship.”

Note that if there is an interaction effect, we know that there is both an age and a sex effect. However, there can also be age or sex effects when the lines are parallel. That's why the tests of total association have 2 d.f.

2.4 Relaxing Linearity Assumption for Continuous Predictors

2.4.1 Avoiding Categorization

Relationships among variables are seldom linear, except in special cases such as when one variable is compared with itself measured at a different time. It is a common belief among practitioners who do not study bias and

efficiency in depth that the presence of non-linearity should be dealt with by chopping continuous variables into intervals. Nothing could be more disastrous.^{13, 14, 17, 45, 82, 185, 187, 215, 294, 300, 379, 446, 465, 521, 533, 559, 597, 646}

Table 2.2 Most Useful Tests for Linear *Age* \times *Sex* Model

Null or Alternative Hypothesis	Mathematical Statement
Effect of age is independent of sex or Effect of sex is independent of age or Age and sex are additive Age effects are parallel	$H_0 : \beta_3 = 0$
Age interacts with sex Age modifies effect of sex Sex modifies effect of age Sex and age are non-additive (synergistic)	$H_a : \beta_3 \neq 0$
Age is not associated with <i>Y</i> Age is associated with <i>Y</i> Age is associated with <i>Y</i> for either Females or males	$H_0 : \beta_1 = \beta_3 = 0$ $H_a : \beta_1 \neq 0$ or $\beta_3 \neq 0$
Sex is not associated with <i>Y</i> Sex is associated with <i>Y</i> Sex is associated with <i>Y</i> for some Value of age	$H_0 : \beta_2 = \beta_3 = 0$ $H_a : \beta_2 \neq 0$ or $\beta_3 \neq 0$
Neither age nor sex is associated with <i>Y</i> Either age or sex is associated with <i>Y</i>	$H_0 : \beta_1 = \beta_2 = \beta_3 = 0$ $H_a : \beta_1 \neq 0$ or $\beta_2 \neq 0$ or $\beta_3 \neq 0$

Problems caused by dichotomization include the following.

1. Estimated values will have reduced precision, and associated tests will have reduced power.
2. Categorization assumes that the relationship between the predictor and the response is flat within intervals; this assumption is far less reasonable than a linearity assumption in most cases.
3. To make a continuous predictor be more accurately modeled when categorization is used, multiple intervals are required. The needed indicator variables will spend more degrees of freedom than will fitting a smooth relationship, hence power and precision will suffer. And because of sample size limitations in the very low and very high range of the variable, the outer intervals (e.g., outer quintiles) will be wide, resulting in significant heterogeneity of subjects within those intervals, and residual confounding.
4. Categorization assumes that there is a discontinuity in response as interval boundaries are crossed. Other than the effect of time (e.g., an instant stock price drop after bad news), there are very few examples in which such discontinuities have been shown to exist.
5. Categorization only seems to yield interpretable estimates such as odds ratios. For example, suppose one computes the odds ratio for stroke for persons with a systolic blood pressure > 160 mmHg compared with persons with a blood

pressure ≤ 160 mmHg. The interpretation of the resulting odds ratio will depend on the exact distribution of blood pressures in the sample (the proportion of subjects > 170 , > 180 , etc.). On the other hand, if blood pressure is modeled as a continuous variable (e.g., using a regression spline, quadratic, or linear effect) one can estimate the ratio of odds for *exact* settings of the predictor, e.g., the odds ratio for 200 mmHg compared with 120 mmHg.

6. Categorization does not condition on full information. When, for example, the risk of stroke is being assessed for a new subject with a known blood pressure (say 162 mmHg), the subject does not report to her physician “my blood pressure exceeds 160” but rather reports 162 mmHg. The risk for this subject will be much lower than that of a subject with a blood pressure of 200 mmHg.
7. If cutpoints are determined in a way that is not blinded to the response variable, calculation of P -values and confidence intervals requires special simulation techniques; ordinary inferential methods are completely invalid. For example, if cutpoints are chosen by trial and error in a way that utilizes the response, even informally, ordinary P -values will be too small and confidence intervals will not have the claimed coverage probabilities. The correct Monte-Carlo simulations must take into account both multiplicities and uncertainty in the choice of cutpoints. For example, if a cutpoint is chosen that minimizes the P -value and the resulting P -value is 0.05, the true type I error can easily be above 0.5³⁰⁰.
8. Likewise, categorization that is not blinded to the response variable results in biased effect estimates^{17, 559}.
9. “Optimal” cutpoints do not replicate over studies. Hollander et al.³⁰⁰ state that “. . . the optimal cutpoint approach has disadvantages. One of these is that in almost every study where this method is applied, another cutpoint will emerge. This makes comparisons across studies extremely difficult or even impossible. Altman et al. point out this problem for studies of the prognostic relevance of the S-phase fraction in breast cancer published in the literature. They identified 19 different cutpoints used in the literature; some of them were solely used because they emerged as the ‘optimal’ cutpoint in a specific data set. In a meta-analysis on the relationship between cathepsin-D content and disease-free survival in node-negative breast cancer patients, 12 studies were included with 12 different cutpoints . . . Interestingly, neither cathepsin-D nor the S-phase fraction are recommended to be used as prognostic markers in breast cancer in the recent update of the American Society of Clinical Oncology.” Giannoni et al.²¹⁵ demonstrated that many claimed “optimal cutpoints” are just the observed median values in the sample, which happens to optimize statistical power for detecting a separation in outcomes and have nothing to do with true outcome thresholds. Disagreements in cutpoints (which are bound to happen whenever one searches for things that do not exist) cause severe interpretation problems. One study may provide an odds ratio for comparing body mass index (BMI) > 30 with BMI ≤ 30 , another for comparing BMI > 28 with BMI ≤ 28 . Neither of these odds ratios has a good definition and the two estimates are not comparable.
10. Cutpoints are arbitrary and manipulatable; cutpoints can be found that can result in both positive and negative associations⁶⁴⁶.
11. If a confounder is adjusted for by categorization, there will be residual confounding that can be explained away by inclusion of the continuous form of the predictor in the model in addition to the categories.

When cutpoints are chosen using Y , categorization represents one of those few times in statistics where both type I and type II errors are elevated.

A scientific quantity is a quantity which can be defined outside of the specifics of the current experiment. The kind of high:low estimates that result from categorizing a continuous variable are not scientific quantities; their interpretation depends on the entire sample distribution of continuous measurements within the chosen intervals.

To summarize problems with categorization it is useful to examine its effective assumptions. Suppose one assumes there is a single cutpoint c for predictor X . Assumptions implicit in seeking or using this cutpoint include (1) the relationship between X and the response Y is discontinuous at $X = c$ and only $X = c$; (2) c is correctly found as the cutpoint; (3) X vs. Y is flat to the left of c ; (4) X vs. Y is flat to the right of c ; (5) the “optimal” cutpoint does not depend on the values of other predictors. Failure to have these assumptions satisfied will result in great error in estimating c (because it doesn’t exist), low predictive accuracy, serious lack of model fit, residual confounding, and overestimation of effects of remaining variables.

A better approach that maximizes power and that only assumes a smooth relationship is to use regression splines for predictors that are not known to predict linearly. Use of flexible parametric approaches such as this allows standard inference techniques (P -values, confidence limits) to be used, as will be described below. Before introducing splines, we consider the simplest approach to allowing for nonlinearity.

2.4.2 Simple Nonlinear Terms

If a continuous predictor is represented, say, as X_1 in the model, the model is assumed to be linear in X_1 . Often, however, the property of Y of interest does not behave linearly in all the predictors. The simplest way to describe a nonlinear effect of X_1 is to include a term for $X_2 = X_1^2$ in the model:

$$C(Y|X_1) = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2. \quad (2.16)$$

If the model is truly linear in X_1 , β_2 will be zero. This model formulation allows one to test H_0 : model is linear in X_1 against H_a : model is quadratic (parabolic) in X_1 by testing $H_0 : \beta_2 = 0$.

Nonlinear effects will frequently not be of a parabolic nature. If a transformation of the predictor is known to induce linearity, that transformation (e.g., $\log(X)$) may be substituted for the predictor. However, often the transformation is not known. Higher powers of X_1 may be included in the model to approximate many types of relationships, but polynomials have some undesirable properties (e.g., undesirable peaks and valleys, and the fit in one region of X can be greatly affected by data in other regions⁴³³) and will not adequately fit many functional forms.¹⁵⁶ For example, polynomials do not adequately fit logarithmic functions or “threshold” effects.

2.4.3 Splines for Estimating Shape of Regression Function and Determining Predictor Transformations

A draftsman's *spline* is a flexible strip of metal or rubber used to draw curves. Spline functions are piecewise polynomials used in curve fitting. That is, they are polynomials within intervals of X that are connected across different intervals of X . Splines have been used, principally in the physical sciences, to approximate a wide variety of functions. The simplest spline function is a linear spline function, a piecewise linear function. Suppose that the x axis is divided into intervals with endpoints at a , b , and c , called *knots*. The linear spline function is given by

$$f(X) = \beta_0 + \beta_1 X + \beta_2(X - a)_+ + \beta_3(X - b)_+ + \beta_4(X - c)_+, \quad (2.17)$$

where

$$\begin{aligned} (u)_+ &= u, u > 0, \\ &0, u \leq 0. \end{aligned} \quad (2.18)$$

The number of knots can vary depending on the amount of available data for fitting the function. The linear spline function can be rewritten as

$$\begin{aligned} f(X) &= \beta_0 + \beta_1 X, & X \leq a \\ &= \beta_0 + \beta_1 X + \beta_2(X - a) & a < X \leq b \\ &= \beta_0 + \beta_1 X + \beta_2(X - a) + \beta_3(X - b) & b < X \leq c \\ &= \beta_0 + \beta_1 X + \beta_2(X - a) \\ &\quad + \beta_3(X - b) + \beta_4(X - c) & c < X. \end{aligned} \quad (2.19)$$

A linear spline is depicted in Figure 2.1.

The general linear regression model can be written assuming only piecewise linearity in X by incorporating constructed variables X_2 , X_3 , and X_4 :

$$C(Y|X) = f(X) = X\beta, \quad (2.20)$$

where $X\beta = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$, and

$$\begin{aligned} X_1 &= X & X_2 &= (X - a)_+ \\ X_3 &= (X - b)_+ & X_4 &= (X - c)_+. \end{aligned} \quad (2.21)$$

By modeling a slope increment for X in an interval $(a, b]$ in terms of $(X - a)_+$, the function is constrained to join ("meet") at the knots. Overall linearity in X can be tested by testing $H_0 : \beta_2 = \beta_3 = \beta_4 = 0$.

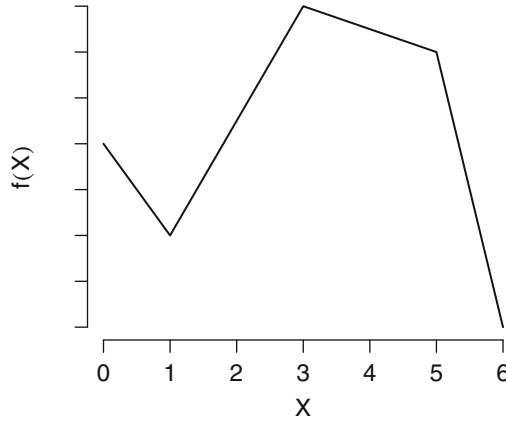


Fig. 2.1 A linear spline function with knots at $a = 1, b = 3, c = 5$.

2.4.4 Cubic Spline Functions

Although the linear spline is simple and can approximate many common relationships, it is not smooth and will not fit highly curved functions well. These problems can be overcome by using piecewise polynomials of order higher than linear. Cubic polynomials have been found to have nice properties with good ability to fit sharply curving shapes. Cubic splines can be made to be smooth at the join points (knots) by forcing the first and second derivatives of the function to agree at the knots. Such a smooth cubic spline function with three knots (a, b, c) is given by

$$\begin{aligned} f(X) &= \beta_0 + \beta_1 X + \beta_2 X^2 + \beta_3 X^3 \\ &\quad + \beta_4 (X - a)_+^3 + \beta_5 (X - b)_+^3 + \beta_6 (X - c)_+^3 \\ &= X\beta \end{aligned} \quad (2.22)$$

with the following constructed variables:

$$\begin{aligned} X_1 &= X & X_2 &= X^2 \\ X_3 &= X^3 & X_4 &= (X - a)_+^3 \\ X_5 &= (X - b)_+^3 & X_6 &= (X - c)_+^3. \end{aligned} \quad (2.23)$$

If the cubic spline function has k knots, the function will require estimating $k + 3$ regression coefficients besides the intercept. See Section 2.4.6 for information on choosing the number and location of knots. 1

There are more numerically stable ways to form a design matrix for cubic spline functions that are based on B-splines instead of the truncated power basis^{152, 575} used here. However, B-splines are more complex and do not allow for extrapolation beyond the outer knots, and the truncated power basis seldom presents estimation problems (see Section 4.6) when modern methods such as the Q-R decomposition are used for matrix inversion. 2

2.4.5 Restricted Cubic Splines

Stone and Koo⁵⁹⁵ have found that cubic spline functions do have a drawback in that they can be poorly behaved in the tails, that is before the first knot and after the last knot. They cite advantages of constraining the function to be linear in the tails. Their restricted cubic spline function (also called *natural splines*) has the additional advantage that only $k - 1$ parameters must be estimated (besides the intercept) as opposed to $k + 3$ parameters with the unrestricted cubic spline. The restricted spline function with k knots t_1, \dots, t_k is given by¹⁵⁶

$$f(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_{k-1} X_{k-1}, \quad (2.24)$$

where $X_1 = X$ and for $j = 1, \dots, k - 2$,

$$\begin{aligned} X_{j+1} = & (X - t_j)_+^3 - (X - t_{k-1})_+^3 (t_k - t_j) / (t_k - t_{k-1}) \\ & + (X - t_k)_+^3 (t_{k-1} - t_j) / (t_k - t_{k-1}). \end{aligned} \quad (2.25)$$

It can be shown that X_j is linear in X for $X \geq t_k$. For numerical behavior and to put all basis functions for X on the same scale, R `Hmisc` and `rms` package functions by default divide the terms in Eq. 2.25 by

$$\tau = (t_k - t_1)^2. \quad (2.26)$$

Figure 2.2 displays the τ -scaled spline component variables X_j for $j = 2, 3, 4$ and $k = 5$ and one set of knots. The left graph magnifies the lower portion of the curves.

```
require(Hmisc)
```

```
x ← rcspline.eval(seq(0,1,.01),
                  knots=seq(.05,.95,length=5), inclx=T)
xm ← x
xm[xm > .0106] ← NA
matplot(x[,1], xm, type="l", ylim=c(0,.01),
        xlab=expression(X), ylab='', lty=1)
matplot(x[,1], x, type="l",
        xlab=expression(X), ylab='', lty=1)
```

Figure 2.3 displays some typical shapes of restricted cubic spline functions with $k = 3, 4, 5$, and 6. These functions were generated using random β .

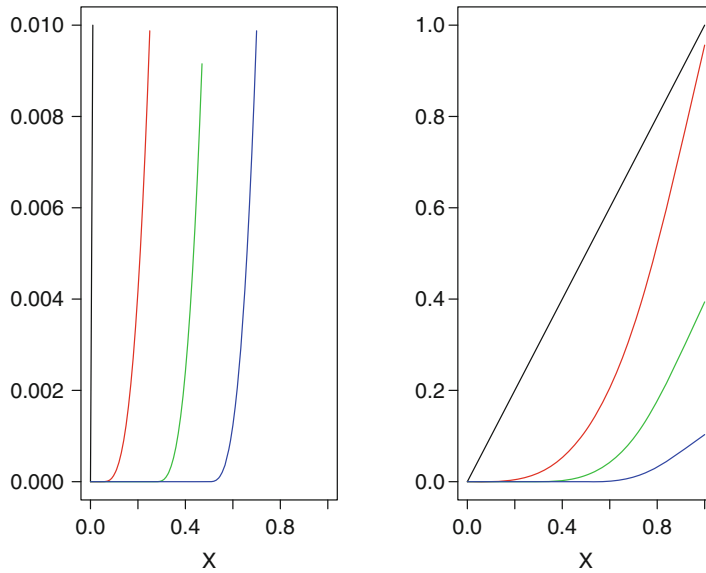


Fig. 2.2 Restricted cubic spline component variables for $k = 5$ and knots at $X = .05, .275, .5, .725, \text{ and } .95$. Nonlinear basis functions are scaled by τ . The left panel is a y -magnification of the right panel. Fitted functions such as those in Figure 2.3 will be linear combinations of these basis functions as long as knots are at the same locations used here.

```
x ← seq(0, 1, length=300)
for(nk in 3:6) {
  set.seed(nk)
  knots ← seq(.05, .95, length=nk)
  xx ← rcspline.eval(x, knots=knots, inclx=T)
  for(i in 1 : (nk - 1))
    xx[,i] ← (xx[,i] - min(xx[,i])) /
             (max(xx[,i]) - min(xx[,i]))
  for(i in 1 : 20) {
    beta ← 2*runif(nk-1) - 1
    xbeta ← xx %*% beta + 2 * runif(1) - 1
    xbeta ← (xbeta - min(xbeta)) /
            (max(xbeta) - min(xbeta))
    if(i == 1) {
      plot(x, xbeta, type="l", lty=1,
           xlab=expression(X), ylab='', bty="n")
      title(sub=paste(nk,"knots"), adj=0, cex=.75)
      for(j in 1 : nk)
        arrows(knots[j], .04, knots[j], -.03,
              angle=20, length=.07, lwd=1.5)
    }
    else lines(x, xbeta, col=i)
  }
}
```

Once $\beta_0, \dots, \beta_{k-1}$ are estimated, the restricted cubic spline can be restated in the form

$$f(X) = \beta_0 + \beta_1 X + \beta_2 (X - t_1)_+^3 + \beta_3 (X - t_2)_+^3 + \dots + \beta_{k+1} (X - t_k)_+^3 \quad (2.27)$$

by dividing $\beta_2, \dots, \beta_{k-1}$ by τ (Eq. 2.26) and computing

$$\begin{aligned} \beta_k &= [\beta_2(t_1 - t_k) + \beta_3(t_2 - t_k) + \dots + \beta_{k-1}(t_{k-2} - t_k)] / (t_k - t_{k-1}) \quad (2.28) \\ \beta_{k+1} &= [\beta_2(t_1 - t_{k-1}) + \beta_3(t_2 - t_{k-1}) + \dots + \beta_{k-1}(t_{k-2} - t_{k-1})] / (t_{k-1} - t_k). \end{aligned}$$

A test of linearity in X can be obtained by testing

$$H_0 : \beta_2 = \beta_3 = \dots = \beta_{k-1} = 0. \quad (2.29)$$

4 The truncated power basis for restricted cubic splines does allow for rational (i.e., linear) extrapolation beyond the outer knots. However, when the outer knots are in the tails of the data, extrapolation can still be dangerous.

When nonlinear terms in Equation 2.25 are normalized, for example, by dividing them by the square of the difference in the outer knots to make all terms have units of X , the ordinary truncated power basis has no numerical difficulties when modern matrix algebra software is used.

2.4.6 Choosing Number and Position of Knots

We have assumed that the locations of the knots are specified in advance; that is, the knot locations are not treated as free parameters to be estimated. If knots were free parameters, the fitted function would have more flexibility but at the cost of instability of estimates, statistical inference problems, and inability to use standard regression modeling software for estimating regression parameters.

How then does the analyst pre-assign knot locations? If the regression relationship were described by prior experience, pre-specification of knot locations would be easy. For example, if a function were known to change curvature at $X = a$, a knot could be placed at a . However, in most situations there is no way to pre-specify knots. Fortunately, Stone⁵⁹³ has found that the location of knots in a restricted cubic spline model is not very crucial in most situations; the fit depends much more on the choice of k , the number of knots. Placing knots at fixed quantiles (percentiles) of a predictor's marginal distribution is a good approach in most datasets. This ensures that enough points are available in each interval, and also guards against letting outliers overly influence knot placement. Recommended equally spaced quantiles are shown in Table 2.3.

5

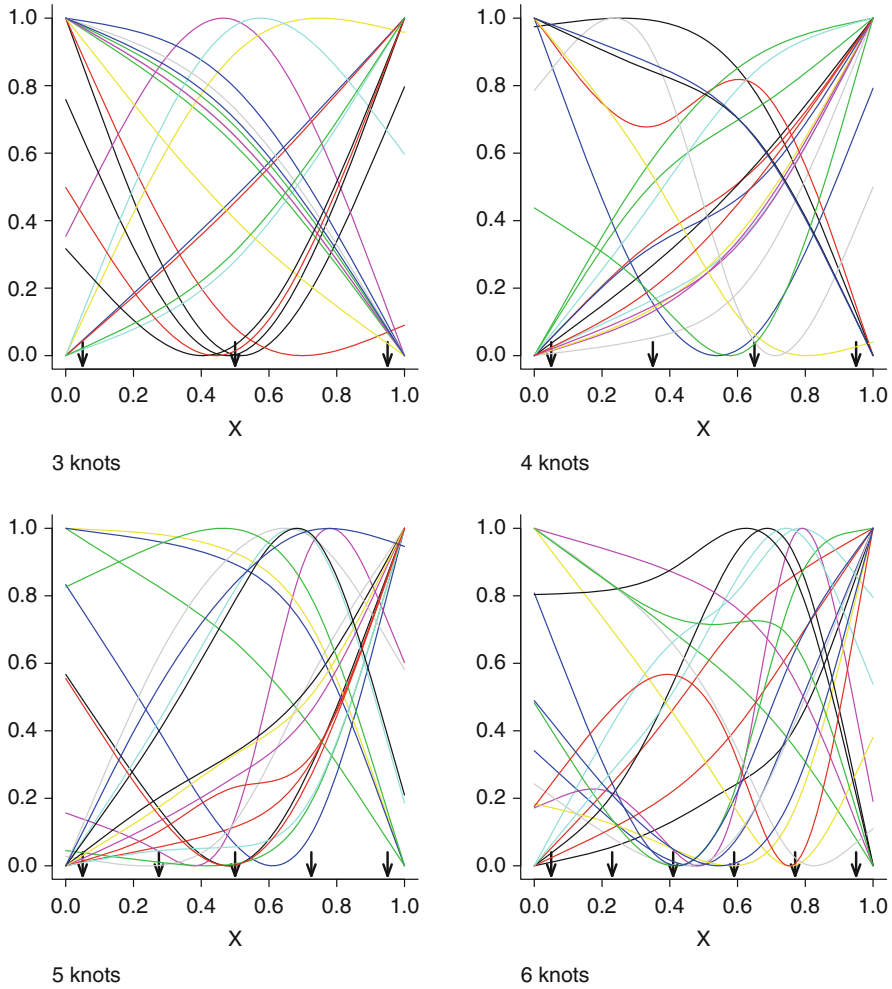


Fig. 2.3 Some typical restricted cubic spline functions for $k = 3, 4, 5, 6$. The y -axis is $X\beta$. Arrows indicate knots. These curves were derived by randomly choosing values of β subject to standard deviations of fitted functions being normalized.

Table 2.3 Default quantiles for knots

k	Quantiles						
3	.10	.5	.90				
4	.05	.35	.65	.95			
5	.05	.275	.5	.725	.95		
6	.05	.23	.41	.59	.77	.95	
7	.025	.1833	.3417	.5	.6583	.8167	.975

The principal reason for using less extreme default quantiles for $k = 3$ and more extreme ones for $k = 7$ is that one usually uses $k = 3$ for small sample sizes and $k = 7$ for large samples. When the sample size is less than 100, the outer quantiles should be replaced by the fifth smallest and fifth largest data points, respectively.⁵⁹⁵ What about the choice of k ? The flexibility of possible fits must be tempered by the sample size available to estimate the unknown parameters. Stone⁵⁹³ has found that more than 5 knots are seldom required in a restricted cubic spline model. The principal decision then is between $k = 3, 4$, or 5. For many datasets, $k = 4$ offers an adequate fit of the model and is a good compromise between flexibility and loss of precision caused by overfitting a small sample. When the sample size is large (e.g., $n \geq 100$ with a continuous uncensored response variable), $k = 5$ is a good choice. Small samples (< 30 , say) may require the use of $k = 3$. Akaike's information criterion (AIC, Section 9.8.1) can be used for a data-based choice of k . The value of k maximizing the model likelihood ratio $\chi^2 - 2k$ would be the best "for the money" using AIC.

The analyst may wish to devote more knots to variables that are thought to be more important, and risk lack of fit for less important variables. In this way the total number of estimated parameters can be controlled (Section 4.1).

2.4.7 Nonparametric Regression

One of the most important results of an analysis is the estimation of the tendency (trend) of how X relates to Y . This trend is useful in its own right and it may be sufficient for obtaining predicted values in some situations, but trend estimates can also be used to guide formal regression modeling (by suggesting predictor variable transformations) and to check model assumptions.

Nonparametric smoothers are excellent tools for determining the shape of the relationship between a predictor and the response. The standard nonparametric smoothers work when one is interested in assessing one continuous predictor at a time and when the property of the response that *should* be linearly related to the predictor is a standard measure of central tendency. For example, when $C(Y)$ is $E(Y)$ or $\Pr[Y = 1]$, standard smoothers are useful, but when $C(Y)$ is a measure of variability or a rate (instantaneous risk), or when Y is only incompletely measured for some subjects (e.g., Y is censored for some subjects), simple smoothers will not work.

The oldest and simplest nonparametric smoother is the moving average. Suppose that the data consist of the points $X = 1, 2, 3, 5$, and 8, with the corresponding Y values 2.1, 3.8, 5.7, 11.1, and 17.2. To smooth the relationship we could estimate $E(Y|X = 2)$ by $(2.1 + 3.8 + 5.7)/3$ and $E(Y|X = (2 + 3 + 5)/3)$ by $(3.8 + 5.7 + 11.1)/3$. Note that overlap is fine; that is one point may be contained in two sets that are averaged. You can immediately see that the

simple moving average has a problem in estimating $E(Y)$ at the outer values of X . The estimates are quite sensitive to the choice of the number of points (or interval width) to use in “binning” the data.

A moving least squares linear regression smoother is far superior to a moving flat line smoother (moving average). Cleveland’s¹¹¹ moving linear regression smoother *loess* has become the most popular smoother. To obtain the smoothed value of Y at $X = x$, we take all the data having X values within a suitable interval about x . Then a linear regression is fitted to all of these points, and the predicted value from this regression at $X = x$ is taken as the estimate of $E(Y|X = x)$. Actually, *loess* uses weighted least squares estimates, which is why it is called a *locally weighted least squares* method. The weights are chosen so that points near $X = x$ are given the most weight^b in the calculation of the slope and intercept. Surprisingly, a good default choice for the interval about x is an interval containing 2/3 of the data points! The weighting function is devised so that points near the extremes of this interval receive almost no weight in the calculation of the slope and intercept.

Because *loess* uses a moving straight line rather than a moving flat one, it provides much better behavior at the extremes of the X s. For example, one can fit a straight line to the first three data points and then obtain the predicted value at the lowest X , which takes into account that this X is not the middle of the three X s.

loess obtains smoothed values for $E(Y)$ at each observed value of X . Estimates for other X s are obtained by linear interpolation.

The *loess* algorithm has another component. After making an initial estimate of the trend line, *loess* can look for outliers off this trend. It can then delete or down-weight those apparent outliers to obtain a more robust trend estimate. Now, different points will appear to be outliers with respect to this second trend estimate. The new set of outliers is taken into account and another trend line is derived. By default, the process stops after these three iterations. *loess* works exceptionally well for binary Y as long as the iterations that look for outliers are not done, that is only one iteration is performed.

For a single X , Friedman’s “super smoother”²⁰⁷ is another efficient and flexible nonparametric trend estimator. For both *loess* and the super smoother the amount of smoothing can be controlled by the analyst. Hastie and Tibshirani²⁷⁵ provided an excellent description of smoothing methods and developed a generalized additive model for multiple X s, in which each continuous predictor is fitted with a nonparametric smoother (see Chapter 16). Interactions are not allowed. Cleveland et al.⁹⁶ have extended two-dimensional smoothers to multiple dimensions without assuming additivity. Their *local regression model* is feasible for up to four or so predictors. Local regression models are extremely flexible, allowing parts of the model to be

6

^b This weight is not to be confused with the regression coefficient; rather the weights are w_1, w_2, \dots, w_n and the fitting criterion is $\sum_i^n w_i (Y_i - \hat{Y}_i)^2$.

parametrically specified, and allowing the analyst to choose the amount of smoothing or the effective number of degrees of freedom of the fit.

Smoothing splines are related to nonparametric smoothers. Here a knot is placed at every data point, but a penalized likelihood is maximized to derive the smoothed estimates. Gray^{237,238} developed a general method that is halfway between smoothing splines and regression splines. He pre-specified, say, 10 fixed knots, but uses a penalized likelihood for estimation. This allows

7

the analyst to control the effective number of degrees of freedom used. Besides using smoothers to estimate regression relationships, smoothers are valuable for examining trends in residual plots. See Sections 14.6 and 21.2 for examples.

2.4.8 Advantages of Regression Splines over Other Methods

There are several advantages of regression splines:²⁷¹

1. Parametric splines are piecewise polynomials and can be fitted using any existing regression program after the constructed predictors are computed. Spline regression is equally suitable to multiple linear regression, survival models, and logistic models for discrete outcomes.
2. Regression coefficients for the spline function are estimated using standard techniques (maximum likelihood or least squares), and statistical inferences can readily be drawn. Formal tests of no overall association, linearity, and additivity can readily be constructed. Confidence limits for the estimated regression function are derived by standard theory.
3. The fitted spline function directly estimates the transformation that a predictor should receive to yield linearity in $C(Y|X)$. The fitted spline transformation sometimes suggests a simple transformation (e.g., square root) of a predictor that can be used if one is not concerned about the proper number of degrees of freedom for testing association of the predictor with the response.
4. The spline function can be used to represent the predictor in the final model. Nonparametric methods do not yield a prediction equation.
5. Splines can be extended to non-additive models (see below). Multidimensional nonparametric estimators often require burdensome computations.

2.5 Recursive Partitioning: Tree-Based Models

Breiman et al.⁶⁹ have developed an essentially model-free approach called *classification and regression trees* (CART), a form of recursive partitioning.

For some implementations of CART, we say “essentially” model-free since a model-based statistic is sometimes chosen as a splitting criterion. The essence of recursive partitioning is as follows.

1. Find the predictor so that the best possible binary split on that predictor has a larger value of some statistical criterion than any other split on any other predictor. For ordinal and continuous predictors, the split is of the form $X < c$ versus $X \geq c$. For polytomous predictors, the split involves finding the best separation of categories, without preserving order.
2. Within each previously formed subset, find the best predictor and best split that maximizes the criterion in the subset of observations passing the previous split.
3. Proceed in like fashion until fewer than k observations remain to be split, where k is typically 20 to 100.
4. Obtain predicted values using a statistic that summarizes each terminal node (e.g., mean or proportion).
5. Prune the tree backward so that a tree with the same number of nodes developed on 0.9 of the data validates best on the remaining 0.1 of the data (average over the 10 cross-validations). Alternatively, shrink the node estimates toward the mean, using a progressively stronger shrinkage factor, until the best cross-validation results.

8

Tree models have the advantage of not requiring any functional form for the predictors and of not assuming additivity of predictors (i.e., recursive partitioning can identify complex interactions). Trees can deal with missing data flexibly. They have the disadvantages of not utilizing continuous variables effectively and of overfitting in three directions: searching for best predictors, for best splits, and searching multiple times. The penalty for the extreme amount of data searching required by recursive partitioning surfaces when the tree does not cross-validate optimally until it is pruned all the way back to two or three splits. Thus reliable trees are often not very discriminating.

9

Tree models are especially useful in messy situations or settings in which overfitting is not so problematic, such as confounder adjustment using propensity scores¹¹⁷ or in missing value imputation. A major advantage of tree modeling is savings of analyst time, but this is offset by the underfitting needed to make trees validate.

2.6 Multiple Degree of Freedom Tests of Association

When a factor is a linear or binary term in the regression model, the test of association for that factor with the response involves testing only a single regression parameter. Nominal factors and predictors that are represented as a quadratic or spline function require multiple regression parameters to be

tested simultaneously in order to assess association with the response. For a nominal factor having k levels, the overall ANOVA-type test with $k - 1$ d.f. tests whether there are any differences in responses between the k categories. It is recommended that this test be done before attempting to interpret individual parameter estimates. If the overall test is not significant, it can be dangerous to rely on individual pairwise comparisons because the type I error will be increased. Likewise, for a continuous predictor for which linearity is not assumed, all terms involving the predictor should be tested simultaneously to check whether the factor is associated with the outcome. This test should precede the test for linearity and should usually precede the attempt to eliminate nonlinear terms. For example, in the model

$$C(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_2^2, \quad (2.30)$$

one should test $H_0 : \beta_2 = \beta_3 = 0$ with 2 d.f. to assess association between X_2 and outcome. In the five-knot restricted cubic spline model

$$C(Y|X) = \beta_0 + \beta_1 X + \beta_2 X' + \beta_3 X'' + \beta_4 X''', \quad (2.31)$$

the hypothesis $H_0 : \beta_1 = \dots = \beta_4 = 0$ should be tested with 4 d.f. to assess whether there is any association between X and Y . If this 4 d.f. test is insignificant, it is dangerous to interpret the shape of the fitted spline function because the hypothesis that the overall function is flat has not been rejected.

A dilemma arises when an overall test of association, say one having 4 d.f., is insignificant, the 3 d.f. test for linearity is insignificant, but the 1 d.f. test for linear association, after deleting nonlinear terms, becomes significant. Had the test for linearity been borderline significant, it would not have been warranted to drop these terms in order to test for a linear association. But with the evidence for nonlinearity not very great, one could attempt to test for association with 1 d.f. This however is not fully justified, because the 1 d.f. test statistic does not have a χ^2 distribution with 1 d.f. since pretesting was done. The original 4 d.f. test statistic does have a χ^2 distribution with 4 d.f. because it was for a pre-specified test.

For quadratic regression, Grambsch and O'Brien²³⁴ showed that the 2 d.f. test of association is nearly optimal when pretesting is done, even when the true relationship is linear. They considered an ordinary regression model $E(Y|X) = \beta_0 + \beta_1 X + \beta_2 X^2$ and studied tests of association between X and Y . The strategy they studied was as follows. First, fit the quadratic model and obtain the partial test of $H_0 : \beta_2 = 0$, that is the test of linearity. If this partial F -test is significant at the $\alpha = 0.05$ level, report as the final test of association between X and Y the 2 d.f. F -test of $H_0 : \beta_1 = \beta_2 = 0$. If the test of linearity is insignificant, the model is refitted without the quadratic term and the test of association is then a 1 d.f. test, $H_0 : \beta_1 = 0 | \beta_2 = 0$. Grambsch and O'Brien demonstrated that the type I error from this two-stage test is greater than the stated α , and in fact a fairly accurate P -value can be obtained if it is computed from an F distribution with 2 numerator

d.f. even when testing at the second stage. This is because in the original 2 d.f. test of association, the 1 d.f. corresponding to the nonlinear effect is deleted if the nonlinear effect is very small; that is one is retaining the most significant part of the 2 d.f. F statistic.

If we use a 2 d.f. F critical value to assess the X effect even when X^2 is not in the model, it is clear that the two-stage approach can only lose power and hence it has no advantage whatsoever. That is because the sum of squares due to regression from the quadratic model is greater than the sum of squares computed from the linear model.

2.7 Assessment of Model Fit

2.7.1 Regression Assumptions

In this section, the regression part of the model is isolated, and methods are described for validating the regression assumptions or modifying the model to meet the assumptions. The general linear regression model is

$$C(Y|X) = X\beta = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k. \quad (2.32)$$

The assumptions of linearity and additivity need to be verified. We begin with a special case of the general model,

$$C(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2, \quad (2.33)$$

where X_1 is binary and X_2 is continuous. One needs to verify that the property of the response $C(Y)$ is related to X_1 and X_2 according to Figure 2.4.

There are several methods for checking the fit of this model. The first method below is based on critiquing the simple model, and the other methods directly “estimate” the model.

1. Fit the simple linear additive model and critically examine residual plots for evidence of systematic patterns. For least squares fits one can compute estimated residuals $e = Y - X\hat{\beta}$ and box plots of e stratified by X_1 and scatterplots of e versus X_1 and \hat{Y} with trend curves. If one is assuming constant conditional variance of Y , the spread of the residual distribution against each of the variables can be checked at the same time. If the normality assumption is needed (i.e., if significance tests or confidence limits are used), the distribution of e can be compared with a normal distribution with mean zero. **Advantage:** Simplicity. **Disadvantages:** Standard residuals can only be computed for continuous uncensored response variables. The judgment of non-randomness is largely subjective, it is difficult to detect interaction, and if interaction is present it is difficult to check any of the other assumptions. Unless trend lines are added to plots, pat-

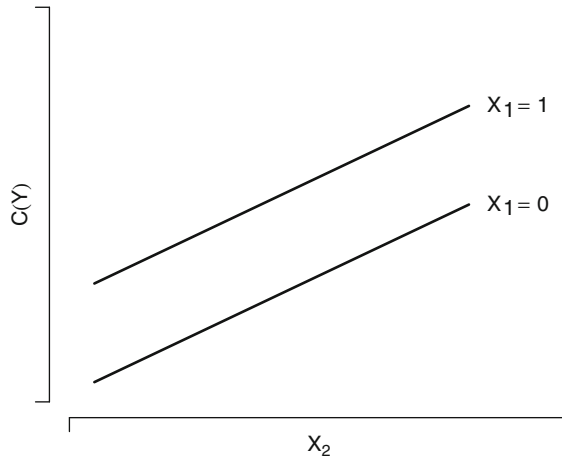


Fig. 2.4 Regression assumptions for one binary and one continuous predictor

terns may be difficult to discern if the sample size is very large. Detecting patterns in residuals does not always inform the analyst of what corrective action to take, although partial residual plots can be used to estimate the needed transformations if interaction is absent.

2. Make a scatterplot of Y versus X_2 using different symbols according to values of X_1 . **Advantages:** Simplicity, and one can sometimes see all regression patterns including interaction. **Disadvantages:** Scatterplots cannot be drawn for binary, categorical, or censored Y . Patterns are difficult to see if relationships are weak or if the sample size is very large.
3. Stratify the sample by X_1 and quantile groups (e.g., deciles) of X_2 . Within each $X_1 \times X_2$ stratum an estimate of $C(Y|X_1, X_2)$ is computed. If X_1 is continuous, the same method can be used after grouping X_1 into quantile groups. **Advantages:** Simplicity, ability to see interaction patterns, can handle censored Y if care is taken. **Disadvantages:** Subgrouping requires relatively large sample sizes and does not use continuous factors effectively as it does not attempt any interpolation. The ordering of quantile groups is not utilized by the procedure. Subgroup estimates have low precision (see p. 488 for an example). Each stratum must contain enough information to allow trends to be apparent above noise in the data. The method of grouping chosen (e.g., deciles vs. quintiles vs. rounding) can alter the shape of the plot.
4. Fit a nonparametric smoother separately for levels of X_1 (Section 2.4.7) relating X_2 to Y . **Advantages:** All regression aspects of the model can be summarized efficiently with minimal assumptions. **Disadvantages:** Does not easily apply to censored Y , and does not easily handle multiple predictors.

5. Fit a flexible parametric model that allows for most of the departures from the linear additive model that you wish to entertain. **Advantages:** One framework is used for examining the model assumptions, fitting the model, and drawing formal inference. Degrees of freedom are well defined and all aspects of statistical inference “work as advertised.” **Disadvantages:** Complexity, and it is generally difficult to allow for interactions when assessing patterns of effects.

The first four methods each have the disadvantage that if confidence limits or formal inferences are desired it is difficult to know how many degrees of freedom were effectively used so that, for example, confidence limits will have the stated coverage probability. For method five, the restricted cubic spline function is an excellent tool for estimating the true relationship between X_2 and $C(Y)$ for continuous variables without assuming linearity. By fitting a model containing X_2 expanded into $k - 1$ terms, where k is the number of knots, one can obtain an estimate of the function of X_2 that could be used linearly in the model:

$$\begin{aligned}\hat{C}(Y|X) &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X'_2 + \hat{\beta}_4 X''_2 \\ &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{f}(X_2),\end{aligned}\tag{2.34}$$

where

$$\hat{f}(X_2) = \hat{\beta}_2 X_2 + \hat{\beta}_3 X'_2 + \hat{\beta}_4 X''_2,\tag{2.35}$$

and X'_2 and X''_2 are constructed spline variables (when $k = 4$) as described previously. We call $\hat{f}(X_2)$ the spline-estimated transformation of X_2 . Plotting the estimated spline function $\hat{f}(X_2)$ versus X_2 will generally shed light on how the effect of X_2 should be modeled. If the sample is sufficiently large, the spline function can be fitted separately for $X_1 = 0$ and $X_1 = 1$, allowing detection of even unusual interaction patterns. A formal test of linearity in X_2 is obtained by testing $H_0 : \beta_3 = \beta_4 = 0$, using a computationally efficient score test, for example (Section 9.2.3).

If the model is nonlinear in X_2 , either a transformation suggested by the spline function plot (e.g., $\log(X_2)$) or the spline function itself (by placing X_2 , X'_2 , and X''_2 simultaneously in any model fitted) may be used to describe X_2 in the model. If a tentative transformation of X_2 is specified, say $g(X_2)$, the adequacy of this transformation can be tested by expanding $g(X_2)$ in a spline function and testing for linearity. If one is concerned only with prediction and not with statistical inference, one can attempt to find a simplifying transformation for a predictor by plotting $g(X_2)$ against $\hat{f}(X_2)$ (the estimated spline transformation) for a variety of g , seeking a linearizing transformation of X_2 . When there are nominal or binary predictors in the model in addition to the continuous predictors, it should be noted that there are no shape assumptions to verify for the binary/nominal predictors. One need only test for interactions between these predictors and the others.

If the model contains more than one continuous predictor, all may be expanded with spline functions in order to test linearity or to describe nonlinear relationships. If one did desire to assess simultaneously, for example, the linearity of predictors X_2 and X_3 in the presence of a linear or binary predictor X_1 , the model could be specified as

$$\begin{aligned} C(Y|X) = & \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_2' + \beta_4 X_2'' \\ & + \beta_5 X_3 + \beta_6 X_3' + \beta_7 X_3'', \end{aligned} \quad (2.36)$$

where X_2' , X_2'' , X_3' , and X_3'' represent components of four knot restricted cubic spline functions.

The test of linearity for X_2 (with 2 d.f.) is $H_0 : \beta_3 = \beta_4 = 0$. The overall test of linearity for X_2 and X_3 is $H_0 : \beta_3 = \beta_4 = \beta_6 = \beta_7 = 0$, with 4 d.f. But as described further in Section 4.1, even though there are many reasons for allowing relationships to be nonlinear, there are reasons for not testing the nonlinear components for significance, as this might tempt the analyst to simplify the model thus distorting inference.²³⁴ Testing for linearity is usually best done to justify to non-statisticians the need for complexity to explain or predict outcomes.

2.7.2 Modeling and Testing Complex Interactions

For testing interaction between X_1 and X_2 (after a needed transformation may have been applied), often a product term (e.g., $X_1 X_2$) can be added to the model and its coefficient tested. A more general simultaneous test of linearity and lack of interaction for a two-variable model in which one variable is binary (or is assumed linear) is obtained by fitting the model

$$\begin{aligned} C(Y|X) = & \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_2' + \beta_4 X_2'' \\ & + \beta_5 X_1 X_2 + \beta_6 X_1 X_2' + \beta_7 X_1 X_2'' \end{aligned} \quad (2.37)$$

and testing $H_0 : \beta_3 = \dots = \beta_7 = 0$. This formulation allows the shape of the X_2 effect to be completely different for each level of X_1 . There is virtually no departure from linearity and additivity that cannot be detected from this expanded model formulation if the number of knots is adequate and X_1 is binary. For binary logistic models, this method is equivalent to fitting two separate spline regressions in X_2 .

Interactions can be complex when all variables are continuous. An approximate approach is to reduce the variables to two transformed variables, in which case interaction may sometimes be approximated by a single product of the two new variables. A disadvantage of this approach is that the estimates of the transformations for the two variables will be different depending

on whether interaction terms are adjusted for when estimating “main effects.” A good compromise method involves fitting interactions of the form $X_1f(X_2)$ and $X_2g(X_1)$:

$$\begin{aligned}
 C(Y|X) = & \beta_0 + \beta_1X_1 + \beta_2X_1' + \beta_3X_1'' \\
 & + \beta_4X_2 + \beta_5X_2' + \beta_6X_2'' \\
 & + \beta_7X_1X_2 + \beta_8X_1X_2' + \beta_9X_1X_2'' \\
 & + \beta_{10}X_2X_1' + \beta_{11}X_2X_1''
 \end{aligned} \tag{2.38}$$

(for $k = 4$ knots for both variables). The test of additivity is $H_0 : \beta_7 = \beta_8 = \dots = \beta_{11} = 0$ with 5 d.f. A test of lack of fit for the simple product interaction with X_2 is $H_0 : \beta_8 = \beta_9 = 0$, and a test of lack of fit for the simple product interaction with X_1 is $H_0 : \beta_{10} = \beta_{11} = 0$.

A general way to model and test interactions, although one requiring a larger number of parameters to be estimated, is based on modeling the $X_1 \times X_2 \times Y$ relationship with a smooth three-dimensional surface. A cubic spline surface can be constructed by covering the $X_1 - X_2$ plane with a grid and fitting a patch-wise cubic polynomial in two variables. The grid is $(u_i, v_j), i = 1, \dots, k, j = 1, \dots, k$, where knots for X_1 are (u_1, \dots, u_k) and knots for X_2 are (v_1, \dots, v_k) . The number of parameters can be reduced by constraining the surface to be of the form $aX_1 + bX_2 + cX_1X_2$ in the lower left and upper right corners of the plane. The resulting restricted cubic spline surface is described by a multiple regression model containing spline expansions in X_1 and X_2 and all cross-products of the restricted cubic spline components (e.g., X_1X_2'). If the same number of knots k is used for both predictors, the number of interaction terms is $(k - 1)^2$. Examples of various ways of modeling interaction are given in Chapter 10. Spline functions made up of cross-products of all terms of individual spline functions are called *tensor splines*.^{50, 274}

11

The presence of more than two predictors increases the complexity of tests for interactions because of the number of two-way interactions and because of the possibility of interaction effects of order higher than two. For example, in a model containing age, sex, and diabetes, the important interaction could be that older male diabetics have an exaggerated risk. However, higher-order interactions are often ignored unless specified a priori based on knowledge of the subject matter. Indeed, the number of two-way interactions alone is often too large to allow testing them all with reasonable power while controlling multiple comparison problems. Often, the only two-way interactions we can afford to test are those that were thought to be important before examining the data. A good approach is to test for all such pre-specified interaction effects with a single global (pooled) test. Then, unless interactions involving only one of the predictors are of special interest, one can either drop all interactions or retain all of them.

For some problems a reasonable approach is, for each predictor separately, to test simultaneously the joint importance of all interactions involving that predictor. For p predictors this results in p tests each with $p - 1$ degrees of freedom. The multiple comparison problem would then be reduced from $p(p - 1)/2$ tests (if all two-way interactions were tested individually) to p tests.

In the fields of biostatistics and epidemiology, some types of interactions that have consistently been found to be important in predicting outcomes and thus may be pre-specified are the following.

1. Interactions between treatment and the severity of disease being treated. Patients with little disease can receive little benefit.
2. Interactions involving age and risk factors. Older subjects are generally less affected by risk factors. They had to have been robust to survive to their current age with risk factors present.
3. Interactions involving age and type of disease. Some diseases are incurable and have the same prognosis regardless of age. Others are treatable or have less effect on younger patients.
4. Interactions between a measurement and the state of a subject during a measurement. Respiration rate measured during sleep may have greater predictive value and thus have a steeper slope versus outcome than respiration rate measured during activity.
5. Interaction between menopausal status and treatment or risk factors.
6. Interactions between race and disease.
7. Interactions between calendar time and treatment. Some treatments have learning curves causing secular trends in the associations.
8. Interactions between month of the year and other predictors, due to seasonal effects.
9. Interaction between the quality and quantity of a symptom, for example, daily frequency of chest pain \times severity of a typical pain episode.
10. Interactions between study center and treatment.

12

2.7.3 Fitting Ordinal Predictors

For the case of an ordinal predictor, spline functions are not useful unless there are so many categories that in essence the variable is continuous. When the number of categories k is small (three to five, say), the variable is usually modeled as a polytomous factor using indicator variables or equivalently as one linear term and $k - 2$ indicators. The latter coding facilitates testing for linearity. For more categories, it may be reasonable to stratify the data by levels of the variable and to compute summary statistics (e.g., logit proportions for a logistic model) or to examine regression coefficients associated with indicator variables over categories. Then one can attempt to summarize the pattern with a linear or some other simple trend. Later hypothesis tests

must take into account this data-driven scoring (by using > 1 d.f., for example), but the scoring can save degrees of freedom when testing for interaction with other factors. In one dataset, the number of comorbid diseases was used to summarize the risk of a set of diseases that was too large to model. By plotting the logit of the proportion of deaths versus the number of diseases, it was clear that the square of the number of diseases would properly score the variables.

Sometimes it is useful to code an ordinal predictor with $k - 1$ indicator variables of the form $[X \geq v_j]$, where $j = 2, \dots, k$ and $[h]$ is 1 if h is true, 0 otherwise.⁶⁴⁸ Although a test of linearity does not arise immediately from this coding, the regression coefficients are interpreted as amounts of change from the previous category. A test of whether the last m categories can be combined with the category $k - m$ does follow easily from this coding.

2.7.4 *Distributional Assumptions*

The general linear regression model is stated as $C(Y|X) = X\beta$ to highlight its regression assumptions. For logistic regression models for binary or nominal responses, there is no distributional assumption if simple random sampling is used and subjects' responses are independent. That is, the binary logistic model and all of its assumptions are contained in the expression $\text{logit}\{Y = 1|X\} = X\beta$. For ordinary multiple regression with constant variance σ^2 , we usually assume that $Y - X\beta$ is normally distributed with mean 0 and variance σ^2 . This assumption can be checked by estimating β with $\hat{\beta}$ and plotting the overall distribution of the residuals $Y - X\hat{\beta}$, the residuals against \hat{Y} , and the residuals against each X . For the latter two, the residuals should be normally distributed within each neighborhood of \hat{Y} or X . A weaker requirement is that the overall distribution of residuals is normal; this will be satisfied if all of the stratified residual distributions are normal. Note a hidden assumption in both models, namely, that there are no omitted predictors. Other models, such as the Weibull survival model or the Cox¹³² proportional hazards model, also have distributional assumptions that are not fully specified by $C(Y|X) = X\beta$. However, regression and distributional assumptions of some of these models are encapsulated by

$$C(Y|X) = C(Y = y|X) = d(y) + X\beta \quad (2.39)$$

for some choice of C . Here $C(Y = y|X)$ is a property of the response Y evaluated at $Y = y$, given the predictor values X , and $d(y)$ is a component of the distribution of Y . For the Cox proportional hazards model, $C(Y = y|X)$ can be written as the log of the hazard of the event at time y , or equivalently as the log of the $-\log$ of the survival probability at time y , and $d(y)$ can be thought of as a log hazard function for a "standard" subject.

If we evaluated the property $C(Y = y|X)$ at predictor values X^1 and X^2 , the difference in properties is

$$\begin{aligned} C(Y = y|X^1) - C(Y = y|X^2) &= d(y) + X^1\beta & (2.40) \\ &= [d(y) + X^2\beta] \\ &= (X^1 - X^2)\beta, \end{aligned}$$

which is independent of y . One way to verify part of the distributional assumption is to estimate $C(Y = y|X^1)$ and $C(Y = y|X^2)$ for set values of X^1 and X^2 using a method that does not make the assumption, and to plot $C(Y = y|X^1) - C(Y = y|X^2)$ versus y . This function should be flat if the distributional assumption holds. The assumption can be tested formally if $d(y)$ can be generalized to be a function of X as well as y . A test of whether $d(y|X)$ depends on X is a test of one part of the distributional assumption. For example, writing $d(y|X) = d(y) + X\Gamma \log(y)$ where

$$X\Gamma = \Gamma_1 X_1 + \Gamma_2 X_2 + \dots + \Gamma_k X_k \quad (2.41)$$

and testing $H_0 : \Gamma_1 = \dots = \Gamma_k = 0$ is one way to test whether $d(y|X)$ depends on X . For semiparametric models such as the Cox proportional hazards model, the only distributional assumption is the one stated above, namely, that the difference in properties between two subjects depends only on the difference in the predictors between the two subjects. Other, parametric, models assume in addition that the property $C(Y = y|X)$ has a specific shape as a function of y , that is that $d(y)$ has a specific functional form. For example, the Weibull survival model has a specific assumption regarding the shape of the hazard or survival distribution as a function of y .

Assessments of distributional assumptions are best understood by applying these methods to individual models as is demonstrated in later chapters.

2.8 Further Reading

- [1] References [152, 575, 578] have more information about cubic splines.
- [2] See Smith⁵⁷⁸ for a good overview of spline functions.
- [3] More material about natural splines may be found in de Boor¹⁵². McNeil et al.⁴⁵¹ discuss the overall smoothness of natural splines in terms of the integral of the square of the second derivative of the regression function, over the range of the data. Govindarajulu et al.²³⁰ compared restricted cubic splines, penalized splines, and fractional polynomial⁵³² fits and found that the first two methods agreed with each other more than with estimated fractional polynomials.
- [4] A tutorial on restricted cubic splines is in [271].
- [5] Durrleman and Simon¹⁶⁸ provide examples in which knots are allowed to be estimated as free parameters, jointly with the regression coefficients. They found that even though the “optimal” knots were often far from a priori knot locations, the model fits were virtually identical.

- [6] Contrast Hastie and Tibshirani's generalized nonparametric additive models²⁷⁵ with Stone and Koo's⁵⁹⁵ additive model in which each continuous predictor is represented with a restricted cubic spline function.
- [7] Gray^{237, 238} provided some comparisons with ordinary regression splines, but he compared penalized regression splines with non-restricted splines with only two knots. Two knots were chosen so as to limit the degrees of freedom needed by the regression spline method to a reasonable number. Gray argued that regression splines are sensitive to knot locations, and he is correct when only two knots are allowed and no linear tail restrictions are imposed. Two knots also prevent the (ordinary maximum likelihood) fit from utilizing some local behavior of the regression relationship. For penalized likelihood estimation using B-splines, Gray²³⁸ provided extensive simulation studies of type I and II error for testing association in which the true regression function, number of knots, and amount of likelihood penalization were varied. He studied both normal regression and Cox regression.
- [8] Breiman et al.'s original CART method⁶⁹ used the Gini criterion for splitting. Later work has used log-likelihoods.¹⁰⁹ Segal,⁵⁶² LeBlanc and Crowley,³⁸⁹ and Ciampi et al.^{107, 108} and Keleş and Segal³⁴² have extended recursive partitioning to censored survival data using the log-rank statistic as the criterion. Zhang⁶⁸² extended tree models to handle multivariate binary responses. Schmoor et al.⁵⁵⁶ used a more general splitting criterion that is useful in therapeutic trials, namely, a Cox test for main and interacting effects. Davis and Anderson¹⁴⁹ used an exponential survival model as the basis for tree construction. Ahn and Loh⁷ developed a Cox proportional hazards model adaptation of recursive partitioning along with bootstrap and cross-validation-based methods to protect against "over-splitting." The Cox-based regression tree methods of Ciampi et al.¹⁰⁷ have a unique feature that allows for construction of "treatment interaction trees" with hierarchical adjustment for baseline variables. Zhang et al.⁶⁸³ provided a new method for handling missing predictor values that is simpler than using surrogate splits. See [34, 140, 270, 629] for examples using recursive partitioning for binary responses in which the prediction trees did not validate well.
- [9] ^{443, 629} discuss other problems with tree models.
- [10] For ordinary linear models, the regression estimates are the same as obtained with separate fits, but standard errors are different (since a pooled standard error is used for the combined fit). For Cox¹³² regression, separate fits can be slightly different since each subset would use a separate ranking of Y .
- [11] Gray's penalized fixed-knot regression splines can be useful for estimating joint effects of two continuous variables while allowing the analyst to control the effective number of degrees of freedom in the fit [237, 238, Section 3.2]. When Y is a non-censored variable, the local regression model of Cleveland et al.,⁹⁶ a multidimensional scatterplot smoother mentioned in Section 2.4.7, provides a good graphical assessment of the joint effects of several predictors so that the forms of interactions can be chosen. See Wang et al.⁶⁵³ and Gustafson²⁴⁸ for several other flexible approaches to analyzing interactions among continuous variables.
- [12] Study site by treatment interaction is often the interaction that is worried about the most in multi-center randomized clinical trials, because regulatory agencies are concerned with consistency of treatment effects over study centers. However, this type of interaction is usually the weakest and is difficult to assess when there are many centers due to the number of interaction parameters to estimate. Schemper⁵⁴⁵ discusses various types of interactions and a general nonparametric test for interaction.

2.9 Problems

For problems 1 to 3, state each model statistically, identifying each predictor with one or more component variables. Identify and interpret each regression parameter except for coefficients of nonlinear terms in spline functions. State each hypothesis below as a formal statistical hypothesis involving the proper parameters, and give the (numerator) degrees of freedom of the test. State alternative hypotheses carefully with respect to unions or intersections of conditions and list the type of alternatives to the null hypothesis that the test is designed to detect.^c

1. A property of Y such as the mean is linear in age and blood pressure and there may be an interaction between the two predictors. Test H_0 : there is no interaction between age and blood pressure. Also test H_0 : blood pressure is not associated with Y (in any fashion). State the effect of blood pressure as a function of age, and the effect of age as a function of blood pressure.
2. Consider a linear additive model involving three treatments (control, drug Z, and drug Q) and one continuous adjustment variable, age. Test H_0 : treatment group is not associated with response, adjusted for age. Also test H_0 : response for drug Z has the same property as the response for drug Q, adjusted for age.
3. Consider models each with two predictors, temperature and white blood count (WBC), for which temperature is always assumed to be linearly related to the appropriate property of the response, and WBC may or may not be linear (depending on the particular model you formulate for each question). Test:
 - a. H_0 : WBC is not associated with the response versus H_a : WBC is linearly associated with the property of the response.
 - b. H_0 : WBC is not associated with Y versus H_a : WBC is quadratically associated with Y . Also write down the formal test of linearity against this quadratic alternative.
 - c. H_0 : WBC is not associated with Y versus H_a : WBC related to the property of the response through a smooth spline function; for example, for WBC the model requires the variables WBC, WBC', and WBC'' where WBC' and WBC'' represent nonlinear components (if there are four knots in a restricted cubic spline function). Also write down the formal test of linearity against this spline function alternative.
 - d. Test for a lack of fit (combined nonlinearity or non-additivity) in an overall model that takes the form of an interaction between temperature and WBC, allowing WBC to be modeled with a smooth spline function.
4. For a fitted model $Y = a + bX + cX^2$ derive the estimate of the effect on Y of changing X from x_1 to x_2 .

^c In other words, under what assumptions does the test have maximum power?

5. In “The Class of 1988: A Statistical Portrait,” the College Board reported mean SAT scores for each state. Use an ordinary least squares multiple regression model to study the mean verbal SAT score as a function of the percentage of students taking the test in each state. Provide plots of fitted functions and defend your choice of the “best” fit. Make sure the shape of the chosen fit agrees with what you know about the variables. Add the raw data points to plots.
- Fit a linear spline function with a knot at $X = 50\%$. Plot the data and the fitted function and do a formal test for linearity and a test for association between X and Y . Give a detailed interpretation of the estimated coefficients in the linear spline model, and use the partial t -test to test linearity in this model.
 - Fit a restricted cubic spline function with knots at $X = 6, 12, 58,$ and 68% (not percentile).^d Plot the fitted function and do a formal test of association between X and Y . Do two tests of linearity that test the same hypothesis:
 - by using a *contrast* to simultaneously test the correct set of coefficients against zero (done by the `anova` function in `rms`);^e
 - by comparing the R^2 from the complex model with that from a simple linear model using a partial F -test.
 Explain why the tests of linearity have the d.f. they have.
 - Using subject matter knowledge, pick a final model (from among the previous models or using another one) that makes sense.

The data are found in Table 2.4 and may be created in R using the `sat.r` code on the RMS course web site.

- Derive the formulas for the restricted cubic spline component variables without cubing or squaring any terms.
- Prove that each component variable is linear in X when $X \geq t_k$, the last knot, using general principles and not algebra or calculus. Derive an expression for the restricted spline regression function when $X \geq t_k$.
- Consider a two-stage procedure in which one tests for linearity of the effect of a predictor X on a property of the response $C(Y|X)$ against a quadratic alternative. If the two-tailed test of linearity is significant at the α level, a two d.f. test of association between X and Y is done. If the test for linearity is not significant, the square term is dropped and a linear model is fitted. The test of association between X and Y is then (apparently) a one d.f. test.
 - Write a formal expression for the test statistic for association.

^d Note: To pre-specify knots for restricted cubic spline functions, use something like `rms(predictor, c(t1,t2,t3,t4))`, where the knot locations are `t1`, `t2`, `t3`, `t4`.

^e Note that `anova` in `rms` computes all needed test statistics from a single model fit object.

- b. Write an expression for the nominal P -value for testing association using this strategy.
- c. Write an expression for the actual P -value or alternatively for the type-I error if using a fixed critical value for the test of association.
- d. For the same two-stage strategy consider an estimate of the effect on $C(Y|X)$ of increasing X from a to b . Write a brief symbolic algorithm for deriving a true two-sided $1 - \alpha$ confidence interval for the $b : a$ effect (difference in $C(Y)$) using the bootstrap.

Table 2.4 SAT data from the College Board, 1988

% Taking SAT (X)	Mean Verbal Score (Y)	% Taking SAT (X)	Mean Verbal Score (Y)
4	482	24	440
5	498	29	460
5	513	37	448
6	498	43	441
6	511	44	424
7	479	45	417
9	480	49	422
9	483	50	441
10	475	52	408
10	476	55	412
10	487	57	400
10	494	58	401
12	474	59	430
12	478	60	433
13	457	62	433
13	485	63	404
14	451	63	424
14	471	63	430
14	473	64	431
16	467	64	437
17	470	68	446
18	464	69	424
20	471	72	420
22	455	73	432
23	452	81	436

Chapter 3

Missing Data

3.1 Types of Missing Data

There are missing data in the majority of datasets one is likely to encounter. Before discussing some of the problems of analyzing data in which some variables are missing for some subjects, we define some nomenclature. 1

Missing completely at random (MCAR)

Data are missing for reasons that are unrelated to any characteristics or responses for the subject, including the value of the missing value, were it to be known. Examples include missing laboratory measurements because of a dropped test tube (if it was not dropped because of knowledge of any measurements), a study that ran out of funds before some subjects could return for follow-up visits, and a survey in which a subject omitted her response to a question for reasons unrelated to the response she would have made or to any other of her characteristics.

Missing at random (MAR)

Data are not missing at random, but the probability that a value is missing depends on values of variables that were actually measured. As an example, consider a survey in which females are less likely to provide their personal income in general (but the likelihood of responding is independent of her actual income). If we know the sex of every subject and have income levels for some of the females, unbiased sex-specific income estimates can be made. That is because the incomes we do have for some of the females are a random sample of all females' incomes. Another way of saying that a variable is MAR

is that given the values of other available variables, subjects having missing values are only randomly different from other subjects.⁵³⁵ Or to paraphrase Greenland and Finkle,²⁴² for MAR the missingness of a covariable cannot depend on unobserved covariable values; for example whether a predictor is observed cannot depend on another predictor when the latter is missing but it can depend on the latter when it is observed. MAR and MCAR data are also called *ignorable* non-responses.

Informative missing (IM)

The tendency for a variable to be missing is a function of data that are not available, including the case when data tend to be missing if their true values are systematically higher or lower. An example is when subjects with lower income levels or very high incomes are less likely to provide their personal income in an interview. IM is also called nonignorable non-response and missing not at random (MNAR).

IM is the most difficult type of missing data to handle. In many cases, there is no fix for IM nor is there a way to use the data to test for the existence of IM. External considerations must dictate the choice of missing data models, and there are few clues for specifying a model under IM. MCAR is the easiest case to handle. Our ability to correctly analyze MAR data depends on the availability of other variables (the sex of the subject in the example above). Most of the methods available for dealing with missing data assume the data are MAR. Fortunately, even though the MAR assumption is not testable, it may hold approximately if enough variables are included in the imputation models²⁵⁶.

3.2 Prelude to Modeling

No matter whether one deletes incomplete cases, carefully imputes (estimates) missing data, or uses a full maximum likelihood or Bayesian techniques to incorporate partial data, it is beneficial to characterize patterns of missingness using exploratory data analysis techniques. These techniques include binary logistic models and recursive partitioning for predicting the probability that a given variable is missing. Patterns of missingness should be reported to help readers understand the limitations of incomplete data. If you do decide to use imputation, it is also important to describe how variables are simultaneously missing. A cluster analysis of missing value status of all the variables is useful here. This can uncover cases where imputation is not as effective. For example, if the only variable moderately related to diastolic blood pressure is systolic pressure, but both pressures are missing on the same subjects, systolic pressure cannot be used to estimate diastolic blood pressure. R

functions `naclus` and `naplot` in the `Hmisc` package (see p. 142) can help detect how variables are simultaneously missing. Recursive partitioning (regression tree) algorithms (see Section 2.5) are invaluable for describing which kinds of subjects are missing on a variable. Logistic regression is also an excellent tool for this purpose. A later example (p. 302) demonstrates these procedures.

It can also be helpful to explore the distribution of non-missing Y by the number of missing variables in X (including zero, i.e., complete cases on X).

3.3 Missing Values for Different Types of Response Variables

When the response variable Y is collected serially but some subjects drop out of the study before completion, there are many ways of dealing with partial information^{42, 412, 480} including multiple imputation in phases,³⁸¹ or efficiently analyzing all available serial data using a full likelihood model. When Y is the time until an event, there are actually no missing values of Y but follow-up will be curtailed for some subjects. That leaves the case where the response is completely measured once.

It is common practice to discard subjects having missing Y . Before doing so, at minimum an analysis should be done to characterize the tendency for Y to be missing, as just described. For example, logistic regression or recursive partitioning can be used to predict whether Y is missing and to test for systematic tendencies as opposed to Y being missing completely at random. In many models, though, more efficient and less biased estimates of regression coefficients can be made by also utilizing observations missing on Y that are non-missing on X . Hence there is a definite place for imputation of Y . von Hippel⁶⁴⁵ found advantages of using all variables to impute all others, and once imputation is finished, discarding those observations having missing Y . However if missing Y values are MCAR, up-front deletion of cases having missing Y may sometimes be preferred, as imputation requires correct specification of the imputation model.

2

3.4 Problems with Simple Alternatives to Imputation

Incomplete predictor information is a very common missing data problem. Statistical software packages use casewise deletion in handling missing predictors; that is, any subject having *any* predictor or Y missing will be excluded from a regression analysis. Casewise deletion results in regression coefficient estimates that can be terribly biased, imprecise, or both³⁵³. First consider an example where bias is the problem. Suppose that the response is death and

the predictors are age, sex, and blood pressure, and that age and sex were recorded for every subject. Suppose that blood pressure was not measured for a fraction of 0.10 of the subjects, and the most common reason for not obtaining a blood pressure was that the subject was about to die. Deletion of these very sick patients will cause a major bias (downward) in the model's intercept parameter. In general, casewise deletion will bias the estimate of the model's intercept parameter (as well as others) when the probability of a case being incomplete is related to Y and not just to X [422, Example 3.3]. van der Heijden et al.⁶²⁸ discuss how complete case analysis (casewise deletion) usually assumes MCAR.

3

Now consider an example in which casewise deletion of incomplete records is inefficient. The inefficiency comes from the reduction of sample size, which causes standard errors to increase,¹⁶² confidence intervals to widen, and power of tests of association and tests of lack of fit to decrease. Suppose that the response is the presence of coronary artery disease and the predictors are age, sex, LDL cholesterol, HDL cholesterol, blood pressure, triglyceride, and smoking status. Suppose that age, sex, and smoking are recorded for all subjects, but that LDL is missing in 0.18 of the subjects, HDL is missing in 0.20, and triglyceride is missing in 0.21. Assume that all missing data are MCAR and that all of the subjects missing LDL are also missing HDL and that overall 0.28 of the subjects have one or more predictors missing and hence would be excluded from the analysis. If total cholesterol were known on every subject, even though it does not appear in the model, it (along perhaps with age and sex) can be used to estimate (*impute*) LDL and HDL cholesterol and triglyceride, perhaps using regression equations from other studies. Doing the analysis on a "filled in" dataset will result in more precise estimates because the sample size would then include the other 0.28 of the subjects.

In general, observations should only be discarded if the MCAR assumption is justified, there is a rarely missing predictor of overriding importance that cannot be reliably imputed from other information, or if the fraction of observations excluded is very small and the original sample size is large. Even then, there is no advantage of such deletion other than saving analyst time. If a predictor is MAR but its missingness depends on Y , casewise deletion is biased.

The first blood pressure example points out why it can be dangerous to handle missing values by adding a dummy variable to the model. Many analysts would set missing blood pressures to a constant (it doesn't matter which constant) and add a variable to the model such as `is.na(blood.pressure)` in R notation. The coefficient for the latter dummy variable will be quite large in the earlier example, and the model will appear to have great ability to predict death. This is because some of the left-hand side of the model contaminates the right-hand side; that is, `is.na(blood.pressure)` is correlated with death. For categorical variables, another common practice is to add a new category to denote missing, adding one more degree of freedom to the

4

predictor and changing its meaning.^a Jones³²⁶, Allison [12, pp. 9–11], Donders et al.¹⁶¹, Knol et al.³⁵³ and van der Heijden et al.⁶²⁸ describe why both of these missing-indicator methods are invalid even when MCAR holds.

5

3.5 Strategies for Developing an Imputation Model

Except in special circumstances that usually involve only very simple models, the primary alternative to deleting incomplete observations is imputation of the missing values. Many non-statisticians find the notion of estimating data distasteful, but the way to think about imputation of missing values is that “making up” data is better than discarding valuable data. It is especially distressing to have to delete subjects who are missing on an adjustment variable when a major variable of interest is not missing. So one goal of imputation is to use as much information as possible for examining any one predictor’s adjusted association with Y . The overall goal of imputation is to preserve the information and meaning of the non-missing data.

At this point the analyst must make some decisions about the information to use in computing predicted values for missing values.

1. Imputation of missing values for one of the variables can ignore all other information. Missing values can be filled in by sampling non-missing values of the variable, or by using a constant such as the median or mean non-missing value.
2. Imputation algorithms can be based only on external information not otherwise used in the model for Y in addition to variables included in later modeling. For example, family income can be imputed on the basis of location of residence when such information is to remain confidential for other aspects of the analysis or when such information would require too many degrees of freedom to be spent in the ultimate response model.
3. Imputations can be derived by only analyzing interrelationships among the X s.
4. Imputations can use relationships among the X s and between X and Y .
5. Imputations can use X , Y , and auxiliary variables not in the model predicting Y .
6. Imputations can take into account the reason for non-response if known.

The model to estimate the missing values in a sometimes-missing (target) variable should include all variables that are either

^a This may work if values are “missing” because of “not applicable”, e.g. one has a measure of marital happiness, dichotomized as high or low, but the sample contains some unmarried people. One could have a 3-category variable with values high, low, and unmarried (Paul Allison, IMPUTE e-mail list, 4Jul09).

1. related to the missing data mechanism;
2. have distributions that differ between subjects that have the target variable missing and those that have it measured;
3. are associated with the target variable when it is not missing; or
4. are included in the final response model⁴³.

The imputation and analysis (response) models should be “congenial” or the imputation model should be more general than the response model or make well-founded assumptions²⁵⁶.

When a variable, say X_j , is to be included as a predictor of Y , and X_j is sometimes missing, ignoring the relationship between X_j and Y for those observations for which both are known will bias regression coefficients for X_j toward zero in the outcome model.⁴²¹ On the other hand, using Y to singly impute X_j using a conditional mean will cause a large inflation in the apparent importance of X_j in the final model. In other words, when the missing X_j are replaced with a mean that is conditional on Y without a random component, this will result in a falsely strong relationship between the imputed X_j values and Y .

At first glance it might seem that using Y to impute one or more of the X s, even with allowance for the correct amount of random variation, would result in a circular analysis in which the importance of the X s will be exaggerated. But the relationship between X and Y in the subset of imputed observations will only be as strong as the associations between X and Y that are evidenced by the non-missing data. In other words, regression coefficients estimated from a dataset that is completed by imputation will not in general be biased high as long as the imputed values have similar variation as non-missing data values.

The next important decision about developing imputation algorithms is the choice of how missing values are estimated.

1. Missings can be estimated using single “best guesses” (e.g., predicted conditional expected values or means) based on relationships between non-missing values. This is called single imputation of conditional means.
2. Missing X_j (or Y) can be estimated using single individual predicted values, where by predicted value we mean a random variable value from the whole conditional distribution of X_j . If one uses ordinary multiple regression to estimate X_j from Y and the other X s, a random residual would be added to the predicted mean value. If assuming a normal distribution for X_j conditional on the other data, such a residual could be computed by a Gaussian random number generator given an estimate of the residual standard deviation. If normality is not assumed, the residual could be a randomly chosen residual from the actual computed residuals. When m missing values need imputation for X_j , the residuals could be sampled with replacement from the entire vector of residuals as in the bootstrap. Better still according to Rubin and Schenker⁵³⁵ would be to use the “approximate Bayesian bootstrap” which involves sampling n residuals with

- replacement from the original n estimated residuals (from observations not missing on X_j), then sampling m residuals with replacement from the first sampled set.
3. More than one random predicted value (as just defined) can be generated for each missing value. This process is called *multiple imputation* and it has many advantages over the other methods in general. This is discussed in Section 3.8.
 4. Matching methods can be used to obtain random draws of other subject's values to replace missing values. Nearest neighbor matching can be used to select a subject that is "close" to the subject in need of imputation, on the basis of a series of variables. This method requires the analyst to make decisions about what constitutes "closeness." To simplify the matching process into a single dimension, Little⁴²⁰ proposed the *predictive mean matching* method where matching is done on the basis of predicted values from a regression model for predicting the sometimes-missing variable (section 3.7). According to Little, in large samples predictive mean matching may be more robust to model misspecification than the method of adding a random residual to the subject's predicted value, but because of difficulties in finding matches the random residual method may be better in smaller samples. The random residual method may be easier to use when multiple imputations are needed, but care must be taken to create the correct degree of uncertainty in residuals.

What if X_j needs to be imputed for some subjects based on other variables that themselves may be missing on the same subjects missing on X_j ? This is a place where recursive partitioning with "surrogate splits" in case of missing predictors may be a good method for developing imputations (see Section 2.5 and p. 142). If using regression to estimate missing values, an algorithm to cycle through all sometimes-missing variables for multiple iterations may perform well. This algorithm is used by the R `transcan` function described in Section 4.7.4 as well as the to-be-described `aregImpute` function. First, all missing values are initialized to medians (modes for categorical variables). Then every time missing values are estimated for a certain variable, those estimates are inserted the next time the variable is used to predict other sometimes-missing variables.

If you want to assess the importance of a specific predictor that is frequently missing, it is a good idea to perform a sensitivity analysis in which all observations containing imputed values for that predictor are temporarily deleted. The test based on a model that included the imputed values may be diluted by the imputation or it may test the wrong hypothesis, especially if Y is not used in imputing X .

Little argues for down-weighting observations containing imputations, to obtain a more accurate variance-covariance matrix. For the ordinary linear model, the weights have been worked out for some cases [421, p. 1231].

3.6 Single Conditional Mean Imputation

For a continuous or binary X that is unrelated to all other predictor variables, the mean or median may be substituted for missing values without much loss of efficiency,¹⁶² although regression coefficients will be biased low since Y was not utilized in the imputation. When the variable of interest is related to the other X s, it is far more efficient to use an individual predictive model for each X based on the other variables.^{79,525,612} The “best guess” imputation method fills in missings with predicted expected values using the multivariable imputation model based on non-missing data^b. It is true that conditional means are the best estimates of unknown values, but except perhaps for binary logistic regression^{621,623} their use will result in biased estimates and very biased (low) variance estimates. The latter problem arises from the reduced variability of imputed values [174, p. 464].

Tree-based models (Section 2.5) may be very useful for imputation since they do not require linearity or additivity assumptions, although such models often have poor discrimination when they don’t overfit. When a continuous X being imputed needs to be non-monotonically transformed to best relate it to the other X s (e.g., blood pressure vs. heart rate), trees and ordinary regression are inadequate. Here a general transformation modeling procedure (Section 4.7) may be needed.

Schemper et al.^{551,553} proposed imputing missing binary covariables by predicted probabilities. For categorical sometimes-missing variables, imputation models can be derived using polytomous logistic regression or a classification tree method. For missing values, the most likely value for each subject (from the series of predicted probabilities from the logistic or recursive partitioning model) can be substituted to avoid creating a new category that is falsely highly correlated with Y . For an ordinal X , the predicted mean value (possibly rounded to the nearest actual data value) or median value from an ordinal logistic model is sometimes useful.

8

3.7 Predictive Mean Matching

In *predictive mean matching*⁴²² (PMM), one replaces a missing (NA) value for the target variable being imputed with the actual value from a donor observation. Donors are identified by matching in only one dimension, namely the predicted value (e.g., predicted mean) of the target. Key considerations are how to

^b Predictors of the target variable include all the other X s along with auxiliary variables that are not included in the final outcome model, as long as they precede the variable being imputed in the causal chain (unlike with multiple imputation).

1. model the target when it is not NA
2. match donors on predicted values
3. avoid overuse of “good” donors to disallow excessive ties in imputed data
4. account for all uncertainties (section 3.8).

The predictive model for each target variable uses any outcome variables, all predictors in the final outcome model, plus any needed auxiliary variables. The modeling method should be flexible, not assuming linearity. Many methods will suffice; parametric additive models are often good choices. Beauties of PMM include the lack of need for distributional assumptions (as no residuals are calculated), and predicted values need only be monotonically related to real predicted values^c

In the original PMM method the donor for an NA was the complete observation whose predicted target was closest to the predicted value of the target from all complete observations^d. This approach can result in some donors being used repeatedly. This can be addressed by sampling from a multinomial distribution, where the probabilities are scaled distances of all potential donors’ predictions to the predicted value y^* of the missing target. Tukey’s tricube function (used in loess) is a good weighting function, implemented in the `Hmisc aregImpute` function:

$$\begin{aligned} w_i &= (1 - \min(d_i/s, 1))^3, \\ d_i &= |\hat{y}_i - y^*| \\ s &= 0.2 \times \text{mean}|\hat{y}_i - y^*|. \end{aligned} \tag{3.1}$$

s above is a good default scale factor, and the w_i are scaled so that $\sum w_i = 1$.

3.8 Multiple Imputation

Imputing missing values and then doing an ordinary analysis as if the imputed values were real measurements is usually better than excluding subjects with incomplete data. However, ordinary formulas for standard errors and other statistics are invalid unless imputation is taken into account.⁶⁵¹ Methods for properly accounting for having incomplete data can be complex. The bootstrap (described later) is an easy method to implement, but the computations can be slow^e.

^c Thus when modeling binary or categorical targets one can frequently take least squares shortcuts in place of maximum likelihood for binary, ordinal, or multinomial logistic models.

^d [662](#) discusses an alternative method based on choosing a donor observation at random from the q closest matches ($q = 3$, for example).

^e To use the bootstrap to correctly estimate variances of regression coefficients, one must repeat the imputation process and the model fitting perhaps 100 times using a

Multiple imputation uses random draws from the conditional distribution of the target variable given the other variables (and any additional information that is relevant)^{85, 417, 421, 536}. The additional information used to predict the missing values can contain any variables that are potentially predictive, including variables measured in the future; the causal chain is not relevant.^{421, 463} When a regression model is used for imputation, the process involves adding a random residual to the “best guess” for missing values, to yield the same conditional variance as the original variable. Methods for estimating residuals were listed in Section 3.5. To properly account for variability due to unknown values, the imputation is repeated M times, where $M \geq 3$. Each repetition results in a “completed” dataset that is analyzed using the standard method. Parameter estimates are averaged over these multiple imputations to obtain better estimates than those from single imputation. The variance–covariance matrix of the averaged parameter estimates, adjusted for variability due to imputation, is estimated using⁴²²

$$V = M^{-1} \sum_i^M V_i + \frac{M+1}{M} B, \quad (3.2)$$

where V_i is the ordinary complete data estimate of the variance–covariance matrix for the model parameters from the i th imputation, and B is the between-imputation sample variance–covariance matrix, the diagonal entries of which are the ordinary sample variances of the M parameter estimates.

After running `aregImpute` (or `MICE`) you can run the `Hmisc` packages’s `fit.mult.impute` function to fit the chosen model separately for each artificially completed dataset corresponding to each imputation. After `fit.mult.impute` fits all of the models, it averages the sets of regression coefficients and computes variance and covariance estimates that are adjusted for imputation (using Eq. 3.2).

White and Royston⁶⁶¹ provide a method for multiply imputing missing covariate values using censored survival time data in the context of the Cox proportional hazards model.

White et al.⁶⁶² recommend choosing the number of imputations M so that the key inferential statistics are very reproducible should the imputation analysis be repeated. They suggest the use of $100f$ imputations when f is the fraction of cases that are incomplete. See also [85, Section 2.7] and²³². Extreme amount of missing data does not prevent one from using multiple imputation, because alternatives are worse³²¹. Horton and Lipsitz³⁰² also have a good overview of multiple imputation and a review of several software packages that implement PMM.

Caution: Multiple imputation methods can generate imputations having very reasonable distributions but still not having the property that final

resampling procedure^{174, 566} (see Section 5.2). Still, the bootstrap can estimate the right variance for the wrong parameter estimates if the imputations are not done correctly.

response model regression coefficients have nominal confidence interval coverage. Among other things, it is worth checking that imputations generate the correct collinearities among covariates.

3.8.1 *The `aregImpute` and Other Chained Equations Approaches*

A flexible approach to multiple imputation that handles a wide variety of target variables to be imputed and allows for multiple variables to be missing on the same subject is the chained equation method. With a chained equations approach, each target variable is predicted by a regression model conditional on all other variables in the model, plus other variables. An iterative process cycles through all target variables to impute all missing values⁶²⁷. This approach is used in the MICE algorithm (multiple imputation using chained equations) implemented in R and other systems. The chained equation method does not attempt to use the full Bayesian multivariate model for all target variables, which makes it more flexible and easy to use but leaves it open to creating improper imputations, e.g., imputing conflicting values for different target variables. However, simulation studies⁶²⁷ so far have demonstrated very good performance of imputation based on chained equations in non-complex situations.

The `aregImpute` algorithm⁴⁶³ takes all aspects of uncertainty into account using the bootstrap while using the same estimation procedures as `transcan` (section 4.7). Different bootstrap resamples used for each imputation by fitting a flexible additive model on a sample with replacement from the original data. This model is used to predict all of the original missing and non-missing values for the target variable for the current imputation. `aregImpute` uses flexible parametric additive regression spline models to predict target variables. There is an option to allow target variables to be optimally transformed, even non-monotonically (but this can overfit). The function implements regression imputation based on adding random residuals to predicted means, but its real value lies in implementing a wide variety of PMM algorithms.

The default method used by `aregImpute` is (weighted) PMM so that no residuals or distributional assumptions are required. The default PMM matching used is van Buuren’s “Type 1” matching [85, Section 3.4.2] to capture the right amount of uncertainty. Here one computes predicted values for missing values using a regression fit on the bootstrap sample, and finds donor observations by matching those predictions to predictions from potential donors using the regression fit from the original sample of complete observations. When a predictor of the target variable is missing, it is first imputed from its last imputation when it was a target variable. The first 3 iterations

Table 3.1 Summary of Methods for Dealing with Missing Values

Method	Deletion	Single	Multiple
Allows nonrandom missing	–	x	x
Reduces sample size	x	–	–
Apparent S.E. of $\hat{\beta}$ too low	–	x	–
Increases real S.E. of $\hat{\beta}$	x	–	–
$\hat{\beta}$ biased	if not MCAR	x	–

of the process are ignored (“burn-in”). `aregImpute` seems to perform as well as MICE but runs significantly faster and allows for nonlinear relationships.

Here is an example using the R `Hmisc` and `rms` packages.

```
a ← aregImpute(~ age + sex + bp + death +
               heart.attack.before.death,
               data=mydata, n.impute=5)
f ← fit.mult.impute(death ~ rcs(age,3) + sex +
                   rcs(bp,5), lrm, a, data=mydata)
```

3.9 Diagnostics

One diagnostic that can be helpful in assessing the MCAR assumption is to compare the distribution of non-missing Y for those subjects having complete X with those having incomplete X . On the other hand, Yucel and Zaslavsky⁶⁸¹ developed a diagnostic that is useful for checking the imputations themselves. In solving a problem related to imputing binary variables using continuous data models, they proposed a simple approach. Suppose we were interested in the reasonableness of imputed values for a sometimes-missing predictor X_j . Duplicate the entire dataset, but in the duplicated observations set all values of X_j to missing. Develop imputed values for the missing values of X_j , and in the observations of the duplicated portion of the dataset corresponding to originally non-missing values of X_j , compare the distribution of imputed X_j with the original values of X_j .

3.10 Summary and Rough Guidelines

Table 3.1 summarizes the advantages and disadvantages of three methods of dealing with missing data. Here “Single” refers to single conditional mean imputation (which cannot utilize Y) and “Multiple” refers to multiple random-draw imputation (which can incorporate Y).

The following contains crude guidelines. Simulation studies are needed to refine the recommendations. Here f refers to the proportion of observations having *any* variables missing.

$f < 0.03$: It doesn't matter very much how you impute missings or whether you adjust variance of regression coefficient estimates for having imputed data in this case. For continuous variables imputing missings with the median non-missing value is adequate; for categorical predictors the most frequent category can be used. Complete case analysis is also an option here. Multiple imputation may be needed to check that the simple approach "worked."

$f \geq 0.03$: Use multiple imputation with number of imputations equal to $\max(5, 100f)$. Fewer imputations may be possible with very large sample sizes. Type 1 predictive mean matching is usually preferred, with weighted selection of donors. Account for imputation in estimating the covariance matrix for final parameter estimates. Use the t distribution instead of the Gaussian distribution for tests and confidence intervals, if possible, using the estimated d.f. for the parameter estimates.

Multiple predictors frequently missing: More imputations may be required. Perform a "sensitivity to order" analysis by creating multiple imputations using different orderings of sometimes missing variables. It may be beneficial to place the variable with the highest number of NAs first so that initialization of other missing variables to medians will have less impact.

It is important to note that the reasons for missing data are more important determinants of how missing values should be handled than is the quantity of missing values.

If the main interest is prediction and not interpretation or inference about individual effects, it is worth trying a simple imputation (e.g., median or normal value substitution) to see if the resulting model predicts the response almost as well as one developed after using customized imputation. But it is not appropriate to use the dummy variable or extra category method, because these methods steal information from Y and bias all β s. Clark and Altman¹¹⁰ presented a nice example of the use of multiple imputation for developing a prognostic model. Marshall et al.⁴⁴² developed a useful method for obtaining predictions on future observations when some of the needed predictors are unavailable. Their method uses an approximate re-fit of the original model for available predictors only, utilizing only the coefficient estimates and covariance matrix from the original fit. Little and An⁴¹⁸ also have an excellent review of imputation methods and developed several approximate formulas for understanding properties of various estimators. They also developed a method combining imputation of missing values with propensity score modeling of the probability of missingness.

3.11 Further Reading

- [1] These types of missing data are well described in an excellent review article on missing data by Schafer and Graham⁵⁴². A good introductory article on missing data and imputation is by Donders et al.¹⁶¹ and a good overview of multiple imputation is by White et al.⁶⁶² and Harel and Zhou²⁵⁶. Paul Allison's booklet¹² and van Buuren's book⁸⁵ are also excellent practical treatments.
- [2] Crawford et al.¹³⁸ give an example where responses are not MCAR for which deleting subjects with missing responses resulted in a biased estimate of the response distribution. They found that multiple imputation of the response resulted in much improved estimates. Wood et al.⁶⁷³ have a good review of how missing response data are typically handled in randomized trial reports, with recommendations for improvements. Barnes et al.⁴² have a good overview of imputation methods and a comparison of bias and confidence interval coverage for the methods when applied to longitudinal data with a small number of subjects. Twist et al.⁶¹⁷ found instability in using multiple imputation of longitudinal data, and advantages of using instead full likelihood models.
- [3] See van Buuren et al.⁶²⁶ for an example in which subjects having missing baseline blood pressure had shorter survival time. Joseph et al.³²⁷ provide examples demonstrating difficulties with casewise deletion and single imputation, and comment on the robustness of multiple imputation methods to violations of assumptions.
- [4] Another problem with the missingness indicator approach arises when more than one predictor is missing and these predictors are missing on almost the same subjects. The missingness indicator variables will be collinear; that is impossible to disentangle.³²⁶
- [5] See [623, pp. 2645–2646] for several problems with the “missing category” approach. A clear example is in¹⁶¹ where covariates X_1, X_2 have true $\beta_1 = 1, \beta_2 = 0$ and X_1 is MCAR. Adding a missingness indicator for X_1 as a covariate resulted in $\hat{\beta}_1 = 0.55, \hat{\beta}_2 = 0.51$ because in the missing observations the constant X_1 was uncorrelated with X_2 . D'Agostino and Rubin¹⁴⁶ developed methods for propensity score modeling that allow for missing data. They mentioned that extra categories may be added to allow for missing data in propensity models and that adding indicator variables describing patterns of missingness will also allow the analyst to match on missingness patterns when comparing non-randomly assigned treatments.
- [6] Harel and Zhou²⁵⁶ and Siddique⁵⁶⁹ discuss the approximate Bayesian bootstrap further.
- [7] Kalton and Kasprzyk³³² proposed a hybrid approach to imputation in which missing values are imputed with the predicted value for the subject plus the residual from the subject having the closest predicted value to the subject being imputed.
- [8] Miller et al.⁴⁵⁸ studied the effect of ignoring imputation when conditional mean fill-in methods are used, and showed how to formalize such methods using linear models.
- [9] Meng⁴⁵⁵ argues against always separating imputation from final analysis, and in favor of sometimes incorporating weights into the process.
- [10] van Buuren et al.⁶²⁶ presented an excellent case study in multiple imputation in the context of survival analysis. Barzi and Woodward⁴³ present a nice review of multiple imputation with detailed comparison of results (point estimates and confidence limits for the effect of the sometimes-missing predictor) for various imputation methods. Barnard and Rubin⁴¹ derived an estimate of the d.f. associated with the imputation-adjusted variance matrix for use in a t -distribution

approximation for hypothesis tests about imputation-averaged coefficient estimates. When d.f. is not very large, the t approximation will result in more accurate P -values than using a normal approximation that we use with Wald statistics after inserting Equation 3.2 as the variance matrix.

- [11] Little and An⁴¹⁸ present imputation methods based on flexible additive regression models using penalized cubic splines. Horton and Kleinman³⁰¹ compare several software packages for handling missing data and have comparisons of results with that of `aregImpute`. Moons et al.⁴⁶³ compared `aregImpute` with MICE.
- [12] He and Zaslavsky²⁸⁰ formalized the duplication approach to imputation diagnostics.
- [13] A good general reference on missing data is Little and Rubin,⁴²² and Volume 16, Nos. 1 to 3 of *Statistics in Medicine*, a large issue devoted to incomplete covariable data. Vach⁶²⁰ is an excellent text describing properties of various methods of dealing with missing data in binary logistic regression (see also [621,622,624]). These references show how to use maximum likelihood to explicitly model the missing data process. Little and Rubin show how imputation can be avoided if the analyst is willing to assume a multivariate distribution for the joint distribution of X and Y . Since X usually contains a strange mixture of binary, polytomous, and continuous but highly skewed predictors, it is unlikely that this approach will work optimally in many problems. That's the reason the imputation approach is emphasized. See Rubin⁵³⁶ for a comprehensive source on multiple imputation. See Little,⁴¹⁹ Vach and Blettner,⁶²³ Rubin and Schenker,⁵³⁵ Zhou et al.,⁶⁸⁸ Greenland and Finkle,²⁴² and Hunsberger et al.³¹³ for excellent reviews of missing data problems and approaches to solving them. Reilly and Pepe have a nice comparison of the “hot-deck” imputation method with a maximum likelihood-based method.⁵²³ White and Carlin⁶⁶⁰ studied bias of multiple imputation vs. complete case analysis.

3.12 Problems

The SUPPORT Study (Study to Understand Prognoses Preferences Outcomes and Risks of Treatments) was a five-hospital study of 10,000 critically ill hospitalized adults^{f352}. Patients were followed for in-hospital outcomes and for long-term survival. We analyze 35 variables and a random sample of 1000 patients from the study.

1. Explore the variables and patterns of missing data in the SUPPORT dataset.
 - a. Print univariable summaries of all variables. Make a plot (showing all variables on one page) that describes especially the continuous variables.
 - b. Make a plot showing the extent of missing data and tendencies for some variables to be missing on the same patients. Functions in the `Hmisc` package may be useful.

^f The dataset is on the book's dataset wiki and may be automatically fetched over the internet and loaded using the `Hmisc` package's command `getHdata(support)`.

- c. Total hospital costs (variable `totcst`) were estimated from hospital-specific Medicare cost-to-charge ratios. Characterize what kind of patients have missing `totcst`. For this characterization use the following patient descriptors: `age`, `sex`, `dzgroup`, `num.co`, `edu`, `income`, `scoma`, `meanbp`, `hrt`, `resp`, `temp`.
2. Prepare for later development of a model to predict costs by developing reliable imputations for missing costs. Remove the observation having zero `totcst`.^g
 - a. The cost estimates are not available on 105 patients. Total hospital charges (bills) are available on all but 25 patients. Relate these two variables to each other with an eye toward using `charges` to predict `totcst` when `totcst` is missing. Make graphs that will tell whether linear regression or linear regression after taking logs of both variables is better.
 - b. Impute missing total hospital costs in `SUPPORT` based on a regression model relating charges to costs, when charges are available. You may want to use a statement like the following in R:

```
support <- transform(support,
                     totcst = ifelse(is.na(totcst),
                                     (expression_in_charges), totcst))
```

If in the previous problem you felt that the relationship between costs and charges should be based on taking logs of both variables, the “expression in charges” above may look something like `exp(intercept + slope * log(charges))`, where constants are inserted for `intercept` and `slope`.

- c. Compute the likely error in approximating total cost using charges by computing the median absolute difference between predicted and observed total costs in the patients having both variables available. If you used a log transformation, also compute the median absolute percent error in imputing total costs by anti-logging the absolute difference in predicted logs.
3. State briefly why single conditional median^h imputation is OK here.
4. Use `transcan` to develop single imputations for total cost, commenting on the strength of the model fitted by `transcan` as well as how strongly each variable can be predicted from all the others.
5. Use predictive mean matching to multiply impute cost 10 times per missing observation. Describe graphically the distributions of imputed values and briefly compare these to distributions of non-imputed values. State in a

^g You can use the R command `subset(support, is.na(totcst) | totcst > 0)`. The `is.na` condition tells R that it is permissible to include observations having missing `totcst` without setting all columns of such observations to `NA`.

^h We are anti-logging predicted log costs and we assume log cost has a symmetric distribution

simple way what the sample variance of multiple imputations for a single observation of a continuous predictor is approximating.

- Using the multiple imputed values, develop an overall least squares model for total cost (using the log transformation) making optimal use of partial information, with variances computed so as to take imputation (except for cost) into account. The model should use the predictors in Problem 1 and should not assume linearity in any predictor but should assume additivity. Interpret one of the resulting ratios of imputation-corrected variance to apparent variance and explain why ratios greater than one do not mean that imputation is inefficient.

Chapter 4

Multivariable Modeling Strategies

Chapter 2 dealt with aspects of modeling such as transformations of predictors, relaxing linearity assumptions, modeling interactions, and examining lack of fit. Chapter 3 dealt with missing data, focusing on utilization of incomplete predictor information. All of these areas are important in the overall scheme of model development, and they cannot be separated from what is to follow. In this chapter we concern ourselves with issues related to the whole model, with emphasis on deciding on the amount of complexity to allow in the model and on dealing with large numbers of predictors. The chapter concludes with three default modeling strategies depending on whether the goal is prediction, estimation, or hypothesis testing.

1

There are many choices to be made when deciding upon a global modeling strategy, including choice between

- parametric and nonparametric procedures
- parsimony and complexity
- parsimony and good discrimination ability
- interpretable models and black boxes.

This chapter addresses some of these issues. One general theme of what follows is the idea that in statistical inference when a method is capable of worsening performance of an estimator or inferential quantity (i.e., when the method is not systematically biased in one's favor), the analyst is allowed to benefit from the method. Variable selection is an example where the analysis is systematically tilted in one's favor by directly selecting variables on the basis of P -values of interest, and all elements of the final result (including regression coefficients and P -values) are biased. On the other hand, the next section is an example of the “capitalize on the benefit when it works, and the method may hurt” approach because one may reduce the complexity of an apparently weak predictor by removing its most important component—

nonlinear effects—from how the predictor is expressed in the model. The method hides tests of nonlinearity that would systematically bias the final result.

The book’s web site contains a number of simulation studies and references to others that support the advocated approaches.

4.1 Prespecification of Predictor Complexity Without Later Simplification

There are rare occasions in which one actually expects a relationship to be linear. For example, one might predict mean arterial blood pressure at two months after beginning drug administration using as baseline variables the pretreatment mean blood pressure and other variables. In this case one expects the pretreatment blood pressure to linearly relate to follow-up blood pressure, and modeling is simple^a. In the vast majority of studies, however, there is every reason to suppose that all relationships involving nonbinary predictors are nonlinear. In these cases, the only reason to represent predictors linearly in the model is that there is insufficient information in the sample to allow us to reliably fit nonlinear relationships.^b

Supposing that nonlinearities are entertained, analysts often use scatter diagrams or descriptive statistics to decide how to represent variables in a model. The result will often be an adequately fitting model, but confidence limits will be too narrow, P -values too small, R^2 too large, and calibration too good to be true. The reason is that the “phantom d.f.” that represented potential complexities in the model that were dismissed during the subjective assessments are forgotten in computing standard errors, P -values, and R^2_{adj} . The same problem is created when one entertains several transformations (log, $\sqrt{\quad}$, etc.) and uses the data to see which one fits best, or when one tries to simplify a spline fit to a simple transformation.

An approach that solves this problem is to prespecify the complexity with which each predictor is represented in the model, without later simplification of the model. The amount of complexity (e.g., number of knots in spline functions or order of ordinary polynomials) one can afford to fit is roughly related to the “effective sample size.” It is also very reasonable to allow for greater complexity for predictors that are thought to be more powerfully related to Y . For example, errors in estimating the curvature of a regression function are consequential in predicting Y only when the regression is somewhere steep. Once the analyst decides to include a predictor in every model, it is fair to

^a Even then, the two blood pressures may need to be transformed to meet distributional assumptions.

^b Shrinkage (penalized estimation) is a general solution (see Section 4.5). One can always use complex models that are “penalized towards simplicity,” with the amount of penalization being greater for smaller sample sizes.

use general measures of association to quantify the predictive potential for a variable. For example, if a predictor has a low rank correlation with the response, it will not “pay” to devote many degrees of freedom to that predictor in a spline function having many knots. On the other hand, a potent predictor (with a high rank correlation) not known to act linearly might be assigned five knots if the sample size allows.

When the effective sample size available is sufficiently large so that a saturated main effects model may be fitted, a good approach to gauging predictive potential is the following.

- Let all continuous predictors be represented as restricted cubic splines with k knots, where k is the maximum number of knots the analyst entertains for the current problem.
- Let all categorical predictors retain their original categories except for pooling of very low prevalence categories (e.g., ones containing < 6 observations).
- Fit this general main effects model.
- Compute the partial χ^2 statistic for testing the association of each predictor with the response, adjusted for all other predictors. In the case of ordinary regression, convert partial F statistics to χ^2 statistics or partial R^2 values.
- Make corrections for chance associations to “level the playing field” for predictors having greatly varying d.f., e.g., subtract the d.f. from the partial χ^2 (the expected value of χ_p^2 is p under H_0).
- Make certain that tests of nonlinearity are not revealed as this would bias the analyst.
- Sort the partial association statistics in descending order.

Commands in the `rms` package can be used to plot only what is needed. Here is an example for a logistic model.

```
f <- lrm(y ~ sex + race + rcs(age,5) + rcs(weight,5) +
        rcs(height,5) + rcs(blood.pressure,5))
plot(anova(f))
```

This approach, and the rank correlation approach about to be discussed, do not require the analyst to really prespecify predictor complexity, so how are they not biased in our favor? There are two reasons: the analyst has already agreed to retain the variable in the model even if the strength of the association is very low, and the assessment of association does not reveal the degree of nonlinearity of the predictor to allow the analyst to “tweak” the number of knots or to discard nonlinear terms. Any predictive ability a variable might have may be concentrated in its nonlinear effects, so using the total association measure for a predictor to save degrees of freedom by restricting the variable to be linear may result in no predictive ability. Likewise, a low association measure between a categorical variable and Y might lead the analyst to collapse some of the categories based on their frequencies. This often helps, but sometimes the categories that are so combined are the

ones that are most different from one another. So if using partial tests or rank correlation to reduce degrees of freedom can harm the model, one might argue that it is fair to allow this strategy to also benefit the analysis.

When collinearities or confounding are not problematic, a quicker approach based on pairwise measures of association can be useful. This approach will not have numerical problems (e.g., singular covariance matrix). When Y is binary or continuous (but not censored), a good general-purpose measure of association that is useful in making decisions about the number of parameters to devote to a predictor is an extension of Spearman's ρ rank correlation. This is the ordinary R^2 from predicting the rank of Y based on the rank of X and the square of the rank of X . This ρ^2 will detect not only nonlinear relationships (as will ordinary Spearman ρ) but some non-monotonic ones as well. It is important that the ordinary Spearman ρ not be computed, as this would tempt the analyst to simplify the regression function (towards monotonicity) if the generalized ρ^2 does not significantly exceed the square of the ordinary Spearman ρ . For categorical predictors, ranks are not squared but instead the predictor is represented by a series of dummy variables. The resulting ρ^2 is related to the Kruskal–Wallis test. See p. 460 for an example. Note that bivariable correlations can be misleading if marginal relationships vary greatly from ones obtained after adjusting for other predictors.

Once one expands a predictor into linear and nonlinear terms and estimates the coefficients, the best way to understand the relationship between predictors and response is to graph this estimated relationship^c. If the plot appears almost linear or the test of nonlinearity is very insignificant there is a temptation to simplify the model. The Grambsch and O'Brien result described in Section 2.6 demonstrates why this is a bad idea.

From the above discussion a general principle emerges. Whenever the response variable is informally or formally linked, in an unmasked fashion, to particular parameters that may be deleted from the model, special adjustments must be made in P -values, standard errors, test statistics, and confidence limits, in order for these statistics to have the correct interpretation. Examples of strategies that are improper without special adjustments (e.g., using the bootstrap) include examining a frequency table or scatterplot to decide that an association is too weak for the predictor to be included in the model at all or to decide that the relationship appears so linear that all nonlinear terms should be omitted. It is also valuable to consider the reverse situation; that is, one posits a simple model and then additional analysis or outside subject matter information makes the analyst want to generalize the model. Once the model is generalized (e.g., nonlinear terms are added), the test of association can be recomputed using multiple d.f. So another general principle is that when one makes the model more complex, the d.f. properly increases and the new test statistics for association have the claimed

^c One can also perform a joint test of all parameters associated with nonlinear effects. This can be useful in demonstrating to the reader that some complexity was actually needed.

distribution. Thus moving from simple to more complex models presents no problems other than conservatism if the new complex components are truly unnecessary.

4.2 Checking Assumptions of Multiple Predictors Simultaneously

Before developing a multivariable model one must decide whether the assumptions of each continuous predictor can be verified by ignoring the effects of all other potential predictors. In some cases, the shape of the relationship between a predictor and the property of response will be different if an adjustment is made for other correlated factors when deriving regression estimates. Also, failure to adjust for an important factor can frequently alter the nature of the distribution of Y . Occasionally, however, it is unwieldy to deal simultaneously with all predictors at each stage in the analysis, and instead the regression function shapes are assessed separately for each continuous predictor.

4.3 Variable Selection

The material covered to this point dealt with a prespecified list of variables to be included in the regression model. For reasons of developing a concise model or because of a fear of collinearity or of a false belief that it is not legitimate to include “insignificant” regression coefficients when presenting results to the intended audience, stepwise variable selection is very commonly employed. Variable selection is used when the analyst is faced with a series of potential predictors but does not have (or use) the necessary subject matter knowledge to enable her to prespecify the “important” variables to include in the model. But using Y to compute P -values to decide which variables to include is similar to using Y to decide how to pool treatments in a five-treatment randomized trial, and then testing for global treatment differences using fewer than four degrees of freedom.

Stepwise variable selection has been a very popular technique for many years, but if this procedure had just been proposed as a statistical method, it would most likely be rejected because it violates every principle of statistical estimation and hypothesis testing. Here is a summary of the problems with this method.

1. It yields R^2 values that are biased high.
2. The ordinary F and χ^2 test statistics do not have the claimed distribution^{d, 234}. Variable selection is based on methods (e.g., F tests for nested models) that were intended to be used to test only prespecified hypotheses.
3. The method yields standard errors of regression coefficient estimates that are biased low and confidence intervals for effects and predicted values that are falsely narrow.¹⁶
4. It yields P -values that are too small (i.e., there are severe multiple comparison problems) and that do not have the proper meaning, and the proper correction for them is a very difficult problem.
5. It provides regression coefficients that are biased high in absolute value and need shrinkage. Even if only a single predictor were being analyzed and one only reported the regression coefficient for that predictor if its association with Y were “statistically significant,” the estimate of the regression coefficient $\hat{\beta}$ is biased (too large in absolute value). To put this in symbols for the case where we obtain a positive association ($\hat{\beta} > 0$), $E(\hat{\beta}|P < 0.05, \hat{\beta} > 0) > \beta$.¹⁰⁰
6. In observational studies, variable selection to determine confounders for adjustment results in residual confounding²⁴¹.
7. Rather than solving problems caused by collinearity, variable selection is made arbitrary by collinearity.
8. It allows us to not think about the problem.

The problems of P -value-based variable selection are exacerbated when the analyst (as she so often does) interprets the final model as if it were prespecified. Copas and Long¹²⁵ stated one of the most serious problems with stepwise modeling eloquently when they said, “The choice of the variables to be included depends on estimated regression coefficients rather than their true values, and so X_j is more likely to be included if its regression coefficient is over-estimated than if its regression coefficient is underestimated.” Derksen and Keselman¹⁵⁵ studied stepwise variable selection, backward elimination, and forward selection, with these conclusions:

1. “The degree of correlation between the predictor variables affected the frequency with which authentic predictor variables found their way into the final model.
2. The number of candidate predictor variables affected the number of noise variables that gained entry to the model.
3. The size of the sample was of little practical importance in determining the number of authentic variables contained in the final model.

^d Lockhart et al.⁴²⁵ provide an example with $n = 100$ and 10 orthogonal predictors where all true β s are zero. The test statistic for the first variable to enter has type I error of 0.39 when the nominal α is set to 0.05, in line with what one would expect with multiple testing using $1 - 0.95^{10} = 0.40$.

4. The population multiple coefficient of determination could be faithfully estimated by adopting a statistic that is adjusted by the total number of candidate predictor variables rather than the number of variables in the final model.”

They found that variables selected for the final model represented noise 0.20 to 0.74 of the time and that the final model usually contained less than half of the actual number of authentic predictors. Hence there are many reasons for using methods such as full-model fits or data reduction, instead of using any stepwise variable selection algorithm.

If stepwise selection must be used, a global test of no regression should be made before proceeding, simultaneously testing all candidate predictors and having degrees of freedom equal to the number of candidate variables (plus any nonlinear or interaction terms). If this global test is not significant, selection of individually significant predictors is usually not warranted.

The method generally used for such variable selection is forward selection of the most significant candidate or backward elimination of the least significant predictor in the model. One of the recommended stopping rules is based on the “residual χ^2 ” with degrees of freedom equal to the number of candidate variables remaining at the current step. The residual χ^2 can be tested for significance (if one is able to forget that because of variable selection this statistic does not have a χ^2 distribution), or the stopping rule can be based on Akaike’s information criterion (AIC³³), here residual $\chi^2 - 2 \times$ d.f.²⁵⁷ Of course, use of more insight from knowledge of the subject matter will generally improve the modeling process substantially. It must be remembered that no currently available stopping rule was developed for data-driven variable selection. Stopping rules such as AIC or Mallows’ C_p are intended for comparing a limited number of *prespecified* models [66, Section 1.3]^{34e}.

If the analyst insists on basing the stopping rule on P -values, the optimum (in terms of predictive accuracy) α to use in deciding which variables to include in the model is $\alpha = 1.0$ unless there are a few powerful variables and several completely irrelevant variables. A reasonable α that does allow for deletion of *some* variables is $\alpha = 0.5$.⁵⁸⁹ These values are far from the traditional choices of $\alpha = 0.05$ or 0.10 .

^e AIC works successfully when the models being entertained are on a progression defined by a single parameter, e.g. a common shrinkage coefficient or the single number of knots to be used by *all* continuous predictors. AIC can also work when the model that is best by AIC is much better than the runner-up so that if the process were bootstrapped the same model would almost always be found. When used for one variable at a time variable selection. AIC is just a restatement of the P -value, and as such, doesn’t solve the severe problems with stepwise variable selection other than forcing us to use slightly more sensible α values. Burnham and Anderson⁸⁴ recommend selection based on AIC for a limited number of theoretically well-founded models. Some statisticians try to deal with multiplicity problems caused by stepwise variable selection by making α smaller than 0.05. This increases bias by giving variables whose effects are estimated with error a greater relative chance of being selected. Variable selection does not compete well with shrinkage methods that simultaneously model all potential predictors.

4

5

Even though forward stepwise variable selection is the most commonly used method, the step-down method is preferred for the following reasons.

6

1. It usually performs better than forward stepwise methods, especially when collinearity is present.⁴³⁷
2. It makes one examine a full model fit, which is the only fit providing accurate standard errors, error mean square, and P -values.
3. The method of Lawless and Singhal³⁸⁵ allows extremely efficient step-down modeling using Wald statistics, in the context of any fit from least squares or maximum likelihood. This method requires passing through the data matrix only to get the initial full fit.

For a given dataset, bootstrapping (Efron et al.^{150,172,177,178}) can help decide between using full and reduced models. Bootstrapping can be done on the whole model and compared with bootstrapped estimates of predictive accuracy based on stepwise variable selection for each resample. Unless most predictors are either very significant or clearly unimportant, the full model usually outperforms the reduced model.

Full model fits have the advantage of providing meaningful confidence intervals using standard formulas. Altman and Andersen¹⁶ gave an example in which the lengths of confidence intervals of predicted survival probabilities were 60% longer when bootstrapping was used to estimate the simultaneous effects of variability caused by variable selection and coefficient estimation, as compared with confidence intervals computed ignoring how a “final” model came to be. On the other hand, models developed on full fits after data reduction will be optimum in many cases.

7

8

In some cases you may want to use the full model for prediction and variable selection for a “best bet” parsimonious list of independently important predictors. This could be accompanied by a list of variables selected in 50 bootstrap samples to demonstrate the imprecision in the “best bet.”

Sauerbrei and Schumacher⁵⁴¹ present a method to use bootstrapping to actually select the set of variables. However, there are a number of drawbacks to this approach³⁵:

1. The choice of an α cutoff for determining whether a variable is retained in a given bootstrap sample is arbitrary.
2. The choice of a cutoff for the proportion of bootstrap samples for which a variable is retained, in order to include that variable in the final model, is somewhat arbitrary.
3. Selection from among a set of correlated predictors is arbitrary, and all highly correlated predictors may have a low bootstrap selection frequency. It may be the case that none of them will be selected for the final model even though when considered individually each of them may be highly significant.

4. By using the bootstrap to choose variables, one must use the double bootstrap to resample the entire modeling process in order to validate the model and to derive reliable confidence intervals. This may be computationally prohibitive.
5. The bootstrap did not improve upon traditional backward stepdown variable selection. Both methods fail at identifying the “correct” variables.

For some applications the list of variables selected may be stabilized by grouping variables according to subject matter considerations or empirical correlations and testing each related group with a multiple degree of freedom test. Then the entire group may be kept or deleted and, if desired, groups that are retained can be summarized into a single variable or the most accurately measured variable within the group can replace the group. See Section 4.7 for more on this.

Kass and Raftery³³⁷ showed that Bayes factors have several advantages in variable selection, including the selection of less complex models that may agree better with subject matter knowledge. However, as in the case with more traditional stopping rules, the final model may still have regression coefficients that are too large. This problem is solved by Tibshirani’s *lasso* method,^{608,609} which is a penalized estimation technique in which the estimated regression coefficients are constrained so that the sum of their scaled absolute values falls below some constant k chosen by cross-validation. This kind of constraint forces some regression coefficient estimates to be exactly zero, thus achieving variable selection while shrinking the remaining coefficients toward zero to reflect the overfitting caused by data-based model selection.

A final problem with variable selection is illustrated by comparing this approach with the sensible way many economists develop regression models. Economists frequently use the strategy of deleting only those variables that are “insignificant” and whose regression coefficients have a nonsensible direction. Standard variable selection on the other hand yields biologically implausible findings in many cases by setting certain regression coefficients exactly to zero. In a study of survival time for patients with heart failure, for example, it would be implausible that patients having a specific symptom live exactly as long as those without the symptom just because the symptom’s regression coefficient was “insignificant.” The lasso method shares this difficulty with ordinary variable selection methods and with any method that in the Bayesian context places nonzero prior probability on β being *exactly* zero.

Many papers claim that there were insufficient data to allow for multivariable modeling, so they did “univariable screening” wherein only “significant” variables (i.e., those that are separately significantly associated with Y) were entered into the model.^f This is just a forward stepwise variable selection in

^f This is akin to doing a t -test to compare the two treatments (out of 10, say) that are apparently most different from each other.

which insignificant variables from the first step are not reanalyzed in later steps. Univariable screening is thus even worse than stepwise modeling as it can miss important variables that are only important after adjusting for other variables.⁵⁹⁸ Overall, neither univariable screening nor stepwise variable selection in any way solves the problem of “too many variables, too few subjects,” and they cause severe biases in the resulting multivariable model fits while losing valuable predictive information from deleting marginally significant variables.

10

The online course notes contain a simple simulation study of stepwise selection using R.

4.4 Sample Size, Overfitting, and Limits on Number of Predictors

When a model is fitted that is too complex, that is, has too many free parameters to estimate for the amount of information in the data, the worth of the model (e.g., R^2) will be exaggerated and future observed values will not agree with predicted values. In this situation, *overfitting* is said to be present, and some of the findings of the analysis come from fitting noise and not just signal, or finding spurious associations between X and Y . In this section general guidelines for preventing overfitting are given. Here we concern ourselves with the *reliability* or *calibration* of a model, meaning the ability of the model to predict future observations as well as it appeared to predict the responses at hand. For now we avoid judging whether the model is adequate for the task, but restrict our attention to the likelihood that the model has significantly overfitted the data.

11

In typical low signal-to-noise ratio situations[§], model validations on independent datasets have found the minimum training sample size for which the fitted model has an independently validated predictive discrimination that equals the apparent discrimination seen with in training sample. Similar validation experiments have considered the margin of error in estimating an absolute quantity such as event probability. Studies such as^{268,270,577} have shown that in many situations a fitted regression model is likely to be reliable when the number of predictors (or *candidate* predictors if using variable selection) p is less than $m/10$ or $m/20$, where m is the “limiting sample size” given in Table 4.1. A good average requirement is $p < \frac{m}{15}$. For example, Smith et al.⁵⁷⁷ found in one series of simulations that the expected error in Cox model predicted five-year survival probabilities was below 0.05 when $p < m/20$ for “average” subjects and below 0.10 when $p < m/20$ for “sick”

12

[§] These are situations where the true R^2 is low, unlike tightly controlled experiments and mechanistic models where signal:noise ratios can be quite high. In those situations, many parameters can be estimated from small samples, and the $\frac{m}{15}$ rule of thumb can be significantly relaxed.

Table 4.1 Limiting Sample Sizes for Various Response Variables

Type of Response Variable	Limiting Sample Size m
Continuous	n (total sample size)
Binary	$\min(n_1, n_2)$ ^h
Ordinal (k categories)	$n - \frac{1}{n^2} \sum_{i=1}^k n_i^3$ ⁱ
Failure (survival) time	number of failures ^j

subjects, where m is the number of deaths. For “average” subjects, $m/10$ was adequate for preventing expected errors > 0.1 . **Note:** The number of non-intercept parameters in the model (p) is usually greater than the number of predictors. Narrowly distributed predictor variables (e.g., if all subjects’ ages are between 30 and 45 or only 5% of subjects are female) will require even higher sample sizes. Note that the number of candidate variables must include all variables screened for association with the response, including nonlinear terms and interactions. Instead of relying on the rules of thumb in the table, the shrinkage factor estimate presented in the next section can be used to guide the analyst in determining how many d.f. to model (see p. 87).

Rules of thumb such as the 15:1 rule do not consider that a certain minimum sample size is needed just to estimate basic parameters such as an intercept or residual variance. This is dealt with in upcoming topics about specific models. For the case of ordinary linear regression, estimation of the residual variance is central. All standard errors, P -values, confidence intervals, and R^2 depend on having a precise estimate of σ^2 . The one-sample problem of estimating a mean, which is equivalent to a linear model containing only an intercept, is the easiest case when estimating σ^2 . When a sample of size n is drawn from a normal distribution, a $1 - \alpha$ two-sided confidence interval for the unknown population variance σ^2 is given by

$$\frac{n-1}{\chi_{1-\alpha/2, n-1}^2} s^2 < \sigma^2 < \frac{n-1}{\chi_{\alpha/2, n-1}^2} s^2, \quad (4.1)$$

^h See [487]. If one considers the power of a two-sample binomial test compared with a Wilcoxon test if the response could be made continuous and the proportional odds assumption holds, the effective sample size for a binary response is $3n_1n_2/n \approx 3 \min(n_1, n_2)$ if n_1/n is near 0 or 1 [664, Eq. 10, 15]. Here n_1 and n_2 are the marginal frequencies of the two response levels.

ⁱ Based on the power of a proportional odds model two-sample test when the marginal cell sizes for the response are n_1, \dots, n_k , compared with all cell sizes equal to unity (response is continuous) [664, Eq. 3]. If all cell sizes are equal, the relative efficiency of having k response categories compared with a continuous response is $1 - 1/k^2$ [664, Eq. 14]; for example, a five-level response is almost as efficient as a continuous one if proportional odds holds across category cutoffs.

^j This is approximate, as the effective sample size may sometimes be boosted somewhat by censored observations, especially for non-proportional hazards methods such as Wilcoxon-type tests.⁴⁹

where s^2 is the sample variance and $\chi_{\alpha, n-1}^2$ is the α critical value of the χ^2 distribution with $n - 1$ degrees of freedom. We take the fold-change or multiplicative margin of error (MMOE) for estimating σ to be

$$\sqrt{\max\left(\frac{\chi_{1-\alpha/2, n-1}^2}{n-1}, \frac{n-1}{\chi_{\alpha/2, n-1}^2}\right)} \quad (4.2)$$

To achieve a MMOE of no worse than 1.2 with 0.95 confidence when estimating σ requires a sample size of 70 subjects.

The linear model case is useful for examining $n : p$ ratio another way. As discussed in the next section, R_{adj}^2 is a nearly unbiased estimate of R^2 , i.e., is not inflated by overfitting if the value used for p is “honest”, i.e., includes all variables screened. We can ask the question “for a given R^2 , what ratio of $n : p$ is required so that R_{adj}^2 does not drop by more than a certain relative or absolute amount from the value of R^2 ?” This assessment takes into account that higher signal:noise ratios allow fitting more variables. For example, with

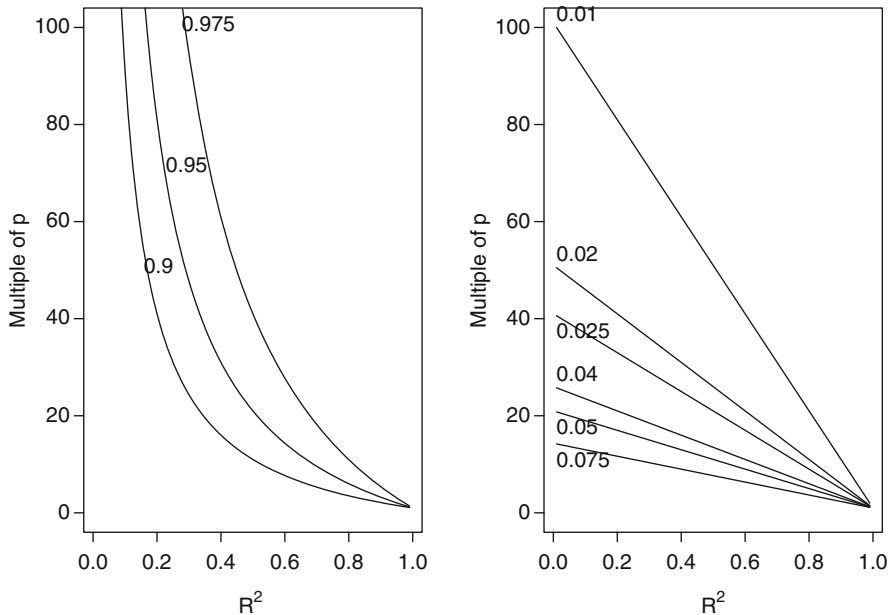


Fig. 4.1 Multiple of p that n must be to achieve a relative drop from R^2 to R_{adj}^2 by the indicated relative factor (left panel, 3 factors) or absolute difference (right panel, 6 decrements)

low R^2 a 100:1 ratio of $n : p$ may be required to prevent R^2 from dropping by more than $\frac{1}{10}$ or by an absolute amount of 0.01. A 15:1 rule would prevent R^2 from dropping by more than 0.075 for low R^2 (Figure 4.1).

4.5 Shrinkage

The term *shrinkage* is used in regression modeling to denote two ideas. The first meaning relates to the slope of a *calibration plot*, which is a plot of observed responses against predicted responses^k. When a dataset is used to fit the model parameters as well as to obtain the calibration plot, the usual estimation process will force the slope of observed versus predicted values to be one. When, however, parameter estimates are derived from one dataset and then applied to predict outcomes on an independent dataset, overfitting will cause the slope of the calibration plot (i.e., the *shrinkage factor*) to be less than one, a result of *regression to the mean*. Typically, low predictions will be too low and high predictions too high. Predictions near the mean predicted value will usually be quite accurate. The second meaning of *shrinkage* is a statistical estimation method that preshrinks regression coefficients towards zero so that the calibration plot for new data will not need shrinkage as its calibration slope will be one.

We turn first to shrinkage as an adverse result of traditional modeling. In ordinary linear regression, we know that all of the coefficient estimates are exactly unbiased estimates of the true effect when the model fits. Isn't the existence of shrinkage and overfitting implying that there is some kind of bias in the parameter estimates? The answer is no because each separate coefficient has the desired expectation. The problem lies in how we use the coefficients. We tend not to pick out coefficients at random for interpretation but we tend to highlight very small and very large coefficients.

A simple example may suffice. Consider a clinical trial with 10 randomly assigned treatments such that the patient responses for each treatment are normally distributed. We can do an ANOVA by fitting a multiple regression model with an intercept and nine dummy variables. The intercept is an unbiased estimate of the mean response for patients on the first treatment, and each of the other coefficients is an unbiased estimate of the difference in mean response between the treatment in question and the first treatment. $\hat{\beta}_0 + \hat{\beta}_1$ is an unbiased estimate of the mean response for patients on the second treatment. But if we plotted the predicted mean response for patients against the observed responses from new data, the slope of this calibration plot would typically be smaller than one. This is because in making this plot we are not picking coefficients at random but we are sorting the coefficients into ascending order. The treatment group having the lowest sample mean response will usually have a higher mean in the future, and the treatment group having the highest sample mean response will typically have a lower mean in the future. The sample mean of the group having the highest sample mean is *not* an unbiased estimate of its population mean.

^k An even more stringent assessment is obtained by stratifying calibration curves by predictor settings.

As an illustration, let us draw 20 samples of size $n = 50$ from a uniform distribution for which the true mean is 0.5. Figure 4.2 displays the 20 means sorted into ascending order, similar to plotting Y versus $\hat{Y} = X\hat{\beta}$ based on least squares after sorting by $X\hat{\beta}$. Bias in the very lowest and highest estimates is evident.

```
set.seed(123)
n <- 50
y <- runif(20*n)
group <- rep(1:20, each=n)
ybar <- tapply(y, group, mean)
ybar <- sort(ybar)
plot(1:20, ybar, type='n', axes=FALSE, ylim=c(.3,.7),
     xlab='Group', ylab='Group Mean')
lines(1:20, ybar)
points(1:20, ybar, pch=20, cex=.5)
axis(2)
axis(1, at=1:20, labels=FALSE)
for(j in 1:20) axis(1, at=j, labels=names(ybar)[j])
abline(h=.5, col=gray(.85))
```

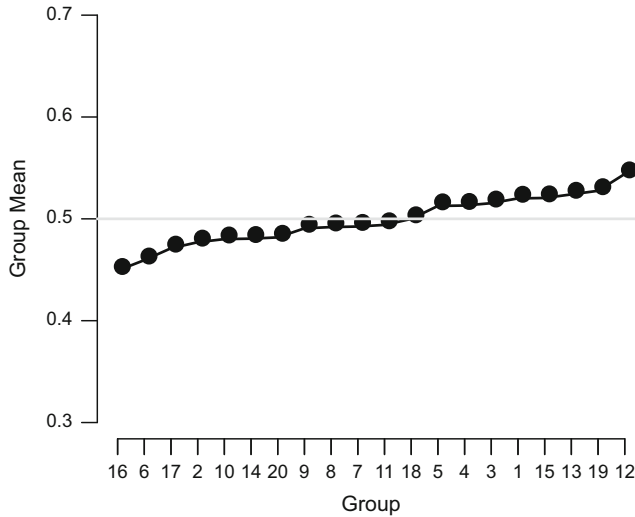


Fig. 4.2 Sorted means from 20 samples of size 50 from a uniform $[0, 1]$ distribution. The reference line at 0.5 depicts the true population value of all of the means.

When we want to highlight a treatment that is not chosen at random (or a priori), the data-based selection of that treatment needs to be compensated for in the estimation process.¹ It is well known that the use of shrinkage

¹ It is interesting that researchers are quite comfortable with adjusting P -values for post hoc selection of comparisons using, for example, the Bonferroni inequality, but they do not realize that post hoc selection of comparisons also biases point estimates.

methods such as the James–Stein estimator to pull treatment means toward the grand mean over all treatments results in estimates of treatment-specific means that are far superior to ordinary stratified means.¹⁷⁶

Turning from a cell means model to the general case where predicted values are general linear combinations $X\hat{\beta}$, the slope γ of properly transformed responses Y against $X\hat{\beta}$ (sorted into ascending order) will be less than one on new data. Estimation of the shrinkage coefficient γ allows quantification of the amount of overfitting present, and it allows one to estimate the likelihood that the model will reliably predict new observations. van Houwelingen and le Cessie [633, Eq. 77] provided a heuristic shrinkage estimate that has worked well in several examples:

$$\hat{\gamma} = \frac{\text{model } \chi^2 - p}{\text{model } \chi^2}, \quad (4.3)$$

where p is the total degrees of freedom for the predictors and $\text{model } \chi^2$ is the likelihood ratio χ^2 statistic for testing the joint influence of all predictors simultaneously (see Section 9.3.1). For ordinary linear models, van Houwelingen and le Cessie proposed a shrinkage factor $\hat{\gamma}$ that can be shown to equal $\frac{n-p-1}{n-1} \frac{R_{\text{adj}}^2}{R^2}$, where the adjusted R^2 is given by

$$R_{\text{adj}}^2 = 1 - (1 - R^2) \frac{n-1}{n-p-1}. \quad (4.4)$$

For such linear models with an intercept β_0 , the shrunk estimate of β is

$$\begin{aligned} \hat{\beta}_0^s &= (1 - \hat{\gamma})\bar{Y} + \hat{\gamma}\hat{\beta}_0 \\ \hat{\beta}_j^s &= \hat{\gamma}\hat{\beta}_j, j = 1, \dots, p, \end{aligned} \quad (4.5)$$

where \bar{Y} is the mean of the response vector. Again, when stepwise fitting is used, the p in these equations is much closer to the number of *candidate* degrees of freedom rather than the number in the “final” model. See Section 5.3 for methods of estimating γ using the bootstrap (p. 115) or cross-validation.

Now turn to the second usage of the term *shrinkage*. Just as clothing is sometimes preshrunk so that it will not shrink further once it is purchased, better calibrated predictions result when shrinkage is built into the estimation process in the first place. The object of shrinking regression coefficient estimates is to obtain a shrinkage coefficient of $\gamma = 1$ on new data. Thus by somewhat discounting $\hat{\beta}$ we make the model underfitted on the data at hand (i.e., apparent $\gamma < 1$) so that on new data extremely low or high predictions are correct.

Ridge regression^{388, 633} is one technique for placing restrictions on the parameter estimates that results in shrinkage. A *ridge parameter* must be chosen to control the amount of shrinkage. Penalized maximum likelihood estimation,^{237, 272, 388, 639} a generalization of ridge regression, is a general shrinkage

procedure. A method such as cross-validation or optimization of a modified AIC must be used to choose an optimal penalty factor. An advantage of penalized estimation is that one can differentially penalize the more complex components of the model such as nonlinear or interaction effects. A drawback of ridge regression and penalized maximum likelihood is that the final model is difficult to validate unbiasedly since the optimal amount of shrinkage is usually determined by examining the entire dataset. Penalization is one of the best ways to approach the “too many variables, too little data” problem. See Section 9.10 for details.

4.6 Collinearity

When at least one of the predictors can be predicted well from the other predictors, the standard errors of the regression coefficient estimates can be inflated and corresponding tests have reduced power.²¹⁷ In stepwise variable selection, collinearity can cause predictors to compete and make the selection of “important” variables arbitrary. Collinearity makes it difficult to estimate and interpret a particular regression coefficient because the data have little information about the effect of changing one variable while holding another (highly correlated) variable constant [101, Chap. 9]. However, collinearity does not affect the joint influence of highly correlated variables when tested simultaneously. Therefore, once groups of highly correlated predictors are identified, the problem can be rectified by testing the contribution of an entire set with a multiple d.f. test rather than attempting to interpret the coefficient or one d.f. test for a single predictor.

Collinearity does not affect predictions made on the same dataset used to estimate the model parameters or on new data that have the same degree of collinearity as the original data [470, pp. 379–381] as long as extreme extrapolation is not attempted. Consider as two predictors the total and LDL cholesterols that are highly correlated. If predictions are made at the same combinations of total and LDL cholesterol that occurred in the training data, no problem will arise. However, if one makes a prediction at an inconsistent combination of these two variables, the predictions may be inaccurate and have high standard errors.

When the ordinary truncated power basis is used to derive component variables for fitting linear and cubic splines, as was described earlier, the component variables can be very collinear. It is very unlikely that this will result in any problems, however, as the component variables are connected algebraically. Thus it is not possible for a combination of, for example, x and $\max(x - 10, 0)$ to be inconsistent with each other. Collinearity problems are then more likely to result from partially redundant subsets of predictors as in the cholesterol example above.

One way to quantify collinearity is with *variance inflation factors* or *VIF*, which in ordinary least squares are diagonals of the inverse of the $X'X$ matrix scaled to have unit variance (except that a column of 1s is retained corresponding to the intercept). Note that some authors compute VIF from the correlation matrix form of the design matrix, omitting the intercept. VIF_i is $1/(1 - R_i^2)$ where R_i^2 is the squared multiple correlation coefficient between column i and the remaining columns of the design matrix. For models that are fitted with maximum likelihood estimation, the information matrix is scaled to correlation form, and VIF is the diagonal of the inverse of this scaled matrix.^{147, 654} Then the VIF are similar to those from a weighted correlation matrix of the original columns in the design matrix. Note that indexes such as VIF are not very informative as some variables are algebraically connected to each other.

16

The SAS `VARCLUS` procedure⁵³⁹ and R `varclus` function can identify collinear predictors. Summarizing collinear variables using a summary score is more powerful and stable than arbitrary selection of one variable in a group of collinear variables (see the next section).

17

4.7 Data Reduction

The sample size need not be as large as shown in Table 4.1 if the model is to be validated independently and if you don't care that the model may fail to validate. However, it is likely that the model will be overfitted and will not validate if the sample size does not meet the guidelines. Use of data reduction methods before model development is strongly recommended if the conditions in Table 4.1 are not satisfied, and if shrinkage is not incorporated into parameter estimation. Methods such as shrinkage and data reduction reduce the effective d.f. of the model, making it more likely for the model to validate on future data. Data reduction is aimed at reducing the number of parameters to estimate in the model, without distorting statistical inference for the parameters. This is accomplished by ignoring Y during data reduction. Manipulations of X in unsupervised learning may result in a loss of information for predicting Y , but when the information loss is small, the gain in power and reduction of overfitting more than offset the loss.

Some available data reduction methods are given below.

1. Use the literature to eliminate unimportant variables.
2. Eliminate variables whose distributions are too narrow.
3. Eliminate candidate predictors that are missing in a large number of subjects, especially if those same predictors are likely to be missing for future applications of the model.
4. Use a statistical data reduction method such as incomplete principal component regression, nonlinear generalizations of principal components such

as principal surfaces, sliced inverse regression, variable clustering, or ordinary cluster analysis on a measure of similarity between variables.

See Chapters 8 and 14 for detailed case studies in data reduction.

4.7.1 Redundancy Analysis

There are many approaches to data reduction. One rigorous approach involves removing predictors that are easily predicted from other predictors, using flexible parametric additive regression models. This approach is unlikely to result in a major reduction in the number of regression coefficients to estimate against Y , but will usually provide insights useful for later data reduction over and above the insights given by methods based on pairwise correlations instead of multiple R^2 .

The `Hmisc redun` function implements the following redundancy checking algorithm.

- Expand each continuous predictor into restricted cubic spline basis functions. Expand categorical predictors into dummy variables.
- Use OLS to predict each predictor with all component terms of all remaining predictors (similar to what the `Hmisc transcan` function does). When the predictor is expanded into multiple terms, use the first canonical variate^m.
- Remove the predictor that can be predicted from the remaining set with the highest adjusted or regular R^2 .
- Predict all remaining predictors from their complement.
- Continue in like fashion until no variable still in the list of predictors can be predicted with an R^2 or adjusted R^2 greater than a specified threshold or until dropping the variable with the highest R^2 (adjusted or ordinary) would cause a variable that was dropped earlier to no longer be predicted at the threshold from the now smaller list of predictors.

Special consideration must be given to categorical predictors. One way to consider a categorical variable redundant is if a linear combination of dummy variables representing it can be predicted from a linear combination of other variables. For example, if there were 4 cities in the data and each city's rainfall was also present as a variable, with virtually the same rainfall reported for all observations for a city, city would be redundant given rainfall (or vice-versa). If two cities had the same rainfall, 'city' might be declared redundant even though tied cities might be deemed non-redundant in another setting. A second, more stringent way to check for redundancy of a categorical predictor is to ascertain whether all dummy variables created from the predictor are individually redundant. The `redun` function implements both approaches.

Examples of use of `redun` are given in two case studies.

^m There is an option to force continuous variables to be linear when they are being predicted.

4.7.2 Variable Clustering

Although the use of subject matter knowledge is usually preferred, statistical clustering techniques can be useful in determining independent dimensions that are described by the entire list of candidate predictors. Once each dimension is scored (see below), the task of regression modeling is simplified, and one quits trying to separate the effects of factors that are measuring the same phenomenon. One type of variable clustering⁵³⁹ is based on a type of oblique-rotation principal component (PC) analysis that attempts to separate variables so that the first PC of each group is representative of that group (the first PC is the linear combination of variables having maximum variance subject to normalization constraints on the coefficients^{142,144}). Another approach, that of doing a hierarchical cluster analysis on an appropriate similarity matrix (such as squared correlations) will often yield the same results. For either approach, it is often advisable to use robust (e.g., rank-based) measures for continuous variables if they are skewed, as skewed variables can greatly affect ordinary correlation coefficients. Pairwise deletion of missing values is also advisable for this procedure—casewise deletion can result in a small biased sample.

20

When variables are not monotonically related to each other, Pearson or Spearman squared correlations can miss important associations and thus are not always good similarity measures. A general and robust similarity measure is Hoeffding's D ,²⁹⁵ which for two variables X and Y is a measure of the agreement between $F(x, y)$ and $G(x)H(y)$, where G, H are marginal cumulative distribution functions and F is the joint CDF. The D statistic will detect a wide variety of dependencies between two variables.

See pp. 330 and 458 for examples of variable clustering.

21

4.7.3 Transformation and Scaling Variables Without Using Y

Scaling techniques often allow the analyst to reduce the number of parameters to fit by estimating transformations for each predictor using only information about associations with other predictors. It may be advisable to cluster variables before scaling so that patterns are derived only from variables that are related. For purely categorical predictors, methods such as correspondence analysis (see, for example, [108, 139, 239, 391, 456]) can be useful for data reduction. Often one can use these techniques to scale multiple dummy variables into a few dimensions. For mixtures of categorical and continuous predictors, qualitative principal component analysis such as the *maximum total variance* (MTV) method of Young et al.^{456, 680} is useful. For the special case of representing a series of variables with one PC, the MTV method is quite easy to implement.

1. Compute PC_1 , the first PC of the variables to reduce X_1, \dots, X_q using the correlation matrix of X s.
2. Use ordinary linear regression to predict PC_1 on the basis of functions of the X s, such as restricted cubic spline functions for continuous X s or a series of dummy variables for polytomous X s. The expansion of each X_j is regressed separately on PC_1 .
3. These separately fitted regressions specify the working transformations of each X .
4. Recompute PC_1 by doing a PC analysis on the transformed X s (predicted values from the fits).
5. Repeat steps 2 to 4 until the proportion of variation explained by PC_1 reaches a plateau. This typically requires three to four iterations.

A transformation procedure that is similar to MTV is the maximum generalized variance (MGV) method due to Sarle [368, pp. 1267–1268]. MGV involves predicting each variable from (the current transformations of) all the other variables. When predicting variable i , that variable is represented as a set of linear and nonlinear terms (e.g., spline components). Analysis of canonical variates²⁷⁹ can be used to find the linear combination of terms for X_i (i.e., find a new transformation for X_i) and the linear combination of the current transformations of all other variables (representing each variable as a single, transformed, variable) such that these two linear combinations have maximum correlation. (For example, if there are only two variables X_1 and X_2 represented as quadratic polynomials, solve for a, b, c, d such that $aX_1 + bX_1^2$ has maximum correlation with $cX_2 + dX_2^2$.) The process is repeated until the transformations converge. The goal of MGV is to transform each variable so that it is most similar to predictions from the other transformed variables. MGV does not use PCs (so one need not precede the analysis by variable clustering), but once all variables have been transformed, you may want to summarize them with the first PC.

The SAS `PRINQUAL` procedure of Kuhfeld³⁶⁸ implements the MTV and MGV methods, and allows for very flexible transformations of the predictors, including monotonic splines and ordinary cubic splines.

A very flexible automatic procedure for transforming each predictor in turn, based on all remaining predictors, is the ACE (alternating conditional expectation) procedure of Breiman and Friedman.⁶⁸ Like SAS `PROC PRINQUAL`, ACE handles monotonically restricted transformations and categorical variables. It fits transformations by maximizing R^2 between one variable and a set of variables. It automatically transforms all variables, using the “super smoother”²⁰⁷ for continuous variables. Unfortunately, ACE does not handle missing values. See Chapter 16 for more about ACE.

It must be noted that at best these automatic transformation procedures generally find only *marginal* transformations, not transformations of each predictor adjusted for the effects of all other predictors. When adjusted transformations differ markedly from marginal transformations, only joint modeling of all predictors (and the response) will find the correct transformations.

Once transformations are estimated using only predictor information, the adequacy of each predictor's transformation can be checked by graphical methods, by nonparametric smooths of transformed X_j versus Y , or by expanding the transformed X_j using a spline function. This approach of checking that transformations are optimal with respect to Y uses the response data, but it accepts the initial transformations unless they are significantly inadequate. If the sample size is low, or if PC_1 for the group of variables used in deriving the transformations is deemed an adequate summary of those variables, that PC_1 can be used in modeling. In that way, data reduction is accomplished two ways: by not using Y to estimate multiple coefficients for a single predictor, and by reducing related variables into a single score, after transforming them. See Chapter 8 for a detailed example of these scaling techniques.

4.7.4 Simultaneous Transformation and Imputation

As mentioned in Chapter 3 (p. 52) if transformations are complex or non-monotonic, ordinary imputation models may not work. SAS PROC PRINQUAL implemented a method for simultaneously imputing missing values while solving for transformations. Unfortunately, the imputation procedure frequently converges to imputed values that are outside the allowable range of the data. This problem is more likely when multiple variables are missing on the same subjects, since the transformation algorithm may simply separate missings and nonmissings into clusters.

A simple modification of the MGCV algorithm of PRINQUAL that simultaneously imputes missing values without these problems is implemented in the R function `transcan`. Imputed values are initialized to medians of continuous variables and the most frequent category of categorical variables. For continuous variables, transformations are initialized to linear functions. For categorical ones, transformations may be initialized to the identify function, to dummy variables indicating whether the observation has the most prevalent categorical value, or to random numbers. Then when using canonical variates to transform each variable in turn, observations that are missing on the current "dependent" variable are excluded from consideration, although missing values for the current set of "predictors" are imputed. Transformed variables are normalized to have mean 0 and standard deviation 1. Although categorical variables are scored using the first canonical variate, `transcan` has an option to use recursive partitioning to obtain imputed values on the original scale (Section 2.5) for these variables. It defaults to imputing categorical variables using the category whose predicted canonical score is closest to the predicted score.

`transcan` uses restricted cubic splines to model continuous variables. It does not implement monotonicity constraints. `transcan` automatically constrains

imputed values (both on transformed and original scales) to be in the same range as non-imputed ones. This adds much stability to the resulting estimates although it can result in a boundary effect. Also, imputed values can optionally be shrunk using Eq. 4.5 to avoid overfitting when developing the imputation models. Optionally, missing values can be set to specified constants rather than estimating them. These constants are ignored during the transformation-estimation phaseⁿ. This technique has proved to be helpful when, for example, a laboratory test is not ordered because a physician thinks the patient has returned to normal with respect to the lab parameter measured by the test. In that case, it's better to use a normal lab value for missings.

The transformation and imputation information created by `transcan` may be used to transform/impute variables in datasets not used to develop the transformation and imputation formulas. There is also an R function to create R functions that compute the final transformed values of each predictor given input values on the original scale.

As an example of non-monotonic transformation and imputation, consider a sample of 1000 hospitalized patients from the SUPPORT^o study.³⁵² Two mean arterial blood pressure measurements were set to missing.

```
require(Hmisc)
getHdata(support) # Get data frame from web site
heart.rate      <- support$hrt
blood.pressure  <- support$meanbp
blood.pressure [400:401]
```

```
Mean Arterial Blood Pressure Day 3
[1] 151 136
```

```
blood.pressure[400:401] <- NA # Create two missings
d <- data.frame(heart.rate, blood.pressure)
par(pch=46) # Figure 4.3
w <- transcan(~ heart.rate + blood.pressure, transformed=TRUE,
              imputed=TRUE, show.na=TRUE, data=d)
```

```
Convergence criterion:2.901 0.035
```

```
0.007
Convergence in 4 iterations
R2 achieved in predicting each variable:

  heart.rate blood.pressure
    0.259      0.259

Adjusted R2:

  heart.rate blood.pressure
    0.254      0.253
```

ⁿ If one were to estimate transformations without removing observations that had these constants inserted for the current Y -variable, the resulting transformations would likely have a spike at $Y = \text{imputation constant}$.

^o Study to Understand Prognoses Preferences Outcomes and Risks of Treatments

```
w$imputed$blood.pressure
```

```
      400      401
132.4057 109.7741
```

```
t <- w$transformed
spe <- round(c(spearman(heart.rate, blood.pressure),
              spearman(t[, 'heart.rate'],
                       t[, 'blood.pressure'])), 2)
```

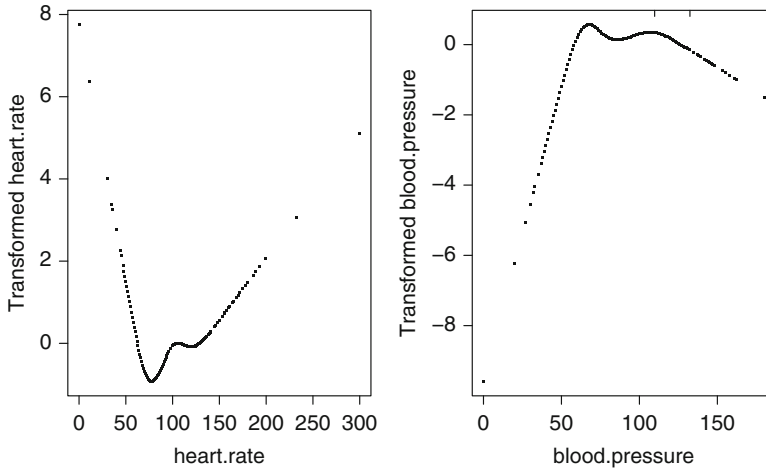


Fig. 4.3 Transformations fitted using `transcan`. Tick marks indicate the two imputed values for blood pressure.

```
plot(heart.rate, blood.pressure) # Figure 4.4
plot(t[, 'heart.rate'], t[, 'blood.pressure'],
     xlab='Transformed hr', ylab='Transformed bp')
```

Spearman's rank correlation ρ between pairs of heart rate and blood pressure was -0.02 , because these variables each require U -shaped transformations. Using restricted cubic splines with five knots placed at default quantiles, `transcan` provided the transformations shown in Figure 4.3. Correlation between transformed variables is $\rho = -0.13$. The fitted transformations are similar to those obtained from relating these two variables to time until death.

4.7.5 Simple Scoring of Variable Clusters

If a subset of the predictors is a series of related dichotomous variables, a simpler data reduction strategy is sometimes employed. First, construct two

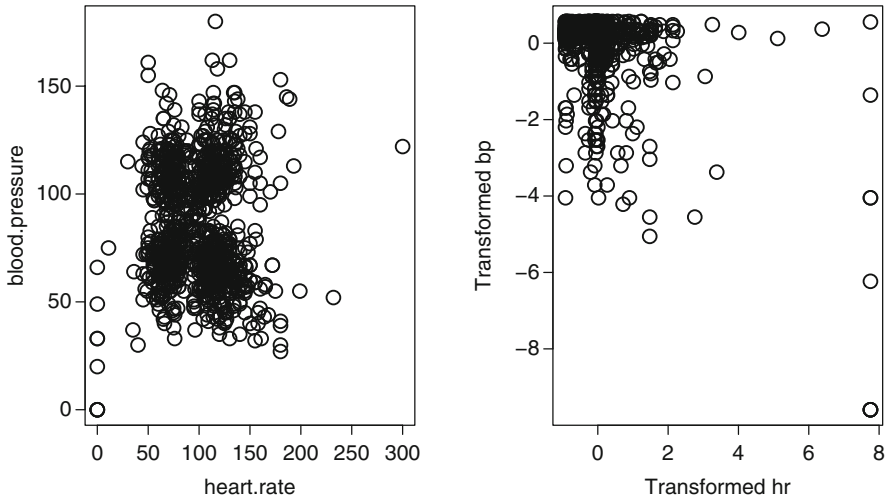


Fig. 4.4 The lower left plot contains raw data (Spearman $\rho = -0.02$); the lower right is a scatterplot of the corresponding transformed values ($\rho = -0.13$). Data courtesy of the SUPPORT study³⁵².

new predictors representing whether any of the factors is positive and a count of the number of positive factors. For the ordinal count of the number of positive factors, score the summary variable to satisfy linearity assumptions as discussed previously. For the more powerful predictor of the two summary measures, test for adequacy of scoring by using all dichotomous variables as candidate predictors after adjusting for the new summary variable. A residual χ^2 statistic can be used to test whether the summary variable adequately captures the predictive information of the series of binary predictors.^P This statistic will have degrees of freedom equal to one less than the number of binary predictors when testing for adequacy of the summary count (and hence will have low power when there are many predictors). Stratification by the summary score and examination of responses over cells can be used to suggest a transformation on the score.

Another approach to scoring a series of related dichotomous predictors is to have “experts” assign severity points to each condition and then to either sum these points or use a hierarchical rule that scores according to the condition with the highest points (see Section 14.3 for an example). The latter has the advantage of being easy to implement for field use. The adequacy of either type of scoring can be checked using tests of linearity in a regression model^Q.

^P Whether this statistic should be used to change the model is problematic in view of model uncertainty.

^Q The R function `score.binary` in the `Hmisc` package (see Section 6.2) assists in computing a summary variable from the series of binary conditions.

4.7.6 *Simplifying Cluster Scores*

If a variable cluster contains many individual predictors, parsimony may sometimes be achieved by predicting the cluster score from a subset of its components (using linear regression or CART (Section 2.5), for example). Then a new cluster score is created and the response model is rerun with the new score in the place of the original one. If one constituent variable has a very high R^2 in predicting the original cluster score, the single variable may sometimes be substituted for the cluster score in refitting the model without loss of predictive discrimination. 22

Sometimes it may be desired to simplify a variable cluster by asking the question “which variables in the cluster are really the predictive ones?,” even though this approach will usually cause true predictive discrimination to suffer. For clusters that are retained after limited step-down modeling, the entire list of variables can be used as candidate predictors and the step-down process repeated. All variables contained in clusters that were not selected initially are ignored. A fair way to validate such two-stage models is to use a resampling method (Section 5.3) with scores for deleted clusters as candidate variables for each resample, along with all the individual variables in the clusters the analyst really wants to retain. A method called *battery reduction* can be used to delete variables from clusters by determining if a subset of the variables can explain most of the variance explained by PC_1 (see [142, Chapter 12] and⁴⁴⁵). This approach does not require examination of associations with Y . Battery reduction can also be used to find a set of individual variables that capture much of the information in the first k principal components. 23

4.7.7 *How Much Data Reduction Is Necessary?*

In addition to using the sample size to degrees of freedom ratio as a rough guide to how much data reduction to do before model fitting, the heuristic shrinkage estimate in Equation 4.3 can also be informative. First, fit a full model with all candidate variables, nonlinear terms, and hypothesized interactions. Let p denote the number of parameters in this model, aside from any intercepts. Let LR denote the log likelihood ratio χ^2 for this full model. The estimated shrinkage is $(LR - p)/LR$. If this falls below 0.9, for example, we may be concerned with the lack of calibration the model may experience on new data. Either a shrunken estimator or data reduction is needed. A reduced model may have acceptable calibration if associations with Y are not used to reduce the predictors.

A simple method, with an assumption, can be used to estimate the target number of total regression degrees of freedom q in the model. In a “best case,” the variables removed to arrive at the reduced model would have no association with Y . The expected value of the χ^2 statistic for testing those

variables would then be $p - q$. The shrinkage for the reduced model is then on average $[\text{LR} - (p - q) - q]/[\text{LR} - (p - q)]$. Setting this ratio to be ≥ 0.9 and solving for q gives $q \leq (\text{LR} - p)/9$. Therefore, reduction of dimensionality down to q degrees of freedom would be expected to achieve $< 10\%$ shrinkage. With these assumptions, there is no hope that a reduced model would have acceptable calibration unless $\text{LR} > p + 9$. If the information explained by the omitted variables is less than one would expect by chance (e.g., their total χ^2 is extremely small), a reduced model could still be beneficial, as long as the conservative bound $(\text{LR} - q)/\text{LR} \geq 0.9$ or $q \leq \text{LR}/10$ were achieved. This conservative bound assumes that no χ^2 is lost by the reduction, that is that the final model $\chi^2 \approx \text{LR}$. This is unlikely in practice. Had the $p - q$ omitted variables had a larger χ^2 of $2(p - q)$ (the break-even point for AIC), q must be $\leq (\text{LR} - 2p)/8$.

As an example, suppose that a binary logistic model is being developed from a sample containing 45 events on 150 subjects. The 10:1 rule suggests we can analyze 4.5 degrees of freedom. The analyst wishes to analyze age, sex, and 10 other variables. It is not known whether interaction between age and sex exists, and whether age is linear. A restricted cubic spline is fitted with four knots, and a linear interaction is allowed between age and sex. These two variables then need $3 + 1 + 1 = 5$ degrees of freedom. The other 10 variables are assumed to be linear and to not interact with themselves or age and sex. There is a total of 15 d.f. The full model with 15 d.f. has $\text{LR} = 50$. Expected shrinkage from this model is $(50 - 15)/50 = 0.7$. Since $\text{LR} > 15 + 9 = 24$, some reduction might yield a better validating model. Reduction to $q = (50 - 15)/9 \approx 4$ d.f. would be necessary, assuming the reduced LR is about $50 - (15 - 4) = 39$. In this case the 10:1 rule yields about the same value for q . The analyst may be forced to assume that age is linear, modeling 3 d.f. for age and sex. The other 10 variables would have to be reduced to a single variable using principal components or another scaling technique. The AIC-based calculation yields a maximum of 2.5 d.f.

If the goal of the analysis is to make a series of hypothesis tests (adjusting P -values for multiple comparisons) instead of to predict future responses, the full model would have to be used.

A summary of the various data reduction methods is given in Figure 4.5.

When principal component analysis or related methods are used for data reduction, the model may be harder to describe since internal coefficients are “hidden.” R code on p. 141 shows how an ordinary linear model fit can be used in conjunction with a logistic model fit based on principal components to draw a nomogram with axes for all predictors.

Fig. 4.5 Summary of Some Data Reduction Methods

Goals	Reasons	Methods
Group predictors so that each group represents a single dimension that can be summarized with a single score	<ul style="list-style-type: none"> • ↓ d.f. arising from multiple predictors • Make PC_1 more reasonable summary 	Variable clustering <ul style="list-style-type: none"> • Subject matter knowledge • Group predictors to maximize proportion of variance explained by PC_1 of each group • Hierarchical clustering using a matrix of similarity measures between predictors
Transform predictors	<ul style="list-style-type: none"> • ↓ d.f. due to nonlinear and dummy variable components • Allows predictors to be optimally combined • Make PC_1 more reasonable summary • Use in customized model for imputing missing values on each predictor 	<ul style="list-style-type: none"> • Maximum total variance on a group of related predictors • Canonical variates on the total set of predictors
Score a group of predictors	↓ d.f. for group to unity	<ul style="list-style-type: none"> • PC_1 • Simple point scores
Multiple dimensional scoring of all predictors	↓ d.f. for all predictors combined	Principal components $1, 2, \dots, k, k < p$ computed from all transformed predictors

4.8 Other Approaches to Predictive Modeling

The approaches recommended in this text are

- fitting fully pre-specified models without deletion of “insignificant” predictors
- using data reduction methods (masked to Y) to reduce the dimensionality of the predictors and then fitting the number of parameters the data’s information content can support

- using shrinkage (penalized estimation) to fit a large model without worrying about the sample size.

Data reduction approaches covered in the last section can yield very interpretable, stable models, but there are many decisions to be made when using a two-stage (reduction/model fitting) approach. Newer single stage approaches are evolving. These new approaches, listed on the text's web site, handle continuous predictors well, unlike recursive partitioning.

When data reduction is not required, generalized additive models^{277,674} should also be considered.

4.9 Overly Influential Observations

Every observation should influence the fit of a regression model. It can be disheartening, however, if a significant treatment effect or the shape of a regression effect rests on one or two observations. Overly influential observations also lead to increased variance of predicted values, especially when variances are estimated by bootstrapping after taking variable selection into account. In some cases, overly influential observations can cause one to abandon a model, “change” the data, or get more data. Observations can be *overly influential* for several major reasons.

1. The most common reason is having too few observations for the complexity of the model being fitted. Remedies for this have been discussed in Sections 4.7 and 4.3.
2. Data transcription or data entry errors can ruin a model fit.
3. Extreme values of the predictor variables can have a great impact, even when these values are validated for accuracy. Sometimes the analyst may deem a subject so atypical of other subjects in the study that deletion of the case is warranted. On other occasions, it is beneficial to truncate measurements where the data density ends. In one dataset of 4000 patients and 2000 deaths, white blood count (WBC) ranged from 500 to 100,000 with .05 and .95 quantiles of 2755 and 26,700, respectively. Predictions from a linear spline function of WBC were sensitive to $WBC > 60,000$, for which there were 16 patients. There were 46 patients with $WBC > 40,000$. Predictions were found to be more stable when WBC was truncated at 40,000, that is, setting WBC to 40,000 if $WBC > 40,000$.
4. Observations containing disagreements between the predictors and the response can influence the fit. Such disagreements should not lead to discarding the observations unless the predictor or response values are erroneous as in Reason 3, or the analysis is made conditional on observations being unlike the influential ones. In one example a single extreme predictor value in a sample of size 8000 that was not on a straight line relationship with

the other (X, Y) pairs caused a χ^2 of 36 for testing nonlinearity of the predictor. Remember that an imperfectly fitting model is a fact of life, and discarding the observations can inflate the model's predictive accuracy. On rare occasions, such lack of fit may lead the analyst to make changes in the model's structure, but ordinarily this is best done from the "ground up" using formal tests of lack of fit (e.g., a test of linearity or interaction).

Influential observations of the second and third kinds can often be detected by careful quality control of the data. Statistical measures can also be helpful. The most common measures that apply to a variety of regression models are *leverage*, DFBETAS, DFFIT, and DFFITS.

Leverage measures the capacity of an observation to be influential due to having extreme predictor values. Such an observation is not *necessarily* influential. To compute leverage in ordinary least squares, we define the *hat matrix* H given by

$$H = X(X'X)^{-1}X'. \quad (4.6)$$

H is the matrix that when multiplied by the response vector gives the predicted values, so it measures how an observation estimates its own predicted response. The diagonals h_{ii} of H are the leverage measures and they are not influenced by Y . It has been suggested⁴⁷ that $h_{ii} > 2(p+1)/n$ signal a high leverage point, where p is the number of columns in the design matrix X aside from the intercept and n is the number of observations. Some believe that the distribution of h_{ii} should be examined for values that are higher than typical.

DFBETAS is the change in the vector of regression coefficient estimates upon deletion of each observation in turn, scaled by their standard errors.⁴⁷ Since DFBETAS encompasses an effect for each predictor's coefficient, DFBETAS allows the analyst to isolate the problem better than some of the other measures. DFFIT is the change in the predicted $X\beta$ when the observation is dropped, and DFFITS is DFFIT standardized by the standard error of the estimate of $X\beta$. In both cases, the standard error used for normalization is recomputed each time an observation is omitted. Some classify an observation as overly influential when $|DFFITS| > 2\sqrt{(p+1)/(n-p-1)}$, while others prefer to examine the entire distribution of DFFITS to identify "outliers".⁴⁷

Section 10.7 discusses influence measures for the logistic model, which requires maximum likelihood estimation. These measures require the use of special residuals and information matrices (in place of $X'X$).

If truly influential observations are identified using these indexes, careful thought is needed to decide how (or whether) to deal with them. Most important, there is no substitute for careful examination of the dataset before doing any analyses.⁹⁹ Spence and Garrison [581, p. 16] feel that

Although the identification of aberrations receives considerable attention in most modern statistical courses, the emphasis sometimes seems to be on disposing of embarrassing data by searching for sources of technical error or

minimizing the influence of inconvenient data by the application of resistant methods. Working scientists often find the most interesting aspect of the analysis inheres in the lack of fit rather than the fit itself.

4.10 Comparing Two Models

Frequently one wants to choose between two competing models on the basis of a common set of observations. The methods that follow assume that the performance of the models is evaluated on a sample not used to develop either one. In this case, predicted values from the model can usually be considered as a single new variable for comparison with responses in the new dataset. These methods listed below will also work if the models are compared using the same set of data used to fit each one, as long as both models have the same effective number of (candidate or actual) parameters. This requirement prevents us from rewarding a model just because it overfits the training sample (see Section 9.8.1 for a method comparing two models of differing complexity). The methods can also be enhanced using bootstrapping or cross-validation on a single sample to get a fair comparison when the playing field is not level, for example, when one model had more opportunity for fitting or overfitting the responses.

Some of the criteria for choosing one model over the other are

1. calibration (e.g., one model is well-calibrated and the other is not),
2. discrimination,
3. face validity,
4. measurement errors in required predictors,
5. use of continuous predictors (which are usually better defined than categorical ones),
6. omission of “insignificant” variables that nonetheless make sense as risk factors,
7. simplicity (although this is less important with the availability of computers), and
8. lack of fit for specific types of subjects.

Items 3 through 7 require subjective judgment, so we focus on the other aspects. If the purpose of the models is only to rank-order subjects, calibration is not an issue. Otherwise, a model having poor calibration can be dismissed outright. Given that the two models have similar calibration, discrimination should be examined critically. Various statistical indexes can quantify discrimination ability (e.g., R^2 , model χ^2 , Somers' D_{xy} , Spearman's ρ , area under ROC curve—see Section 10.8). Rank measures (D_{xy} , ρ , ROC area) only measure how well predicted values can rank-order responses. For example, predicted probabilities of 0.01 and 0.99 for a pair of subjects are no better than probabilities of 0.2 and 0.8 using rank measures, if the first subject had

a lower response value than the second. Therefore, rank measures such as ROC area (c index), although fine for describing a given model, may not be very sensitive in choosing between two models^{118,488,493}. This is especially true when the models are strong, as it is easier to move a rank correlation from 0.6 to 0.7 than it is to move it from 0.9 to 1.0. Measures such as R^2 and the model χ^2 statistic (calculated from the predicted and observed responses) are more sensitive. Still, one may not know how to interpret the added utility of a model that boosts the R^2 from 0.80 to 0.81.

Again given that both models are equally well calibrated, discrimination can be studied more simply by examining the distribution of predicted values \hat{Y} . Suppose that the predicted value is the probability that a subject dies. Then high-resolution histograms of the predicted risk distributions for the two models can be very revealing. If one model assigns 0.02 of the sample to a risk of dying above 0.9 while the other model assigns 0.08 of the sample to the high risk group, the second model is more discriminating. The worth of a model can be judged by how far it goes out on a limb while still maintaining good calibration.

25

Frequently, one model will have a similar discrimination index to another model, but the likelihood ratio χ^2 statistic is meaningfully greater for one. Assuming corrections have been made for complexity, the model with the higher χ^2 usually has a better fit for *some* subjects, although not necessarily for the *average* subject. A crude plot of predictions from the first model against predictions from the second, possibly stratified by Y , can help describe the differences in the models. More specific analyses will determine the characteristics of subjects where the differences are greatest. Large differences may be caused by an omitted, underweighted, or improperly transformed predictor, among other reasons. In one example, two models for predicting hospital mortality in critically ill patients had the same discrimination index (to two decimal places). For the relatively small subset of patients with extremely low white blood counts or serum albumin, the model that treated these factors as continuous variables provided predictions that were very much different from a model that did not.

When comparing predictions for two models that may not be calibrated (from overfitting, e.g.), the two sets of predictions may be shrunk so as to not give credit for overfitting (see Equation 4.3).

Sometimes one wishes to compare two models that used the response variable differently, a much more difficult problem. For example, an investigator may want to choose between a survival model that used time as a continuous variable, and a binary logistic model for dead/alive at six months. Here, other considerations are also important (see Section 17.1). A model that predicts dead/alive at six months does not use the response variable effectively, and it provides no information on the chance of dying within three months.

When one or both of the models is fitted using least squares, it is useful to compare them using an error measure that was not used as the optimization criterion, such as mean absolute error or median absolute error. Mean

and median absolute errors are excellent measures for judging the value of a model developed without transforming the response to a model fitted after transforming Y , then back-transforming to get predictions.

26

4.11 Improving the Practice of Multivariable Prediction

Standards for published predictive modeling and feature selection in high-dimensional problems are not very high. There are several things that a good analyst can do to improve the situation.

1. Insist on validation of predictive models and discoveries, using rigorous internal validation based on resampling or using external validation.
2. Show collaborators that split-sample validation is not appropriate unless the number of subjects is huge
 - This can be demonstrated by splitting the data more than once and seeing volatile results, and by calculating a confidence interval for the predictive accuracy in the test dataset and showing that it is very wide.
3. Run a simulation study with no real associations and show that associations are easy to find if a dangerous data mining procedure is used. Alternately, analyze the collaborator's data after randomly permuting the Y vector and show some "positive" findings.
4. Show that alternative explanations are easy to posit. For example:
 - The importance of a risk factor may disappear if 5 "unimportant" risk factors are added back to the model
 - Omitted main effects can explain away apparent interactions.
 - Perform a *uniqueness analysis*: attempt to predict the predicted values from a model derived by data torture from all of the features not used in the model. If one can obtain $R^2 = 0.85$ in predicting the "winning" feature signature (predicted values) from the "losing" features, the "winning" pattern is not unique and may be unreliable.

4.12 Summary: Possible Modeling Strategies

Some possible global modeling strategies are to

- Use a method known not to work well (e.g., stepwise variable selection without penalization; recursive partitioning resulting in a single tree), document how poorly the model performs (e.g. using the bootstrap), and use the model anyway
- Develop a black box model that performs poorly and is difficult to interpret (e.g., does not incorporate penalization)

- Develop a black box model that performs well and is difficult to interpret
- Develop interpretable approximations to the black box
- Develop an interpretable model (e.g. give priority to additive effects) that performs well and is likely to perform equally well on future data from the same stream.

As stated in the Preface, the strategy emphasized in this text, stemming from the last philosophy, is to decide how many degrees of freedom can be “spent,” where they should be spent, and then to spend them. If statistical tests or confidence limits are required, later reconsideration of how d.f. are spent is not usually recommended. In what follows some default strategies are elaborated. These strategies are far from failsafe, but they should allow the reader to develop a strategy that is tailored to a particular problem. At the least these default strategies are concrete enough to be criticized so that statisticians can devise better ones.

4.12.1 Developing Predictive Models

The following strategy is generic although it is aimed principally at the development of accurate predictive models.

1. Assemble as much accurate pertinent data as possible, with wide distributions for predictor values. For survival time data, follow-up must be sufficient to capture enough events as well as the clinically meaningful phases if dealing with a chronic process.
2. Formulate good hypotheses that lead to specification of relevant candidate predictors and possible interactions. Don't use Y (either informally using graphs, descriptive statistics, or tables, or formally using hypothesis tests or estimates of effects such as odds ratios) in devising the list of candidate predictors.
3. If there are missing Y values on a small fraction of the subjects but Y can be reliably substituted by a surrogate response, use the surrogate to replace the missing values. Characterize tendencies for Y to be missing using, for example, recursive partitioning or binary logistic regression. Depending on the model used, even the information on X for observations with missing Y can be used to improve precision of $\hat{\beta}$, so multiple imputation of Y can sometimes be effective. Otherwise, discard observations having missing Y .
4. Impute missing X s if the fraction of observations with any missing X s is not tiny. Characterize observations that had to be discarded. Special imputation models may be needed if a continuous X needs a non-monotonic transformation (p. 52). These models can simultaneously impute missing values while determining transformations. In most cases, multiply impute missing X s based on other X s and Y , and other available information about the missing data mechanism.

5. For each predictor specify the complexity or degree of nonlinearity that should be allowed (see Section 4.1). When prior knowledge does not indicate that a predictor has a linear effect on the property $C(Y|X)$ (the property of the response that *can* be linearly related to X), specify the number of degrees of freedom that should be devoted to the predictor. The d.f. (or number of knots) can be larger when the predictor is thought to be more important in predicting Y or when the sample size is large.
6. If the number of terms fitted or tested in the modeling process (counting nonlinear and cross-product terms) is too large in comparison with the number of outcomes in the sample, use data reduction (ignoring Y) until the number of remaining free variables needing regression coefficients is tolerable. Use the $m/10$ or $m/15$ rule or an estimate of likely shrinkage or overfitting (Section 4.7) as a guide. Transformations determined from the previous step may be used to reduce each predictor into 1 d.f., or the transformed variables may be clustered into highly correlated groups if more data reduction is required. Alternatively, use penalized estimation with the entire set of variables. This will also effectively reduce the total degrees of freedom.²⁷²
7. Use the entire sample in the model development as data are too precious to waste. If steps listed below are too difficult to repeat for each bootstrap or cross-validation sample, hold out test data from **all** model development steps that follow.
8. When you can test for model complexity in a very structured way, you may be able to simplify the model without a great need to penalize the final model for having made this initial look. For example, it can be advisable to test an entire group of variables (e.g., those more expensive to collect) and to either delete or retain the entire group for further modeling, based on a single P -value (especially if the P value is not between 0.05 and 0.2). Another example of structured testing to simplify the “initial” model is making *all* continuous predictors have the same number of knots k , varying k from 0 (linear), 3, 4, 5, . . . , and choosing the value of k that optimizes AIC. A composite test of all nonlinear effects in a model can also be used, and statistical inferences are not invalidated if the global test of nonlinearity yields $P > 0.2$ or so and the analyst deletes all nonlinear terms.
9. Make tests of linearity of effects in the model only to demonstrate to others that such effects are often statistically significant. Don’t remove insignificant effects from the model when tested separately by predictor. Any examination of the response that might result in simplifying the model needs to be accounted for in computing confidence limits and other statistics. It is preferable to retain the complexity that was prespecified in Step 5 regardless of the results of assessments of nonlinearity.

10. Check additivity assumptions by testing prespecified interaction terms. If the global test for additivity is significant or equivocal, all prespecified interactions should be retained in the model. If the test is decisive (e.g., $P > 0.3$), all interaction terms can be omitted, and in all likelihood there is no need to repeat this pooled test for each resample during model validation. In other words, one can assume that had the global interaction test been carried out for each bootstrap resample it would have been insignificant at the 0.05 level more than, say, 0.9 of the time. In this large P -value case the pooled interaction test did not induce an uncertainty in model selection that needed accounting.
11. Check to see if there are overly influential observations.
12. Check distributional assumptions and choose a different model if needed.
13. Do limited backwards step-down variable selection if parsimony is more important than accuracy.⁵⁸² The cost of doing any aggressive variable selection is that the variable selection algorithm must also be included in a resampling procedure to properly validate the model or to compute confidence limits and the like.
14. This is the “final” model.
15. Interpret the model graphically (Section 5.1) and by examining predicted values and using appropriate significance tests without trying to interpret some of the individual model parameters. For collinear predictors obtain pooled tests of association so that competition among variables will not give misleading impressions of their total significance.
16. Validate the final model for calibration and discrimination ability, preferably using bootstrapping (see Section 5.3). Steps 9 to 13 must be repeated for each bootstrap sample, at least approximately. For example, if age was transformed when building the final model, and the transformation was suggested by the data using a fit involving age and age², each bootstrap repetition should include both age variables with a possible step-down from the quadratic to the linear model based on automatic significance testing at each step.
17. Shrink parameter estimates if there is overfitting but no further data reduction is desired, if shrinkage was not built into the estimation process.
18. When missing values were imputed, adjust final variance–covariance matrix for imputation wherever possible (e.g., using bootstrap or multiple imputation). This may affect some of the other results.
19. When all steps of the modeling strategy can be automated, consider using Faraway’s method¹⁸⁶ to penalize for the randomness inherent in the multiple steps.
20. Develop simplifications to the full model by approximating it to any desired degrees of accuracy (Section 5.5).

4.12.2 Developing Models for Effect Estimation

By effect estimation is meant point and interval estimation of differences in properties of the responses between two or more settings of some predictors, or estimating some function of these differences such as the antilog. In ordinary multiple regression with no transformation of Y such differences are absolute estimates. In regression involving $\log(Y)$ or in logistic or proportional hazards models, effect estimation is, at least initially, concerned with estimation of relative effects. As discussed on pp. 4 and 224, estimation of absolute effects for these models must involve accurate prediction of overall response values, so the strategy in the previous section applies.

When estimating differences or relative effects, the bias in the effect estimate, besides being influenced by the study design, is related to how well subject heterogeneity and confounding are taken into account. The variance of the effect estimate is related to the distribution of the variable whose levels are being compared, and, in least squares estimates, to the amount of variation “explained” by the entire set of predictors. Variance of the estimated difference can increase if there is overfitting. So for estimation, the previous strategy largely applies.

The following are differences in the modeling strategy when effect estimation is the goal.

1. There is even less gain from having a parsimonious model than when developing overall predictive models, as estimation is usually done at the time of analysis. Leaving insignificant predictors in the model increases the likelihood that the confidence interval for the effect of interest has the stated coverage. By contrast, overall predictions are conditional on the values of all predictors in the model. The variance of such predictions is increased by the presence of unimportant variables, as predictions are still conditional on the particular values of these variables (Section 5.5.1) and cancellation of terms (which occurs when differences are of interest) does not occur.
2. Careful consideration of inclusion of interactions is still a major consideration for estimation. If a predictor whose effects are of major interest is allowed to interact with one or more other predictors, effect estimates must be conditional on the values of the other predictors and hence have higher variance.
3. A major goal of imputation is to avoid lowering the sample size because of missing values in adjustment variables. If the predictor of interest is the only variable having a substantial number of missing values, multiple imputation is less worthwhile, unless it corrects for a substantial bias caused by deletion of nonrandomly missing data.

4. The analyst need not be very concerned about conserving degrees of freedom devoted to the predictor of interest. The complexity allowed for this variable is usually determined by prior beliefs, with compromises that consider the bias-variance trade-off.
5. If penalized estimation is used, the analyst may wish to not shrink parameter estimates for the predictor of interest.
6. Model validation is not necessary unless the analyst wishes to use it to quantify the degree of overfitting.

4.12.3 Developing Models for Hypothesis Testing

A default strategy for developing a multivariable model that is to be used as a basis for hypothesis testing is almost the same as the strategy used for estimation.

1. There is little concern for parsimony. A full model fit, including insignificant variables, will result in more accurate P -values for tests for the variables of interest.
2. Careful consideration of inclusion of interactions is still a major consideration for hypothesis testing. If one or more predictors interacts with a variable of interest, either separate hypothesis tests are carried out over the levels of the interacting factors, or a combined “main effect + interaction” test is performed. For example, a very well-defined test is whether treatment is effective for *any* race group.
3. If the predictor of interest is the only variable having a substantial number of missing values, multiple imputation is less worthwhile. In some cases, multiple imputation may increase power (e.g., in ordinary multiple regression one can obtain larger degrees of freedom for error) but in others there will be little net gain. However, the test can be biased due to exclusion of nonrandomly missing observations if imputation is not done.
4. As before, the analyst need not be very concerned about conserving degrees of freedom devoted to the predictor of interest. The degrees of freedom allowed for this variable is usually determined by prior beliefs, with careful consideration of the trade-off between bias and power.
5. If penalized estimation is used, the analyst should not shrink parameter estimates for the predictors being tested.
6. Model validation is not necessary unless the analyst wishes to use it to quantify the degree of overfitting. This may shed light on whether there is overadjustment for confounders.

4.13 Further Reading

- [1] Some good general references that address modeling strategies are [216, 269, 476, 590].
- [2] Even though they used a generalized correlation index for *screening* variables and not for transforming them, Hall and Miller²⁴⁹ present a related idea, computing the ordinary R^2 against a cubic spline transformation of each potential predictor.
- [3] Simulation studies are needed to determine the effects of modifying the model based on assessments of “predictor promise.” Although it is unlikely that this strategy will result in regression coefficients that are biased high in absolute value, it may on some occasions result in somewhat optimistic standard errors and a slight elevation in type I error probability. Some simulation results may be found on the Web site. Initial promising findings for least squares models for two uncorrelated predictors indicate that the procedure is conservative in its estimation of σ^2 and in preserving type I error.
- [4] Verweij and Houwelingen⁶⁴⁰ and Shao⁵⁶⁵ describe how cross-validation can be used in formulating a stopping rule. Luo et al.⁴³⁰ developed an approach to tuning forward selection by adding noise to Y .
- [5] Roecker⁵²⁸ compared forward variable selection (FS) and all possible subsets selection (APS) with full model fits in ordinary least squares. APS had a greater tendency to select smaller, less accurate models than FS. Neither selection technique was as accurate as the full model fit unless more than half of the candidate variables was redundant or unnecessary.
- [6] Wiegand⁶⁶⁸ showed that it is not very fruitful to try different stepwise algorithms and then to be comforted by agreements in some of the variables selected. It is easy for different stepwise methods to agree on the wrong set of variables.
- [7] Other results on how variable selection affects inference may be found in Hurvich and Tsai³¹⁶ and Breiman [66, Section 8.1].
- [8] Goring et al.²²⁷ presented an interesting analysis of the huge bias caused by conditioning analyses on statistical significance in a high-dimensional genetics context.
- [9] Steyerberg et al.⁵⁸⁹ have comparisons of smoothly penalized estimators with the lasso and with several stepwise variable selection algorithms.
- [10] See Weiss,⁶⁵⁶ Faraway,¹⁸⁶ and Chatfield¹⁰⁰ for more discussions of the effect of not prespecifying models, for example, dependence of point estimates of effects on the variables used for adjustment.
- [11] Greenland²⁴¹ provides an example in which overfitting a logistic model resulted in far too many predictors with $P < 0.05$.
- [12] See Peduzzi et al.^{486, 487} for studies of the relationship between “events per variable” and types I and II error, accuracy of variance estimates, and accuracy of normal approximations for regression coefficient estimators. Their findings are consistent with those given in the text (but⁶⁴⁴ has a slightly different take). van der Ploeg et al.⁶²⁹ did extensive simulations to determine the events per variable ratio needed to avoid a drop-off (in an independent test sample) in more than 0.01 in the c -index, for a variety of predictive methods. They concluded that support vector machines, neural networks, and random forests needed far more events per variable to achieve freedom from overfitting than does logistic regression, and that recursive partitioning was not competitive. Logistic regression required between 20 and 50 events per variable to avoid overfitting. Different results might have been obtained had the authors used a proper accuracy score.
- [13] Copas [122, Eq. 8.5] adds 2 to the numerator of Equation 4.3 (see also [504, 631]).

- [14] An excellent discussion about such indexes may be found in <http://r.789695.n4.nabble.com/Adjusted-R-squared-formula-in-lm-td4656857.html> where J. Lucke points out that R^2 tends to $\frac{p}{n-1}$ when the population R^2 is zero, but R_{adj}^2 converges to zero.
- [15] Efron [173, Eq. 4.23] and van Houwelingen and le Cessie⁶³³ showed that the average expected optimism in a mean logarithmic quality score for a p -predictor binary logistic model is p/n . Taylor et al.⁶⁰⁰ showed that the ratio of variances for certain quantities is proportional to the ratio of the number of parameters in two models. Copas stated that “Shrinkage can be particularly marked when stepwise fitting is used: the shrinkage is then closer to that expected of the full regression rather than of the subset regression actually fitted.”^{122, 504, 631} Spiegelhalter,⁵⁸² in arguing against variable selection, states that better prediction will often be obtained by fitting all candidate variables in the final model, shrinking the vector of regression coefficient estimates towards zero.
- [16] See Belsley [46, pp. 28–30] for some reservations about using VIF.
- [17] Friedman and Wall²⁰⁸ discuss and provide graphical devices for explaining *suppression* by a predictor not correlated with the response but that is correlated with another predictor. Adjusting for a suppressor variable will increase the predictive discrimination of the model. Meinshausen⁴⁵³ developed a novel hierarchical approach to gauging the importance of collinear predictors.
- [18] For incomplete principal component regression see [101, 119, 120, 142, 144, 320, 325]. See^{396, 686} for sparse principal component analysis methods in which constraints are applied to loadings so that some of them are set to zero. The latter reference provides a principal component method for binary data. See²⁴⁶ for a type of sparse principal component analysis that also encourages loadings to be similar for a group of highly correlated variables and allows for a type of variable clustering. See [390] for principal surfaces. Sliced inverse regression is described in [104, 119, 120, 189, 403, 404]. For material on variable clustering see [142, 144, 268, 441, 539]. A good general reference on cluster analysis is [634, Chapter 11]. de Leeuw and Mair in their R `homals` package [153] have one of the most general approaches to data reduction related to optimal scaling. Their approach includes nonlinear principal component analysis among several other multivariate analyses.
- [19] The redundancy analysis described here is related to *principal variables*⁴⁴⁸ but is faster.
- [20] Meinshausen⁴⁵³ developed a method of testing the importance of competing (collinear) variables using an interesting automatic clustering procedure.
- [21] The R `ClustOfVar` package by Marie Chavent, Vanessa Kuentz, Benoit Liquet, and Jerome Saracco generalizes variable clustering and explicitly handles a mixture of quantitative and categorical predictors. It also implements bootstrap cluster stability analysis.
- [22] Principal components are commonly used to summarize a cluster of variables. Vines⁶⁴³ developed a method to constrain the principal component coefficients to be integers without much loss of explained variability.
- [23] Jolliffe³²⁴ presented a way to discard some of the variables making up principal components. Wang and Gehan⁶⁴⁹ presented a new method for finding subsets of predictors that approximate a set of principal components, and surveyed other methods for simplifying principal components.
- [24] See D’Agostino et al.¹⁴⁴ for excellent examples of variable clustering (including a two-stage approach) and other data reduction techniques using both statistical methods and subject-matter expertise.
- [25] Cook¹¹⁸ and Pencina et al.^{490, 492, 493} present an approach for judging the added value of new variables that is based on evaluating the extent to which the new information moves predicted probabilities higher for subjects having events and lower for subjects not having events. But see^{292, 592}.

- [26] The `Hmisc abs.error.pred` function computes a variety of accuracy measures based on absolute errors.
- [27] Shen et al.⁵⁶⁷ developed an “optimal approximation” method to make correct inferences after model selection.

4.14 Problems

Analyze the SUPPORT dataset (`getHdata(support)`) as directed below to relate selected variables to total cost of the hospitalization. Make sure this response variable is utilized in a way that approximately satisfies the assumptions of normality-based multiple regression so that statistical inferences will be accurate. See problems at the end of Chapters 3 and 7 of the text for more information. Consider as predictors mean arterial blood pressure, heart rate, age, disease group, and coma score.

1. Do an analysis to understand interrelationships among predictors, and find optimal scaling (transformations) that make the predictors better relate to each other (e.g., optimize the variation explained by the first principal component).
2. Do a redundancy analysis of the predictors, using both a less stringent and a more stringent approach to assessing the redundancy of the multiple-level variable disease group.
3. Do an analysis that helps one determine how many d.f. to devote to each predictor.
4. Fit a model, assuming the above predictors act additively, but do not assume linearity for the age and blood pressure effects. Use the truncated power basis for fitting restricted cubic spline functions with 5 knots. Estimate the shrinkage coefficient $\hat{\gamma}$.
5. Make appropriate graphical diagnostics for this model.
6. Test linearity in age, linearity in blood pressure, and linearity in heart rate, and also do a joint test of linearity simultaneously in all three predictors.
7. Expand the model to not assume additivity of age and blood pressure. Use a tensor natural spline or an appropriate restricted tensor spline. If you run into any numerical difficulties, use 4 knots instead of 5. Plot in an interpretable fashion the estimated 3-D relationship between age, blood pressure, and cost for a fixed disease group.
8. Test for additivity of age and blood pressure. Make a joint test for the overall absence of complexity in the model (linearity and additivity simultaneously).

Chapter 5

Describing, Resampling, Validating, and Simplifying the Model

5.1 Describing the Fitted Model

5.1.1 *Interpreting Effects*

Before addressing issues related to describing and interpreting the model and its coefficients, one can never apply too much caution in attempting to interpret results in a causal manner. Regression models are excellent tools for estimating and inferring *associations* between an X and Y given that the “right” variables are in the model. Any ability of a model to provide *causal* inference rests entirely on the faith of the analyst in the experimental design, completeness of the set of variables that are thought to measure confounding and are used for adjustment when the experiment is not randomized, lack of important measurement error, and lastly the goodness of fit of the model.

The first line of attack in interpreting the results of a multivariable analysis is to interpret the model’s parameter estimates. For simple linear, additive models, regression coefficients may be readily interpreted. If there are interactions or nonlinear terms in the model, however, simple interpretations are usually impossible. Many programs ignore this problem, routinely printing such meaningless quantities as the effect of increasing age² by one day while holding age constant. A meaningful age change needs to be chosen, and connections between mathematically related variables must be taken into account. These problems can be solved by relying on predicted values and differences between predicted values.

Even when the model contains no nonlinear effects, it is difficult to compare regression coefficients across predictors having varying scales. Some analysts like to gauge the relative contributions of different predictors on a common scale by multiplying regression coefficients by the standard deviations of the predictors that pertain to them. This does not make sense for nonnormally distributed predictors (and regression models should not need

to make assumptions about the distributions of predictors). When a predictor is binary (e.g., sex), the standard deviation makes no sense as a scaling factor as the scale would depend on the prevalence of the predictor.^a

1

It is more sensible to estimate the change in Y when X_j is changed by an amount that is subject-matter relevant. For binary predictors this is a change from 0 to 1. For many continuous predictors the interquartile range is a reasonable default choice. If the 0.25 and 0.75 quantiles of X_j are g and h , linearity holds, and the estimated coefficient of X_j is b ; $b \times (h - g)$ is the effect of increasing X_j by $h - g$ units, which is a span that contains half of the sample values of X_j .

For the more general case of continuous predictors that are monotonically but not linearly related to Y , a useful point summary is the change in $X\beta$ when the variable changes from its 0.25 quantile to its 0.75 quantile. For models for which $\exp(X\beta)$ is meaningful, antilogging the predicted change in $X\beta$ results in quantities such as interquartile-range odds and hazards ratios. When the variable is involved in interactions, these ratios are estimated separately for various levels of the interacting factors. For categorical predictors, ordinary effects are computed by comparing each level of the predictor with a reference level. See Section 10.10 and Chapter 11 for tabular and graphical examples of this approach.

2

The model can be described using *partial effect plots* by plotting each X against $X\hat{\beta}$ holding other predictors constant. Modified versions of such plots, by nonlinearly rank-transforming the predictor axis, can show the relative importance of a predictor³³⁶.

For an X that interacts with other factors, separate curves are drawn on the same graph, one for each level of the interacting factor.

3

Nomograms^{40, 254, 339, 427} provide excellent graphical depictions of all the variables in the model, in addition to enabling the user to obtain predicted values manually. Nomograms are especially good at helping the user envision interactions. See Section 10.10 and Chapter 11 for examples.

4

5.1.2 Indexes of Model Performance

5.1.2.1 Error Measures

Care must be taken in the choice of accuracy scores to be used in validation. Indexes can be broken down into three main areas.

Central tendency of prediction errors: These measures include mean absolute differences, mean squared differences, and logarithmic scores. An absolute measure is mean $|Y - \hat{Y}|$. The mean squared error is a commonly used and sensitive measure if there are no outliers. For the special case

^a The s.d. of a binary variable is, aside from a multiplier of $\frac{n}{n-1}$, equal to $\sqrt{a(1-a)}$, where a is the proportion of ones.

where Y is binary, such a measure is the Brier score, which is a quadratic proper scoring rule that combines calibration and discrimination^b. The logarithmic proper scoring rules (related to average log-likelihood) is even more sensitive but can be harder to interpret and can be destroyed by a single predicted probability of 0 or 1 that was incorrect.

Discrimination measures: A measure of pure discrimination is a rank correlation of \hat{Y} and Y , including Spearman's ρ , Kendall's τ , and Somers' D_{xy} . When Y is binary, $D_{xy} = 2 \times (c - \frac{1}{2})$ where c is the concordance probability or area under the receiver operating characteristic curve, a linear translation of the Wilcoxon-Mann-Whitney statistic. R^2 is *mostly* a measure of discrimination, and R_{adj}^2 is a good overfitting-corrected measure, if the model is pre-specified. See Section 10.8 for more information about rank-based measures.

Discrimination measures based on variation in \hat{Y} : These include the regression sum of squares and the g -Index (see below).

Calibration measures: These assess absolute prediction accuracy. *Calibration-in-the-large* compares the average \hat{Y} with the average Y . A *high-resolution calibration curve* or *calibration-in-the-small* assesses the absolute forecast accuracy of predictions at individual levels of \hat{Y} . When the calibration curve is linear, this can be summarized by the calibration slope and intercept. A more general approach uses the *loess* nonparametric smoother to estimate the calibration curve³⁷. For any shape of calibration curve, errors can be summarized by quantities such as the maximum absolute calibration error, mean absolute calibration error, and 0.9 quantile of calibration error.

The g -index is a new measure of a model's predictive discrimination based only on $X\hat{\beta} = \hat{Y}$ that applies quite generally. It is based on Gini's mean difference for a variable Z , which is the mean over all possible $i \neq j$ of $|Z_i - Z_j|$. The g -index is an interpretable, robust, and highly efficient measure of variation. For example, when predicting systolic blood pressure, $g = 11\text{mmHg}$ represents a typical difference in \hat{Y} . g is independent of censoring and other complexities. For models in which the anti-log of a difference in \hat{Y} represents meaningful ratios (e.g., odds ratios, hazard ratios, ratio of medians), g_r can be defined as $\exp(g)$. For models in which \hat{Y} can be turned into a probability estimate (e.g., logistic regression), g_p is defined as Gini's mean difference of \hat{P} . These g -indexes represent e.g. "typical" odds ratios, and "typical" risk differences. Partial g indexes can also be defined. More details may be found in the documentation for the R `rms` package's `gIndex` function.

5

^b There are decompositions of the Brier score into discrimination and calibration components.

5.2 The Bootstrap

When one assumes that a random variable Y has a certain population distribution, one can use simulation or analytic derivations to study how a statistical estimator computed from samples from this distribution behaves. For example, when Y has a log-normal distribution, the variance of the sample median for a sample of size n from that distribution can be derived analytically. Alternatively, one can simulate 500 samples of size n from the log-normal distribution, compute the sample median for each sample, and then compute the sample variance of the 500 sample medians. Either case requires knowledge of the population distribution function.

Efron's *bootstrap*^{150, 177, 178} is a general-purpose technique for obtaining estimates of the properties of statistical estimators without making assumptions about the distribution giving rise to the data. Suppose that a random variable Y comes from a cumulative distribution function $F(y) = \text{Prob}\{Y \leq y\}$ and that we have a sample of size n from this unknown distribution, Y_1, Y_2, \dots, Y_n . The basic idea is to repeatedly simulate a sample of size n from F , computing the statistic of interest, and assessing how the statistic behaves over B repetitions. Not having F at our disposal, we can estimate F by the empirical cumulative distribution function

$$F_n(y) = \frac{1}{n} \sum_{i=1}^n [Y_i \leq y]. \quad (5.1)$$

F_n corresponds to a density function that places probability $1/n$ at each observed datapoint (k/n if that point were duplicated k times and its value listed only once).

As an example, consider a random sample of size $n = 30$ from a normal distribution with mean 100 and standard deviation 10. Figure 5.1 shows the population and empirical cumulative distribution functions.

Now pretend that $F_n(y)$ is the original population distribution $F(y)$. Sampling from F_n is equivalent to sampling with replacement from the observed data Y_1, \dots, Y_n . For large n , the expected fraction of original datapoints that are selected for each bootstrap sample is $1 - e^{-1} = 0.632$. Some points are selected twice, some three times, a few four times, and so on. We take B samples of size n with replacement, with B chosen so that the summary measure of the individual statistics is nearly as good as taking $B = \infty$. The bootstrap is based on the fact that the distribution of the *observed* differences between a resampled estimate of a parameter of interest and the original estimate of the parameter from the whole sample tells us about the distribution of *unobservable* differences between the original estimate and the unknown population value of the parameter.

As an example, consider the data (1, 5, 6, 7, 8, 9) and suppose that we would like to obtain a 0.80 confidence interval for the population median, as well as an estimate of the population expected value of the sample median (the latter

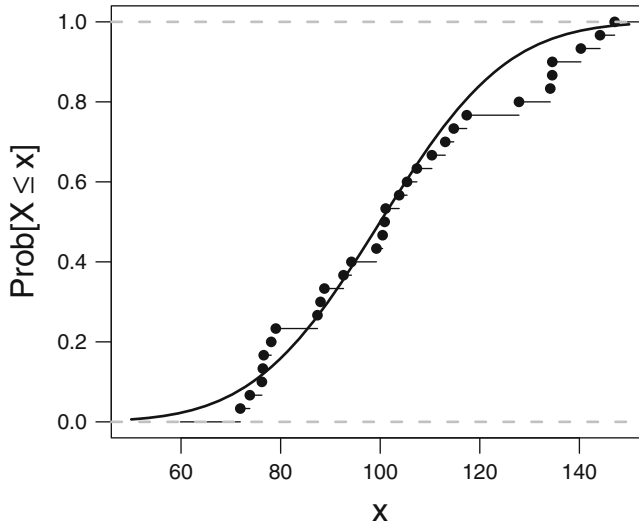


Fig. 5.1 Empirical and population cumulative distribution function

is only used to estimate bias in the sample median). The first 20 bootstrap samples (after sorting data values) and the corresponding sample medians are shown in Table 5.1.

For a given number B of bootstrap samples, our estimates are simply the sample 0.1 and 0.9 quantiles of the sample medians, and the mean of the sample medians. Not knowing how large B should be, we could let B range from, say, 50 to 1000, stopping when we are sure the estimates have converged. In the left plot of Figure 5.2, B varies from 1 to 400 for the mean (10 to 400 for the quantiles). It can be seen that the bootstrap estimate of the population mean of the sample median can be estimated satisfactorily when $B > 50$. For the lower and upper limits of the 0.8 confidence interval for the population median Y , B must be at least 200. For more extreme confidence limits, B must be higher still.

For the final set of 400 sample medians, a histogram (right plot in Figure 5.2) can be used to assess the form of the sampling distribution of the sample median. Here, the distribution is almost normal, although there is a slightly heavy left tail that comes from the data themselves having a heavy left tail. For large samples, sample medians are normally distributed for a wide variety of population distributions. Therefore we could use bootstrapping to estimate the variance of the sample median and then take ± 1.28 standard errors as a 0.80 confidence interval. In other cases (e.g., regression coefficient estimates for certain models), estimates are asymmetrically distributed, and the bootstrap quantiles are better estimates than confidence intervals that are based on a normality assumption. Note that because sample quantiles are more or less restricted to equal one of the values in the sample, the boot-

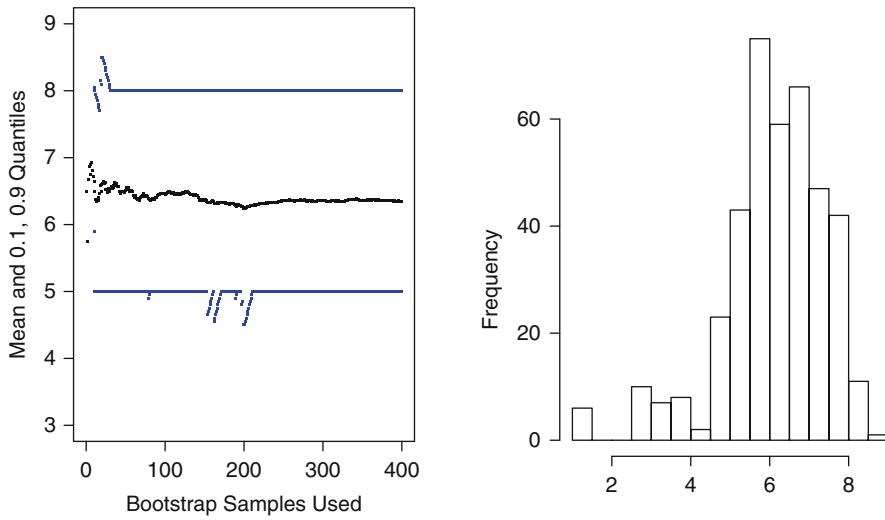


Fig. 5.2 Estimating properties of sample median using the bootstrap

Table 5.1 First 20 bootstrap samples

Bootstrap Sample	Sample Median
1 6 6 7 8 9	6.5
1 5 5 5 6 8	5.0
5 7 8 9 9 9	8.5
7 7 7 8 8 9	7.5
1 5 7 7 9 9	7.0
1 5 6 6 7 8	6.0
7 8 8 8 8 8	8.0
5 5 5 7 9 9	6.0
1 5 5 7 7 9	6.0
1 5 5 7 7 8	6.0
1 1 5 5 7 7	5.0
1 1 5 5 7 8	5.0
1 5 5 7 7 8	6.0
1 5 6 7 8 8	6.5
1 5 6 7 9 9	6.5
6 6 7 7 8 9	7.0
1 5 7 8 8 9	7.5
6 6 8 9 9 9	8.5
1 1 5 5 6 9	5.0
1 6 8 9 9 9	8.5

strap distribution is discrete and can be dependent on a small number of outliers. For this reason, bootstrapping quantiles does not work particularly well for small samples [150, pp. 41–43].

The method just presented for obtaining a nonparametric confidence interval for the population median is called the *bootstrap percentile method*. It is the simplest but not necessarily the best performing bootstrap method. 7

In this text we use the bootstrap primarily for computing statistical estimates that are much different from standard errors and confidence intervals, namely, estimates of model performance.

5.3 Model Validation

5.3.1 Introduction

The surest method to have a model fit the data at hand is to discard much of the data. A p -variable fit to $p + 1$ observations will perfectly predict Y as long as no two observations have the same Y . Such a model will, however, yield predictions that appear almost random with respect to responses on a different dataset. Therefore, unbiased estimates of predictive accuracy are essential.

Model validation is done to ascertain whether predicted values from the model are likely to accurately predict responses on future subjects or subjects not used to develop our model. Three major causes of failure of the model to validate are overfitting, changes in measurement methods/changes in definition of categorical variables, and major changes in subject inclusion criteria. 8

There are two major modes of model validation, *external* and *internal*. The most stringent external validation involves testing a final model developed in one country or setting on subjects in another country or setting at another time. This validation would test whether the data collection instrument was translated into another language properly, whether cultural differences make earlier findings nonapplicable, and whether secular trends have changed associations or base rates. Testing a finished model on new subjects from the same geographic area but from a different institution as subjects used to fit the model is a less stringent form of external validation. The least stringent form of external validation involves using the first m of n observations for model training and using the remaining $n - m$ observations as a test sample. This is very similar to data-splitting (Section 5.3.3). For details about methods for external validation see the R `val.prob` and `val.surv` functions in the `rms` package. 9

Even though external validation is frequently favored by non-statisticians, it is often problematic. Holding back data from the model-fitting phase re-

sults in lower precision and power, and one can increase precision and learn more about geographic or time differences by fitting a unified model to the entire subject series including, for example, country or calendar time as a main effect and/or as an interacting effect. Indeed one could use the following working definition of external validation: validation of a prediction tool using data that were not available when the tool needed to be completed. An alternate definition could be taken as the validation of a prediction tool by an independent research team.

One suggested hierarchy of the quality of various validation methods is as follows, ordered from worst to best.

1. Attempting several validations (internal or external) and reporting only the one that “worked”
2. Reporting apparent performance on the training dataset (no validation)
3. Reporting predictive accuracy on an undersized independent test sample
4. Internal validation using data-splitting where at least one of the training and test samples is not huge and the investigator is not aware of the arbitrariness of variable selection done on a single sample
5. Strong internal validation using 100 repeats of 10-fold cross-validation or several hundred bootstrap resamples, repeating *all* analysis steps involving Y afresh at each re-sample and the arbitrariness of selected “important variables” is reported (if variable selection is used)
6. External validation on a large test sample, done by the original research team
7. Re-analysis by an independent research team using strong internal validation of the original dataset
8. External validation using new test data, done by an independent research team
9. External validation using new test data generated using different instruments/technology, done by an independent research team

Internal validation involves fitting and validating the model by carefully using one series of subjects. One uses the combined dataset in this way to estimate the likely performance of the final model on new subjects, which after all is often of most interest. Most of the remainder of Section 5.3 deals with internal validation.

5.3.2 Which Quantities Should Be Used in Validation?

For ordinary multiple regression models, the R^2 index is a good measure of the model’s predictive ability, especially for the purpose of quantifying drop-off in predictive ability when applying the model to other datasets. R^2 is biased, however. For example, if one used nine predictors to predict outcomes of 10 subjects, $R^2 = 1.0$ but the R^2 that will be achieved on future

subjects will be close to zero. In this case, dramatic overfitting has occurred. The *adjusted* R^2 (Equation 4.4) solves this problem, at least when the model has been completely prespecified and no variables or parameters have been “screened” out of the final model fit. That is, R_{adj}^2 is only valid when p in its formula is honest—when it includes all parameters ever examined (formally or informally, e.g., using graphs or tables) whether these parameters are in the final model or not.

Quite often we need to validate indexes other than R^2 for which adjustments for p have not been created.^c We also need to validate models containing “phantom degrees of freedom” that were screened out earlier, formally or informally. For these purposes, we obtain nearly unbiased estimates of R^2 or other indexes using data splitting, cross-validation, or the bootstrap. The bootstrap provides the most precise estimates.

The g -index is another discrimination measure to validate. But g and R^2 measures only one aspect of predictive ability. In general, there are two major aspects of predictive accuracy that need to be assessed. As discussed in Section 4.5, *calibration* or *reliability* is the ability of the model to make unbiased estimates of outcome. *Discrimination* is the model’s ability to separate subjects’ outcomes. Validation of the model is recommended even when a data reduction technique is used. This is a way to ensure that the model was not overfitted or is otherwise inaccurate.

5.3.3 Data-Splitting

The simplest validation method is one-time *data-splitting*. Here a dataset is split into *training* (model development) and *test* (model validation) samples by a random process with or without balancing distributions of the response and predictor variables in the two samples. In some cases, a chronological split is used so that the validation is prospective. The model’s calibration and discrimination are validated in the test set.

In ordinary least squares, calibration may be assessed by, for example, plotting Y against \hat{Y} . Discrimination here is assessed by R^2 and it is of interest in comparing R^2 in the training sample with that achieved in the test sample. A drop in R^2 indicates overfitting, and the absolute R^2 in the test sample is an unbiased estimate of predictive discrimination. Note that in extremely overfitted models, R^2 in the test set can be negative, since it is computed on “frozen” intercept and regression coefficients using the formula $1 - SSE/SST$, where SSE is the error sum of squares, SST is the total sum

^c For example, in the binary logistic model, there is a generalization of R^2 available, but no adjusted version. For logistic models we often validate other indexes such as the ROC area or rank correlation between predicted probabilities and observed outcomes. We also validate the calibration accuracy of \hat{Y} in predicting Y .

of squares, and SSE can be greater than SST (when predictions are worse than the constant predictor \bar{Y}).

10

To be able to validate predictions from the model over an entire test sample (without validating it separately in particular subsets such as in males and females), the test sample must be large enough to precisely fit a model containing one predictor. For a study with a continuous uncensored response variable, the test sample size should ordinarily be ≥ 100 at a bare minimum. For survival time studies, the test sample should at least be large enough to contain a minimum of 100 outcome events. For binary outcomes, the test sample should contain a bare minimum of 100 subjects in the least frequent outcome category. Once the size of the test sample is determined, the remaining portion of the original sample can be used as a training sample. Even with these test sample sizes, validation of extreme predictions is difficult.

Data-splitting has the advantage of allowing hypothesis tests to be confirmed in the test sample. However, it has the following disadvantages.

1. Data-splitting greatly reduces the sample size for both model development and model testing. Because of this, Roecker⁵²⁸ found this method “appears to be a costly approach, both in terms of predictive accuracy of the fitted model and the precision of our estimate of the accuracy.” Breiman [66, Section 1.3] found that bootstrap validation on the original sample was as efficient as having a separate test sample twice as large³⁶.
2. It requires a larger sample to be held out than cross-validation (see below) to be able to obtain the same precision of the estimate of predictive accuracy.
3. The split may be fortuitous; if the process were repeated with a different split, different assessments of predictive accuracy may be obtained.
4. Data-splitting does not validate the final model, but rather a model developed on only a subset of the data. The training and test sets are recombined for fitting the final model, which is not validated.
5. Data-splitting requires the split before the *first* analysis of the data. With other methods, analyses can proceed in the usual way on the complete dataset. Then, after a “final” model is specified, the modeling process is rerun on multiple resamples from the original data to mimic the process that produced the “final” model.

5.3.4 Improvements on Data-Splitting: Resampling

Bootstrapping, jackknifing, and other resampling plans can be used to obtain nearly unbiased estimates of model performance without sacrificing sample size. These methods work when either the model is completely specified except for the regression coefficients, or all important steps of the modeling process, especially variable selection, are automated. Only then can each

bootstrap replication be a reflection of all sources of variability in modeling. Note that most analyses involve examination of graphs and testing for lack of model fit, with many intermediate decisions by the analyst such as simplification of interactions. These processes are difficult to automate. But variable selection alone is often the greatest source of variability because of multiple comparison problems, so the analyst must go to great lengths to bootstrap or jackknife variable selection.

The ability to study the arbitrariness of how a stepwise variable selection algorithm selects “important” factors is a major benefit of bootstrapping. A useful display is a matrix of blanks and asterisks, where an asterisk is placed in column x of row i if variable x is selected in bootstrap sample i (see p. 263 for an example). If many variables appear to be selected at random, the analyst may want to turn to a data reduction method rather than using stepwise selection (see also [541]).

Cross-validation is a generalization of data-splitting that solves some of the problems of data-splitting. *Leave-out-one cross-validation*,^{565, 633} the limit of cross-validation, is similar to jackknifing.⁶⁷⁵ Here one observation is omitted from the analytical process and the response for that observation is predicted using a model derived from the remaining $n - 1$ observations. The process is repeated n times to obtain an average accuracy. Efron¹⁷² reports that grouped cross-validation is more accurate; here groups of k observations are omitted at a time. Suppose, for example, that 10 groups are used. The original dataset is divided into 10 equal subsets at random. The first 9 subsets are used to develop a model (transformation selection, interaction testing, stepwise variable selection, etc. are all done). The resulting model is assessed for accuracy on the remaining 1/10th of the sample. This process is repeated at least 10 times to get an average of 10 indexes such as R^2 .

11

A drawback of cross-validation is the choice of the number of observations to hold out from each fit. Another is that the number of repetitions needed to achieve accurate estimates of accuracy often exceeds 200. For example, one may have to omit $\frac{1}{10}$ th of the sample 500 times to accurately estimate the index of interest. Thus the sample would need to be split into tenths 50 times. Another possible problem is that cross-validation may not fully represent the variability of variable selection. If 20 subjects are omitted each time from a sample of size 1000, the lists of variables selected from each training sample of size 980 are likely to be much more similar than lists obtained from fitting independent samples of 1000 subjects. Finally, as with data-splitting, cross-validation does not validate the full 1000-subject model.

12

An interesting way to study overfitting could be called the randomization method. Here we ask the question “How well can the response be predicted when we use our best procedure on random responses when the predictive accuracy should be near zero?” The better the fit on random Y , the worse the overfitting. The method takes a random permutation of the response variable and develops a model with optional variable selection based on the original X and permuted Y . Suppose this yields $R^2 = .2$ for the fitted sample. Apply the

fit to the original data to estimate optimism. If overfitting is not a problem, R^2 would be the same for both fits and it will ordinarily be very near zero.

13

5.3.5 Validation Using the Bootstrap

Efron,^{172,173} Efron and Gong,¹⁷⁵ Gong,²²⁴ Efron and Tibshirani,^{177,178} Linnet,⁴¹⁶ and Breiman⁶⁶ describe several bootstrapping procedures for obtaining nearly unbiased estimates of future model performance without holding back data when making the final estimates of model parameters. With the “simple bootstrap” [178, p. 247], one repeatedly fits the model in a bootstrap sample and evaluates the performance of the model on the original sample. The estimate of the likely performance of the final model on future data is estimated by the average of all of the indexes computed on the original sample.

Efron showed that an enhanced bootstrap estimates future model performance more accurately than the simple bootstrap. Instead of estimating an accuracy index directly from averaging indexes computed on the original sample, the enhanced bootstrap uses a slightly more indirect approach by estimating the bias due to overfitting or the “optimism” in the final model fit. After the optimism is estimated, it can be subtracted from the index of accuracy derived from the original sample to obtain a bias-corrected or overfitting-corrected estimate of predictive accuracy. The bootstrap method is as follows. From the original X and Y in the sample of size n , draw a sample with replacement also of size n . Derive a model in the bootstrap sample and apply it without change to the original sample. The accuracy index from the bootstrap sample minus the index computed on the original sample is an estimate of optimism. This process is repeated for 100 or so bootstrap replications to obtain an average optimism, which is subtracted from the final model fit’s apparent accuracy to obtain the overfitting-corrected estimate.

14

Note that bootstrapping validates the *process* that was used to fit the original model (as does cross-validation). It provides an estimate of the *expected value* of the optimism, which when subtracted from the original index, provides an estimate of the *expected* bias-corrected index. If stepwise variable selection is part of the bootstrap process (as it must be if the final model is developed that way), and not all resamples (samples with replacement or training samples in cross-validation) resulted in the same model (which is almost always the case), this internal validation process actually provides an unbiased estimate of the future performance of the *process* used to identify markers and scoring systems; it does not validate a single final model. But resampling does tend to provide good estimates of the future performance of the final model that was selected using the same procedure repeated in the resamples.

15

Note that by drawing samples from X and Y , we are estimating aspects of the *unconditional* distribution of statistical quantities. One could instead draw samples from quantities such as residuals from the model to obtain a distribution that is conditional on X . However, this approach requires that the model be specified correctly, whereas the unconditional bootstrap does not. Also, the unconditional estimators are similar to conditional estimators except for very skewed or very small samples [186, p. 217].

Bootstrapping can be used to estimate the optimism in virtually any index. Besides discrimination indexes such as R^2 , slope and intercept calibration factors can be estimated. When one fits the model $C(Y|X) = X\beta$, and then refits the model $C(Y|X) = \gamma_0 + \gamma_1 X\hat{\beta}$ on the same data, where $\hat{\beta}$ is an estimate of β , $\hat{\gamma}_0$ and $\hat{\gamma}_1$ will necessarily be 0 and 1, respectively. However, when $\hat{\beta}$ is used to predict responses on another dataset, $\hat{\gamma}_1$ may be < 1 if there is overfitting, and $\hat{\gamma}_0$ will be different from zero to compensate. Thus a bootstrap estimate of γ_1 will not only quantify overfitting nicely, but can also be used to shrink predicted values to make them more calibrated (similar to [582]). Efron's optimism bootstrap is used to estimate the optimism in $(0, 1)$ and then (γ_0, γ_1) are estimated by subtracting the optimism in the constant estimator $(0, 1)$. Note that in cross-validation one estimates β with $\hat{\beta}$ from the training sample and fits $C(Y|X) = \gamma X\hat{\beta}$ on the test sample directly. Then the γ estimates are averaged over all test samples. This approach does not require the choice of a parameter that determines the amount of shrinkage as does ridge regression or penalized maximum likelihood estimation; instead one estimates how to make the initial fit well calibrated.^{123,633} However, this approach is not as reliable as building shrinkage into the original estimation process. The latter allows different parameters to be shrunk by different factors.

16

Ordinary bootstrapping can sometimes yield overly optimistic estimates of optimism, that is, may underestimate the amount of overfitting. This is especially true when the ratio of the number of observations to the number of parameters estimated is not large.²⁰⁵ A variation on the bootstrap that improves precision of the assessment is the “.632” method, which Efron found to be optimal in several examples.¹⁷² This method provides a bias-corrected estimate of predictive accuracy by substituting $0.632 \times [\text{apparent accuracy} - \hat{\epsilon}_0]$ for the estimate of optimism, where $\hat{\epsilon}_0$ is a weighted average of accuracies evaluated on observations *omitted* from bootstrap samples [178, Eq.17.25, p. 253].

17

For ordinary least squares, where the genuine per-observation .632 estimator can be used, several simulations revealed close agreement with the modified .632 estimator, even in small, highly overfitted samples. In these overfitted cases, the ordinary bootstrap bias-corrected accuracy estimates were significantly higher than the .632 estimates. Simulations^{259,591} have shown, however, that for most types of indexes of accuracy of binary logistic regression models, Efron's original bootstrap has lower mean squared error than the .632 bootstrap when $n = 200, p = 30$. Bootstrap overfitting-corrected estimates of model performance can be biased in favor of the model. Although

18

Table 5.2 Example validation with and without variable selection

Method	Apparent Rank Correlation of Predicted vs. Observed	Over- Optimism	Bias-Corrected Correlation
Full Model	0.50	0.06	0.44
Stepwise Model	0.47	0.05	0.42

cross-validation is less biased than the bootstrap, Efron¹⁷² showed that it has much higher variance in estimating overfitting-corrected predictive accuracy than bootstrapping. In other words, cross-validation, like data-splitting, can yield significantly different estimates when the entire validation process is repeated.

It is frequently very informative to estimate a measure of predictive accuracy forcing all candidate factors into the fit and then to separately estimate accuracy allowing stepwise variable selection, possibly with different stopping rules. Consistent with Spiegelhalter's proposal to use all factors and then to shrink the coefficients to adjust for overfitting,⁵⁸² the full model fit will outperform the stepwise model more often than not. Even though stepwise modeling has slightly less optimism in predictive discrimination, this improvement is not enough to offset the loss of information from deleting even marginally important variables. Table 5.2 shows a typical scenario. In this example, stepwise modeling lost a possible $0.50 - 0.47 = 0.03$ predictive discrimination. The full model fit will especially be an improvement when

1. the stepwise selection deletes several variables that are almost significant;
2. these marginal variables have *some* real predictive value, even if it's slight; and
3. there is no small set of extremely dominant variables that would be easily found by stepwise selection.

19

Faraway¹⁸⁶ has a fascinating study showing how resampling methods can be used to estimate the distributions of predicted values and of effects of a predictor, adjusting for an automated multistep modeling process. Bootstrapping can be used, for example, to penalize the variance in predicted values for choosing a transformation for Y and for outlier and influential observation deletion, in addition to variable selection. Estimation of the transformation of Y greatly increased the variance in Faraway's examples. Brownstone [77, p. 74] states that "In spite of considerable efforts, theoretical statisticians have been unable to analyze the sampling properties of [usual multistep modeling strategies] under realistic conditions" and concludes that the modeling strategy must be completely specified and then bootstrapped to get consistent estimates of variances and other sampling properties.

20

5.4 Bootstrapping Ranks of Predictors

When the order of importance of predictors is not pre-specified but the researcher attempts to determine that order by assessing multiple associations with Y , the process of selecting “winners” and “losers” is unreliable. The bootstrap can be used to demonstrate the difficulty of this task, by estimating confidence intervals for the ranks of all the predictors. Even though the bootstrap intervals are wide, they actually underestimate the true widths²⁵⁰.

The following example uses simulated data with known ranks of importance of 12 predictors, using an ordinary linear model. The importance metric is the partial χ^2 minus its degrees of freedom, while the true metric is the partial β , as all covariates have $U(0, 1)$ distributions.

```
# Use the plot method for anova, with pl=FALSE to suppress
# actual plotting of chi-square - d.f. for each bootstrap
# repetition. Rank the negative of the adjusted chi-squares
# so that a rank of 1 is assigned to the highest. It is
# important to tell plot.anova.rms not to sort the results,
# or every bootstrap replication would have ranks of 1,2,3,
# ... for the partial test statistics.
require(rms)
n <- 300
set.seed(1)
d <- data.frame(x1=runif(n), x2=runif(n), x3=runif(n),
               x4=runif(n), x5=runif(n), x6=runif(n), x7=runif(n),
               x8=runif(n), x9=runif(n), x10=runif(n), x11=runif(n),
               x12=runif(n))
d$y <- with(d, 1*x1 + 2*x2 + 3*x3 + 4*x4 + 5*x5 + 6*x6 +
            7*x7 + 8*x8 + 9*x9 + 10*x10 + 11*x11 +
            12*x12 + 9*rnorm(n))

f <- ols(y ~ x1+x2+x3+x4+x5+x6+x7+x8+x9+x10+x11+x12, data=d)
B <- 1000
ranks <- matrix(NA, nrow=B, ncol=12)
rankvars <- function(fit)
  rank(plot(anova(fit), sort='none', pl=FALSE))
Rank <- rankvars(f)
for(i in 1:B) {
  j <- sample(1:n, n, TRUE)
  bootfit <- update(f, data=d, subset=j)
  ranks[i,] <- rankvars(bootfit)
}
lim <- t(apply(ranks, 2, quantile, probs=c(.025,.975)))
predictor <- factor(names(Rank), names(Rank))
w <- data.frame(predictor, Rank, lower=lim[,1], upper=lim[,2])
require(ggplot2)
ggplot(w, aes(x=predictor, y=Rank)) + geom_point() +
  coord_flip() + scale_y_continuous(breaks=1:12) +
  geom_errorbar(aes(ymin=lim[,1], ymax=lim[,2]), width=0)
```

With a sample size of $n = 300$ the observed ranks of predictor importance do not coincide with population β s, even when there are no collinearities among

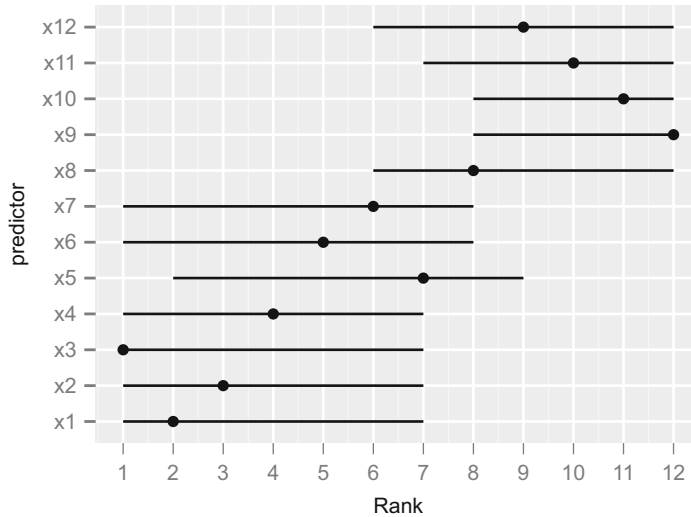


Fig. 5.3 Bootstrap percentile 0.95 confidence limits for ranks of predictors in an OLS model. Ranking is on the basis of partial χ^2 minus d.f. Point estimates are original ranks

the predictors. Confidence intervals are wide; for example the 0.95 confidence interval for the rank of x_7 (which has a true rank of 7) is $[1, 8]$, so we are only confident that x_7 is not one of the 4 most influential predictors. The confidence intervals do include the true ranks in each case (Figure 5.3).

5.5 Simplifying the Final Model by Approximating It

5.5.1 Difficulties Using Full Models

A model that contains all prespecified terms will usually be the one that predicts the most accurately on new data. It is also a model for which confidence limits and statistical tests have the claimed properties. Often, however, this model will not be very parsimonious. The full model may require more predictors than the researchers care to collect in future samples. It also requires predicted values to be conditional on all of the predictors, which can increase the variance of the predictions.

As an example suppose that least squares has been used to fit a model containing several variables including race (with four categories). Race may be an insignificant predictor and may explain a tiny fraction of the observed variation in Y . Yet when predictions are requested, a value for race must be inserted. If the subject is of the majority race, and this race has a majority of,

say 0.75, the variance of the predicted value will not be significantly greater than the variance for a predicted value from a model that excluded race for its list of predictors. If, however, the subject is of a minority race (say “other” with a prevalence of 0.01), the predicted value will have much higher variance. One approach to this problem, that does not require development of a second model, is to ignore the subject’s race and to get a weighted average prediction. That is, we obtain predictions for each of the four races and weight these predictions by the relative frequencies of the four races.^d This weighted average estimates the expected value of Y unconditional on race. It has the advantage of having exactly correct confidence limits when model assumptions are satisfied, because the correct “error term” is being used (one that deducts 3 d.f. for having ever estimated the race effect). In regression models having nonlinear link functions, this process does not yield such a simple interpretation.

When predictors are collinear, their competition results in larger P -values when predictors are (often inappropriately) tested individually. Likewise, confidence intervals for individual effects will be wide and uninterpretable (can other variables really be held constant when one is changed?).

5.5.2 Approximating the Full Model

When the full model contains several predictors that do not appreciably affect the predictions, the above process of “unconditioning” is unwieldy. In the search for a simple solution, the most commonly used procedure for making the model parsimonious is to remove variables on the basis of P -values, but this results in a variety of problems as we have seen. Our approach instead is to consider the full model fit as the “gold standard” model, especially the model from which formal inferences are made. We then proceed to approximate this full model to any desired degree of accuracy. For any approximate model we calculate the accuracy with which it approximates the best model. One goal this process accomplishes is that it provides different degrees of parsimony to different audiences, based on their needs. One investigator may be able to collect only three variables, another one seven. Each investigator will know how much she is giving up by using a subset of the predictors. In approximating the gold standard model it is very important to note that there is nothing gained in removing certain nonlinear terms; gains in parsimony come only from removing entire predictors. Another accomplishment of model approximation is that when the full model has been fitted using

^d Using the `rms` package described in Chapter 6, such estimates and their confidence limits can easily be obtained, using for example `contrast(fit, list(age=50, disease='hypertension', race=levels(race)), type='average', weights=table(race))`.

shrinkage (penalized estimation, Section 9.10), the approximate models will inherit the shrinkage (see Section 14.10 for an example).

Approximating complex models with simpler ones has been used to decode “black boxes” such as artificial neural networks. Recursive partitioning trees (Section 2.5) may sometimes be used in this context. One develops a regression tree to predict the predicted value $X\hat{\beta}$ on the basis of the unique variables in X , using R^2 , the average absolute prediction error, or the maximum absolute prediction error as a stopping rule, for example¹⁸⁴. The user desiring simplicity may use the tree to obtain predicted values, using the first k nodes, with k just large enough to yield a low enough absolute error in predicting the more comprehensive prediction. Overfitting is not a problem as it is when the tree procedure is used to predict the outcome, because (1) given the predictor values the predictions are deterministic and (2) the variable being predicted is a continuous, completely observed variable. Hence the best cross-validating tree approximation will be one with one subject per node. One advantage of the tree-approximation procedure is that data collection on an individual subject whose outcome is being predicted may be abbreviated by measuring only those X s that are used in the top nodes, until the prediction is resolved to within a tolerable error.

When principal component regression is being used, trees can also be used to approximate the components or to make them more interpretable.

Full models may also be approximated using least squares as long as the linear predictor $X\hat{\beta}$ is the target, and not some nonlinear transformation of it such as a logistic model probability. When the original model was fitted using unpenalized least squares, submodels fitted against \hat{Y} will have the same coefficients as if least squares had been used to fit the subset of predictors directly against Y . To see this, note that if X denotes the entire design matrix and T denotes a subset of the columns of X , the coefficient estimates for the full model are $(X'X)^{-1}X'Y$, $\hat{Y} = X(X'X)^{-1}X'Y$, estimates for a reduced model fitted against Y are $(T'T)^{-1}T'Y$, and coefficients fitted against \hat{Y} are $(T'T)T'X(X'X)^{-1}X'Y$ which can be shown to equal $(T'T)^{-1}T'Y$.

When least squares is used for both the full and reduced models, the variance–covariance matrix of the coefficient estimates of the reduced model is $(T'T)^{-1}\sigma^2$, where the residual variance σ^2 is estimated using the *full* model. When σ^2 is estimated by the unbiased estimator using the d.f. from the full model, which provides the only unbiased estimate of σ^2 , the estimated variance–covariance matrix of the reduced model will be appropriate (unlike that from stepwise variable selection) although the bootstrap may be needed to fully take into account the source of variation due to how the approximate model was selected.

So if in the least squares case the approximate model coefficients are identical to coefficients obtained upon fitting the reduced model against Y , how is model approximation any different from stepwise variable selection? There are several differences, in addition to how σ^2 is estimated.

1. When the full model is approximated by a backward step-down procedure against \hat{Y} , the stopping rule is less arbitrary. One stops deleting variables when deleting any further variable would make the approximation inadequate (e.g., the R^2 for predictions from the reduced model against the original \hat{Y} drops below 0.95).
2. Because the stopping rule is different (i.e., is not based on P -values), the approximate model will have a different number of predictors than an ordinary stepwise model.
3. If the original model used penalization, approximate models will inherit the amount of shrinkage used in the full fit.

Typically, though, if one performed ordinary backward step-down against Y using a large cutoff for α (e.g., 0.5), the approximate model would be very similar to the step-down model. The main difference would be the use of a larger estimate of σ^2 and smaller error d.f. than are used for the ordinary step-down approach (an estimate that pretended the final reduced model was prespecified).

When the full model was not fitted using least squares, least squares can still easily be used to approximate the full model. If the coefficient estimates from the full model are $\hat{\beta}$, estimates from the approximate model are matrix contrasts of $\hat{\beta}$, namely, $W\hat{\beta}$, where $W = (T'T)^{-1}T'X$. So the variance-covariance matrix of the reduced coefficient estimates is given by

$$WVW', \quad (5.2)$$

where V is the variance matrix for $\hat{\beta}$. See Section 19.5 for an example. Ambler et al.²¹ studied model simplification using simulation studies based on several clinical datasets, and compared it with ordinary backward stepdown variable selection and with shrinkage methods such as the *lasso* (see Section 4.3). They found that ordinary backwards variable selection can be competitive when there is a large fraction of truly irrelevant predictors (something that can be difficult to know in advance). Paul et al.⁴⁸⁵ found advantages to modeling the response with a complex but reliable approach, and then developing a parsimonious model using the *lasso* or stepwise variable selection against \hat{Y} . See Section 11.7 for a case study in model approximation.

5.6 Further Reading

- [1] Gelman²¹³ argues that continuous variables should be scaled by two standard deviations to make them comparable to binary predictors. However his approach assumes linearity in the predictor effect and assumes the prevalence of the binary predictor is near 0.5. John Fox [202, p. 95] points out that if two predictors are on the same scale and have the same impact (e.g., years of employment and years of education), standardizing the coefficients will make them appear to have different impacts.

- [2] Levine et al.⁴⁰¹ have a compelling argument for graphing effect ratios on a logarithmic scale.
- [3] Hankins²⁵⁴ is a definitive reference on nomograms and has multi-axis examples of historical significance. According to Hankins, Maurice d’Ocagne could be called the inventor of the nomogram, starting with alignment diagrams in 1884 and declaring a new science of “nomography” in 1899. d’Ocagne was at *École des Ponts et Chaussées*, a French civil engineering school. Julien and Hanley³²⁸ have a nice example of adding axes to a nomogram to estimate the absolute effect of a treatment estimated using a Cox proportional hazards model. Kattan and Marasco³³⁹ have several clinical examples and explain advantages to the user of nomograms over “black box” computerized prediction.
- [4] Graham and Clavel²³¹ discuss graphical and tabular ways of obtaining risk estimates. van Gorp et al.⁶³⁰ have a nice example of a score chart for manually obtaining estimates.
- [5] Larsen and Merlo³⁷⁵ developed a similar measure—the median odds ratio. Gönen and Heller²²³ developed a *c*-index that like *g* is a function of the covariate distribution.
- [6] Booth and Sarkar⁶¹ have a nice analysis of the number of bootstrap resamples needed to guarantee with 0.95 confidence that a variance estimate has a sufficiently small relative error. They concentrate on the Monte Carlo simulation error, showing that small errors in variance estimates can lead to important differences in *P*-values. Canty et al.⁹¹ provide a number of diagnostics to check the reliability of bootstrap calculations.
- [7] There are many variations on the basic bootstrap for computing confidence limits.^{150,178} See Booth and Sarkar⁶¹ for useful information about choosing the number of resamples. They report the number of resamples necessary to not appreciably change *P*-values, for example. Booth and Sarkar propose a more conservative number of resamples than others use (e.g., 800 resamples) for estimating variances. Carpenter and Bithell⁹² have an excellent overview of bootstrap confidence intervals, with practical guidance. They also have a good discussion of the unconditional nonparametric bootstrap versus the conditional semiparametric bootstrap.
- [8] Altman and Royston¹⁸ have a good general discussion of what it means to validate a predictive model, including issues related to study design and consideration of uses to which the model will be put.
- [9] An excellent paper on external validation and generalizability is Justice et al.³²⁹. Bleeker et al.⁵⁸ provide an example where internal validation is misleading when compared with a true external validation done using subjects from different centers in a different time period. Vergouwe et al.⁶³⁸ give good guidance about the number of events needed in sample used for external validation of binary logistic models.
- [10] See Picard and Berk⁵⁰⁵ for more about data-splitting.
- [11] In the context of variable selection where one attempts to select the set of variables with nonzero true regression coefficients in an ordinary regression model, Shao⁵⁶⁵ demonstrated that leave-out-one cross-validation selects models that are “too large.” Shao also showed that the number of observations held back for validation should often be larger than the number used to train the model. This is because in this case one is not interested in an accurate model (you fit the whole sample to do that), but an accurate estimate of prediction error is mandatory so as to know which variables to allow into the final model. Shao suggests using a cross-validation strategy in which approximately $n^{3/4}$ observations are used in each training sample and the remaining observations are used in the test sample. A repeated balanced or Monte Carlo splitting approach is used, and accuracy estimates are averaged over $2n$ (for the Monte Carlo method) repeated splits.

- [12] Picard and Cook's Monte Carlo cross-validation procedure⁵⁰⁶ is an improvement over ordinary cross-validation.
- [13] The randomization method is related to Kipnis' "chaotization relevancy principle"³⁴⁸ in which one chooses between two models by measuring how far each is from a nonsense model. Tibshirani and Knight also use a randomization method for estimating the optimism in a model fit.⁶¹¹
- [14] This method used here is a slight change over that presented in [172], where Efron wrote predictive accuracy as a sum of per-observation components (such as 1 if the observation is classified correctly, 0 otherwise). Here we are writing $m \times$ the unitless summary index of predictive accuracy in the place of Efron's sum of m per-observation accuracies [416, p. 613].
- [15] See [633] and [66, Section 4] for insight on the meaning of expected optimism.
- [16] See Copas,¹²³ van Houwelingen and le Cessie [633, p. 1318], Verweij and van Houwelingen,⁶⁴⁰ and others⁶³¹ for other methods of estimating shrinkage coefficients.
- [17] Efron¹⁷² developed the ".632" estimator only for the case where the index being bootstrapped is estimated on a per-observation basis. A natural generalization of this method can be derived by assuming that the accuracy evaluated on observation i that is omitted from a bootstrap sample has the same expectation as the accuracy of any other observation that would be omitted from the sample. The modified estimate of ϵ_0 is then given by

$$\hat{\epsilon}_0 = \sum_{i=1}^B w_i T_i, \quad (5.3)$$

where T_i is the accuracy estimate derived from fitting a model on the i th bootstrap sample and evaluating it on the observations omitted from that bootstrap sample, and w_i are weights derived for the B bootstrap samples:

$$w_i = \frac{1}{n} \sum_{j=1}^n \frac{[\text{bootstrap sample } i \text{ omits observation } j]}{\#\text{bootstrap samples omitting observation } j}. \quad (5.4)$$

Note that $\hat{\epsilon}_0$ is undefined if any observation is included in every bootstrap sample. Increasing B will avoid this problem. This modified ".632" estimator is easy to compute if one assembles the bootstrap sample assignments and computes the w_i before computing the accuracy indexes T_i . For large n , the w_i approach $1/B$ and so $\hat{\epsilon}_0$ becomes equivalent to the accuracy computed on the observations not contained in the bootstrap sample and then averaged over the B repetitions.

- [18] Efron and Tibshirani¹⁷⁹ have reduced the bias of the ".632" estimator further with only a modest increase in its variance. Simulation has, however, shown no advantage of this ".632+" method over the basic optimism bootstrap for most accuracy indexes used in logistic models.
- [19] van Houwelingen and le Cessie⁶³³ have several interesting developments in model validation. See Breiman⁶⁶ for a discussion of the choice of X for which to validate predictions. Steyerberg et al.⁵⁸⁷ present simulations showing the number of bootstrap samples needed to obtain stable estimates of optimism of various accuracy measures. They demonstrate that bootstrap estimates of optimism are nearly unbiased when compared with simulated external estimates. They also discuss problems with precision of estimates of accuracy, especially when using external validation on small samples.
- [20] Blettner and Sauerbrei also demonstrate the variability caused by data-driven analytic decisions.⁵⁹ Chatfield¹⁰⁰ has more results on the effects of using the data to select the model.

5.7 Problem

Perform a simulation study to understand the performance of various internal validation methods for binary logistic models. Modify the R code below in at least two meaningful ways with regard to covariate distribution or number, sample size, true regression coefficients, number of resamples, or number of times certain strategies are averaged. Interpret your findings and give recommendations for best practice for the type of configuration you studied. The R code from this assignment may be downloaded from the RMS course wiki page.

For each of 200 simulations, the code below generates a training sample of 200 observations with p predictors ($p = 15$ or 30) and a binary response. The predictors are independently $U(-0.5, 0.5)$. The response is sampled so as to follow a logistic model where the intercept is zero and all regression coefficients equal 0.5. The “gold standard” is the predictive ability of the fitted model on a test sample containing 50,000 observations generated from the same population model. For each of the 200 simulations, several validation methods are employed to estimate how the training sample model predicts responses in the 50,000 observations. These validation methods involve fitting 40 or 200 models in resamples.

g -fold cross-validation is done using the command `validate(f, method='cross', B=g)` using the `rms` package. This was repeated and averaged using an extra loop, shown below.

For bootstrap methods, `validate(f, method='boot' or '.632', B=40 or B=200)` was used. `method='.632'` does Efron’s “.632” method¹⁷⁹, labeled `632a` in the output. An ad-hoc modification of the `.632` method, `632b` was also done. Here a “bias-corrected” index of accuracy is simply the index evaluated in the observation omitted from the bootstrap resample. The “gold standard” external validations were obtained from the `val.prob` function in the `rms` package. The following indexes of predictive accuracy are used:

D_{xy} : Somers’ rank correlation between predicted probability that $Y = 1$ vs. the binary Y values. This equals $2(C - 0.5)$ where C is the “ROC Area” or concordance probability.

D : Discrimination index — likelihood ratio χ^2 divided by the sample size

U : Unreliability index — unitless index of how far the logit calibration curve intercept and slope are from $(0, 1)$

Q : Logarithmic accuracy score — a scaled version of the log-likelihood achieved by the predictive model

Intercept: Calibration intercept on logit scale

Slope: Calibration slope (slope of predicted log odds vs. true log odds)

Accuracy of the various resampling procedures may be estimated by computing the mean absolute errors and the root mean squared errors of estimates (e.g., of D_{xy} from the bootstrap on the 200 observations) against the “gold standard” (e.g., D_{xy} for the fitted 200-observation model achieved in the 50,000 observations).

```

require(rms)
set.seed(1) # so can reproduce results

n      ← 200          # Size of training sample
reps   ← 200          # Simulations
npop   ← 50000        # Size of validation gold standard sample
methods ← c('Boot 40','Boot 200','632a 40','632a 200',
            '632b 40','632b 200','10-fold x 4','4-fold x 10',
            '10-fold x 20','4-fold x 50')
R ← expand.grid(sim     = 1:reps,
               p       = c(15,30),
               method  = methods)
R$Dxy ← R$Intercept ← R$Slope ← R$D ← R$U ← R$Q ←
R$repmeth ← R$B ← NA
R$n ← n

## Function to do r overall reps of B resamples, averaging to
## get estimates similar to as if r*B resamples were done

val ← function(fit, method, B, r) {
  contains ← function(m) length(grep(m, method)) > 0
  meth ← if(contains('Boot')) 'boot' else
          if(contains('fold')) 'crossvalidation' else
          if(contains('632')) '.632'
  z ← 0
  for(i in 1:r) z ← z + validate(fit, method=meth, B=B)[
    c("Dxy","Intercept","Slope","D","U","Q"),
    'index.corrected']
  z/r
}

```

```

for(p in c(15, 30)) {

  ## For each p create the true betas, the design matrix,
  ## and realizations of binary y in the gold standard
  ## large sample
  Beta ← rep(.5, p) # True betas
  X ← matrix(runif(npop*p), nrow=npop) - 0.5
  LX ← matxv(X, Beta)
  Y ← ifelse(runif(npop) ≤ plogis(LX), 1, 0)

  ## For each simulation create the data matrix and
  ## realizations of y
  for(j in 1:reps) {

    ## Make training sample
    x ← matrix(runif(n*p), nrow=n) - 0.5
    L ← matxv(x, Beta)
    y ← ifelse(runif(n) ≤ plogis(L), 1, 0)
    f ← lrm(y ~ x, x=TRUE, y=TRUE)
    beta ← f$coef
    forecast ← matxv(X, beta)
    ## Validate in population
  }
}

```

```

v ← val.prob(logit=forecast, y=Y, pl=FALSE)[
  c("Dxy", "Intercept", "Slope", "D", "U", "Q")]

for(method in methods) {
  repmeth ← 1
  if(method %in% c('Boot 40', '632a 40', '632b 40'))
    B ← 40
  if(method %in% c('Boot 200', '632a 200', '632b 200'))
    B ← 200
  if(method == '10-fold x 4') {
    B ← 10
    repmeth ← 4
  }
  if(method == '4-fold x 10') {
    B ← 4
    repmeth ← 10
  }
  if(method == '10-fold x 20') {
    B ← 10
    repmeth ← 20
  }
  if(method == '4-fold x 50') {
    B ← 4
    repmeth ← 50
  }

  z ← val(f, method, B, repmeth)
  k ← which(R$sim == j & R$p == p & R$method == method)
  if(length(k) != 1) stop('program logic error')
  R[k, names(z)] ← z - v
  R[k, c('B', 'repmeth')] ← c(B=B, repmeth=repmeth)
} # end over methods
} # end over reps
} # end over p

```

Results are best summarized in a multi-way dot chart. Bootstrap nonparametric percentile 0.95 confidence limits are included.

```

statnames ← names(R)[6:11]
w ← reshape(R, direction='long', varying=list(statnames),
  v.names='x', timevar='stat', times=statnames)
w$p ← paste('p', w$p, sep='')
require(lattice)
s ← with(w, summarize(abs(x), llist(p, method, stat),
  smean.cl.boot, stat.name='mae'))
Dotplot(method ~ Cbind(mae, Lower, Upper) | stat*p, data=s,
  xlab='Mean |error|')
s ← with(w, summarize(x^2, llist(p, method, stat),
  smean.cl.boot, stat.name='mse'))
Dotplot(method ~ Cbind(sqrt(mse), sqrt(Lower), sqrt(Upper)) |
  stat*p, data=s,
  xlab=expression(sqrt(MSE)))

```

Chapter 6

R Software

The methods described in this book are useful in any regression model that involves a linear combination of regression parameters. The software that is described below is useful in the same situations. Functions in R⁵²⁰ allow interaction spline functions as well as a wide variety of predictor parameterizations for any regression function, and facilitate model validation by resampling.

R is the most comprehensive tool for general regression models for the following reasons.

1. It is very easy to write R functions for new models, so R has implemented a wide variety of modern regression models.
2. Designs can be generated for any model. There is no need to find out whether the particular modeling function handles what SAS calls “class” variables—dummy variables are generated automatically when an R `category`, `factor`, `ordered`, or `character` variable is analyzed.
3. A single R object can contain all information needed to test hypotheses and to obtain predicted values for new data.
4. R has superior graphics.
5. *Classes* in R make possible the use of generic function names (e.g., `predict`, `summary`, `anova`) to examine fits from a large set of specific model-fitting functions.

R^{44,601,635} is a high-level object-oriented language for statistical analysis with over six thousand packages and tens of thousands of functions available. The R system^{318,520} is the basis for R software used in this text, centered around the Regression Modeling Strategies (`rms`) package²⁶¹. See the Appendix and the Web site for more information about software implementations.

1

6.1 The R Modeling Language

R has a battery of functions that make up a statistical modeling language.⁹⁶

At the heart of the modeling functions is an R *formula* of the form

```
response ~ terms
```

The `terms` represent additive components of a general linear model. Although variables and functions of variables make up the `terms`, the formula refers to additive combinations; for example, when `terms` is `age + blood.pressure`, it refers to $\beta_1 \times \text{age} + \beta_2 \times \text{blood.pressure}$. Some examples of formulas are below.

```
y ~ age + sex           # age + sex main effects
y ~ age + sex + age:sex # add second-order interaction
y ~ age*sex            # second-order interaction +
                      # all main effects
y ~ (age + sex + pressure)^2
                      # age+sex+pressure+age:sex+age:pressure...
y ~ (age + sex + pressure)^2 - sex:pressure
                      # all main effects and all 2nd order
                      # interactions except sex:pressure
y ~ (age + race)*sex   # age+race+sex+age:sex+race:sex
y ~ treatment*(age*race + age*sex)
                      # no interact. with race,sex
sqrt(y) ~ sex*sqrt(age) + race
# functions, with dummy variables generated if
# race is an R factor (classification) variable
y ~ sex + poly(age,2) # poly makes orthogonal polynomials
race.sex <- interaction(race,sex)
y ~ age + race.sex    # if desire dummy variables for all
                      # combinations of the factors
```

The formula for a regression model is given to a modeling function; for example,

```
lrm(y ~ rcs(x,4))
```

is read “use a logistic regression model to model `y` as a function of `x`, representing `x` by a restricted cubic spline with four default knots.”^a You can use the R function `update` to refit a model with changes to the model terms or the data used to fit it:

```
f <- lrm(y ~ rcs(x,4) + x2 + x3)
f2 <- update(f, subset=sex=="male")
f3 <- update(f, .~-x2)           # remove x2 from model
f4 <- update(f, .~. + rcs(x5,5)) # add rcs(x5,5) to model
f5 <- update(f, y2 ~ .)         # same terms, new response var.
```

^a `lrm` and `rcs` are in the `rms` package.

6.2 User-Contributed Functions

In addition to the many functions that are packaged with R, a wide variety of user-contributed functions is available on the Internet (see the Appendix or Web site for addresses). Two packages of functions used extensively in this text are `Hmisc`²⁰ and `rms` written by the author. The `Hmisc` package contains miscellaneous functions such as `varclus`, `spearman2`, `transcan`, `hoeffd`, `rcspline.eval`, `impute`, `cut2`, `describe`, `sas.get`, `latex`, and several power and sample size calculation functions. The `varclus` function uses the R `hclust` hierarchical clustering function to do variable clustering, and the R `plclust` function to draw dendrograms depicting the clusters. `varclus` offers a choice of three similarity measures (Pearson r^2 , Spearman ρ^2 , and Hoeffding D) and uses pairwise deletion of missing values. `varclus` automatically generates a series of dummy variables for categorical factors. The `Hmisc` `hoeffd` function computes a matrix of Hoeffding D s for a series of variables. The `spearman2` function will do Wilcoxon, Spearman, and Kruskal–Wallis tests and generalizes Spearman’s ρ to detect non-monotonic relationships.

`Hmisc`’s `transcan` function (see Section 4.7) performs a similar function to PROC PRINQUAL in SAS—it uses restricted splines, dummy variables, and canonical variates to transform each of a series of variables while imputing missing values. An option to shrink regression coefficients for the imputation models avoids overfitting for small samples or a large number of predictors. `transcan` can also do multiple imputation and adjust variance–covariance matrices for imputation. See Chapter 8 for an example of using these functions for data reduction.

See the Web site for a list of R functions for correspondence analysis, principal component analysis, and missing data imputation available from other users. Venables and Ripley [635, Chapter 11] provide a nice description of the multivariate methods that are available in R, and they provide several new multivariate analysis functions.

A basic function in `Hmisc` is the `rcspline.eval` function, which creates a design matrix for a restricted (natural) cubic spline using the truncated power basis. Knot locations are optionally estimated using methods described in Section 2.4.6, and two types of normalizations to reduce numerical problems are supported. You can optionally obtain the design matrix for the anti-derivative of the spline function. The `rcspline.restate` function computes the coefficients (after un-normalizing if needed) that translate the restricted cubic spline function to unrestricted form (Equation 2.27). `rcspline.restate` also outputs L^AT_EX and R representations of spline functions in simplified form.

6.3 The rms Package

A package of R functions called `rms` contains several functions that extend R to make the analyses described in this book easy to do. A central function in `rms` is `datadist`, which computes statistical summaries of predictors to automate estimation and plotting of effects. `datadist` exists as a separate function so that the candidate predictors may be summarized once, thus saving time when fitting several models using subsets or different transformations of predictors. If `datadist` is called before model fitting, the distributional summaries are stored with the fit so that the fit is self-contained with respect to later estimation. Alternatively, `datadist` may be called after the fit to create temporary summaries to use as plot ranges and effect intervals, or these ranges may be specified explicitly to `Predict` and `summary` (see below), without ever calling `datadist`. The input to `datadist` may be a data frame, a list of individual predictors, or a combination of the two.

The characteristics saved by `datadist` include the overall range and certain quantiles for continuous variables, and the distinct values for discrete variables (i.e., R factor variables or variables with 10 or fewer unique values). The quantiles and set of distinct values facilitate estimation and plotting, as described later. When a function of a predictor is used (e.g., `po1(pmin(x,50),2)`), the limits saved apply to the innermost variable (here, `x`). When a plot is requested for how `x` relates to the response, the plot will have `x` on the x -axis, not `pmin(x,50)`. The way that defaults are computed can be controlled by the `q.effect` and `q.display` parameters to `datadist`. By default, continuous variables are plotted with ranges determined by the tenth smallest and tenth largest values occurring in the data (if $n < 200$, the 0.05 and 0.95 quantiles are used). The default range for estimating effects such as odds and hazard ratios is the lower and upper quartiles. When a predictor is adjusted to a constant so that the effects of changes in other predictors can be studied, the default constant used is the median for continuous predictors and the most frequent category for factor variables. The R system option `datadist` is used to point to the result returned by the `datadist` function. See the help files for `datadist` for more information.

`rms` fitting functions save detailed information for later prediction, plotting, and testing. `rms` also allows for special restricted interactions and sets the default method of generating contrasts for categorical variables to "`contr.-treatment`", the traditional dummy-variable approach.

`rms` has a special operator `%ia%` in the terms of a formula that allows for restricted interactions. For example, one may specify a model that contains `sex` and a five-knot linear spline for `age`, but restrict the `age × sex` interaction to be linear in `age`. To be able to connect this incomplete interaction with the main effects for later hypothesis testing and estimation, the following formula would be given:

```
y ~ sex + lsp(age, c(20,30,40,50,60)) +
    sex %ia% lsp(age, c(20,30,40,50,60))
```

Table 6.1 rms Fitting Functions

Function	Purpose	Related R Functions
<code>ols</code>	Ordinary least squares linear model	<code>lm</code>
<code>lrm</code>	Binary and ordinal logistic regression model Has options for penalized MLE	<code>glm</code>
<code>orm</code>	Ordinal semi-parametric regression model with several link functions	<code>polr</code> , <code>lrm</code>
<code>psm</code>	Accelerated failure time parametric survival models	<code>survreg</code>
<code>cph</code>	Cox proportional hazards regression	<code>coxph</code>
<code>bj</code>	Buckley–James censored least squares model	<code>survreg</code> , <code>lm</code>
<code>Glm</code>	General linear models	<code>glm</code>
<code>Gls</code>	Generalized least squares	<code>gls</code>
<code>Rq</code>	Quantile regression	<code>rq</code>

The following expression would restrict the $\text{age} \times \text{cholesterol}$ interaction to be of the form $AF(B) + BG(A)$ by removing doubly nonlinear terms.

```
y ~ lsp(age,30) + rcs(cholesterol,4) +
  lsp(age,30) %ia% rcs(cholesterol,4)
```

`rms` has special fitting functions that facilitate many of the procedures described in this book, shown in Table 6.1.

`Glm` is a slight modification of the built-in R `glm` function so that `rms` methods can be run on the resulting fit object. `glm` fits general linear models under a wide variety of distributions of Y . `Gls` is a modification of the `gls` function from the `nlme` package of Pinheiro and Bates⁵⁰⁹, for repeated measures (longitudinal) and spatially correlated data. The `Rq` function is a modification of the `quantreg` package's `rq` function^{356,357}. Functions related to survival analysis make heavy use of Therneau's `survival` package⁴⁸².

You may want to specify to the fitting functions an option for how missing values (`NA`s) are handled. The method for handling missing data in R is to specify an `na.action` function. Some possible `na.actions` are given in Table 6.2. The default `na.action` is `na.delete` when you use `rms`'s fitting functions. An easy way to specify a new default `na.action` is, for example,

```
options(na.action="na.omit")# don't report frequency of NAs
```

before using a fitting function. If you use `na.delete` you can also use the system option `na.detail.response` that makes model fits store information about Y stratified by whether each X is missing. The default descriptive statistics for Y are the sample size and mean. For a survival time response object the sample size and proportion of events are used. Other summary functions can be specified using the `na.fun.response` option.

Table 6.2 Some `na.actions` Used in `rms`

Function Name	Method Used
<code>na.fail</code>	Stop with error message if any missing values present
<code>na.omit</code>	Function to remove observations with any predictors or responses missing
<code>na.delete</code>	Modified version of <code>na.omit</code> to also report on frequency of NAs for each variable

```
options(na.action="na.delete", na.detail.response=TRUE,
        na.fun.response="mystats")
# Just use na.fun.response="quantile" if don't care about n
mystats <- function(y) {
  z <- quantile(y, na.rm=T)
  n <- sum(!is.na(y))
  c(N=n, z)      # elements named N, 0%, 25%, etc.
}
```

When R deletes missing values during the model-fitting procedure, residuals, fitted values, and other quantities stored with the fit will not correspond row-for-row with observations in the original data frame (which retained NAs). This is problematic when, for example, `age` in the dataset is plotted against the residual from the fitted model. Fortunately, for many `na.actions` including `na.delete` and a modified version of `na.omit`, a class of R functions called `naresid` written by Therneau works behind the scenes to put NAs back into residuals, predicted values, and other quantities when the `predict` or `residuals` functions (see below) are used. Thus for some of the `na.actions`, predicted values and residuals will automatically be arranged to match the original data.

Any R function can be used in the terms for formulas given to the fitting function, but if the function represents a transformation that has data-dependent parameters (such as the standard R functions `poly` or `ns`), R will not in general be able to compute predicted values correctly for new observations. For example, the function `ns` that automatically selects knots for a B-spline fit will not be conducive to obtaining predicted values if the knots are kept “secret.” For this reason, a set of functions that keep track of transformation parameters, exists in `rms` for use with the functions highlighted in this book. These are shown in Table 6.3. Of these functions, `asis`, `catg`, `scored`, and `matrx` are almost always called implicitly and are not mentioned by the user. `catg` is usually called explicitly when the variable is a numeric variable to be used as a polytomous factor, and it has not been converted to an R categorical variable using the `factor` function.

Table 6.3 rms Transformation Functions

Function	Purpose	Related R Functions
<code>asis</code>	No post-transformation (seldom used explicitly)	<code>I</code>
<code>rcs</code>	Restricted cubic spline	<code>ns</code>
<code>pol</code>	Polynomial using standard notation	<code>poly</code>
<code>lsp</code>	Linear spline	
<code>catg</code>	Categorical predictor (seldom)	<code>factor</code>
<code>scored</code>	Ordinal categorical variables	<code>ordered</code>
<code>matrx</code>	Keep variables as group for <code>anova</code> and <code>fastbw</code>	<code>matrix</code>
<code>strat</code>	Nonmodeled stratification factors (used for <code>cph</code> only)	<code>strata</code>

These functions can be used with any function of a predictor. For example, to obtain a four-knot cubic spline expansion of the cube root of x , specify `rcs(x^(1/3),4)`.

When the transformation functions are called, they are usually given one or two arguments, such as `rcs(x,5)`. The first argument is the predictor variable or some function of it. The second argument is an optional vector of parameters describing a transformation, for example location or number of knots. Other arguments may be provided.

The `Hmisc` package's `cut2` function is sometimes used to create a categorical variable from a continuous variable x . You can specify the actual interval endpoints (`cuts`), the number of observations to have in each interval on the average (`m`), or the number of quantile groups (`g`). Use, for example, `cuts=c(0,1,2)` to cut into the intervals $[0, 1)$, $[1, 2]$.

A key concept in fitting models in R is that the fitting function returns an object that is an R list. This object contains basic information about the fit (e.g., regression coefficient estimates and covariance matrix, model χ^2) as well as information about how each parameter of the model relates to each factor in the model. Components of the fit object are addressed by, for example, `fit$coef`, `fit$var`, `fit$loglik`. `rms` causes the following information to also be retained in the fit object: the limits for plotting and estimating effects for each factor (if `options(datadist="name")` was in effect), the label for each factor, and a vector of values indicating which parameters associated with a factor are nonlinear (if any). Thus the “fit object” contains all the information needed to get predicted values, plots, odds or hazard ratios, and hypothesis tests, and to do “smart” variable selection that keeps parameters together when they are all associated with the same predictor.

R uses the notion of the *class* of an object. The object-oriented class idea allows one to write a few generic functions that decide which specific functions to call based on the class of the object passed to the generic function. An example is the function for printing the main results of a logistic model.

The `lrm` function returns a fit object of class "lrm". If you specify the R command `print(fit)` (or just `fit` if using R interactively—this invokes `print`), the `print` function invokes the `print.lrm` function to do the actual printing specific to logistic models. To find out which particular methods are implemented for a given generic function, type `methods(generic.name)`.

Generic functions that are used in this book include those in Table 6.4.

Table 6.4 rms Package and R Generic Functions

Function	Purpose	Related Functions
<code>print</code>	Print parameters and statistics of fit	
<code>coef</code>	Fitted regression coefficients	
<code>formula</code>	Formula used in the fit	
<code>specs</code>	Detailed specifications of fit	
<code>vcov</code>	Fetch covariance matrix	
<code>logLik</code>	Fetch maximized log-likelihood	
<code>AIC</code>	Fetch AIC	
<code>lrtest</code>	Likelihood ratio test for two nested models	
<code>univarLR</code>	Compute all univariable LR χ^2	
<code>robcov</code>	Robust covariance matrix estimates	
<code>bootcov</code>	Bootstrap covariance matrix estimates and bootstrap distributions of estimates	
<code>pentrace</code>	Find optimum penalty factors by tracing effective AIC for a grid of penalties	
<code>effective.df</code>	Print effective d.f. for each type of variable in model, for penalized fit or <code>pentrace</code> result	
<code>summary</code>	Summary of effects of predictors	
<code>plot.summary</code>	Plot continuously shaded confidence bars for results of <code>summary</code>	
<code>anova</code>	Wald tests of most meaningful hypotheses	
<code>plot.anova</code>	Graphical depiction of <code>anova</code>	
<code>contrast</code>	General contrasts, C.L., tests	
<code>Predict</code>	Predicted values and confidence limits easily varying a subset of predictors and leaving the rest set at default values	
<code>plot.Predict</code>	Plot the result of <code>Predict</code> using <code>lattice</code>	
<code>ggplot</code>	Plot the result of <code>Predict</code> using <code>ggplot2</code>	
<code>bplot</code>	3-dimensional plot when <code>Predict</code> varied two continuous predictors over a fine grid	
<code>gendata</code>	Easily generate predictor combinations	
<code>predict</code>	Obtain predicted values or design matrix	
<code>fastbw</code>	Fast backward step-down variable selection	<code>step</code>
<code>residuals</code>	(or <code>resid</code>) Residuals, influence stats from fit	
<code>sensuc</code>	Sensitivity analysis for unmeasured confounder	
<code>which.influence</code>	Which observations are overly influential	<code>residuals</code>
<code>latex</code>	L ^A T _E X representation of fitted model	<code>Function</code>

continued on next page

<i>continued from previous page</i>		
Function	Purpose	Related Functions
Function	R function analytic representation of $X\hat{\beta}$ from a fitted regression model	latex
Hazard	R function analytic representation of a fitted hazard function (for psm)	
Survival	R function analytic representation of fitted survival function (for psm , cph)	
ExProb	R function analytic representation of exceedance probabilities for orm	
Quantile	R function analytic representation of fitted function for quantiles of survival time (for psm , cph)	
Mean	R function analytic representation of fitted function for mean survival time or for ordinal logistic	
nomogram	Draws a nomogram for the fitted model	latex , plot
survest	Estimate survival probabilities (psm , cph)	survfit
survplot	Plot survival curves (psm , cph)	plot.survfit
validate	Validate indexes of model fit using resampling	
calibrate	Estimate calibration curve using resampling	val.prob
vif	Variance inflation factors for fitted model	
naresid	Bring elements corresponding to missing data back into predictions and residuals	
naprint	Print summary of missing values	
impute	Impute missing values	transcan

The first argument of the majority of functions is the object returned from the model fitting function. When used with **ols**, **lrm**, **orm**, **psm**, **cph**, **Glm**, **Gls**, **Rq**, **bj**, these functions do the following. **specs** prints the design specifications, for example, number of parameters for each factor, levels of categorical factors, knot locations in splines, and so on. **vcov** returns the variance-covariance matrix for the model. **logLik** retrieves the maximized log-likelihood, whereas **AIC** computes the Akaike Information Criterion for the model on the minus twice log-likelihood scale (with an option to compute it on the χ^2 scale if you specify **type='chisq'**). **lrtest**, when given two fit objects from nested models, computes the likelihood ratio test for the extra variables. **univarLR** computes all univariable likelihood ratio χ^2 statistics, one predictor at a time.

The **robcov** function computes the Huber robust covariance matrix estimate. **bootcov** uses the bootstrap to estimate the covariance matrix of parameter estimates. Both **robcov** and **bootcov** assume that the design matrix and response variable were stored with the fit. They have options to adjust for cluster sampling. Both replace the original variance-covariance matrix with robust estimates and return a new fit object that can be passed to any of the other functions. In that way, robust Wald tests, variable selection, confidence limits, and many other quantities may be computed automatically. The functions do save the old covariance estimates in component **orig.var** of the new fit object. **bootcov** also optionally returns the matrix of parameter estimates over the bootstrap simulations. These estimates can be used to derive bootstrap confidence intervals that don't assume normality or symmetry. Associated with **bootcov** are plotting functions for drawing histogram

and smooth density estimates for bootstrap distributions. `bootcov` also has a feature for deriving approximate nonparametric simultaneous confidence sets. For example, the function can get a simultaneous 0.90 confidence region for the regression effect of age over its entire range.

The `pentrace` function assists in selection of penalty factors for fitting regression models using penalized maximum likelihood estimation (see Section 9.10). Different types of model terms can be penalized by different amounts. For example, one can penalize interaction terms more than main effects. The `effective.df` function prints details about the effective degrees of freedom devoted to each type of model term in a penalized fit.

`summary` prints a summary of the effects of each factor. When `summary` is used to estimate effects (e.g., odds or hazard ratios) for continuous variables, it allows the levels of interacting factors to be easily set, as well as allowing the user to choose the interval for the effect. This method of estimating effects allows for nonlinearity in the predictor. By default, interquartile range effects (differences in $X\hat{\beta}$, odds ratios, hazards ratios, etc.) are printed for continuous factors, and all comparisons with the reference level are made for categorical factors. See the example at the end of the `summary` documentation for a method of quickly computing pairwise treatment effects and confidence intervals for a large series of values of factors that interact with the treatment variable. Saying `plot(summary(fit))` will depict the effects graphically, with bars for a list of confidence levels.

The `anova` function automatically tests most meaningful hypotheses in a design. For example, suppose that age and cholesterol are predictors, and that a general interaction is modeled using a restricted spline surface. `anova` prints Wald statistics for testing linearity of age, linearity of cholesterol, age effect (age + age \times cholesterol interaction), cholesterol effect (cholesterol + age \times cholesterol interaction), linearity of the age \times cholesterol interaction (i.e., adequacy of the simple age \times cholesterol 1 d.f. product), linearity of the interaction in age alone, and linearity of the interaction in cholesterol alone. Joint tests of all interaction terms in the model and all nonlinear terms in the model are also performed. The `plot.anova` function draws a dot chart showing the relative contribution (χ^2 , χ^2 minus d.f., AIC, partial R^2 , P -value, etc.) of each factor in the model.

The `contrast` function is used to obtain general contrasts and corresponding confidence limits and test statistics. This is most useful for testing effects in the presence of interactions (e.g., type II and type III contrasts). See the help file for `contrast` for several examples of how to obtain joint tests of multiple contrasts (see Section 9.3.2) as well as double differences (interaction contrasts).

The `predict` function is used to obtain a variety of values or predicted values from either the data used to fit the model or a new dataset. The `Predict` function is easier to use for most purposes, and has a special `plot` method. The `gendata` function makes it easy to obtain a data frame containing predictor combinations for obtaining selected predicted values.

The `fastbw` function performs a slightly inefficient but numerically stable version of fast backward elimination on factors, using a method based on Lawless and Singhal.³⁸⁵ This method uses the fitted complete model and computes approximate Wald statistics by computing conditional (restricted) maximum likelihood estimates assuming multivariate normality of estimates. It can be used in simulations since it returns indexes of factors retained and dropped:

```
fit ← ols(y ~ x1*x2*x3)
# run, and print results:
fastbw(fit, optional_arguments)
# typically used in simulations:
z ← fastbw(fit, optional_args)
# least squares fit of reduced model:
lm.fit(X[,z$params.kept], Y)
```

`fastbw` deletes factors, not columns of the design matrix. Factors requiring multiple d.f. will be retained or dropped as a group. The function prints the deletion statistics for each variable in turn, and prints approximate parameter estimates for the model after deleting variables. The approximation is better when the number of factors deleted is not large. For `ols`, the approximation is exact.

The `which.influence` function creates a list with a component for each factor in the model. The names of the components are the factor names. Each component contains the observation identifiers of all observations that are “overly influential” with respect to that factor, meaning that $|\mathbf{dfbetas}| > u$ for at least one β_i associated with that factor, for a given u . The default u is `.2`. You must have specified `x=TRUE`, `y=TRUE` in the fitting function to use `which.influence`. The first argument is the fit object, and the second argument is the cutoff u .

The following R program will print the set of predictor values that were very influential for each factor. It assumes that the data frame containing the data used in the fit is called `df`.

```
f ← lrm(y ~ x1 + x2 + ..., data=df, x=TRUE, y=TRUE)
w ← which.influence(f, .4)
nam ← names(w)
for(i in 1:length(nam)) {
  cat("Influential observations for effect of",
      nam[i], "\n")
  print(df[w[[i]],])
}
```

The `latex` function is a generic function available in the `Hmisc` package. It invokes a specific `latex` function for most of the fit objects created by `rms` to create a L^AT_EX algebraic representation of the fitted model for inclusion in a report or viewing on the screen. This representation documents all parameters in the model and the functional form being assumed for Y , and is especially useful for getting a simplified version of restricted cubic spline functions. On

the other hand, the `print` method with optional argument `latex=TRUE` is used to output L^AT_EX code representing the model results in tabular form to the console. This is intended for use with `knitr`⁶⁷⁷ or `Sweave`³⁹⁹.

The `Function` function composes an R function that you can use to evaluate $X\hat{\beta}$ analytically from a fitted regression model. The documentation for `Function` also shows how to use a subsidiary function `sascode` that will (almost) translate such an R function into SAS code for evaluating predicted values in new subjects. Neither `Function` nor `latex` handles third-order interactions.

The `nomogram` function draws a partial nomogram for obtaining predictions from the fitted model manually. It constructs different scales when interactions (up to third-order) are present. The constructed nomogram is not complete, in that point scores are obtained for each predictor and the user must add the point scores manually before reading predicted values on the final axis of the nomogram. The constructed nomogram is useful for interpreting the model fit, especially for non-monotonically transformed predictors (their scales wrap around an axis automatically).

The `vif` function computes variance inflation factors from the covariance matrix of a fitted model, using [147, 654].

The `impute` function is another generic function. It does simple imputation by default. It can also work with the `transcan` function to multiply or singly impute missing values using a flexible additive model.

As an example of using many of the functions, suppose that a categorical variable `treat` has values "a", "b", and "c", an ordinal variable `num.diseases` has values 0,1,2,3,4, and that there are two continuous variables, `age` and `cholesterol`. `age` is fitted with a restricted cubic spline, while `cholesterol` is transformed using the transformation `log(cholesterol+10)`. Cholesterol is missing on three subjects, and we impute these using the overall median cholesterol. We wish to allow for interaction between `treat` and `cholesterol`. The following R program will fit a logistic model, test all effects in the design, estimate effects, and plot estimated transformations. The fit for `num.diseases` really considers the variable to be a five-level categorical variable. The only difference is that a 3 d.f. test of linearity is done to assess whether the variable can be remodeled "asis". Here we also show statements to attach the `rms` package and store predictor characteristics from `datadist`.

```
require(rms) # make new functions available
ddist <- datadist(cholesterol, treat, num.diseases, age)
# Could have used ddist <- datadist(data.frame.name)
options(datadist="ddist") # defines data dist. to rms
cholesterol <- impute(cholesterol)
fit <- lrm(y ~ treat + scored(num.diseases) + rcs(age) +
          log(cholesterol+10) +
          treat:log(cholesterol+10))
describe(y ~ treat + scored(num.diseases) + rcs(age))
# or use describe(formula(fit)) for all variables used in
# fit. describe function (in Hmisc) gets simple statistics
# on variables
# fit <- robcov(fit)# Would make all statistics that follow
```

```

# use a robust covariance matrix
# would need x=TRUE, y=TRUE in lrm()
# Describe the design characteristics
specs(fit)
anova(fit)
anova(fit, treat, cholesterol) # Test these 2 by themselves
plot(anova(fit)) # Summarize anova graphically
summary(fit) # Est. effects; default ranges
plot(summary(fit)) # Graphical display of effects with C.I.
# Specific reference cell and adjustment value:
summary(fit, treat="b", age=60)
# Estimate effect of increasing age: 50->70
summary(fit, age=c(50,70))
# Increase age 50->70, adjust to 60 when estimating
# effects of other factors:
summary(fit, age=c(50,60,70))
# If had not defined datadist, would have to define
# ranges for all variables

# Estimate and test treatment (b-a) effect averaged
# over 3 cholesterol:
contrast(fit, list(treat='b', cholesterol=c(150,200,250)),
         list(treat='a', cholesterol=c(150,200,250)),
         type='average')

p ← Predict(fit, age=seq(20,80,length=100), treat,
            conf.int=FALSE)
plot(p) # Plot relationship between age and
# or ggplot(p) # log odds, separate curve for each
# treat, no C.I.
plot(p, ~ age | treat) # Same but 2 panels
ggplot(p, groups=FALSE)
bplot(Predict(fit, age, cholesterol, np=50))
# 3-dimensional perspective plot for
# age, cholesterol, and log odds
# using default ranges for both
# Plot estimated probabilities instead of log odds:
plot(Predict(fit, num.diseases,
            fun=function(x) 1/(1+exp(-x)),
            conf.int=.9), ylab="Prob")
# Again, if no datadist were defined, would have to tell
# plot all limits
logit ← predict(fit, expand.grid(treat="b", num.dis=1:3,
                                age=c(20,40,60),
                                cholesterol=seq(100,300,length=10)))
# Could obtain list of predictor settings interactively
logit ← predict(fit, gendata(fit, nobs=12))
# An easier approach is
# Predict(fit, treat='b', num.dis=1:3,...)

# Since age doesn't interact with anything, we can quickly
# and interactively try various transformations of age,
# taking the spline function of age as the gold standard.
# We are seeking a linearizing transformation.

```

```

ag ← 10:80
logit ← predict(fit, expand.grid(treat="a", num.dis=0,
                               age=ag,
                               cholesterol=median(cholesterol)),
               type="terms")[,"age"]
# Note: if age interacted with anything, this would be the
# age `main effect' ignoring interaction terms
# Could also use logit ← Predict(f, age=ag, ...)$yhat,
# which allows evaluation of the shape for any level of
# interacting factors. When age does not interact with
# anything, the result from predict(f, ..., type="terms")
# would equal the result from Predict if all other terms
# were ignored

# Could also specify:
# logit ← predict(fit,
#                gendata(fit, age=ag, cholesterol=...))
# Unmentioned variables are set to reference values

plot(ag^.5, logit) # try square root vs. spline transform.
plot(ag^1.5, logit) # try 1.5 power

# Pretty printing of table of estimates and
# summary statistics:
print(fit, latex=TRUE) # print  $\LaTeX$  code to console
latex(fit) # invokes latex.lrm, creates fit.tex
# Draw a nomogram for the model fit
plot(nomogram(fit))

# Compose R function to evaluate linear predictors
# analytically
g ← Function(fit)
g(treat='b', cholesterol=260, age=50)
# Letting num.diseases default to reference value

```

To examine interactions in a simpler way, you may want to group age into tertiles:

```

age.tertile ← cut2(age, g=3)
# For auto ranges later, specify age.tertile to datadist
fit ← lrm(y ~ age.tertile * rcs(cholesterol))

```

Example output from these functions is shown in Chapter 10 and later chapters.

Note that `type="terms"` in `predict` scores each factor in a model with its fitted transformation. This may be used to compute, for example, rank correlation between the response and each transformed factor, pretending it has 1 d.f.

When regression is done on principal components, one may use an ordinary linear model to decode “internal” regression coefficients for helping to understand the final model. Here is an example.

```

require(rms)
dd ← datadist(my.data)
options(datadist='dd')
pcfit ← princomp(~ pain.symptom1 + pain.symptom2 + sign1 +
                 sign2 + sign3 + smoking)
pc2 ← pcfit$scores[,1:2] # first 2 PCs as matrix
logistic.fit ← lrm(death ~ rcs(age,4) + pc2)
predicted.logit ← predict(logistic.fit)
linear.mod ← ols(predicted.logit ~ rcs(age,4) +
                 pain.symptom1 + pain.symptom2 +
                 sign1 + sign2 + sign3 + smoking)
# This model will have R-squared=1
nom ← nomogram(linear.mod, fun=function(x)1/(1+exp(-x)),
               funlabel="Probability of Death")
# can use fun=plogis
plot(nom)
# 7 Axes showing effects of all predictors, plus a reading
# axis converting to predicted probability scale

```

In addition to many of the add-on functions described above, there are several other R functions that validate models. The first, `predab.resample`, is a general-purpose function that is used by functions for specific models described later. `predab.resample` computes estimates of optimism and bias-corrected estimates of a vector of indexes of predictive accuracy, for a model with a specified design matrix, with or without fast backward step-down of predictors. If `bw=TRUE`, `predab.resample` prints a matrix of asterisks showing which factors were selected at each repetition, along with a frequency distribution of the number of factors retained across resamples. The function has an optional parameter that may be specified to force the bootstrap algorithm to do sampling with replacement from clusters rather than from original records, which is useful when each subject has multiple records in the dataset. It also has a parameter that can be used to validate predictions in a subset of the records even though models are refit using all records.

The generic function `validate` invokes `predab.resample` with model-specific fits and measures of accuracy. The function `calibrate` invokes `predab.resample` to estimate bias-corrected model calibration and to plot the calibration curve. Model calibration is estimated at a sequence of predicted values.

6.4 Other Functions

For principal component analysis, R has the `princomp` and `prcomp` functions. Canonical correlations and canonical variates can be easily computed using the `cancor` function. There are many other R functions for examining associations and for fitting models. The `supsmu` function implements Friedman's "super smoother."²⁰⁷ The `lowess` function implements Cleveland's two-dimensional smoother.¹¹¹ The `glm` function will fit general linear models under

a wide variety of distributions of Y . There are functions to fit Hastie and Tibshirani's²⁷⁵ generalized additive model for a variety of distributions. More is said about parametric and nonparametric additive multiple regression functions in Chapter 16. The `loess` function fits a multidimensional scatterplot smoother (the local regression model of Cleveland et al.⁹⁶). `loess` provides approximate test statistics for normal or symmetrically distributed Y :

```
f <- loess(y ~ age * pressure)
plot(f) # cross-sectional plots
ages <- seq(20,70,length=40)
pressures <- seq(80,200,length=40)
pred <- predict(f,
               expand.grid(age=ages, pressure=pressures))
persp(ages, pressures, pred) # 3-D plot
```

`loess` has a large number of options allowing various restrictions to be placed on the fitted surface.

Atkinson and Therneau's `rpart` recursive partitioning package and related functions implement classification and regression trees⁶⁹ algorithms for binary, continuous, and right-censored response variables (assuming an exponential distribution for the latter). `rpart` deals effectively with missing predictor values using surrogate splits. The `rms` package has a `validate` function for `rpart` objects for obtaining cross-validated mean squared errors and Somers' D_{xy} rank correlations (Brier score and ROC areas for probability models).

For displaying which variables tend to be missing on the same subjects, the `Hmisc` `naclus` function can be used (e.g., `plot(naclus(dataframename))` or `naplot(naclus(dataframename))`). For characterizing what type of subjects have NA's on a given predictor (or response) variable, a tree model whose response variable is `is.na(varname)` can be quite useful.

```
require(rpart)
f <- rpart(is.na(cholesterol) ~ age + sex + trig + smoking)
plot(f) # plots the tree
text(f) # labels the tree
```

The `Hmisc` `rcorr.cens` function can compute Somers' D_{xy} rank correlation coefficient and its standard error, for binary or continuous (and possibly right-censored) responses. A simple transformation of D_{xy} yields the c index (generalized ROC area). The `Hmisc` `improveProb` function is useful for comparing two probability models using the methods of Pencina et al.^{490,492,493} in an external validation setting. See also the `rcorrr.cens` function in this context.

6.5 Further Reading

- 1 Harrell and Goldstein²⁶³ list components of statistical languages or packages and compare several popular packages for survival analysis capabilities.
- 2 Imai et al.³¹⁹ have further generalized R as a statistical modeling language.

Chapter 7

Modeling Longitudinal Responses using Generalized Least Squares

In this chapter we consider models for a multivariate response variable represented by serial measurements over time within subject. This setup induces correlations between measurements on the same subject that must be taken into account to have optimal model fits and honest inference. Full likelihood model-based approaches have advantages including (1) optimal handling of imbalanced data and (2) robustness to missing data (dropouts) that occur not completely at random. The three most popular model-based full likelihood approaches are mixed effects models, generalized least squares, and Bayesian hierarchical models. For continuous Y , generalized least squares has a certain elegance, and a case study will demonstrate its use after surveying competing approaches. As OLS is a special case of generalized least squares, the case study is also helpful in developing and interpreting OLS models^a.

Some good references on longitudinal data analysis include^{148, 159, 252, 414, 509, 635, 637}.

7.1 Notation and Data Setup

Suppose there are N independent subjects, with subject i ($i = 1, 2, \dots, N$) having n_i responses measured at times $t_{i1}, t_{i2}, \dots, t_{in_i}$. The response at time t for subject i is denoted by Y_{it} . Suppose that subject i has baseline covariates X_i . Generally the response measured at time $t_{i1} = 0$ is a covariate in X_i instead of being the first measured response Y_{i0} .

For flexible analysis, longitudinal data are usually arranged in a “tall and thin” layout. This allows measurement times to be irregular. In studies com-

^a A case study in OLS—Chapter 7 from the first edition—may be found on the text’s web site.

paring two or more treatments, a response is often measured at baseline (pre-randomization). The analyst has the option to use this measurement as Y_{i0} or as part of X_i . There are many reasons to put initial measurements of Y in X , i.e., to use baseline measurements as baseline .

1

7.2 Model Specification for Effects on $E(Y)$

Longitudinal data can be used to estimate overall means or the mean at the last scheduled follow-up, making maximum use of incomplete records. But the real value of longitudinal data comes from modeling the entire time course. Estimating the time course leads to understanding slopes, shapes, overall trajectories, and periods of treatment effectiveness. With continuous Y one typically specifies the time course by a mean time-response profile. Common representations for such profiles include

- k dummy variables for $k + 1$ unique times (assumes no functional form for time but assumes discrete measurement times and may spend many d.f.)
- $k = 1$ for linear time trend, $g_1(t) = t$
- k -order polynomial in t
- $k + 1$ -knot restricted cubic spline (one linear term, $k - 1$ nonlinear terms)

Suppose the time trend is modeled with k parameters so that the time effect has k d.f. Let the basis functions modeling the time effect be $g_1(t)$, $g_2(t), \dots, g_k(t)$ to allow it to be nonlinear. A model for the time profile without interactions between time and any X is given by

$$E[Y_{it}|X_i] = X_i\beta + \gamma_1g_1(t) + \gamma_2g_2(t) + \dots + \gamma_kg_k(t). \quad (7.1)$$

To allow the slope or shape of the time-response profile to depend on some of the X s we add product terms for desired interaction effects. For example, to allow the mean time trend for subjects in group 1 (reference group) to be arbitrarily different from the time trend for subjects in group 2, have a dummy variable for group 2, a time “main effect” curve with k d.f. and all k products of these time components with the dummy variable for group 2.

Once the right hand side of the model is formulated, predicted values, contrasts, and ANOVAs are obtained just as with a univariate model. For these purposes time is no different than any other covariate except for what is described in the next section.

7.3 Modeling Within-Subject Dependence

Sometimes understanding within-subject correlation patterns is of interest in itself. More commonly, accounting for intra-subject correlation is crucial for inferences to be valid. Some methods of analysis cover up the correlation

pattern while others assume a restrictive form for the pattern. The following table is an attempt to briefly survey available longitudinal analysis methods. LOCF and the summary statistic method are not modeling methods. LOCF is an ad hoc attempt to account for longitudinal dropouts, and summary statistics can convert multivariate responses to univariate ones with few assumptions (other than minimal dropouts), with some information loss.

2

What Methods To Use for Repeated Measurements / Serial Data? ^{ab}

	Repeated Measures ANOVA	GEE	Mixed Effects Model	GLS	LOCF	Summary Statistic ^c
Assumes normality	×		×	×		
Assumes independence of measurements within subject	× ^d	× ^e				
Assumes a correlation structure ^f	×	× ^g	×	×		
Requires same measurement times for all subjects	×				?	
Does not allow smooth modeling of time to save d.f.	×					
Does not allow adjustment for baseline covariates	×					
Does not easily extend to non-continuous Y	×			×		
Loses information by not using intermediate measurements					× ^h	×
Does not allow widely varying # of observations per subject	×	× ⁱ			×	× ^j
Does not allow for subjects to have distinct trajectories ^k	×	×		×	×	
Assumes subject-specific effects are Gaussian			×			
Badly biased if non-random dropouts	?	×			×	
Biased in general					×	
Harder to get tests & CLs			× ^l		× ^m	
Requires large # subjects/clusters		×				
SEs are wrong	× ⁿ				×	
Assumptions are not verifiable in small samples	×	N/A	×	×	×	
Does not extend to complex settings such as time-dependent covariates and dynamic ^o models	×		×	×	×	?

^a Thanks to Charles Berry, Brian Cade, Peter Flom, Bert Gunter, and Leena Choi for valuable input.

^b GEE: generalized estimating equations; GLS: generalized least squares; LOCF: last observation carried forward.

^c E.g., compute within-subject slope, mean, or area under the curve over time. Assumes that the summary measure is an adequate summary of the time profile and assesses the relevant treatment effect.

The most prevalent full modeling approach is mixed effects models in which baseline predictors are fixed effects, and random effects are used to describe subject differences and to induce within-subject correlation. Some disadvantages of mixed effects models are

- The induced correlation structure for Y may be unrealistic if care is not taken in specifying the model.
- Random effects require complex approximations for distributions of test statistics.
- The most commonly used models assume that random effects follow a normal distribution. This assumption may not hold.

It could be argued that an extended linear model (with no random effects) is a logical extension of the univariate OLS model^b. This model, called the generalized least squares or growth curve model^{221, 509, 510}, was developed long before mixed effect models became popular.

We will assume that $Y_{it}|X_i$ has a multivariate normal distribution with mean given above and with variance-covariance matrix V_i , an $n_i \times n_i$ matrix that is a function of t_{i1}, \dots, t_{in_i} . We further assume that the diagonals of V_i are all equal^b. This *extended linear model* has the following assumptions:

- all the assumptions of OLS at a single time point including correct modeling of predictor effects and univariate normality of responses conditional on X

^d Unless one uses the Huynh-Feldt or Greenhouse-Geisser correction

^e For full efficiency, if using the working independence model

^f Or requires the user to specify one

^g For full efficiency of regression coefficient estimates

^h Unless the last observation is missing

ⁱ The cluster sandwich variance estimator used to estimate SEs in GEE does not perform well in this situation, and neither does the working independence model because it does not weight subjects properly.

^j Unless one knows how to properly do a weighted analysis

^k Or uses population averages

^l Unlike GLS, does not use standard maximum likelihood methods yielding simple likelihood ratio χ^2 statistics. Requires high-dimensional integration to marginalize random effects, using complex approximations, and if using SAS, unintuitive d.f. for the various tests.

^m Because there is no correct formula for SE of effects; ordinary SEs are not penalized for imputation and are too small

ⁿ If correction not applied

^o E.g., a model with a predictor that is a lagged value of the response variable

^b E.g., few statisticians use subject random effects for univariate Y . Pinheiro and Bates [509, Section 5.1.2] state that “in some applications, one may wish to avoid incorporating random effects in the model to account for dependence among observations, choosing to use the within-group component A_i to directly model variance-covariance structure of the response.”

^b This procedure can be generalized to allow for heteroscedasticity over time or with respect to X , e.g., males may be allowed to have a different variance than females.

- the distribution of two responses at two different times for the same subject, conditional on X , is bivariate normal with a specified correlation coefficient
- the joint distribution of all n_i responses for the i^{th} subject is multivariate normal with the given correlation pattern (which implies the previous two distributional assumptions)
- responses from two different subjects are uncorrelated.

7.4 Parameter Estimation Procedure

Generalized least squares is like weighted least squares but uses a covariance matrix that is not diagonal. Each subject can have her own shape of V_i due to each subject being measured at a different set of times. This is a maximum likelihood procedure. Newton-Raphson or other trial-and-error methods are used for estimating parameters. For a small number of subjects, there are advantages in using REML (restricted maximum likelihood) instead of ordinary MLE [159, Section 5.3] [509, Chapter 5]²²¹ (especially to get a more unbiased estimate of the covariance matrix).

When imbalances of measurement times are not severe, OLS fitted ignoring subject identifiers may be efficient for estimating β . But OLS standard errors will be too small as they don't take intra-cluster correlation into account. This may be rectified by substituting a covariance matrix estimated using the Huber-White cluster sandwich estimator or from the cluster bootstrap. When imbalances are severe and intra-subject correlations are strong, OLS (or GEE using a working independence model) is not expected to be efficient because it gives equal weight to each observation; a subject contributing two distant observations receives $\frac{1}{5}$ the weight of a subject having 10 tightly-spaced observations.

7.5 Common Correlation Structures

We usually restrict ourselves to *isotropic* correlation structures which assume the correlation between responses within subject at two times depends only on a measure of the distance between the two times, not the individual times. We simplify further and assume it depends on $|t_1 - t_2|^c$. Assume that the correlation coefficient for Y_{it_1} vs. Y_{it_2} conditional on baseline covariates X_i for subject i is $h(|t_1 - t_2|, \rho)$, where ρ is a vector (usually a scalar) set of fundamental correlation parameters. Some commonly used structures when

^c We can speak interchangeably of correlations of residuals within subjects or correlations between responses measured at different times on the same subject, conditional on covariates X .

times are continuous and are not equally spaced [509, Section 5.3.3] are shown below, along with the correlation function names from the R `nlme` package.

Compound symmetry: $h = \rho$ if $t_1 \neq t_2$, 1 if $t_1 = t_2$ `nlme corCompSymm`
 (Essentially what two-way ANOVA assumes)
 Autoregressive-moving average lag 1: $h = \rho^{|t_1 - t_2|} = \rho^s$ `corCAR1`
 where $s = |t_1 - t_2|$
 Exponential: $h = \exp(-s/\rho)$ `corExp`
 Gaussian: $h = \exp[-(s/\rho)^2]$ `corGaus`
 Linear: $h = (1 - s/\rho)[s < \rho]$ `corLin`
 Rational quadratic: $h = 1 - (s/\rho)^2/[1 + (s/\rho)^2]$ `corRatio`
 Spherical: $h = [1 - 1.5(s/\rho) + 0.5(s/\rho)^3][s < \rho]$ `corSpher`
 Linear exponent AR(1): $h = \rho^{d_{min} + \delta \frac{s - d_{min}}{d_{max} - d_{min}}}$, 1 if $t_1 = t_2$ ⁵⁷²

The structures 3–7 use ρ as a scaling parameter, not as something restricted to be in $[0, 1]$

7.6 Checking Model Fit

The constant variance assumption may be checked using typical residual plots. The univariate normality assumption (but not multivariate normality) may be checked using typical Q-Q plots on residuals. For checking the correlation pattern, a *variogram* is a very helpful device based on estimating correlations of all possible pairs of residuals at different time points^d. Pairs of estimates obtained at the same absolute time difference s are pooled. The variogram is a plot with $y = 1 - \hat{h}(s, \rho)$ vs. s on the x -axis, and the theoretical variogram of the correlation model currently being assumed is superimposed.

7.7 Sample Size Considerations

Section 4.4 provided some guidance about sample sizes needed for OLS. A good way to think about sample size adequacy for generalized least squares is to determine the effective number of independent observations that a given configuration of repeated measurements has. For example, if the standard error of an estimate from three measurements on each of 20 subjects is the same as the standard error from 27 subjects measured once, we say that the 20×3 study has an effective sample size of 27, and we equate power from the univariate analysis on n subjects measured once to $\frac{20n}{27}$ subjects measured three times. Faes et al.¹⁸¹ have a nice approach to effective sample sizes with a variety of correlation patterns in longitudinal data. For an AR(1) correlation structure with n equally spaced measurement times on each of N subjects,

^d Variograms can be unstable.

with the correlation between two consecutive times being ρ , the effective sample size is $\frac{n-(n-2)\rho}{1+\rho}N$. Under compound symmetry, the effective size is $\frac{nN}{1+\rho(n-1)}$.

7.8 R Software

The nonlinear mixed effects model package `nlme` of Pinheiro & Bates in R provides many useful functions. For fitting linear models, fitting functions are `lme` for mixed effects models and `gls` for generalized least squares without random effects. The `rms` package has a front-end function `Gls` so that many features of `rms` can be used:

anova: all partial Wald tests, test of linearity, pooled tests
summary: effect estimates (differences in \hat{Y}) and confidence limits
Predict and plot: partial effect plots
nomogram: nomogram
Function: generate R function code for the fitted model
latex: L^AT_EX representation of the fitted model.

In addition, `Gls` has a cluster bootstrap option (hence you do not use `rms`'s `bootcov` for `Gls` fits). When `B` is provided to `Gls()`, bootstrapped regression coefficients and correlation estimates are saved, the former setting up for bootstrap percentile confidence limits^e. The `nlme` package has many graphics and fit-checking functions. Several functions will be demonstrated in the case study.

7.9 Case Study

Consider the dataset in Table 6.9 of Davis [148, pp. 161–163] from a multi-center, randomized controlled trial of botulinum toxin type B (BotB) in patients with cervical dystonia from nine U.S. sites. Patients were randomized to placebo ($N = 36$), 5000 units of BotB ($N = 36$), or 10,000 units of BotB ($N = 37$). The response variable is the total score on the Toronto Western Spasmodic Torticollis Rating Scale (TWSTRS), measuring severity, pain, and disability of cervical dystonia (high scores mean more impairment). TWSTRS is measured at baseline (week 0) and weeks 2, 4, 8, 12, 16 after treatment began. The dataset name on the dataset wiki page is `cdystonia`.

^e To access regular `gls` functions named `anova` (for likelihood ratio tests, AIC, etc.) or `summary` use `anova.gls` or `summary.gls`.

7.9.1 Graphical Exploration of Data

Graphics which follow display raw data as well as quartiles of TWSTRS by time, site, and treatment. A table shows the realized measurement schedule.

```
require(rms)
```

```
getHdata(cdystonia)
attach(cdystonia)

# Construct unique subject ID
uid ← with(cdystonia, factor(paste(site, id)))

# Tabulate patterns of subjects' time points
table(tapply(week, uid,
             function(w) paste(sort(unique(w)), collapse=' ')))
```

	0	0 2 4	0 2 4 12 16	0 2 4 8	0 2 4 8 12
	1	1	3	1	1
0 2 4 8 12 16	0 2 4 8 16	0 2 8 12 16	0 4 8 12 16	0 4 8 16	
94	1	2	4	1	

```
# Plot raw data, superposing subjects
xl ← xlab('Week'); yl ← ylab('TWSTRS-total score')
ggplot(cdystonia, aes(x=week, y=twstrs, color=factor(id))) +
  geom_line() + xl + yl + facet_grid(treat ~ site) +
  guides(color=FALSE) # Fig. 7.1
```

```
# Show quartiles
ggplot(cdystonia, aes(x=week, y=twstrs)) + xl + yl +
  ylim(0, 70) + stat_summary(fun.data="median_hilow",
                             conf.int=0.5, geom='smooth') +
  facet_wrap(~ treat, nrow=2) # Fig. 7.2
```

Next the data are rearranged so that Y_{i0} is a baseline covariate.

```
baseline ← subset(data.frame(cdystonia, uid), week == 0,
                  -week)
baseline ← upData(baseline, rename=c(twstrs='twstrs0'),
                  print=FALSE)
followup ← subset(data.frame(cdystonia, uid), week > 0,
                  c(uid, week, twstrs))
rm(uid)
both ← merge(baseline, followup, by='uid')

dd ← datadist(both)
options(datadist='dd')
```

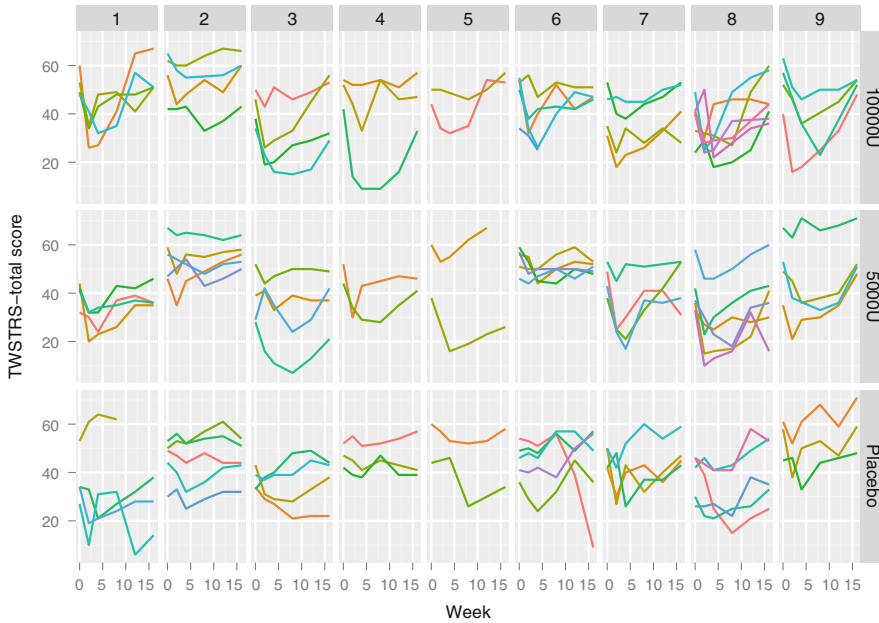


Fig. 7.1 Time profiles for individual subjects, stratified by study site and dose

7.9.2 Using Generalized Least Squares

We stay with baseline adjustment and use a variety of correlation structures, with constant variance. Time is modeled as a restricted cubic spline with 3 knots, because there are only 3 unique interior values of `week`. Below, six correlation patterns are attempted. In general it is better to use scientific knowledge to guide the choice of the correlation structure.

```
require(nlme)
```

```
cp<- list(corCAR1,corExp,corCompSymm,corLin,corGaus,corSpher)
z  <- vector('list',length(cp))
for(k in 1:length(cp)) {
  z[[k]] <- gls(twstrs ~ treat * rcs(week, 3) +
               rcs(twstrs0, 3) + rcs(age, 4) * sex, data=both,
               correlation=cp[[k]](form = ~week | uid))
}
```

```
anova(z[[1]],z[[2]],z[[3]],z[[4]],z[[5]],z[[6]])
```

	Model	df	AIC	BIC	logLik
z[[1]]	1	20	3553.906	3638.357	-1756.953
z[[2]]	2	20	3553.906	3638.357	-1756.953
z[[3]]	3	20	3587.974	3672.426	-1773.987
z[[4]]	4	20	3575.079	3659.531	-1767.540
z[[5]]	5	20	3621.081	3705.532	-1790.540
z[[6]]	6	20	3570.958	3655.409	-1765.479

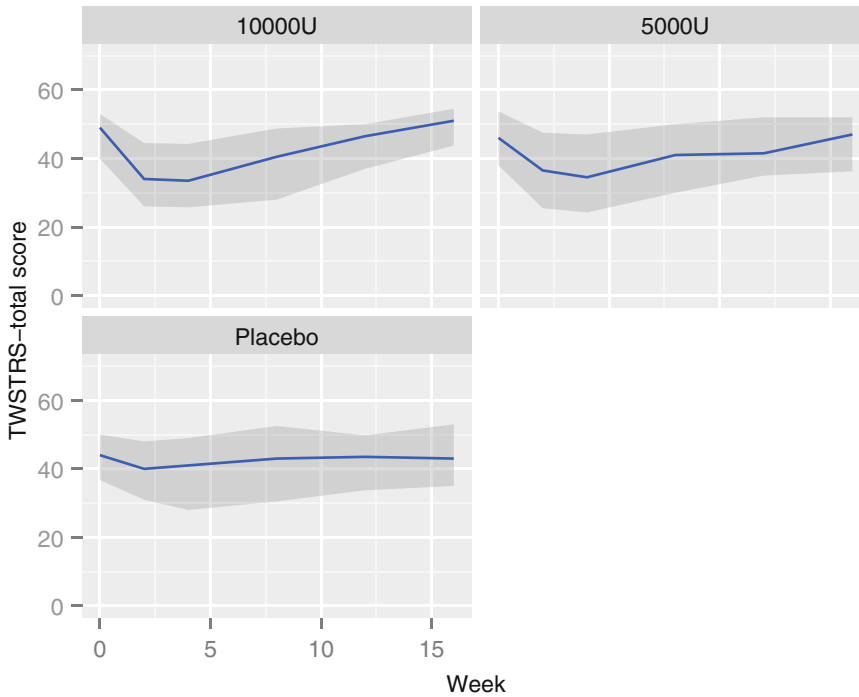


Fig. 7.2 Quartiles of TWSTRS stratified by dose

AIC computed above is set up so that smaller values are best. From this the continuous-time AR1 and exponential structures are tied for the best. For the remainder of the analysis we use `corCAR1`, using `Gls`.

3

```
a <- GlS(twstrs ~ treat * rcs(week, 3) + rcs(twstrs0, 3) +
         rcs(age, 4) * sex, data=both,
         correlation=corCAR1(form=~week | uid))
```

```
print(a, latex=TRUE)
```

Generalized Least Squares Fit by REML

```
Gls(model = twstrs ~ treat * rcs(week, 3) + rcs(twstrs0, 3) +
     rcs(age, 4) * sex, data = both, correlation = corCAR1
     (form = ~week | uid))
```

Obs	522	Log-restricted-likelihood	-1756.95
Clusters	108	Model d.f.	17
g	11.334	σ	8.5917
		d.f.	504

	Coef	S.E.	t	Pr(> t)
Intercept	-0.3093	11.8804	-0.03	0.9792
treat=5000U	0.4344	2.5962	0.17	0.8672
treat=Placebo	7.1433	2.6133	2.73	0.0065
week	0.2879	0.2973	0.97	0.3334
week'	0.7313	0.3078	2.38	0.0179
twstrs0	0.8071	0.1449	5.57	< 0.0001
twstrs0'	0.2129	0.1795	1.19	0.2360
age	-0.1178	0.2346	-0.50	0.6158
age'	0.6968	0.6484	1.07	0.2830
age''	-3.4018	2.5599	-1.33	0.1845
sex=M	24.2802	18.6208	1.30	0.1929
treat=5000U * week	0.0745	0.4221	0.18	0.8599
treat=Placebo * week	-0.1256	0.4243	-0.30	0.7674
treat=5000U * week'	-0.4389	0.4363	-1.01	0.3149
treat=Placebo * week'	-0.6459	0.4381	-1.47	0.1411
age * sex=M	-0.5846	0.4447	-1.31	0.1892
age' * sex=M	1.4652	1.2388	1.18	0.2375
age'' * sex=M	-4.0338	4.8123	-0.84	0.4023

Correlation Structure: Continuous AR(1)

Formula: ~week | uid

Parameter estimate(s):

Phi

0.8666689

$\hat{\rho} = 0.867$, the estimate of the correlation between two measurements taken one week apart on the same subject. The estimated correlation for measurements 10 weeks apart is $0.867^{10} = 0.24$.

```
v <- Variogram(a, form=~ week | uid)
plot(v) # Figure 7.3
```

The empirical variogram is largely in agreement with the pattern dictated by AR(1).

Next check constant variance and normality assumptions.

```
both$resid <- r <- resid(a); both$fitted <- fitted(a)
y1 <- ylab('Residuals')
p1 <- ggplot(both, aes(x=fitted, y=resid)) + geom_point() +
  facet_grid(~ treat) + y1
p2 <- ggplot(both, aes(x=twstrs0, y=resid)) + geom_point()+y1
p3 <- ggplot(both, aes(x=week, y=resid)) + y1 + ylim(-20,20) +
  stat_summary(fun.data="mean_sd1", geom='smooth')
p4 <- ggplot(both, aes(sample=resid)) + stat_qq() +
  geom_abline(intercept=mean(r), slope=sd(r)) + y1
gridExtra::grid.arrange(p1, p2, p3, p4, ncol=2) # Figure 7.4
```

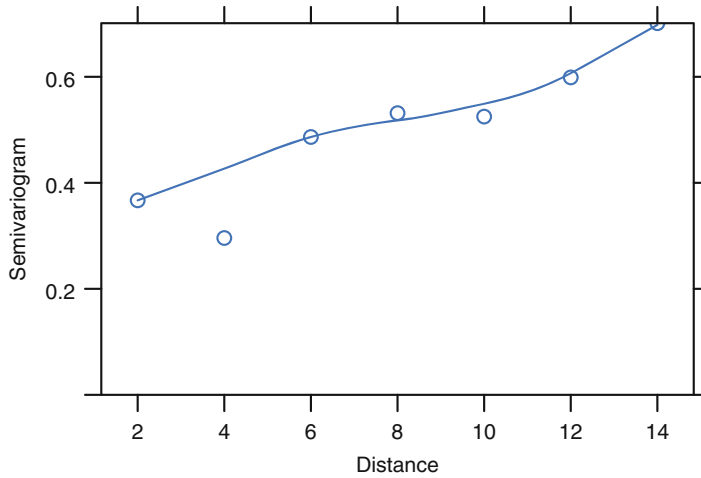


Fig. 7.3 Variogram, with assumed correlation pattern superimposed

These model assumptions appear to be well satisfied, so inferences are likely to be trustworthy if the more subtle multivariate assumptions hold.

Now get hypothesis tests, estimates, and graphically interpret the model.

```
plot(anova(a)) # Figure 7.5
```

```
ylm <- ylim(25, 60)
p1 <- ggplot(Predict(a, week, treat, conf.int=FALSE),
             adj.subtitle=FALSE, legend.position='top') + ylm
p2 <- ggplot(Predict(a, twstrs0), adj.subtitle=FALSE) + ylm
p3 <- ggplot(Predict(a, age, sex), adj.subtitle=FALSE,
             legend.position='top') + ylm
gridExtra::grid.arrange(p1, p2, p3, ncol=2) # Figure 7.6
```

```
latex(summary(a), file='', table.env=FALSE) # Shows for week 8
```

	Low	High	Δ	Effect	S.E.	Lower 0.95	Upper 0.95
week	4	12	8	6.69100	1.10570	4.5238	8.8582
twstrs0	39	53	14	13.55100	0.88618	11.8140	15.2880
age	46	65	19	2.50270	2.05140	-1.5179	6.5234
treat — 5000U:10000U	1	2		0.59167	1.99830	-3.3249	4.5083
treat — Placebo:10000U	1	3		5.49300	2.00430	1.5647	9.4212
sex — M:F	1	2		-1.08500	1.77860	-4.5711	2.4011

```
# To get results for week 8 for a different reference group
# for treatment, use e.g. summary(a, week=4, treat='Placebo')
# Compare low dose with placebo, separately at each time
```

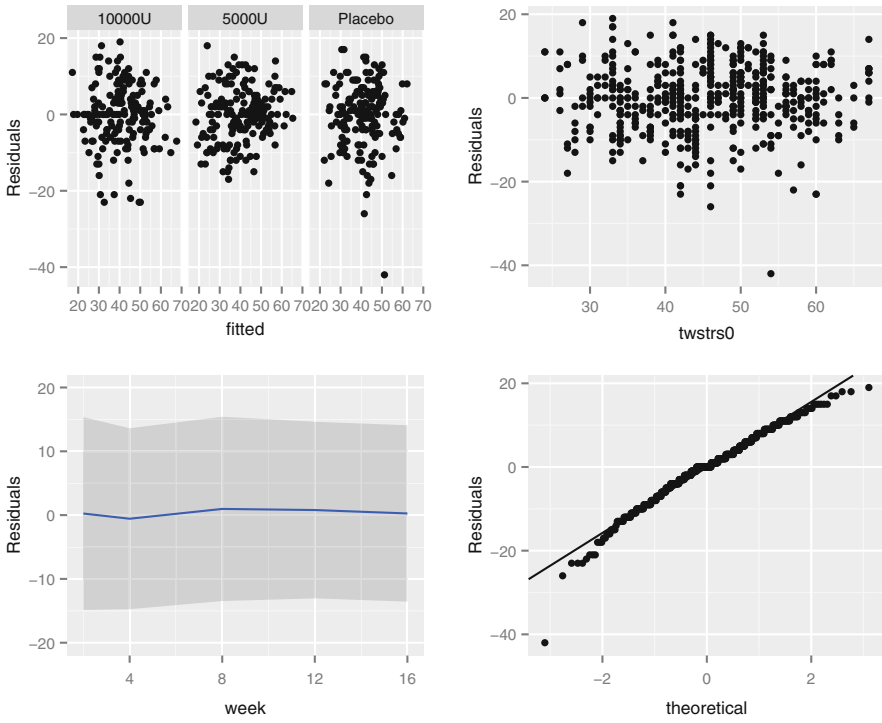


Fig. 7.4 Three residual plots to check for absence of trends in central tendency and in variability. Upper right panel shows the baseline score on the x -axis. Bottom left panel shows the mean $\pm 2 \times SD$. Bottom right panel is the QQ plot for checking normality of residuals from the GLS fit.

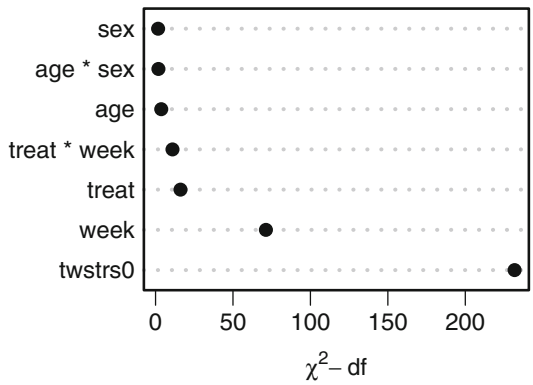


Fig. 7.5 Results of `anova` from generalized least squares fit with continuous time AR1 correlation structure. As expected, the baseline version of Y dominates.

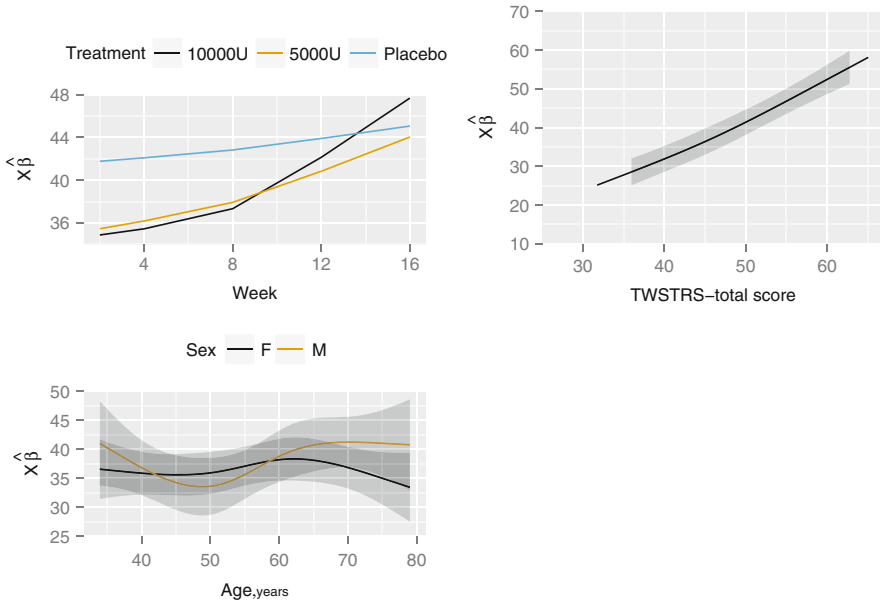


Fig. 7.6 Estimated effects of time, baseline TWSTRS, age, and sex

```
k1 ← contrast(a, list(week=c(2,4,8,12,16), treat='5000U'),
              list(week=c(2,4,8,12,16), treat='Placebo'))
options(width=80)
print(k1, digits=3)
```

	week	twstrs0	age	sex	Contrast	S.E.	Lower	Upper	Z	Pr(> z)
1	2	46	56	F	-6.31	2.10	-10.43	-2.186	-3.00	0.0027
2	4	46	56	F	-5.91	1.82	-9.47	-2.349	-3.25	0.0011
3	8	46	56	F	-4.90	2.01	-8.85	-0.953	-2.43	0.0150
4*	12	46	56	F	-3.07	1.75	-6.49	0.361	-1.75	0.0795
5*	16	46	56	F	-1.02	2.10	-5.14	3.092	-0.49	0.6260

Redundant contrasts are denoted by *

Confidence intervals are 0.95 individual intervals

```
# Compare high dose with placebo
k2 ← contrast(a, list(week=c(2,4,8,12,16), treat='10000U'),
              list(week=c(2,4,8,12,16), treat='Placebo'))
print(k2, digits=3)
```

	week	twstrs0	age	sex	Contrast	S.E.	Lower	Upper	Z	Pr(> z)
1	2	46	56	F	-6.89	2.07	-10.96	-2.83	-3.32	0.0009
2	4	46	56	F	-6.64	1.79	-10.15	-3.13	-3.70	0.0002
3	8	46	56	F	-5.49	2.00	-9.42	-1.56	-2.74	0.0061
4*	12	46	56	F	-1.76	1.74	-5.17	1.65	-1.01	0.3109
5*	16	46	56	F	2.62	2.09	-1.47	6.71	1.25	0.2099

Redundant contrasts are denoted by *

Confidence intervals are 0.95 individual intervals

```

k1 ← as.data.frame(k1[c('week', 'Contrast', 'Lower',
                        'Upper')])
p1 ← ggplot(k1, aes(x=week, y=Contrast)) + geom_point() +
      geom_line() + ylab('Low Dose - Placebo') +
      geom_errorbar(aes(ymin=Lower, ymax=Upper), width=0)
k2 ← as.data.frame(k2[c('week', 'Contrast', 'Lower',
                        'Upper')])
p2 ← ggplot(k2, aes(x=week, y=Contrast)) + geom_point() +
      geom_line() + ylab('High Dose - Placebo') +
      geom_errorbar(aes(ymin=Lower, ymax=Upper), width=0)
gridExtra::grid.arrange(p1, p2, ncol=2) # Figure 7.7

```

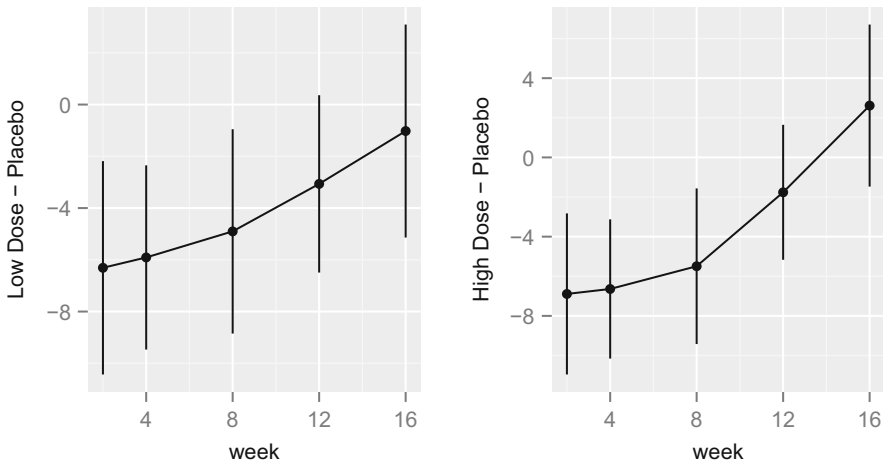


Fig. 7.7 Contrasts and 0.95 confidence limits from GLS fit

Although multiple d.f. tests such as total treatment effects or treatment \times time interaction tests are comprehensive, their increased degrees of freedom can dilute power. In a treatment comparison, treatment contrasts at the last time point (single d.f. tests) are often of major interest. Such contrasts are informed by all the measurements made by all subjects (up until dropout times) when a smooth time trend is assumed. They use appropriate extrapolation past dropout times based on observed trajectories of subjects followed the entire observation period. In agreement with the top left panel of Figure 7.6, Figure 7.7 shows that the treatment, despite causing an early improvement, wears off by 16 weeks at which time no benefit is seen.

A nomogram can be used to obtain predicted values, as well as to better understand the model, just as with a univariate Y .

```

n ← nomogram(a, age=c(seq(20, 80, by=10), 85))
plot(n, cex.axis=.55, cex.var=.8, lmgp=.25) # Figure 7.8

```

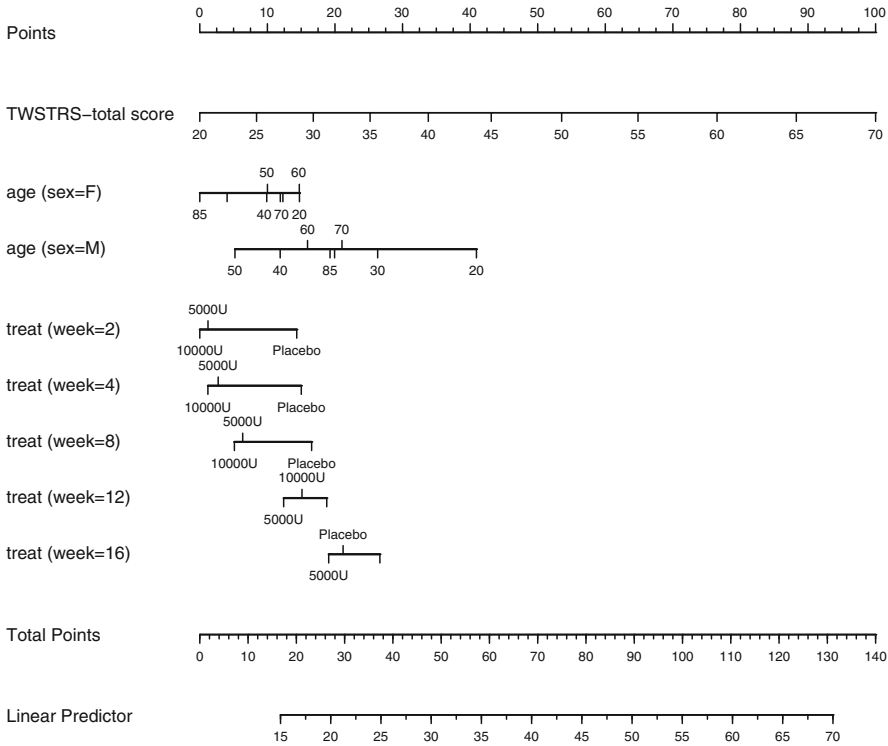


Fig. 7.8 Nomogram from GLS fit. Second axis is the baseline score.

7.10 Further Reading

- [1] Jim Rochon (Rho, Inc., Chapel Hill NC) has the following comments about using the baseline measurement of Y as the first longitudinal response.

For RCTs [randomized clinical trials], I draw a sharp line at the point when the intervention begins. The LHS [left hand side of the model equation] is reserved for something that is a response to treatment. Anything before this point can potentially be included as a covariate in the regression model. This includes the “baseline” value of the outcome variable. Indeed, the best predictor of the outcome at the end of the study is typically where the patient began at the beginning. It drinks up a lot of variability in the outcome; and, the effect of other covariates is typically mediated through this variable.

I treat anything after the intervention begins as an outcome. In the western scientific method, an “effect” must follow the “cause” even if by a split second.

Note that an RCT is different than a cohort study. In a cohort study, “Time 0” is not terribly meaningful. If we want to model, say, the trend over time, it would be legitimate, in my view, to include the “baseline” value on the LHS of that regression model.

Now, even if the intervention, e.g., surgery, has an immediate effect, I would include still reserve the LHS for anything that might legitimately be considered as the response to the intervention. So, if we cleared a blocked artery and then measured the MABP, then that would still be included on the LHS.

Now, it could well be that most of the therapeutic effect occurred by the time that the first repeated measure was taken, and then levels off. Then, a plot of the means would essentially be two parallel lines and the treatment effect is the distance between the lines, i.e., the difference in the intercepts.

If the linear trend from baseline to Time 1 continues beyond Time 1, then the lines will have a common intercept but the slopes will diverge. Then, the treatment effect will be the difference in slopes.

One point to remember is that the estimated intercept is the value at time 0 that we predict from the set of repeated measures post randomization. In the first case above, the model will predict different intercepts even though randomization would suggest that they would start from the same place. This is because we were asleep at the switch and didn't record the "action" from baseline to time 1. In the second case, the model will predict the same intercept values because the linear trend from baseline to time 1 was continued thereafter.

More importantly, there are considerable benefits to including it as a covariate on the RHS. The baseline value tends to be the best predictor of the outcome post-randomization, and this maneuver increases the precision of the estimated treatment effect. Additionally, any other prognostic factors correlated with the outcome variable will also be correlated with the baseline value of that outcome, and this has two important consequences. First, this greatly reduces the need to enter a large number of prognostic factors as covariates in the linear models. Their effect is already mediated through the baseline value of the outcome variable. Secondly, any imbalances across the treatment arms in important prognostic factors will induce an imbalance across the treatment arms in the baseline value of the outcome. Including the baseline value thereby reduces the need to enter these variables as covariates in the linear models.

Stephen Senn⁵⁶³ states that temporally and logically, a "baseline cannot be a *response* to treatment", so baseline and response cannot be modeled in an integrated framework.

... one should focus clearly on 'outcomes' as being the only values that can be influenced by treatment and examine critically any schemes that assume that these are linked in some rigid and deterministic view to 'baseline' values. An alternative tradition sees a baseline as being merely one of a number of measurements capable of improving predictions of outcomes and models it in this way.

The final reason that baseline cannot be modeled as the response at time zero is that many studies have inclusion/exclusion criteria that include cutoffs on the baseline variable yielding a truncated distribution. In general it is not appropriate to model the baseline with the same distributional shape as the follow-up

measurements. Thus the approach recommended by Liang and Zeger⁴⁰⁵ and Liu et al.⁴²³ are problematic^f.

- [2] Gardiner et al.²¹¹ compared several longitudinal data models, especially with regard to assumptions and how regression coefficients are estimated. Peters et al.⁵⁰⁰ have an empirical study confirming that the “use all available data” approach of likelihood-based longitudinal models makes imputation of follow-up measurements unnecessary.
- [3] Keselman et al.³⁴⁷ did a simulation study to study the reliability of AIC for selecting the correct covariance structure in repeated measurement models. In choosing from among 11 structures, AIC selected the correct structure 47% of the time. Gurka et al.²⁴⁷ demonstrated that fixed effects in a mixed effects model can be biased, independent of sample size, when the specified covariate matrix is more restricted than the true one.

^f In addition to this, one of the paper’s conclusions that analysis of covariance is not appropriate if the population means of the baseline variable are not identical in the treatment groups is arguable⁵⁶³. See³⁴⁶ for a discussion of^{f423}.

Chapter 8

Case Study in Data Reduction

Recall that the aim of data reduction is to reduce (without using the outcome) the number of parameters needed in the outcome model. The following case study illustrates these techniques:

1. redundancy analysis;
2. variable clustering;
3. data reduction using principal component analysis (PCA), sparse PCA, and pretransformations;
4. restricted cubic spline fitting using ordinary least squares, in the context of scaling; and
5. scaling/variable transformations using canonical variates and nonparametric additive regression.

8.1 Data

Consider the 506-patient prostate cancer dataset from Byar and Green.⁸⁷ The data are listed in [28, Table 46] and are available in ASCII form from `StatLib` (`lib.stat.cmu.edu`) in the `Datasets` area from this book's Web page. These data were from a randomized trial comparing four treatments for stage 3 and 4 prostate cancer, with almost equal numbers of patients on placebo and each of three doses of estrogen. Four patients had missing values on all of the following variables: `wt`, `pf`, `hx`, `sbp`, `dbp`, `ekg`, `hg`, `bm`; two of these patients were also missing `sz`. These patients are excluded from consideration. The ultimate goal of an analysis of the dataset might be to discover patterns in survival or to do an analysis of covariance to assess the effect of treatment while adjusting for patient heterogeneity. See Chapter 21 for such analyses. The data reductions developed here are general and can be used for a variety of dependent variables.

pf

n missing unique
502 0 4

normal activity (450, 90%), in bed < 50% daytime (37, 7%)
in bed > 50% daytime (13, 3%), confined to bed (2, 0%)

hx : History of Cardiovascular Disease

n missing unique Info Sum Mean
502 0 2 0.73 213 0.4243

sbp : Systolic Blood Pressure/10

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
502 0 18 0.98 14.35 11 12 13 14 16 17 18

Frequency	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	30
%	1	3	14	27	65	74	98	74	72	34	17	12	3	2	3	1	1	1
	0	1	3	5	13	15	20	15	14	7	3	2	1	0	1	0	0	0

dbp : Diastolic Blood Pressure/10

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
502 0 12 0.95 8.149 6 6 7 8 9 10 10

Frequency	4	5	6	7	8	9	10	11	12	13	14	18
%	4	5	43	107	165	94	66	9	5	2	1	1
	1	1	9	21	33	19	13	2	1	0	0	0

ekg

n missing unique
494 8 7

normal (168, 34%), benign (23, 5%)
rhythmic disturb & electrolyte ch (51, 10%)
heart block or conduction def (26, 5%), heart strain (150, 30%)
old MI (75, 15%), recent MI (1, 0%)

hg : Serum Hemoglobin (g/100ml)

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
502 0 91 1 13.45 10.2 10.7 12.3 13.7 14.7 15.8 16.4

lowest : 5.899 7.000 7.199 7.800 8.199
highest: 17.297 17.500 17.598 18.199 21.199

sz: Size of Primary Tumor (cm²)

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
497 5 55 1 14.63 2.0 3.0 5.0 11.0 21.0 32.0 39.2

lowest : 0 1 2 3 4, highest: 54 55 61 62 69

sg : Combined Index of Stage and Hist. Grade

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
491 11 11 0.96 10.31 8 8 9 10 11 13 13

Frequency	5	6	7	8	9	10	11	12	13	14	15
%	3	8	7	67	137	33	114	26	75	5	16
	1	2	1	14	28	7	23	5	15	1	3

```

ap : Serum Prostatic Acid Phosphatase
  n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
502      0      128   1 12.18 0.300 0.300 0.500 0.700 2.975 21.689 38.470

lowest :  0.09999  0.19998  0.29999  0.39996  0.50000
highest: 316.00000 353.50000 367.00000 596.00000 999.87500
    
```

```

bm : Bone Metastases
  n missing unique Info Sum Mean
502      0       2 0.41  82 0.1633
    
```

`stage` is defined by `ap` as well as X-ray results. Of the patients in stage 3, 0.92 have `ap` \leq 0.8. Of those in stage 4, 0.93 have `ap` $>$ 0.8. Since `stage` can be predicted almost certainly from `ap`, we do not consider `stage` in some of the analyses.

8.2 How Many Parameters Can Be Estimated?

There are 354 deaths among the 502 patients. If predicting survival time were of major interest, we could develop a reliable model if no more than about $354/15 = 24$ parameters were *examined* against Y in unpenalized modeling. Suppose that a full model with no interactions is fitted and that linearity is not assumed for any continuous predictors. Assuming `age` is almost linear, we could fit a restricted cubic spline function with three knots. For the other continuous variables, let us use five knots. For categorical predictors, the maximum number of degrees of freedom needed would be one fewer than the number of categories. For `pf` we could lump the last two categories since the last category has only 2 patients. Likewise, we could combine the last two levels of `ekg`. Table 8.1 lists the candidate predictors with the maximum number of parameters we consider for each.

Table 8.1 Degrees of freedom needed for predictors

Predictor:	rx	age	wt	pf	hx	sbp	dbp	ekg	hg	sz	sg	ap	bm
# Parameters:	3	2	4	2	1	4	4	5	4	4	4	4	1

8.3 Redundancy Analysis

As described in Section 4.7.1, it is occasionally useful to do a rigorous redundancy analysis on a set of potential predictors. Let us run the algorithm discussed there, on the set of predictors we are considering. We will use a low threshold (0.3) for R^2 for demonstration purposes.

```
# Allow only 1 d.f. for three of the predictors
prostate <-
  transform(prostate,
            ekg.norm = 1*(ekg %in% c("normal","benign")),
            rxn = as.numeric(rx),
            pfn = as.numeric(pf))
# Force pfn, rxn to be linear because of difficulty of placing
# knots with so many ties in the data
# Note: all incomplete cases are deleted (inefficient)
redun(~ stage + I(rxn) + age + wt + I(pfn) + hx +
      sbp + dbp + ekg.norm + hg + sz + sg + ap + bm,
      r2=.3, type='adjusted', data=prostate)
```

Redundancy Analysis

```
redun(formula = ~stage + I(rxn) + age + wt + I(pfn) + hx +
      sbp + dbp + ekg.norm + hg + sz + sg + ap + bm,
      data = prostate, r2 = 0.3, type = "adjusted")
```

n: 483 p: 14 nk: 3

Number of NAs: 19

Frequencies of Missing Values Due to Each Variable

	stage	I(rxn)	age	wt	I(pfn)	hx	sbp
dbp	0	0	1	2	0	0	0
0							
ekg.norm		hg	sz	sg	ap	bm	
	0	0	5	11	0	0	

Transformation of target variables forced to be linear

R^2 cutoff: 0.3 Type: adjusted

R^2 with which each variable can be predicted from all other variables:

	stage	I(rxn)	age	wt	I(pfn)	hx	sbp
dbp	0.658	0.000	0.073	0.111	0.156	0.062	0.452
0.417							
ekg.norm		hg	sz	sg	ap	bm	
	0.055	0.146	0.192	0.540	0.147	0.391	

Rendundant variables:

stage sbp bm sg

Predicted from variables:

I(rxn) age wt I(pfn) hx dbp ekg.norm hg sz ap

	Variable Deleted	R^2	R^2 after later deletions		
1	stage	0.658	0.658	0.646	0.494
2	sbp	0.452		0.453	0.455
3	bm	0.374			0.367
4	sg	0.342			

By any reasonable criterion on R^2 , none of the predictors is redundant. `stage` can be predicted with an $R^2 = 0.658$ from the other 13 variables, but only with $R^2 = 0.493$ after deletion of 3 variables later declared to be “redundant.”

8.4 Variable Clustering

From Table 8.1, the total number of parameters is 42, so some data reduction should be considered. We resist the temptation to take the “easy way out” using stepwise variable selection so that we can achieve a more stable modeling process and obtain unbiased standard errors. Before using a variable clustering procedure, note that `ap` is extremely skewed. To handle skewness, we use Spearman rank correlations for continuous variables (later we transform each variable using `transcan`, which will allow ordinary correlation coefficients to be used). After classifying `ekg` as “normal/benign” versus everything else, the Spearman correlations are plotted below.

1

```
x ← with(prostate ,
         cbind(stage, rx, age, wt, pf, hx, sbp, dbp,
              ekg.norm, hg, sz, sg, ap, bm))
# If no missing data, could use cor(apply(x, 2, rank))
r ← rcorr(x, type="spearman")$r      # rcorr in Hmisc
maxabsr ← max(abs(r[row(r) != col(r)]))
```

```
p ← nrow(r)
plot(c(-.35,p+.5),c(.5,p+.25), type='n', axes=FALSE,
     xlab='',ylab='') # Figure 8.1
v ← dimnames(r)[[1]]
text(rep(.5,p), 1:p, v, adj=1)
for(i in 1:(p-1)) {
  for(j in (i+1):p) {
    lines(c(i,i),c(j,j+r[i,j]/maxabsr/2),
          lwd=3, lend='butt')
    lines(c(i-.2,i+.2),c(j,j), lwd=1, col=gray(.7))
  }
  text(i, i, v[i], srt=-45, adj=0)
}
```

We perform a hierarchical cluster analysis based on a similarity matrix that contains pairwise Hoeffding D statistics.²⁹⁵ D will detect nonmonotonic associations.

```
vc ← varclus(~ stage + rxn + age + wt + pfn + hx +
             sbp + dbp + ekg.norm + hg + sz + sg + ap + bm,
             sim='hoeffding', data=prostate)
plot(vc) # Figure 8.2
```

We combine `sbp` and `dbp`, and tentatively combine `ap`, `sg`, `sz`, and `bm`.

8.5 Transformation and Single Imputation Using `transcan`

Now we turn to the scoring of the predictors to potentially reduce the number of regression parameters that are needed later by doing away with the need for

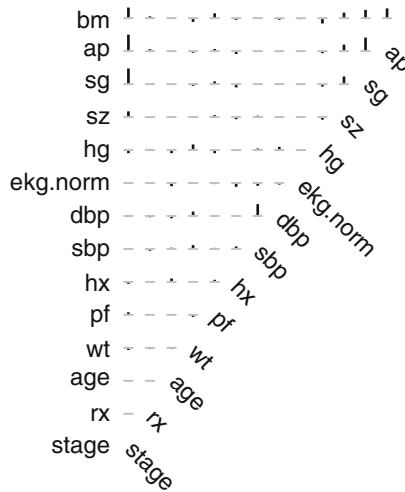


Fig. 8.1 Matrix of Spearman ρ rank correlation coefficients between predictors. Horizontal gray scale lines correspond to $\rho = 0$. The tallest bar corresponds to $|\rho| = 0.78$.

nonlinear terms and multiple dummy variables. The R `Hmisc` package `transcan` function defaults to using a maximum generalized variance method³⁶⁸ that incorporates canonical variates to optimally transform both sides of a multiple regression model. Each predictor is treated in turn as a variable being predicted, and all variables are expanded into restricted cubic splines (for continuous variables) or dummy variables (for categorical ones).

```
# Combine 2 levels of ekg (one had freq. 1)
levels(prostate$ekg)[levels(prostate$ekg) %in%
                    c('old MI', 'recent MI')] ← 'MI'

prostate$pf.coded ← as.integer(prostate$pf)
```

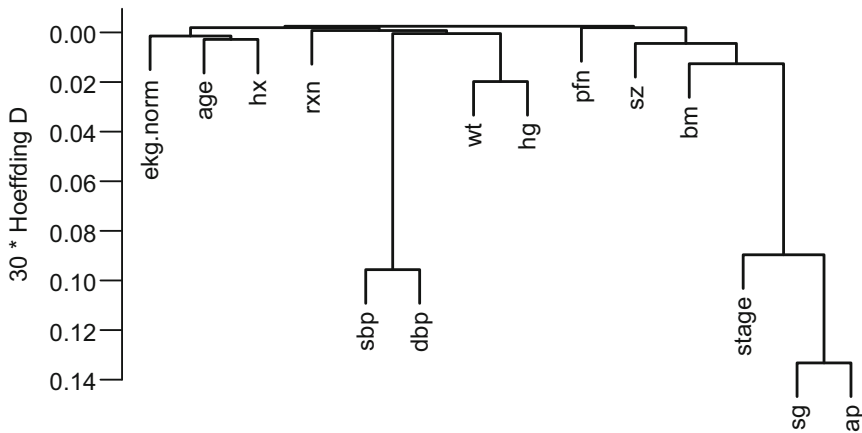


Fig. 8.2 Hierarchical clustering using Hoeffding's D as a similarity measure. Dummy variables were used for the categorical variable `ekg`. Some of the dummy variables cluster together since they are by definition negatively correlated.

```
# make a numeric version; combine last 2 levels of original
levels(prostate$pf) ← levels(prostate$pf)[c(1,2,3,3)]

ptrans ←
  transcan(~ sz + sg + ap + sbp + dbp +
           age + wt + hg + ekg + pf + bm + hx, imputed=TRUE,
           transformed=TRUE, trantab=TRUE, pl=FALSE,
           show.na=TRUE, data=prostate, frac=.1, pr=FALSE)
summary(ptrans, digits=4)
```

```
transcan(x = ~sz + sg + ap + sbp + dbp + age + wt + hg + ekg +
         pf + bm + hx, imputed = TRUE, trantab = TRUE, transformed = TRUE,
         pr = FALSE, pl = FALSE, show.na = TRUE, data = prostate,
         frac = 0.1)

Iterations: 8

R2 achieved in predicting each variable:
      sz      sg      ap      sbp      dbp      age      wt      hg      ekg      pf      bm      hx
0.207 0.556 0.573 0.498 0.485 0.095 0.122 0.158 0.092 0.113 0.349 0.108

Adjusted R2:
      sz      sg      ap      sbp      dbp      age      wt      hg      ekg      pf      bm      hx
0.180 0.541 0.559 0.481 0.468 0.065 0.093 0.129 0.059 0.086 0.331 0.083

Coefficients of canonical variates for predicting each (row) variable
      sz      sg      ap      sbp      dbp      age      wt      hg      ekg      pf      bm
sz      0.66      0.20      0.33      0.33      -0.01      -0.01      0.11      0.11      0.03      -0.36
sg      0.23      0.84      0.08      0.07      -0.02      0.01      -0.01      -0.07      0.02      -0.20
ap      0.07      0.80      -0.11      -0.05      0.03      -0.02      0.01      -0.01      0.00      -0.83
sbp     0.13      0.10      -0.14      -0.94      0.14      -0.09      0.03      0.10      0.10      -0.03
dbp     0.13      0.09      -0.06      -0.98      0.14      0.07      0.05      0.03      0.04      0.03
age     -0.02      -0.06      0.18      0.58      0.57      0.14      0.46      0.43      -0.03      1.05
wt      -0.02      0.06      -0.08      -0.31      0.23      0.12      0.51      -0.06      0.21      -1.09
hg      0.13      -0.02      0.03      0.09      0.15      0.33      0.43      -0.02      0.24      -1.53
ekg     0.20      -0.38      0.10      0.42      0.12      0.41      -0.04      -0.04      0.15      -0.42
pf      0.04      0.08      0.02      0.36      0.14      -0.03      0.22      0.29      0.13      -1.75
bm     -0.02      -0.03      -0.13      0.00      0.00      0.03      -0.04      -0.06      -0.01      -0.06
```



```

hx      0.04  0.05 -0.01 -0.04  0.00 -0.06  0.02 -0.01 -0.09 -0.04 -0.05
  hx
sz      0.34
sg      0.14
ap     -0.03
sbp    -0.14
dbp    -0.01
age    -0.76
wt      0.27
hg     -0.12
ekg    -1.23
pf     -0.46
bm     -0.02
hx

```

Summary of imputed values

```

sz
  n missing  unique  Info  Mean
  5         0         4  0.95 12.86
6 (2, 40%), 7.416 (1, 20%), 20.18 (1, 20%), 24.69 (1, 20%)
sg
  n missing  unique  Info  Mean  .05  .10  .25  .50
 11         0         10    1  10.1  6.900  7.289  7.697  10.270
 .75        .90        .95
10.560  15.000  15.000
 6.511  7.289  7.394  8 10.25 10.27 10.32 10.39 10.73 15
Frequency  1  1  1  1  1  1  1  1  1  2
%          9  9  9  9  9  9  9  9  18
age
  n missing  unique  Info  Mean
  1         0         1    0  71.65
wt
  n missing  unique  Info  Mean
  2         0         2    1  97.77
91.24 (1, 50%), 104.3 (1, 50%)
ekg
  n missing  unique  Info  Mean
  8         0         4    0.9  2.625
1 (3, 38%), 3 (3, 38%), 4 (1, 12%), 5 (1, 12%)

```

Starting estimates for imputed values:

```

  sz  sg  ap  sbp  dbp  age  wt  hg  ekg  pf  bm  hx
11.0 10.0 0.7 14.0  8.0 73.0 98.0 13.7  1.0  1.0  0.0  0.0

```

```

ggplot(ptrans, scale=TRUE) +
  theme(axis.text.x=element_text(size=6)) # Figure 8.3

```

The plotted output is shown in Figure 8.3. Note that at face value the transformation of `ap` was derived in a circular manner, since the combined index of stage and histologic grade, `sg`, uses in its stage component a cutoff on `ap`. However, if `sg` is omitted from consideration, the resulting transformation for `ap` does not change appreciably. Note that `bm` and `hx` are represented as binary variables, so their coefficients in the table of canonical variable coefficients are on a different scale. For the variables that were actually transformed, the coefficients are for standardized transformed variables (mean 0, variance 1). From examining the R^2 s, `age`, `wt`, `ekg`, `pf`, and `hx` are not strongly related to other variables. Imputations for `age`, `wt`, `ekg` are thus relying more on the median or modal values from the marginal distributions. From the coefficients of first (standardized) canonical variates, `sbp` is predicted almost solely from `dbp`; `bm` is predicted mainly from `ap`, `hg`, and `pf`.

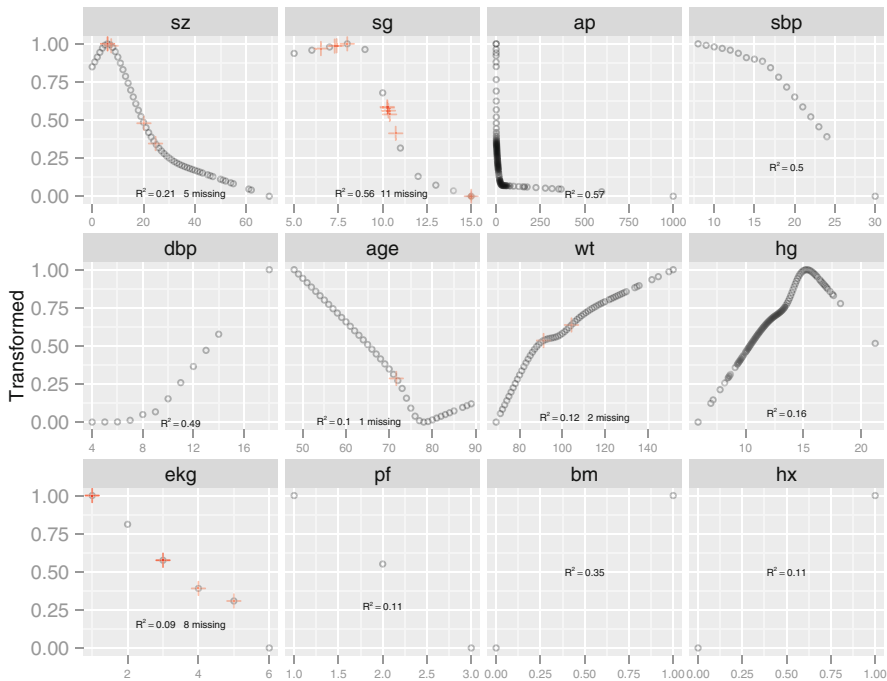


Fig. 8.3 Simultaneous transformation and single imputation of all candidate predictors using *transcan*. Imputed values are shown as red plus signs. Transformed values are arbitrarily scaled to $[0, 1]$.

8.6 Data Reduction Using Principal Components

The first PC, PC_1 , is the linear combination of standardized variables having maximum variance. PC_2 is the linear combination of predictors having the second largest variance such that PC_2 is orthogonal to (uncorrelated with) PC_1 . If there are p raw variables, the first k PCs, where $k < p$, will explain only part of the variation in the whole system of p variables unless one or more of the original variables is exactly a linear combination of the remaining variables. Note that it is common to scale and center variables to have mean zero and variance 1 before computing PCs.

The response variable (here, time until death due to any cause) is not examined during data reduction, so that if PCs are selected by variance explained in the X -space and not by variation explained in Y , one needn't correct for model uncertainty or multiple comparisons.

PCA results in data reduction when the analyst uses only a subset of the p possible PCs in predicting Y . This is called *incomplete principal component regression*. When one sequentially enters PCs into a predictive model in a strict pre-specified order (i.e., by descending amounts of variance explained

for the system of p variables), model uncertainty requiring bootstrap adjustment is minimized. In contrast, model uncertainty associated with stepwise regression (driven by associations with Y) is massive.

For the prostate dataset, consider PCs on raw candidate predictors, expanding polytomous factors using dummy variables. The R function `princomp` is used, after singly imputing missing raw values using `transcan`'s optimal additive nonlinear models. In this series of analyses we ignore the treatment variable, `rx`.

```
# Impute all missing values in all variables given to transcan
imputed <- impute(ptrans, data=prostate, list.out=TRUE)
```

Imputed missing values with the following frequencies
and stored them in variables with their original names:

```
sz  sg  age  wt  ekg
5   11  1   2   8
```

```
imputed <- as.data.frame(imputed)

# Compute principal components on imputed data.
# Create a design matrix from ekg categories
Ekg <- model.matrix(~ ekg, data=imputed)[, -1]
# Use correlation matrix
pfn <- prostate$pfn
prin.raw <- princomp(~ sz + sg + ap + sbp + dbp + age +
                    wt + hg + Ekg + pfn + bm + hx,
                    cor=TRUE, data=imputed)

plot(prin.raw, type='lines', main='', ylim=c(0,3)) #Figure 8.4
# Add cumulative fraction of variance explained
addscree <- function(x, npcs=min(10, length(x$sdev)),
                    plotv=FALSE,
                    col=1, offset=.8, adj=0, pr=FALSE) {
  vars <- x$sdev^2
  cumv <- cumsum(vars)/sum(vars)
  if(pr) print(cumv)
  text(1:npcs, vars[1:npcs] + offset*par('cxy')[2],
       as.character(round(cumv[1:npcs], 2)),
       srt=45, adj=adj, cex=.65, xpd=NA, col=col)
  if(plotv) lines(1:npcs, vars[1:npcs], type='b', col=col)
}
addscree(prin.raw)
prin.trans <- princomp(ptrans$transformed, cor=TRUE)
addscree(prin.trans, npcs=10, plotv=TRUE, col='red',
         offset=-.8, adj=1)
```

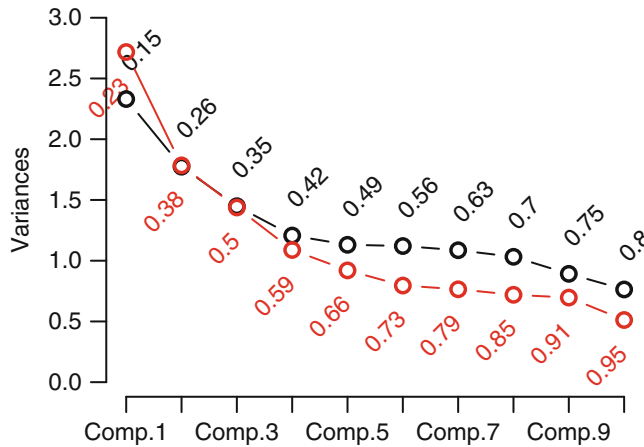


Fig. 8.4 Variance of the system of raw predictors (black) explained by individual principal components (lines) along with cumulative proportion of variance explained (text), and variance explained by components computed on `transcan`-transformed variables (red)

The resulting plot shown in Figure 8.4 is called a “scree” plot [325, pp. 96–99, 104, 106]. It shows the variation explained by the first k principal components as k increases all the way to 16 parameters (no data reduction). It requires 10 of the 16 possible components to explain > 0.8 of the variance, and the first 5 components explain 0.49 of the variance of the system. Two of the 16 dimensions are almost totally redundant.

After repeating this process when transforming all predictors via `transcan`, we have only 12 degrees of freedom for the 12 predictors. The variance explained is depicted in Figure 8.4 in red. It requires at least 9 of the 12 possible components to explain ≥ 0.9 of the variance, and the first 5 components explain 0.66 of the variance as opposed to 0.49 for untransformed variables.

Let us see how the PCs “explain” the times until death using the Cox regression¹³² function from `rms`, `cph`, described in Chapter 20. In what follows we vary the number of components used in the Cox models from 1 to all 16, computing the AIC for each model. AIC is related to model log likelihood penalized for number of parameters estimated, and lower is better. For reference, the AIC of the model using all of the original predictors, and the AIC of a full additive spline model are shown as horizontal lines.

```
require(rms)
```

```
S <- with(prostate, Surv(dtime, status != "alive"))
# two-column response var.

pcs <- prin.raw$scores          # pick off all PCs
aic <- numeric(16)
for(i in 1:16) {
```

```

  ps ← pcs[,1:i]
  aic[i] ← AIC(cph(S ~ ps))
} # Figure 8.5
plot(1:16, aic, xlab='Number of Components Used',
     ylab='AIC', type='l', ylim=c(3950,4000))
f ← cph(S ~ sz + sg + log(ap) + sbp + dbp + age + wt + hg +
        ekg + pf + bm + hx, data=imputed)
abline(h=AIC(f), col='blue')
f ← cph(S ~ rcs(sz,5) + rcs(sg,5) + rcs(log(ap),5) +
        rcs(sbp,5) + rcs(dbp,5) + rcs(age,3) + rcs(wt,5) +
        rcs(hg,5) + ekg + pf + bm + hx,
        tol=1e-14, data=imputed)

```

```
abline(h=AIC(f), col='blue', lty=2)
```

For the money, the first 5 components adequately summarizes all variables, if linearly transformed, and the full linear model is no better than this. The model allowing all continuous predictors to be nonlinear is not worth its added degrees of freedom.

Next check the performance of a model derived from cluster scores of transformed variables.

```

# Compute PC1 on a subset of transcan-transformed predictors
pco ← function(v) {
  f ← princomp(ptrans$transformed[,v], cor=TRUE)
  vars ← f$sdev^2
  cat('Fraction of variance explained by PC1:',
      round(vars[1]/sum(vars),2), '\n')
  f$scores[,1]
}
tumor ← pco(c('sz','sg','ap','bm'))

```

```
Fraction of variance explained by PC1: 0.59
```

```
bp ← pco(c('sbp','dbp'))
```

```
Fraction of variance explained by PC1: 0.84
```

```
cardiac ← pco(c('hx','ekg'))
```

```
Fraction of variance explained by PC1: 0.61
```

```

# Get transformed individual variables that are not clustered
other ← ptrans$transformed[,c('hg','age','pf','wt')]
f ← cph(S ~ tumor + bp + cardiac + other) # other is matrix
AIC(f)

```

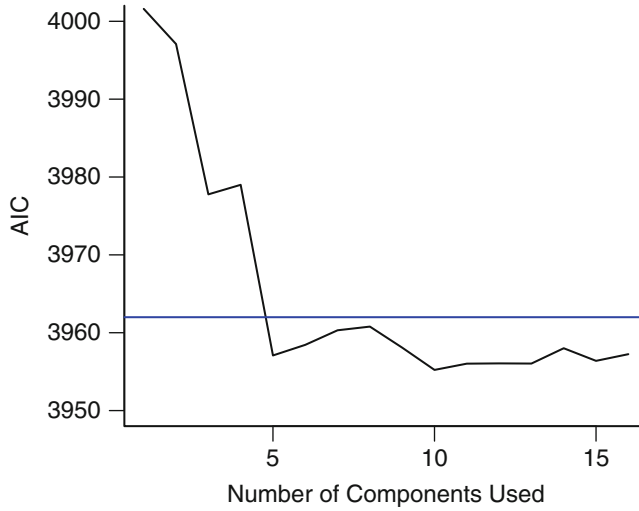


Fig. 8.5 AIC of Cox models fitted with progressively more principal components. The solid blue line depicts the AIC of the model with all original covariates. The dotted blue line is positioned at the AIC of the full spline model.

```
[1] 3954.393
```

```
print(f, latex=TRUE, long=FALSE, title='')
```

		Model Tests		Discrimination Indexes	
Obs	502	LR χ^2	81.11	R^2	0.149
Events	354	d.f.	7	D_{xy}	0.286
Center	0	$\Pr(> \chi^2)$	0.0000	g	0.562
		Score χ^2	86.81	g_r	1.755
		$\Pr(> \chi^2)$	0.0000		

	Coef	S.E.	Wald Z	$\Pr(> Z)$
tumor	-0.1723	0.0367	-4.69	< 0.0001
bp	-0.0251	0.0424	-0.59	0.5528
cardiac	-0.2513	0.0516	-4.87	< 0.0001
hg	-0.1407	0.0554	-2.54	0.0111
age	-0.1034	0.0579	-1.79	0.0739
pf	-0.0933	0.0487	-1.92	0.0551
wt	-0.0910	0.0555	-1.64	0.1012

The `tumor` and `cardiac` clusters seem to dominate prediction of mortality, and the AIC of the model built from cluster scores of transformed variables compares favorably with other models (Figure 8.5).

8.6.1 Sparse Principal Components

A disadvantage of principal components is that every predictor receives a nonzero weight for every component, so many coefficients are involved even through the effective degrees of freedom with respect to the response model are reduced. *Sparse principal components*⁶⁷² uses a penalty function to reduce the magnitude of the loadings variables receive in the components. If an L1 penalty is used (as with the *lasso*), some loadings are shrunk to zero, resulting in some simplicity. Sparse principal components combines some elements of variable clustering, scoring of variables within clusters, and redundancy analysis.

Filzmoser, Fritz, and Kalcher¹⁹¹ have written a nice R package `pcaPP` for doing sparse PC analysis.^a The following example uses the prostate data again. To allow for nonlinear transformations and to score the `ekg` variable in the prostate dataset down to a scalar, we use the `transcan`-transformed predictors as inputs.

```
require(pcaPP)
```

```
s <- sPCAgrid(ptrans$transformed, k=10, method='sd',
              center=mean, scale=sd, scores=TRUE,
              maxiter=10)
plot(s, type='lines', main='', ylim=c(0,3)) # Figure 8.6
addscree(s)
s$loadings # These loadings are on the orig. transcan scale
```

```
Loadings:
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8 Comp.9 Comp.10
sz    0.248
sg    0.620
ap    0.634
sbp   -0.707
dbp    0.707
age           1.000
wt           1.000
hg           1.000
ekg           1.000
pf           1.000
bm   -0.391
hx           1.000
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7 Comp.8
SS loadings 1.000 1.000 1.000 1.000 1.000 1.000 1.000 1.000
Proportion Var 0.083 0.083 0.083 0.083 0.083 0.083 0.083 0.083
Cumulative Var 0.083 0.167 0.250 0.333 0.417 0.500 0.583 0.667
      Comp.9 Comp.10
SS loadings 1.000 1.000
Proportion Var 0.083 0.083
Cumulative Var 0.750 0.833
```

Only nonzero loadings are shown. The first sparse PC is the tumor cluster used above, and the second is the blood pressure cluster. Let us see how well incomplete sparse principal component regression predicts time until death.

^a The `spca` package is a new sparse PC package that should also be considered.

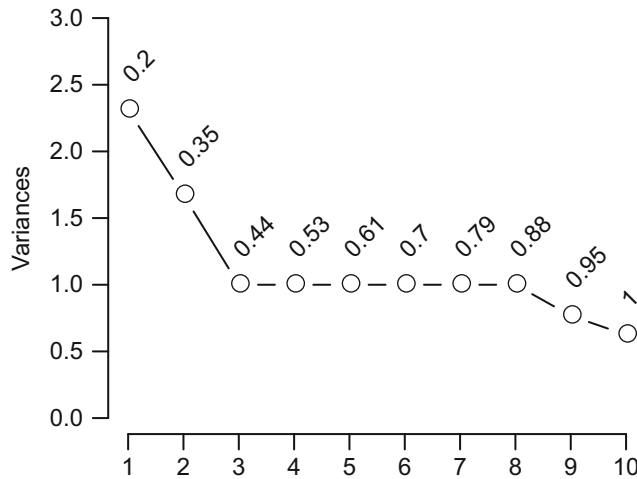


Fig. 8.6 Variance explained by individual sparse principal components (lines) along with cumulative proportion of variance explained (text)

```

pcs ← s$scores          # pick off sparse PCs
aic ← numeric(10)
for(i in 1:10) {
  ps ← pcs[,1:i]
  aic[i] ← AIC(cph(S ~ ps))
} # Figure 8.7
plot(1:10, aic, xlab='Number of Components Used',
     ylab='AIC', type='l', ylim=c(3950,4000))

```

More components are required to optimize AIC than were seen in Figure 8.5, but a model built from 6–8 sparse PCs performed as well as the other models.

8.7 Transformation Using Nonparametric Smoothers

The ACE nonparametric additive regression method of Breiman and Friedman⁶⁸ transforms both the left-hand-side variable and all the right-hand-side variables so as to optimize R^2 . ACE can be used to transform the predictors using the R `ace` function in the `acepack` package, called by the `transace` function in the `Hmisc` package. `transace` does not impute data but merely does casewise deletion of missing values. Here `transace` is run after single imputation by `transcan`. `binary` is used to tell `transace` which variables not to try to predict (because they need no transformation). Several predictors are restricted to be monotonically transformed.

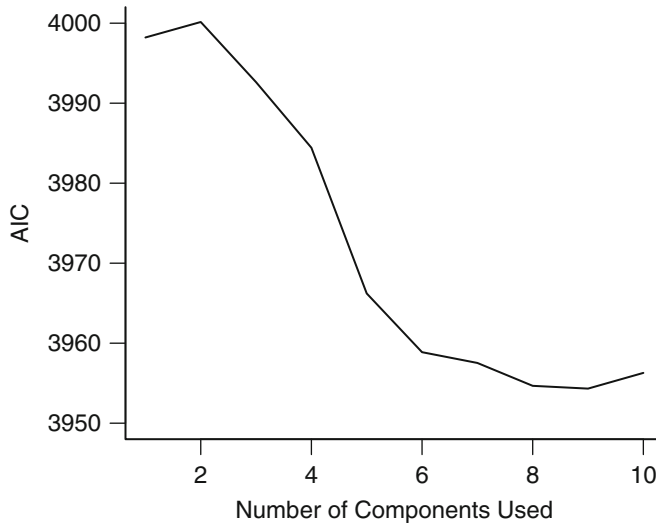


Fig. 8.7 Performance of sparse principal components in Cox models

```
x ← with(imputed,
         cbind(sz, sg, ap, sbp, dbp, age, wt, hg, ekg, pf,
              bm, hx))
monotonic ← c("sz", "sg", "ap", "sbp", "dbp", "age", "pf")
transace(x, monotonic, # Figure 8.8
         categorical="ekg", binary=c("bm", "hx"))
```

R^2 achieved in predicting each variable:

sz	sg	ap	sbp	dbp	age	wt
0.2265824	0.5762743	0.5717747	0.4823852	0.4580924	0.1514527	0.1732244
hg	ekg	pf	bm	hx		
0.2001008	0.1110709	0.1778705	NA	NA		

Except for `ekg`, `age`, and for arbitrary sign reversals, the transformations in Figure 8.8 determined using `transace` were similar to those in Figure 8.3. The `transcan` transformation for `ekg` makes more sense.

8.8 Further Reading

- 1 Sauerbrei and Schumacher⁵⁴¹ used the bootstrap to demonstrate the variability of a standard variable selection procedure for the prostate cancer dataset.
- 2 Schemper and Heinze⁵⁵¹ used logistic models to impute dichotomizations of the predictors for this dataset.

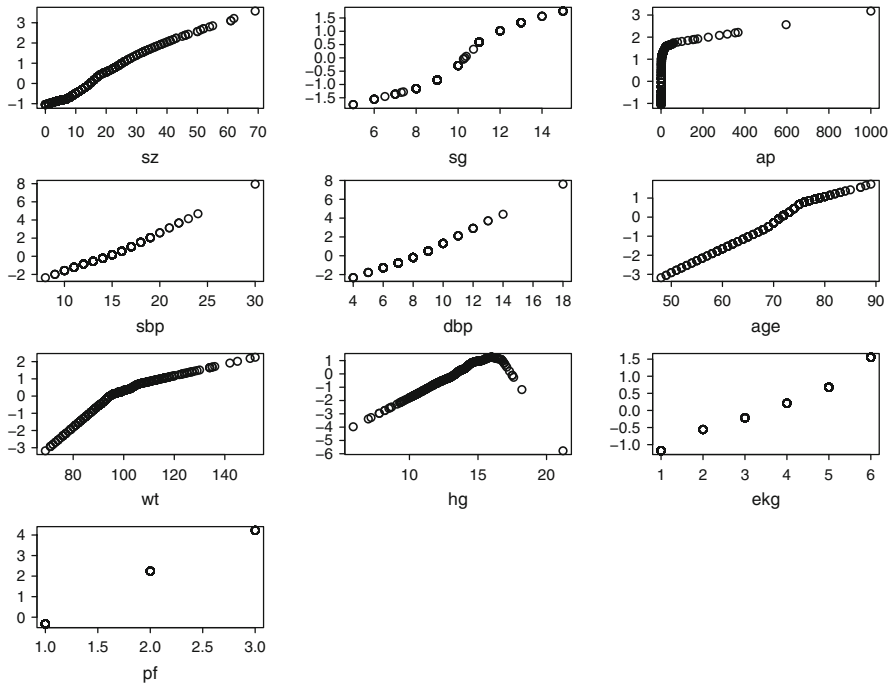


Fig. 8.8 Simultaneous transformation of all variables using ACE.

8.9 Problems

The Mayo Clinic conducted a randomized trial in primary biliary cirrhosis (PBC) of the liver between January 1974 and May 1984, to compare D-penicillamine with placebo. The drug was found to be ineffective [197, p. 2], and the trial was done before liver transplantation was common, so this trial constitutes a natural history study for PBC. Followup continued through July, 1986. For the 19 patients that did undergo transplant, followup time was censored (`status=0`) at the day of transplant. 312 patients were randomized, and another 106 patients were entered into a registry. The nonrandomized patients have most of their laboratory values missing, except for bilirubin, albumin, and prothrombin time. 28 randomized patients had both serum cholesterol and triglycerides missing. The data, which consist of clinical, biochemical, serologic, and histologic information, are listed in [197, pp. 359–375]. The PBC data are discussed and analyzed in [197, pp. 2–7, 102–104, 153–162], [158], [7] (a tree-based analysis which on its p. 480 mentions some possible lack of fit of the earlier analyses), and [361]. The data are stored in the datasets web site so may be accessed using the `Hmisc::getHdata` function with argument `pbc`. Use only the data on randomized patients for all analyses. For Problems 1–6, ignore followup time, status, and drug.

1. Do an initial variable clustering based on ranks, using pairwise deletion of missing data. Comment on the potential for one-dimensional summaries of subsets of variables being adequate summaries of prognostic information.
2. `cholesterol`, `triglycerides`, `platelets`, and `copper` are missing on some patients. Impute them using a method you recommend. Use some or all of the remaining predictors and possibly the outcome. Provide a correlation coefficient describing the usefulness of each imputation model. Provide the actual imputed values, specifying observation numbers. For all later analyses, use imputed values for missing values.
3. Perform a scaling/transformation analysis to better measure how the predictors interrelate and to possibly pretransform some of them. Use `transcan` or `ACE`. Repeat the variable clustering using the transformed scores and Pearson correlation or using an oblique rotation principal component analysis. Determine if the correlation structure (or variance explained by the first principal component) indicates whether it is possible to summarize multiple variables into single scores.
4. Do a principal component analysis of all transformed variables simultaneously. Make a graph of the number of components versus the cumulative proportion of explained variation. Repeat this for laboratory variables alone.
5. Repeat the overall PCA using sparse principal components. Pay attention to how best to solve for sparse components, e.g., consider the `lambda` parameter in `sPCAgrid`.
6. How well can variables (lab and otherwise) that are routinely collected (on nonrandomized patients) capture the information (variation) of the variables that are often missing? It would be helpful to explore the strength of interrelationships by
 - a. correlating two PC_1 s obtained from untransformed variables,
 - b. correlating two PC_1 s obtained from transformed variables,
 - c. correlating the best linear combination of one set of variables with the best linear combination of the other set, and
 - d. doing the same on transformed variables.

For this problem consider only complete cases, and transform the 5 non-numeric categorical predictors to binary 0–1 variables.

7. Consider the patients having complete data who were randomized to placebo. Consider only models that are linear in all the covariates.
 - a. Fit a survival model to predict time of death using the following covariates: `bili`, `albumin`, `stage`, `protime`, `age`, `alk.phos`, `sgot`, `chol`, `trig`, `platelet`, `copper`.
 - b. Perform an ordinary principal component analysis. Fit the survival model using only the first 3 PCs. Compare the likelihood ratio χ^2 and AIC with that of the model using the original variables.

- c. Considering the PCs are fixed, use the bootstrap to estimate the 0.95 confidence interval of the inter-quartile-range age effect on the original scale, and the same type of confidence interval for the coefficient of PC_1 .
- d. Now accounting for uncertainty in the PCs, compute the same two confidence intervals. Compare and interpret the two sets. Take into account the fact that PCs are not unique to within a sign change.

R programming hints for this exercise are found on the course web site.

Chapter 9

Overview of Maximum Likelihood Estimation

9.1 General Notions—Simple Cases

In ordinary least squares multiple regression, the objective in fitting a model is to find the values of the unknown parameters that minimize the sum of squared errors of prediction. When the response variable is non-normal, polytomous, or not observed completely, one needs a more general objective function to optimize.

Maximum likelihood (ML) estimation is a general technique for estimating parameters and drawing statistical inferences in a variety of situations, especially nonstandard ones. Before laying out the method in general, ML estimation is illustrated with a standard situation, the one-sample binomial problem. Here, independent binary responses are observed and one wishes to draw inferences about an unknown parameter, the probability of an event in a population.

Suppose that in a population of individuals, each individual has the same probability P that an event occurs. We could also say that the event has already been observed, so that P is the prevalence of some condition in the population. For each individual, let $Y = 1$ denote the occurrence of the event and $Y = 0$ denote nonoccurrence. Then $\text{Prob}\{Y = 1\} = P$ for each individual. Suppose that a random sample of size 3 from the population is drawn and that the first individual had $Y = 1$, the second had $Y = 0$, and the third had $Y = 1$. The respective probabilities of these outcomes are P , $1 - P$, and P . The joint probability of observing the independent events $Y = 1, 0, 1$ is $P(1 - P)P = P^2(1 - P)$. Now the value of P is unknown, but we can solve for the value of P that makes the observed data $(Y = 1, 0, 1)$ *most likely to have occurred*. In this case, the value of P that maximizes $P^2(1 - P)$ is $P = 2/3$. This value for P is the *maximum likelihood estimate (MLE)* of the population probability.

Let us now study the situation of independent binary trials in general. Let the sample size be n and the observed responses be Y_1, Y_2, \dots, Y_n . The joint probability of observing the data is given by

$$L = \prod_{i=1}^n P^{Y_i} (1 - P)^{1 - Y_i}. \quad (9.1)$$

Now let s denote the sum of the Y 's or the number of times that the event occurred ($Y_i = 1$), that is the number of "successes." The number of non-occurrences ("failures") is $n - s$. The likelihood of the data can be simplified to

$$L = P^s (1 - P)^{n - s}. \quad (9.2)$$

It is easier to work with the *log likelihood function*, which also has desirable statistical properties. For the one-sample binary response problem, the log likelihood is

$$\log L = s \log(P) + (n - s) \log(1 - P). \quad (9.3)$$

The MLE of P is that value of P that maximizes L or $\log L$. Since $\log L$ is a smooth function of P , its maximum value can be found by finding the point at which $\log L$ has a slope of 0. The slope or first derivative of $\log L$, with respect to P , is

$$U(P) = \partial \log L / \partial P = s/P - (n - s)/(1 - P). \quad (9.4)$$

The first derivative of the log likelihood function with respect to the parameter(s), here $U(P)$, is called the *score function*. Equating this function to zero requires that $s/P = (n - s)/(1 - P)$. Multiplying both sides of the equation by $P(1 - P)$ yields $s(1 - P) = (n - s)P$ or that $s = (n - s)P + sP = nP$. Thus the MLE of P is $p = s/n$.

Another important function is called the *Fisher information* about the unknown parameters. The information function is the expected value of the negative of the curvature in $\log L$, which is the negative of the slope of the slope as a function of the parameter, or the negative of the second derivative of $\log L$. Motivation for consideration of the Fisher information is as follows. If the log likelihood function has a distinct peak, the sample provides information that allows one to readily discriminate between a good parameter estimate (the location of the obvious peak) and a bad one. In such a case the MLE will have good precision or small variance. If on the other hand the likelihood function is relatively flat, almost any estimate will do and the chosen estimate will have poor precision or large variance. The degree of peakedness of a function at a given point is the speed with which the slope is changing at that point, that is, the slope of the slope or second derivative of the function at that point.

Here, the information is

$$\begin{aligned}
 I(P) &= E\{-\partial^2 \log L / \partial P^2\} \\
 &= E\{s/P^2 + (n - s)/(1 - P)^2\} \\
 &= nP/P^2 + n(1 - P)/(1 - P)^2 = n/[P(1 - P)].
 \end{aligned}
 \tag{9.5}$$

We estimate the information by substituting the MLE of P into $I(P)$, yielding $I(p) = n/[p(1 - p)]$.

Figures 9.1, 9.2, and 9.3 depict, respectively, $\log L$, $U(P)$, and $I(P)$, all as a function of P . Three combinations of n and s were used in each graph. These combinations correspond to $p = .5, .6, \text{ and } .6$, respectively.

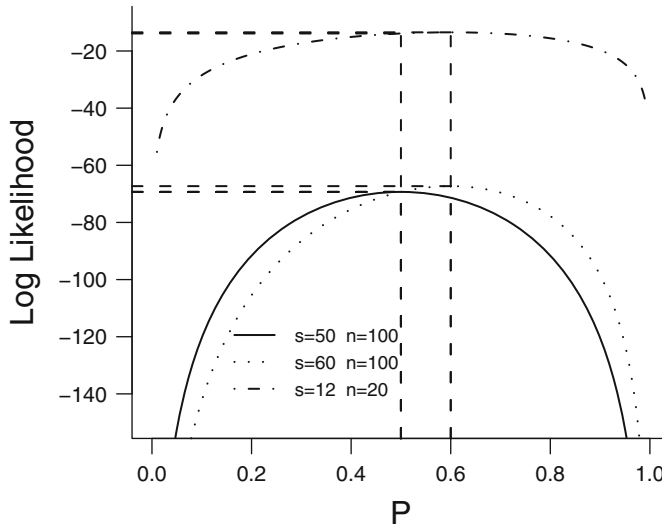


Fig. 9.1 log likelihood functions for three one-sample binomial problems

In each case it can be seen that the value of P that makes the data most likely to have occurred (the value that maximizes L or $\log L$) is p given above. Also, the score function (slope of $\log L$) is zero at $P = p$. Note that the information function $I(P)$ is highest for P approaching 0 or 1 and is lowest for P near .5, where there is maximum uncertainty about P . Note also that while $\log L$ has the same shape for the $s = 60$ and $s = 12$ curves in Figure 9.1, the range of $\log L$ is much greater for the larger sample size. Figures 9.2 and 9.3 show that the larger sample size produces a sharper likelihood. In other words, with larger n , one can zero in on the true value of P with more precision.

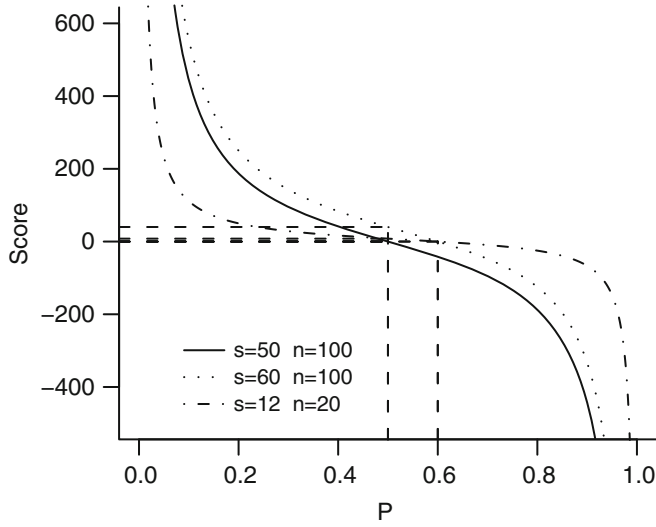


Fig. 9.2 Score functions $(\partial L / \partial P)$

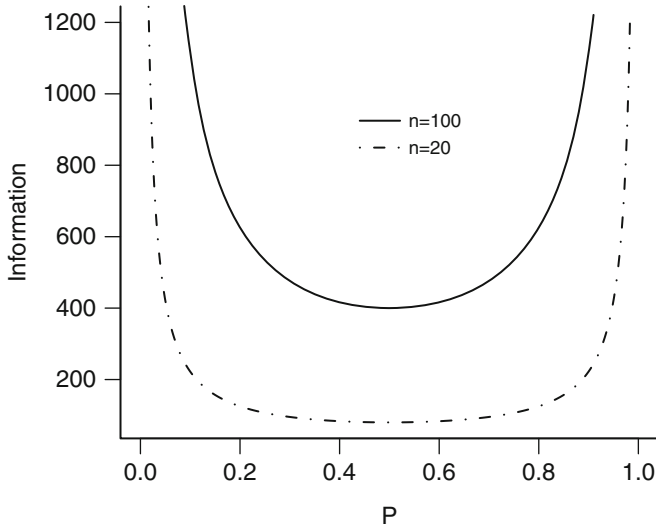


Fig. 9.3 Information functions $(-\partial^2 \log L / \partial P^2)$

In this binary response one-sample example let us now turn to inference about the parameter P . First, we turn to the estimation of the variance of the MLE, p . An estimate of this variance is given by the inverse of the information at $P = p$:

$$\text{Var}(p) = I(p)^{-1} = p(1-p)/n. \quad (9.6)$$

Note that the variance is smallest when the information is greatest ($p = 0$ or 1).

The variance estimate forms a basis for confidence limits on the unknown parameter. For large n , the MLE p is approximately normally distributed with expected value (mean) P and variance $P(1-P)/n$. Since $p(1-p)$ is a consistent estimate of $P(1-P)/n$, it follows that $p \pm z[p(1-p)/n]^{1/2}$ is an approximate $1 - \alpha$ confidence interval for P if z is the $1 - \alpha/2$ critical value of the standard normal distribution.

9.2 Hypothesis Tests

Now let us turn to hypothesis tests about the unknown population parameter $P - H_0 : P = P_0$. There are three kinds of statistical tests that arise from likelihood theory.

9.2.1 Likelihood Ratio Test

This test statistic is the ratio of the likelihood at the hypothesized parameter values to the likelihood of the data at the maximum (i.e., at parameter values = MLEs). It turns out that $-2 \times$ the log of this likelihood ratio has desirable statistical properties. The likelihood ratio test statistic is given by

$$\begin{aligned} LR &= -2 \log(L \text{ at } H_0 / L \text{ at MLEs}) \\ &= -2(\log L \text{ at } H_0) - [-2(\log L \text{ at MLEs})]. \end{aligned} \quad (9.7)$$

The LR statistic, for large enough samples, has approximately a χ^2 distribution with degrees of freedom equal to the number of parameters estimated, if the null hypothesis is “simple,” that is, doesn’t involve any unknown parameters. Here LR has 1 d.f.

The value of $\log L$ at H_0 is

$$\log L(H_0) = s \log(P_0) + (n - s) \log(1 - P_0). \quad (9.8)$$

The maximum value of $\log L$ (at MLEs) is

$$\log L(P = p) = s \log(p) + (n - s) \log(1 - p). \quad (9.9)$$

For the hypothesis $H_0 : P = P_0$, the test statistic is

$$LR = -2\{s \log(P_0/p) + (n - s) \log[(1 - P_0)/(1 - p)]\}. \quad (9.10)$$

Note that when p happens to equal P_0 , $LR = 0$. When p is far from P_0 , LR will be large. Suppose that $P_0 = 1/2$, so that H_0 is $P = 1/2$. For $n = 100$, $s = 50$, $LR = 0$. For $n = 100$, $s = 60$,

$$LR = -2\{60 \log(.5/.6) + 40 \log(.5/.4)\} = 4.03. \quad (9.11)$$

For $n = 20$, $s = 12$,

$$LR = -2\{12 \log(.5/.6) + 8 \log(.5/.4)\} = .81 = 4.03/5. \quad (9.12)$$

Therefore, even though the best estimate of P is the same for these two cases, the test statistic is more impressive when the sample size is five times larger.

9.2.2 Wald Test

The Wald test statistic is a generalization of a t - or z -statistic. It is a function of the difference in the MLE and its hypothesized value, normalized by an estimate of the standard deviation of the MLE. Here the statistic is

$$W = [p - P_0]^2 / [p(1 - p)/n]. \quad (9.13)$$

For large enough n , W is distributed as χ^2 with 1 d.f. For $n = 100$, $s = 50$, $W = 0$. For the other samples, W is, respectively, 4.17 and 0.83 (note $0.83 = 4.17/5$).

Many statistical packages treat \sqrt{W} as having a t distribution instead of a normal distribution. As pointed out by Gould,²²⁸ there is no basis for this outside of ordinary linear models^a.

9.2.3 Score Test

If the MLE happens to equal the hypothesized value P_0 , P_0 maximizes the likelihood and so $U(P_0) = 0$. Rao's score statistic measures how far from zero the score function is when evaluated at the null hypothesis. The score function

^a In linear regression, a t distribution is used to penalize for the fact that the variance of $Y|X$ is estimated. In models such as the logistic model, there is no separate variance parameter to estimate. Gould has done simulations that show that the normal distribution provides more accurate P -values than the t for binary logistic regression.

(slope or first derivative of $\log L$) is normalized by the information (curvature or second derivative of $-\log L$). The test statistic for our example is

$$S = U(P_0)^2/I(P_0), \quad (9.14)$$

which formally does not involve the MLE, p . The statistic can be simplified as follows.

$$\begin{aligned} U(P_0) &= s/P_0 - (n-s)/(1-P_0) \\ I(P_0) &= s/P_0^2 + (n-s)/(1-P_0)^2 \\ S &= (s-nP_0)^2/[nP_0(1-P_0)] = n(p-P_0)^2/[P_0(1-P_0)]. \end{aligned} \quad (9.15)$$

Note that the numerator of S involves $s-nP_0$, the difference between the observed number of successes and the number of successes expected under H_0 .

As with the other two test statistics, $S = 0$ for the first sample. For the last two samples S is, respectively, 4 and $.8 = 4/5$.

1

9.2.4 Normal Distribution—One Sample

Suppose that a sample of size n is taken from a population for a random variable Y that is known to be normally distributed with unknown mean μ and variance σ^2 . Denote the observed values of the random variable by Y_1, Y_2, \dots, Y_n . Now unlike the binary response case ($Y = 0$ or 1), we cannot use the notion of the probability that Y equals an observed value. This is because Y is continuous and the probability that it will take on a given value is zero. We substitute the *density function* for the probability. The density at a point y is the limit as d approaches zero of

$$\text{Prob}\{y < Y \leq y + d\}/d = [F(y+d) - F(y)]/d, \quad (9.16)$$

where $F(y)$ is the normal cumulative distribution function (for a mean of μ and variance of σ^2). The limit of the right-hand side of the above equation as d approaches zero is $f(y)$, the density function of a normal distribution with mean μ and variance σ^2 . This density function is

$$f(y) = (2\pi\sigma^2)^{-1/2} \exp\{-(y-\mu)^2/2\sigma^2\}. \quad (9.17)$$

The likelihood of observing the observed sample values is the joint density of the Y s. The log likelihood function here is a function of two unknowns, μ and σ^2 .

$$\log L = -.5n \log(2\pi\sigma^2) - .5 \sum_{i=1}^n (Y_i - \mu)^2/\sigma^2. \quad (9.18)$$

It can be shown that the value of μ that maximizes $\log L$ is the value that minimizes the sum of squared deviations about μ , which is the sample mean \bar{Y} . The MLE of σ^2 is

$$s^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 / n. \quad (9.19)$$

Recall that the sample variance uses $n - 1$ instead of n in the denominator. It can be shown that the expected value of the MLE of σ^2 , s^2 , is $[(n - 1)/n]\sigma^2$; in other words, s^2 is too small by a factor of $(n - 1)/n$ on the average. The sample variance is unbiased, but being unbiased does not necessarily make it a better estimator. The MLE has greater precision (smaller mean squared error) in many cases.

9.3 General Case

Suppose we need to estimate a vector of unknown parameters $B = \{B_1, B_2, \dots, B_p\}$ from a sample of size n based on observations Y_1, \dots, Y_n . Denote the probability or density function of the random variable Y for the i th observation by $f_i(y; B)$. The likelihood for the i th observation is $L_i(B) = f_i(Y_i; B)$. In the one-sample binary response case, recall that $L_i(B) = L_i(P) = P^{Y_i}[1 - P]^{1 - Y_i}$. The likelihood function, or joint likelihood of the sample, is given by

$$L(B) = \prod_{i=1}^n f_i(Y_i; B). \quad (9.20)$$

The log likelihood function is

$$\log L(B) = \sum_{i=1}^n \log L_i(B). \quad (9.21)$$

The MLE of B is that value of the vector B that maximizes $\log L(B)$ as a function of B . In general, the solution for B requires iterative trial-and-error methods as outlined later. Denote the MLE of B as $b = \{b_1, \dots, b_p\}$. The *score vector* is the vector of first derivatives of $\log L(B)$ with respect to B_1, \dots, B_p :

$$\begin{aligned} U(B) &= \{\partial/\partial B_1 \log L(B), \dots, \partial/\partial B_p \log L(B)\} \\ &= (\partial/\partial B) \log L(B). \end{aligned} \quad (9.22)$$

The Fisher *information matrix* is the $p \times p$ matrix whose elements are the negative of the expectation of all second partial derivatives of $\log L(B)$:

$$\begin{aligned} I^*(B) &= -\{E[(\partial^2 \log L(B)/\partial B_j \partial B_k)]\}_{p \times p} \\ &= -E\{(\partial^2/\partial B \partial B') \log L(B)\}. \end{aligned} \quad (9.23)$$

The *observed information matrix* $I(B)$ is $I^*(B)$ without taking the expectation. In other words, observed values remain in the second derivatives:

$$I(B) = -(\partial^2/\partial B\partial B') \log L(B). \quad (9.24)$$

This information matrix is often estimated from the sample using the *estimated observed information* $I(b)$, by inserting b , the MLE of B , into the formula for $I(B)$.

Under suitable conditions, which are satisfied for most situations likely to be encountered, the MLE b for large samples is an optimal estimator (has as great a chance of being close to the true parameter as all other types of estimators) and has an approximate multivariate normal distribution with mean vector B and variance-covariance matrix $I^{*-1}(B)$, where C^{-1} denotes the inverse of the matrix C . (C^{-1} is the matrix such that $C^{-1}C$ is the identity matrix, a matrix with ones on the diagonal and zeros elsewhere. If C is a 1×1 matrix, $C^{-1} = 1/C$.) A consistent estimator of the variance-covariance matrix is given by the matrix V , obtained by inserting b for B in $I(B)$: $V = I^{-1}(b)$.

9.3.1 Global Test Statistics

Suppose we wish to test the null hypothesis $H_0 : B = B^0$. The likelihood ratio test statistic is

$$\begin{aligned} LR &= -2 \log(L \text{ at } H_0 / L \text{ at MLEs}) \\ &= -2[\log L(B^0) - \log L(b)]. \end{aligned} \quad (9.25)$$

The corresponding Wald test statistic, using the estimated observed information matrix, is

$$W = (b - B^0)' I(b) (b - B^0) = (b - B^0)' V^{-1} (b - B^0). \quad (9.26)$$

(A quadratic form $a'Va$ is a matrix generalization of a^2V .) Note that if the number of estimated parameters is $p = 1$, W reduces to $(b - B^0)^2/V$, which is the square of a z - or t -type statistic (estimate - hypothesized value divided by estimated standard deviation of estimate).

The score statistic for H_0 is

$$S = U'(B^0) I^{-1}(B^0) U(B^0). \quad (9.27)$$

Note that as before, S does not require solving for the MLE. For large samples, LR , W , and S have a χ^2 distribution with p d.f. under suitable conditions.

9.3.2 Testing a Subset of the Parameters

Let $B = \{B_1, B_2\}$ and suppose that we wish to test $H_0 : B_1 = B_1^0$. We are treating B_2 as a nuisance parameter. For example, we may want to test whether blood pressure and cholesterol are risk factors after adjusting for confounders age and sex. In that case B_1 is the pair of regression coefficients for blood pressure and cholesterol and B_2 is the pair of coefficients for age and sex. B_2 must be estimated to allow adjustment for age and sex, although B_2 is a nuisance parameter and is not of primary interest.

Let the number of parameters of interest be k so that B_1 is a vector of length k . Let the number of “nuisance” or “adjustment” parameters be q , the length of B_2 (note $k + q = p$).

Let b_2^* be the MLE of B_2 under the restriction that $B_1 = B_1^0$. Then the likelihood ratio statistic is

$$LR = -2[\log L \text{ at } H_0 - \log L \text{ at MLE}]. \quad (9.28)$$

Now $\log L$ at H_0 is more complex than before because H_0 involves an unknown nuisance parameter B_2 that must be estimated. $\log L$ at H_0 is the maximum of the likelihood function for any value of B_2 but subject to the condition that $B_1 = B_1^0$. Thus

$$LR = -2[\log L(B_1^0, b_2^*) - \log L(b)], \quad (9.29)$$

where as before b is the overall MLE of B . Note that LR requires maximizing two log likelihood functions. The first component of LR is a restricted maximum likelihood and the second component is the overall or unrestricted maximum.

LR is often computed by examining successively more complex models in a stepwise fashion and calculating the increment in likelihood ratio χ^2 in the overall model. The $LR \chi^2$ for testing $H_0 : B_2 = 0$ when B_1 is not in the model is

$$LR(H_0 : B_2 = 0 | B_1 = 0) = -2[\log L(0, 0) - \log L(0, b_2^*)]. \quad (9.30)$$

Here we are specifying that B_1 is not in the model by setting $B_1 = B_1^0 = 0$, and we are testing $H_0 : B_2 = 0$. (We are also ignoring nuisance parameters such as an intercept term in the test for $B_2 = 0$.)

The $LR \chi^2$ for testing $H_0 : B_1 = B_2 = 0$ is given by

$$LR(H_0 : B_1 = B_2 = 0) = -2[\log L(0, 0) - \log L(b)]. \quad (9.31)$$

Subtracting $LR \chi^2$ for the smaller model from that of the larger model yields

$$\begin{aligned} & -2[\log L(0, 0) - \log L(b)] - [-2[\log L(0, 0) - \log L(0, b_2^*)]] \\ = & \qquad \qquad \qquad -2[\log L(0, b_2^*) - \log L(b)], \end{aligned} \quad (9.32)$$

which is the same as above (letting $B_1^0 = 0$).

Table 9.1 Example tests

Variables (Parameters) in Model	$LR \chi^2$	Number of Parameters
Intercept, age	1000	2
Intercept, age, age ²	1010	3
Intercept, age, age ² , sex	1013	4

For example, suppose successively larger models yield the $LR \chi^2$ s in Table 9.1. The $LR \chi^2$ for testing for linearity in age (not adjusting for sex) against quadratic alternatives is $1010 - 1000 = 10$ with 1 d.f. The $LR \chi^2$ for testing the added information provided by sex, adjusting for a quadratic effect of age, is $1013 - 1010 = 3$ with 1 d.f. The $LR \chi^2$ for testing the joint importance of sex and the nonlinear (quadratic) effect of age is $1013 - 1000 = 13$ with 2 d.f.

To derive the Wald statistic for testing $H_0 : B_1 = B_1^0$ with B_2 being a nuisance parameter, let the MLE b be partitioned into $b = \{b_1, b_2\}$. We can likewise partition the estimated variance–covariance matrix V into

$$V = \begin{bmatrix} V_{11} & V_{12} \\ V'_{12} & V_{22} \end{bmatrix}. \quad (9.33)$$

The Wald statistic is

$$W = (b_1 - B_1^0)' V_{11}^{-1} (b_1 - B_1^0), \quad (9.34)$$

which when $k = 1$ reduces to (estimate – hypothesized value)²/ estimated variance, with the estimates adjusted for the parameters in B_2 .

The score statistic for testing $H_0 : B_1 = B_1^0$ does not require solving for the full set of unknown parameters. Only the MLEs of B_2 must be computed, under the restriction that $B_1 = B_1^0$. This restricted MLE is b_2^* from above. Let $U(B_1^0, b_2^*)$ denote the vector of first derivatives of $\log L$ with respect to all parameters in B , evaluated at the hypothesized parameter values B_1^0 for the first k parameters and at the restricted MLE b_2^* for the last q parameters. (Since the last q estimates are MLEs, the last q elements of U are zero, so the formulas that follow simplify.) Let $I(B_1^0, b_2^*)$ be the observed information matrix evaluated at the same values of B as is U . The score statistic for testing $H_0 : B_1 = B_1^0$ is

$$S = U'(B_1^0, b_2^*) I^{-1}(B_1^0, b_2^*) U(B_1^0, b_2^*). \quad (9.35)$$

Under suitable conditions, the distribution of LR, W , and S can be adequately approximated by a χ^2 distribution with k d.f.

9.3.3 Tests Based on Contrasts

Wald tests are also done by setting up a general linear contrast. $H_0 : CB = 0$ is tested by a Wald statistic of the form

$$W = (Cb)'(CVC')^{-1}(Cb), \quad (9.36)$$

where C is a contrast matrix that “picks off” the proper elements of B . The contrasts can be much more general by allowing elements of C to be other than zero and one. For the normal linear model, W is converted to an F -statistic by dividing by the rank r of C (normally the number of rows in C), yielding a statistic with an F -distribution with r numerator degrees of freedom.

Many interesting contrasts are tested by forming differences in predicted values. By forming more contrasts than are really needed, one can develop a surprisingly flexible approach to hypothesis testing using predicted values. This has the major advantage of not requiring the analyst to account for how the predictors are coded. Suppose that one wanted to assess the difference in two vectors of predicted values, $X_1b - X_2b = (X_1 - X_2)b = \Delta b$ to test $H_0 : \Delta B = 0$, where $\Delta = X_1 - X_2$. The covariance matrix for Δb is given by

$$\text{var}(\Delta b) = \Delta V \Delta'. \quad (9.37)$$

Let r be the rank of $\text{var}(\Delta b)$, i.e., the number of non-linearly-dependent (non-redundant) differences of predicted values of Δ . The value of r and the rows of Δ that are not redundant may easily be determined using the QR decomposition as done by the R function `qr`^b. The χ^2 statistic with r degrees of freedom (or F -statistic upon dividing the statistic by r) may be obtained by computing $\Delta^* V^* \Delta^{*'} where Δ^* is the subset of elements of Δ corresponding to non-redundant contrasts and V^* is the corresponding sub-matrix of V .$

The “difference in predictions” approach can be used to compare means in a 30 year old male with a 40 year old female^c. But the true utility of the approach is most obvious when the contrast involves multiple nonlinear terms for a single predictor, e.g., a spline function. To test for a difference in two curves, one can compare predictions at one predictor value against predictions at a series of values with at least one value that pertains to each basis function. Points can be placed between every pair of knots and beyond the outer knots, or just obtain predictions at 100 equally spaced X -values.

^b For example, in a 3-treatment comparison one could examine contrasts between treatments A and B, A and C, and B and C by obtaining predicted values for those treatments, even though only two differences are required.

^c The `rms` command could be `contrast(fit, list(sex='male',age=30), list(sex='female',age=40))` where all other predictors are set to medians or modes.

Suppose that there are three treatment groups (A, B, C) interacting with a cubic spline function of X . If one wants to test the multiple degree of freedom hypothesis that the profile for X is the same for treatment A and B vs. the alternative hypothesis that there is a difference between A and B for at least one value of X , one can compare predicted values at treatment A and a vector of X values against predicted values at treatment B and the same vector of X values. If the X relationship is linear, any two X values will suffice, and if X is quadratic, any three points will suffice. It would be difficult to test complex hypotheses involving only 2 of 3 treatments using other methods.

The `contrast` function in `rms` can estimate a wide variety of contrasts and make joint tests involving them, automatically computing the number of non-linearly-dependent contrasts as the test's degrees of freedom. See its help file for several examples.

9.3.4 Which Test Statistics to Use When

At this point, one may ask why three types of test statistics are needed. The answer lies in the statistical properties of the three tests as well as in computational expense in different situations. From the standpoint of statistical properties, LR is the best statistic, followed by S and W . The major statistical problem with W is that it is sensitive to problems in the estimated variance–covariance matrix in the full model. For some models, most notably the logistic regression model,²⁷⁸ the variance–covariance estimates can be too large as the effects in the model become very strong, resulting in values of W that are too small (or significance levels that are too large). W is also sensitive to the way the parameter appears in the model. For example, a test of $H_0 : \log \text{ odds ratio} = 0$ will yield a different value of W than will $H_0 : \text{odds ratio} = 1$.

Relative computational efficiency of the three types of tests is also an issue. Computation of LR and W requires estimating all p unknown parameters, and in addition LR requires re-estimating the last q parameters under that restriction that the first k parameters = B_1^0 . Therefore, when one is contemplating whether a set of parameters should be added to a model, the score test is the easiest test to carry out. For example, if one were interested in testing all two-way interactions among 4 predictors, the score test statistic for $H_0 : \text{“no interactions present”}$ could be computed without estimating the $4 \times 3/2 = 6$ interaction effects. S would also be appealing for testing linearity of effects in a model—the nonlinear spline terms could be tested for significance after adjusting for the linear effects (with estimation of only the linear effects). Only parameters for linear effects must be estimated to compute S , resulting in fewer numerical problems such as lack of convergence of the Newton–Raphson algorithm.

Table 9.2 Choice of test statistics

Type of Test	Recommended Test Statistic
Global association	LR (S for large no. parameters)
Partial association	W (LR or S if problem with W)
Lack of fit, 1 d.f.	W or S
Lack of fit, > 1 d.f.	S
Inclusion of additional predictors	S

The Wald tests are very easy to make after all the parameters in a model have been estimated. Wald tests are thus appealing in a multiple regression setup when one wants to test whether a given predictor or set of predictors is “significant.” A score test would require re-estimating the regression coefficients under the restriction that the parameters of interest equal zero.

Likelihood ratio tests are used often for testing the global hypothesis that no effects are significant, as the log likelihood evaluated at the MLEs is already available from fitting the model and the log likelihood evaluated at a “null model” (e.g., a model containing only an intercept) is often easy to compute. Likelihood ratio tests should also be used when the validity of a Wald test is in question as in the example cited above.

Table 9.2 summarizes recommendations for choice of test statistics for various situations.

9.3.5 Example: Binomial—Comparing Two Proportions

Suppose that a binary random variable Y_1 represents responses for population 1 and Y_2 represents responses for population 2. Let $P_i = \text{Prob}\{Y_i = 1\}$ and assume that a random sample has been drawn from each population with respective sample sizes n_1 and n_2 . The sample values are denoted by $Y_{i1}, \dots, Y_{in_i}, i = 1$ or 2. Let

$$s_1 = \sum_{j=1}^{n_1} Y_{1j} \quad s_2 = \sum_{j=1}^{n_2} Y_{2j}, \quad (9.38)$$

the respective observed number of “successes” in the two samples. Let us test the null hypothesis $H_0 : P_1 = P_2$ based on the two samples.

The likelihood function is

$$L = \prod_{i=1}^2 \prod_{j=1}^{n_i} P_i^{Y_{ij}} (1 - P_i)^{1 - Y_{ij}}$$

$$= \prod_{i=1}^2 P_i^{s_i} (1 - P_i)^{n_i - s_i} \quad (9.39)$$

$$\log L = \sum_{i=1}^2 \{s_i \log(P_i) + (n_i - s_i) \log(1 - P_i)\}. \quad (9.40)$$

Under H_0 , $P_1 = P_2 = P$, so

$$\log L(H_0) = s \log(P) + (n - s) \log(1 - P), \quad (9.41)$$

where $s = s_1 + s_2$, $n = n_1 + n_2$. The (restricted) MLE of this common P is $p = s/n$ and $\log L$ at this value is $s \log(p) + (n - s) \log(1 - p)$.

Since the original unrestricted log likelihood function contains two terms with separate parameters, the two parts may be maximized separately giving MLEs

$$p_1 = s_1/n_1 \quad \text{and} \quad p_2 = s_2/n_2. \quad (9.42)$$

$\log L$ evaluated at these (unrestricted) MLEs is

$$\begin{aligned} \log L &= s_1 \log(p_1) + (n_1 - s_1) \log(1 - p_1) \\ &\quad + s_2 \log(p_2) + (n_2 - s_2) \log(1 - p_2). \end{aligned} \quad (9.43)$$

The likelihood ratio statistic for testing $H_0 : P_1 = P_2$ is then

$$\begin{aligned} LR &= -2\{s \log(p) + (n - s) \log(1 - p) \\ &\quad - [s_1 \log(p_1) + (n_1 - s_1) \log(1 - p_1) \\ &\quad + s_2 \log(p_2) + (n_2 - s_2) \log(1 - p_2)]\}. \end{aligned} \quad (9.44)$$

This statistic for large enough n_1 and n_2 has a χ^2 distribution with 1 d.f. since the null hypothesis involves the estimation of one fewer parameter than does the unrestricted case. This LR statistic is the likelihood ratio χ^2 statistic for a 2×2 contingency table. It can be shown that the corresponding score statistic is equivalent to the Pearson χ^2 statistic. The better LR statistic can be used routinely over the Pearson χ^2 for testing hypotheses in contingency tables.

9.4 Iterative ML Estimation

In most cases, one cannot explicitly solve for MLEs but must use trial-and-error numerical methods to solve for parameter values B that maximize $\log L(B)$ or yield a score vector $U(B) = 0$. One of the fastest and most applicable methods for maximizing a function is the Newton–Raphson method, which is based on approximating $U(B)$ by a linear function of B in a small

region. A starting estimate b^0 of the MLE b is made. The linear approximation (a first-order Taylor series approximation)

$$U(b) = U(b^0) - I(b^0)(b - b^0) \quad (9.45)$$

is equated to 0 and solved by b yielding

$$b = b^0 + I^{-1}(b^0)U(b^0). \quad (9.46)$$

The process is continued in like fashion. At the i th step the next estimate is obtained from the previous estimate using the formula

$$b^{i+1} = b^i + I^{-1}(b^i)U(b^i). \quad (9.47)$$

If the log likelihood actually worsened at b^{i+1} , “step halving” is used; b^{i+1} is replaced with $(b^i + b^{i+1})/2$. Further step halving is done if the log likelihood still is worse than the log likelihood at b^i , after which the original iterative strategy is resumed. The Newton–Raphson iterations continue until the -2 log likelihood changes by only a small amount over the previous iteration (say .025). The reasoning behind this stopping rule is that estimates of B that change the -2 log likelihood by less than this amount do not affect statistical inference since -2 log likelihood is on the χ^2 scale.

3

9.5 Robust Estimation of the Covariance Matrix

The estimator for the covariance matrix of b found in Section 9.3 assumes that the model is correctly specified in terms of distribution, regression assumptions, and independence assumptions. The model may be incorrect in a variety of ways such as non-independence (e.g., repeated measurements within subjects), lack of fit (e.g., omitted covariable, incorrect covariable transformation, omitted interaction), and distributional (e.g., Y has a Γ distribution instead of a normal distribution). Variances and covariances, and hence confidence intervals and Wald tests, will be incorrect when these assumptions are violated.

For the case in which the observations are independent and identically distributed but other assumptions are possibly violated, Huber³¹² provided a covariance matrix estimator that is consistent. His “sandwich” estimator is given by

$$H = I^{-1}(b) \left[\sum_{i=1}^n U_i U_i' \right] I^{-1}(b), \quad (9.48)$$

where $I(b)$ is the observed information matrix (Equation 9.24) and U_i is the vector of derivatives, with respect to all parameters, of the log likelihood component for the i th observation (assuming the log likelihood can be partitioned into per-observation contributions). For the normal multiple linear regression case, H was derived by White.⁶⁵⁹

$$(X'X)^{-1} \left[\sum_{i=1}^n (Y_i - X_i b)^2 X_i X_i' \right] (X'X)^{-1}, \quad (9.49)$$

where X is the design matrix (including an intercept if appropriate) and X_i is the vector of predictors (including an intercept) for the i th observation. This covariance estimator allows for any pattern of variances of $Y|X$ across observations. Note that even though H improves the bias of the covariance matrix of b , it may actually have larger mean squared error than the ordinary estimate in some cases due to increased variance.^{164, 529}

4

When observations are dependent within clusters, and the number of observations within clusters is very small in comparison to the total sample size, a simple adjustment to Equation 9.48 can be used to derive appropriate covariance matrix estimates (see Lin [407, p. 2237], Rogers,⁵²⁹ and Lee et al. [393, Eq. 5.1, p. 246]). One merely accumulates sums of elements of U within clusters before computing cross-product terms:

$$H_c = I^{-1}(b) \left[\sum_{i=1}^c \left\{ \left(\sum_{j=1}^{n_i} U_{ij} \right) \left(\sum_{j=1}^{n_i} U_{ij} \right)' \right\} \right] I^{-1}(b), \quad (9.50)$$

where c is the number of clusters, n_i is the number of observations in the i th cluster, U_{ij} is the contribution of the j th observation within the i th cluster to the score vector, and $I(b)$ is computed as before ignoring clusters. For a model such as the Cox model which has no per-observation score contributions, special score residuals^{393, 407, 410, 605} are used for U .

Bootstrapping can also be used to derive robust covariance matrix estimates^{177, 178} in many cases, especially if covariances of b that are not conditional on X are appropriate. One merely generates approximately 200 samples with replacement from the original dataset, computes 200 sets of parameter estimates, and computes the sample covariance matrix of these parameter estimates. Sampling with replacement from entire clusters can be used to derive variance estimates in the presence of intracluster correlation.¹⁸⁸ Bootstrap estimates of the conditional variance–covariance matrix given X are harder to obtain and depend on the model assumptions being satisfied. The simpler unconditional estimates may be more appropriate for many non-experimental studies where one may desire to “penalize” for the X being random variables. It is interesting that these unconditional estimates may be very difficult to obtain parametrically, since a multivariate distribution may need to be assumed for X .

5

The previous discussion addresses the use of a “working independence model” with clustered data. Here one estimates regression coefficients assuming independence of all records (observations). Then a sandwich or bootstrap method is used to increase standard errors to reflect some redundancy in the correlated observations. The parameter estimates will often be consistent estimates of the true parameter values, but they may be inefficient for certain cluster or correlation structures.

6

The `rms` package's `robcov` function computes the Huber robust covariance matrix estimator, and the `bootcov` function computes the bootstrap covariance estimator. Both of these functions allow for clustering.

9.6 Wald, Score, and Likelihood-Based Confidence Intervals

A $1 - \alpha$ confidence interval for a parameter β_i is the set of all values β_i^0 that if hypothesized would be accepted in a test of $H_0 : \beta_i = \beta_i^0$ at the α level. What test should form the basis for the confidence interval? The Wald test is most frequently used because of its simplicity. A two-sided $1 - \alpha$ confidence interval is $b_i \pm z_{1-\alpha/2}s$, where z is the critical value from the normal distribution and s is the estimated standard error of the parameter estimate b_i .^d The problem with s discussed in Section 9.3.4 points out that Wald statistics may not always be a good basis. Wald-based confidence intervals are also symmetric even though the coverage probability may not be.¹⁶⁰ Score- and LR-based confidence limits have definite advantages. When Wald-type confidence intervals are appropriate, the analyst may consider insertion of robust covariance estimates (Section 9.5) into the confidence interval formulas (note that adjustments for heterogeneity and correlated observations are not available for score and LR statistics).

7

Wald- (asymptotic normality) based statistics are convenient for deriving confidence intervals for linear or more complex combinations of the model's parameters. As in Equation 9.36, the variance-covariance matrix of Cb , where C is an appropriate matrix and b is the vector of parameter estimates, is $CV C'$, where V is the variance matrix of b . In regression models we commonly substitute a vector of predictors (and optional intercept) for C to obtain the variance of the linear predictor Xb as

$$\text{var}(Xb) = XVX'. \quad (9.51)$$

See Section 9.3.3 for related information.

^d This is the basis for confidence limits computed by the R `rms` package's `Predict`, `summary`, and `contrast` functions. When the `robcov` function has been used to replace the information-matrix-based covariance matrix with a Huber robust covariance estimate with an optional cluster sampling correction, the functions are using a "robust" Wald statistic basis. When the `bootcov` function has been used to replace the model fit's covariance matrix with a bootstrap unconditional covariance matrix estimate, the two functions are computing confidence limits based on a normal distribution but using more nonparametric covariance estimates.

9.6.1 Simultaneous Wald Confidence Regions

The confidence intervals just discussed are *pointwise* confidence intervals. For OLS regression there are methods for computing confidence intervals with exact *simultaneous* confidence coverage for multiple estimates³⁷⁴. There are approximate methods for simultaneous confidence limits for all models for which the vector of estimates b is approximately multivariately normally distributed. The method of Hothorn et al.³⁰⁷ is quite general; in their R package `multcomp`'s `glht` function, the user can specify any contrast matrix over which the individual confidence limits will be simultaneous. A special case of a contrast matrix is the design matrix X itself, resulting in simultaneous confidence bands for any number of predicted values. An example is shown in Figure 9.5. See Section 9.3.3 for a good use for simultaneous contrasts.

9.7 Bootstrap Confidence Regions

A more nonparametric method for computing confidence intervals for functions of the vector of parameters B can be based on bootstrap percentile confidence limits. For each sample with replacement from the original dataset, one computes the MLE of B , b , and then the quantity of interest $g(b)$. Then the g s are sorted and the desired quantiles are computed. At least 1000 bootstrap samples will be needed for accurate assessment of outer confidence limits. This method is suitable for obtaining pointwise confidence bands for a nonlinear regression function, say, the relationship between age and the log odds of disease. At each of 100 age values the predicted logits are computed for each bootstrap sample. Then separately for each age point the 0.025 and 0.975 quantiles of 1000 estimates of the logit are computed to derive a 0.95 confidence band. Other more complex bootstrap schemes will achieve somewhat greater accuracy of confidence interval coverage,¹⁷⁸ and as described in Section 9.5 one can use variations on the basic bootstrap in which the predictors are considered fixed and/or cluster sampling is taken into account. The R function `bootcov` in the `rms` package bootstraps model fits to obtain unconditional (with respect to predictors) bootstrap distributions with or without cluster sampling. `bootcov` stores the matrix of bootstrap regression coefficients so that the bootstrapped quantities of interest can be computed in one sweep of the coefficient matrix once bootstrapping is completed.

For many regression models, the `rms` package's `Predict`, `summary`, and `contrast` functions make it easy to compute pointwise bootstrap confidence intervals in a variety of contexts. As an example, consider 200 simulated x values from a log-normal distribution and simulate binary y from a true population binary logistic model given by

8

9

$$\text{Prob}(Y = 1|X = x) = \frac{1}{1 + \exp[-(1 + x/2)]}. \quad (9.52)$$

Not knowing the true model, a quadratic logistic model is fitted. The R code needed to generate the data and fit the model is given below.

```
require(rms)
```

```
n <- 200
set.seed(15)
x1 <- rnorm(n)
logit <- x1/2
y <- ifelse(runif(n) <= plogis(logit), 1, 0)
dd <- datadist(x1); options(datadist='dd')
f <- lrm(y ~ pol(x1,2), x=TRUE, y=TRUE)
print(f, latex=TRUE)
```

Logistic Regression Model

```
lrm(formula = y ~ pol(x1, 2), x = TRUE, y = TRUE)
```

	Model Likelihood Ratio Test	Discrimination Indexes	Rank Discrim. Indexes
Obs	200 LR χ^2 16.37	R^2 0.105	C 0.642
0	97 d.f.	g 0.680	D_{xy} 0.285
1	103 Pr(> χ^2) 0.0003	g_r 1.973	γ 0.286
max $ \frac{\partial \log L}{\partial \beta} $	3×10^{-9}	g_p 0.156	τ_a 0.143
		Brier 0.231	

	Coef	S.E.	Wald Z	Pr(> Z)
Intercept	-0.0842	0.1823	-0.46	0.6441
x1	0.5902	0.1580	3.74	0.0002
x1 ²	0.1557	0.1136	1.37	0.1708

```
latex(anova(f), file='', table.env=FALSE)
```

	χ^2	d.f.	P
x1	13.99	2	0.0009
<i>Nonlinear</i>	1.88	1	0.1708
TOTAL	13.99	2	0.0009

The `bootcov` function is used to draw 1000 resamples to obtain bootstrap estimates of the covariance matrix of the regression coefficients as well as to save the 1000×3 matrix of regression coefficients. Then, because individual regression coefficients for x do not tell us much, we summarize the

x -effect by computing the effect (on the logit scale) of increasing x from 1 to 5. We first compute bootstrap nonparametric percentile confidence intervals the long way. The 1000 bootstrap estimates of the log odds ratio are computed easily using a single matrix multiplication with the difference in predictions approach, multiplying the difference in two design matrices, and we obtain the bootstrap estimate of the standard error of the log odds ratio by computing the sample standard deviation of the 1000 values^e. Bootstrap percentile confidence limits are just sample quantiles from the bootstrapped log odds ratios.

```
# Get 2-row design matrix for obtaining predicted values
# for x = 1 and 5
X ← cbind(Intercept=1,
          predict(f, data.frame(x1=c(1,5)), type='x'))
Xdif ← X[2,,drop=FALSE] - X[1,,drop=FALSE]
Xdif
```

```
Intercept pol(x1, 2)x1 pol(x1, 2)x1^2
2          0          4          24
```

```
b ← bootcov(f, B=1000)
boot.log.odds.ratio ← b$boot.Coeff %*% t(Xdif)
sd(boot.log.odds.ratio)
```

```
[1] 2.752103
```

```
# This is the same as from summary(b, x=c(1,5)) as summary
# uses the bootstrap covariance matrix
summary(b, x1=c(1,5))[1, 'S.E. ']
```

```
[1] 2.752103
```

```
# Compare this s.d. with one from information matrix
summary(f, x1=c(1,5))[1, 'S.E. ']
```

```
[1] 2.988373
```

```
# Compute percentiles of bootstrap odds ratio
exp(quantile(boot.log.odds.ratio, c(.025, .975)))
```

```
2.5%          97.5%
2.795032e+00 2.067146e+05
```

```
# Automatic:
summary(b, x1=c(1,5))[' Odds Ratio ',]
```

^e As indicated below, this standard deviation can also be obtained by using the `summary` function on the object returned by `bootcov`, as `bootcov` returns a fit object like one from `lrm` except with the bootstrap covariance matrix substituted for the information-based one.

	Low	High	Diff.	Effect	S.E.
	1.000000e+00	5.000000e+00	4.000000e+00	4.443932e+02	NA
	Lower 0.95	Upper 0.95	Type		
	2.795032e+00	2.067146e+05	2.000000e+00		

```
print(contrast(b, list(x1=5), list(x1=1), fun=exp))
```

	Contrast	S.E.	Lower	Upper	Z	Pr(> z)
11	6.09671	2.752103	1.027843	12.23909	2.22	0.0267

Confidence intervals are 0.95 bootstrap nonparametric percentile intervals

```
# Figure 9.4
hist(boot.log.odds.ratio, nclass=100, xlab='log(OR)',
     main='')
```

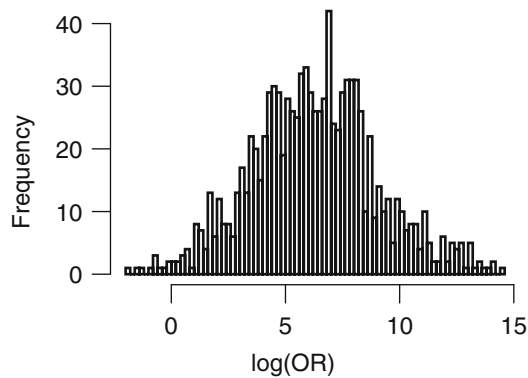


Fig. 9.4 Distribution of 1000 bootstrap $x=1:5$ log odds ratios

Figure 9.4 shows the distribution of log odds ratios.

Now consider confidence bands for the true log odds that $y = 1$, across a sequence of x values. The `Predict` function automatically calculates point-by-point bootstrap percentiles, basic bootstrap, or BCa²⁰³ confidence limits when the fit has passed through `bootcov`. Simultaneous Wald-based confidence intervals³⁰⁷ and Wald intervals substituting the bootstrap covariance matrix estimator are added to the plot when `Predict` calls the `multcomp` package (Figure 9.5).

```
x1s <- seq(0, 5, length=100)
pwald <- Predict(f, x1=x1s)
psand <- Predict(robcov(f), x1=x1s)
pbootcov <- Predict(b, x1=x1s, usebootcoef=FALSE)
pbootnp <- Predict(b, x1=x1s)
pbootbca <- Predict(b, x1=x1s, boot.type='bca')
pbootbas <- Predict(b, x1=x1s, boot.type='basic')
psimult <- Predict(b, x1=x1s, conf.type='simultaneous')
```

```

z <- rbind('Boot percentile' = pbootnp,
          'Robust sandwich'  = psand,
          'Boot BCa'        = pbootbca,
          'Boot covariance+Wald' = pbootcov,
          Wald              = pwald,
          'Boot basic'      = pbootbas,
          Simultaneous      = psimult)

z$class <- ifelse(z$.set. %in% c('Boot percentile', 'Boot bca',
                               'Boot basic'), 'Other', 'Wald')
ggplot(z, groups=c('.set.', 'class'),
       conf='line', ylim=c(-1, 9), legend.label=FALSE)

```

See Problems at chapter's end for a worrisome investigation of bootstrap confidence interval coverage using simulation. It appears that when the model's log odds distribution is not symmetric and includes very high or very low probabilities, neither the bootstrap percentile nor the bootstrap BCa intervals have good coverage, while the basic bootstrap and ordinary Wald intervals are fairly accurate^f. It is difficult in general to know when to trust the bootstrap for logistic and perhaps other models when computing confidence intervals, and the simulation problem suggests that the basic bootstrap should be used more frequently. Similarly, the distribution of bootstrap effect estimates can be suspect. Asymmetry in this distribution does not imply that the true sampling distribution is asymmetric or that the percentile intervals are preferred.

9.8 Further Use of the Log Likelihood

9.8.1 Rating Two Models, Penalizing for Complexity

Suppose that from a single sample two competing models were developed. Let the respective $-2 \log$ likelihoods for these models be denoted by L_1 and L_2 , and let p_1 and p_2 denote the number of parameters estimated in each model. Suppose that $L_1 < L_2$. It may be tempting to rate model one as the “best” fitting or “best” predicting model. That model may provide a better fit for the data at hand, but if it required many more parameters to be estimated, it may not be better “for the money.” If both models were applied to a new sample, model one's overfitting of the original dataset may actually result in a worse fit on the new dataset.

^f Limited simulations using the conditional bootstrap and Firth's penalized likelihood²⁸¹ did not show significant improvement in confidence interval coverage.

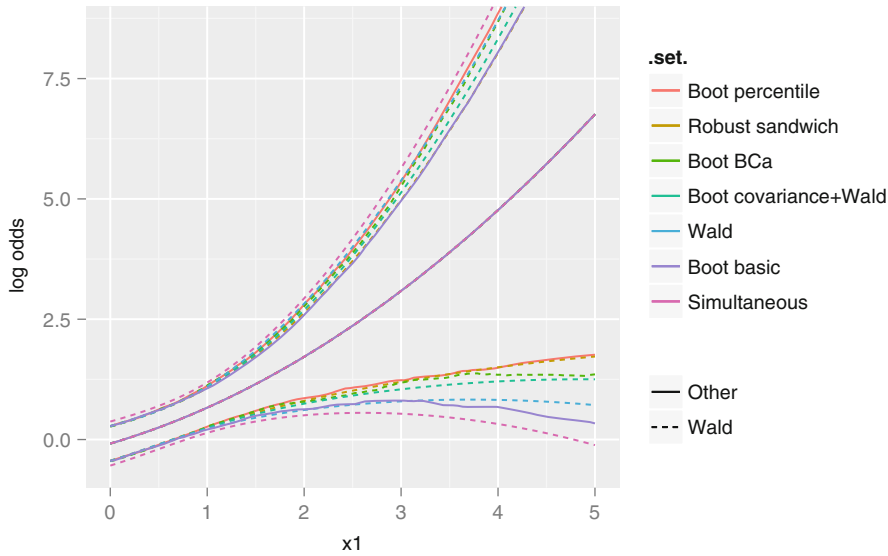


Fig. 9.5 Predicted log odds and confidence bands for seven types of confidence intervals. Seven categories are ordered top to bottom corresponding to order of lower confidence bands at $x_1=5$. Dotted lines are for Wald-type methods that yield symmetric confidence intervals and assume normality of point estimators.

Akaike’s information criterion (AIC^{33,359,633}) provides a method for penalizing the log likelihood achieved by a given model for its complexity to obtain a more unbiased assessment of the model’s worth. The penalty is to subtract the number of parameters estimated from the log likelihood, or equivalently to add twice the number of parameters to the $-2 \log$ likelihood. The penalized log likelihood is analogous to Mallows’ C_p in ordinary multiple regression. AIC would choose the model by comparing $L_1 + 2p_1$ to $L_2 + 2p_2$ and picking the model with the lower value. We often use AIC in “adjusted χ^2 ” form:

10

$$\text{AIC} = \text{LR } \chi^2 - 2p. \tag{9.53}$$

Breiman [66, Section 1.3] and Chatfield [100, Section 4] discuss the fallacy of AIC and C_p for selecting from a series of non-prespecified models.

11

9.8.2 Testing Whether One Model Is Better than Another

One way to test whether one model (A) is better than another (B) is to embed both models in a more general model ($A + B$). Then a $LR \chi^2$ test

can be done to test whether A is better than B by changing the hypothesis to test whether A adds predictive information to B ($H_0 : A + B > B$) and whether B adds information to A ($H_0 : A + B > A$). The approach of testing $A > B$ via testing $A + B > B$ and $A + B > A$ is especially useful for selecting from competing predictors such as a multivariable model and a subjective assessor.^{131, 264, 395, 669}

Note that $LR \chi^2$ for $H_0 : A + B > B$ minus $LR \chi^2$ for $H_0 : A + B > A$ equals $LR \chi^2$ for $H_0 : A$ has no predictive information minus $LR \chi^2$ for $H_0 : B$ has no predictive information,⁶⁶⁵ the difference in $LR \chi^2$ for testing each model (set of variables) separately. This gives further support to the use of two separately computed Akaike's information criteria for rating the two sets of variables.

12

See Section 9.8.4 for an example.

9.8.3 Unitless Index of Predictive Ability

The global likelihood ratio test for regression is useful for determining whether any predictor is associated with the response. If the sample is large enough, even weak associations can be "statistically significant." Even though a likelihood ratio test does not shed light on a model's predictive strength, the log likelihood (L.L.) can still be useful here. Consider the following L.L.s:

Best (lowest) possible -2 L.L.:

$L^* = -2$ L.L. for a hypothetical model that perfectly predicts the outcome.

-2 L.L. achieved:

$L = -2$ L.L. for the fitted model.

Worst -2 L.L.:

$L^0 = -2$ L.L. for a model that has no predictive information.

The last -2 L.L., for a "no information" model, is the -2 L.L. under the null hypothesis that all regression coefficients except for intercepts are zero. A "no information" model often contains only an intercept and some distributional parameters (a variance, for example).

13

The quantity $L^0 - L$ is LR , the log likelihood ratio statistic for testing the global null hypothesis that no predictors are related to the response. It is also the -2 log likelihood "explained" by the model. The best (lowest) -2 L.L. is L^* , so the amount of L.L. that is capable of being explained by the model is $L^0 - L^*$. The fraction of -2 L.L. explained that was capable of being explained is

$$(L^0 - L)/(L^0 - L^*) = LR/(L^0 - L^*). \quad (9.54)$$

The fraction of log likelihood explained is analogous to R^2 in an ordinary linear model, although Korn and Simon^{365,366} provide a much more precise notion.

Akaike's information criterion can be used to penalize this measure of association for the number of parameters estimated (p , say) to transform this unitless measure of association into a quantity that is analogous to the adjusted R^2 or Mallows' C_p in ordinary linear regression. We let R denote the square root of such a penalized fraction of log likelihood explained. R is defined by

$$R^2 = (LR - 2p)/(L_0 - L^*). \quad (9.55)$$

The R index can be used to assess how well the model compares with a "perfect" model, as well as to judge whether a more complex model has predictive strength that justifies its additional parameters. Had p been used in Equation 9.55 rather than $2p$, R^2 is negative if the log likelihood explained is less than what one would expect by chance. R will be the square root of $1 - 2p/(L_0 - L^*)$ if the model perfectly predicts the response. This upper limit will be near one if the sample size is large.

Partial R indexes can also be defined by substituting the -2 L.L. explained for a given factor in place of that for the entire model, LR . The "penalty factor" p becomes one. This index R_{partial} is defined by

$$R_{\text{partial}}^2 = (LR_{\text{partial}} - 2)/(L_0 - L^*), \quad (9.56)$$

which is the (penalized) fraction of -2 log likelihood explained by the predictor. Here LR_{partial} is the log likelihood ratio statistic for testing whether the predictor is associated with the response, after adjustment for the other predictors. Since such likelihood ratio statistics are tedious to compute, the 1 d.f. Wald χ^2 can be substituted for the LR statistic (keeping in mind that difficulties with the Wald statistic can arise).

Liu and Dyer⁴²⁴ and Cox and Wermuth¹³⁶ point out difficulties with the R^2 measure for binary logistic models. Cox and Snell¹³⁵ and Magee⁴³² used other analogies to derive other R^2 measures that may have better properties. For a sample of size n and a Wald statistic for testing overall association, they defined

$$\begin{aligned} R_W^2 &= \frac{W}{n + W} \\ R_{LR}^2 &= 1 - \exp(-LR/n) \\ &= 1 - \lambda^{2/n}, \end{aligned} \quad (9.57)$$

where λ is the null model likelihood divided by the fitted model likelihood. In the case of ordinary least squares with normality both of the above indexes are equal to the traditional R^2 . R_{LR}^2 is equivalent to Maddala's index [431, Eq. 2.44]. Cragg and Uhler¹³⁷ and Nagelkerke⁴⁷¹ suggested dividing R_{LR}^2 by

its maximum attainable value

$$R_{\max}^2 = 1 - \exp(-L^0/n) \quad (9.58)$$

to derive R_N^2 which ranges from 0 to 1. This is the form of the R^2 index we use throughout.

For penalizing for overfitting, see Verweij and van Houwelingen⁶⁴⁰ for an overfitting-corrected R^2 that uses a cross-validated likelihood.

14

9.8.4 Unitless Index of Adequacy of a Subset of Predictors

Log likelihoods are also useful for quantifying the predictive information contained in a subset of the predictors compared with the information contained in the entire set of predictors.²⁶⁴ Let LR again denote the -2 log likelihood ratio statistic for testing the joint significance of the full set of predictors. Let LR^s denote the -2 log likelihood ratio statistic for testing the importance of the subset of predictors of interest, excluding the other predictors from the model. A measure of adequacy of the subset for predicting the response is given by

$$A = LR^s / LR. \quad (9.59)$$

A is then the proportion of log likelihood explained by the subset with reference to the log likelihood explained by the entire set. When $A = 1$, the subset contains all the predictive information found in the whole set of predictors; that is, the subset is adequate by itself and the additional predictors contain no independent information. When $A = 0$, the subset contains no predictive information by itself.

Califf et al.⁸⁹ used the A index to quantify the adequacy (with respect to prognosis) of two competing sets of predictors that each describe the extent of coronary artery disease. The response variable was time until cardiovascular death and the statistical model used was the Cox¹³² proportional hazards model. Some of their results are reproduced in Table 9.3. A chance-corrected adequacy measure could be derived by squaring the ratio of the R -index for the subset to the R -index for the whole set. A formal test of superiority of $X_1 =$ maximum % stenosis over $X_2 =$ jeopardy score can be obtained by testing whether X_1 adds to X_2 ($LR \chi^2 = 57.5 - 42.6 = 14.9$) and whether X_2 adds to X_1 ($LR \chi^2 = 57.5 - 51.8 = 5.7$). X_1 adds more to X_2 (14.9) than X_2 adds to X_1 (5.7). The difference $14.9 - 5.7 = 9.2$ equals the difference in single factor χ^2 ($51.8 - 42.6$)⁶⁶⁵.

15

Table 9.3 Completing prognostic markers

Predictors Used	LR	χ^2 Adequacy
Coronary jeopardy score	42.6	0.74
Maximum % stenosis in each artery	51.8	0.90
Combined	57.5	1.00

9.9 Weighted Maximum Likelihood Estimation

It is commonly the case that data elements represent combinations of values that pertain to a set of individuals. This occurs, for example, when unique combinations of X and Y are determined from a massive dataset, along with the frequency of occurrence of each combination, for the purpose of reducing the size of the dataset to analyze. For the i th combination we have a *case weight* w_i that is a positive integer representing a frequency. Assuming that observations represented by combination i are independent, the likelihood needed to represent all w_i observations is computed simply by multiplying all of the likelihood elements (each having value L_i), yielding a total likelihood contribution for combination i of $L_i^{w_i}$ or a log likelihood contribution of $w_i \log L_i$. To obtain a likelihood for the entire dataset one computes the product over all combinations. The total log likelihood is $\sum w_i \log L_i$. As an example, the weighted likelihood that would be used to fit a weighted logistic regression model is given by

$$L = \prod_{i=1}^n P_i^{w_i Y_i} (1 - P_i)^{w_i (1 - Y_i)}, \quad (9.60)$$

where there are n combinations, $\sum_{i=1}^n w_i > n$, and P_i is $\text{Prob}[Y_i = 1 | X_i]$ as dictated by the model. Note that in general the correct likelihood function cannot be obtained by weighting the data and using an unweighted likelihood.

By a small leap one can obtain weighted maximum likelihood estimates from the above method even if the weights do not represent frequencies or even integers, as long as the weights are non-negative. Non-frequency weights are commonly used in sample surveys to adjust estimates back to better represent a target population when some types of subjects have been over-sampled from that population. Analysts should beware of possible losses in efficiency when obtaining weighted estimates in sample surveys.^{363, 364} Making the regression estimates conditional on sampling strata by including strata as covariables may be preferable to re-weighting the strata. If weighted estimates must be obtained, the weighted likelihood function is generally valid for obtaining properly weighted parameter estimates. However, the variance-covariance matrix obtained by inverting the information matrix from the weighted likelihood will not be correct in general. For one thing, the sum of the weights may be far from the number of subjects in the sample. A rough

approximation to the variance–covariance matrix may be obtained by first multiplying each weight by $n/\sum w_i$ and then computing the weighted information matrix, where n is the number of actual subjects in the sample.

16

9.10 Penalized Maximum Likelihood Estimation

Maximizing the log likelihood provides the best fit to the dataset at hand, but this can also result in fitting noise in the data. For example, a categorical predictor with 20 levels can produce extreme estimates for some of the 19 regression parameters, especially for the small cells (see Section 4.5). A shrinkage approach will often result in regression coefficient estimates that while biased are lower in mean squared error and hence are more likely to be close to the true unknown parameter values. Ridge regression is one approach to shrinkage, but a more general and better developed approach is penalized maximum likelihood estimation,^{237, 388, 639, 641} which is really a special case of Bayesian modeling with a Gaussian prior. Letting L denote the usual likelihood function and λ be a penalty factor, we maximize the penalized log likelihood given by

17

$$\log L - \frac{1}{2}\lambda \sum_{i=1}^p (s_i\beta_i)^2, \quad (9.61)$$

where s_1, s_2, \dots, s_p are scale factors chosen to make $s_i\beta_i$ unitless. Most authors standardize the data first and do not have scale factors in the equation, but Equation 9.61 has the advantage of allowing estimation of β on the original scale of the data. The usual methods (e.g., Newton–Raphson) are used to maximize 9.61.

The choice of the scaling constants has received far too little attention in the ridge regression and penalized MLE literature. It is common to use the standard deviation of each column of the design matrix to scale the corresponding parameter. For models containing nothing but continuous variables that enter the regression linearly, this is usually a reasonable approach. For continuous variables represented with multiple terms (one of which is linear), it is not always reasonable to scale each nonlinear term with its own standard deviation. For dummy variables, scaling using the standard deviation ($\sqrt{d(1-d)}$, where d is the mean of the dummy variable, i.e., the fraction of observations in that cell) is problematic since this will result in high prevalence cells getting more shrinkage than low prevalence ones because the high prevalence cells will dominate the penalty function.

18

An advantage of the formulation in Equation 9.61 is that one can assign scale constants of zero for parameters for which no shrinkage is desired.^{237, 639} For example, one may have prior beliefs that a linear additive model will fit the data. In that case, nonlinear and non-additive terms may be penalized.

For a categorical predictor having c levels, users of ridge regression often do not recognize that the amount of shrinkage and the predicted values from the fitted model depend on how the design matrix is coded. For example, one will get different predictions depending on which cell is chosen as the reference cell when constructing dummy variables. The setup in Equation 9.61 has the same problem. For example, if for a three-category factor we use category 1 as the reference cell and have parameters β_2 and β_3 , the unscaled penalty function is $\beta_2^2 + \beta_3^2$. If category 3 were used as the reference cell instead, the penalty would be $\beta_3^2 + (\beta_2 - \beta_3)^2$. To get around this problem, Verweij and van Houwelingen⁶³⁹ proposed using the penalty function $\sum_i^c (\beta_i - \bar{\beta})^2$, where $\bar{\beta}$ is the mean of all c β s. This causes shrinkage of all parameters toward the mean parameter value. Letting the first category be the reference cell, we use $c - 1$ dummy variables and define $\beta_1 \equiv 0$. For the case $c = 3$ the sum of squares is $2[\beta_2^2 + \beta_3^2 - \beta_2\beta_3]/3$. For $c = 2$ the penalty is $\beta_2^2/2$. If no scale constant is used, this is the same as scaling β_2 with $\sqrt{2} \times$ the standard deviation of a binary dummy variable with prevalence of 0.5.

The sum of squares can be written in matrix form as $[\beta_2, \dots, \beta_c]'(A - B)[\beta_2, \dots, \beta_c]$, where A is a $c - 1 \times c - 1$ identity matrix and B is a $c - 1 \times c - 1$ matrix all of whose elements are $\frac{1}{c}$.

For general penalty functions such as that just described, the penalized log likelihood can be generalized to

$$\log L - \frac{1}{2}\lambda\beta'P\beta. \quad (9.62)$$

For purposes of using the Newton–Raphson procedure, the first derivative of the penalty function with respect to β is $-\lambda P\beta$, and the negative of the second derivative is λP .

Another problem in penalized estimation is how the choice of λ is made. Many authors use cross-validation. A limited number of simulation studies in binary logistic regression modeling has shown that for each λ being considered, at least 10-fold cross-validation must be done so as to obtain a reasonable estimate of predictive accuracy. Even then, a smoother²⁰⁷ (“super smoother”) must be used on the $(\lambda, \text{accuracy})$ pairs to allow location of the optimum value unless one is careful in choosing the initial sub-samples and uses these same splits throughout. Simulation studies have shown that a modified AIC is not only much quicker to compute (since it requires no cross-validation) but performs better at finding a good value of λ (see below).

For a given λ , the effective number of parameters being estimated is reduced because of shrinkage. Gray [237, Eq. 2.9] and others estimate the effective degrees of freedom by computing the expected value of a global Wald statistic for testing association, when the null hypothesis of no association is true. The d.f. is equal to

$$\text{trace}[I(\hat{\beta}^P)V(\hat{\beta}^P)], \quad (9.63)$$

where $\hat{\beta}^P$ is the penalized MLE (the parameters that maximize Equation 9.61), I is the information matrix computed from ignoring the penalty function, and V is the covariance matrix computed by inverting the information matrix that included the second derivatives with respect to β in the penalty function.

22

Gray [237, Eq. 2.6] states that a better estimate of the variance–covariance matrix for $\hat{\beta}^P$ than $V(\hat{\beta}^P)$ is

$$V^* = V(\hat{\beta}^P)I(\hat{\beta}^P)V(\hat{\beta}^P). \quad (9.64)$$

Therneau (personal communication, 2000) has found in a limited number of simulation studies that V^* underestimates the true variances, and that a better estimate of the variance–covariance matrix is simply $V(\hat{\beta}^P)$, assuming that the model is correctly specified. This is the covariance matrix used by default in the `rms` package (the user can request that the sandwich estimator be used instead) and is in fact the one Gray used for Wald tests.

Penalization will bias estimates of β , so hypothesis tests and confidence intervals using $\hat{\beta}^P$ may not have a simple interpretation. The same problem arises in score and likelihood ratio tests. So far, penalization is better understood in pure prediction mode unless Bayesian methods are used.

Equation 9.63 can be used to derive a modified AIC (see [639, Eq. 6] and [641, Eq. 7]) on the model χ^2 scale:

$$\text{LR } \chi^2 - 2 \times \text{effective d.f.}, \quad (9.65)$$

where $\text{LR } \chi^2$ is the likelihood ratio χ^2 for the penalized model, but ignoring the penalty function. If a variety of λ are tried and one plots the (λ, AIC) pairs, the λ that maximizes AIC will often be a good choice, that is, it is likely to be near the value of λ that maximizes predictive accuracy on a future dataset[§].

Note that if one does penalized maximum likelihood estimation where a set of variables being penalized has a negative value for the unpenalized $\chi^2 - 2 \times \text{d.f.}$, the value of λ that will optimize the overall model AIC will be ∞ .

As an example, consider some simulated data ($n = 100$) with one predictor in which the true model is $Y = X_1 + \epsilon$, where ϵ has a standard normal distribution and so does X_1 . We use a series of penalties (found by trial and error) that give rise to sensible effective d.f., and fit penalized restricted cubic spline functions with five knots. We penalize two ways: all terms in the model including the coefficient of X_1 , which in reality needs no penalty; and only the nonlinear terms. The following R program, in conjunction with the `rms` package, does the job.

[§] Several examples from simulated datasets have shown that using BIC to choose a penalty results in far too much shrinkage.

```

set.seed(191)
x1 ← rnorm(100)
y ← x1 + rnorm(100)
pens ← df ← aic ← c(0,.07,.5,2,6,15,60)
all ← nl ← list()

for(penalize in 1:2) {
  for(i in 1:length(pens)) {
    f ← ols(y ~ rcs(x1,5), penalty=
            list(simple=if(penalize==1)pens[i] else 0,
                  nonlinear=pens[i]))
    df[i] ← f$stats['d.f.']
    aic[i] ← AIC(f)
    nam ← paste(if(penalize == 1) 'all' else 'nl',
                ' penalty:', pens[i], sep='')
    nam ← as.character(pens[i])
    p ← Predict(f, x1=seq(-2.5, 2.5, length=100),
                conf.int=FALSE)
    if(penalize == 1) all[[nam]] ← p else nl[[nam]] ← p
  }
  print(rbind(df=df, aic=aic))
}

```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]
df	4.0000	3.213591	2.706069	2.30273	2.029282	1.822758
aic	270.6653	269.154045	268.222855	267.56594	267.288988	267.552915
		[,7]				
df	1.513609					
aic	270.805033					
	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]
df	4.0000	3.219149	2.728126	2.344807	2.109741	1.960863
aic	270.6653	269.167108	268.287933	267.718681	267.441197	267.347475
		[,7]				
df	1.684421					
aic	267.892073					

```

all ← do.call('rbind', all); all$type ← 'Penalize All'
nl ← do.call('rbind', nl); nl$type ← 'Penalize Nonlinear'
both ← as.data.frame(rbind.data.frame(all, nl))
both$Penalty ← both$.set.
ggplot(both, aes(x=x1, y=yhat, color=Penalty)) + geom_line() +
  geom_abline(col=gray(.7)) + facet_grid(~ type)
# Figure 9.6

```

The left panel in Figure 9.6 corresponds to `penalty = list(simple=a, nonlinear=a)` in the R program, meaning that all parameters except the intercept are shrunk by the same amount `a` (this would be more appropriate had there been multiple predictors). As effective d.f. get smaller (penalty factor gets larger), the regression fits get flatter (too flat for the largest penalties) and confidence bands get narrower. The right graph corresponds to `penalty=list(simple=0, nonlinear=a)`, causing only the cubic spline terms that are nonlinear in X_1 to be shrunk. As the amount of shrinkage increases (d.f. lowered), the fits become more linear and closer to the true regression line (longer dotted line). Again, confidence intervals become smaller.

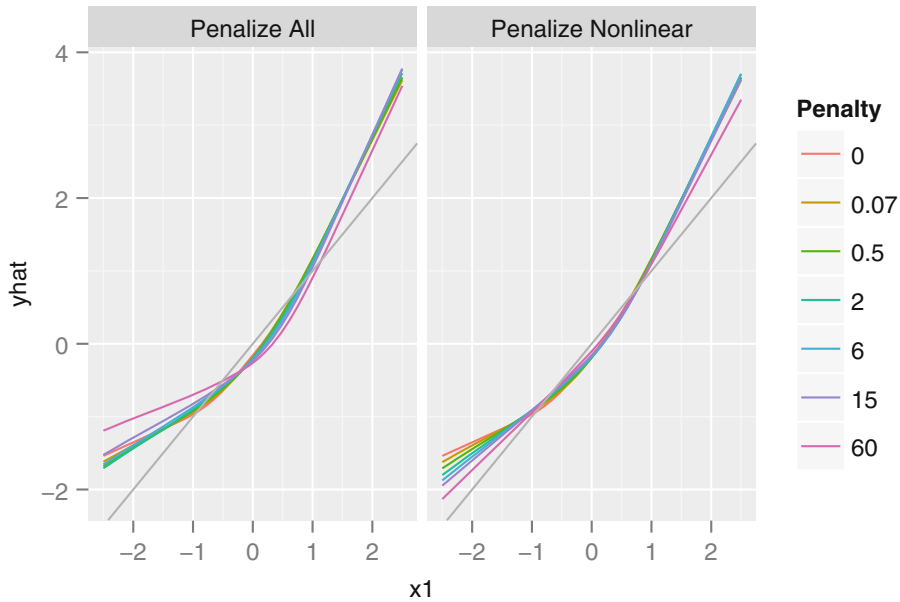


Fig. 9.6 Penalized least squares estimates for an unnecessary five-knot restricted cubic spline function. In the left graph all parameters (except the intercept) are penalized. The effective d.f. are 4, 3.21, 2.71, 2.30, 2.03, 1.82, and 1.51. In the right graph, only parameters associated with nonlinear functions of X_1 are penalized. The effective d.f. are 4, 3.22, 2.73, 2.34, 2.11, 1.96, and 1.68.

9.11 Further Reading

- 1 Boos⁶⁰ has some nice generalizations of the score test. Morgan et al.⁴⁶⁴ show how score test χ^2 statistics may negative unless the *expected* information matrix is used.
- 2 See Marubini and Valsecchi [444, pp. 164–169] for an excellent description of the relationship between the three types of test statistics.
- 3 References [115, 507] have good descriptions of methods used to maximize $\log L$.
- 4 As Long and Ervin⁴²⁶ argue, for small sample sizes, the usual Huber–White covariance estimator should not be used because there the residuals do not have constant variance even under homoscedasticity. They showed that a simple correction due to Efron and others can result in substantially better estimates. Lin and Wei,⁴¹⁰ Binder,⁵⁵ and Lin⁴⁰⁷ have applied the Huber estimator to the Cox¹³² survival model. Freedman²⁰⁶ questioned the use of sandwich estimators because they are often used to obtain the right variances on the wrong parameters when the model doesn't fit. He also has some excellent background information.
- 5 Feng et al.¹⁸⁸ showed that in the case of cluster correlations arising from repeated measurement data with Gaussian errors, the cluster bootstrap performs excellently even when the number of observations per cluster is large and the number of subjects is small. Xiao and Abrahamowicz⁶⁷⁶ compared the cluster bootstrap with a two-stage cluster bootstrap in the context of the Cox model.

- [6] Graubard and Korn²³⁵ and Fitzmaurice¹⁹⁵ describe the kinds of situations in which the working independence model can be trusted.
- [7] Minkin,⁴⁶⁰ Alho,¹¹ Doganaksoy and Schmee,¹⁶⁰ and Meeker and Escobar⁴⁵² discuss the need for LR and score-based confidence intervals. Alho found that score-based intervals are usually more tedious to compute, and provided useful algorithms for the computation of either type of interval (see also [452] and [444, p. 167]). Score and LR intervals require iterative computations and have to deal with the fact that when one parameter is changed (e.g., b_i is restricted to be zero), all other parameter estimates change. DiCiccio and Efron¹⁵⁷ provide a method for very accurate confidence intervals for exponential families that requires a modest amount of additional computation. Venzon and Moolgavkar provide an efficient general method for computing LR-based intervals.⁶³⁶ Brazzale and Davison⁶⁵ developed some promising and feasible ways to make unconditional likelihood-based inferences more accurate in small samples.
- [8] Carpenter and Bithell⁹² have an excellent overview of several variations on the bootstrap for obtaining confidence limits.
- [9] Tibshirani and Knight⁶¹⁰ developed an easy to program approach for deriving simultaneous confidence sets that is likely to be useful for getting simultaneous confidence regions for the entire vector of model parameters, for population values for an entire sequence of predictor values, and for a set of regression effects (e.g., interquartile-range odds ratios for age for both sexes). The basic idea is that during the, say, 1000 bootstrap repetitions one stores the $-2 \log$ likelihood for each model fit, being careful to compute the likelihood at the current bootstrap parameter estimates but with respect to the *original* data matrix, not the bootstrap sample of the data matrix. To obtain an approximate simultaneous 0.95 confidence set one computes the 0.95 quantile of the $-2 \log$ likelihood values and determines which vectors of parameter estimates correspond to $-2 \log$ likelihoods that are at least as small as the 0.95 quantile of all $-2 \log$ likelihoods. Once the qualifying parameter estimates are found, the quantities of interest are computed from those parameter estimates and an outer envelope of those quantities is found. Computations are facilitated with the `rms` package `confplot` function.
- [10] van Houwelingen and le Cessie [633, Eq. 52] showed, consistent with AIC, that the average optimism in a mean logarithmic (minus log likelihood) quality score for logistic models is p/n .
- [11] Schwarz⁵⁶⁰ derived a different penalty using large-sample Bayesian properties of competing models. His Bayesian Information Criterion (BIC) chooses the model having the lowest value of $L + 1/2p \log n$ or the highest value of $\text{LR } \chi^2 - p \log n$. Kass and Raftery have done several studies of BIC.³³⁷ Smith and Spiegelhalter⁵⁷⁶ and Laud and Ibrahim³⁷⁷ discussed other useful generalizations of likelihood penalties. Zheng and Loh⁶⁸⁵ studied several penalty measures, and found that AIC does not penalize enough for overfitting in the ordinary regression case. Kass and Raftery [337, p. 790] provide a nice review of this topic, stating that “AIC picks the correct model asymptotically if the complexity of the true model grows with sample size” and that “AIC selects models that are too big even when the sample size is large.” But they also cite other papers that show the existence of cases where AIC can work better than BIC. According to Buckland et al.,⁸⁰ BIC “assumes that a true model exists and is low-dimensional.” Hurvich and Tsai^{314, 315} made an improvement in AIC that resulted in much better model selection for small n . They defined the corrected AIC as

$$\text{AIC}_C = \text{LR } \chi^2 - 2p \left[1 + \frac{p+1}{n-p-1} \right]. \quad (9.66)$$

In [314] they contrast asymptotically efficient model selection with AIC when the true model has infinitely many parameters with improvements using other indexes such as AIC_C when the model is finite.

One difficulty in applying the Schwarz, AIC_C , and related criteria is that with censored or binary responses it is not clear that the actual sample size n should be used in the formula.

- [12] Goldstein,²²² Willan et al.,⁶⁶⁹ and Royston and Thompson⁵³⁴ have nice discussions on comparing non-nested regression models. Schemper's method⁵⁴⁹ is useful for testing whether a set of variables provides significantly greater information (using an R^2 measure) than another set of variables.
- [13] van Houwelingen and le Cessie [633, Eq. 22] recommended using $L/2$ (also called the Kullback–Leibler error rate) as a quality index.
- [14] Schemper⁵⁴⁹ provides a bootstrap technique for testing for significant differences between correlated R^2 measures. Mittlböck and Schemper,⁴⁶¹ Schemper and Stare,⁵⁵⁴ Korn and Simon,^{365, 366} Menard,⁴⁵⁴ and Zheng and Agresti⁶⁸⁴ have excellent discussions about the pros and cons of various indexes of the predictive value of a model.
- [15] Al-Radi et al.¹⁰ presented another analysis comparing competing predictors using the adequacy index and a receiver operating characteristic curve area approach based on a test for whether one predictor has a higher probability of being “more concordant” than another.
- [16] [55, 97, 409] provide good variance–covariance estimators from a weighted maximum likelihood analysis.
- [17] Huang and Harrington³¹⁰ developed penalized partial likelihood estimates for Cox models and provided useful background information and theoretical results about improvements in mean squared errors of regression estimates. They used a bootstrap error estimate for selection of the penalty parameter.
- [18] Sardy⁵³⁸ proposes that the square roots of the diagonals of the inverse of the covariance matrix for the predictors be used for scaling rather than the standard deviations.
- [19] Park and Hastie⁴⁸³ and articles referenced therein describe how quadratic penalized logistic regression automatically sets coefficient estimates for empty cells to zero and forces the sum of k coefficients for a k -level categorical predictor to equal zero.
- [20] Greenland²⁴¹ has a nice discussion of the relationship between penalized maximum likelihood estimation and mixed effects models. He cautions against estimating the shrinkage parameter.
- [21] See³¹⁰ for a bootstrap approach to selection of λ .
- [22] Verweij and van Houwelingen [639, Eq. 4] derived another expression for d.f., but it requires more computation and did not perform any better than Equation 9.63 in choosing λ in several examples tested.
- [23] See van Houwelingen and Thorogood⁶³¹ for an approximate empirical Bayes approach to shrinkage. See Tibshirani⁶⁰⁸ for the use of a non-smooth penalty function that results in variable selection as well as shrinkage (see Section 4.3). Verweij and van Houwelingen⁶⁴⁰ used a “cross-validated likelihood” based on leave-out-one estimates to penalize for overfitting. Wang and Taylor⁶⁵² presented some methods for carrying out hypothesis tests and computing confidence limits under penalization. Moons et al.⁴⁶² presented a case study of penalized estimation and discussed the advantages of penalization.

Table 9.4 Likelihood ratio global test statistics

Variables in Model	$LR \chi^2$
age	100
sex	108
age, sex	111
age ²	60
age, age ²	102
age, age ² , sex	115

9.12 Problems

1. A sample of size 100 from a normal distribution with unknown mean and standard deviation (μ and σ) yielded the following log likelihood values when computed at two values of μ .

$$\log L(\mu = 10, \sigma = 5) = -800$$

$$\log L(\mu = 20, \sigma = 5) = -820.$$

What do you know about μ ? What do you know about \bar{Y} ?

2. Several regression models were considered for predicting a response. $LR \chi^2$ (corrected for the intercept) for models containing various combinations of variables are found in Table 9.4. Compute all possible meaningful $LR \chi^2$. For each, state the d.f. and an approximate P -value. State which $LR \chi^2$ involving only one variable is not very meaningful.
3. For each problem below, rank Wald, score, and LR statistics by overall statistical properties and then by computational convenience.
 - a. A forward stepwise variable selection (to be later accounted for with the bootstrap) is desired to determine a concise model that contains most of the independent information in all potential predictors.
 - b. A test of independent association of each variable in a given model (each variable adjusted for the effects of all other variables in the given model) is to be obtained.
 - c. A model that contains only additive effects is fitted. A large number of potential interaction terms are to be tested using a global (multiple d.f.) test.
4. Consider a univariate saturated model in 3 treatments (A, B, C) that is quadratic in age. Write out the model with all the β s, and write in detail the contrast for comparing treatment B with treatment C for 30 year olds. Sketch out the same contrast using the “difference in predictions” approach without simplification.

5. Simulate a binary logistic model for $n = 300$ with an average fraction of events somewhere between 0.15 and 0.3. Use 5 continuous covariates and assume the model is everywhere linear. Fit an unpenalized model, then solve for the optimum quadratic penalty λ . Relate the resulting effective d.f. to the 15:1 rule of thumb, and compute the heuristic shrinkage coefficient $\hat{\gamma}$ for the unpenalized model and for the optimally penalized model, inserting the effective d.f. for the number of non-intercept parameters in the model.
6. For a similar setup as the binary logistic model simulation in Section 9.7, do a Monte Carlo simulation to determine the coverage probabilities for ordinary Wald and for three types of bootstrap confidence intervals for the true $x=5$ to $x=1$ log odds ratio. In addition, consider the Wald-type confidence interval arising from the sandwich covariance estimator. Estimate the non-coverage probabilities in both tails. Use a sample size $n = 200$ with the single predictor x_1 having a standard log-normal distribution, and the true model being $\text{logit}(Y = 1) = 1 + x_1/2$. Determine whether increasing the sample size relieves any problem you observed. Some R code for this simulation is on the web site.

Chapter 10

Binary Logistic Regression

10.1 Model

Binary responses are commonly studied in many fields. Examples include the presence or absence of a particular disease, death during surgery, or a consumer purchasing a product. Often one wishes to study how a set of predictor variables X is related to a dichotomous response variable Y . The predictors may describe such quantities as treatment assignment, dosage, risk factors, and calendar time.

1

For convenience we define the response to be $Y = 0$ or 1 , with $Y = 1$ denoting the occurrence of the event of interest. Often a dichotomous outcome can be studied by calculating certain proportions, for example, the proportion of deaths among females and the proportion among males. However, in many situations, there are multiple descriptors, or one or more of the descriptors are continuous. Without a statistical model, studying patterns such as the relationship between age and occurrence of a disease, for example, would require the creation of arbitrary age groups to allow estimation of disease prevalence as a function of age.

Letting X denote the vector of predictors $\{X_1, X_2, \dots, X_k\}$, a first attempt at modeling the response might use the ordinary linear regression model

$$E\{Y|X\} = X\beta, \quad (10.1)$$

since the expectation of a binary variable Y is $\text{Prob}\{Y = 1\}$. However, such a model by definition cannot fit the data over the whole range of the predictors since a purely linear model $E\{Y|X\} = \text{Prob}\{Y = 1|X\} = X\beta$ can allow $\text{Prob}\{Y = 1\}$ to exceed 1 or fall below 0. The statistical model that is generally preferred for the analysis of binary responses is instead the binary logistic regression model, stated in terms of the probability that $Y = 1$ given X , the values of the predictors:

$$\text{Prob}\{Y = 1|X\} = [1 + \exp(-X\beta)]^{-1}. \quad (10.2)$$

As before, $X\beta$ stands for $\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$. The binary logistic regression model was developed primarily by Cox¹²⁹ and Walker and Duncan.⁶⁴⁷ The regression parameters β are estimated by the method of maximum likelihood (see below).

The function

$$P = [1 + \exp(-x)]^{-1} \quad (10.3)$$

is called the logistic function. This function is plotted in Figure 10.1 for x varying from -4 to $+4$. This function has an unlimited range for x while P is restricted to range from 0 to 1.

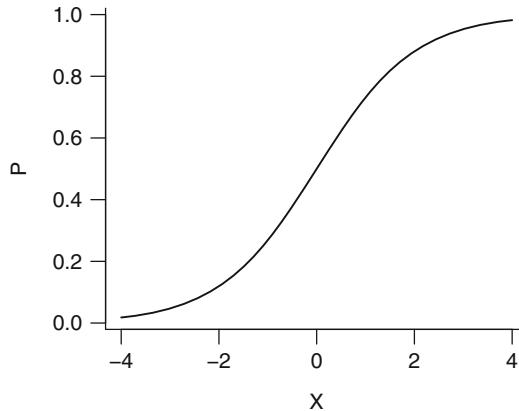


Fig. 10.1 Logistic function

For future derivations it is useful to express x in terms of P . Solving the equation above for x by using

$$1 - P = \exp(-x)/[1 + \exp(-x)] \quad (10.4)$$

yields the inverse of the logistic function:

$$x = \log[P/(1 - P)] = \log[\text{odds that } Y = 1 \text{ occurs}] = \text{logit}\{Y = 1\}. \quad (10.5)$$

Other methods that have been used to analyze binary response data include the probit model, which writes P in terms of the cumulative normal distribution, and discriminant analysis. Probit regression, although assuming a similar shape to the logistic function for the regression relationship between $X\beta$ and $\text{Prob}\{Y = 1\}$, involves more cumbersome calculations, and there is no natural interpretation of its regression parameters. In the past, discriminant analysis has been the predominant method since it is the simplest computationally. However, it makes more assumptions than logistic regression. The model used in discriminant analysis is stated in terms of the

distribution of X given the outcome group Y , even though one is seldom interested in the distribution of the predictors per se. The discriminant model has to be inverted using Bayes' rule to derive the quantity of primary interest, $\text{Prob}\{Y = 1\}$. By contrast, the logistic model is a *direct probability model* since it is stated in terms of $\text{Prob}\{Y = 1|X\}$. Since the distribution of a binary random variable Y is completely defined by the true probability that $Y = 1$ and since the model makes no assumption about the distribution of the predictors, the logistic model makes no distributional assumptions whatsoever.

10.1.1 Model Assumptions and Interpretation of Parameters

Since the logistic model is a direct probability model, its only assumptions relate to the form of the regression equation. Regression assumptions are verifiable, unlike the assumption of multivariate normality made by discriminant analysis. The logistic model assumptions are most easily understood by transforming $\text{Prob}\{Y = 1\}$ to make a model that is linear in $X\beta$:

$$\begin{aligned}\text{logit}\{Y = 1|X\} &= \text{logit}(P) = \log[P/(1 - P)] \\ &= X\beta,\end{aligned}\tag{10.6}$$

where $P = \text{Prob}\{Y = 1|X\}$. Thus the model is a linear regression model in the log odds that $Y = 1$ since $\text{logit}(P)$ is a weighted sum of the X s. If all effects are additive (i.e., no interactions are present), the model assumes that for every predictor X_j ,

$$\begin{aligned}\text{logit}\{Y = 1|X\} &= \beta_0 + \beta_1 X_1 + \dots + \beta_j X_j + \dots + \beta_k X_k \\ &= \beta_j X_j + C,\end{aligned}\tag{10.7}$$

where if all other factors are held constant, C is a constant given by

$$C = \beta_0 + \beta_1 X_1 + \dots + \beta_{j-1} X_{j-1} + \beta_{j+1} X_{j+1} + \dots + \beta_k X_k.\tag{10.8}$$

The parameter β_j is then the change in the log odds per unit change in X_j if X_j represents a single factor that is linear and does not interact with other factors and if all other factors are held constant. Instead of writing this relationship in terms of log odds, it could just as easily be written in terms of the odds that $Y = 1$:

$$\text{odds}\{Y = 1|X\} = \exp(X\beta),\tag{10.9}$$

and if all factors other than X_j are held constant,

$$\text{odds}\{Y = 1|X\} = \exp(\beta_j X_j + C) = \exp(\beta_j X_j) \exp(C). \quad (10.10)$$

The regression parameters can also be written in terms of *odds ratios*. The odds that $Y = 1$ when X_j is increased by d , divided by the odds at X_j is

$$\begin{aligned} \frac{\text{odds}\{Y = 1|X_1, X_2, \dots, X_j + d, \dots, X_k\}}{\text{odds}\{Y = 1|X_1, X_2, \dots, X_j, \dots, X_k\}} &= \frac{\exp[\beta_j(X_j + d)] \exp(C)}{[\exp(\beta_j X_j) \exp(C)]} \\ &= \exp[\beta_j X_j + \beta_j d - \beta_j X_j] = \exp(\beta_j d). \end{aligned} \quad (10.11)$$

Thus the effect of increasing X_j by d is to increase the odds that $Y = 1$ by a factor of $\exp(\beta_j d)$, or to increase the log odds that $Y = 1$ by an increment of $\beta_j d$. In general, the ratio of the odds of response for an individual with predictor variable values X^* compared with an individual with predictors X is

$$\begin{aligned} X^* : X \text{ odds ratio} &= \exp(X^* \beta) / \exp(X \beta) \\ &= \exp[(X^* - X) \beta]. \end{aligned} \quad (10.12)$$

Now consider some special cases of the logistic multiple regression model. If there is only one predictor X and that predictor is binary, the model can be written

$$\begin{aligned} \text{logit}\{Y = 1|X = 0\} &= \beta_0 \\ \text{logit}\{Y = 1|X = 1\} &= \beta_0 + \beta_1. \end{aligned} \quad (10.13)$$

Here β_0 is the log odds of $Y = 1$ when $X = 0$. By subtracting the two equations above, it can be seen that β_1 is the difference in the log odds when $X = 1$ as compared with $X = 0$, which is equivalent to the log of the ratio of the odds when $X = 1$ compared with the odds when $X = 0$. The quantity $\exp(\beta_1)$ is the odds ratio for $X = 1$ compared with $X = 0$. Letting $P^0 = \text{Prob}\{Y = 1|X = 0\}$ and $P^1 = \text{Prob}\{Y = 1|X = 1\}$, the regression parameters are interpreted by

$$\begin{aligned} \beta_0 &= \text{logit}(P^0) = \log[P^0/(1 - P^0)] \\ \beta_1 &= \text{logit}(P^1) - \text{logit}(P^0) \\ &= \log[P^1/(1 - P^1)] - \log[P^0/(1 - P^0)] \\ &= \log\{[P^1/(1 - P^1)]/[P^0/(1 - P^0)]\}. \end{aligned} \quad (10.14)$$

Since there are only two quantities to model and two free parameters, there is no way that this two-sample model can't fit; the model in this case is essentially fitting two cell proportions. Similarly, if there are $g - 1$ dummy indicator X s representing g groups, the ANOVA-type logistic model must always fit.

If there is one continuous predictor X , the model is

$$\text{logit}\{Y = 1|X\} = \beta_0 + \beta_1 X, \quad (10.15)$$

and without further modification (e.g., taking log transformation of the predictor), the model assumes a straight line in the log odds, or that an increase in X by one unit increases the odds by a factor of $\exp(\beta_1)$.

Now consider the simplest analysis of covariance model in which there are two treatments (indicated by $X_1 = 0$ or 1) and one continuous covariable (X_2). The simplest logistic model for this setup is

$$\text{logit}\{Y = 1|X\} = \beta_0 + \beta_1 X_1 + \beta_2 X_2, \quad (10.16)$$

which can be written also as

$$\begin{aligned} \text{logit}\{Y = 1|X_1 = 0, X_2\} &= \beta_0 + \beta_2 X_2 \\ \text{logit}\{Y = 1|X_1 = 1, X_2\} &= \beta_0 + \beta_1 + \beta_2 X_2. \end{aligned} \quad (10.17)$$

The $X_1 = 1 : X_1 = 0$ odds ratio is $\exp(\beta_1)$, independent of X_2 . The odds ratio for a one-unit increase in X_2 is $\exp(\beta_2)$, independent of X_1 .

This model, with no term for a possible interaction between treatment and covariable, assumes that for each treatment the relationship between X_2 and log odds is linear, and that the lines have equal slope; that is, they are parallel. Assuming linearity in X_2 , the only way that this model can fail is for the two slopes to differ. Thus, the only assumptions that need verification are linearity and lack of interaction between X_1 and X_2 .

To adapt the model to allow or test for interaction, we write

$$\text{logit}\{Y = 1|X\} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3, \quad (10.18)$$

where the derived variable X_3 is defined to be $X_1 X_2$. The test for lack of interaction (equal slopes) is $H_0 : \beta_3 = 0$. The model can be amplified as

$$\begin{aligned} \text{logit}\{Y = 1|X_1 = 0, X_2\} &= \beta_0 + \beta_2 X_2 \\ \text{logit}\{Y = 1|X_1 = 1, X_2\} &= \beta_0 + \beta_1 + \beta_2 X_2 + \beta_3 X_2 \\ &= \beta'_0 + \beta'_2 X_2, \end{aligned} \quad (10.19)$$

Table 10.1 Effect of an odds ratio of two on various risks

Without Risk Factor		With Risk Factor	
Probability	Odds	Odds	Probability
.2	.25	.5	.33
.5	1	2	.67
.8	4	8	.89
.9	9	18	.95
.98	49	98	.99

where $\beta'_0 = \beta_0 + \beta_1$ and $\beta'_2 = \beta_2 + \beta_3$. The model with interaction is therefore equivalent to fitting two separate logistic models with X_2 as the only predictor, one model for each treatment group. Here the $X_1 = 1 : X_1 = 0$ odds ratio is $\exp(\beta_1 + \beta_3 X_2)$.

10.1.2 Odds Ratio, Risk Ratio, and Risk Difference

As discussed above, the logistic model quantifies the effect of a predictor in terms of an odds ratio or log odds ratio. An odds ratio is a natural description of an effect in a probability model since an odds ratio *can* be constant. For example, suppose that a given risk factor doubles the odds of disease. Table 10.1 shows the effect of the risk factor for various levels of initial risk.

Since odds have an unlimited range, any positive odds ratio will still yield a valid probability. If one attempted to describe an effect by a risk ratio, the effect can only occur over a limited range of risk (probability). For example, a risk ratio of 2 can only apply to risks below .5; above that point the risk ratio must diminish. (Risk ratios are similar to odds ratios if the risk is small.) Risk differences have the same difficulty; the risk difference cannot be constant and must depend on the initial risk. Odds ratios, on the other hand, can describe an effect over the entire range of risk. An odds ratio can, for example, describe the effect of a treatment independently of covariables affecting risk. Figure 10.2 depicts the relationship between risk of a subject without the risk factor and the increase in risk for a variety of relative increases (odds ratios). It demonstrates how absolute risk increase is a function of the baseline risk. Risk increase will also be a function of factors that interact with the risk factor, that is, factors that modify its relative effect. Once a model is developed for estimating $\text{Prob}\{Y = 1|X\}$, this model can easily be used to estimate the absolute risk increase as a function of baseline risk factors as well as interacting factors. Let X_1 be a binary risk factor and let $A = \{X_2, \dots, X_p\}$ be the other factors (which for convenience we assume do not interact with X_1). Then the estimate of $\text{Prob}\{Y = 1|X_1 = 1, A\} - \text{Prob}\{Y = 1|X_1 = 0, A\}$ is

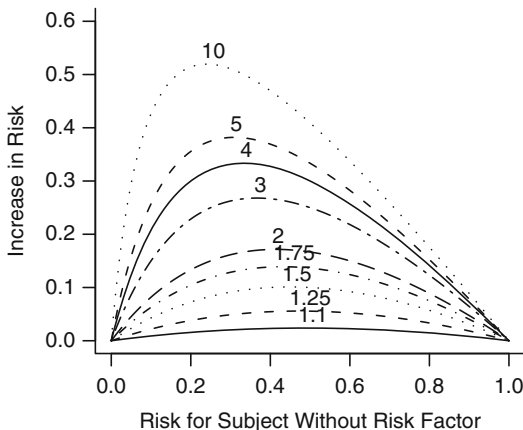


Fig. 10.2 Absolute benefit as a function of risk of the event in a control subject and the relative effect (odds ratio) of the risk factor. The odds ratios are given for each curve.

Table 10.2 Example binary response data

Females	Age:	37	39	39	42	47	48	48	52	53	55	56	57	58	58	60	64	65	68	68	70	
	Response:	0	0	0	0	0	0	1	0	0	0	0	0	0	1	0	0	1	1	1	1	1
Males	Age:	34	38	40	40	41	43	43	43	44	46	47	48	48	50	50	52	55	60	61	61	61
	Response:	1	1	0	0	0	1	1	1	0	0	1	1	1	0	1	1	1	1	1	1	1

$$\begin{aligned}
 & \frac{1}{1 + \exp -[\hat{\beta}_0 + \hat{\beta}_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p]} \\
 & - \frac{1}{1 + \exp -[\hat{\beta}_0 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p]} \qquad (10.20) \\
 & = \frac{1}{1 + (\frac{1-\hat{R}}{\hat{R}}) \exp(-\hat{\beta}_1)} - \hat{R},
 \end{aligned}$$

where \hat{R} is the estimate of the baseline risk, $\text{Prob}\{Y = 1|X_1 = 0\}$. The risk difference estimate can be plotted against \hat{R} or against levels of variables in A to display absolute risk increase against overall risk (Figure 10.2) or against specific subject characteristics.

4

10.1.3 Detailed Example

Consider the data in Table 10.2. A graph of the data, along with a fitted logistic model (described later), appears in Figure 10.3. The graph also displays proportions of responses obtained by stratifying the data by sex and

age group ($< 45, 45 - 54, \geq 55$). The age points on the abscissa for these groups are the overall mean ages in the three age intervals (40.2, 49.1, and 61.1, respectively).

```
require(rms)

getHdata(sex.age.response)
d <- sex.age.response
dd <- datadist(d); options(datadist='dd')
f <- lrm(response ~ sex + age, data=d)
fasr <- f # Save for later
w <- function(...)
  with(d, {
    m <- sex=='male'
    f <- sex=='female'
    lpoints(age[f], response[f], pch=1)
    lpoints(age[m], response[m], pch=2)
    af <- cut2(age, c(45,55), levels.mean=TRUE)
    prop <- tapply(response, list(af, sex), mean,
                   na.rm=TRUE)
    agem <- as.numeric(row.names(prop))
    lpoints(agem, prop[, 'female'],
            pch=4, cex=1.3, col='green')
    lpoints(agem, prop[, 'male'],
            pch=5, cex=1.3, col='green')
    x <- rep(62, 4); y <- seq(.25, .1, length=4)
    lpoints(x, y, pch=c(1, 2, 4, 5),
            col=rep(c('blue', 'green'), each=2))
    ltext(x+5, y,
          c('F Observed', 'M Observed',
            'F Proportion', 'M Proportion'), cex=.8)
  }) # Figure 10.3

plot(Predict(f, age=seq(34, 70, length=200), sex, fun=plogis),
     ylab='Pr[response]', ylim=c(-.02, 1.02), addpanel=w)
ltx <- function(fit) latex(fit, inline=TRUE, columns=54,
                           file='', after='$.', digits=3,
                           size='Ssize', before='$X\\hat{\\beta}$')
ltx(f)
```

$$X\hat{\beta} = -9.84 + 3.49[\text{male}] + 0.158 \text{ age.}$$

Descriptive statistics for assessing the association between sex and response, age group and response, and age group and response stratified by sex are found below. Corresponding fitted logistic models, with sex coded as 0 = female, 1 = male are also given. Models were fitted first with sex as the only predictor, then with age as the (continuous) predictor, then with sex and age simultaneously. First consider the relationship between sex and response, ignoring the effect of age.

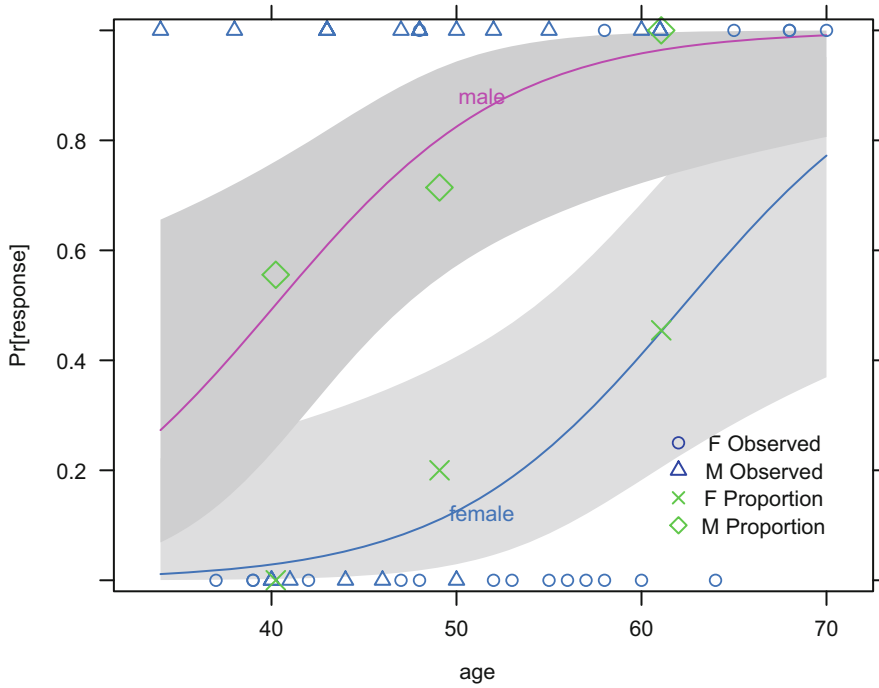


Fig. 10.3 Data, subgroup proportions, and fitted logistic model, with 0.95 pointwise confidence bands

sex	response		Total	Odds/Log
	0	1		
F	14 70.00	6 30.00	20	6/14=.429 -.847
M	6 30.00	14 70.00	20	14/6=2.33 .847
Total	20	20	40	

M:F odds ratio = (14/6)/(6/14) = 5.44, log=1.695

Statistics for sex × response

Statistic	d.f.	Value	P
χ^2	1	6.400	0.011
Likelihood Ratio χ^2	1	6.583	0.010
Parameter Estimate	Std Err	Wald χ^2	P
β_0	-0.8473	0.4880	3.0152
β_1	1.6946	0.6901	6.0305 0.0141

Note that the estimate of β_0 , $\hat{\beta}_0$ is the log odds for females and that $\hat{\beta}_1$ is the log odds (M:F) ratio. $\hat{\beta}_0 + \hat{\beta}_1 = .847$, the log odds for males. The likelihood ratio test for H_0 : no effect of sex on probability of response is obtained as follows.

$$\begin{aligned} \text{Log likelihood } (\beta_1 = 0) &: -27.727 \\ \text{Log likelihood (max)} &: -24.435 \\ \text{LR } \chi^2(H_0 : \beta_1 = 0) &: -2(-27.727 - -24.435) = 6.584. \end{aligned}$$

(Note the agreement of the LR χ^2 with the contingency table likelihood ratio χ^2 , and compare 6.584 with the Wald statistic 6.03.)

Next, consider the relationship between age and response, ignoring sex.

age	response			
Frequency	0	1	Total	Odds/Log
Row Pct				
<45	8	5	13	5/8=.625
	61.5	38.4		-.47
45-54	6	6	12	6/6=1
	50.0	50.0		0
55+	6	9	15	9/6=1.5
	40.0	60.0		.405
Total	20	20	40	

$$55+ : <45 \text{ odds ratio} = (9/6)/(5/8) = 2.4, \log=.875$$

Parameter	Estimate	Std Err	Wald χ^2	P
β_0	-2.7338	1.8375	2.2134	0.1368
β_1	0.0540	0.0358	2.2763	0.1314

The estimate of β_1 is in rough agreement with that obtained from the frequency table. The 55+ : < 45 log odds ratio is .875, and since the respective mean ages in the 55+ and <45 age groups are 61.1 and 40.2, an estimate of the log odds ratio increase per year is $.875/(61.1 - 40.2) = .875/20.9 = .042$.

The likelihood ratio test for H_0 : no association between age and response is obtained as follows.

$$\begin{aligned} \text{Log likelihood } (\beta_1 = 0) &: -27.727 \\ \text{Log likelihood (max)} &: -26.511 \\ \text{LR } \chi^2(H_0 : \beta_1 = 0) &: -2(-27.727 - -26.511) = 2.432. \end{aligned}$$

(Compare 2.432 with the Wald statistic 2.28.)

Next we consider the simultaneous association of age and sex with response.

sex=F

age	response		Total
	0	1	
Frequency			
Row Pct			
<45	4	0	4
	100.0	0.0	
45-54	4	1	5
	80.0	20.0	
55+	6	5	11
	54.6	45.4	
Total	14	6	20

sex=M

age	response		Total
	0	1	
Frequency			
Row Pct			
<45	4	5	9
	44.4	55.6	
45-54	2	5	7
	28.6	71.4	
55+	0	4	4
	0.0	100.0	
Total	6	14	20

A logistic model for relating sex and age simultaneously to response is given below.

Parameter	Estimate	Std Err	Wald χ^2	P
β_0	-9.8429	3.6758	7.1706	0.0074
β_1 (sex)	3.4898	1.1992	8.4693	0.0036
β_2 (age)	0.1581	0.0616	6.5756	0.0103

Likelihood ratio tests are obtained from the information below.

Log likelihood ($\beta_1 = 0, \beta_2 = 0$)	: -27.727
Log likelihood (max)	: -19.458
Log likelihood ($\beta_1 = 0$)	: -26.511
Log likelihood ($\beta_2 = 0$)	: -24.435
LR χ^2 ($H_0 : \beta_1 = \beta_2 = 0$)	: $-2(-27.727 - -19.458) = 16.538$
LR χ^2 ($H_0 : \beta_1 = 0$) sex age	: $-2(-26.511 - -19.458) = 14.106$
LR χ^2 ($H_0 : \beta_2 = 0$) age sex	: $-2(-24.435 - -19.458) = 9.954$

The 14.1 should be compared with the Wald statistic of 8.47, and 9.954 should be compared with 6.58. The fitted logistic model is plotted separately

for females and males in Figure 10.3. The fitted model is

$$\text{logit}\{\text{Response} = 1|\text{sex,age}\} = -9.84 + 3.49 \times \text{sex} + .158 \times \text{age}, \quad (10.21)$$

where as before sex = 0 for females, 1 for males. For example, for a 40-year-old female, the predicted logit is $-9.84 + .158(40) = -3.52$. The predicted probability of a response is $1/[1 + \exp(3.52)] = .029$. For a 40-year-old male, the predicted logit is $-9.84 + 3.49 + .158(40) = -.03$, with a probability of .492.

10.1.4 Design Formulations

The logistic multiple regression model can incorporate the same designs as can ordinary linear regression. An analysis of variance (ANOVA) model for a treatment with k levels can be formulated with $k - 1$ dummy variables. This logistic model is equivalent to a $2 \times k$ contingency table. An analysis of covariance logistic model is simply an ANOVA model augmented with covariables used for adjustment.

One unique design that is interesting to consider in the context of logistic models is a simultaneous comparison of multiple factors between two groups. Suppose, for example, that in a randomized trial with two treatments one wished to test whether any of 10 baseline characteristics are mal-distributed between the two groups. If the 10 factors are continuous, one could perform a two-sample Wilcoxon–Mann–Whitney test or a t -test for each factor (if each is normally distributed). However, this procedure would result in multiple comparison problems and would also not be able to detect the combined effect of small differences across all the factors. A better procedure would be a multivariate test. The Hotelling T^2 test is designed for just this situation. It is a k -variable extension of the one-variable unpaired t -test. The T^2 test, like discriminant analysis, assumes multivariate normality of the k factors. This assumption is especially tenuous when some of the factors are polytomous. A better alternative is the global test of no regression from the logistic model. This test is valid because it can be shown that $H_0 : \text{mean } X$ is the same for both groups ($= H_0 : \text{mean } X$ does not depend on group $= H_0 : \text{mean } X | \text{group} = \text{constant}$) is true if and only if $H_0 : \text{Prob}\{\text{group}|X\} = \text{constant}$. Thus k factors can be tested simultaneously for differences between the two groups using the binary logistic model, which has far fewer assumptions than does the Hotelling T^2 test. The logistic global test of no regression (with k d.f.) would be expected to have greater power if there is non-normality. Since the logistic model makes no assumption regarding the distribution of the descriptor variables, it can easily test for simultaneous group differences involving a mixture of continuous, binary, and nominal variables. In observational studies, such

models for treatment received or exposure (propensity score models) hold great promise for adjusting for confounding.^{117, 380, 526, 530, 531}

5

O'Brien⁴⁷⁹ has developed a general test for comparing group 1 with group 2 for a single measurement. His test detects location and scale differences by fitting a logistic model for $\text{Prob}\{\text{Group } 2\}$ using X and X^2 as predictors.

For a randomized study where adjustment for confounding is seldom necessary, adjusting for covariables using a binary logistic model results in *increases* in standard errors of regression coefficients.⁵²⁷ This is the opposite of what happens in linear regression where there is an unknown variance parameter that is estimated using the residual squared error. Fortunately, adjusting for covariables using logistic regression, by accounting for subject heterogeneity, will result in larger regression coefficients even for a randomized treatment variable. The increase in estimated regression coefficients more than offsets the increase in standard error.^{284, 285, 527, 588}

10.2 Estimation

10.2.1 Maximum Likelihood Estimates

The parameters in the logistic regression model are estimated using the maximum likelihood (ML) method. The method is based on the same principles as the one-sample proportion example described in Section 9.1. The difference is that the general logistic model is not a single sample or a two-sample problem. The probability of response for the i th subject depends on a particular set of predictors X_i , and in fact the list of predictors may not be the same for any two subjects. Denoting the response and probability of response of the i th subject by Y_i and P_i , respectively, the model states that

$$P_i = \text{Prob}\{Y_i = 1|X_i\} = [1 + \exp(-X_i\beta)]^{-1}. \quad (10.22)$$

The likelihood of an observed response Y_i given predictors X_i and the unknown parameters β is

$$P_i^{Y_i} [1 - P_i]^{1 - Y_i}. \quad (10.23)$$

The joint likelihood of all responses Y_1, Y_2, \dots, Y_n is the product of these likelihoods for $i = 1, \dots, n$. The likelihood and log likelihood functions are rewritten by using the definition of P_i above to allow them to be recognized as a function of the unknown parameters β . Except in simple special cases (such as the k -sample problem in which all X s are dummy variables), the ML estimates (MLE) of β cannot be written explicitly. The Newton–Raphson method described in Section 9.4 is usually used to solve iteratively for the list of values β that maximize the log likelihood. The MLEs are denoted by

$\hat{\beta}$. The inverse of the estimated observed information matrix is taken as the estimate of the variance–covariance matrix of $\hat{\beta}$.

Under $H_0 : \beta_1 = \beta_2 = \dots = \beta_k = 0$, the intercept parameter β_0 can be estimated explicitly and the log likelihood under this global null hypothesis can be computed explicitly. Under the global null hypothesis, $P_i = P = [1 + \exp(-\beta_0)]^{-1}$ and the MLE of P is $\hat{P} = s/n$ where s is the number of responses and n is the sample size. The MLE of β_0 is $\hat{\beta}_0 = \text{logit}(\hat{P})$. The log likelihood under this null hypothesis is

6

$$\begin{aligned} & s \log(\hat{P}) + (n - s) \log(1 - \hat{P}) \\ &= s \log(s/n) + (n - s) \log[(n - s)/n] \\ &= s \log s + (n - s) \log(n - s) - n \log(n). \end{aligned} \tag{10.24}$$

10.2.2 Estimation of Odds Ratios and Probabilities

Once β is estimated, one can estimate any log odds, odds, or odds ratios. The MLE of the $X_j + 1 : X_j$ log odds ratio is $\hat{\beta}_j$, and the estimate of the $X_j + d : X_j$ log odds ratio is $\hat{\beta}_j d$, all other predictors remaining constant (assuming the absence of interactions and nonlinearities involving X_j). For large enough samples, the MLEs are normally distributed with variances that are consistently estimated from the estimated variance–covariance matrix. Letting z denote the $1 - \alpha/2$ critical value of the standard normal distribution, a two-sided $1 - \alpha$ confidence interval for the log odds ratio for a one-unit increase in X_j is $[\hat{\beta}_j - zs, \hat{\beta}_j + zs]$, where s is the estimated standard error of $\hat{\beta}_j$. (Note that for $\alpha = .05$, i.e., for a 95% confidence interval, $z = 1.96$.)

A theorem in statistics states that the MLE of a function of a parameter is that same function of the MLE of the parameter. Thus the MLE of the $X_j + 1 : X_j$ odds ratio is $\exp(\hat{\beta}_j)$. Also, if a $1 - \alpha$ confidence interval of a parameter β is $[c, d]$ and $f(u)$ is a one-to-one function, a $1 - \alpha$ confidence interval of $f(\beta)$ is $[f(c), f(d)]$. Thus a $1 - \alpha$ confidence interval for the $X_j + 1 : X_j$ odds ratio is $\exp[\hat{\beta}_j \pm zs]$. Note that while the confidence interval for β_j is symmetric about $\hat{\beta}_j$, the confidence interval for $\exp(\beta_j)$ is not. By the same theorem just used, the MLE of $P_i = \text{Prob}\{Y_i = 1 | X_i\}$ is

$$\hat{P}_i = [1 + \exp(-X_i \hat{\beta})]^{-1}. \tag{10.25}$$

A confidence interval for P_i could be derived by computing the standard error of \hat{P}_i , yielding a symmetric confidence interval. However, such an interval would have the disadvantage that its endpoints could fall below zero or exceed one. A better approach uses the fact that for large samples $X \hat{\beta}$ is approximately normally distributed. An estimate of the variance of $X \hat{\beta}$ in matrix notation is XVX' where V is the estimated variance–covariance

matrix of $\hat{\beta}$ (see Equation 9.51). This variance is the sum of all variances and covariances of $\hat{\beta}$ weighted by squares and products of the predictors. The estimated standard error of $X\hat{\beta}$, s , is the square root of this variance estimate. A $1 - \alpha$ confidence interval for P_i is then

$$\{1 + \exp[-(X_i\hat{\beta} \pm zs)]\}^{-1}. \quad (10.26)$$

7

10.2.3 Minimum Sample Size Requirement

Suppose there were no covariates, so that the only parameter in the model is the intercept. What is the sample size required to allow the estimate of the intercept to be precise enough so that the predicted probability is within 0.1 of the true probability with 0.95 confidence, when the true intercept is in the neighborhood of zero? The answer is $n=96$. What if there were one covariate, and it was binary with a prevalence of $\frac{1}{2}$? One would need 96 subjects with $X = 0$ and 96 with $X = 1$ to have an upper bound on the margin of error for estimating $\text{Prob}\{Y = 1|X = x\}$ not exceed 0.1 for either value of x^a .

Now consider a very simple single continuous predictor case in which X has a normal distribution with mean zero and standard deviation σ , with the true $\text{Prob}\{Y = 1|X = x\} = [1 + \exp(-x)]^{-1}$. The expected number of events is $\frac{n}{2}^b$. The following simulation answers the question “What should n be so that the expected maximum absolute error (over $x \in [-1.5, 1.5]$) in \hat{P} is less than ϵ ?”

```

sigmas  <- c(.5, .75, 1, 1.25, 1.5, 1.75, 2, 2.5, 3, 4)
ns      <- seq(25, 300, by=25)
nsim    <- 1000
xs      <- seq(-1.5, 1.5, length=200)
pactual <- plogis(xs)

dn <- list(sigma=format(sigmas), n=format(ns))
maxerr <- N1 <- array(NA, c(length(sigmas), length(ns)), dn)
require(rms)

i <- 0
for(s in sigmas) {
  i <- i + 1
  j <- 0
  for(n in ns) {

```

^a The general formula for the sample size required to achieve a margin of error of δ in estimating a true probability of θ at the 0.95 confidence level is $n = (\frac{1.96}{\delta})^2 \times \theta(1-\theta)$. Set $\theta = \frac{1}{2}$ (intercept=0) for the worst case.

^b The R code can easily be modified for other event frequencies, or the minimum of the number of events and non-events for a dataset at hand can be compared with $\frac{n}{2}$ in this simulation. An average maximum absolute error of 0.05 corresponds roughly to a half-width of the 0.95 confidence interval of 0.1.


```

j ← j + 1
n1 ← maxe ← 0
for(k in 1:nsim) {
  x ← rnorm(n, 0, s)
  P ← plogis(x)
  y ← ifelse(runif(n) ≤ P, 1, 0)
  n1 ← n1 + sum(y)
  beta ← lrm.fit(x, y)$coefficients
  phat ← plogis(beta[1] + beta[2] * xs)
  maxe ← maxe + max(abs(phat - pactual))
}
n1 ← n1/nsim
maxe ← maxe/nsim
maxerr[i,j] ← maxe
N1[i,j] ← n1
}
}
xrange ← range(xs)
simerr ← llist(N1, maxerr, sigmas, ns, nsim, xrange)

maxe ← reShape(maxerr)
# Figure 10.4
xYplot(maxerr ~ n, groups=sigma, data=maxe,
        ylab=expression(paste('Average Maximum ',
                               abs(hat(P) - P))),
        type='l', lty=rep(1:2, 5), label.curve=FALSE,
        abline=list(h=c(.15, .1, .05), col=gray(.85)))
Key(.8, .68, other=list(cex=.7,
                        title=expression(~~~~~sigma)))

```

10.3 Test Statistics

The likelihood ratio, score, and Wald statistics discussed earlier can be used to test any hypothesis in the logistic model. The likelihood ratio test is generally preferred. When true parameters are near the null values all three statistics usually agree. The Wald test has a significant drawback when the true parameter value is very far from the null value. In such case the standard error estimate becomes too large. As $\hat{\beta}_j$ increases from 0, the Wald test statistic for $H_0 : \beta_j = 0$ becomes larger, but after a certain point it becomes smaller. The statistic will eventually drop to zero if $\hat{\beta}_j$ becomes infinite.²⁷⁸ Infinite estimates can occur in the logistic model especially when there is a binary predictor whose mean is near 0 or 1. Wald statistics are especially problematic in this case. For example, if 10 out of 20 males had a disease and 5 out of 5 females had the disease, the female : male odds ratio is infinite and so is the logistic regression coefficient for sex. If such a situation occurs, the likelihood ratio or score statistic should be used instead of the Wald statistic.

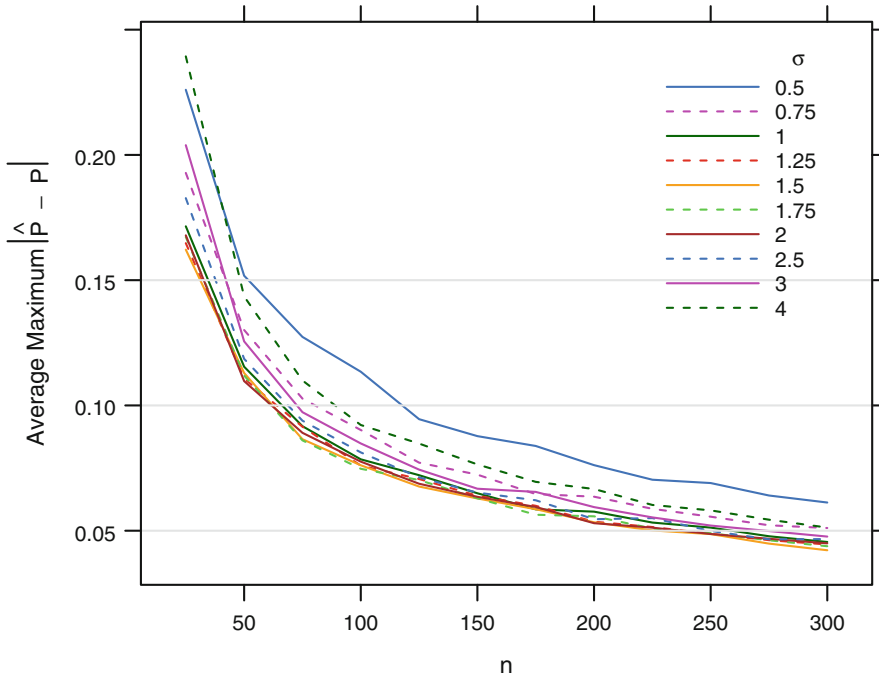


Fig. 10.4 Simulated expected maximum error in estimating probabilities for $x \in [-1.5, 1.5]$ with a single normally distributed X with mean zero

For k -sample (ANOVA-type) logistic models, logistic model statistics are equivalent to contingency table χ^2 statistics. As exemplified in the logistic model relating sex to response described previously, the global likelihood ratio statistic for all dummy variables in a k -sample model is identical to the contingency table (k -sample binomial) likelihood ratio χ^2 statistic. The score statistic for this same situation turns out to be identical to the $k - 1$ degrees of freedom Pearson χ^2 for a $k \times 2$ table.

As mentioned in Section 2.6, it can be dangerous to interpret individual parameters, make pairwise treatment comparisons, or test linearity if the overall test of association for a factor represented by multiple parameters is insignificant.

10.4 Residuals

Several types of residuals can be computed for binary logistic model fits. Many of these residuals are used to examine the influence of individual observations on the fit. The *partial residual* can be used for directly assessing how each

predictor should be transformed. For the i th observation, the partial residual for the j th element of X is defined by

$$r_{ij} = \hat{\beta}_j X_{ij} + \frac{Y_i - \hat{P}_i}{\hat{P}_i(1 - \hat{P}_i)}, \quad (10.27)$$

where X_{ij} is the value of the j th variable in the i th observation, Y_i is the corresponding value of the response, and \hat{P}_i is the predicted probability that $Y_i = 1$. A smooth plot (using, e.g., loess) of X_{ij} against r_{ij} will provide an estimate of how X_j should be transformed, adjusting for the other X s (using their current transformations). Typically one tentatively models X_j linearly and checks the smoothed plot for linearity. A U -shaped relationship in this plot, for example, indicates that a squared term or spline function needs to be added for X_j . This approach does assume additivity of predictors.

9

10.5 Assessment of Model Fit

As the logistic regression model makes no distributional assumptions, only the assumptions of linearity and additivity need to be verified (in addition to the usual assumptions about independence of observations and inclusion of important covariables). In ordinary linear regression there is no global test for lack of model fit unless there are replicate observations at various settings of X . This is because ordinary regression entails estimation of a separate variance parameter σ^2 . In logistic regression there are global tests for goodness of fit. Unfortunately, some of the most frequently used ones are inappropriate. For example, it is common to see a deviance test of goodness of fit based on the “residual” log likelihood, with P -values obtained from a χ^2 distribution with $n - p$ d.f. This P -value is inappropriate since the deviance does not have an asymptotic χ^2 distribution, due to the facts that the number of parameters estimated is increasing at the same rate as n and the expected cell frequencies are far below five (by definition).

Hosmer and Lemeshow³⁰⁴ have developed a commonly used test for goodness of fit for binary logistic models based on grouping into deciles of predicted probability and performing an ordinary χ^2 test for the mean predicted probability against the observed fraction of events (using 8 d.f. to account for evaluating fit on the model development sample). The Hosmer–Lemeshow test is dependent on the choice of how predictions are grouped³⁰³ and it is not clear that the choice of the number of groups should be independent of n . Hosmer et al.³⁰³ have compared a number of global goodness of fit tests for binary logistic regression. They concluded that the simple unweighted sum of squares test of Copas¹²⁴ as modified by le Cessie and van Houwelingen³⁸⁷ is as

good as any. They used a normal Z -test for the sum of squared errors ($n \times B$, where B is the Brier index in Equation 10.35). This test takes into account the fact that one cannot obtain a χ^2 distribution for the sum of squares. It also takes into account the estimation of β . It is not yet clear for which types of lack of fit this test has reasonable power. Returning to the external validation case where uncertainty of β does not need to be accounted for, Stallard⁵⁸⁴ has further documented the lack of power of the original Hosmer-Lemeshow test and found more power with a logarithmic scoring rule (deviance test) and a χ^2 test that, unlike the simple unweighted sum of squares test, weights each squared error by dividing it by $\hat{P}_i(1 - \hat{P}_i)$. A scaled χ^2 distribution seemed to provide the best approximation to the null distribution of the test statistics.

More power for detecting lack of fit is expected to be obtained from testing specific alternatives to the model. In the model

$$\text{logit}\{Y = 1|X\} = \beta_0 + \beta_1 X_1 + \beta_2 X_2, \quad (10.28)$$

where X_1 is binary and X_2 is continuous, one needs to verify that the log odds is related to X_1 and X_2 according to Figure 10.5.

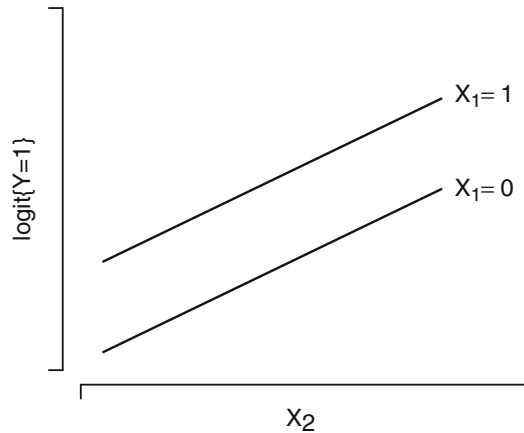


Fig. 10.5 Logistic regression assumptions for one binary and one continuous predictor

The simplest method for validating that the data are consistent with the no-interaction linear model involves stratifying the sample by X_1 and quantile groups (e.g., deciles) of X_2 .²⁶⁵ Within each stratum the proportion of responses \hat{P} is computed and the log odds calculated from $\log[\hat{P}/(1 - \hat{P})]$. The number of quantile groups should be such that there are at least 20 (and perhaps many more) subjects in each $X_1 \times X_2$ group. Otherwise, probabilities cannot be estimated precisely enough to allow trends to be seen above “noise” in the data. Since at least 3 X_2 groups must be formed to allow assessment of linearity, the total sample size must be at least $2 \times 3 \times 20 = 120$ for this method to work at all.

Figure 10.6 demonstrates this method for a large sample size of 3504 subjects stratified by sex and deciles of age. Linearity is apparent for males while there is evidence for slight interaction between age and sex since the age trend for females appears curved.

```
getHdata(acath)
acath$sex ← factor(acath$sex, 0:1, c('male','female'))
dd ← datadist(acath); options(datadist='dd')
f ← lrm(sigdz ~ rcs(age, 4) * sex, data=acath)

w ← function(...)
  with(acath, {
    plsmo(age, sigdz, group=sex, fun=qlogis, lty='dotted',
          add=TRUE, grid=TRUE)
    af ← cut2(age, g=10, levels.mean=TRUE)
    prop ← qlogis(tapply(sigdz, list(af, sex), mean,
                             na.rm=TRUE))
    agem ← as.numeric(row.names(prop))
    lpoints(agem, prop[, 'female'], pch=4, col='green')
    lpoints(agem, prop[, 'male'], pch=2, col='green')
  }) # Figure 10.6
plot(Predict(f, age, sex), ylim=c(-2,4), addpanel=w,
     label.curve=list(offset=unit(0.5, 'cm')))
```

The subgrouping method requires relatively large sample sizes and does not use continuous factors effectively. The ordering of values is not used at all between intervals, and the estimate of the relationship for a continuous variable has little resolution. Also, the method of grouping chosen (e.g., deciles vs. quintiles vs. rounding) can alter the shape of the plot.

In this dataset with only two variables, it is efficient to use a nonparametric smoother for age, separately for males and females. Nonparametric smoothers, such as `loess`¹¹¹ used here, work well for binary response variables (see Section 2.4.7); the logit transformation is made on the smoothed probability estimates. The smoothed estimates are shown in Figure 10.6.

10

When there are several predictors, the restricted cubic spline function is better for estimating the true relationship between X_2 and $\text{logit}\{Y = 1\}$ for continuous variables without assuming linearity. By fitting a model containing X_2 expanded into $k - 1$ terms, where k is the number of knots, one can obtain an estimate of the transformation of X_2 as discussed in Section 2.4:

$$\begin{aligned} \text{logit}\{Y = 1|X\} &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \hat{\beta}_3 X_2' + \hat{\beta}_4 X_2'' \\ &= \hat{\beta}_0 + \hat{\beta}_1 X_1 + f(X_2), \end{aligned} \quad (10.29)$$

where X_2' and X_2'' are constructed spline variables (when $k = 4$). Plotting the estimated spline function $f(X_2)$ versus X_2 will estimate how the effect of X_2 should be modeled. If the sample is sufficiently large, the spline function can be fitted separately for $X_1 = 0$ and $X_1 = 1$, allowing detection of even unusual interaction patterns. A formal test of linearity in X_2 is obtained by testing $H_0 : \beta_3 = \beta_4 = 0$.

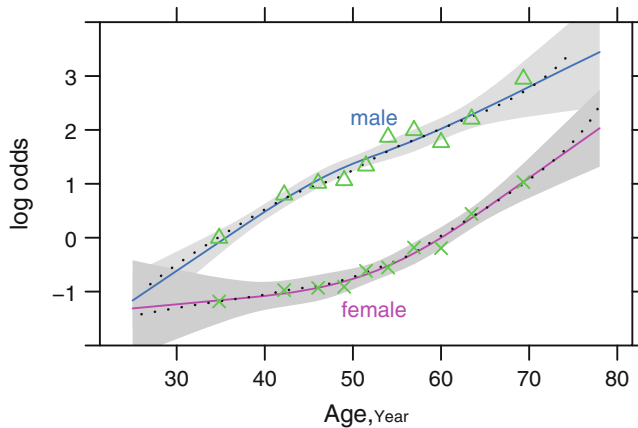


Fig. 10.6 Logit proportions of significant coronary artery disease by sex and deciles of age for $n=3504$ patients, with spline fits (smooth curves). Spline fits are for $k = 4$ knots at age= 36, 48, 56, and 68 years, and interaction between age and sex is allowed. Shaded bands are pointwise 0.95 confidence limits for predicted log odds. Smooth nonparametric estimates are shown as dotted curves. Data courtesy of the Duke Cardiovascular Disease Databank.

For testing interaction between X_1 and X_2 , a product term (e.g., X_1X_2) can be added to the model and its coefficient tested. A more general simultaneous test of linearity and lack of interaction for a two-variable model in which one variable is binary (or is assumed linear) is obtained by fitting the model

$$\begin{aligned} \text{logit}\{Y = 1|X\} = & \beta_0 + \beta_1X_1 + \beta_2X_2 + \beta_3X_1' + \beta_4X_2'' \\ & + \beta_5X_1X_2 + \beta_6X_1X_2' + \beta_7X_1X_2'' \end{aligned} \quad (10.30)$$

and testing $H_0 : \beta_3 = \dots = \beta_7 = 0$. This formulation allows the shape of the X_2 effect to be completely different for each level of X_1 . There is virtually no departure from linearity and additivity that cannot be detected from this expanded model formulation. The most computationally efficient test for lack of fit is the score test (e.g., X_1 and X_2 are forced into a tentative model and the remaining variables are candidates). Figure 10.6 also depicts a fitted spline logistic model with $k = 4$, allowing for general interaction between age and sex as parameterized above. The fitted function, after expanding the restricted cubic spline function for simplicity (see Equation 2.27), is given above. Note the good agreement between the empirical estimates of log odds and the spline fits and nonparametric estimates in this large dataset.

An analysis of log likelihood for this model and various sub-models is found in Table 10.3. The χ^2 for global tests is corrected for the intercept and the degrees of freedom does not include the intercept.

Table 10.3 LR χ^2 tests for coronary artery disease risk

Model / Hypothesis	Likelihood Ratio χ^2	d.f.	P	Formula
a: sex, age (linear, no interaction)	766.0	2		
b: sex, age, age \times sex	768.2	3		
c: sex, spline in age	769.4	4		
d: sex, spline in age, interaction	782.5	7		
H_0 : no age \times sex interaction given linearity	2.2	1	.14	$(b - a)$
H_0 : age linear no interaction	3.4	2	.18	$(c - a)$
H_0 : age linear, no interaction	16.6	5	.005	$(d - a)$
H_0 : age linear, product form interaction	14.4	4	.006	$(d - b)$
H_0 : no interaction, allowing for nonlinearity in age	13.1	3	.004	$(d - c)$

Table 10.4 AIC on χ^2 scale by number of knots

k	Model χ^2	AIC
0	99.23	97.23
3	112.69	108.69
4	121.30	115.30
5	123.51	115.51
6	124.41	114.51

This analysis confirms the first impression from the graph, namely, that age \times sex interaction is present but it is not of the form of a simple product between age and sex (change in slope). In the context of a linear age effect, there is no significant product interaction effect ($P = .14$). Without allowing for interaction, there is no significant nonlinear effect of age ($P = .18$). However, the general test of lack of fit with 5 d.f. indicates a significant departure from the linear additive model ($P = .005$).

In Figure 10.7, data from 2332 patients who underwent cardiac catheterization at Duke University Medical Center and were found to have significant ($\geq 75\%$) diameter narrowing of at least one major coronary artery were analyzed (the dataset is available from the Web site). The relationship between the time from the onset of symptoms of coronary artery disease (e.g., angina, myocardial infarction) to the probability that the patient has severe (three-vessel disease or left main disease—`tvdlm`) coronary disease was of interest. There were 1129 patients with `tvdlm`. A logistic model was used with the duration of symptoms appearing as a restricted cubic spline function with $k = 3, 4, 5$, and 6 equally spaced knots in terms of quantiles between .05 and .95. The best fit for the number of parameters was chosen using Akaike's information criterion (AIC), computed in Table 10.4 as the model likelihood

ratio χ^2 minus twice the number of parameters in the model aside from the intercept. The linear model is denoted $k = 0$.

```
dz ← subset(acath, sigdz==1)
dd ← datadist(dz)
```

```
f ← lrm(tvdlm ~ rcs(cad.dur, 5), data=dz)
w ← function(...)
  with(dz, {
    plsmo(cad.dur, tvdlm, fun=qlogis, add=TRUE,
          grid=TRUE, lty='dotted')
    x ← cut2(cad.dur, g=15, levels.mean=TRUE)
    prop ← qlogis(tapply(tvdlm, x, mean, na.rm=TRUE))
    xm ← as.numeric(names(prop))
    lpoints(xm, prop, pch=2, col='green')
  }) # Figure 10.7
plot(Predict(f, cad.dur), addpanel=w)
```

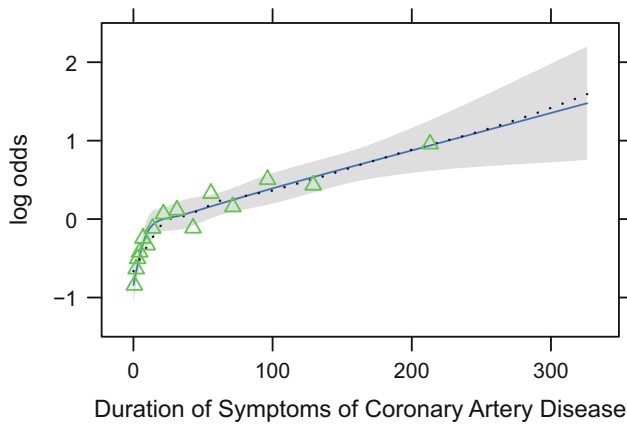


Fig. 10.7 Estimated relationship between duration of symptoms and the log odds of severe coronary artery disease for $k = 5$. Knots are marked with arrows. Solid line is spline fit; dotted line is a nonparametric loess estimate.

Figure 10.7 displays the spline fit for $k = 5$. The triangles represent subgroup estimates obtained by dividing the sample into groups of 150 patients. For example, the leftmost triangle represents the logit of the proportion of `tvdlm` in the 150 patients with the shortest duration of symptoms, versus the mean duration in that group. A Wald test of linearity, with 3 d.f., showed highly significant nonlinearity ($\chi^2 = 23.92$ with 3 d.f.). The plot of the spline transformation suggests a log transformation, and when $\log(\text{duration of symptoms in months} + 1)$ was fitted in a logistic model, the log likelihood of the model (119.33 with 1 d.f.) was virtually as good as the spline model (123.51 with 4 d.f.); the corresponding Akaike information criteria (on the χ^2 scale) are 117.33 and 115.51. To check for adequacy in the log transformation,

a five-knot restricted cubic spline function was fitted to $\log_{10}(\text{months} + 1)$, as displayed in Figure 10.8. There is some evidence for lack of fit on the right, but the Wald χ^2 for testing linearity yields $P = .27$.

```
f <- lrm(tvd1m ~ log10(cad.dur + 1), data=dz)
w <- function(...)
  with(dz, {
    x <- cut2(cad.dur, m=150, levels.mean=TRUE)
    prop <- tapply(tvd1m, x, mean, na.rm=TRUE)
    xm <- as.numeric(names(prop))
    lpoints(xm, prop, pch=2, col='green')
  })
# Figure 10.8
plot(Predict(f, cad.dur, fun=plogis), ylab='P',
     ylim=c(.2, .8), addpanel=w)
```

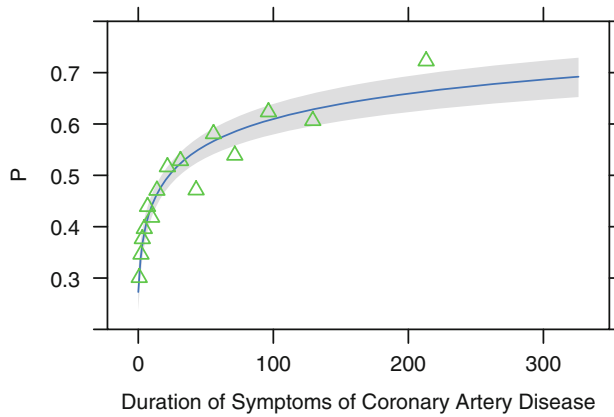


Fig. 10.8 Fitted linear logistic model in $\log_{10}(\text{duration} + 1)$, with subgroup estimates using groups of 150 patients. Fitted equation is $\text{logit}(\text{tvd1m}) = -.9809 + .7122 \log_{10}(\text{months} + 1)$.

If the model contains two continuous predictors, they may both be expanded with spline functions in order to test linearity or to describe nonlinear relationships. Testing interaction is more difficult here. If X_1 is continuous, one might temporarily group X_1 into quantile groups. Consider the subset of 2258 (1490 with disease) of the 3504 patients used in Figure 10.6 who have serum cholesterol measured. A logistic model for predicting significant coronary disease was fitted with age in tertiles (modeled with two dummy variables), sex, age \times sex interaction, four-knot restricted cubic spline in cholesterol, and age tertile \times cholesterol interaction. Except for the sex adjustment this model is equivalent to fitting three separate spline functions in cholesterol, one for each age tertile. The fitted model is shown in Figure 10.9 for cholesterol and age tertile against logit of significant disease. Significant age \times cholesterol interaction is apparent from the figure and is suggested by

the Wald χ^2 statistic (10.03) that follows. Note that the test for linearity of the interaction with respect to cholesterol is very insignificant ($\chi^2 = 2.40$ on 4 d.f.), but we retain it for now. The fitted function is

```
acath ← transform(acath,
                  cholesterol = choleste,
                  age.tertile = cut2(age,g=3),
                  sx = as.integer(acath$sex) - 1)
# sx for loess, need to code as numeric
dd ← datadist(acath); options(datadist='dd')

# First model stratifies age into tertiles to get more
# empirical estimates of age x cholesterol interaction

f ← lrm(sigdz ~ age.tertile*(sex + rcs(cholesterol,4)),
        data=acath)
print(f, latex=TRUE)
```

Logistic Regression Model

```
lrm(formula = sigdz ~ age.tertile * (sex + rcs(cholesterol, 4)),
    data = acath)
```

Frequencies of Missing Values Due to Each Variable

```
sigdz age.tertile          sex cholesterol
      0           0           0           1246
```

	Model Likelihood Ratio Test	Discrimination Indexes	Rank Discrim. Indexes
Obs	2258	$LR \chi^2$ 533.52	R^2 0.291
0	768	d.f. 14	C 0.780
1	1490	$Pr(> \chi^2) < 0.0001$	D_{xy} 0.560
$\max \left \frac{\partial \log L}{\partial \beta} \right $	2×10^{-8}	g_r 3.729	γ 0.562
		g_p 0.252	τ_a 0.251
		Brier 0.173	

	Coef	S.E.	Wald Z	$Pr(> Z)$
Intercept	-0.4155	1.0987	-0.38	0.7053
age.tertile=[49,58]	0.8781	1.7337	0.51	0.6125
age.tertile=[58,82]	4.7861	1.8143	2.64	0.0083
sex=female	-1.6123	0.1751	-9.21	< 0.0001
cholesterol	0.0029	0.0060	0.48	0.6347
cholesterol'	0.0384	0.0242	1.59	0.1126
cholesterol''	-0.1148	0.0768	-1.49	0.1350
age.tertile=[49,58] * sex=female	-0.7900	0.2537	-3.11	0.0018
age.tertile=[58,82] * sex=female	-0.4530	0.2978	-1.52	0.1283
age.tertile=[49,58] * cholesterol	0.0011	0.0095	0.11	0.9093

	Coef	S.E.	Wald Z	Pr(> Z)
age.tertile=[58,82] * cholesterol	-0.0158	0.0099	-1.59	0.1111
age.tertile=[49,58] * cholesterol'	-0.0183	0.0365	-0.50	0.6162
age.tertile=[58,82] * cholesterol'	0.0127	0.0406	0.31	0.7550
age.tertile=[49,58] * cholesterol''	0.0582	0.1140	0.51	0.6095
age.tertile=[58,82] * cholesterol''	-0.0092	0.1301	-0.07	0.9436

```
ltx(f)
```

$$X\hat{\beta} = -0.415 + 0.878[\text{age.tertile} \in [49, 58]] + 4.79[\text{age.tertile} \in [58, 82]] - 1.61[\text{female}] + 0.00287\text{cholesterol} + 1.52 \times 10^{-6}(\text{cholesterol} - 160)_+^3 - 4.53 \times 10^{-6}(\text{cholesterol} - 208)_+^3 + 3.44 \times 10^{-6}(\text{cholesterol} - 243)_+^3 - 4.28 \times 10^{-7}(\text{cholesterol} - 319)_+^3 + [\text{female}][-0.79[\text{age.tertile} \in [49, 58]] - 0.453[\text{age.tertile} \in [58, 82]]] + [\text{age.tertile} \in [49, 58]][0.00108\text{cholesterol} - 7.23 \times 10^{-7}(\text{cholesterol} - 160)_+^3 + 2.3 \times 10^{-6}(\text{cholesterol} - 208)_+^3 - 1.84 \times 10^{-6}(\text{cholesterol} - 243)_+^3 + 2.69 \times 10^{-7}(\text{cholesterol} - 319)_+^3] + [\text{age.tertile} \in [58, 82]][-0.0158\text{cholesterol} + 5 \times 10^{-7}(\text{cholesterol} - 160)_+^3 - 3.64 \times 10^{-7}(\text{cholesterol} - 208)_+^3 - 5.15 \times 10^{-7}(\text{cholesterol} - 243)_+^3 + 3.78 \times 10^{-7}(\text{cholesterol} - 319)_+^3].$$

```
# Table 10.5:
```

```
latex(anova(f), file='', size='smaller',
      caption='Crudely categorizing age into tertiles',
      label='tab:anova-tertiles')
```

```
y1 <- c(-1, 5)
```

```
plot(Predict(f, cholesterol, age.tertile),
     adj.subtitle=FALSE, ylim=y1) # Figure 10.9
```

Table 10.5 Crudely categorizing age into tertiles

	χ^2	d.f.	P
age.tertile (Factor+Higher Order Factors)	120.74	10	< 0.0001
<i>All Interactions</i>	21.87	8	0.0052
sex (Factor+Higher Order Factors)	329.54	3	< 0.0001
<i>All Interactions</i>	9.78	2	0.0075
cholesterol (Factor+Higher Order Factors)	93.75	9	< 0.0001
<i>All Interactions</i>	10.03	6	0.1235
<i>Nonlinear (Factor+Higher Order Factors)</i>	9.96	6	0.1263
age.tertile × sex (Factor+Higher Order Factors)	9.78	2	0.0075
age.tertile × cholesterol (Factor+Higher Order Factors)	10.03	6	0.1235
<i>Nonlinear</i>	2.62	4	0.6237
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	2.62	4	0.6237
TOTAL NONLINEAR	9.96	6	0.1263
TOTAL INTERACTION	21.87	8	0.0052
TOTAL NONLINEAR + INTERACTION	29.67	10	0.0010
TOTAL	410.75	14	< 0.0001

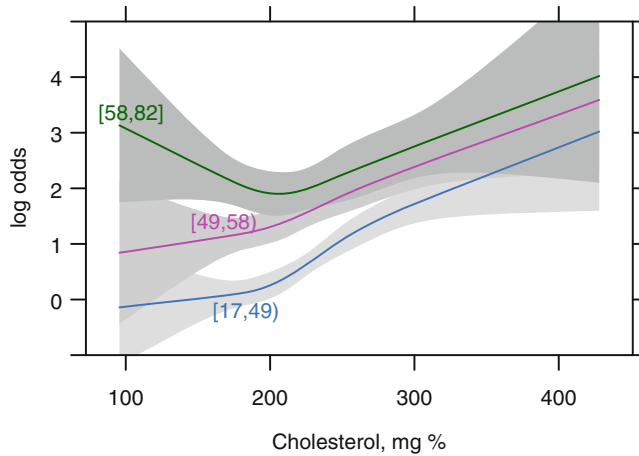


Fig. 10.9 Log odds of significant coronary artery disease modeling age with two dummy variables

Before fitting a parametric model that allows interaction between age and cholesterol, let us use the local regression model of Cleveland et al.⁹⁶ discussed in Section 2.4.7. This nonparametric smoothing method is not meant to handle binary Y , but it can still provide useful graphical displays in the binary case. Figure 10.10 depicts the fit from a local regression model predicting $Y = 1 =$ significant coronary artery disease. Predictors are sex (modeled parametrically with a dummy variable), age, and cholesterol, the last two fitted nonparametrically. The effect of not explicitly modeling a probability is seen in the figure, as the predicted probabilities exceeded 1. Because of this we do not take the logit transformation but leave the predicted values in raw form. However, the overall shape is in agreement with Figure 10.10.

```
# Re-do model with continuous age
f <- loess(sigmoid ~ age * (sx + cholesterol), data=acath,
          parametric="sx", drop.square="sx")
ages <- seq(25, 75, length=40)
chols <- seq(100, 400, length=40)
g <- expand.grid(cholesterol=chols, age=ages, sx=0)
# drop sex dimension of grid since held to 1 value
p <- drop(predict(f, g))
p[p < 0.001] <- 0.001
p[p > 0.999] <- 0.999
zl <- c(-3, 6) # Figure 10.10
wireframe(qlogis(p) ~ cholesterol*age,
          xlab=list(rot=30), ylab=list(rot=-40),
          zlab=list(label='log odds', rot=90), zlim=zl,
          scales = list(arrows = FALSE), data=g)
```

Chapter 2 discussed linear splines, which can be used to construct linear spline surfaces by adding all cross-products of the linear variables and spline terms in the model. With a sufficient number of knots for each predictor, the linear spline surface can fit a wide variety of patterns. However, it requires

a large number of parameters to be estimated. For the age–sex–cholesterol example, a linear spline surface is fitted for age and cholesterol, and a sex \times age spline interaction is also allowed. Figure 10.11 shows a fit that placed knots at quartiles of the two continuous variables^c. The algebraic form of the fitted model is shown below.

```
f ← lrm(sigdz ~ lsp(age, c(46, 52, 59)) *
        (sex + lsp(cholesterol, c(196, 224, 259))),
        data=acath)
ltx(f)
```

$$\begin{aligned}
 X\hat{\beta} = & -1.83 + 0.0232 \text{ age} + 0.0759(\text{age} - 46)_+ - 0.0025(\text{age} - 52)_+ + \\
 & 2.27(\text{age} - 59)_+ + 3.02[\text{female}] - 0.0177 \text{ cholesterol} + 0.114(\text{cholesterol} - 196)_+ - \\
 & 0.131(\text{cholesterol} - 224)_+ + 0.0651(\text{cholesterol} - 259)_+ + [\text{female}][-0.112 \text{ age} + \\
 & 0.0852(\text{age} - 46)_+ - 0.0302(\text{age} - 52)_+ + 0.176(\text{age} - 59)_+] + \text{age} \\
 & [0.000577 \text{ cholesterol} - 0.00286(\text{cholesterol} - 196)_+ + 0.00382(\text{cholesterol} - \\
 & 224)_+ - 0.00205(\text{cholesterol} - 259)_+] + (\text{age} - 46)_+ [-0.000936 \text{ cholesterol} + \\
 & 0.00643(\text{cholesterol} - 196)_+ - 0.0115(\text{cholesterol} - 224)_+ + 0.00756(\text{cholesterol} - \\
 & 259)_+] + (\text{age} - 52)_+ [0.000433 \text{ cholesterol} - 0.0037(\text{cholesterol} - 196)_+ + \\
 & 0.00815(\text{cholesterol} - 224)_+ - 0.00715(\text{cholesterol} - 259)_+] + (\text{age} - 59)_+ \\
 & [-0.0124 \text{ cholesterol} + 0.015(\text{cholesterol} - 196)_+ - 0.0067(\text{cholesterol} - 224)_+ + \\
 & 0.00752(\text{cholesterol} - 259)_+].
 \end{aligned}$$

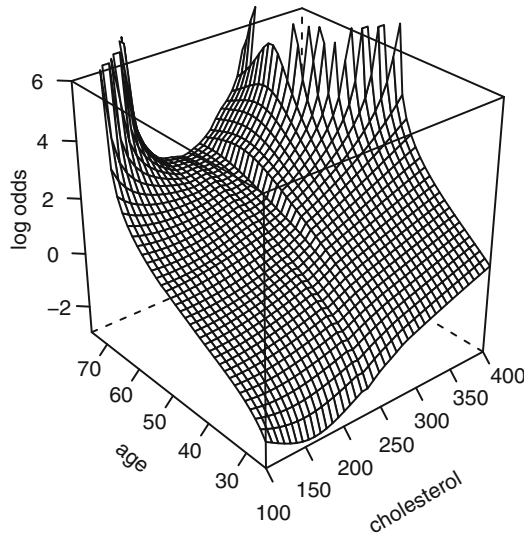


Fig. 10.10 Local regression fit for the logit of the probability of significant coronary disease vs. age and cholesterol for males, based on the `loess` function.

^c In the wireframe plots that follow, predictions for cholesterol–age combinations for which fewer than 5 exterior points exist are not shown, so as to not extrapolate to regions not supported by at least five points beyond the data perimeter.

```
latex(anova(f), caption='Linear spline surface', file='',
      size='smaller', label='tab:anova-lsp') # Table 10.6
```

```
perim ← with(acath,
             perimeter(cholesterol, age, xinc=20, n=5))
zl ← c(-2, 4) # Figure 10.11
bplot(Predict(f, cholesterol, age, np=40), perim=perim,
      lfun=wireframe, zlim=zl, adj.subtitle=FALSE)
```

Table 10.6 Linear spline surface

	χ^2	d.f.	P
age (Factor+Higher Order Factors)	164.17	24	< 0.0001
<i>All Interactions</i>	42.28	20	0.0025
<i>Nonlinear (Factor+Higher Order Factors)</i>	25.21	18	0.1192
sex (Factor+Higher Order Factors)	343.80	5	< 0.0001
<i>All Interactions</i>	23.90	4	0.0001
cholesterol (Factor+Higher Order Factors)	100.13	20	< 0.0001
<i>All Interactions</i>	16.27	16	0.4341
<i>Nonlinear (Factor+Higher Order Factors)</i>	16.35	15	0.3595
age × sex (Factor+Higher Order Factors)	23.90	4	0.0001
<i>Nonlinear</i>	12.97	3	0.0047
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	12.97	3	0.0047
age × cholesterol (Factor+Higher Order Factors)	16.27	16	0.4341
<i>Nonlinear</i>	11.45	15	0.7204
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	11.45	15	0.7204
<i>f(A,B) vs. Af(B) + Bg(A)</i>	9.38	9	0.4033
<i>Nonlinear Interaction in age vs. Af(B)</i>	9.99	12	0.6167
<i>Nonlinear Interaction in cholesterol vs. Bg(A)</i>	10.75	12	0.5503
TOTAL NONLINEAR	33.22	24	0.0995
TOTAL INTERACTION	42.28	20	0.0025
TOTAL NONLINEAR + INTERACTION	49.03	26	0.0041
TOTAL	449.26	29	< 0.0001

Chapter 2 also discussed a tensor spline extension of the restricted cubic spline model to fit a smooth function of two predictors, $f(X_1, X_2)$. Since this function allows for general interaction between X_1 and X_2 , the two-variable cubic spline is a powerful tool for displaying and testing interaction, assuming the sample size warrants estimating $2(k-1) + (k-1)^2$ parameters for a rectangular grid of $k \times k$ knots. Unlike the linear spline surface, the cubic surface is smooth. It also requires fewer parameters in most situations. The general cubic model with $k = 4$ (ignoring the sex effect here) is

$$\begin{aligned}
 & \beta_0 + \beta_1 X_1 + \beta_2 X_1' + \beta_3 X_1'' + \beta_4 X_2 + \beta_5 X_2' + \beta_6 X_2'' + \beta_7 X_1 X_2 \\
 + & \quad \beta_8 X_1 X_2' + \beta_9 X_1 X_2'' + \beta_{10} X_1' X_2 + \beta_{11} X_1' X_2' \\
 + & \quad + \beta_{12} X_1' X_2'' + \beta_{13} X_1'' X_2 + \beta_{14} X_1'' X_2' + \beta_{15} X_1'' X_2'',
 \end{aligned} \tag{10.31}$$

where $X'_1, X''_1, X'_2,$ and X''_2 are restricted cubic spline component variables for X_1 and X_2 for $k = 4$. A general test of interaction with 9 d.f. is $H_0 : \beta_7 = \dots = \beta_{15} = 0$. A test of adequacy of a simple product form interaction is $H_0 : \beta_8 = \dots = \beta_{15} = 0$ with 8 d.f. A 13 d.f. test of linearity and additivity is $H_0 : \beta_2 = \beta_3 = \beta_5 = \beta_6 = \beta_7 = \beta_8 = \beta_9 = \beta_{10} = \beta_{11} = \beta_{12} = \beta_{13} = \beta_{14} = \beta_{15} = 0$.

Figure 10.12 depicts the fit of this model. There is excellent agreement with Figures 10.9 and 10.11, including an increased (but probably insignificant) risk with low cholesterol for age ≥ 57 .

```
f ← lrm(sigdz ~ rcs(age,4)*(sex + rcs(cholesterol,4)),
      data=acath, tol=1e-11)
ltx(f)
```

$$X\hat{\beta} = -6.41 + 0.166\text{age} - 0.00067(\text{age} - 36)_+^3 + 0.00543(\text{age} - 48)_+^3 - 0.00727(\text{age} - 56)_+^3 + 0.00251(\text{age} - 68)_+^3 + 2.87[\text{female}] + 0.00979\text{cholesterol} + 1.96 \times 10^{-6}(\text{cholesterol} - 160)_+^3 - 7.16 \times 10^{-6}(\text{cholesterol} - 208)_+^3 + 6.35 \times 10^{-6}(\text{cholesterol} - 243)_+^3 - 1.16 \times 10^{-6}(\text{cholesterol} - 319)_+^3 + [\text{female}] [-0.109\text{age} + 7.52 \times 10^{-5}(\text{age} - 36)_+^3 + 0.00015(\text{age} - 48)_+^3 - 0.00045(\text{age} - 56)_+^3 + 0.000225(\text{age} - 68)_+^3] + \text{age} [-0.00028\text{cholesterol} + 2.68 \times 10^{-9}(\text{cholesterol} - 160)_+^3 + 3.03 \times 10^{-8}(\text{cholesterol} - 208)_+^3 - 4.99 \times 10^{-8}(\text{cholesterol} - 243)_+^3 + 1.69 \times 10^{-8}(\text{cholesterol} - 319)_+^3] + \text{age}' [0.00341\text{cholesterol} - 4.02 \times 10^{-7}(\text{cholesterol} - 160)_+^3 + 9.71 \times 10^{-7}(\text{cholesterol} - 208)_+^3 - 5.79 \times 10^{-7}(\text{cholesterol} - 243)_+^3 + 8.79 \times 10^{-9}(\text{cholesterol} - 319)_+^3] + \text{age}'' [-0.029\text{cholesterol} + 3.04 \times 10^{-6}(\text{cholesterol} -$$

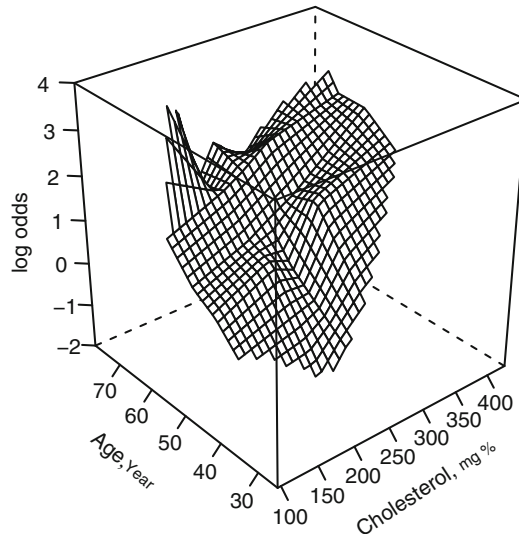


Fig. 10.11 Linear spline surface for males, with knots for age at 46, 52, 59 and knots for cholesterol at 196, 224, and 259 (quartiles).

$$160)_+^3 - 7.34 \times 10^{-6}(\text{cholesterol} - 208)_+^3 + 4.36 \times 10^{-6}(\text{cholesterol} - 243)_+^3 - 5.82 \times 10^{-8}(\text{cholesterol} - 319)_+^3].$$

```
latex(anova(f), caption='Cubic spline surface', file='',
      size='smaller', label='tab:anova-rcs') #Table 10.7
```

```
# Figure 10.12:
bplot(Predict(f, cholesterol, age, np=40), perim=perim,
      lfun=wireframe, zlim=z1, adj.subtitle=FALSE)
```

Table 10.7 Cubic spline surface

	χ^2	d.f.	P
age (Factor+Higher Order Factors)	165.23	15	< 0.0001
<i>All Interactions</i>	37.32	12	0.0002
<i>Nonlinear (Factor+Higher Order Factors)</i>	21.01	10	0.0210
sex (Factor+Higher Order Factors)	343.67	4	< 0.0001
<i>All Interactions</i>	23.31	3	< 0.0001
cholesterol (Factor+Higher Order Factors)	97.50	12	< 0.0001
<i>All Interactions</i>	12.95	9	0.1649
<i>Nonlinear (Factor+Higher Order Factors)</i>	13.62	8	0.0923
age × sex (Factor+Higher Order Factors)	23.31	3	< 0.0001
<i>Nonlinear</i>	13.37	2	0.0013
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	13.37	2	0.0013
age × cholesterol (Factor+Higher Order Factors)	12.95	9	0.1649
<i>Nonlinear</i>	7.27	8	0.5078
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	7.27	8	0.5078
<i>f(A,B) vs. Af(B) + Bg(A)</i>	5.41	4	0.2480
<i>Nonlinear Interaction in age vs. Af(B)</i>	6.44	6	0.3753
<i>Nonlinear Interaction in cholesterol vs. Bg(A)</i>	6.27	6	0.3931
TOTAL NONLINEAR	29.22	14	0.0097
TOTAL INTERACTION	37.32	12	0.0002
TOTAL NONLINEAR + INTERACTION	45.41	16	0.0001
TOTAL	450.88	19	< 0.0001

Statistics for testing age × cholesterol components of this fit are above. None of the nonlinear interaction components is significant, but we again retain them.

The general interaction model can be restricted to be of the form

$$f(X_1, X_2) = f_1(X_1) + f_2(X_2) + X_1g_2(X_2) + X_2g_1(X_1) \tag{10.32}$$

by removing the parameters $\beta_{11}, \beta_{12}, \beta_{14}$, and β_{15} from the model. The previous table of Wald statistics included a test of adequacy of this reduced form ($\chi^2 = 5.41$ on 4 d.f., $P = .248$). The resulting fit is in Figure 10.13.

```
f ← lrm(sigdz ~ sex*rca(age,4) + rca(cholesterol,4) +
      rca(age,4) %ia% rca(cholesterol,4), data=acath)
latex(anova(f), file='', size='smaller',
      caption='Singly nonlinear cubic spline surface',
      label='tab:anova-ria') #Table 10.8
```

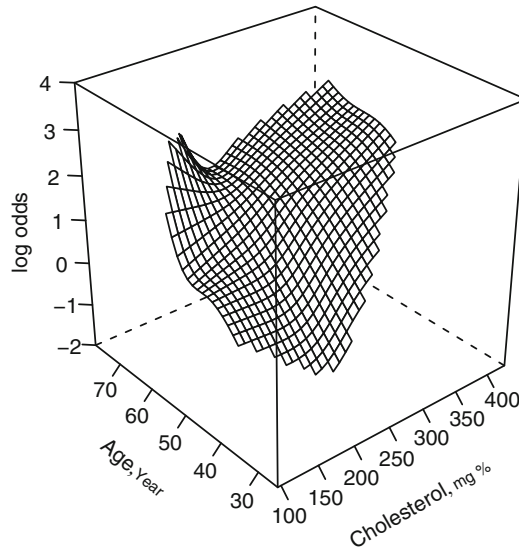



Fig. 10.12 Restricted cubic spline surface in two variables, each with $k = 4$ knots

Table 10.8 Singly nonlinear cubic spline surface

	χ^2	d.f.	P
sex (Factor+Higher Order Factors)	343.42	4	< 0.0001
<i>All Interactions</i>	24.05	3	< 0.0001
age (Factor+Higher Order Factors)	169.35	11	< 0.0001
<i>All Interactions</i>	34.80	8	< 0.0001
<i>Nonlinear (Factor+Higher Order Factors)</i>	16.55	6	0.0111
cholesterol (Factor+Higher Order Factors)	93.62	8	< 0.0001
<i>All Interactions</i>	10.83	5	0.0548
<i>Nonlinear (Factor+Higher Order Factors)</i>	10.87	4	0.0281
age \times cholesterol (Factor+Higher Order Factors)	10.83	5	0.0548
<i>Nonlinear</i>	3.12	4	0.5372
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	3.12	4	0.5372
<i>Nonlinear Interaction in age vs. Af(B)</i>	1.60	2	0.4496
<i>Nonlinear Interaction in cholesterol vs. Bg(A)</i>	1.64	2	0.4400
sex \times age (Factor+Higher Order Factors)	24.05	3	< 0.0001
<i>Nonlinear</i>	13.58	2	0.0011
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	13.58	2	0.0011
TOTAL NONLINEAR	27.89	10	0.0019
TOTAL INTERACTION	34.80	8	< 0.0001
TOTAL NONLINEAR + INTERACTION	45.45	12	< 0.0001
TOTAL	453.10	15	< 0.0001

Figure 10.13:

```
bplot(Predict(f, cholesterol, age, np=40),
      perim=perim,
      lfun=wireframe, zlim=z1, adj.subtitle=FALSE)
ltx(f)
```

Table 10.9 Linear interaction surface

	χ^2	d.f.	P
age (Factor+Higher Order Factors)	167.83	7	< 0.0001
<i>All Interactions</i>	31.03	4	< 0.0001
<i>Nonlinear (Factor+Higher Order Factors)</i>	14.58	4	0.0057
sex (Factor+Higher Order Factors)	345.88	4	< 0.0001
<i>All Interactions</i>	22.30	3	0.0001
cholesterol (Factor+Higher Order Factors)	89.37	4	< 0.0001
<i>All Interactions</i>	7.99	1	0.0047
<i>Nonlinear</i>	10.65	2	0.0049
age \times cholesterol (Factor+Higher Order Factors)	7.99	1	0.0047
age \times sex (Factor+Higher Order Factors)	22.30	3	0.0001
<i>Nonlinear</i>	12.06	2	0.0024
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	12.06	2	0.0024
TOTAL NONLINEAR	25.72	6	0.0003
TOTAL INTERACTION	31.03	4	< 0.0001
TOTAL NONLINEAR + INTERACTION	43.59	8	< 0.0001
TOTAL	452.75	11	< 0.0001

$$X\hat{\beta} = -7.2 + 2.96[\text{female}] + 0.164\text{age} + 7.23 \times 10^{-5}(\text{age} - 36)_+^3 - 0.000106(\text{age} - 48)_+^3 - 1.63 \times 10^{-5}(\text{age} - 56)_+^3 + 4.99 \times 10^{-5}(\text{age} - 68)_+^3 + 0.0148\text{cholesterol} + 1.21 \times 10^{-6}(\text{cholesterol} - 160)_+^3 - 5.5 \times 10^{-6}(\text{cholesterol} - 208)_+^3 + 5.5 \times 10^{-6}(\text{cholesterol} - 243)_+^3 - 1.21 \times 10^{-6}(\text{cholesterol} - 319)_+^3 + \text{age}[-0.00029\text{cholesterol} + 9.28 \times 10^{-9}(\text{cholesterol} - 160)_+^3 + 1.7 \times 10^{-8}(\text{cholesterol} - 208)_+^3 - 4.43 \times 10^{-8}(\text{cholesterol} - 243)_+^3 + 1.79 \times 10^{-8}(\text{cholesterol} - 319)_+^3] + \text{cholesterol}[2.3 \times 10^{-7}(\text{age} - 36)_+^3 + 4.21 \times 10^{-7}(\text{age} - 48)_+^3 - 1.31 \times 10^{-6}(\text{age} - 56)_+^3 + 6.64 \times 10^{-7}(\text{age} - 68)_+^3] + [\text{female}][-0.111\text{age} + 8.03 \times 10^{-5}(\text{age} - 36)_+^3 + 0.000135(\text{age} - 48)_+^3 - 0.00044(\text{age} - 56)_+^3 + 0.000224(\text{age} - 68)_+^3].$$

The fit is similar to the former one except that the climb in risk for low-cholesterol older subjects is less pronounced. The test for nonlinear interaction is now more concentrated ($P = .54$ with 4 d.f.). Figure 10.14 accordingly depicts a fit that allows age and cholesterol to have nonlinear main effects, but restricts the interaction to be a product between (untransformed) age and cholesterol. The function agrees substantially with the previous fit.

```
f ← lrm(sigdz ~ rcs(age,4)*sex + rcs(cholesterol,4) +
      age %ia% cholesterol, data=acath)
latex(anova(f), caption='Linear interaction surface', file='',
      size='smaller', label='tab:anova-lia') #Table 10.9
```

```
# Figure 10.14:
bplot(Predict(f, cholesterol, age, np=40), perim=perim,
      lfun=wireframe, zlim=z1, adj.subtitle=FALSE)
f.linia ← f # save linear interaction fit for later
ltx(f)
```

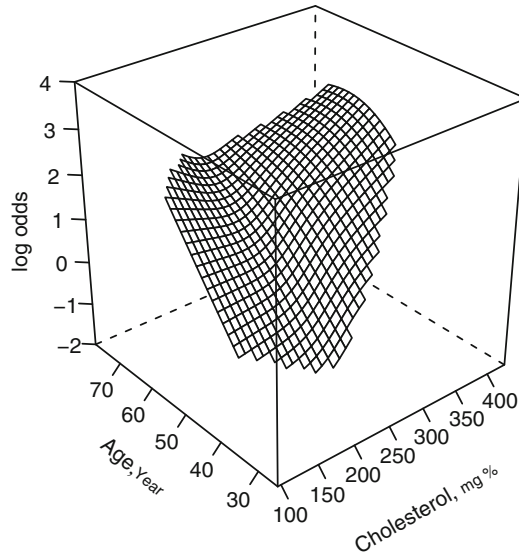


Fig. 10.13 Restricted cubic spline fit with age \times spline(cholesterol) and cholesterol \times spline(age)

$$X\hat{\beta} = -7.36 + 0.182\text{age} - 5.18 \times 10^{-5}(\text{age} - 36)_+^3 + 8.45 \times 10^{-5}(\text{age} - 48)_+^3 - 2.91 \times 10^{-6}(\text{age} - 56)_+^3 - 2.99 \times 10^{-5}(\text{age} - 68)_+^3 + 2.8[\text{female}] + 0.0139\text{cholesterol} + 1.76 \times 10^{-6}(\text{cholesterol} - 160)_+^3 - 4.88 \times 10^{-6}(\text{cholesterol} - 208)_+^3 + 3.45 \times 10^{-6}(\text{cholesterol} - 243)_+^3 - 3.26 \times 10^{-7}(\text{cholesterol} - 319)_+^3 - 0.00034 \text{age} \times \text{cholesterol} + [\text{female}][-0.107\text{age} + 7.71 \times 10^{-5}(\text{age} - 36)_+^3 + 0.000115(\text{age} - 48)_+^3 - 0.000398(\text{age} - 56)_+^3 + 0.000205(\text{age} - 68)_+^3].$$

The Wald test for age \times cholesterol interaction yields $\chi^2 = 7.99$ with 1 d.f., $P = .005$. These analyses favor the nonlinear model with simple product interaction in Figure 10.14 as best representing the relationships among cholesterol, age, and probability of prognostically severe coronary artery disease. A nomogram depicting this model is shown in Figure 10.21.

Using this simple product interaction model, Figure 10.15 displays predicted cholesterol effects at the mean age within each age tertile. Substantial agreement with Figure 10.9 is apparent.

```
# Make estimates of cholesterol effects for mean age in
# tertiles corresponding to initial analysis
mean.age <-
  with(acath,
    as.vector(tapply(age, age.tertile, mean, na.rm=TRUE)))
plot(Predict(f, cholesterol, age=round(mean.age, 2),
  sex="male"),
  adj.subtitle=FALSE, ylim=y1) #3 curves, Figure 10.15
```

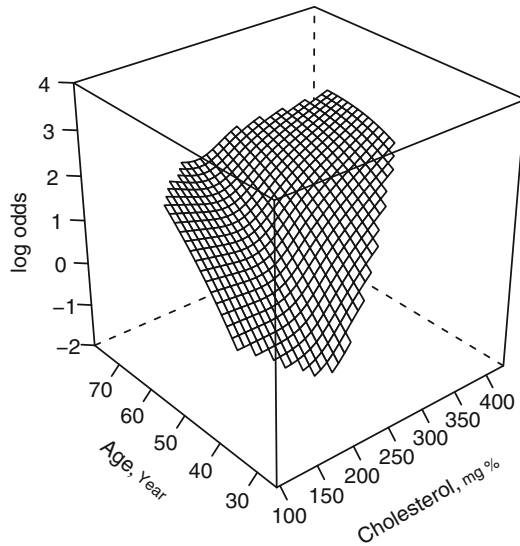


Fig. 10.14 Spline fit with nonlinear effects of cholesterol and age and a simple product interaction

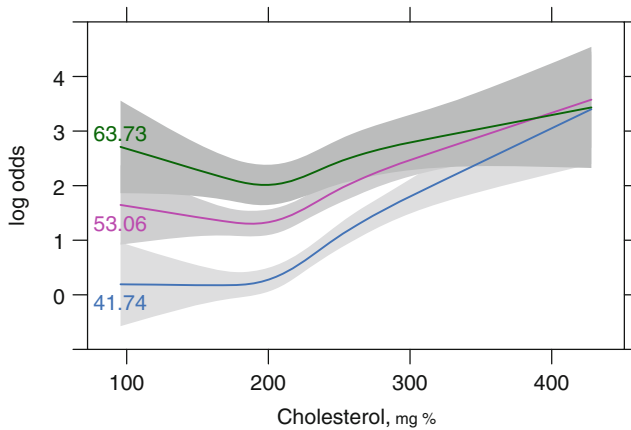


Fig. 10.15 Predictions from linear interaction model with mean age in tertiles indicated.

The partial residuals discussed in Section 10.4 can be used to check logistic model fit (although it may be difficult to deal with interactions). As an example, reconsider the “duration of symptoms” fit in Figure 10.7. Figure 10.16 displays “loess smoothed” and raw partial residuals for the original and log-transformed variable. The latter provides a more linear relationship, especially where the data are most dense.

Table 10.10 Merits of Methods for Checking Logistic Model Assumptions

Method	Choice Required	Assumes Additivity	Uses Ordering of X	Low Variance	Good Resolution on X
Stratification	Intervals				
Smoother on X_1 stratifying on X_2	Bandwidth		x (not on X_2)	x (if min. strat.)	x (X_1)
Smooth partial residual plot	Bandwidth	x	x	x	x
Spline model for all X s	Knots	x	x	x	x

```
f <- lrm(tvd1m ~ cad.dur, data=dz, x=TRUE, y=TRUE)
resid(f, "partial", pl="loess", xlim=c(0,250), ylim=c(-3,3))
scat1d(dz$cad.dur)
log.cad.dur <- log10(dz$cad.dur + 1)
f <- lrm(tvd1m ~ log.cad.dur, data=dz, x=TRUE, y=TRUE)
resid(f, "partial", pl="loess", ylim=c(-3,3))
scat1d(log.cad.dur) # Figure 10.16
```

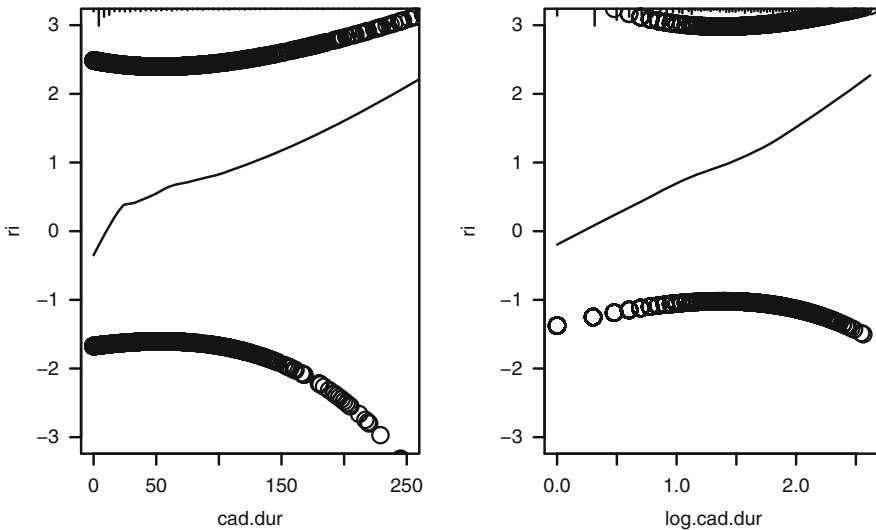


Fig. 10.16 Partial residuals for duration and $\log_{10}(\text{duration}+1)$. Data density shown at top of each plot.

Table 10.10 summarizes the relative merits of stratification, nonparametric smoothers, and regression splines for determining or checking binary logistic model fits.

10.6 Collinearity

The variance inflation factors (VIFs) discussed in Section 4.6 can apply to any regression fit.^{147,654} These VIFs allow the analyst to isolate which variable(s) are responsible for highly correlated parameter estimates. Recall that, in general, collinearity is not a large problem compared with nonlinearity and overfitting.

10.7 Overly Influential Observations

Pregibon⁵¹¹ developed a number of regression diagnostics that apply to the family of regression models of which logistic regression is a member. Influence statistics based on the “leave-out-one” method use an approximation to avoid having to refit the model n times for n observations. This approximation uses the fit and covariance matrix at the last iteration and assumes that the “weights” in the weighted least squares fit can be kept constant, yielding a computationally feasible one-step estimate of the leave-out-one regression coefficients.

Hosmer and Lemeshow [305, pp. 149–170] discuss many diagnostics for logistic regression and show how the final fit can be used in any least squares program that provides diagnostics. A new dependent variable to be used in that way is

$$Z_i = X\hat{\beta} + \frac{Y_i - \hat{P}_i}{V_i}, \quad (10.33)$$

where $V_i = \hat{P}_i(1 - \hat{P}_i)$, and $\hat{P}_i = [1 + \exp -X\hat{\beta}]^{-1}$ is the predicted probability that $Y_i = 1$. The $V_i, i = 1, 2, \dots, n$ are used as weights in an ordinary weighted least squares fit of X against Z . This least squares fit will provide regression coefficients identical to b . The new standard errors will be off from the actual logistic model ones by a constant.

As discussed in Section 4.9, the standardized change in the regression coefficients upon leaving out each observation in turn (DFBETAS) is one of the most useful diagnostics, as these can pinpoint which observations are influential on each part of the model. After carefully modeling predictor transformations, there should be no lack of fit due to improper transformations. However, as the white blood count example in Section 4.9 indicates, it is commonly the case that extreme predictor values can still have too much influence on the estimates of coefficients involving that predictor.

In the age–sex–response example of Section 10.1.3, both DFBETAS and DFFITS identified the same influential observations. The observation given by age = 48 sex = female response = 1 was influential for both age and sex, while the observation age = 34 sex = male response = 1 was influential for age and the observation age = 50 sex = male response = 0 was influential for sex. It can readily be seen from Figure 10.3 that these points do not fit the overall trends in the data. However, as these data were simulated from a

Table 10.11 Example influence statistics

Females				Males			
DFBETAS			DFFITS	DFBETAS			DFFITS
Intercept	Age	Sex		Intercept	Age	Sex	
0.0	0.0	0.0	0	0.5	-0.5	-0.2	2
0.0	0.0	0.0	0	0.2	-0.3	0.0	1
0.0	0.0	0.0	0	-0.1	0.1	0.0	-1
0.0	0.0	0.0	0	-0.1	0.1	0.0	-1
-0.1	0.1	0.1	0	-0.1	0.1	-0.1	-1
-0.1	0.1	0.1	0	0.0	0.0	0.1	0
0.7	-0.7	-0.8	3	0.0	0.0	0.1	0
-0.1	0.1	0.1	0	0.0	0.0	0.1	0
-0.1	0.1	0.1	0	0.0	0.0	-0.2	-1
-0.1	0.1	0.1	0	0.1	-0.1	-0.2	-1
-0.1	0.1	0.1	0	0.0	0.0	0.1	0
-0.1	0.0	0.1	0	-0.1	0.1	0.1	0
-0.1	0.0	0.1	0	-0.1	0.1	0.1	0
0.1	0.0	-0.2	1	0.3	-0.3	-0.4	-2
0.0	0.0	0.1	-1	-0.1	0.1	0.1	0
0.1	-0.2	0.0	-1	-0.1	0.1	0.1	0
-0.1	0.2	0.0	1	-0.1	0.1	0.1	0
-0.2	0.2	0.0	1	0.0	0.0	0.0	0
-0.2	0.2	0.0	1	0.0	0.0	0.0	0
-0.2	0.2	0.1	1	0.0	0.0	0.0	0

population model that is truly linear in age and additive in age and sex, the apparent influential observations are just random occurrences. It is unwise to assume that in real data all points will agree with overall trends. Removal of such points would bias the results, making the model apparently more predictive than it will be prospectively. See Table 10.11.

[11]

```
f ← update(fasr, x=TRUE, y=TRUE)
which.influence(f, .4) # Table 10.11
```

10.8 Quantifying Predictive Ability

The test statistics discussed above allow one to test whether a factor or set of factors is related to the response. If the sample is sufficiently large, a factor that grades risk from .01 to .02 may be a significant risk factor. However, that factor is not very useful in predicting the response for an individual subject. There is controversy regarding the appropriateness of R^2 from ordinary least squares in this setting.^{136,424} The generalized R_N^2 index of Nagelkerke⁴⁷¹ and Cragg and Uhler¹³⁷, Maddala⁴³¹, and Magee⁴³² described in Section 9.8.3 can be useful for quantifying the predictive strength of a model:

[12]

$$R_N^2 = \frac{1 - \exp(-LR/n)}{1 - \exp(-L^0/n)}, \quad (10.34)$$

where LR is the global log likelihood ratio statistic for testing the importance of all p predictors in the model and L^0 is the -2 log likelihood for the null model.

Tjur⁶¹³ coined the term “coefficient of discrimination” D , defined as the average \hat{P} when $Y = 1$ minus the average \hat{P} when $Y = 0$, and showed how it ties in with sum of squares–based R^2 measures. D has many advantages as an index of predictive power^d.

Linnet⁴¹⁶ advocates quadratic and logarithmic probability scoring rules for measuring predictive performance for probability models. Linnet shows how to bootstrap such measures to get bias-corrected estimates and how to use bootstrapping to compare two correlated scores. The quadratic scoring rule is Brier’s score, frequently used in judging meteorologic forecasts^{30,73}:

$$B = \frac{1}{n} \sum_{i=1}^n (\hat{P}_i - Y_i)^2, \quad (10.35)$$

where \hat{P}_i is the predicted probability and Y_i the corresponding observed response for the i th observation.

A unitless index of the strength of the rank correlation between predicted probability of response and actual response is a more interpretable measure of the fitted model’s predictive discrimination. One such index is the probability of concordance, c , between predicted probability and response. The c index, which is derived from the Wilcoxon–Mann–Whitney two-sample rank test, is computed by taking all possible pairs of subjects such that one subject responded and the other did not. The index is the proportion of such pairs with the responder having a higher predicted probability of response than the nonresponder.

Bamber³⁹ and Hanley and McNeil²⁵⁵ have shown that c is identical to a widely used measure of diagnostic discrimination, the area under a “receiver operating characteristic” (ROC) curve. A value of c of .5 indicates random predictions, and a value of 1 indicates perfect prediction (i.e., perfect separation of responders and nonresponders). A model having c greater than roughly .8 has some utility in predicting the responses of individual subjects. The concordance index is also related to another widely used index, Somers’ D_{xy} rank correlation⁵⁷⁹ between predicted probabilities and observed responses, by the identity

$$D_{xy} = 2(c - .5). \quad (10.36)$$

D_{xy} is the difference between concordance and discordance probabilities. When $D_{xy} = 0$, the model is making random predictions. When $D_{xy} = 1$,

^d Note that D and B (below) and other indexes not related to c (below) do not work well in case-control studies because of their reliance on absolute probability estimates.

13

14

the predictions are perfectly discriminating. These rank-based indexes have the advantage of being insensitive to the prevalence of positive responses.

15

A commonly used measure of predictive ability for binary logistic models is the fraction of correctly classified responses. Here one chooses a cutoff on the predicted probability of a positive response and then predicts that a response will be positive if the predicted probability exceeds this cutoff. There are a number of reasons why this measure should be avoided.

1. It's highly dependent on the cutpoint chosen for a "positive" prediction.
2. You can add a highly significant variable to the model and have the percentage classified correctly actually decrease. Classification error is a very insensitive and statistically inefficient measure^{264, 633} since if the threshold for "positive" is, say 0.75, a prediction of 0.99 rates the same as one of 0.751.
3. It gets away from the purpose of fitting a logistic model. A logistic model is a model for the probability of an event, not a model for the occurrence of the event. For example, suppose that the event we are predicting is the probability of being struck by lightning. Without having any data, we would predict that you won't get struck by lightning. However, you might develop an interesting model that discovers real risk factors that yield probabilities of being struck that range from 0.000000001 to 0.001.
4. If you make a classification rule from a probability model, you are being presumptuous. Suppose that a model is developed to assist physicians in diagnosing a disease. Physicians sometimes profess to desiring a binary decision model, but if given a probability they will rightfully apply different thresholds for treating different patients or for ordering other diagnostic tests. Even though the age of the patient may be a strong predictor of the probability of disease, the physician will often use a lower threshold of disease likelihood for treating a young patient. This usage is above and beyond how age affects the likelihood.
5. If a disease were present in only 0.02 of the population, one could be 0.98 accurate in diagnosing the disease by ruling that everyone is disease-free, i.e., by avoiding predictors. The proportion classified correctly fails to take the difficulty of the task into account.
6. van Houwelingen and le Cessie⁶³³ demonstrated a peculiar property that occurs when you try to obtain an honest estimate of classification error using cross-validation. The cross-validated error rate corrects the apparent error rate only if the predicted probability is exactly $1/2$ or is $1/2 \pm 1/(2n)$. The cross-validation estimate of optimism is "zero for n even and negligibly small for n odd." Better measures of error rate such as the Brier score and logarithmic scoring rule do not have this problem. They also have the nice property of being maximized when the predicted probabilities are the population probabilities.⁴¹⁶

16

10.9 Validating the Fitted Model

The major cause of unreliable models is overfitting the data. The methods described in Section 5.3 can be used to assess the accuracy of models fairly. If a sample has been held out and never used to study associations with the response, indexes of predictive accuracy can now be estimated using that sample. More efficient is cross-validation, and bootstrapping is the most efficient validation procedure. As discussed earlier, bootstrapping does not require holding out any data, since all aspects of model development (stepwise variable selection, tests of linearity, estimation of coefficients, etc.) are re-validated on samples taken with replacement from the whole sample.

Cox¹³⁰ proposed and Harrell and Lee²⁶⁷ and Miller et al.⁴⁵⁷ further developed the idea of fitting a new binary logistic model to a new sample to estimate the relationship between the predicted probability and the observed outcome in that sample. This fit provides a simple calibration equation that can be used to quantify unreliability (lack of calibration) and to calibrate the predictions for future use. This logistic calibration also leads to indexes of unreliability (U), discrimination (D), and overall quality ($Q = D - U$) which are derived from likelihood ratio tests²⁶⁷. Q is a logarithmic scoring rule, which can be compared with Brier's index (Equation 10.35). See [633] for many more ideas.

With bootstrapping we do not have a separate validation sample for assessing calibration, but we can estimate the overoptimism in assuming that the final model needs no calibration, that is, it has overall intercept=0 and slope=1. As discussed in Section 5.3, refitting the model

$$P_c = \text{Prob}\{Y = 1|X\hat{\beta}\} = [1 + \exp -(\gamma_0 + \gamma_1 X\hat{\beta})]^{-1} \quad (10.37)$$

(where P_c denotes the calibrated probability and the original predicted probability is $\hat{P} = [1 + \exp(-X\hat{\beta})]^{-1}$) in the original sample will always result in $\gamma = (\gamma_0, \gamma_1) = (0, 1)$, since a logistic model will always "fit" the training sample when assessed overall. We thus estimate γ by using Efron's¹⁷² method to estimate the overoptimism in $(0, 1)$ to obtain bias-corrected estimates of the true calibration. Simulations have shown this method produces an efficient estimate of γ .²⁵⁹

More stringent calibration checks can be made by running separate calibrations for different covariate levels. Smooth nonparametric curves described in Section 10.11 are more flexible than the linear-logit calibration method just described.

A good set of indexes to estimate for summarizing a model validation is the c or D_{xy} indexes and measures of calibration. In addition, the overoptimism in the indexes may be reported to quantify the amount of overfitting present. The estimate of γ can be used to draw a calibration curve by plotting \hat{P} on the x -axis and $\hat{P}_c = [1 + \exp -(\gamma_0 + \gamma_1 L)]^{-1}$ on the y -axis, where $L = \text{logit}(\hat{P})$.^{130,267} An easily interpreted index of unreliability, E_{max} , follows immediately from this calibration model:

$$E_{max}(a, b) = \max_{a \leq \hat{P} \leq b} |\hat{P} - \hat{P}_c|, \quad (10.38)$$

the maximum error in predicted probabilities over the range $a \leq \hat{P} \leq b$. In some cases, we would compute the maximum absolute difference in predicted and calibrated probabilities over the entire interval, that is, use $E_{max}(0, 1)$. The null hypothesis $H_0 : E_{max}(0, 1) = 0$ can easily be tested by testing $H_0 : \gamma_0 = 0, \gamma_1 = 1$ as above. Since E_{max} does not weight the discrepancies by the actual distribution of predictions, it may be preferable to compute the average absolute discrepancy over the actual distribution of predictions (or to use a mean squared error, incorporating the same calibration function).

If stepwise variable selection is being done, a matrix depicting which factors are selected at each bootstrap sample will shed light on how arbitrary is the selection of “significant” factors. See Section 5.3 for reasons to compare full and stepwise model fits.

As an example using bootstrapping to validate the calibration and discrimination of a model, consider the data in Section 10.1.3. Using 150 samples with replacement, we first validate the additive model with age and sex forced into every model. The optimism-corrected discrimination and calibration statistics produced by `validate` (see Section 10.11) are in the table below.

```
d ← sex.age.response
dd ← datadist(d); options(datadist='dd')
f ← lrm(response ~ sex + age, data=d, x=TRUE, y=TRUE)
set.seed(3) # for reproducibility
v1 ← validate(f, B=150)
```

```
latex(v1,
      caption='Bootstrap Validation, 2 Predictors Without
Stepdown', digits=2, size='Ssize', file='')
```

Bootstrap Validation, 2 Predictors Without Stepdown					
Index	Original Training Sample	Test Sample	Optimism Corrected		n
D_{xy}	0.70	0.70	0.67	0.04	0.66 150
R^2	0.45	0.48	0.43	0.05	0.40 150
Intercept	0.00	0.00	0.01	-0.01	0.01 150
Slope	1.00	1.00	0.91	0.09	0.91 150
E_{max}	0.00	0.00	0.02	0.02	0.02 150
D	0.39	0.44	0.36	0.07	0.32 150
U	-0.05	-0.05	0.04	-0.09	0.04 150
Q	0.44	0.49	0.32	0.16	0.28 150
B	0.16	0.15	0.18	-0.03	0.19 150
g	2.10	2.49	1.97	0.52	1.58 150
g_p	0.35	0.35	0.34	0.01	0.34 150

Now we incorporate variable selection. The variables selected in the first 10 bootstrap replications are shown below. The apparent Somers' D_{xy} is 0.7, and the bias-corrected D_{xy} is 0.66. The slope shrinkage factor is 0.91. The maximum absolute error in predicted probability is estimated to be 0.02.

We next allow for step-down variable selection at each resample. For illustration purposes only, we use a suboptimal stopping rule based on significance of *individual* variables at the $\alpha = 0.10$ level. Of the 150 repetitions, both age and sex were selected in 137, and neither variable was selected in 3 samples. The validation statistics are in the table below.

```
v2 ← validate(f, B=150, bw=TRUE,
             rule='p', sls=.1, type='individual')
```

```
latex(v2,
      caption='Bootstrap Validation, 2 Predictors with Stepdown',
      digits=2, B=15, file='', size='Ssize')
```

Bootstrap Validation, 2 Predictors with Stepdown					
Index	Original Training Sample	Test Sample	Optimism	Corrected	n
D_{xy}	0.70	0.70	0.64	0.07	0.63 150
R^2	0.45	0.49	0.41	0.09	0.37 150
Intercept	0.00	0.00	-0.04	0.04	-0.04 150
Slope	1.00	1.00	0.84	0.16	0.84 150
E_{\max}	0.00	0.00	0.05	0.05	0.05 150
D	0.39	0.45	0.34	0.11	0.28 150
U	-0.05	-0.05	0.06	-0.11	0.06 150
Q	0.44	0.50	0.28	0.22	0.22 150
B	0.16	0.14	0.18	-0.04	0.20 150
g	2.10	2.60	1.88	0.72	1.38 150
g_p	0.35	0.35	0.33	0.02	0.33 150

Factors Retained in Backwards Elimination
First 15 Resamples


```
function(x) paste(v[1:2][x], collapse=', ')
table(paste(as, ' ', nx, 'Xs'))
```

	0 Xs		1 Xs	age	2 Xs	age, sex	0 Xs
	50		3		1		34
age, sex	1 Xs	age, sex	2 Xs	age, sex	3 Xs	age, sex	4 Xs
	17		11		7		1
sex	0 Xs	sex	1 Xs				
	12		3				

```
latex(v3,
caption='Bootstrap Validation with 5 Noise Variables and
Stepdown', digits=2, B=15, size='Ssize', file='')
```

Bootstrap Validation with 5 Noise Variables and Stepdown

Index	Original Training Sample	Test Sample	Optimism	Corrected	n
D_{xy}	0.70	0.47	0.38	0.09	0.60 139
R^2	0.45	0.34	0.23	0.11	0.34 139
Intercept	0.00	0.00	0.03	-0.03	0.03 139
Slope	1.00	1.00	0.78	0.22	0.78 139
E_{max}	0.00	0.00	0.06	0.06	0.06 139
D	0.39	0.31	0.18	0.13	0.26 139
U	-0.05	-0.05	0.07	-0.12	0.07 139
Q	0.44	0.36	0.11	0.25	0.19 139
B	0.16	0.17	0.22	-0.04	0.20 139
g	2.10	1.81	1.06	0.75	1.36 139
g_p	0.35	0.23	0.19	0.04	0.31 139

Factors Retained in Backwards Elimination

First 15 Resamples

age	sex	x1	x2	x3	x4	x5
•	•			•	•	•
•	•	•				•
•	•					
•	•				•	•
•	•	•				
•	•					
•	•			•		
•	•					
•	•			•		

Frequencies of Numbers of Factors Retained

0	1	2	3	4	5	6
50	15	37	18	11	7	1

Using step-down variable selection with the same stopping rule as before, the “final” model on the original sample correctly deleted x_1, \dots, x_5 . Of the 150 bootstrap repetitions, 11 samples yielded a singularity or non-convergence either in the full-model fit or after step-down variable selection. Of the 139 successful repetitions, the frequencies of the number of factors selected, as well as the frequency of variable combinations selected, are shown above. Validation statistics are also shown above.

Figure 10.17 depicts the calibration (reliability) curves for the three strategies using the corrected intercept and slope estimates in the above tables as γ_0 and γ_1 , and the logistic calibration model $P_c = [1 + \exp -(\gamma_0 + \gamma_1 L)]^{-1}$, where P_c is the “actual” or calibrated probability, L is $\text{logit}(\hat{P})$, and \hat{P} is the predicted probability. The shape of the calibration curves (driven by slopes < 1) is typical of overfitting—low predicted probabilities are too low and high predicted probabilities are too high. Predictions near the overall prevalence of the outcome tend to be calibrated even when overfitting is present.

```
g <- function(v) v[c('Intercept', 'Slope'), 'index.corrected']
k <- rbind(g(v1), g(v2), g(v3))
co <- c(2,5,4,1)
plot(0, 0, ylim=c(0,1), xlim=c(0,1),
     xlab="Predicted Probability",
     ylab="Estimated Actual Probability", type="n")
legend(.45, .35, c("age, sex", "age, sex stepdown",
                  "age, sex, x1-x5", "ideal"),
       lty=1, col=co, cex=.8, bty="n")
probs <- seq(0, 1, length=200); L <- qlogis(probs)
for(i in 1:3) {
  P <- plogis(k[i, 'Intercept'] + k[i, 'Slope'] * L)
  lines(probs, P, col=co[i], lwd=1)
}
abline(a=0, b=1, col=co[4], lwd=1) # Figure 10.17
```

“Honest” calibration curves may also be estimated using nonparametric smoothers in conjunction with bootstrapping and cross-validation (see Section 10.11).

10.10 Describing the Fitted Model

Once the proper variables have been modeled and all model assumptions have been met, the analyst needs to present and interpret the fitted model. There are at least three ways to proceed. The coefficients in the model may be interpreted. For each variable, the change in log odds for a sensible change in the variable value (e.g., interquartile range) may be computed. Also, the odds

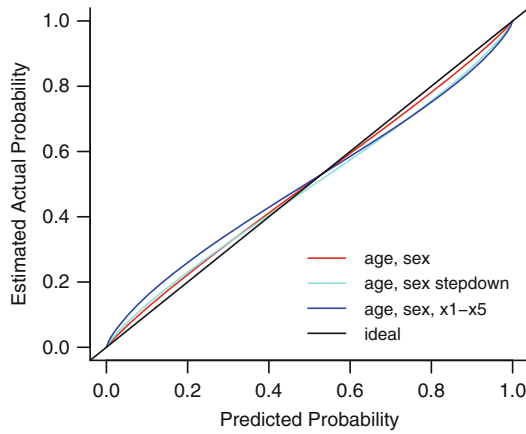


Fig. 10.17 Estimated logistic calibration (reliability) curves obtained by bootstrapping three modeling strategies.

Table 10.12 Effects Response : sigdz

	Low	High	Δ	Effect	S.E.	Lower 0.95	Upper 0.95
age	46	59	13	0.90629	0.18381	0.546030	1.26650
<i>Odds Ratio</i>	46	59	13	2.47510		1.726400	3.54860
cholesterol	196	259	63	0.75479	0.13642	0.487410	1.02220
<i>Odds Ratio</i>	196	259	63	2.12720		1.628100	2.77920
sex — female:male	1	2		-2.42970	0.14839	-2.720600	-2.13890
<i>Odds Ratio</i>	1	2		0.08806		0.065837	0.11778

ratio or factor by which the odds increases for a certain change in a predictor, holding all other predictors constant, may be displayed. Table 10.12 contains such summary statistics for the linear age \times cholesterol interaction surface fit described in Section 10.5.

```
s ← summary(f.linia) # Table 10.12
latex(s, file='', size='Ssize',
      label='tab:lrn-cholexage-confbar')
```

```
plot(s) # Figure 10.18
```

The outer quartiles of age are 46 and 59 years, so the “half-sample” odds ratio for age is 2.47, with 0.95 confidence interval [1.63, 3.74] when sex is male and cholesterol is set to its median. The effect of increasing cholesterol from 196 (its lower quartile) to 259 (its upper quartile) is to increase the log odds by 0.79 or to increase the odds by a factor of 2.21. Since there are interactions allowed between age and sex and between age and cholesterol, each odds ratio in the above table depends on the setting of at least one other factor. The

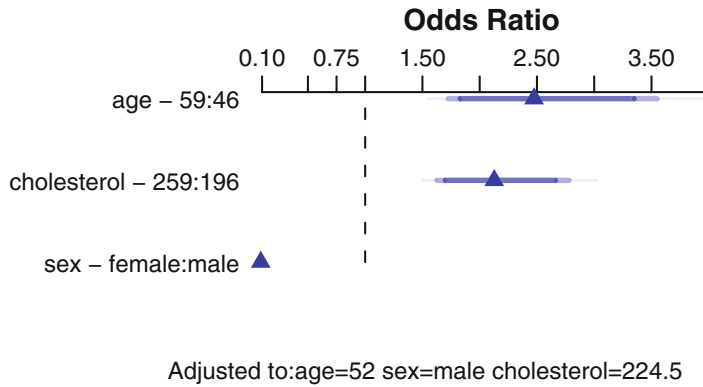


Fig. 10.18 Odds ratios and confidence bars, using quartiles of age and cholesterol for assessing their effects on the odds of coronary disease

results are shown graphically in Figure 10.18. The shaded confidence bars show various levels of confidence and do not pin the analyst down to, say, the 0.95 level.

For those used to thinking in terms of odds or log odds, the preceding description may be sufficient. Many prefer instead to interpret the model in terms of predicted probabilities instead of odds. If the model contains only a single predictor (even if several spline terms are required to represent that predictor), one may simply plot the predictor against the predicted response. Such a plot is shown in Figure 10.19 which depicts the fitted relationship between age of diagnosis and the probability of acute bacterial meningitis (ABM) as opposed to acute viral meningitis (AVM), based on an analysis of 422 cases from Duke University Medical Center.⁵⁸⁰ The data may be found on the web site. A linear spline function with knots at 1, 2, and 22 years was used to model this relationship.

When the model contains more than one predictor, one may graph the predictor against log odds, and barring interactions, the shape of this relationship will be independent of the level of the other predictors. When displaying the model on what is usually a more interpretable scale, the probability scale, a difficulty arises in that unlike log odds the relationship between one predictor and the probability of response depends on the levels of all other factors. For example, in the model

$$\text{Prob}\{Y = 1|X\} = \{1 + \exp[-(\beta_0 + \beta_1 X_1 + \beta_2 X_2)]\}^{-1} \quad (10.39)$$

there is no way to factor out X_1 when examining the relationship between X_2 and the probability of a response. For the two-predictor case one can plot X_2 versus predicted probability for each level of X_1 . When it is uncertain whether to include an interaction in this model, consider presenting graphs for two models (with and without interaction terms included) as was done in [658].

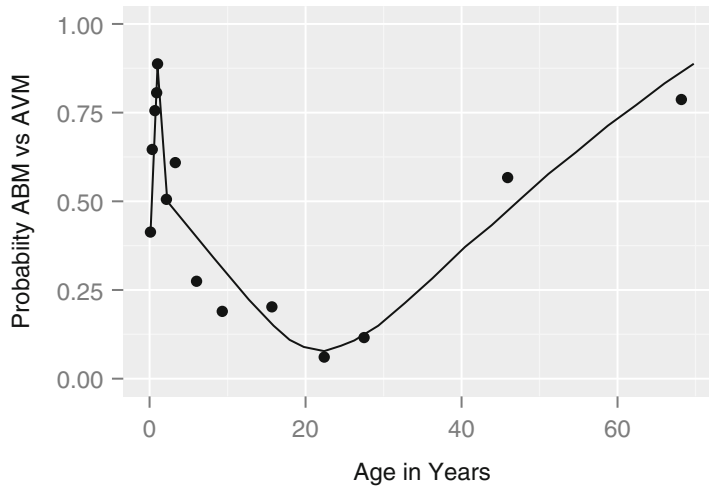


Fig. 10.19 Linear spline fit for probability of bacterial versus viral meningitis as a function of age at onset⁵⁸⁰. Points are simple proportions by age quantile groups.

When three factors are present, one could draw a separate graph for each level of X_3 , a separate curve on each graph for each level of X_1 , and vary X_2 on the x -axis. Instead of this, or if more than three factors are present, a good way to display the results may be to plot “adjusted probability estimates” as a function of one predictor, adjusting all other factors to constants such as the mean. For example, one could display a graph relating serum cholesterol to probability of myocardial infarction or death, holding age constant at 55, sex at 1 (male), and systolic blood pressure at 120 mmHg.

The final method for displaying the relationship between several predictors and probability of response is to construct a nomogram.^{40,254} A nomogram not only sheds light on how the effect of one predictor on the probability of response depends on the levels of other factors, but it allows one to quickly estimate the probability of response for individual subjects. The nomogram in Figure 10.20 allows one to predict the probability of acute bacterial meningitis (given the patient has either viral or bacterial meningitis) using the same sample as in Figure 10.19. Here there are four continuous predictor values, none of which are linearly related to log odds of bacterial meningitis: age at admission (expressed as a linear spline function), month of admission (expressed as $|\text{month} - 8|$), cerebrospinal fluid glucose/blood glucose ratio (linear effect truncated at .6; that is, the effect is the glucose ratio if it is $\leq .6$, and .6 if it exceeded .6), and the cube root of the total number of polymorphonuclear leukocytes in the cerebrospinal fluid.

17

The model associated with Figure 10.14 is depicted in what could be called a “precision nomogram” in Figure 10.21. Discrete cholesterol levels were required because of the interaction between two continuous variables.

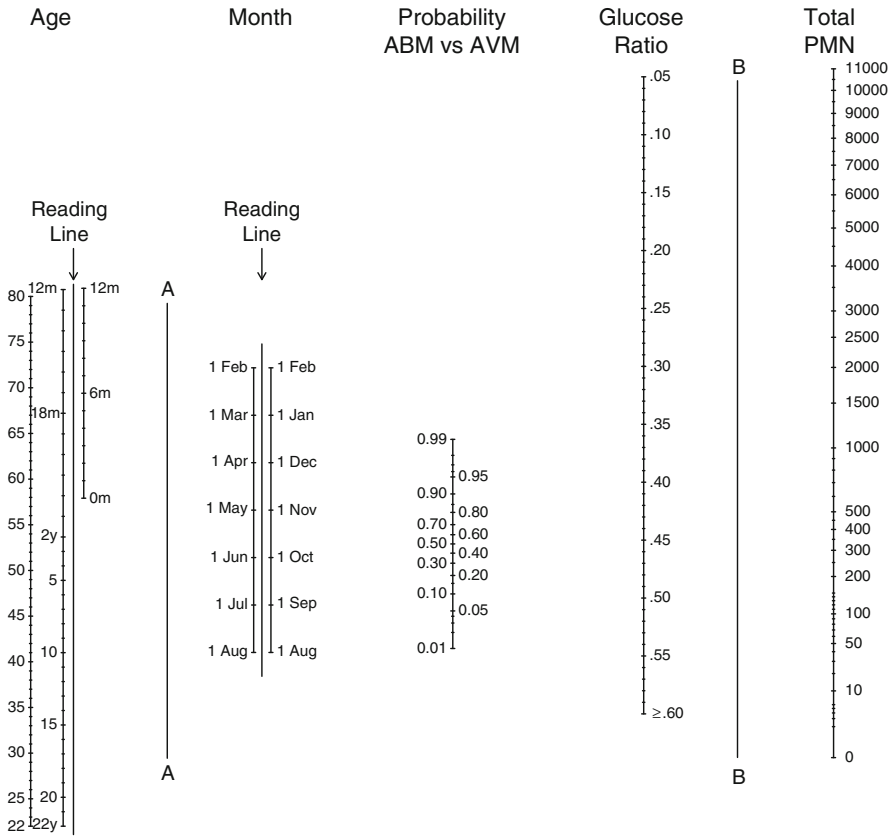


Fig. 10.20 Nomogram for estimating probability of bacterial (ABM) versus viral (AVM) meningitis. Step 1, place ruler on reading lines for patient's age and month of presentation and mark intersection with line A; step 2, place ruler on values for glucose ratio and total polymorphonuclear leukocyte (PMN) count in cerebrospinal fluid and mark intersection with line B; step 3, use ruler to join marks on lines A and B, then read off the probability of ABM versus AVM.⁵⁸⁰

```
# Draw a nomogram that shows examples of confidence intervals
nom ← nomogram(f.linia, cholesterol=seq(150, 400, by=50),
  interact=list(age=seq(30, 70, by=10)),
  lp.at=seq(-2, 3.5, by=.5),
  conf.int=TRUE, conf.lp="all",
  fun=function(x)1/(1+exp(-x)), # or plogis
  funlabel="Probability of CAD",
  fun.at=c(seq(.1, .9, by=.1), .95, .99)
) # Figure 10.21
plot(nom, col.grid = gray(c(0.8, 0.95)),
  varname.label=FALSE, ia.space=1, xfrac=.46, lmgp=.2)
```

10.11 R Functions

The general R statistical modeling functions⁹⁶ described in Section 6.2 work with the author's `lrm` function for fitting binary and ordinal logistic regression models. `lrm` has several options for doing penalized maximum likelihood estimation, with special treatment of categorical predictors so as to shrink all estimates (including the reference cell) to the mean. The following example fits a logistic model containing predictors `age`, `blood.pressure`, and `sex`, with `age` fitted with a smooth five-knot restricted cubic spline function and a different shape of the age relationship for males and females.

18

```
fit ← lrm(death ~ blood.pressure + sex * rcs(age,5))
anova(fit)
plot(Predict(fit, age, sex))
```

The `pentrace` function makes it easy to check the effects of a sequence of penalties. The following code fits an unpenalized model and plots the AIC and Schwarz BIC for a variety of penalties so that approximately the best cross-validating model can be chosen (and so we can learn how the penalty relates to the effective degrees of freedom). Here we elect to only penalize the nonlinear or non-additive parts of the model.

```
f ← lrm(death ~ rcs(age,5)*treatment + lsp(sbp,c(120,140)),
        x=TRUE, y=TRUE)
plot(pentrace(f,
             penalty=list(nonlinear=seq(.25,10,by=.25))) )
```

See Sections 9.8.1 and 9.10 for more information.

19

The `residuals` function for `lrm` and the `which.influence` function can be used to check predictor transformations as well as to analyze overly influential observations in binary logistic regression. See Figure 10.16 for one application. The `residuals` function will also perform the unweighted sum of squares test for global goodness of fit described in Section 10.5.

The `validate` function when used on an object created by `lrm` does resampling validation of a logistic regression model, with or without backward step-down variable deletion. It provides bias-corrected Somers' D_{xy} rank correlation, R_N^2 index, the intercept and slope of an overall logistic calibration equation, the maximum absolute difference in predicted and calibrated probabilities E_{max} , the discrimination index D [(model L.R. $\chi^2 - 1$)/ n], the unreliability index U = (difference in -2 log likelihood between uncalibrated $X\beta$ and $X\beta$ with overall intercept and slope calibrated to test sample)/ n , and the overall quality index $Q = D - U$.²⁶⁷ The "corrected" slope can be thought of as a shrinkage factor that takes overfitting into account. See `predab.resample` in Section 6.2 for the list of resampling methods.

The `calibrate` function produces bootstrapped or cross-validated calibration curves for logistic and linear models. The "apparent" calibration accuracy is estimated using a nonparametric smoother relating predicted probabilities

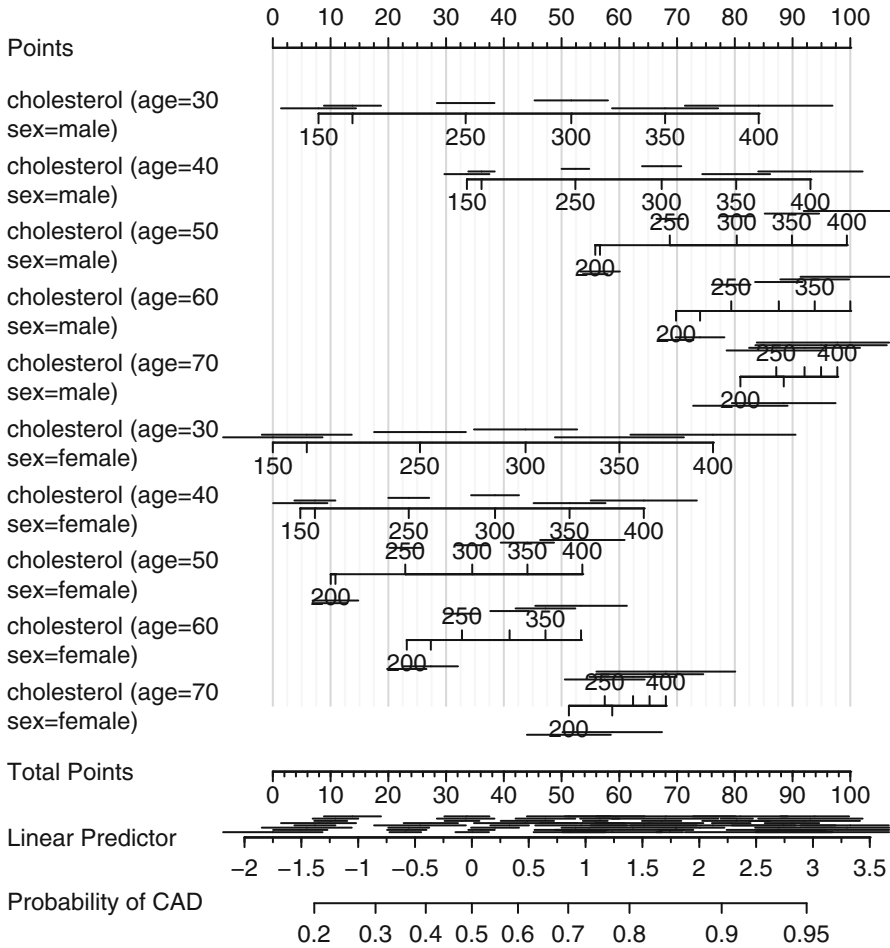


Fig. 10.21 Nomogram relating age, sex, and cholesterol to the log odds and to the probability of significant coronary artery disease. Select one axis corresponding to sex and to age $\in \{30, 40, 50, 60, 70\}$. There is linear interaction between age and sex and between age and cholesterol. 0.70 and 0.90 confidence intervals are shown (0.90 in gray). Note that for the “Linear Predictor” scale there are various lengths of confidence intervals near the same value of $X\hat{\beta}$, demonstrating that the standard error of $X\hat{\beta}$ depends on the individual X values. Also note that confidence intervals corresponding to smaller patient groups (e.g., females) are wider.

to observed binary outcomes. The nonparametric estimate is evaluated at a sequence of predicted probability levels. Then the distances from the 45° line are compared with the differences when the current model is evaluated back on the whole sample (or omitted sample for cross-validation). The differences in the differences are estimates of overoptimism. After averaging over many replications, the predicted-value-specific differences are then subtracted from

the apparent differences and an adjusted calibration curve is obtained. Unlike `validate`, `calibrate` does not assume a linear logistic calibration. For an example, see the end of Chapter 11. `calibrate` will print the mean absolute calibration error, the 0.9 quantile of the absolute error, and the mean squared error, all over the observed distribution of predicted values.

The `val.prob` function is used to compute measures of discrimination and calibration of predicted probabilities for a separate sample from the one used to derive the probability estimates. Thus `val.prob` is used in external validation and data-splitting. The function computes similar indexes as `validate` plus the Brier score and a statistic for testing for unreliability or $H_0 : \gamma_0 = 0, \gamma_1 = 1$.

In the following example, a logistic model is fitted on 100 observations simulated from the actual model given by

$$\text{Prob}\{Y = 1 | X_1, X_2, X_3\} = [1 + \exp[-(-1 + 2X_1)]]^{-1}, \quad (10.40)$$

where X_1 is a random uniform $[0, 1]$ variable. Hence X_2 and X_3 are irrelevant. After fitting a linear additive model in X_1, X_2 , and X_3 , the coefficients are used to predict $\text{Prob}\{Y = 1\}$ on a separate sample of 100 observations.

```
set.seed(13)
n <- 200
x1 <- runif(n)
x2 <- runif(n)
x3 <- runif(n)
logit <- 2*(x1-.5)
P <- 1/(1+exp(-logit))
y <- ifelse(runif(n) <= P, 1, 0)
d <- data.frame(x1, x2, x3, y)
f <- lrm(y ~ x1 + x2 + x3, subset=1:100)
phat <- predict(f, d[101:200,], type='fitted')
# Figure 10.22
v <- val.prob(phat, y[101:200], m=20, cex=.5)
```

The output is shown in Figure 10.22.

The R built-in function `glm`, a very general modeling function, can fit binary logistic models. The response variable *must* be coded 0/1 for `glm` to work. `Glm` is a slight modification of the built-in `glm` function in the `rms` package that allows fits to use `rms` methods. This facilitates Poisson and several other types of regression analysis.

10.12 Further Reading

- [1] See [590] for modeling strategies specific to binary logistic regression.
- [2] See [632] for a nice review of logistic modeling. Agresti⁶ is an excellent source for categorical Y in general.
- [3] Not only does discriminant analysis assume the same regression model as logistic regression, but it also assumes that the predictors are each normally distributed and that jointly the predictors have a multivariate normal distribution. These assumptions are unlikely to be met in practice, especially when

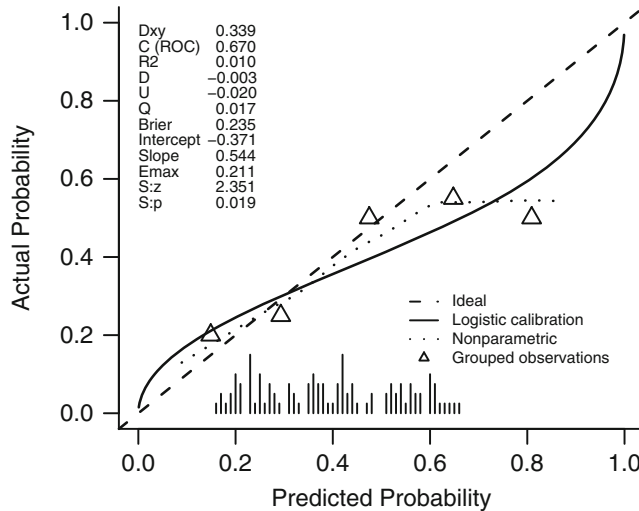


Fig. 10.22 Validation of a logistic model in a test sample of size $n = 100$. The calibrated risk distribution (histogram of logistic-calibrated probabilities) is shown.

one of the predictors is a discrete variable such as sex group. When discriminant analysis assumptions are violated, logistic regression yields more accurate estimates.^{251,514} Even when discriminant analysis is optimal (i.e., when all its assumptions are satisfied) logistic regression is virtually as accurate as the discriminant model.²⁶⁴

- [4] See [573] for a review of measures of effect for binary outcomes.
- [5] Cepedaet al.⁹⁵ found that propensity adjustment is better than covariate adjustment with logistic models when the number of events per variable is less than 8.
- [6] Pregibon⁵¹² developed a modification of the log likelihood function that when maximized results in a fit that is resistant to overly influential and outlying observations.
- [7] See Hosmer and Lemeshow³⁰⁶ for methods of testing for a difference in the observed event proportion and the predicted event probability (average of predicted probabilities) for a group of heterogeneous subjects.
- [8] See Hosmer and Lemeshow,³⁰⁵ Kay and Little,³⁴¹ and Collett [115, Chap. 5]. Landwehr et al.³⁷³ proposed the partial residual (see also Fowlkes¹⁹⁹).
- [9] See Berk and Booth⁵¹ for other partial-like residuals.
- [10] See [341] for an example comparing a smoothing method with a parametric logistic model fit.
- [11] See Collett [115, Chap. 5] and Pregibon⁵¹² for more information about influence statistics. Pregibon's resistant estimator of β handles overly influential *groups* of observations and allows one to estimate the weight that an observation contributed to the fit after making the fit robust. Observations receiving low weight are partially ignored but are not deleted.
- [12] Buyse⁸⁶ showed that in the case of a single categorical predictor, the ordinary R^2 has a ready interpretation in terms of variance explained for binary responses. Menard⁴⁵⁴ studied various indexes for binary logistic regression. He criticized R_N^2 for being too dependent on the proportion of observations with $Y = 1$. Hu et al.³⁰⁹ further studied the properties of variance-based R^2 measures for binary responses. Tjur⁶¹³ has a nice discussion discrimination graphics

and sum of squares–based R^2 measures for binary logistic regression, as well as a good discussion of “separation” and infinite regression coefficients. Sums of squares are approximated various ways.

- [13] Very little work has been done on developing adjusted R^2 measures in logistic regression and other non-linear model setups. Liao and McGee⁴⁰⁶ developed one adjusted R^2 measure for binary logistic regression, but it uses simulation to adjust for the bias of overfitting. One might as well use the bootstrap to adjust any of the indexes discussed in this section.
- [14] [123, 633] have more pertinent discussion of probability accuracy scores.
- [15] Copas¹²¹ demonstrated how ROC areas can be misleading when applied to different responses having greatly different prevalences. He proposed another approach, the logit rank plot. Newsom⁴⁷³ is an excellent reference on D_{xy} . Newsom⁴⁷⁴ developed several generalizations to D_{xy} including a stratified version, and discussed the jackknife variance estimator for them. ROC areas are not very useful for comparing two models^{118, 493} (but see⁴⁹⁰).
- [16] Gneiting and Raftery²¹⁹ have an excellent review of proper scoring rules. Hand²⁵³ contains much information about assessing classification accuracy. Mittlböck and Schemper⁴⁶¹ have an excellent review of indexes of explained variation for binary logistic models. See also Korn and Simon³⁶⁶ and Zheng and Agresti.⁶⁸⁴
- [17] Pryor et al.⁵¹⁵ presented nomograms for a 10-variable logistic model. One of the variables was sex, which interacted with some of the other variables. Evaluation of predicted probabilities was simplified by the construction of separate nomograms for females and males. Seven terms for discrete predictors were collapsed into one weighted point score axis in the nomograms, and age by risk factor interactions were captured by having four age scales.
- [18] Moons et al.⁴⁶² presents a case study in penalized binary logistic regression modeling.
- [19] The `rcspline.plot` function in the `Hmisc` R package does not allow for interactions as does `lrm`, but it can provide detailed output for checking spline fits. This function plots the estimated spline regression and confidence limits, placing summary statistics on the graph. If there are no adjustment variables, `rcspline.plot` can also plot two alternative estimates of the regression function: proportions or logit proportions on grouped data, and a nonparametric estimate. The nonparametric regression estimate is based on smoothing the binary responses and taking the logit transformation of the smoothed estimates, if desired. The smoothing uses the “super smoother” of Friedman²⁰⁷ implemented in the R function `supsmu`.

10.13 Problems

1. Consider the age–sex–response example in Section 10.1.3. This dataset is available from the text’s web site in the Datasets area.
 - a. Duplicate the analyses done in Section 10.1.3.
 - b. For the model containing both age and sex, test H_0 : logit response is linear in age versus H_a : logit response is quadratic in age. Use the best test statistic.
 - c. Using a Wald test, test H_0 : no age \times sex interaction. Interpret all parameters in the model.

- d. Plot the estimated logit response as a function of age and sex, with and without fitting an interaction term.
 - e. Perform a likelihood ratio test of H_0 : the model containing only age and sex is adequate versus H_a : model is inadequate. Here, “inadequate” may mean nonlinearity (quadratic) in age or presence of an interaction.
 - f. Assuming no interaction is present, test H_0 : model is linear in age versus H_a : model is nonlinear in age. Allow “nonlinear” to be more general than quadratic. (Hint: use a restricted cubic spline function with knots at age=39, 45, 55, 64 years.)
 - g. Plot age against the estimated spline transformation of age (the transformation that would make age fit linearly). You can set the sex and intercept terms to anything you choose. Also plot $\text{Prob}\{\text{response} = 1 \mid \text{age, sex}\}$ from this fitted restricted cubic spline logistic model.
2. Consider a binary logistic regression model using the following predictors: age (years), sex, race (white, African-American, Hispanic, Oriental, other), blood pressure (mmHg). The fitted model is given by

$$\begin{aligned} \text{logit Prob}[Y = 1|X] = X\hat{\beta} = & -1.36 + .03(\text{race} = \text{African-American}) \\ & - .04(\text{race} = \text{hispanic}) + .05(\text{race} = \text{oriental}) - .06(\text{race} = \text{other}) \\ & + .07|\text{blood pressure} - 110| + .3(\text{sex} = \text{male}) - .1\text{age} + .002\text{age}^2 + \\ & (\text{sex} = \text{male})[.05\text{age} - .003\text{age}^2]. \end{aligned}$$

- a. Compute the predicted logit (log odds) that $Y = 1$ for a 50-year-old female Hispanic with a blood pressure of 90 mmHg. Also compute the odds that $Y = 1$ ($\text{Prob}[Y = 1]/\text{Prob}[Y = 0]$) and the estimated probability that $Y = 1$.
- b. Estimate odds ratios for each nonwhite race compared with the reference group (white), holding all other predictors constant. Why can you estimate the relative effect of race for all types of subjects without specifying their characteristics?
- c. Compute the odds ratio for a blood pressure of 120 mmHg compared with a blood pressure of 105, holding age first to 30 years and then to 40 years.
- d. Compute the odds ratio for a blood pressure of 120 mmHg compared with a blood pressure of 105, all other variables held to unspecified constants. Why is this relative effect meaningful without knowing the subject’s age, race, or sex?
- e. Compute the estimated risk difference in changing blood pressure from 105 mmHg to 120 mmHg, first for age = 30 then for age = 40, for a white female. Why does the risk difference depend on age?
- f. Compute the relative odds for males compared with females, for age = 50 and other variables held constant.
- g. Same as the previous question but for females : males instead of males : females.
- h. Compute the odds ratio resulting from increasing age from 50 to 55 for males, and then for females, other variables held constant. What is wrong with the following question: What is the relative effect of changing age by one year?

Chapter 11

Case Study in Binary Logistic Regression, Model Selection and Approximation: Predicting Cause of Death

11.1 Overview

This chapter contains a case study on developing, describing, and validating a binary logistic regression model. In addition, the following methods are exemplified:

1. Data reduction using incomplete linear and nonlinear principal components
2. Use of AIC to choose from five modeling variations, deciding which is best for the number of parameters
3. Model simplification using stepwise variable selection and approximation of the full model
4. The relationship between the degree of approximation and the degree of predictive discrimination loss
5. Bootstrap validation that includes penalization for model uncertainty (variable selection) and that demonstrates a loss of predictive discrimination over the full model even when compensating for overfitting the full model.

The data reduction and pre-transformation methods used here were discussed in more detail in Chapter 8. Single imputation will be used because of the limited quantity of missing data.

11.2 Background

Consider the randomized trial of estrogen for treatment of prostate cancer⁸⁷ described in Chapter 8. In this trial, larger doses of estrogen reduced the effect of prostate cancer but at the cost of increased risk of cardiovascular death.

Kay³⁴⁰ did a formal analysis of the competing risks for cancer, cardiovascular, and other deaths. It can also be quite informative to study how treatment and baseline variables relate to the cause of death for those patients who died.³⁷⁶ We subset the original dataset of those patients dying from prostate cancer ($n = 130$), heart or vascular disease ($n = 96$), or cerebrovascular disease ($n = 31$). Our goal is to predict cardiovascular–cerebrovascular death (cvd , $n = 127$) given the patient died from either cvd or prostate cancer. Of interest is whether the time to death has an effect on the cause of death, and whether the importance of certain variables depends on the time of death.

11.3 Data Transformations and Single Imputation

In R, first obtain the desired subset of the data and do some preliminary calculations such as combining an infrequent category with the next category, and dichotomizing ekg for use in ordinary principal components (PCs).

```
require(rms)
```

```
getHdata(prostate)
prostate <-
  within(prostate, {
    levels(ekg)[levels(ekg) %in%
      c('old MI', 'recent MI')] <- 'MI'
    ekg.norm <- 1*(ekg %in% c('normal', 'benign'))
    levels(ekg) <- abbreviate(levels(ekg))
    pfn <- as.numeric(pf)
    levels(pf) <- levels(pf)[c(1,2,3,3)]
    cvd <- status %in% c("dead - heart or vascular",
      "dead - cerebrovascular")
    rxn = as.numeric(rx) })
# Use transcan to compute optimal pre-transformations
ptrans <- # See Figure 8.3
  transcan(~ sz + sg + ap + sbp + dbp +
    age + wt + hg + ekg + pf + bm + hx + dtime + rx,
    imputed=TRUE, transformed=TRUE,
    data=prostate, pl=FALSE, pr=FALSE)
# Use transcan single imputations
imp <- impute(ptrans, data=prostate, list.out=TRUE)
```

Imputed missing values with the following frequencies
and stored them in variables with their original names:

```
sz  sg  age  wt  ekg
5   11  1    2    8
```

```
NAvars <- all.vars(~ sz + sg + age + wt + ekg)
for(x in NAvars) prostate[[x]] <- imp[[x]]
subset <- prostate$status %in% c("dead - heart or vascular",
```

```

      "dead - cerebrovascular", "dead - prostatic ca")
trans <- ptrans$transformed[subset,]
psub  <- prostate[subset,]

```

11.4 Regression on Original Variables, Principal Components and Pretransformations

We first examine the performance of data reduction in predicting the cause of death, similar to what we did for survival time in Section 8.6. The first analyses assess how well PCs (on raw and transformed variables) predict the cause of death.

There are 127 *cvd*s. We use the 15:1 rule of thumb discussed on P. 72 to justify using the first 8 PCs. *ap* is log-transformed because of its extreme distribution.

```

# Function to compute the first k PCs
ipc <- function(x, k=1, ...)
  princomp(x, ..., cor=TRUE)$scores[,1:k]
# Compute the first 8 PCs on raw variables then on
# transformed ones
pc8 <- ipc(~ sz + sg + log(ap) + sbp + dbp + age +
           wt + hg + ekg.norm + pfn + bm + hx + rxn + dtime,
           data=psub, k=8)
f8 <- lrm(cvd ~ pc8, data=psub)
pc8t <- ipc(trans, k=8)
f8t <- lrm(cvd ~ pc8t, data=psub)
# Fit binary logistic model on original variables
f <- lrm(cvd ~ sz + sg + log(ap) + sbp + dbp + age +
         wt + hg + ekg + pf + bm + hx + rx + dtime, data=psub)
# Expand continuous variables using splines
g <- lrm(cvd ~ rcs(sz,4) + rcs(sg,4) + rcs(log(ap),4) +
         rcs(sbp,4) + rcs(dbp,4) + rcs(age,4) + rcs(wt,4) +
         rcs(hg,4) + ekg + pf + bm + hx + rx + rcs(dtime,4),
         data=psub)
# Fit binary logistic model on individual transformed var.
h <- lrm(cvd ~ trans, data=psub)

```

The five approaches to modeling the outcome are compared using AIC (where smaller is better).

```
c(f8=AIC(f8), f8t=AIC(f8t), f=AIC(f), g=AIC(g), h=AIC(h))
```

f8	f8t	f	g	h
257.6573	254.5172	255.8545	263.8413	254.5317

Based on AIC, the more traditional model fitted to the raw data and assuming linearity for all the continuous predictors has only a slight chance of producing worse cross-validated predictive accuracy than other methods.

The chances are also good that effect estimates from this simple model will have competitive mean squared errors.

11.5 Description of Fitted Model

Here we describe the simple all-linear full model. Summary statistics and a Wald-ANOVA table are below, followed by partial effects plots with pointwise confidence bands, and odds ratios over default ranges of predictors.

```
print(f, latex=TRUE)
```

Logistic Regression Model

```
lrm(formula = cvd ~ sz + sg + log(ap) + sbp + dbp + age + wt +
    hg + ekg + pf + bm + hx + rx + dtime, data = psub)
```

		Model Likelihood Ratio Test	Discrimination Indexes	Rank Discrim. Indexes
Obs	257	LR χ^2 144.39	R^2 0.573	C 0.893
FALSE	130	d.f. 21	g 2.688	D_{xy} 0.786
TRUE	127	$\Pr(> \chi^2) < 0.0001$	g_r 14.701	γ 0.787
max $ \frac{\partial \log L}{\partial \beta} $		6×10^{-11}	g_p 0.394	τ_a 0.395
			Brier 0.133	

	Coef	S.E.	Wald Z	$\Pr(> Z)$
Intercept	-4.5130	3.2210	-1.40	0.1612
sz	-0.0640	0.0168	-3.80	0.0001
sg	-0.2967	0.1149	-2.58	0.0098
ap	-0.3927	0.1411	-2.78	0.0054
sbp	-0.0572	0.0890	-0.64	0.5201
dbp	0.3917	0.1629	2.40	0.0162
age	0.0926	0.0286	3.23	0.0012
wt	-0.0177	0.0140	-1.26	0.2069
hg	0.0860	0.0925	0.93	0.3524
ekg=bngn	1.0781	0.8793	1.23	0.2202
ekg=rd&ec	-0.1929	0.6318	-0.31	0.7601
ekg=hbocd	-1.3679	0.8279	-1.65	0.0985
ekg=hrts	0.4365	0.4582	0.95	0.3407
ekg=MI	0.3039	0.5618	0.54	0.5886
pf=in bed < 50% daytime	0.9604	0.6956	1.38	0.1673
pf=in bed > 50% daytime	-2.3232	1.2464	-1.86	0.0623
bm	0.1456	0.5067	0.29	0.7738
hx	1.0913	0.3782	2.89	0.0039

	Coef	S.E.	Wald Z	Pr(> Z)
rx=0.2 mg estrogen	-0.3022	0.4908	-0.62	0.5381
rx=1.0 mg estrogen	0.7526	0.5272	1.43	0.1534
rx=5.0 mg estrogen	0.6868	0.5043	1.36	0.1733
dtime	-0.0136	0.0107	-1.27	0.2040

```
an ← anova(f)
latex(an, file='', table.env=FALSE)
```

	χ^2	d.f.	P
sz	14.42	1	0.0001
sg	6.67	1	0.0098
ap	7.74	1	0.0054
sbp	0.41	1	0.5201
dbp	5.78	1	0.0162
age	10.45	1	0.0012
wt	1.59	1	0.2069
hg	0.86	1	0.3524
ekg	6.76	5	0.2391
pf	5.52	2	0.0632
bm	0.08	1	0.7738
hx	8.33	1	0.0039
rx	5.72	3	0.1260
dtime	1.61	1	0.2040
TOTAL	66.87	21	< 0.0001

```
plot(an) # Figure 11.1
s ← f$stats
gamma.hat ← (s['Model L.R. '] - s['d.f. '])/s['Model L.R. ']
```

```
dd ← datadist(psub); options(datadist='dd')
ggplot(Predict(f), sepdiscrete='vertical', vnames='names',
        rdata=psub,
        histSpike.opts=list(frac=function(f) .1*f/max(f) ))
# Figure 11.2
```

```
plot(summary(f), log=TRUE) # Figure 11.3
```

The van Houwelingen–Le Cessie heuristic shrinkage estimate (Equation 4.3) is $\hat{\gamma} = 0.85$, indicating that this model will validate on new data about 15% worse than on this dataset.

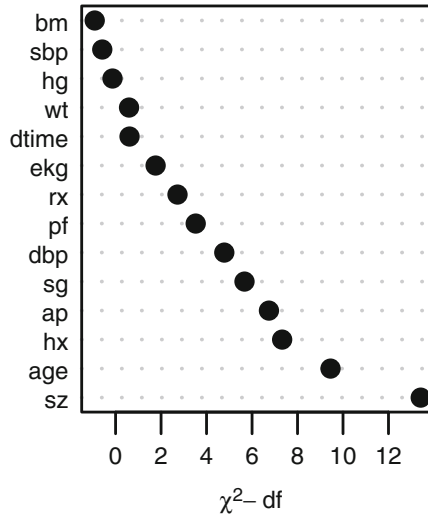


Fig. 11.1 Ranking of apparent importance of predictors of cause of death

11.6 Backwards Step-Down

Now use fast backward step-down (with total residual AIC as the stopping rule) to identify the variables that explain the bulk of the cause of death. Later validation will take this screening of variables into account. The greatly reduced model results in a simple nomogram.

```
fastbw(f)
```

Deleted	Chi-Sq	d.f.	P	Residual	d.f.	P	AIC
ekg	6.76	5	0.2391	6.76	5	0.2391	-3.24
bm	0.09	1	0.7639	6.85	6	0.3349	-5.15
hg	0.38	1	0.5378	7.23	7	0.4053	-6.77
sbp	0.48	1	0.4881	7.71	8	0.4622	-8.29
wt	1.11	1	0.2932	8.82	9	0.4544	-9.18
dtime	1.47	1	0.2253	10.29	10	0.4158	-9.71
rx	5.65	3	0.1302	15.93	13	0.2528	-10.07
pf	4.78	2	0.0915	20.71	15	0.1462	-9.29
sg	4.28	1	0.0385	25.00	16	0.0698	-7.00
dbp	5.84	1	0.0157	30.83	17	0.0209	-3.17

Approximate Estimates after Deleting Factors

	Coef	S.E.	Wald Z	P
Intercept	-3.74986	1.82887	-2.050	0.0403286
sz	-0.04862	0.01532	-3.174	0.0015013
ap	-0.40694	0.11117	-3.660	0.0002518
age	0.06000	0.02562	2.342	0.0191701
hx	0.86969	0.34339	2.533	0.0113198

Factors in Final Model

```
[1] sz ap age hx
```

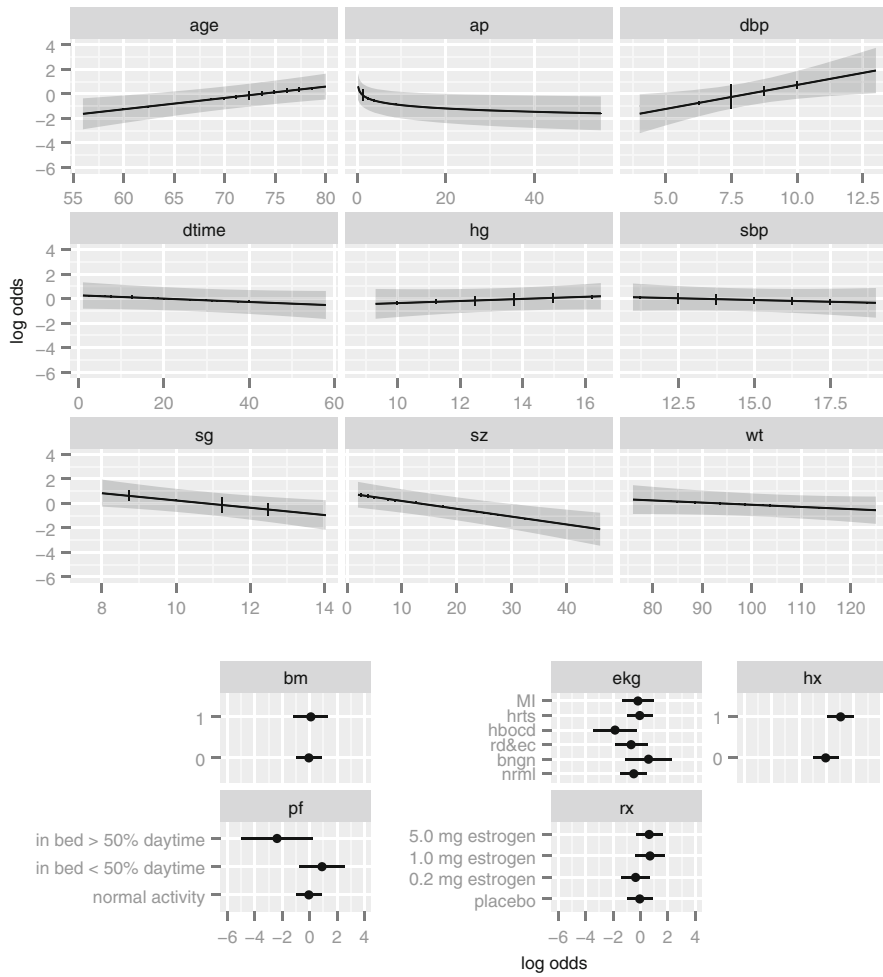


Fig. 11.2 Partial effects (log odds scale) in full model for cause of death, along with vertical line segments showing the raw data distribution of predictors

```
fred ← lrm(cvd ~ sz + log(ap) + age + hx, data=psub)
latex(fred, file='')
```

$$\text{Prob}\{\text{cvd}\} = \frac{1}{1 + \exp(-X\hat{\beta})}, \text{ where}$$

$$X\hat{\beta} = -5.009276 - 0.05510121 \text{ sz} - 0.509185 \log(\text{ap}) + 0.0788052 \text{ age} + 1.070601 \text{ hx}$$

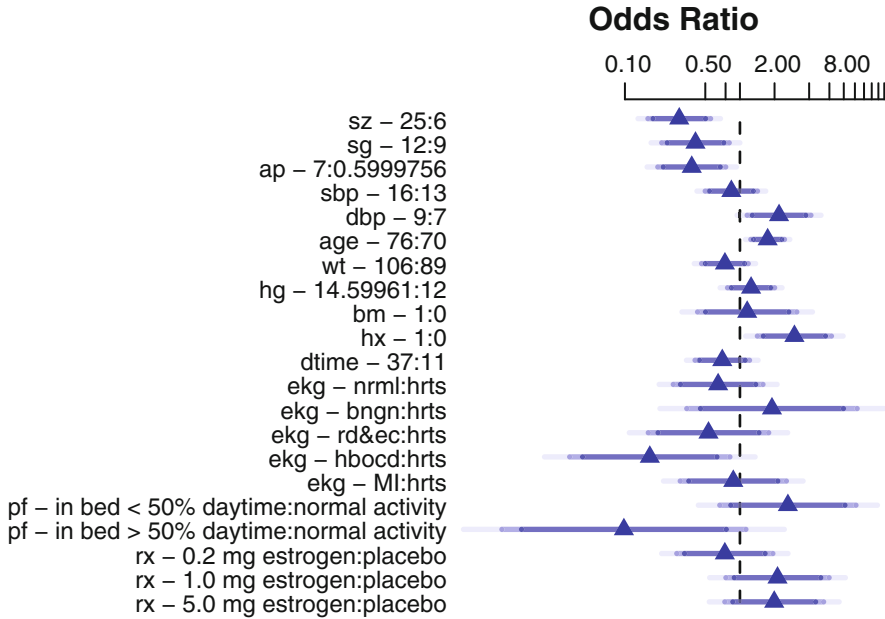


Fig. 11.3 Interquartile-range odds ratios for continuous predictors and simple odds ratios for categorical predictors. Numbers at left are upper quartile : lower quartile or current group : reference group. The bars represent 0.9, 0.95, 0.99 confidence limits. The intervals are drawn on the log odds ratio scale and labeled on the odds ratio scale. Ranges are on the original scale.

```
nom ← nomogram(fred, ap=c(.1, .5, 1, 5, 10, 50),
               fun=plgis, funlabel="Probability",
               fun.at=c(.01, .05, .1, .25, .5, .75, .9, .95, .99))
plot(nom, xfrac=.45) # Figure 11.4
```

It is readily seen from this model that patients with a history of heart disease, and patients with less extensive prostate cancer are those more likely to die from *cvd* rather than from cancer. But beware that it is easy to over-interpret findings when using unpenalized estimation, and confidence intervals are too narrow. Let us use the bootstrap to study the uncertainty in the selection of variables and to penalize for this uncertainty when estimating predictive performance of the model. The variables selected in the first 20 bootstrap resamples are shown, making it obvious that the set of “significant” variables, i.e., the final model, is somewhat arbitrary.

```
f ← update(f, x=TRUE, y=TRUE)
v ← validate(f, B=200, bw=TRUE)
```

```
latex(v, B=20, digits=3)
```

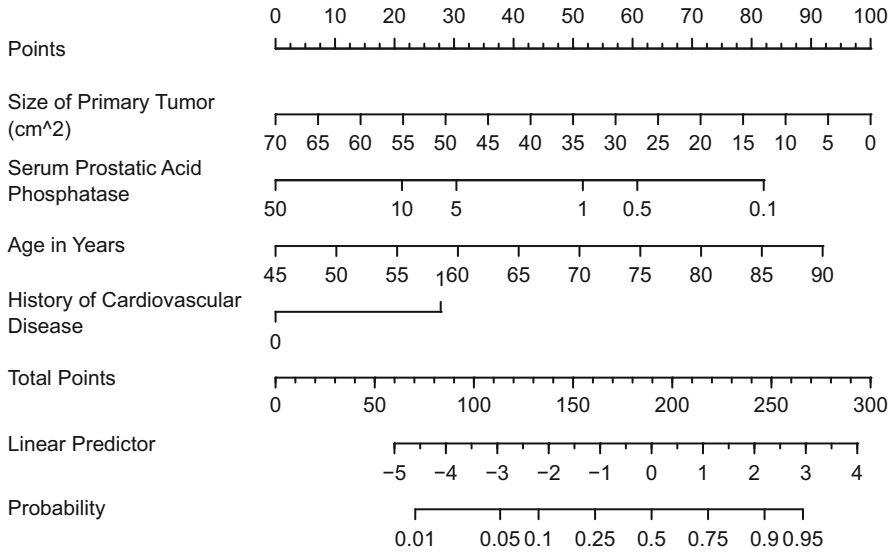


Fig. 11.4 Nomogram calculating $X\hat{\beta}$ and \hat{P} for *cvd* as the cause of death, using the step-down model. For each predictor, read the points assigned on the 0–100 scale and add these points. Read the result on the Total Points scale and then read the corresponding predictions below it.

Index	Original Sample	Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.682	0.713	0.643	0.071	0.611	200
R^2	0.439	0.481	0.393	0.088	0.351	200
Intercept	0.000	0.000	-0.006	0.006	-0.006	200
Slope	1.000	1.000	0.811	0.189	0.811	200
E_{\max}	0.000	0.000	0.048	0.048	0.048	200
D	0.395	0.449	0.346	0.102	0.293	200
U	-0.008	-0.008	0.018	-0.026	0.018	200
Q	0.403	0.456	0.329	0.128	0.275	200
B	0.162	0.151	0.174	-0.022	0.184	200
g	1.932	2.213	1.756	0.457	1.475	200
g_p	0.341	0.355	0.320	0.035	0.306	200

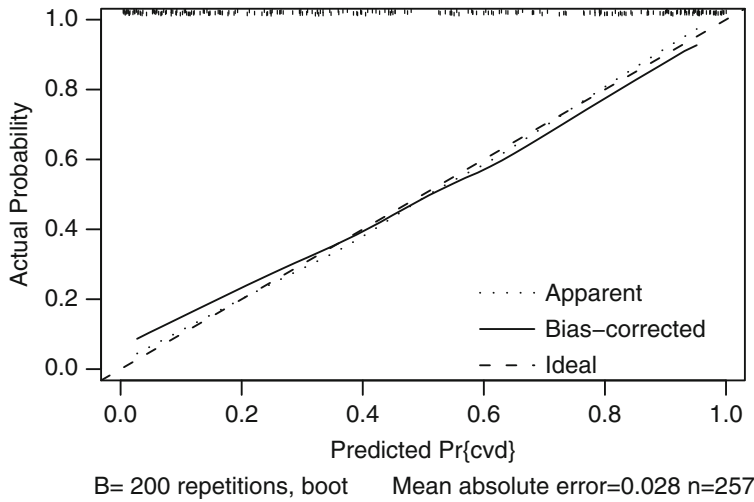


Fig. 11.5 Bootstrap overfitting-corrected calibration curve estimate for the backwards step-down cause of death logistic model, along with a rug plot showing the distribution of predicted risks. The smooth nonparametric calibration estimator (*loess*) is used.

Index	Original Training Sample	Test Sample	Optimism	Corrected Index	<i>n</i>
D_{xy}	0.786	0.833	0.738	0.095	0.691 200
R^2	0.573	0.641	0.501	0.140	0.433 200
Intercept	0.000	0.000	-0.013	0.013	-0.013 200
Slope	1.000	1.000	0.690	0.310	0.690 200
E_{\max}	0.000	0.000	0.085	0.085	0.085 200
D	0.558	0.653	0.468	0.185	0.373 200
U	-0.008	-0.008	0.051	-0.058	0.051 200
Q	0.566	0.661	0.417	0.244	0.322 200
B	0.133	0.115	0.150	-0.035	0.168 200
g	2.688	3.464	2.355	1.108	1.579 200
g_p	0.394	0.416	0.366	0.050	0.344 200

Compared to the validation of the full model, the step-down model has less optimism, but it started with a smaller D_{xy} due to loss of information from removing moderately important variables. The improvement in optimism was not enough to offset the effect of eliminating variables. If shrinkage were used with the full model, it would have better calibration and discrimination than the reduced model, since shrinkage does not diminish D_{xy} . Thus stepwise variable selection failed at delivering excellent predictive discrimination.

Finally, compare previous results with a bootstrap validation of a step-down model using a better significance level for a variable to stay in the

model ($\alpha = 0.5$,⁵⁸⁹) and using individual approximate Wald tests rather than tests combining all deleted variables.

```
v5 ← validate(f, bw=TRUE, sls=0.5, type='individual', B=200)
```

```

Backwards Step-down - Original Model

Deleted Chi-Sq d.f. P      Residual d.f. P      AIC
ekg      6.76  5    0.2391  6.76  5    0.2391  -3.24
bm       0.09  1    0.7639  6.85  6    0.3349  -5.15
hg       0.38  1    0.5378  7.23  7    0.4053  -6.77
sbp     0.48  1    0.4881  7.71  8    0.4622  -8.29
wt      1.11  1    0.2932  8.82  9    0.4544  -9.18
dtime   1.47  1    0.2253 10.29 10    0.4158  -9.71
rx      5.65  3    0.1302 15.93 13    0.2528 -10.07

Approximate Estimates after Deleting Factors

              Coef      S.E. Wald Z      P
Intercept   -4.86308  2.67292  -1.819  0.068852
sz          -0.05063  0.01581  -3.202  0.001366
sg          -0.28038  0.11014  -2.546  0.010903
ap          -0.24838  0.12369  -2.008  0.044629
dbp         0.28288  0.13036   2.170  0.030008
age         0.08502  0.02690   3.161  0.001572
pf=in bed < 50% daytime  0.81151  0.66376   1.223  0.221485
pf=in bed > 50% daytime -2.19885  1.21212  -1.814  0.069670
hx          0.87834  0.35203   2.495  0.012592

Factors in Final Model

[1] sz  sg  ap  dbp age pf  hx

```

```
latex(v5, digits=3, B=0)
```

Index	Original Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.739	0.801	0.716	0.085	0.654 200
R^2	0.517	0.598	0.481	0.117	0.400 200
Intercept	0.000	0.000	-0.008	0.008	-0.008 200
Slope	1.000	1.000	0.745	0.255	0.745 200
E_{\max}	0.000	0.000	0.067	0.067	0.067 200
D	0.486	0.593	0.444	0.149	0.337 200
U	-0.008	-0.008	0.033	-0.040	0.033 200
Q	0.494	0.601	0.411	0.190	0.304 200
B	0.147	0.125	0.156	-0.030	0.177 200
g	2.351	2.958	2.175	0.784	1.567 200
g_p	0.372	0.401	0.358	0.043	0.330 200

The performance statistics are midway between the full model and the smaller stepwise model.

11.7 Model Approximation

Frequently a better approach than stepwise variable selection is to approximate the full model, using its estimates of precision, as discussed in Section 5.5. Stepwise variable selection as well as regression trees are useful for making the approximations, and the sacrifice in predictive accuracy is always apparent.

We begin by computing the “gold standard” linear predictor from the full model fit ($R^2 = 1.0$), then running backwards step-down OLS regression to approximate it.

```
lp ← predict(f) # Compute linear predictor from full model
# Insert sigma=1 as otherwise sigma=0 will cause problems
a ← ols(lp ~ sz + sg + log(ap) + sbp + dbp + age + wt +
        hg + ekg + pf + bm + hx + rx + dtime, sigma=1,
        data=psub)
# Specify silly stopping criterion to remove all variables
s ← fastbw(a, aics=10000)
betas ← s$Coefficients # matrix, rows=iterations
X ← cbind(1, f$x) # design matrix
# Compute the series of approximations to lp
ap ← X %*% t(betas)
# For each approx. compute approximation R^2 and ratio of
# likelihood ratio chi-square for approximate model to that
# of original model
m ← ncol(ap) - 1 # all but intercept-only model
r2 ← frac ← numeric(m)
fullchisq ← f$stats['Model L.R. ']
for(i in 1:m) {
  lpa ← ap[,i]
  r2[i] ← cor(lpa, lp)^2
  fapprox ← lrm(cvd ~ lpa, data=psub)
  frac[i] ← fapprox$stats['Model L.R. '] / fullchisq
} # Figure 11.6:
plot(r2, frac, type='b',
     xlab=expression(paste('Approximation ', R^2)),
     ylab=expression(paste('Fraction of ',
                           chi^2, ' Preserved')))
abline(h=.95, col=gray(.83)); abline(v=.95, col=gray(.83))
abline(a=0, b=1, col=gray(.83))
```

After 6 deletions, slightly more than 0.05 of both the LR χ^2 and the approximation R^2 are lost (see Figure 11.6). Therefore we take as our approximate model the one that removed 6 predictors. The equation for this model is below, and its nomogram is in Figure 11.7.

```
fapprox ← ols(lp ~ sz + sg + log(ap) + age + ekg + pf + hx +
              rx, data=psub)
fapprox$stats['R2'] # as a check
```

R2 0.9453396

```
latex(fapprox, file='')
```

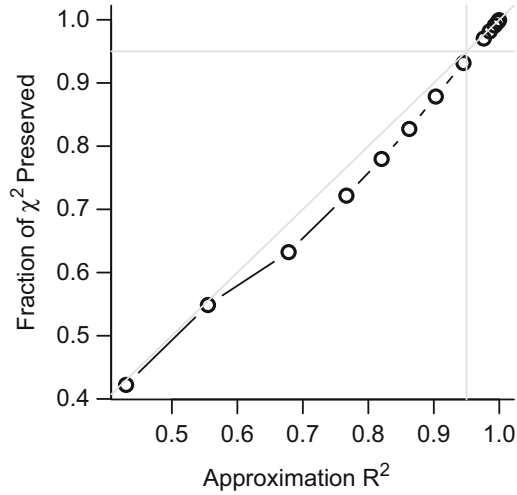


Fig. 11.6 Fraction of explainable variation (full model LR χ^2) in *cvd* that was explained by approximate models, along with approximation accuracy (*x*-axis)

$$E(lp) = X\beta, \text{ where}$$

$$\begin{aligned}
 X\hat{\beta} = & -2.868303 - 0.06233241 \text{ sz} - 0.3157901 \text{ sg} - 0.3834479 \log(\text{ap}) + 0.09089393 \text{ age} \\
 & + 1.396922[\text{bngn}] + 0.06275034[\text{rd\&ec}] - 1.24892[\text{hbocd}] + 0.6511938[\text{hrts}] \\
 & + 0.3236771[\text{MI}] \\
 & + 1.116028[\text{in bed} < 50\% \text{ daytime}] - 2.436734[\text{in bed} > 50\% \text{ daytime}] \\
 & + 1.05316 \text{ hx} \\
 & - 0.3888534[0.2 \text{ mg estrogen}] + 0.6920495[1.0 \text{ mg estrogen}] \\
 & + 0.7834498[5.0 \text{ mg estrogen}]
 \end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise.

```

nom ← nomogram(fapprox, ap=c(.1, .5, 1, 5, 10, 20, 30, 40),
               fun=plogis, funlabel="Probability",
               lp.at=(-5):4,
               fun.lp.at=qlogis(c(.01, .05, .25, .5, .75, .95, .99)))
plot(nom, xfrac=.45) # Figure 11.7

```

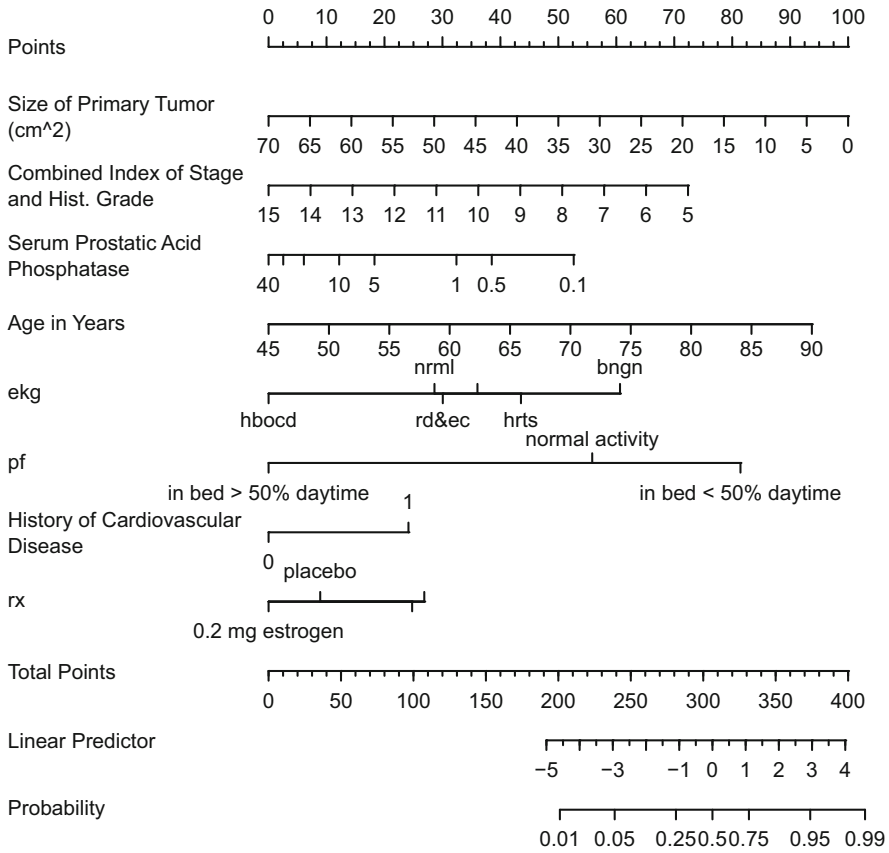


Fig. 11.7 Nomogram for predicting the probability of cvd based on the approximate model

Chapter 12

Logistic Model Case Study 2: Survival of Titanic Passengers

This case study demonstrates the development of a binary logistic regression model to describe patterns of survival in passengers on the *Titanic*, based on passenger age, sex, ticket class, and the number of family members accompanying each passenger. Nonparametric regression is also used. Since many of the passengers had missing ages, multiple imputation is used so that the complete information on the other variables can be efficiently utilized. Titanic passenger data were gathered by many researchers. Primary references are the *Encyclopedia Titanica* at www.encyclopedia-titanica.org and Eaton and Haas.¹⁶⁹ Titanic survival patterns have been analyzed previously^{151, 296, 571} but without incorporation of individual passenger ages. Thomas Cason while a University of Virginia student compiled and interpreted the data from the World Wide Web. One thousand three hundred nine of the passengers are represented in the dataset, which is available from this text's Web site under the name `titanic3`. An early analysis of Titanic data may be found in Bron⁷⁵.

12.1 Descriptive Statistics

First we obtain basic descriptive statistics on key variables.


```
require(rms)

getHdata(titanic3)      # get dataset from web site
# List of names of variables to analyze
v <- c('pclass', 'survived', 'age', 'sex', 'sibsp', 'parch')
t3 <- titanic3[, v]
units(t3$age) <- 'years'
latex(describe(t3), file='')
```

t3
6 Variables 1309 Observations

pclass
n missing unique
1309 0 3
1st (323, 25%), 2nd (277, 21%), 3rd (709, 54%)

survived : Survived
n missing unique Info Sum Mean
1309 0 2 0.71 500 0.382

age : Age [years] 
n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
1046 263 98 1 29.88 5 14 21 28 39 50 57
lowest : 0.1667 0.3333 0.4167 0.6667 0.7500
highest: 70.5000 71.0000 74.0000 76.0000 80.0000

sex
n missing unique
1309 0 2
female (466, 36%), male (843, 64%)

sibsp : Number of Siblings/Spouses Aboard
n missing unique Info Mean
1309 0 7 0.67 0.4989
Frequency 0 1 2 3 4 5 8
891 319 42 20 22 6 9
% 68 24 3 2 2 0 1

parch : Number of Parents/Children Aboard
n missing unique Info Mean
1309 0 8 0.55 0.385
Frequency 0 1 2 3 4 5 6 9
1002 170 113 8 6 6 2 2
% 77 13 9 1 0 0 0

Next, we obtain access to the needed variables and observations, and save data distribution characteristics for plotting and for computing predictor effects. There are not many passengers having more than 3 siblings or spouses or more than 3 children, so we truncate two variables at 3 for the purpose of estimating stratified survival probabilities.

```
dd <- datadist(t3)
# describe distributions of variables to rms
options(datadist='dd')
s <- summary(survived ~ age + sex + pclass +
             cut2(sibsp,0:3) + cut2(parch,0:3), data=t3)
plot(s, main='', subtitles=FALSE) # Figure 12.1
```

Note the large number of missing ages. Also note the strong effects of sex and passenger class on the probability of surviving. The age effect does not appear to be very strong, because as we show later, much of the effect is restricted to

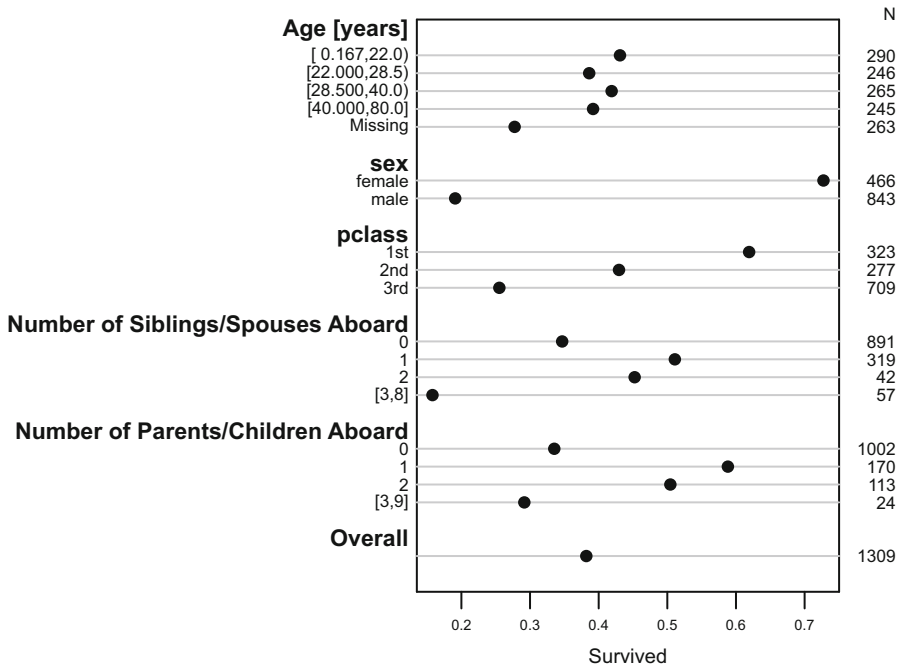


Fig. 12.1 Univariable summaries of Titanic survival

age < 21 years for one of the sexes. The effects of the last two variables are unclear as the estimated proportions are not monotonic in the values of these descriptors. Although some of the cell sizes are small, we can show four-way empirical relationships with the fraction of surviving passengers by creating four cells for `sibsp` × `parch` combinations and by creating two age groups. We suppress proportions based on fewer than 25 passengers in a cell. Results are shown in Figure 12.2.

```
tn <- transform(t3,
  agec = ifelse(age < 21, 'child', 'adult'),
  sibsp = ifelse(sibsp == 0, 'no sib/sp', 'sib/sp'),
  parch = ifelse(parch == 0, 'no par/child', 'par/child'))

g <- function(y) if(length(y) < 25) NA else mean(y)
s <- with(tn, summarize(survived,
  llist(agec, sex, pclass, sibsp, parch), g))
# llist, summarize in Hmisc package
# Figure 12.2:
ggplot(subset(s, agec != 'NA'),
  aes(x=survived, y=pclass, shape=sex)) +
  geom_point() + facet_grid(agec ~ sibsp * parch) +
  xlab('Proportion Surviving') + ylab('Passenger Class') +
  scale_x_continuous(breaks=c(0, .5, 1))
```

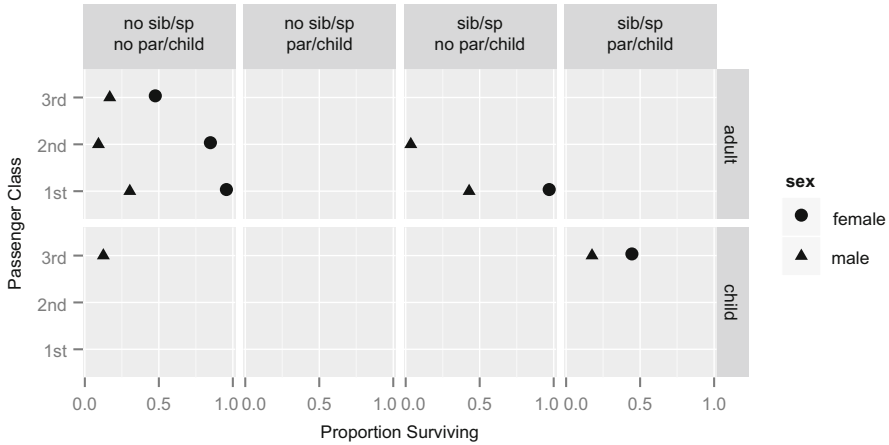


Fig. 12.2 Multi-way summary of Titanic survival

Note that none of the effects of `sibsp` or `parch` for common passenger groups appear strong on an absolute risk scale.

12.2 Exploring Trends with Nonparametric Regression

As described in Section 2.4.7, the `loess` smoother has excellent performance when the response is binary, as long as outlier detection is turned off. Here we use a `ggplot2` add-on function `histSpikeg` in the `Hmisc` package to obtain and plot the `loess` fit and age distribution. `histSpikeg` uses the “no iteration” option for the R `lowess` function when the response is binary.

```
# Figure 12.3
b ← scale_size_discrete(range=c(.1, .85))
yl ← ylab(NULL)
p1 ← ggplot(t3, aes(x=age, y=survived)) +
  histSpikeg(survived ~ age, lowess=TRUE, data=t3) +
  ylim(0,1) + yl
p2 ← ggplot(t3, aes(x=age, y=survived, color=sex)) +
  histSpikeg(survived ~ age + sex, lowess=TRUE,
    data=t3) + ylim(0,1) + yl
p3 ← ggplot(t3, aes(x=age, y=survived, size=pclass)) +
  histSpikeg(survived ~ age + pclass, lowess=TRUE,
    data=t3) + b + ylim(0,1) + yl
p4 ← ggplot(t3, aes(x=age, y=survived, color=sex,
  size=pclass)) +
  histSpikeg(survived ~ age + sex + pclass,
    lowess=TRUE, data=t3) +
  b + ylim(0,1) + yl
gridExtra::grid.arrange(p1, p2, p3, p4, ncol=2) # combine 4
```

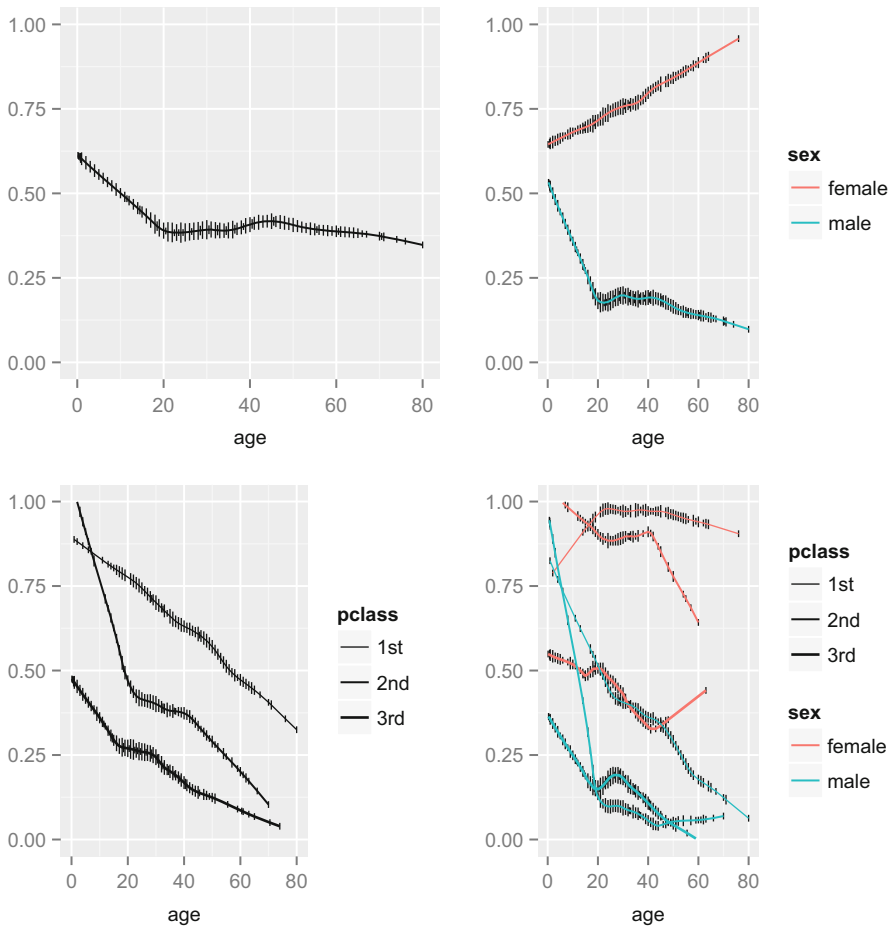


Fig. 12.3 Nonparametric regression (*loess*) estimates of the relationship between age and the probability of surviving the Titanic, with tick marks depicting the age distribution. The top left panel shows unstratified estimates of the probability of survival. Other panels show nonparametric estimates by various stratifications.

Figure 12.3 shows much of the story of passenger survival patterns. “Women and children first” seems to be true except for women in third class. It is interesting that there is no real cutoff for who is considered a child. For men, the younger the greater chance of surviving. The interpretation of the effects of the “number of relatives”-type variables will be more difficult, as their definitions are a function of age. Figure 12.4 shows these relationships.

```
# Figure 12.4
top <- theme(legend.position='top')
p1 <- ggplot(t3, aes(x=age, y=survived, color=cut2(sibsp,
  0:2))) + stat_plsmo() + b + ylim(0,1) + y1 + top +
  scale_color_discrete(name='siblings/spouses')
```

```
p2 ← ggplot(t3, aes(x=age, y=survived, color=cut2(parch,
0:2))) + stat_plsmo() + b + ylim(0,1) + y1 + top +
  scale_color_discrete(name='parents/children')
gridExtra::grid.arrange(p1, p2, ncol=2)
```

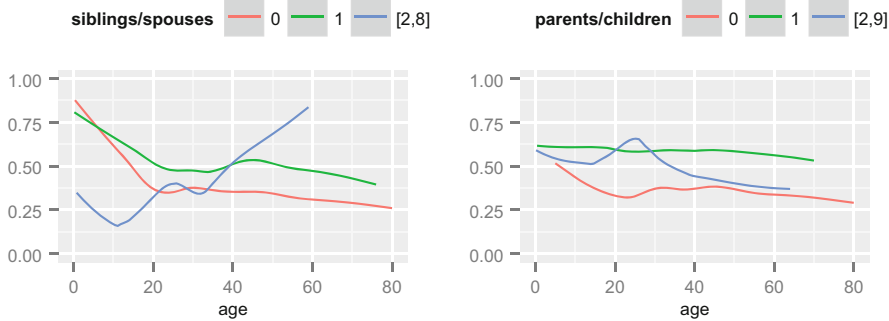


Fig. 12.4 Relationship between age and survival stratified by the number of siblings or spouses on board (left panel) or by the number of parents or children of the passenger on board (right panel).

12.3 Binary Logistic Model With Casewise Deletion of Missing Values

What follows is the standard analysis based on eliminating observations having any missing data. We develop an initial somewhat saturated logistic model, allowing for a flexible nonlinear age effect that can differ in shape for all six $\text{sex} \times \text{class}$ strata. The `sibsp` and `parch` variables do not have sufficiently dispersed distributions to allow for us to model them nonlinearly. Also, there are too few passengers with nonzero values of these two variables in $\text{sex} \times \text{pclass} \times \text{age}$ strata to allow us to model complex interactions involving them. The meaning of these variables does depend on the passenger's age, so we consider only age interactions involving `sibsp` and `parch`.

```
f1 ← lrm(survived ~ sex*pclass*rCs(age,5) +
  rCs(age,5)*(sibsp + parch), data=t3) # Table 12.1
latex(anova(f1), file='', label='titanic-anova3',
  size='small')
```

Three-way interactions are clearly insignificant ($P = 0.4$) in Table 12.1. So is `parch` ($P = 0.6$ for testing the combined main effect + interaction effects for `parch`, i.e., whether `parch` is important for any age). These effects would be deleted in almost all bootstrap resamples had we bootstrapped a variable selection procedure using $\alpha = 0.1$ for retention of terms, so we can safely ignore these terms for future steps. The model not containing those terms

Table 12.1 Wald Statistics for survived

	χ^2	d.f.	<i>P</i>
sex (Factor+Higher Order Factors)	187.15	15	< 0.0001
<i>All Interactions</i>	59.74	14	< 0.0001
pclass (Factor+Higher Order Factors)	100.10	20	< 0.0001
<i>All Interactions</i>	46.51	18	0.0003
age (Factor+Higher Order Factors)	56.20	32	0.0052
<i>All Interactions</i>	34.57	28	0.1826
<i>Nonlinear (Factor+Higher Order Factors)</i>	28.66	24	0.2331
sibsp (Factor+Higher Order Factors)	19.67	5	0.0014
<i>All Interactions</i>	12.13	4	0.0164
parch (Factor+Higher Order Factors)	3.51	5	0.6217
<i>All Interactions</i>	3.51	4	0.4761
sex × pclass (Factor+Higher Order Factors)	42.43	10	< 0.0001
sex × age (Factor+Higher Order Factors)	15.89	12	0.1962
<i>Nonlinear (Factor+Higher Order Factors)</i>	14.47	9	0.1066
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	4.17	3	0.2441
pclass × age (Factor+Higher Order Factors)	13.47	16	0.6385
<i>Nonlinear (Factor+Higher Order Factors)</i>	12.92	12	0.3749
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	6.88	6	0.3324
age × sibsp (Factor+Higher Order Factors)	12.13	4	0.0164
<i>Nonlinear</i>	1.76	3	0.6235
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	1.76	3	0.6235
age × parch (Factor+Higher Order Factors)	3.51	4	0.4761
<i>Nonlinear</i>	1.80	3	0.6147
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	1.80	3	0.6147
sex × pclass × age (Factor+Higher Order Factors)	8.34	8	0.4006
<i>Nonlinear</i>	7.74	6	0.2581
TOTAL NONLINEAR	28.66	24	0.2331
TOTAL INTERACTION	75.61	30	< 0.0001
TOTAL NONLINEAR + INTERACTION	79.49	33	< 0.0001
TOTAL	241.93	39	< 0.0001

is fitted below. The $\wedge 2$ in the model formula means to expand the terms in parentheses to include all main effects and second-order interactions.

```
f <- lrm(survived ~ (sex + pclass + rcs(age,5))^2 +
          rcs(age,5)*sibsp, data=t3)
print(f, latex=TRUE)
```

Logistic Regression Model

```
lrm(formula = survived ~ (sex + pclass + rcs(age, 5))^2
    + rcs(age, 5) * sibsp, data = t3)
```

Frequencies of Missing Values Due to Each Variable

survived	sex	pclass	age	sibsp
0	0	0	263	0

		Model Likelihood Ratio Test		Discrimination Indexes		Rank Discrim. Indexes	
Obs	1046	LR χ^2	553.87	R^2	0.555	C	0.878
0	619	d.f.	26	g	2.427	D_{xy}	0.756
1	427	Pr($> \chi^2$) < 0.0001		g_r	11.325	γ	0.758
max $ \frac{\partial \log L}{\partial \beta} $		6×10^{-6}		g_p	0.365	τ_a	0.366
				Brier	0.130		

	Coef	S.E.	Wald Z	Pr(> Z)
Intercept	3.3075	1.8427	1.79	0.0727
sex=male	-1.1478	1.0878	-1.06	0.2914
pclass=2nd	6.7309	3.9617	1.70	0.0893
pclass=3rd	-1.6437	1.8299	-0.90	0.3691
age	0.0886	0.1346	0.66	0.5102
age'	-0.7410	0.6513	-1.14	0.2552
age''	4.9264	4.0047	1.23	0.2186
age'''	-6.6129	5.4100	-1.22	0.2216
sibsp	-1.0446	0.3441	-3.04	0.0024
sex=male * pclass=2nd	-0.7682	0.7083	-1.08	0.2781
sex=male * pclass=3rd	2.1520	0.6214	3.46	0.0005
sex=male * age	-0.2191	0.0722	-3.04	0.0024
sex=male * age'	1.0842	0.3886	2.79	0.0053
sex=male * age''	-6.5578	2.6511	-2.47	0.0134
sex=male * age'''	8.3716	3.8532	2.17	0.0298
pclass=2nd * age	-0.5446	0.2653	-2.05	0.0401
pclass=3rd * age	-0.1634	0.1308	-1.25	0.2118
pclass=2nd * age'	1.9156	1.0189	1.88	0.0601
pclass=3rd * age'	0.8205	0.6091	1.35	0.1780
pclass=2nd * age''	-8.9545	5.5027	-1.63	0.1037
pclass=3rd * age''	-5.4276	3.6475	-1.49	0.1367
pclass=2nd * age'''	9.3926	6.9559	1.35	0.1769
pclass=3rd * age'''	7.5403	4.8519	1.55	0.1202
age * sibsp	0.0357	0.0340	1.05	0.2933
age' * sibsp	-0.0467	0.2213	-0.21	0.8330
age'' * sibsp	0.5574	1.6680	0.33	0.7382
age''' * sibsp	-1.1937	2.5711	-0.46	0.6425

```

latex(anova(f), file=' ', label='titanic-anova2 ', size='small ')
#12.2
    
```

This is a very powerful model (ROC area = $c = 0.88$); the survival patterns are easy to detect. The Wald ANOVA in Table 12.2 indicates especially strong sex and pclass effects ($\chi^2 = 199$ and 109, respectively). There is a very strong

Table 12.2 Wald Statistics for survived

	χ^2	d.f.	<i>P</i>
sex (Factor+Higher Order Factors)	199.42	7	< 0.0001
<i>All Interactions</i>	56.14	6	< 0.0001
pclass (Factor+Higher Order Factors)	108.73	12	< 0.0001
<i>All Interactions</i>	42.83	10	< 0.0001
age (Factor+Higher Order Factors)	47.04	20	0.0006
<i>All Interactions</i>	24.51	16	0.0789
<i>Nonlinear (Factor+Higher Order Factors)</i>	22.72	15	0.0902
sibsp (Factor+Higher Order Factors)	19.95	5	0.0013
<i>All Interactions</i>	10.99	4	0.0267
sex × pclass (Factor+Higher Order Factors)	35.40	2	< 0.0001
sex × age (Factor+Higher Order Factors)	10.08	4	0.0391
<i>Nonlinear</i>	8.17	3	0.0426
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	8.17	3	0.0426
pclass × age (Factor+Higher Order Factors)	6.86	8	0.5516
<i>Nonlinear</i>	6.11	6	0.4113
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	6.11	6	0.4113
age × sibsp (Factor+Higher Order Factors)	10.99	4	0.0267
<i>Nonlinear</i>	1.81	3	0.6134
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	1.81	3	0.6134
TOTAL NONLINEAR	22.72	15	0.0902
TOTAL INTERACTION	67.58	18	< 0.0001
TOTAL NONLINEAR + INTERACTION	70.68	21	< 0.0001
TOTAL	253.18	26	< 0.0001

sex × pclass interaction and a strong age × sibsp interaction, considering the strength of sibsp overall.

Let us examine the shapes of predictor effects. With so many interactions in the model we need to obtain predicted values at least for all combinations of sex and pclass. For sibsp we consider only two of its possible values.

```
p ← Predict(f, age, sex, pclass, sibsp=0, fun=plogis)
ggplot(p) # Fig. 12.5
```

Note the agreement between the lower right-hand panel of Figure 12.3 with Figure 12.5. This results from our use of similar flexibility in the parametric and nonparametric approaches (and similar effective degrees of freedom). The estimated effect of sibsp as a function of age is shown in Figure 12.6.

```
ggplot(Predict(f, sibsp, age=c(10,15,20,50), conf.int=FALSE))
## Figure 12.6
```

Note that children having many siblings apparently had lower survival. Married adults had slightly higher survival than unmarried ones.

There will never be another Titanic, so we do not need to validate the model for prospective use. But we use the bootstrap to validate the model anyway, in an effort to detect whether it is overfitting the data. We do not penalize the calculations that follow for having examined the effect of parch or

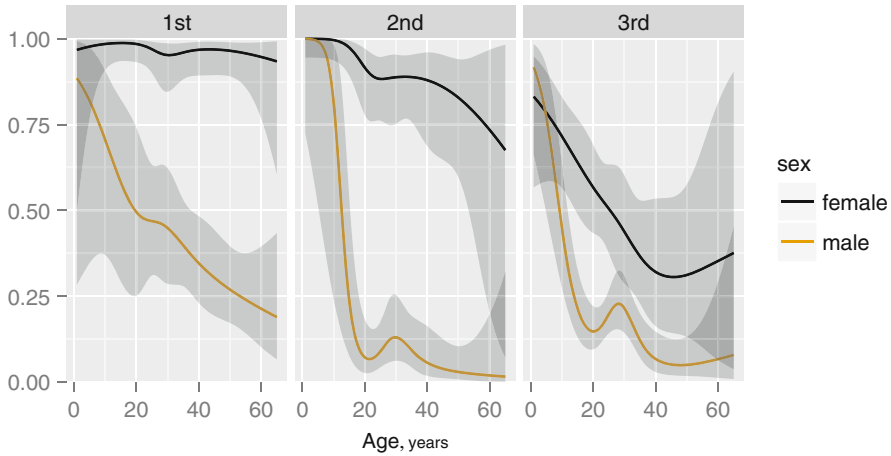


Fig. 12.5 Effects of predictors on probability of survival of Titanic passengers, estimated for zero siblings or spouses

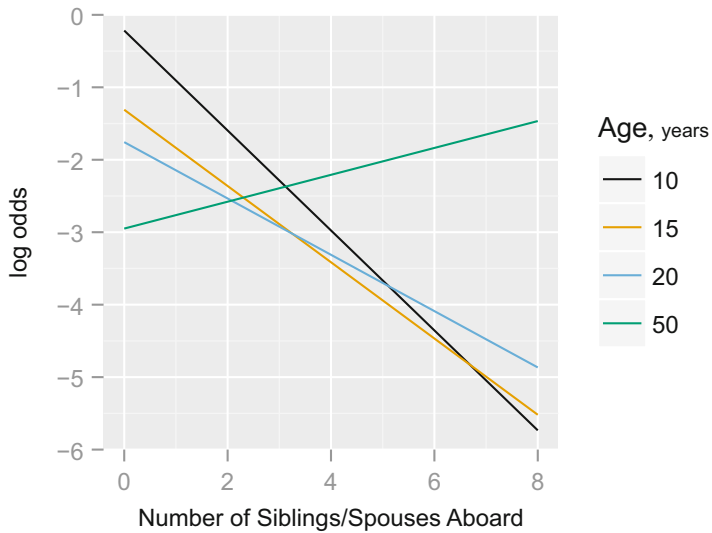


Fig. 12.6 Effect of number of siblings and spouses on the log odds of surviving, for third class males

for testing three-way interactions, in the belief that these tests would replicate well.

```
f ← update(f, x=TRUE, y=TRUE)
# x=TRUE, y=TRUE adds raw data to fit object so can bootstrap
set.seed(131) # so can replicate re-samples
latex(validate(f, B=200), digits=2, size='Ssize')
```

Index	Original Training Sample	Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.76	0.77	0.74	0.03	0.72	200
R^2	0.55	0.58	0.53	0.05	0.50	200
Intercept	0.00	0.00	-0.08	0.08	-0.08	200
Slope	1.00	1.00	0.87	0.13	0.87	200
E_{\max}	0.00	0.00	0.05	0.05	0.05	200
D	0.53	0.56	0.50	0.06	0.46	200
U	0.00	0.00	0.01	-0.01	0.01	200
Q	0.53	0.56	0.49	0.07	0.46	200
B	0.13	0.13	0.13	-0.01	0.14	200
g	2.43	2.75	2.37	0.37	2.05	200
g_p	0.37	0.37	0.35	0.02	0.35	200

```
cal ← calibrate(f, B=200) # Figure 12.7
plot(cal, subtitles=FALSE)
```

```
n=1046 Mean absolute error=0.009 Mean squared error=0.00012
0.9 Quantile of absolute error=0.017
```

The output of `validate` indicates minor overfitting. Overfitting would have been worse had the risk factors not been so strong. The closeness of the calibration curve to the 45° line in Figure 12.7 demonstrates excellent validation on an absolute probability scale. But the extent of missing data casts some doubt on the validity of this model, and on the efficiency of its parameter estimates.

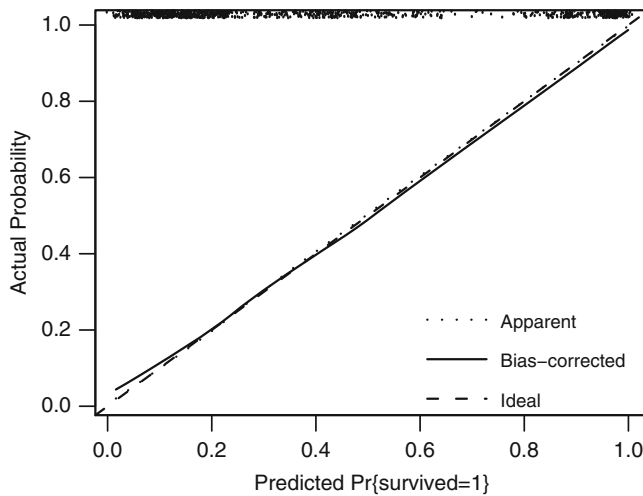


Fig. 12.7 Bootstrap overfitting-corrected loess nonparametric calibration curve for casewise deletion model

12.4 Examining Missing Data Patterns

The first step to dealing with missing data is understanding the patterns of missing values. To do this we use the `Hmisc` library's `naclus` and `naplot` functions, and the recursive partitioning library of Atkinson and Therneau. Below `naclus` tells us which variables tend to be missing on the same persons, and it computes the proportion of missing values for each variable. The `rpart` function derives a tree to predict which types of passengers tended to have age missing.

```
na.patterns <- naclus(titanic3)
require(rpart)      # Recursive partitioning package
```

```
who.na <- rpart(is.na(age) ~ sex + pclass + survived +
               sibsp + parch, data=titanic3, minbucket=15)
naplot(na.patterns, 'na per var')
plot(who.na, margin=.1); text(who.na) # Figure 12.8
plot(na.patterns)
```

We see in Figure 12.8 that age tends to be missing on the same passengers as the body bag identifier, and that it is missing in only 0.09 of first or second class passengers. The category of passengers having the highest fraction of missing ages is third class passengers having no parents or children on board. Below we use `Hmisc`'s `summary.formula` function to plot simple descriptive statistics on the fraction of missing ages, stratified by other variables. We see that without adjusting for other variables, age is slightly more missing on nonsurviving passengers.

```
plot(summary(is.na(age) ~ sex + pclass + survived +
            sibsp + parch, data=t3)) # Figure 12.9
```

Let us derive a logistic model to predict missingness of age, to see if the survival bias maintains after adjustment for the other variables.

```
m <- lrm(is.na(age) ~ sex * pclass + survived + sibsp + parch,
        data=t3)
print(m, latex=TRUE, needspace='2in')
```

Logistic Regression Model

```
lrm(formula = is.na(age) ~ sex * pclass + survived + sibsp +
    parch, data = t3)
```

		Model Likelihood Ratio Test	Discrimination Indexes	Rank Discrim. Indexes
Obs	1309	LR χ^2 114.99	R^2 0.133	C 0.703
FALSE	1046	d.f. 8	g 1.015	D_{xy} 0.406
TRUE	263	$\Pr(> \chi^2) < 0.0001$	g_r 2.759	γ 0.452
$\max \left \frac{\partial \log L}{\partial \beta} \right $	5×10^{-6}		g_p 0.126	τ_a 0.131
			Brier 0.148	

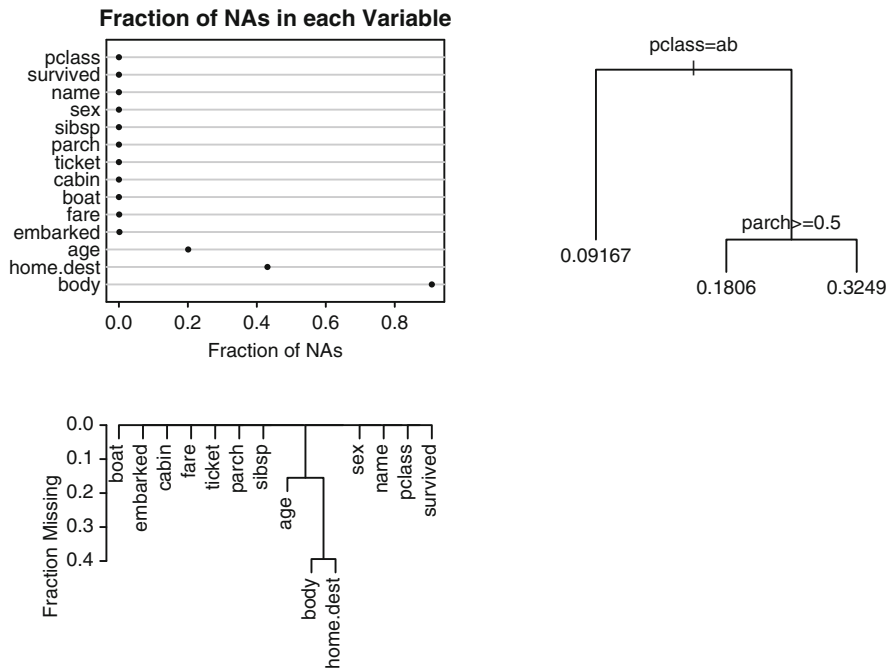


Fig. 12.8 Patterns of missing data. Upper left panel shows the fraction of observations missing on each predictor. Lower panel depicts a hierarchical cluster analysis of missingness combinations. The similarity measure shown on the Y-axis is the fraction of observations for which both variables are missing. Right panel shows the result of recursive partitioning for predicting `is.na(age)`. The `rpart` function found only strong patterns according to passenger class.

	Coef	S.E.	Wald Z	Pr(> Z)
Intercept	-2.2030	0.3641	-6.05	< 0.0001
sex=male	0.6440	0.3953	1.63	0.1033
pclass=2nd	-1.0079	0.6658	-1.51	0.1300
pclass=3rd	1.6124	0.3596	4.48	< 0.0001
survived	-0.1806	0.1828	-0.99	0.3232
sibsp	0.0435	0.0737	0.59	0.5548
parch	-0.3526	0.1253	-2.81	0.0049
sex=male * pclass=2nd	0.1347	0.7545	0.18	0.8583
sex=male * pclass=3rd	-0.8563	0.4214	-2.03	0.0422

```

latex(anova(m), file='', label='titanic-anova.na')
# Table 12.3
    
```

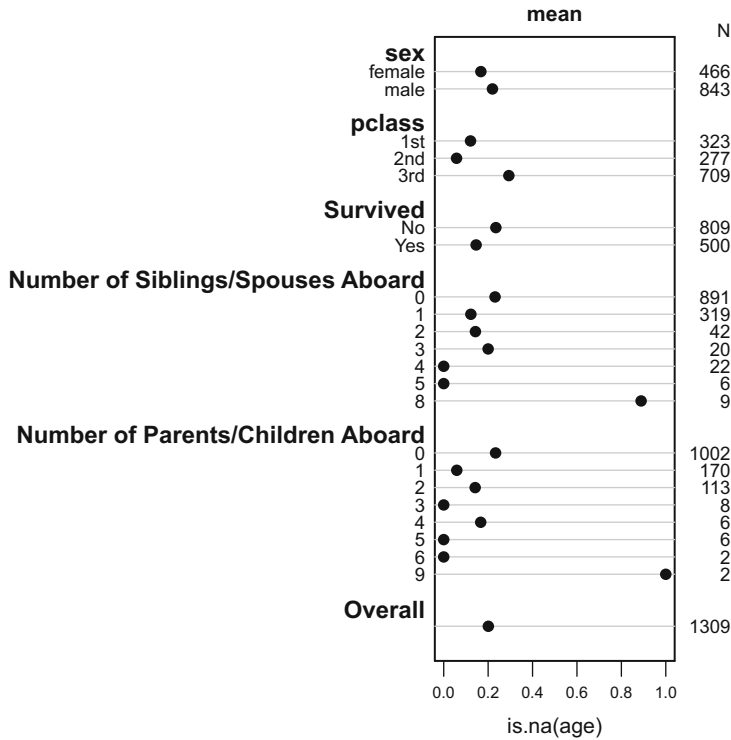


Fig. 12.9 Univariable descriptions of proportion of passengers with missing age

Fortunately, after controlling for other variables, Table 12.3 provides evidence that nonsurviving passengers are no more likely to have age missing. The only important predictors of missingness are `pclass` and `parch` (the more parents or children the passenger has on board, the less likely age was to be missing).

12.5 Multiple Imputation

Multiple imputation is expected to reduce bias in estimates as well as to provide an estimate of the variance–covariance matrix of $\hat{\beta}$ penalized for imputation. With multiple imputation, survival status can be used to impute missing ages, so the age relationship will not be as attenuated as with single conditional mean imputation. `aregImpute` The following uses the `Hmisc` package `aregImpute` function to do predictive mean matching, using van Buuren’s “Type 1” matching [85, Section 3.4.2] in conjunction with bootstrapping to incorporate all uncertainties, in the context of smooth additive imputation

Table 12.3 Wald Statistics for `is.na(age)`

	χ^2	d.f.	<i>P</i>
sex (Factor+Higher Order Factors)	5.61	3	0.1324
<i>All Interactions</i>	5.58	2	0.0614
pclass (Factor+Higher Order Factors)	68.43	4	< 0.0001
<i>All Interactions</i>	5.58	2	0.0614
survived	0.98	1	0.3232
sibsp	0.35	1	0.5548
parch	7.92	1	0.0049
sex × pclass (Factor+Higher Order Factors)	5.58	2	0.0614
TOTAL	82.90	8	< 0.0001

models. Sampling of donors is handled by distance weighting to yield better distributions of imputed values. By default, `aregImpute` does not transform `age` when it is being predicted from the other variables. Four knots are used to transform `age` when used to impute other variables (not needed here as no other missings were present in the variables of interest). Since the fraction of observations with missing `age` is $\frac{263}{1309} = 0.2$ we use 20 imputations.

```
set.seed(17)           # so can reproduce random aspects
mi ← aregImpute(~ age + sex + pclass +
               sibsp + parch + survived,
               data=t3, n.impute=20, nk=4, pr=FALSE)
```

```
mi
```

Multiple Imputation using Bootstrap and PMM

```
aregImpute(formula = ~age + sex + pclass + sibsp + parch + survived,
            data = t3, n.impute = 20, nk = 4, pr = FALSE)
```

```
n: 1309          p: 6      Imputations: 20          nk: 4
```

Number of NAs:

```
   age      sex  pclass  sibsp  parch survived
263      0      0      0      0      0      0
```

```
      type d.f.
age      s    1
sex      c    1
pclass   c    2
sibsp    s    2
parch    s    2
survived l    1
```

Transformation of Target Variables Forced to be Linear

R-squares for Predicting Non-Missing Values for Each Variable
Using Last Imputations of Predictors

```
age
0.295
```

```
# Print the first 10 imputations for the first 10 passengers
# having missing age
mi$imputed$age[1:10, 1:10]
```

	[,1]	[,2]	[,3]	[,4]	[,5]	[,6]	[,7]	[,8]	[,9]	[,10]
16	40	49	24	29	60.0	58	64	36	50	61
38	33	45	40	49	80.0	2	38	38	36	53
41	29	24	19	31	40.0	60	64	42	30	65
47	40	42	29	48	36.0	46	64	30	38	42
60	52	40	22	31	38.0	22	19	24	40	33
70	16	14	23	23	18.0	24	19	27	59	23
71	30	62	57	30	42.0	31	64	40	40	63
75	43	23	36	61	45.5	58	64	27	24	50
81	44	57	47	31	45.0	30	64	62	39	67
107	52	18	24	62	32.5	38	64	47	19	23

```
plot(mi)
Ecdf(t3$age, add=TRUE, col='gray', lwd=2,
      subtitles=FALSE)#Fig. 12.10
```

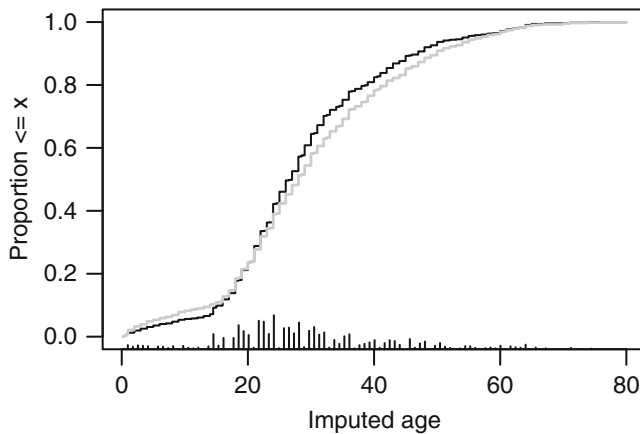


Fig. 12.10 Distributions of imputed and actual ages for the Titanic dataset. Imputed values are in black and actual ages in gray.

We now fit logistic models for five completed datasets. The `fit.mult.impute` function fits five models and examines the within- and between-imputation variances to compute an imputation-corrected variance-covariance matrix that is stored in the fit object `f.mi`. `fit.mult.impute` will also average the five $\hat{\beta}$ vectors, storing the result in `f.mi$coefficients`. The function also prints the ratio of imputation-corrected variances to average ordinary variances.

```
f.mi <- fit.mult.impute(
  survived ~ (sex + pclass + rcs(age,5))^2 +
  rcs(age,5)*sibsp,
```


Table 12.4 Wald Statistics for survived

	χ^2	d.f.	<i>P</i>
sex (Factor+Higher Order Factors)	240.42	7	< 0.0001
<i>All Interactions</i>	54.56	6	< 0.0001
pclass (Factor+Higher Order Factors)	114.21	12	< 0.0001
<i>All Interactions</i>	36.43	10	0.0001
age (Factor+Higher Order Factors)	50.37	20	0.0002
<i>All Interactions</i>	25.88	16	0.0557
<i>Nonlinear (Factor+Higher Order Factors)</i>	24.21	15	0.0616
sibsp (Factor+Higher Order Factors)	24.22	5	0.0002
<i>All Interactions</i>	12.86	4	0.0120
sex × pclass (Factor+Higher Order Factors)	30.99	2	< 0.0001
sex × age (Factor+Higher Order Factors)	11.38	4	0.0226
<i>Nonlinear</i>	8.15	3	0.0430
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	8.15	3	0.0430
pclass × age (Factor+Higher Order Factors)	5.30	8	0.7246
<i>Nonlinear</i>	4.63	6	0.5918
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	4.63	6	0.5918
age × sibsp (Factor+Higher Order Factors)	12.86	4	0.0120
<i>Nonlinear</i>	1.84	3	0.6058
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	1.84	3	0.6058
TOTAL NONLINEAR	24.21	15	0.0616
TOTAL INTERACTION	67.12	18	< 0.0001
TOTAL NONLINEAR + INTERACTION	70.99	21	< 0.0001
TOTAL	298.78	26	< 0.0001

```

lrm, mi, data=t3, pr=FALSE)
latex(anova(f.mi), file='', label='titanic-anova.mi ',
      size='small') # Table 12.4

```

The Wald χ^2 for age is reduced by accounting for imputation but is increased (by a lesser amount) by using patterns of association with survival status to impute missing age. The Wald tests are all adjusted for multiple imputation. Now examine the fitted age relationship using multiple imputation vs. casewise deletion.

```

p1 ← Predict(f, age, pclass, sex, sibsp=0, fun=plogis)
p2 ← Predict(f.mi, age, pclass, sex, sibsp=0, fun=plogis)
p ← rbind('Casewise Deletion'=p1, 'Multiple Imputation'=p2)
ggplot(p, groups='sex', ylab='Probability of Surviving')
# Figure 12.11

```

12.6 Summarizing the Fitted Model

In this section we depict the model fitted using multiple imputation, by computing odds ratios and by showing various predicted values. For age, the odds ratio for an increase from 1 year old to 30 years old is computed, instead of the default odds ratio based on outer quartiles of age. The estimated odds

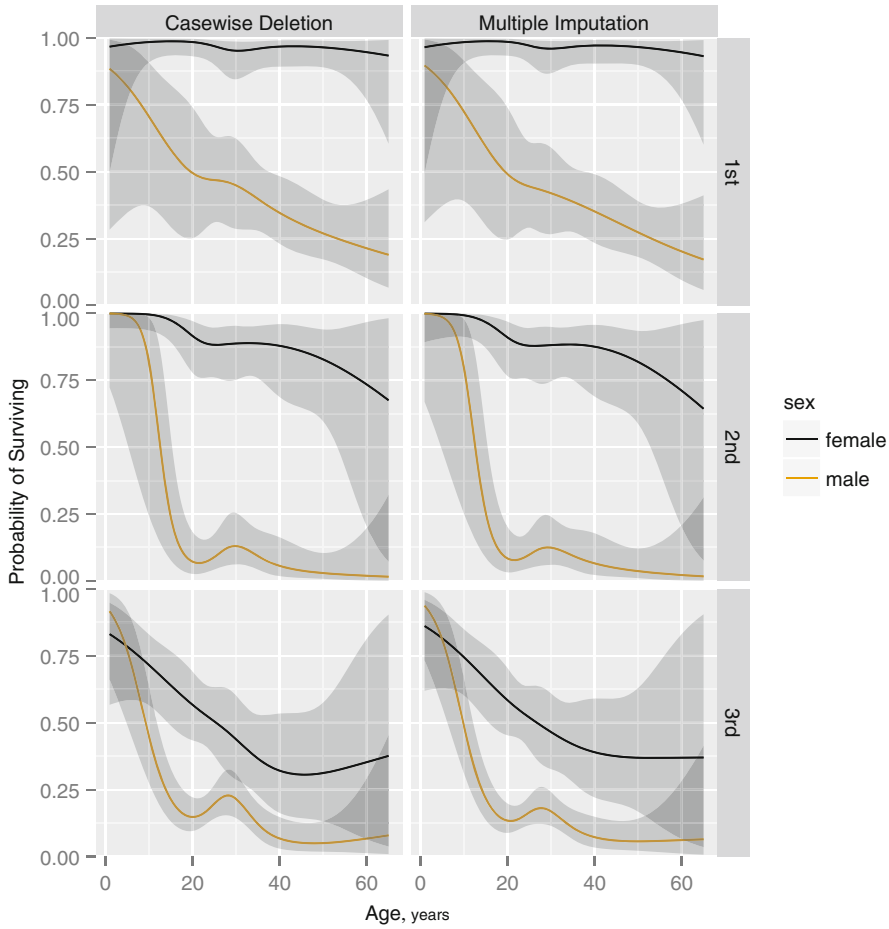


Fig. 12.11 Predicted probability of survival for males from fit using casewise deletion again (top) and multiple random draw imputation (bottom). Both sets of predictions are for `sibsp=0`.

ratios are very dependent on the levels of interacting factors, so Figure 12.12 depicts only one of many patterns.

```
# Get predicted values for certain types of passengers
s <- summary(f.mi, age=c(1,30), sibsp=0:1)
# override default ranges for 3 variables
plot(s, log=TRUE, main='') # Figure 12.12
```

Now compute estimated probabilities of survival for a variety of settings of the predictors.

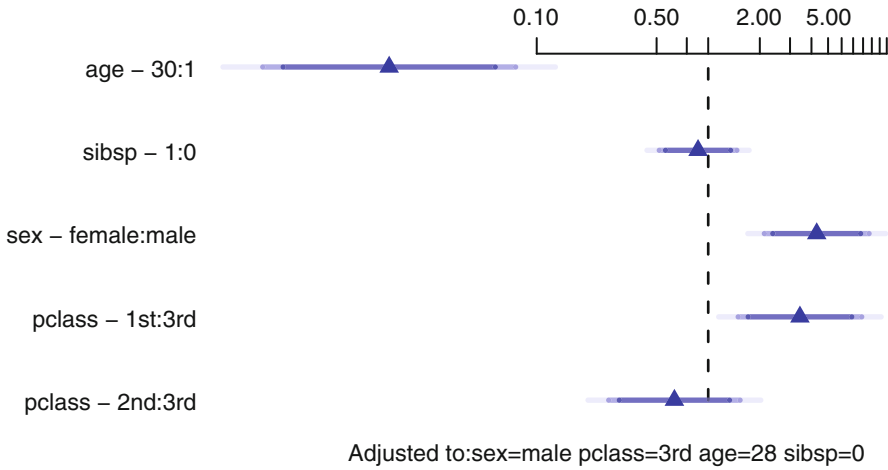


Fig. 12.12 Odds ratios for some predictor settings

```
phat <- predict(f.mi,
               combos <-
                 expand.grid(age=c(2,21,50), sex=levels(t3$sex),
                             pclass=levels(t3$pclass),
                             sibsp=0), type='fitted')
# Can also use Predict(f.mi, age=c(2,21,50), sex, pclass,
#                       sibsp=0, fun=plogis)$yhat
options(digits=1)
data.frame(combos, phat)
```

	age	sex	pclass	sibsp	phat
1	2	female	1st	0	0.97
2	21	female	1st	0	0.98
3	50	female	1st	0	0.97
4	2	male	1st	0	0.88
5	21	male	1st	0	0.48
6	50	male	1st	0	0.27
7	2	female	2nd	0	1.00
8	21	female	2nd	0	0.90
9	50	female	2nd	0	0.82
10	2	male	2nd	0	1.00
11	21	male	2nd	0	0.08
12	50	male	2nd	0	0.04
13	2	female	3rd	0	0.85
14	21	female	3rd	0	0.57
15	50	female	3rd	0	0.37
16	2	male	3rd	0	0.91
17	21	male	3rd	0	0.13
18	50	male	3rd	0	0.06

```
options(digits=5)
```

We can also get predicted values by creating an R function that will evaluate the model on demand.

```

pred.logit ← Function(f.mi)
# Note: if don't define sibsp to pred.logit, defaults to 0
# normally just type the function name to see its body
latex(pred.logit, file='', type='Sinput', size='small',
      width.cutoff=49)

```

```

pred.logit ← function (sex = "male", pclass = "3rd",
  age = 28, sibsp = 0)
{
  3.2427671 - 0.95431809 * (sex == "male") + 5.4086505 *
    (pclass == "2nd") - 1.3378623 * (pclass ==
    "3rd") + 0.091162649 * age - 0.00031204327 *
    pmax(age - 6, 0)^3 + 0.0021750413 * pmax(age -
    21, 0)^3 - 0.0027627032 * pmax(age - 27, 0)^3 +
    0.0009805137 * pmax(age - 36, 0)^3 - 8.0808484e-05 *
    pmax(age - 55.8, 0)^3 - 1.1567976 * sibsp +
    (sex == "male") * (-0.46061284 * (pclass ==
    "2nd") + 2.0406523 * (pclass == "3rd")) +
    (sex == "male") * (-0.22469066 * age + 0.00043708296 *
    pmax(age - 6, 0)^3 - 0.0026505136 * pmax(age -
    21, 0)^3 + 0.0031201404 * pmax(age - 27,
    0)^3 - 0.00097923749 * pmax(age - 36,
    0)^3 + 7.2527708e-05 * pmax(age - 55.8,
    0)^3) + (pclass == "2nd") * (-0.46144083 *
    age + 0.00070194849 * pmax(age - 6, 0)^3 -
    0.0034726662 * pmax(age - 21, 0)^3 + 0.0035255387 *
    pmax(age - 27, 0)^3 - 0.0007900891 * pmax(age -
    36, 0)^3 + 3.5268151e-05 * pmax(age - 55.8,
    0)^3) + (pclass == "3rd") * (-0.17513289 *
    age + 0.00035283358 * pmax(age - 6, 0)^3 -
    0.0023049372 * pmax(age - 21, 0)^3 + 0.0028978962 *
    pmax(age - 27, 0)^3 - 0.00105145 * pmax(age -
    36, 0)^3 + 0.00010565735 * pmax(age - 55.8,
    0)^3) + sibsp * (0.040830773 * age - 1.5627772e-05 *
    pmax(age - 6, 0)^3 + 0.00012790256 * pmax(age -
    21, 0)^3 - 0.00025039385 * pmax(age - 27,
    0)^3 + 0.00017871701 * pmax(age - 36, 0)^3 -
    4.0597949e-05 * pmax(age - 55.8, 0)^3)
}

```

```

# Run the newly created function
plogis(pred.logit(age=c(2,21,50), sex='male', pclass='3rd'))

```

```
[1] 0.914817 0.132640 0.056248
```

A nomogram could be used to obtain predicted values manually, but this is not feasible when so many interaction terms are present.

Chapter 13

Ordinal Logistic Regression

13.1 Background

Many medical and epidemiologic studies incorporate an ordinal response variable. In some cases an ordinal response Y represents levels of a standard measurement scale such as severity of pain (none, mild, moderate, severe). In other cases, ordinal responses are constructed by specifying a hierarchy of separate endpoints. For example, clinicians may specify an ordering of the severity of several component events and assign patients to the worst event present from among none, heart attack, disabling stroke, and death. Still another use of ordinal response methods is the application of rank-based methods to continuous responses so as to obtain robust inferences. For example, the proportional odds model described later allows for a continuous Y and is really a generalization of the Wilcoxon–Mann–Whitney rank test. Thus the semiparametric proportional odds model is a direct competitor of ordinary linear models.

There are many variations of logistic models used for predicting an ordinal response variable Y . All of them have the advantage that they do not assume a spacing between levels of Y . In other words, the same regression coefficients and P -values result from an analysis of a response variable having levels 0, 1, 2 when the levels are recoded 0, 1, 20. Thus ordinal models use only the rank-ordering of values of Y .

In this chapter we consider two of the most popular ordinal logistic models, the proportional odds (PO) form of an ordinal logistic model⁶⁴⁷ and the forward continuation ratio (CR) ordinal logistic model.¹⁹⁰ Chapter 15 deals with a wider variety of ordinal models with emphasis on analysis of continuous Y .

1

13.2 Ordinality Assumption

A basic assumption of all commonly used ordinal regression models is that the response variable behaves in an ordinal fashion with respect to each predictor. Assuming that a predictor X is linearly related to the log odds of some appropriate event, a simple way to check for ordinality is to plot the mean of X stratified by levels of Y . These means should be in a consistent order. If for many of the X s, two adjacent categories of Y do not distinguish the means, that is evidence that those levels of Y should be pooled.

One can also estimate the mean or expected value of $X|Y = j$ ($E(X|Y = j)$) given that the ordinal model assumptions hold. This is a useful tool for checking those assumptions, at least in an unadjusted fashion. For simplicity, assume that X is discrete, and let $P_{jx} = \Pr(Y = j|X = x)$ be the probability that $Y = j$ given $X = x$ that is dictated from the model being fitted, with X being the only predictor in the model. Then

$$\begin{aligned} \Pr(X = x|Y = j) &= \Pr(Y = j|X = x)\Pr(X = x)/\Pr(Y = j) \\ E(X|Y = j) &= \sum_x xP_{jx} \Pr(X = x)/\Pr(Y = j), \end{aligned} \quad (13.1)$$

and the expectation can be estimated by

$$\hat{E}(X|Y = j) = \sum_x x\hat{P}_{jx}f_x/g_j, \quad (13.2)$$

where \hat{P}_{jx} denotes the estimate of P_{jx} from the fitted one-predictor model (for inner values of Y in the PO models, these probabilities are differences between terms given by Equation 13.4 below), f_x is the frequency of $X = x$ in the sample of size n , and g_j is the frequency of $Y = j$ in the sample. This estimate can be computed conveniently without grouping the data by X . For n subjects let the n values of X be x_1, x_2, \dots, x_n . Then

$$\hat{E}(X|Y = j) = \sum_{i=1}^n x_i\hat{P}_{jx_i}/g_j. \quad (13.3)$$

Note that if one were to compute differences between conditional means of X and the conditional means of X given PO, and if furthermore the means were conditioned on $Y \geq j$ instead of $Y = j$, the result would be proportional to means of score residuals defined later in Equation 13.6.

13.3 Proportional Odds Model

13.3.1 Model

The most commonly used ordinal logistic model was described in Walker and Duncan⁶⁴⁷ and later called the *proportional odds (PO) model* by McCullagh.⁴⁴⁹ The PO model is best stated as follows, for a response variable having levels $0, 1, 2, \dots, k$:

$$\Pr[Y \geq j|X] = \frac{1}{1 + \exp[-(\alpha_j + X\beta)]}, \quad (13.4)$$

where $j = 1, 2, \dots, k$. Some authors write the model in terms of $Y \leq j$. Our formulation makes the model coefficients consistent with the binary logistic model. There are k intercepts (α s). For fixed j , the model is an ordinary logistic model for the event $Y \geq j$. By using a common vector of regression coefficients β connecting probabilities for varying j , the PO model allows for parsimonious modeling of the distribution of Y .

There is a nice connection between the PO model and the Wilcoxon–Mann–Whitney two-sample test: when there is a single predictor X_1 that is binary, the numerator of the score test for testing $H_0 : \beta_1 = 0$ is proportional to the two-sample test statistic [664, pp. 2258–2259].

2

13.3.2 Assumptions and Interpretation of Parameters

There is an implicit assumption in the PO model that the regression coefficients (β) are independent of j , the cutoff level for Y . One could say that there is no $X \times Y$ interaction if PO holds. For a specific Y -cutoff j , the model has the same assumptions as the binary logistic model (Section 10.1.1). That is, the model in its simplest form assumes the log odds that $Y \geq j$ is linearly related to each X and that there is no interaction between the X s.

In designing clinical studies, one sometimes hears the statement that an ordinal outcome should be avoided since statistical tests of patterns of those outcomes are hard to interpret. In fact, one interprets effects in the PO model using ordinary odds ratios. The difference is that a single odds ratio is assumed to apply equally to *all* events $Y \geq j, j = 1, 2, \dots, k$. If linearity and additivity hold, the $X_m + 1 : X_m$ odds ratio for $Y \geq j$ is $\exp(\beta_m)$, whatever the cutoff j .

The proportional hazards assumption is frequently violated, just as the assumptions of normality of residuals with equal variance in ordinary regression are frequently violated, but the PO model can still be useful and powerful in this situation. As stated by Senn and Julious⁵⁶⁴,

Clearly, the dependence of the proportional odds model on the assumption of proportionality can be over-stressed. Suppose that two different statisticians would cut the same three-point scale at different cut points. It is hard to see how anybody who could accept either dichotomy could object to the compromise answer produced by the proportional odds model.

Sometimes it helps in interpreting the model to estimate the mean Y as a function of one or more predictors, even though this assumes a spacing for the Y -levels.

3

13.3.3 Estimation

The PO model is fitted using MLE on a somewhat complex likelihood function that is dependent on differences in logistic model probabilities. The estimation process forces the α s to be in descending order.

13.3.4 Residuals

Schoenfeld residuals⁵⁵⁷ are very effective²³³ in checking the proportional hazards assumption in the Cox¹³² survival model. For the PO model one could analogously compute each subject's contribution to the first derivative of the log likelihood function with respect to β_m , average them separately by levels of Y , and examine trends in the residual plots as in Section 20.6.2. A few examples have shown that such plots are usually hard to interpret. Easily interpreted score residual plots for the PO model can be constructed, however, by using the fitted PO model to predict a series of binary events $Y \geq j, j = 1, 2, \dots, k$, using the corresponding predicted probabilities

$$\hat{P}_{ij} = \frac{1}{1 + \exp[-(\hat{\alpha}_j + X_i \hat{\beta})]}, \quad (13.5)$$

where X_i stands for a vector of predictors for subject i . Then, after forming an indicator variable for the event currently being predicted ($[Y_i \geq j]$), one computes the score (first derivative) components U_{im} from an ordinary binary logistic model:

$$U_{im} = X_{im}([Y_i \geq j] - \hat{P}_{ij}), \quad (13.6)$$

for the subject i and predictor m . Then, for each column of U , plot the mean $\bar{U}_{.m}$ and confidence limits, with Y (i.e., j) on the x -axis. For each predictor the trend against j should be flat if PO holds. ^aIn binary logistic regression, *partial residuals* are very useful as they allow the analyst to fit linear effects

^a If $\hat{\beta}$ were derived from separate binary fits, all $\bar{U}_{.m} \equiv 0$.

for all the predictors but then to nonparametrically estimate the true transformation that each predictor requires (Section 10.4). The partial residual is defined as follows, for the i th subject and m th predictor variable.^{115, 373}

$$r_{im} = \hat{\beta}_m X_{im} + \frac{Y_i - \hat{P}_i}{\hat{P}_i(1 - \hat{P}_i)}, \quad (13.7)$$

where

$$\hat{P}_i = \frac{1}{1 + \exp[-(\alpha + X_i \hat{\beta})]}. \quad (13.8)$$

A smoothed plot (e.g., using the moving linear regression algorithm in `loess`¹¹¹) of X_{im} against r_{im} provides a nonparametric estimate of how X_m relates to the log relative odds that $Y = 1|X_m$. For ordinal Y , we just need to compute binary model partial residuals for all cutoffs j :

$$r_{im} = \hat{\beta}_m X_{im} + \frac{[Y_i \geq j] - \hat{P}_{ij}}{\hat{P}_{ij}(1 - \hat{P}_{ij})}, \quad (13.9)$$

then to make a plot for each m showing smoothed partial residual curves for all j , looking for similar shapes and slopes for a given predictor for all j . Each curve provides an estimate of how X_m relates to the relative log odds that $Y \geq j$. Since partial residuals allow examination of predictor transformations (linearity) while simultaneously allowing examination of PO (parallelism), partial residual plots are generally preferred over score residual plots for ordinal models.

Li and Shepherd⁴⁰² have a residual for ordinal models that serves for the entire range of Y without the need to consider cutoffs. Their residual is useful for checking functional form of predictors but not the proportional odds assumption.

13.3.5 Assessment of Model Fit

Peterson and Harrell⁵⁰² developed score and likelihood ratio tests for testing the PO assumption. The score test is used in the SAS PROC LOGISTIC,⁵⁴⁰ but its extreme anti-conservatism in many cases can make it unreliable.⁵⁰²

For determining whether the PO assumption is likely to be satisfied for each predictor separately, there are several graphics that are useful. One is the graph comparing means of $X|Y$ with and without assuming PO, as described in Section 13.2 (see Figure 14.2 for an example). Another is the simple method of stratifying on each predictor and computing the logits of all proportions of the form $Y \geq j, j = 1, 2, \dots, k$. When proportional odds holds, the differences in logits between different values of j should be the same at all levels of X ,

because the model dictates that $\text{logit}(Y \geq j|X) - \text{logit}(Y \geq i|X) = \alpha_j - \alpha_i$, for any constant X . An example of this is in Figure 13.1.

```
require(Hmisc)

getHdata(support)
sfdm <- as.integer(support$sfdm2) - 1
sf <- function(y)
  c('Y ≥ 1'=qlogis(mean(y ≥ 1)), 'Y ≥ 2'=qlogis(mean(y ≥ 2)),
    'Y ≥ 3'=qlogis(mean(y ≥ 3)))
s <- summary(sfdm ~ adlsc + sex + age + meanbp, fun=sf,
             data=support)
plot(s, which=1:3, pch=1:3, xlab='logit', vnames='names',
     main='', width.factor=1.5) # Figure 13.1
```

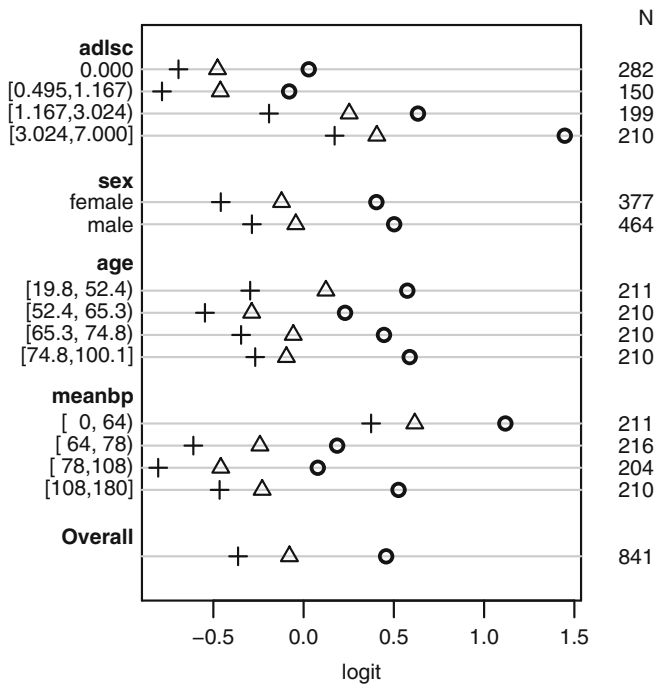


Fig. 13.1 Checking PO assumption separately for a series of predictors. The circle, triangle, and plus sign correspond to $Y \geq 1, 2, 3$, respectively. PO is checked by examining the vertical constancy of distances between any two of these three symbols. Response variable is the severe functional disability scale `sfdm2` from the 1000-patient SUPPORT dataset, with the last two categories combined because of low frequency of coma/intubation.

When Y is continuous or almost continuous and X is discrete, the PO model assumes that the logit of the cumulative distribution function of Y is parallel

across categories of X . The corresponding, more rigid, assumptions of the ordinary linear model (here, parametric ANOVA) are parallelism and linearity if the normal inverse cumulative distribution function across categories of X . As an example consider the web site's `diabetes` dataset, where we consider the distribution of log glycohemoglobin across subjects' body frames.

```
getHdata(diabetes)
a <- Ecdf(~ log(glyhb), group=frame, fun=qnorm,
          xlab='log(HbA1c)', label.curves=FALSE, data=diabetes,
          ylab=expression(paste(Phi^-1, (F[n](x)))) # Fig. 13.2
b <- Ecdf(~ log(glyhb), group=frame, fun=qlogis,
          xlab='log(HbA1c)', label.curves=list(keys='lines'),
          data=diabetes, ylab=expression(logit(F[n](x))))
print(a, more=TRUE, split=c(1,1,2,1))
print(b, split=c(2,1,2,1))
```

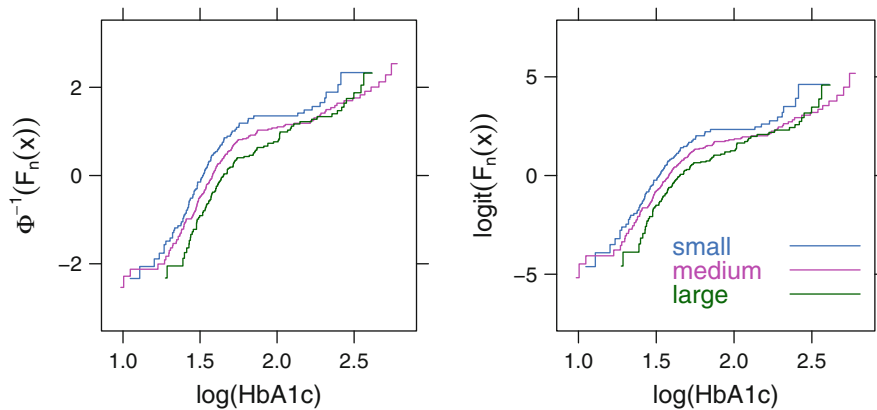


Fig. 13.2 Transformed empirical cumulative distribution functions stratified by body frame in the `diabetes` dataset. Left panel: checking all assumptions of the parametric ANOVA. Right panel: checking all assumptions of the PO model (here, Kruskal–Wallis test).

One could conclude the right panel of Figure 13.2 displays more parallelism than the left panel displays linearity, so the assumptions of the PO model are better satisfied than the assumptions of the ordinary linear model.

Chapter 14 has many examples of graphics for assessing fit of PO models. Regarding assessment of linearity and additivity assumptions, splines, partial residual plots, and interaction tests are among the best tools. Fagerland and Hosmer¹⁸² have a good review of goodness-of-fit tests for the PO model.

13.3.6 Quantifying Predictive Ability

The R_N^2 coefficient is really computed from the model LR χ^2 (χ^2 added to a model containing only the k intercept parameters) to describe the model's predictive power. The Somers' D_{xy} rank correlation between $X\hat{\beta}$ and Y is an easily interpreted measure of predictive discrimination. Since it is a rank measure, it does not matter which intercept α is used in the calculation. The probability of concordance, c , is also a useful measure. Here one takes all possible pairs of subjects having differing Y values and computes the fraction of such pairs for which the values of $X\hat{\beta}$ are in the same direction as the two Y values. c could be called a generalized ROC area in this setting. As before, $D_{xy} = 2(c - 0.5)$. Note that D_{xy} , c , and the Brier score B can easily be computed for various dichotomizations of Y , to investigate predictive ability in more detail.

13.3.7 Describing the Fitted Model

As discussed in Section 5.1, models are best described by computing predicted values or differences in predicted values. For PO models there are four and sometimes five types of relevant predictions:

1. $\text{logit}[Y \geq j|X]$, i.e., the linear predictor
2. $\text{Prob}[Y \geq j|X]$
3. $\text{Prob}[Y = j|X]$
4. Quantiles of $Y|X$ (e.g., the median^b)
5. $E(Y|X)$ if Y is interval scaled.

For the first two quantities above a good default choice for j is the middle category. Partial effect plots are as simple to draw for PO models as they are for binary logistic models. Other useful graphics, as before, are odds ratio charts and nomograms. For the latter, an axis displaying the predicted mean makes the model more interpretable, under scaling assumptions on Y .

13.3.8 Validating the Fitted Model

The PO model is validated much the same way as the binary logistic model (see Section 10.9). For estimating an overfitting-corrected calibration curve (Section 10.11) one estimates $\text{Pr}(Y \geq j|X)$ using one j at a time.

^b If Y does not have very many levels, the median will be a discontinuous function of X and may not be satisfactory.

13.3.9 R Functions

The `rms` package's `lrm` and `orm` functions fit the PO model directly, assuming that the levels of the response variable (e.g., the `levels` of a `factor` variable) are listed in the proper order. `lrm` is intended to be used for the case where the number of unique values of Y are less than a few dozen whereas `orm` handles the continuous Y case efficiently, as well as allowing for links other than the logit. See Chapter 15 for more information.

If the response is numeric, `lrm` assumes the numeric codes properly order the responses. If it is a character vector and is not a `factor`, `lrm` assumes the correct ordering is alphabetic. Of course `ordered` variables in R are appropriate response variables for ordinal regression. The `predict` function (`predict.lrm`) can compute all the quantities listed in Section 13.3.7 except for quantiles.

The R functions `popower` and `posamsize` (in the `Hmisc` package) compute power and sample size estimates for ordinal responses using the proportional odds model.

The function `plot.xmean.ordinaly` in `rms` computes and graphs the quantities described in Section 13.2. It plots simple Y -stratified means overlaid with $\hat{E}(X|Y = j)$, with j on the x -axis. The \hat{E} s are computed for both PO and continuation ratio ordinal logistic models. The `Hmisc` package's `summary.formula` function is also useful for assessing the PO assumption (Figure 13.1). Generic `rms` functions such as `validate`, `calibrate`, and `nomogram` work with PO model fits from `lrm` as long as the analyst specifies which intercept(s) to use. `rms` has a special function generator `Mean` for constructing an easy-to-use function for getting the predicted mean Y from a PO model. This is handy with `plot` and `nomogram`. If the fit has been run through the `bootcov` function, it is easy to use the `Predict` function to estimate bootstrap confidence limits for predicted means.

13.4 Continuation Ratio Model

13.4.1 Model

Unlike the PO model, which is based on *cumulative* probabilities, the continuation ratio (CR) model is based on *conditional* probabilities. The (forward) CR model^{31, 52, 190} is stated as follows for $Y = 0, \dots, k$.

$$\begin{aligned} \Pr(Y = j|Y \geq j, X) &= \frac{1}{1 + \exp[-(\theta_j + X\gamma)]} \\ \text{logit}(Y = 0|Y \geq 0, X) &= \text{logit}(Y = 0|X) \\ &= \theta_0 + X\gamma \end{aligned} \tag{13.10}$$

$$\begin{aligned} \text{logit}(Y = 1|Y \geq 1, X) &= \theta_1 + X\gamma \\ &\dots \\ \text{logit}(Y = k - 1|Y \geq k - 1, X) &= \theta_{k-1} + X\gamma. \end{aligned}$$

The CR model has been said to be likely to fit ordinal responses when subjects have to “pass through” one category to get to the next. The CR model is a discrete version of the Cox proportional hazards model. The discrete hazard function is defined as $\Pr(Y = j|Y \geq j)$.

13.4.2 Assumptions and Interpretation of Parameters

The CR model assumes that the vector of regression coefficients, γ , is the same regardless of which conditional probability is being computed.

One could say that there is no $X \times$ condition interaction if the CR model holds. For a specific condition $Y \geq j$, the model has the same assumptions as the binary logistic model (Section 10.1.1). That is, the model in its simplest form assumes that the log odds that $Y = j$ conditional on $Y \geq j$ is linearly related to each X and that there is no interaction between the X s.

A single odds ratio is assumed to apply equally to *all* conditions $Y \geq j$, $j = 0, 1, 2, \dots, k - 1$. If linearity and additivity hold, the $X_m + 1 : X_m$ odds ratio for $Y = j$ is $\exp(\beta_m)$, whatever the conditioning event $Y \geq j$.

To compute $\Pr(Y > 0|X)$ from the CR model, one only needs to take one minus $\Pr(Y = 0|X)$. To compute other unconditional probabilities from the CR model, one must multiply the conditional probabilities. For example, $\Pr(Y > 1|X) = \Pr(Y > 1|X, Y \geq 1) \times \Pr(Y \geq 1|X) = [1 - \Pr(Y = 1|Y \geq 1, X)][1 - \Pr(Y = 0|X)] = [1 - 1/(1 + \exp[-(\theta_1 + X\gamma)])][1 - 1/(1 + \exp[-(\theta_0 + X\gamma)])]$.

13.4.3 Estimation

Armstrong and Sloan³¹ and Berridge and Whitehead⁵² showed how the CR model can be fitted using an ordinary binary logistic model likelihood function, after certain rows of the X matrix are duplicated and a new binary Y vector is constructed. For each subject, one constructs separate records by considering successive conditions $Y \geq 0, Y \geq 1, \dots, Y \geq k - 1$ for a response variable with values $0, 1, \dots, k$. The binary response for each applicable condition or “cohort” is set to 1 if the subject failed at the current “cohort” or “risk set,” that is, if $Y = j$ where the cohort being considered is $Y \geq j$. The constructed cohort variable is carried along with the new X and Y . This variable is considered to be categorical and its coefficients are fitted by adding $k - 1$ dummy variables to the binary logistic model. For ease of computation,

the CR model is restated as follows, with the first cohort used as the reference cell.

$$\Pr(Y = j|Y \geq j, X) = \frac{1}{1 + \exp[-(\alpha + \theta_j + X\gamma)]}. \quad (13.11)$$

Here α is an overall intercept, $\theta_0 \equiv 0$, and $\theta_1, \dots, \theta_{k-1}$ are increments from α .

13.4.4 Residuals

To check CR model assumptions, binary logistic model partial residuals are again valuable. We separately fit a sequence of binary logistic models using a series of binary events and the corresponding applicable (increasingly small) subsets of subjects, and plot smoothed partial residuals against X for all of the binary events. Parallelism in these plots indicates that the CR model's constant γ assumptions are satisfied.

13.4.5 Assessment of Model Fit

The partial residual plots just described are very useful for checking the constant slope assumption of the CR model. The next section shows how to test this assumption formally. Linearity can be assessed visually using the smoothed partial residual plot, and interactions between predictors can be tested as usual.

13.4.6 Extended CR Model

The PO model has been extended by Peterson and Harrell⁵⁰² to allow for unequal slopes for some or all of the X s for some or all levels of Y . This partial PO model requires specialized software. The CR model can be extended more easily. In R notation, the ordinary CR model is specified as

```
y ~ cohort + X1 + X2 + X3 + ...
```

with `cohort` denoting a polytomous variable. The CR model can be extended to allow for some or all of the β s to change with the cohort or Y -cutoff.³¹ Suppose that non-constant slope is allowed for `x1` and `x2`. The R notation for the extended model would be

```
y ~ cohort*(X1 + X2) + X3
```

The extended CR model is a discrete version of the Cox survival model with time-dependent covariables.

There is nothing about the CR model that makes it fit a given dataset better than other ordinal models such as the PO model. The real benefit of the CR model is that using standard binary logistic model software one can flexibly specify how the equal-slopes assumption can be relaxed.

13.4.7 Role of Penalization in Extended CR Model

As demonstrated in the upcoming case study, penalized MLE is invaluable in allowing the model to be extended into an unequal-slopes model insofar as the information content in the data will support. Faraway¹⁸⁶ has demonstrated how all data-driven steps of the modeling process increase the real variance in “final” parameter estimates, when one estimates variances without assuming that the final model was prespecified. For ordinal regression modeling, the most important modeling steps are (1) choice of predictor variables, (2) selecting or modeling predictor transformations, and (3) allowance for unequal slopes across Y -cutoffs (i.e., non-PO or non-CR). Regarding Steps (2) and (3) one is tempted to rely on graphical methods such as residual plots to make detours in the strategy, but it is very difficult to estimate variances or to properly penalize assessments of predictive accuracy for subjective modeling decisions. Regarding (1), shrinkage has been proven to work better than stepwise variable selection when one is attempting to build a main-effects model. Choosing a shrinkage factor is a well-defined, smooth, and often a unique process as opposed to binary decisions on whether variables are “in” or “out” of the model. Likewise, instead of using arbitrary subjective (residual plots) or objective (χ^2 due to `cohort` \times covariable interactions, i.e., non-constant covariable effects), shrinkage can systematically allow model enhancements insofar as the information content in the data will support, through the use of differential penalization. Shrinkage is a solution to the dilemma faced when the analyst attempts to choose between a parsimonious model and a more complex one that fits the data. Penalization does not require the analyst to make a binary decision, and it is a process that can be validated using the bootstrap.

13.4.8 Validating the Fitted Model

Validation of statistical indexes such as D_{xy} and model calibration is done using techniques discussed previously, except that certain problems must be addressed. First, when using the bootstrap, the resampling must take into account the existence of multiple records per subject that were created to use

the binary logistic likelihood trick. That is, sampling should be done with replacement from *subjects* rather than *records*. Second, the analyst must isolate which event to predict. This is because when observations are expanded in order to use a binary logistic likelihood function to fit the CR model, several different events are being predicted simultaneously. Somers' D_{xy} could be computed by relating $X\hat{\gamma}$ (ignoring intercepts) to the ordinal Y , but other indexes are not defined so easily. The simplest approach here would be to validate a single prediction for $\Pr(Y = j|Y \geq j, X)$, for example. The simplest event to predict is $\Pr(Y = 0|X)$, as this would just require subsetting on all observations in the first cohort level in the validation sample. It would also be easy to validate any one of the later conditional probabilities. The validation functions described in the next section allow for such subsetting, as well as handling the cluster sampling. Specialized calculations would be needed to validate an unconditional probability such as $\Pr(Y \geq 2|X)$.

13.4.9 R Functions

The `cr.setup` function in `rms` returns a list of vectors useful in constructing a dataset used to trick a binary logistic function such as `lrm` into fitting CR models. The `subs` vector in this list contains observation numbers in the original data, some of which are repeated. Here is an example.

```
u <- cr.setup(Y)           # Y=original ordinal response
attach(mydata[u$subs,])   # mydata is the original dataset
                           # mydata[i,] subscripts input data,
                           # using duplicate values of i for
                           # repeats
y       <- u$y             # constructed binary responses
cohort  <- u$cohort       # cohort or risk set categories
f <- lrm(y ~ cohort*age + sex)
```

Since the `lrm` and `pentrace` functions have the capability to penalize different parts of the model by different amounts, they are valuable for fitting extended CR models in which the `cohort` \times predictor interactions are allowed to be only as important as the information content in the data will support. Simple main effects can be unpenalized or slightly penalized as desired.

The `validate` and `calibrate` functions for `lrm` allow specification of subject identifiers when using the bootstrap, so the samples can be constructed with replacement from the original subjects. In other words, cluster sampling is done from the expanded records. This is handled internally by the `predab.resample` function. These functions also allow one to specify a subset of the records to use in the validation, which makes it especially easy to validate the part of the model used to predict $\Pr(Y = 0|X)$.

The `plot.xmean.ordinaly` function is useful for checking the CR assumption for single predictors, as described earlier.

13.5 Further Reading

- ① See^{5, 25, 26, 31, 32, 52, 63, 64, 113, 126, 240, 245, 276, 354, 449, 502, 561, 664, 679} for some excellent background references, applications, and extensions to the ordinal models.⁶⁶³ and⁴²⁸ demonstrate how to model ordinal outcomes with repeated measurements within subject using random effects in Bayesian models. The first to develop an ordinal regression model were Aitchison and Silvey⁸.
- ② Some analysts feel that combining categories improves the performance of test statistics when fitting PO models when sample sizes are small and cells are sparse. Murad et al.⁴⁶⁹ demonstrated that this causes more problems, because it results in overly conservative Wald tests.
- ③ Anderson and Philips [26, p. 29] proposed methods for constructing properly spaced response values given a fitted PO model.
- ④ The simplest demonstration of this is to consider a model in which there is a single predictor that is totally independent of a nine-level response Y , so PO *must* hold. A PO model is fitted in SAS using:

```
DATA test;
DO i=1 to 50;
  y=FLOOR(RANUNI(151)*9);
  x=RANNOR(5);
  OUTPUT;
  END;
PROC LOGISTIC; MODEL y=x;
```

The score test for PO was $\chi^2 = 56$ on 7 d.f., $P < 0.0001$. This problem results from some small cell sizes in the distribution of Y .⁵⁰² The P -value for testing the regression effect for X was 0.76.

- ⑤ The R `glmnetcr` package by Kellie Archer provides a different way to fit continuation ratio models.
- ⑥ Bender and Benner⁴⁸ have some examples using the precursor of the `rms` package for fitting and assessing the goodness of fit of ordinal logistic regression models.

13.6 Problems

Test for the association between disease group and total hospital cost in SUPPORT, without imputing any missing costs (exclude the one patient having zero cost).

1. Use the Kruskal–Wallis rank test.
2. Use the proportional odds ordinal logistic model generalization of the Wilcoxon–Mann–Whitney Kruskal–Wallis Spearman test. Group total cost into 20 quantile groups so that only 19 intercepts will need to be in the model, not one less than the number of subjects (this would have taken the program too long to fit the model). Use the likelihood ratio χ^2 for this and later steps.
3. Use a binary logistic model to test for association between disease group and whether total cost exceeds the median of total cost. In other words, group total cost into two quantile groups and use this binary variable as the response. What is wrong with this approach?

4. Instead of using only two cost groups, group cost into 3, 4, 5, 6, 8, 10, and 12 quantile groups. Describe the relationship between the number of intervals used to approximate the continuous response variable and the efficiency of the analysis. How many intervals of total cost, assuming that the ordering of the different intervals is used in the analysis, are required to avoid losing significant information in this continuous variable?
5. If you were selecting one of the rank-based tests for testing the association between disease and cost, which of any of the tests considered would you choose?
6. Why do all of the tests you did have the same number of degrees of freedom for the hypothesis of no association between `dzgroup` and `totcst`?
7. What is the advantage of a rank-based test over a parametric test based on $\log(\text{cost})$?
8. Show that for a two-sample problem, the numerator of the score test for comparing the two groups using a proportional odds model is exactly the numerator of the Wilcoxon-Mann-Whitney two-sample rank-sum test.

Chapter 14

Case Study in Ordinal Regression, Data Reduction, and Penalization

This case study is taken from Harrell et al.²⁷² which described a World Health Organization study⁴³⁹ in which vital signs and a large number of clinical signs and symptoms were used to develop a predictive model for an ordinal response. This response consists of laboratory assessments of diagnosis and severity of illness related to pneumonia, meningitis, and sepsis. Much of the modeling strategy given in Chapter 4 was used to develop the model, with additional emphasis on penalized maximum likelihood estimation (Section 9.10). The following laboratory data are used in the response: cerebrospinal fluid (CSF) culture from a lumbar puncture (LP), blood culture (BC), arterial oxygen saturation (SaO_2 , a measure of lung dysfunction), and chest X-ray (CXR). The sample consisted of 4552 infants aged 90 days or less.

This case study covers these topics:

1. definition of the ordinal response (Section 14.1);
2. scoring and clustering of clinical signs (Section 14.2);
3. testing adequacy of weights specified by subject-matter specialists and assessing the utility of various scoring schemes using a tentative ordinal logistic model (Section 14.3);
4. assessing the basic ordinality assumptions and examining the proportional odds and continuation ratio (PO and CR) assumptions separately for each predictor (Section 14.4);
5. deriving a tentative PO model using cluster scores and regression splines (Section 14.5);
6. using residual plots to check PO, CR, and linearity assumptions (Section 14.6);
7. examining the fit of a CR model (Section 14.7);
8. utilizing an extended CR model to allow some or all of the regression coefficients to vary with cutoffs of the response level as well as to provide formal tests of constant slopes (Section 14.8);

Table 14.1 Ordinal Outcome Scale

Outcome Level Y	Definition	n	Fraction in Outcome Level		
			BC, CXR Indicated ($n = 2398$)	Not Indicated ($n = 1979$)	Random Sample ($n = 175$)
0	None of the below	3551	0.63	0.96	0.91
1	$90\% \leq SaO_2 < 95\%$ or CXR+	490	0.17	0.04 ^a	0.05
2	BC+ or CSF+ or $SaO_2 < 90\%$	511	0.21	0.00 ^b	0.03

^a SaO_2 was measured but CXR was not done

^b Assumed zero since neither BC nor LP were done.

9. using penalized maximum likelihood estimation to improve accuracy (Section 14.9);
10. approximating the full model by a sub-model and drawing a nomogram on the basis of the sub-model (Section 14.10); and
11. validating the ordinal model using the bootstrap (Section 14.11).

14.1 Response Variable

To be a candidate for BC and CXR, an infant had to have a clinical indication for one of the three diseases, according to prespecified criteria in the study protocol ($n = 2398$). Blood work-up (but not necessarily LP) and CXR was also done on a random sample intended to be 10% of infants having no signs or symptoms suggestive of infection ($n = 175$). Infants with signs suggestive of meningitis had LP done. All 4552 infants received a full physical exam and standardized pulse oximetry to measure SaO_2 . The vast majority of infants getting CXR had the X-rays interpreted by three independent radiologists.

The analyses that follow are not corrected for verification bias⁶⁸⁷ with respect to BC, LP, and CXR, but Section 14.1 has some data describing the extent of the problem, and the problem is reduced by conditioning on a large number of covariates.

Patients were assigned to the worst qualifying outcome category. Table 14.1 shows the definition of the ordinal outcome variable Y and shows the distribution of Y by the lab work-up strategy.

The effect of verification bias is a false negative fraction of 0.03 for $Y = 2$, from comparing the detection fraction of zero for $Y = 2$ in the “Not Indicated” group with the observed positive fraction of 0.03 in the random sample that was fully worked up. The extent of verification bias in $Y = 1$ is $0.05 - 0.04 = 0.01$. These biases are ignored in this analysis.

14.2 Variable Clustering

Forty-seven clinical signs were collected for each infant. Most questionnaire items were scored as a single variable using equally spaced codes, with 0 to 3 representing, for example, sign not present, mild, moderate, severe. The resulting list of clinical signs with their abbreviations is given in Table 14.2. The signs are organized into clusters as discussed later.

Table 14.2 Clinical Signs

Cluster Name	Sign Abbreviation	Name of Sign	Values
bul.conv	abb	bulging fontanel	0-1
	convul	hx convulsion	0-1
hydration	abk	sunken fontanel	0-1
	hdi	hx diarrhoea	0-1
	deh	dehydrated	0-2
	stu	skin turgor	0-2
	dcp	digital capillary refill	0-2
drowsy	hcl	less activity	0-1
	qcr	quality of crying	0-2
	csd	drowsy state	0-2
	slpm	sleeping more	0-1
	wake	wakes less easily	0-1
	aro	arousal	0-2
	mvm	amount of movement	0-2
agitated	hcm	crying more	0-1
	slpl	sleeping less	0-1
	con	consolability	0-2
	csa	agitated state	0-1
crying	hcm	crying more	0-1
	hcs	crying less	0-1
	qcr	quality of crying	0-2
	smi2	smiling ability \times age > 42 days	0-2
reffort	nfl	nasal flaring	0-3
	lcw	lower chest in-drawing	0-3
	gru	grunting	0-2
	ccy	central cyanosis	0-1
stop.breath	hap	hx stop breathing	0-1
	apn	apnea	0-1
ausc	whz	wheezing	0-1
	coh	cough heard	0-1
	crs	crepitation	0-2
hxprob	hfb	fast breathing	0-1
	hdb	difficulty breathing	0-1
	hlt	mother report resp. problems	none, chest, other
feeding	hfa	hx abnormal feeding	0-3
	absu	sucking ability	0-2
	afe	drinking ability	0-2
labor	chi	previous child died	0-1
	fde	fever at delivery	0-1
	ldy	days in labor	1-9
	twb	water broke	0-1
abdominal	adb	abdominal distension	0-4
	jau	jaundice	0-1
	omph	omphalitis	0-1
fever.ill	illd	age-adjusted no. days ill	
	hfe	hx fever	0-1
pustular	conj	conjunctivitis	0-1
	oto	otoscopy impression	0-2
	puskin	pustular skin rash	0-1

Table 14.3 Clinician Combinations, Rankings, and Scorings of Signs

Cluster	Combined/Ranked Signs in Order of Severity	Weights
bul.conv	abb \cup convul	0-1
drowsy	hcl, qcr>0, csd>0 \cup slpm \cup wake, aro>0, mvm>0	0-5
agitated	hcm, slpl, con=1, csa, con=2	0, 1, 2, 7, 8, 10
reffort	nfl>0, lcw>1, gru=1, gru=2, ccy	0-5
ausc	whz, coh, crs>0	0-3
feeding	hfa=1, hfa=2, hfa=3, absu=1 \cup afe=1, absu=2 \cup afe=2	0-5
abdominal	jau \cup abd>0 \cup omph	0-1

for analyzing the principal components were to see if some of the clusters could be removed from consideration so that the clinicians would not spend time developing scoring rules for them. Let us “peek” at Y to assist in scoring clusters at this point, but to do so in a very structured way that does not involve the examination of a large number of individual coefficients.

To judge any cluster scoring scheme, we must pick a tentative outcome model. For this purpose we chose the PO model. By using the 14 PC_1 s corresponding to the 14 clusters, the fitted PO model had a likelihood ratio (LR) χ^2 of 1155 with 14 d.f., and the predictive discrimination of the clusters was quantified by a Somers’ D_{xy} rank correlation between $X\hat{\beta}$ and Y of 0.596. The following clusters were not statistically important predictors and we assumed that the lack of importance of the PC_1 s in predicting Y (adjusted for the other PC_1 s) justified a conclusion that no sign within that cluster was clinically important in predicting Y : hydration, hxprob, pustular, crying, fever.ill, stop.breath, labor. This list was identified using a backward step-down procedure on the full model. The total Wald χ^2 for these seven PC_1 s was 22.4 ($P = 0.002$). The reduced model had LR $\chi^2 = 1133$ with 7 d.f., $D_{xy} = 0.591$. The bootstrap validation in Section 14.11 penalizes for examining all candidate predictors.

The clinicians were asked to rank the clinical severity of signs within each potentially important cluster. During this step, the clinicians also ranked severity levels of some of the component signs, and some cluster scores were simplified, especially when the signs within a cluster occurred infrequently. The clinicians also assessed whether the severity points or weights should be equally spaced, assigning unequally spaced weights for one cluster (agitated). The resulting rankings and sign combinations are shown in Table 14.3. The signs or sign combinations separated by a comma are treated as separate categories, whereas some signs were unioned (“or”-ed) when the clinicians deemed them equally important. As an example, if an additive cluster score was to be used for drowsy, the scorings would be 0 = none present, 1 = hcl, 2 = qcr>0, 3 = csd>0 or slpm or wake, 4 = aro>0, 5 = mvm>0 and the scores would be added.

Table 14.4 Predictive information of various cluster scoring strategies. AIC is on the likelihood ratio χ^2 scale.

Scoring Method	LR χ^2	d.f.	AIC
PC_1 of each cluster	1133	7	1119
Union of all signs	1045	7	1031
Union of higher categories	1123	7	1109
Hierarchical (worst sign)	1194	7	1180
Additive, equal weights	1155	7	1141
Additive using clinician weights	1183	7	1169
Hierarchical, data-driven weights	1227	25	1177

This table reflects some data reduction already (unioning some signs and selection of levels of ordinal signs) but more reduction is needed. Even after signs are ranked within a cluster, there are various ways of assigning the cluster scores. We investigated six methods. We started with the purely statistical approach of using PC_1 to summarize each cluster. Second, all sign combinations within a cluster were unioned to represent a 0/1 cluster score. Third, only sign combinations thought by the clinicians to be severe were unioned, resulting in $\text{drowsy}=\text{aro}>0$ or $\text{mvm}>0$, $\text{agitated}=\text{csa}$ or $\text{con}=2$, $\text{reffort}=\text{lcw}>1$ or $\text{gru}>0$ or ccy , $\text{ausc}=\text{crs}>0$, and $\text{feeding}=\text{absu}>0$ or $\text{afe}>0$. For clusters that are not scored 0/1 in Table 14.3, the fourth summarization method was a hierarchical one that used the weight of the worst applicable category as the cluster score. For example, if $\text{aro}=1$ but $\text{mvm}=0$, drowsy would be scored as 4. The fifth method counted the number of positive signs in the cluster. The sixth method summed the weights of all signs or sign combinations present. Finally, the worst sign combination present was again used as in the second method, but the points assigned to the category were data-driven ones obtained by using extra dummy variables. This provided an assessment of the adequacy of the clinician-specified weights. By comparing rows 4 and 7 in Table 14.4 we see that response data-driven sign weights have a slightly worse AIC, indicating that the number of extra β parameters estimated was not justified by the improvement in χ^2 . The hierarchical method, using the clinicians' weights, performed quite well. The only cluster with inadequate clinician weights was ausc —see below. The PC_1 method, without any guidance, performed well, as in²⁶⁸. The only reasons not to use it are that it requires a coefficient for every sign in the cluster and the coefficients are not translatable into simple scores such as 0, 1, . . .

Representation of clusters by a simple union of selected signs or of all signs is inadequate, but otherwise the choice of methods is not very important in terms of explaining variation in Y . We chose the fourth method, a hierarchical severity point assignment (using weights that were prespecified by the clinicians), for its ease of use and of handling missing component variables (in most cases) and potential for speeding up the clinical exam (examining

to detect more important signs first). Because of what was learned regarding the relationship between `ausc` and Y , we modified the `ausc` cluster score by redefining it as `ausc=crs>0` (crepitations present). Note that neither the “tweaking” of `ausc` nor the examination of the seven scoring methods displayed in Table 14.4 is taken into account in the model validation.

14.4 Assessing Ordinality of Y for each X , and Unadjusted Checking of PO and CR Assumptions

Section 13.2 described a graphical method for assessing the ordinality assumption for Y separately with respect to each X , and for assessing PO and CR assumptions individually. Figure 14.2 is an example of such displays. For this dataset we expect strongly nonlinear effects for `temp`, `rr`, and `hrrat`, so for those predictors we plot the mean absolute differences from suitable “normal” values as an approximate solution.

```
Sc <- transform(Sc,
                ausc = 1 * (ausc == 3),
                bul.conv = 1 * (bul.conv == 'TRUE'),
                abdominal = 1 * (abdominal == 'TRUE'))
plot.xmean.ordinality(Y ~ age + abs(temp-37) + abs(rr-60) +
                      abs(hrrat-125) + waz + bul.conv + drowsy +
                      agitated + reffort + ausc + feeding +
                      abdominal, data=Sc, cr=TRUE,
                      subn=FALSE, cex.points=.65) # Figure 14.2
```

The plot is shown in Figure 14.2. Y does not seem to operate in an ordinal fashion with respect to `age`, `|rr-60|`, or `ausc`. For the other variables, ordinality holds, and PO holds reasonably well for the other variables. For heart rate, the PO assumption appears to be satisfied perfectly. CR model assumptions appear to be more tenuous than PO assumptions, when one variable at a time is fitted.

14.5 A Tentative Full Proportional Odds Model

Based on what was determined in Section 14.3, the original list of 47 signs was reduced to seven predictors: two unions of signs (`bul.conv`, `abdominal`), one single sign (`ausc`), and four “worst category” point assignments (`drowsy`, `agitated`, `reffort`, `feeding`). Seven clusters were dropped for the time being because of weak associations with Y . Such a limited use of variable selection reduces the severe problems inherent with that technique.

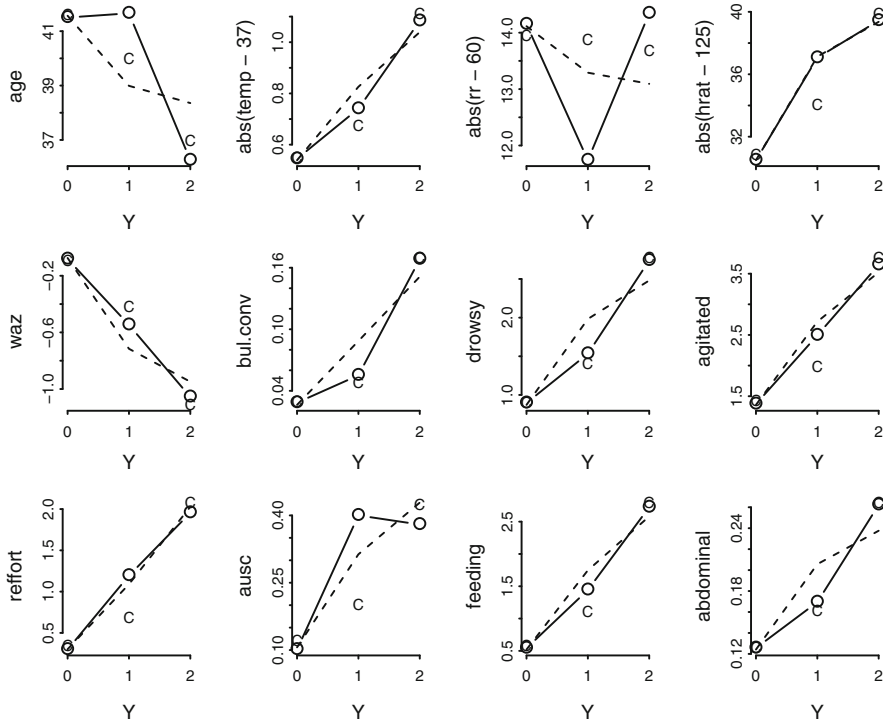


Fig. 14.2 Examination of the ordinality of Y for each predictor by assessing how varying Y relate to the mean X , and whether the trend is monotonic. Solid lines connect the simple stratified means, and dashed lines connect the estimated expected value of $X|Y = j$ given that PO holds. Estimated expected values from the CR model are marked with Cs.

At this point in model development add to the model `age` and vital signs: `temp` (temperature), `rr` (respiratory rate), `hrat` (heart rate), and `waz`, weight-for-age Z -score. Since `age` was expected to modify the interpretation of `temp`, `rr`, and `hrat`, and interactions between continuous variables would be difficult to use in the field, we categorized `age` into three intervals: 0–6 days ($n = 302$), 7–59 days ($n = 3042$), and 60–90 days ($n = 1208$).^a

```
Sc$ageg ← cut2(Sc$age, c(7, 60))
```

The new variables `temp`, `rr`, `hrat`, `waz` were missing in, respectively, $n = 13$, 11 , 147 , and 20 infants. Since the three vital sign variables are somewhat correlated with each other, customized single imputation models were developed to impute all the missing values without assuming linearity or even monotonicity of any of the regressions.

^a These age intervals were also found to adequately capture most of the interaction effects.

```
vsign.trans ← transcan(~ temp + hrat + rr, data=Sc,
                       imputed=TRUE, pl=FALSE)
```

```
Convergence criterion:2.222 0.643 0.191 0.056 0.016
Convergence in 6 iterations
R2 achieved in predicting each variable:
```

```
temp hrat rr
0.168 0.160 0.066
```

```
Adjusted R2:
```

```
temp hrat rr
0.167 0.159 0.064
```

```
Sc ← transform(Sc,
               temp = impute(vsign.trans, temp),
               hrat = impute(vsign.trans, hrat),
               rr = impute(vsign.trans, rr))
```

After `transcan` estimated optimal restricted cubic spline transformations, `temp` could be predicted with adjusted $R^2 = 0.17$ from `hrat` and `rr`, `hrat` could be predicted with adjusted $R^2 = 0.16$ from `temp` and `rr`, and `rr` could be predicted with adjusted R^2 of only 0.06. The first two R^2 , while not large, mean that customized imputations are more efficient than imputing with constants. Imputations on `rr` were closer to the median `rr` of 48/minute as compared with the other two vital signs whose imputations have more variation. In a similar manner, `waz` was imputed using `age`, birth weight, head circumference, body length, and prematurity (adjusted R^2 for predicting `waz` from the others was 0.74). The continuous predictors `temp`, `hrat`, `rr` were not assumed to linearly relate to the log odds that $Y \geq j$. Restricted cubic spline functions with five knots for `temp`, `rr` and four knots for `hrat`, `waz` were used to model the effects of these variables:

```
f1 ← lrm(Y ~ ageg*(rcs(temp,5)+rcs(rr,5)+rcs(hrat,4)) +
         rcs(waz,4) + bul.conv + drowsy + agitated +
         reffort + ausc + feeding + abdominal,
         data=Sc, x=TRUE, y=TRUE)
# x=TRUE, y=TRUE used by resid() below
print(f1, latex=TRUE, coefs=5)
```

Logistic Regression Model

```
lrm(formula = Y ~ ageg * (rcs(temp, 5) + rcs(rr, 5) + rcs(hrat,
4)) + rcs(waz, 4) + bul.conv + drowsy + agitated + reffort +
    ausc + feeding + abdominal, data = Sc, x = TRUE, y = TRUE)
```

		Model Likelihood Ratio Test		Discrimination Indexes		Rank Discrim. Indexes	
Obs	4552	LR χ^2	1393.18	R^2	0.355	C	0.826
0	3551	d.f.	45	g	1.485	D_{xy}	0.653
1	490	$\Pr(> \chi^2) < 0.0001$		g_r	4.414	γ	0.654
2	511			g_p	0.225	τ_a	0.240
$\max \left \frac{\partial \log L}{\partial \beta} \right $		2×10^{-6}		Brier	0.120		

	Coef	S.E.	Wald Z	$\Pr(> Z)$
$y \geq 1$	0.0653	7.6563	0.01	0.9932
$y \geq 2$	-1.0646	7.6563	-0.14	0.8894
ageg=[7,60)	9.5590	9.9071	0.96	0.3346
ageg=[60,90]	29.1376	15.8915	1.83	0.0667
temp	-0.0694	0.2160	-0.32	0.7480
...				

Wald tests of nonlinearity and interaction are shown in Table 14.5.

```

latex(anova(f1), file='', label='ordinal-anova.f1',
caption='Wald statistics from the proportional odds model',
size='smaller') # Table 14.5
    
```

The bottom four lines of the table are the most important. First, there is strong evidence that some associations with Y exist (45 d.f. test) and very strong evidence of nonlinearity in one of the vital signs or in `waz` (26 d.f. test). There is moderately strong evidence for an interaction effect somewhere in the model (22 d.f. test). We see that the grouped age variable `ageg` is predictive of Y , but mainly as an effect modifier for `rr`, and `hrat`. `temp` is extremely nonlinear, and `rr` is moderately so. `hrat`, a difficult variable to measure reliably in young infants, is perhaps not important enough ($\chi^2 = 19,9$ d.f.) to keep in the final model.

14.6 Residual Plots

Section 13.3.4 defined binary logistic score residuals for isolating the PO assumption in an ordinal model. For the tentative PO model, score residuals for four of the variables were plotted using

```

resid(f1, 'score.binary', pl=TRUE, which=c(17,18,20,21))
## Figure 14.3
    
```

The result is shown in Figure 14.3. We see strong evidence of non-PO for `ausc` and moderate evidence for `drowsy` and `bul.conv`, in agreement with Figure 14.2.

Table 14.5 Wald statistics from the proportional odds model

	χ^2	d.f.	<i>P</i>
ageg (Factor+Higher Order Factors)	41.49	24	0.0147
<i>All Interactions</i>	40.48	22	0.0095
temp (Factor+Higher Order Factors)	37.08	12	0.0002
<i>All Interactions</i>	6.77	8	0.5617
<i>Nonlinear (Factor+Higher Order Factors)</i>	31.08	9	0.0003
rr (Factor+Higher Order Factors)	81.16	12	< 0.0001
<i>All Interactions</i>	27.37	8	0.0006
<i>Nonlinear (Factor+Higher Order Factors)</i>	27.36	9	0.0012
hrat (Factor+Higher Order Factors)	19.00	9	0.0252
<i>All Interactions</i>	8.83	6	0.1836
<i>Nonlinear (Factor+Higher Order Factors)</i>	7.35	6	0.2901
waz	35.82	3	< 0.0001
<i>Nonlinear</i>	13.21	2	0.0014
bul.conv	12.16	1	0.0005
drowsy	17.79	1	< 0.0001
agitated	8.25	1	0.0041
reffort	63.39	1	< 0.0001
ausc	105.82	1	< 0.0001
feeding	30.38	1	< 0.0001
abdominal	0.74	1	0.3895
ageg × temp (Factor+Higher Order Factors)	6.77	8	0.5617
<i>Nonlinear</i>	6.40	6	0.3801
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	6.40	6	0.3801
ageg × rr (Factor+Higher Order Factors)	27.37	8	0.0006
<i>Nonlinear</i>	14.85	6	0.0214
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	14.85	6	0.0214
ageg × hrat (Factor+Higher Order Factors)	8.83	6	0.1836
<i>Nonlinear</i>	2.42	4	0.6587
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	2.42	4	0.6587
TOTAL NONLINEAR	78.20	26	< 0.0001
TOTAL INTERACTION	40.48	22	0.0095
TOTAL NONLINEAR + INTERACTION	96.31	32	< 0.0001
TOTAL	1073.78	45	< 0.0001

Partial residuals computed separately for each *Y*-cutoff (Section 13.3.4) are the most useful residuals for ordinal models as they simultaneously check linearity, find needed transformations, and check PO. In Figure 14.4, smoothed partial residual plots were obtained for all predictors, after first fitting a simple model in which every predictor was assumed to operate linearly. Interactions were temporarily ignored and *age* was used as a continuous variable.

```
f2 ← lrm(Y ~ age + temp + rr + hrat + waz +
        bul.conv + drowsy + agitated + reffort + ausc +
        feeding + abdominal, data=Sc, x=TRUE, y=TRUE)
resid(f2, 'partial', pl=TRUE, label.curves=FALSE) # Figure 14.4
```

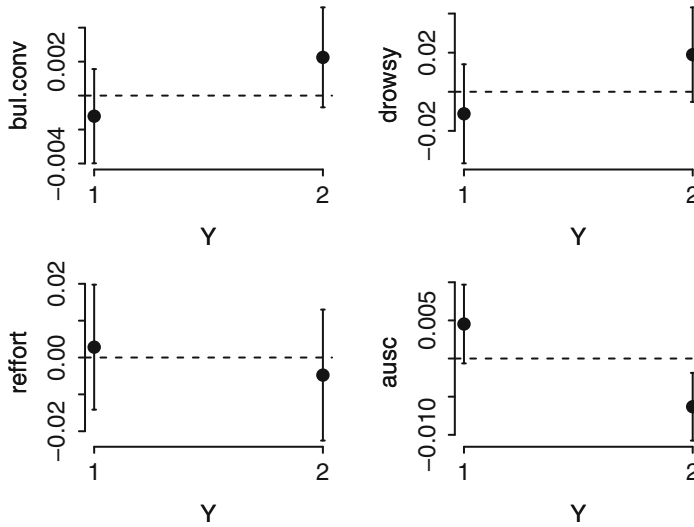


Fig. 14.3 Binary logistic model score residuals for binary events derived from two cutoffs of the ordinal response Y . Note that the mean residuals, marked with closed circles, correspond closely to differences between solid and dashed lines at $Y = 1, 2$ in Figure 14.2. Score residual assessments for spline-expanded variables such as `rr` would have required one plot per d.f.

The degree of non-parallelism generally agreed with the degree of non-flatness in Figure 14.3 and with the other score residual plots that were not shown. The partial residuals show that `temp` is highly nonlinear and that it is much more useful in predicting $Y = 2$. For the cluster scores, the linearity assumption appears reasonable, except possibly for `drowsy`. Other nonlinear effects are taken into account using splines as before (except for `age`, which is categorized).

A model can have significant lack of fit with respect to some of the predictors and still yield quite accurate predictions. To see if that is the case for this PO model, we computed predicted probabilities of $Y = 2$ for all infants from the model and compared these with predictions from a customized binary logistic model derived to predict $\Pr(Y = 2)$. The mean absolute difference in predicted probabilities between the two models is only 0.02, but the 0.90 quantile of that difference is 0.059. For high-risk infants, discrepancies of 0.2 were common. Therefore we elected to consider a different model.

14.7 Graphical Assessment of Fit of CR Model

In order to take a first look at the fit of a CR model, let us consider the two binary events that need to be predicted, and assess linearity and paral-

lelism over Y -cutoffs. Here we fit a sequence of binary fits and then use the `plot.lrm.partial` function, which assembles partial residuals for a sequence of fits and constructs one graph per predictor.

```
cr0 <- lrm(Y==0 ~ age + temp + rr + hrat + waz +
          bul.conv + drowsy + agitated + reffort + ausc +
          feeding + abdominal, data=Sc, x=TRUE, y=TRUE)
# Use the update function to save repeating model right-
# hand side. An indicator variable for Y=1 is the
# response variable below
cr1 <- update(cr0, Y==1 ~ ., subset=Y ≥ 1)
plot.lrm.partial(cr0, cr1, center=TRUE) # Figure 14.5
```

The output is in Figure 14.5. There is not much more parallelism here than in Figure 14.4. For the two most important predictors, `ausc` and `rr`, there are strongly differing effects for the different events being predicted (e.g., $Y = 0$ or $Y = 1|Y \geq 1$). As is often the case, there is no one constant β model that satisfies assumptions with respect to all predictors simultaneously, especially

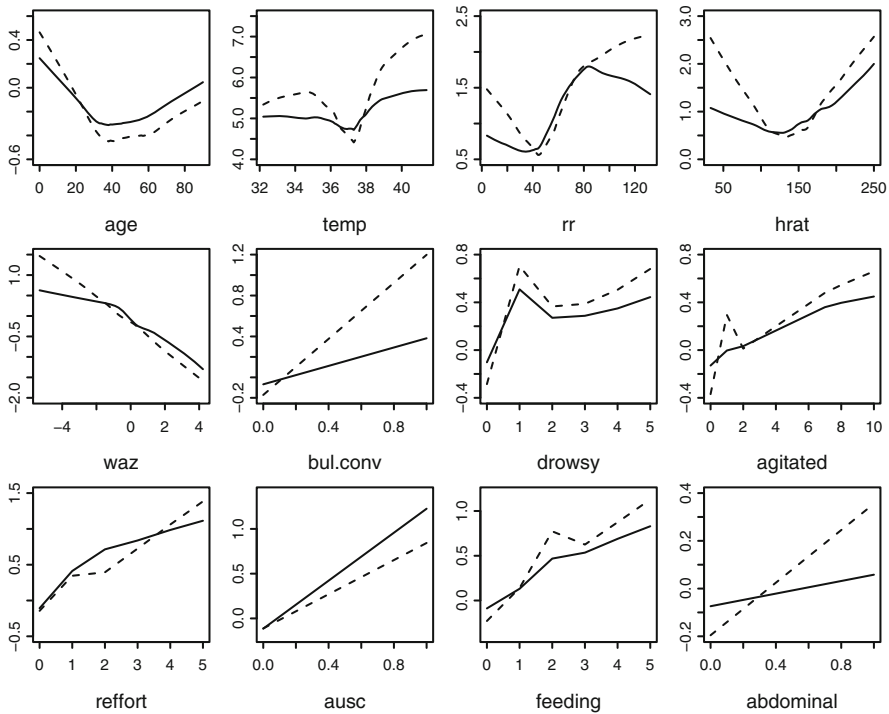


Fig. 14.4 Smoothed partial residuals corresponding to two cutoffs of Y , from a model in which all predictors were assumed to operate linearly and additively. The smoothed curves estimate the actual predictor transformations needed, and parallelism relates to the PO assumption. Solid lines denote $Y \geq 1$ while dashed lines denote $Y \geq 2$.

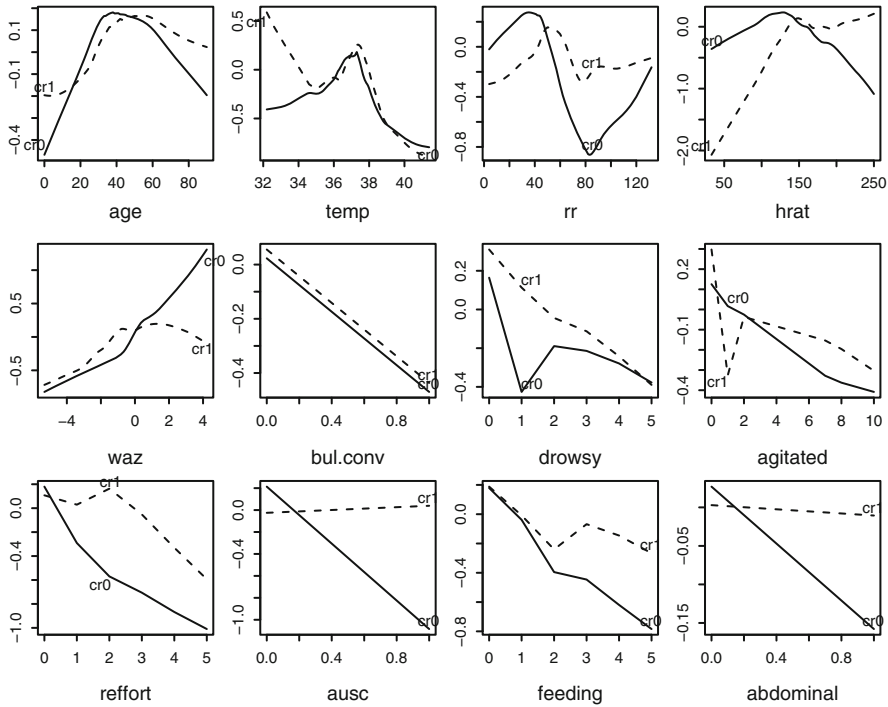


Fig. 14.5 loess smoothed partial residual plots for binary models that are components of an ordinal continuation ratio model. Solid lines correspond to a model for $Y = 0$, and dotted lines correspond to a model for $Y = 1|Y \geq 1$.

when there is evidence for non-ordinality for `ausc` in Figure 14.2. The CR model will need to be generalized to adequately fit this dataset.

14.8 Extended Continuation Ratio Model

The CR model in its ordinary form has no advantage over the PO model for this dataset. But Section 13.4.6 discussed how the CR model can easily be extended to relax any of its assumptions. First we use the `cr.setup` function to set up the data for fitting a CR model using the binary logistic trick.

```
u <- cr.setup(Y)
Sc.expanded <- Sc[u$subs, ]
y <- u$y
cohort <- u$cohort
```

Here the `cohort` variable has values 'all', 'Y>=1' corresponding to the conditioning events in Equation 13.10. Once the data frame is expanded to include the different risk cohorts, vectors such as `age` are lengthened (to 5553 records). Now we fit a fully extended CR model that makes no equal slopes assumptions; that is, the model *has* to fit Y assuming the covariables are linear and additive. At this point, we omit `hrat` but add back all variables that were deleted by examining their association with Y . Recall that most of these seven cluster scores were summarized using PC_1 . Adding back “insignificant” variables will allow us to validate the model fairly using the bootstrap, as well as to obtain confidence intervals that are not falsely narrow.¹⁶

```
full <-
  lrm(y ~ cohort*(ageg*(rcs(temp,5) + rcs(rr,5)) +
    rcs(waz,4) + bul.conv + drowsy + agitated + reffort +
    ausc + feeding + abdominal + hydration + hxprob +
    pustular + crying + fever.ill + stop.breath + labor),
    data=Sc.expanded, x=TRUE, y=TRUE)
# x=TRUE, y=TRUE are for pentrace, validate, calibrate below
perf <- function(fit) { # model performance for Y=0
  pr <- predict(fit, type='fitted')[cohort == 'all']
  s <- round(somers2(pr, y[cohort == 'all']), 3)
  pr <- 1 - pr # Predict Prob[Y > 0] instead of Prob[Y = 0]
  f <- round(c(mean(pr < .05), mean(pr > .25),
    mean(pr > .5)), 2)
  f <- paste(f[1], ', ', f[2], ', and ', f[3], '.', sep='')
  list(somers=s, fractions=f)
}
perf.unpen <- perf(full)
print(full, latex=TRUE, coefs=5)
```

Logistic Regression Model

```
lrm(formula = y ~ cohort * (ageg * (rcs(temp, 5) +
  rcs(rr, 5)) + rcs(waz, 4) + bul.conv + drowsy +
  agitated + reffort + ausc + feeding + abdominal +
  hydration + hxprob + pustular + crying + fever.ill +
  stop.breath + labor), data = Sc.expanded, x = TRUE,
  y = TRUE)
```

		Model Likelihood Ratio Test	Discrimination Indexes	Rank Discrim. Indexes
Obs	5553	LR χ^2 1824.33	R^2 0.406	C 0.843
0	1512	d.f. 87	g 1.677	D_{xy} 0.685
1	4041	$\Pr(> \chi^2) < 0.0001$	g_r 5.350	γ 0.687
max $ \frac{\partial \log L}{\partial \beta} 8 \times 10^{-7}$			g_p 0.269	τ_a 0.272
			Brier 0.135	

Table 14.6 Wald statistics for `cohort` in the CR model

	χ^2	d.f.	P
cohort (Factor+Higher Order Factors)	199.47	44	< 0.0001
<i>All Interactions</i>	172.12	43	< 0.0001
TOTAL	199.47	44	< 0.0001

	Coef	S.E.	Wald Z	Pr(> Z)
Intercept	1.3966	9.0827	0.15	0.8778
cohort= $Y \geq 1$	1.5077	14.6443	0.10	0.9180
ageg=[7,60)	-9.3715	11.4104	-0.82	0.4115
ageg=[60,90]	-26.4502	17.2188	-1.54	0.1245
temp	-0.0049	0.2551	-0.02	0.9846
...				

```
latex(anova(full, cohort), file='', # Table 14.6
      caption='Wald statistics for \\co{cohort} in the CR model',
      size='smaller[2]', label='ordinal-anova.cohort')
```

```
an ← anova(full, india=FALSE, indnl=FALSE)
```

```
latex(an, file='', label='ordinal-anova.full',
      caption='Wald statistics for the continuation ratio model.
      Interactions with \\co{cohort} assess non-proportional
      hazards', caption.lot='Wald statistics for $Y$ in the
      continuation ratio model',
      size='smaller[2]') # Table 14.7
```

This model has LR $\chi^2 = 1824$ with 87 d.f. Wald statistics are in Tables 14.6 and 14.7. The global test of the constant slopes assumption in the CR model (test of all interactions involving `cohort`) has Wald $\chi^2 = 172$ with 43 d.f., $P < 0.0001$. Consistent with Figure 14.5, the formal tests indicate that `ausc` is the biggest violator, followed by `waz` and `rr`.

14.9 Penalized Estimation

We know that the CR model must be extended to fit these data adequately. If the model is fully extended to allow for all `cohort` \times predictor interactions, we have not gained any precision or power in using an ordinal model over using a polytomous logistic model. Therefore we seek some restrictions on the model's parameters. The `lrm` and `pentrace` functions allow for differing λ for shrinking different types of terms in the model. Here we do a grid search to determine the optimum penalty for simple main effect (non-interaction) terms and the penalty for interaction terms, most of which are terms interacting with `cohort`

Table 14.7 Wald statistics for the continuation ratio model. Interactions with cohort assess non-proportional hazards

	χ^2	d.f.	P
cohort	199.47	44	< 0.0001
ageg	48.89	36	0.0742
temp	59.37	24	0.0001
rr	93.77	24	< 0.0001
waz	39.69	6	< 0.0001
bul.conv	10.80	2	0.0045
drowsy	15.19	2	0.0005
agitated	13.55	2	0.0011
reffort	51.85	2	< 0.0001
ausc	109.80	2	< 0.0001
feeding	27.47	2	< 0.0001
abdominal	1.78	2	0.4106
hydration	4.47	2	0.1069
hxprob	6.62	2	0.0364
pustular	3.03	2	0.2194
crying	1.55	2	0.4604
fever.ill	3.63	2	0.1630
stop.breath	5.34	2	0.0693
labor	5.35	2	0.0690
ageg \times temp	8.18	16	0.9432
ageg \times rr	38.11	16	0.0015
cohort \times ageg	14.88	18	0.6701
cohort \times temp	8.77	12	0.7225
cohort \times rr	19.67	12	0.0736
cohort \times waz	9.04	3	0.0288
cohort \times bul.conv	0.33	1	0.5658
cohort \times drowsy	0.57	1	0.4489
cohort \times agitated	0.55	1	0.4593
cohort \times reffort	2.29	1	0.1298
cohort \times ausc	38.11	1	< 0.0001
cohort \times feeding	2.48	1	0.1152
cohort \times abdominal	0.09	1	0.7696
cohort \times hydration	0.53	1	0.4682
cohort \times hxprob	2.54	1	0.1109
cohort \times pustular	2.40	1	0.1210
cohort \times crying	0.39	1	0.5310
cohort \times fever.ill	3.17	1	0.0749
cohort \times stop.breath	2.99	1	0.0839
cohort \times labor	0.05	1	0.8309
cohort \times ageg \times temp	2.22	8	0.9736
cohort \times ageg \times rr	10.22	8	0.2500
TOTAL NONLINEAR	93.36	40	< 0.0001
TOTAL INTERACTION	203.10	59	< 0.0001
TOTAL NONLINEAR + INTERACTION	257.70	67	< 0.0001
TOTAL	1211.73	87	< 0.0001

to allow for unequal slopes. The following code uses `pentrace` on the full extended CR model fit to find the optimum penalty factors. All combinations of the `simple` and `interaction` λ s for which the interaction penalty \geq the penalty for the simple parameters are examined.

```
d ← options(digits=4)
pentrace(full,
         list(simple=c(0, .025, .05, .075, .1),
              interaction=c(0, 10, 50, 100, 125, 150)))
```

Best penalty:

simple	interaction	df				
0.05	125	49.75				
simple	interaction	df	aic	bic	aic.c	
0.000	0	87.00	1650	1074	1648	
0.000	10	60.63	1671	1269	1669	
0.025	10	60.11	1672	1274	1670	
0.050	10	59.80	1672	1276	1670	
0.075	10	59.58	1671	1277	1670	
0.100	10	59.42	1671	1278	1670	
0.000	50	54.64	1671	1309	1670	
0.025	50	54.14	1672	1313	1671	
0.050	50	53.83	1672	1316	1671	
0.075	50	53.62	1672	1317	1671	
0.100	50	53.46	1672	1318	1671	
0.000	100	51.61	1672	1330	1671	
0.025	100	51.11	1673	1334	1672	
0.050	100	50.81	1673	1336	1672	
0.075	100	50.60	1672	1337	1671	
0.100	100	50.44	1672	1338	1671	
0.000	125	50.55	1672	1337	1671	
0.025	125	50.05	1673	1341	1672	
0.050	125	49.75	1673	1343	1672	
0.075	125	49.54	1672	1344	1672	
0.100	125	49.39	1672	1345	1671	
0.000	150	49.65	1672	1343	1671	
0.025	150	49.15	1672	1347	1672	
0.050	150	48.85	1673	1349	1672	
0.075	150	48.64	1672	1350	1671	
0.100	150	48.49	1672	1351	1671	

```
options(d)
```

We see that shrinkage from 87 d.f. down to 49.75 effective d.f. results in an improvement in χ^2 -scaled AIC of 23. The optimum penalty factors were 0.05 for simple terms and 125 for interaction terms.

Let us now store a penalized version of the full fit, find where the effective d.f. were reduced, and compute χ^2 for each factor in the model. We take the effective d.f. for a collection of model parameters to be the sum of the

diagonals of the matrix product defined underneath Gray's Equation 2.9²³⁷ that correspond to those parameters.

```
full.pen ←
  update(full,
    penalty=list(simple=.05, interaction=125))
print(full.pen, latex=TRUE, coefs=FALSE)
```

Logistic Regression Model

```
lrm(formula = y ~ cohort * (age * (rcs(temp, 5) + rcs(rr, 5)) +
  rcs(waz, 4) + bul.conv + drowsy + agitated + reffort + ausc +
  feeding + abdominal + hydration + hxprob + pustular + crying +
  fever.ill + stop.breath + labor), data = Sc.expanded, x = TRUE,
  y = TRUE, penalty = list(simple = 0.05, interaction = 125))
```

Penalty factors

```
simple nonlinear interaction nonlinear.interaction
0.05      0.05      125      125
```

		Model Likelihood	Discrimination	Rank Discrim.
		Ratio Test	Indexes	Indexes
Obs	5553	LR χ^2 1772.11	R^2 0.392	C 0.840
0	1512	d.f. 49.75	g 1.594	D_{xy} 0.679
1	4041	$\Pr(> \chi^2) < 0.0001$	g_r 4.924	γ 0.681
$\max \left \frac{\partial \log L}{\partial \beta} \right $	1×10^{-7}	Penalty 21.48	g_p 0.263	τ_a 0.269
			Brier 0.136	

```
effective.df(full.pen)
```

Original and Effective Degrees of Freedom

	Original	Penalized
All	87	49.75
Simple Terms	20	19.98
Interaction or Nonlinear	67	29.77
Nonlinear	40	16.82
Interaction	59	22.57
Nonlinear Interaction	32	9.62

```
## Compute discrimination for Y=0 vs. Y>0
perf.pen ← perf(full.pen) # Figure 14.6
# Exclude interactions and cohort effects from plot
plot(anova(full.pen), cex.labels=0.75, rm.ia=TRUE,
  rm.other='cohort (Factor+Higher Order Factors)')
```

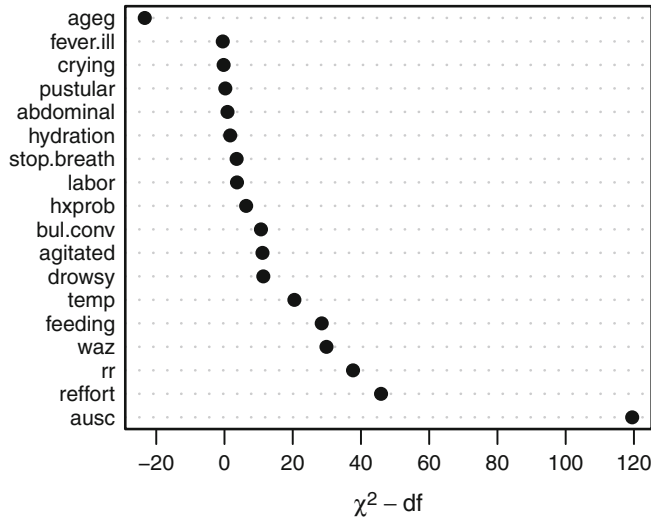


Fig. 14.6 Importance of predictors in full penalized model, as judged by partial Wald χ^2 minus the predictor d.f. The Wald χ^2 values for each line in the dot plot include contributions from all higher-order effects. Interaction effects by themselves have been removed as has the `cohort` effect.

This will be the final model except for the model used in Section 14.10. The model has LR $\chi^2 = 1772$. The output of `effective.df` shows that non-interaction terms have barely been penalized, and coefficients of interaction terms have been shrunk from 59 d.f. to effectively 22.6 d.f. Predictive discrimination was assessed by computing the Somers' D_{xy} rank correlation between $X\hat{\beta}$ and whether $Y = 0$, in the subset of records for which $Y = 0$ is what was being predicted. Here $D_{xy} = 0.672$, and the ROC area is 0.838 (the unpenalized model had an apparent $D_{xy} = 0.676$). To summarize in another way the effectiveness of this model in screening infants for risks of any abnormality, the fraction of infants with predicted probabilities that $Y > 0$ being < 0.05 , > 0.25 , and > 0.5 are, respectively, 0.1, 0.28, and 0.14. `anova` output is plotted in Figure 14.6 to give a snapshot of the importance of the various predictors. The Wald statistics used here are computed on a variance-covariance matrix which is adjusted for penalization (using Gray Equation 2.6²³⁷ before it was determined that the sandwich covariance estimator performs less well than the inverse of the penalized information matrix—see p. 211).

The full equation for the fitted model is below. Only the part of the equation used for predicting $\Pr(Y = 0)$ is shown, other than an intercept for $Y \geq 1$ that does not apply when $Y = 0$.

```
latex(full.pen, which=1:21, file='')
```

$$\begin{aligned}
X\hat{\beta} = & \\
& -1.337435[Y \geq 1] \\
& +0.1074525[\text{ageg} \in [7, 60]] + 0.1971287[\text{ageg} \in [60, 90]] \\
& +0.1978706\text{temp} + 0.1091831(\text{temp} - 36.19998)_+^3 - 2.833442(\text{temp} - 37)_+^3 \\
& +5.07114(\text{temp} - 37.29999)_+^3 - 2.507527(\text{temp} - 37.69998)_+^3 \\
& +0.1606456(\text{temp} - 39)_+^3 \\
& +0.02090741\text{rr} - 6.336873 \times 10^{-5}(\text{rr} - 32)_+^3 + 8.405441 \times 10^{-5}(\text{rr} - 42)_+^3 \\
& +6.152416 \times 10^{-5}(\text{rr} - 49)_+^3 - 0.0001018105(\text{rr} - 59)_+^3 + 1.960063 \times 10^{-5}(\text{rr} - 76)_+^3 \\
& -0.07589699\text{waz} + 0.02508918(\text{waz} + 2.9)_+^3 - 0.1185068(\text{waz} + 0.75)_+^3 \\
& +0.1225752(\text{waz} - 0.28)_+^3 - 0.02915754(\text{waz} - 1.73)_+^3 - 0.4418073 \text{bul.conv} \\
& -0.08185088 \text{drowsy} - 0.05327209 \text{agitated} - 0.2304409 \text{reffort} \\
& -1.158604 \text{ausc} - 0.1599588 \text{feeding} - 0.1608684 \text{abdominal} \\
& -0.05409718 \text{hydration} + 0.08086387 \text{hxprob} + 0.007519746 \text{pustular} \\
& +0.04712091 \text{crying} + 0.004298725 \text{fever.ill} - 0.3519033 \text{stop.breath} \\
& +0.06863879 \text{labor} \\
& +[\text{ageg} \in [7, 60]][6.499592 \times 10^{-5} \text{temp} - 0.00279976(\text{temp} - 36.19998)_+^3 \\
& -0.008691166(\text{temp} - 37)_+^3 - 0.004987871(\text{temp} - 37.29999)_+^3 \\
& +0.0259236(\text{temp} - 37.69998)_+^3 - 0.009444801(\text{temp} - 39)_+^3] \\
& +[\text{ageg} \in [60, 90]][0.0001320368\text{temp} - 0.00182639(\text{temp} - 36.19998)_+^3 \\
& -0.01640406(\text{temp} - 37)_+^3 - 0.0476041(\text{temp} - 37.29999)_+^3 \\
& +0.09142148(\text{temp} - 37.69998)_+^3 - 0.02558693(\text{temp} - 39)_+^3] \\
& +[\text{ageg} \in [7, 60]][-0.0009437598\text{rr} - 1.044673 \times 10^{-6}(\text{rr} - 32)_+^3 \\
& -1.670499 \times 10^{-6}(\text{rr} - 42)_+^3 - 5.189082 \times 10^{-6}(\text{rr} - 49)_+^3 + 1.428634 \times 10^{-5}(\text{rr} - 59)_+^3 \\
& -6.382087 \times 10^{-6}(\text{rr} - 76)_+^3] \\
& +[\text{ageg} \in [60, 90]][-0.001920811\text{rr} - 5.52134 \times 10^{-6}(\text{rr} - 32)_+^3 \\
& -8.628392 \times 10^{-6}(\text{rr} - 42)_+^3 - 4.147347 \times 10^{-6}(\text{rr} - 49)_+^3 + 3.813427 \times 10^{-5}(\text{rr} - 59)_+^3 \\
& -1.98372 \times 10^{-5}(\text{rr} - 76)_+^3]
\end{aligned}$$

where $[c] = 1$ if subject is in group c , 0 otherwise; $(x)_+ = x$ if $x > 0$, 0 otherwise.

Now consider displays of the shapes of effects of the predictors. For the continuous variables `temp` and `rr` that interact with age group, we show the effects for all three age groups separately for each Y cutoff. All effects have been centered so that the log odds at the median predictor value is zero when `cohort='all'`, so these plots actually show log odds relative to reference values. The patterns in Figures 14.9 and 14.8 are in agreement with those in Figure 14.5.


```

yl ← c(-3, 1)      # put all plots on common y-axis scale

# Plot predictors that interact with another predictor
# Vary ageg over all age groups, then vary temp over its
# default range (10th smallest to 10th largest values in
# data). Make a separate plot for each 'cohort'
# ref.zero centers effects using median x

dd ← datadist(Sc.expanded); dd ← datadist(dd, cohort)
options(datadist='dd')

p1 ← Predict(full.pen, temp, ageg, cohort,
             ref.zero=TRUE, conf.int=FALSE)
p2 ← Predict(full.pen, rr, ageg, cohort,
             ref.zero=TRUE, conf.int=FALSE)
p ← rbind(temp=p1, rr=p2) # Figure 14.7:
source(paste('http://biostat.mc.vanderbilt.edu/wiki/pub/Main',
             'RConfiguration/graphicsSet.r', sep='/'))
ggplot(p, ~ cohort, groups='ageg', varypred=TRUE,
       ylim=yl, layout=c(2, 1), legend.position=c(.85,.8),
       addlayer=ltheme(width=3, height=3, text=2.5, title=2.5),
       adj.subtitle=FALSE) # ltheme defined with source()

```

```

# For each predictor that only interacts with cohort, show
# the differing effects of the predictor for predicting
# Pr(Y=0) and Pr(Y=1 given Y exceeds 0) on the same graph

dd$limits['Adjust to','cohort'] ← 'Y ≥ 1'
v ← Cs(waz, bul.conv, drowsy, agitated, reffort, ausc,
       feeding, abdominal, hydration, hxprob, pustular,
       crying)
yeq1 ← Predict(full.pen, name=v, ref.zero=TRUE)
yl ← c(-1.5, 1.5)
ggplot(yeq1, ylim=yl, sepdiscrete='vertical') # Figure 14.8

```

```

dd$limits['Adjust to','cohort'] ← 'all' # original default
all ← Predict(full.pen, name=v, ref.zero=TRUE)
ggplot(all, ylim=yl, sepdiscrete='vertical') # Figure 14.9

```

1

14.10 Using Approximations to Simplify the Model

Parsimonious models can be developed by approximating predictions from the model to any desired level of accuracy. Let $\hat{L} = X\hat{\beta}$ denote the predicted log odds from the full penalized ordinal model, including multiple records for subjects with $Y > 0$. Then we can use a variety of techniques to approximate \hat{L} from a subset of the predictors (in their raw form). With this approach one can immediately see what is lost over the full model by computing, for

example, the mean absolute error in predicting \hat{L} . Another advantage to full model approximation is that shrinkage used in computing \hat{L} is inherited by any model that predicts \hat{L} . In contrast, the usual stepwise methods result in $\hat{\beta}$ that are too large since the final coefficients are estimated as if the model structure were prespecified.

2

CART would be particularly useful as a model approximator as it would result in a prediction tree that would be easy for health workers to use.

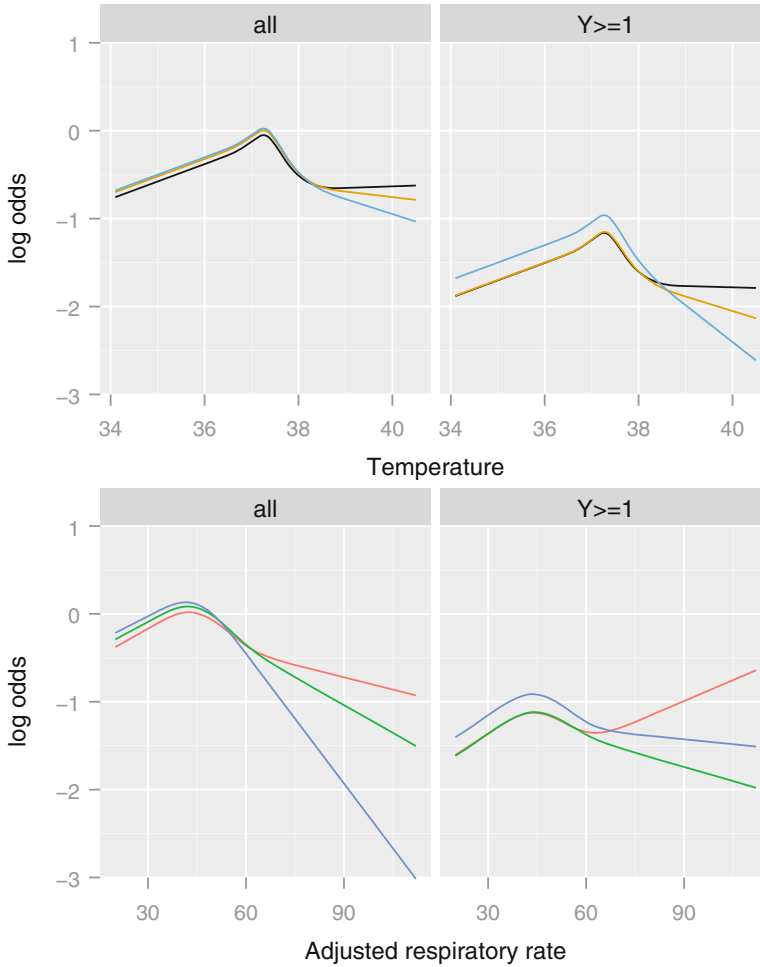


Fig. 14.7 Centered effects of predictors on the log odds, showing the effects of two predictors with interaction effects for the age intervals noted. The title **all** refers to the prediction of $Y = 0|Y \geq 0$, that is, $Y = 0$. **Y>=1** refers to predicting the probability of $Y = 1|Y \geq 1$.

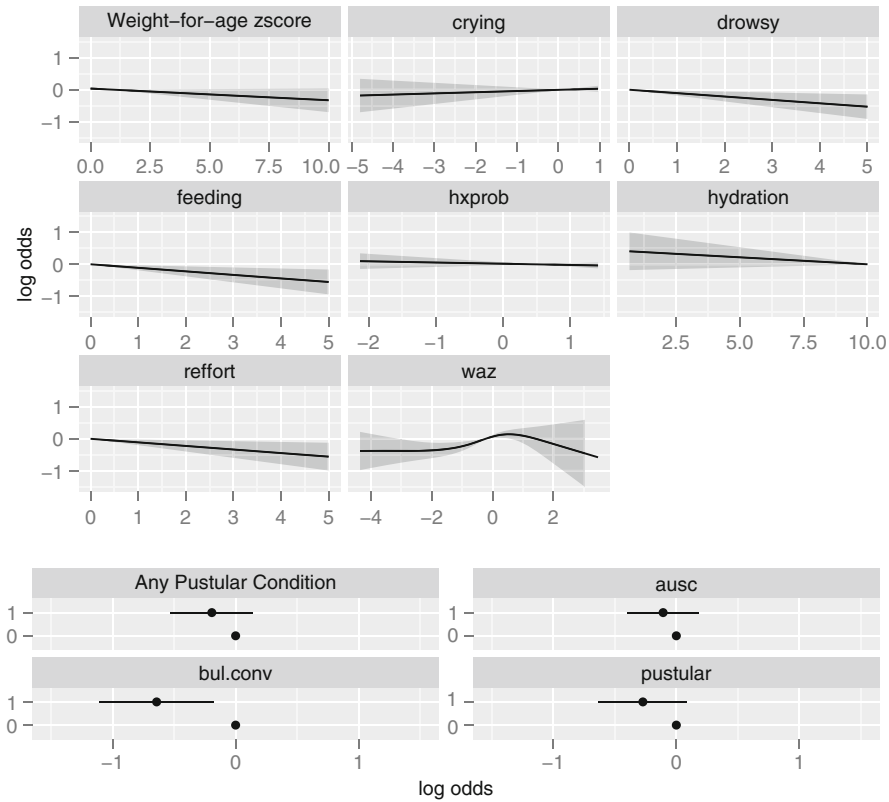


Fig. 14.8 Centered effects of predictors on the log odds, for predicting $Y = 1|Y \geq 1$

Unfortunately, a 50-node CART was required to predict \hat{L} with an $R^2 \geq 0.9$, and the mean absolute error in the predicted logit was still 0.4. This will happen when the model contains many important continuous variables.

Let's approximate the full model using its important components, by using a step-down technique predicting \hat{L} from all of the component variables using ordinary least squares. In using step-down with the least squares function `ols` in `rms` there is a problem when the initial $R^2 = 1.0$ as in that case the estimate of $\sigma = 0$. This can be circumvented by specifying an arbitrary nonzero value of σ to `ols` (here 1.0), as we are not using the variance-covariance matrix from `ols` anyway. Since `cohort` interacts with the predictors, separate approximations can be developed for each level of Y . For this example we approximate the log odds that $Y = 0$ using the cohort of patients used for determining $Y = 0$, that is, $Y \geq 0$ or `cohort='all'`.

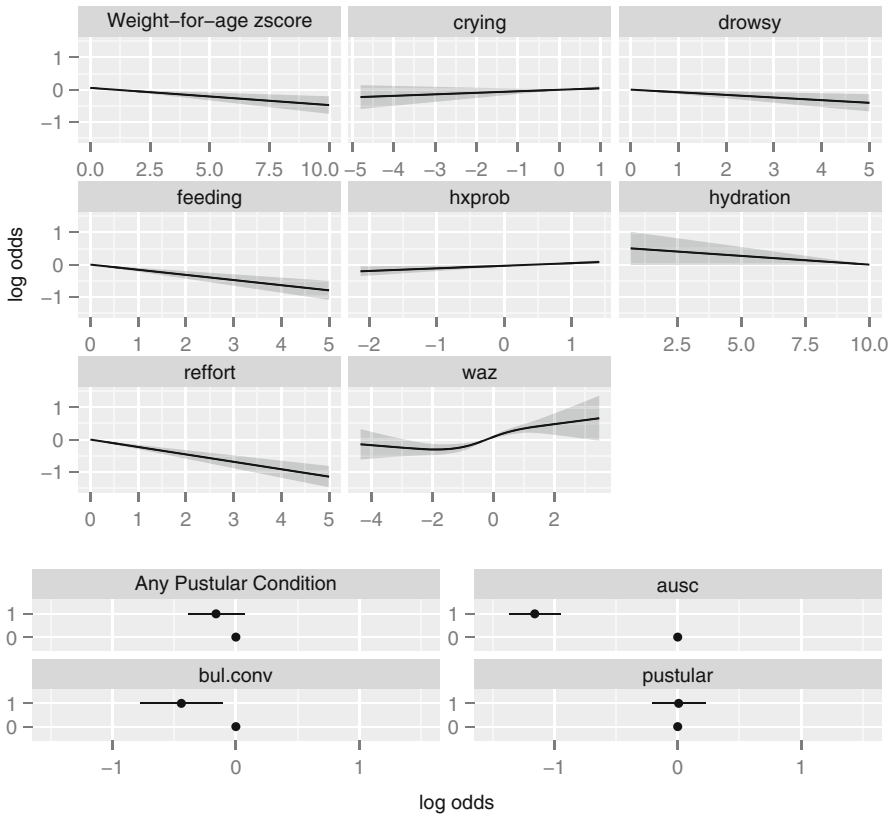


Fig. 14.9 Centered effects of predictors on the log odds, for predicting $Y \geq 1$. No plot was made for the fever.ill, stop.breath. or labor cluster scores.

```

plogit ← predict(full.pen)
f ← ols(plogit ~ ageg*(rcs(temp,5) + rcs(rr,5)) +
        rcs(waz,4) + bul.conv + drowsy + agitated +
        reffort + ausc + feeding + abdominal + hydration +
        hxprob + pustular + crying + fever.ill +
        stop.breath + labor,
        subset=cohort=='all', data=Sc.expanded, sigma=1)

# Do fast backward stepdown
w ← options(width=120)
fastbw(f, aics=1e10)

```

Deleted	Chi-Sq	d.f.	P	Residual	d.f.	P	AIC	R2
ageg * temp	1.87	8	0.9848	1.87	8	0.9848	-14.13	1.000
ageg	0.05	2	0.9740	1.92	10	0.9969	-18.08	1.000
pustular	0.02	1	0.8778	1.94	11	0.9987	-20.06	1.000
fever.ill	0.08	1	0.7828	2.02	12	0.9994	-21.98	1.000
crying	9.47	1	0.0021	11.49	13	0.5698	-14.51	0.999
abdominal	12.66	1	0.0004	24.15	14	0.0440	-3.85	0.997
rr	17.90	4	0.0013	42.05	18	0.0011	6.05	0.995
hydration	13.21	1	0.0003	55.26	19	0.0000	17.26	0.993
labor	23.48	1	0.0000	78.74	20	0.0000	38.74	0.990
stop.breath	33.40	1	0.0000	112.14	21	0.0000	70.14	0.986
bul.conv	51.53	1	0.0000	163.67	22	0.0000	119.67	0.980
agitated	63.66	1	0.0000	227.33	23	0.0000	181.33	0.972
hxprob	84.16	1	0.0000	311.49	24	0.0000	263.49	0.962
drowsy	109.86	1	0.0000	421.35	25	0.0000	371.35	0.948
temp	295.67	4	0.0000	717.01	29	0.0000	659.01	0.911
waz	368.86	3	0.0000	1085.87	32	0.0000	1021.87	0.866
reffort	449.83	1	0.0000	1535.70	33	0.0000	1469.70	0.810
ageg * rr	751.19	8	0.0000	2286.90	41	0.0000	2204.90	0.717
ausc	1906.82	1	0.0000	4193.72	42	0.0000	4109.72	0.482
feeding	3900.33	1	0.0000	8094.04	43	0.0000	8008.04	0.000

Approximate Estimates after Deleting Factors

	Coef	S.E.	Wald	Z	P
[1,]	1.617	0.01482	109.1	0	

Factors in Final Model

None

```
options(w)
# 1e10 causes all variables to eventually be
# deleted so can see most important ones in order

# Fit an approximation to the full penalized model using
# most important variables
full.approx <-
  ols(plogit ~ rcs(temp,5) + ageg*rcs(rr,5) +
      rcs(waz,4) + bul.conv + drowsy + reffort +
      ausc + feeding,
      subset=cohort=='all', data=Sc.expanded)
p <- predict(full.approx)
abserr <- mean(abs(p - plogit[cohort == 'all']))
Dxy <- somers2(p, y[cohort == 'all'])['Dxy']
```

The approximate model had R^2 against the full penalized model of 0.972, and the mean absolute error in predicting \hat{L}_{xy} was 0.17. The D_{xy} rank correlation between the approximate model's predicted logit and the binary event $Y = 0$

is 0.665 as compared with the full model's $D_{xy} = 0.672$. See Section 19.5 for an example of computing correct estimates of variance of the parameters in an approximate model.

Next turn to diagramming this model approximation so that all predicted values can be computed without the use of a computer. We draw a type of nomogram that converts each effect in the model to a 0 to 100 scale which is just proportional to the log odds. These points are added across predictors to derive the “Total Points,” which are converted to \hat{L} and then to predicted probabilities. For the interaction between `rr` and `ageg`, `rms`'s `nomogram` function automatically constructs three `rr` axes—only one is added into the total point score for a given subject. Here we draw a nomogram for predicting the probability that $Y > 0$, which is $1 - \Pr(Y = 0)$. This probability is derived by negating $\hat{\beta}$ and $X\hat{\beta}$ in the model derived to predict $\Pr(Y = 0)$.

```
f <- full.approx
f$coefficients <- -f$coefficients
f$linear.predictors <- -f$linear.predictors

n <- nomogram(f,
              temp=32:41, rr=seq(20,120,by=10),
              waz=seq(-1.5,2,by=.5),
              fun=plogis, funlabel='Pr(Y>0)',
              fun.at=c(.02,.05,seq(.1,.9,by=.1),.95,.98))
# Print n to see point tables
plot(n, lmgp=.2, cex.axis=.6) # Figure 14.10
newsobject <-
  data.frame(ageg='[ 0, 7)', rr=30, temp=39, waz=0, drowsy=5,
             reffort=2, bul.conv=0, ausc=0, feeding=0)
xb <- predict(f, newsobject)
```

The nomogram is shown in Figure 14.10. As an example in using the nomogram, a six-day-old infant gets approximately 9 points for having a respiration rate of 30/minute, 19 points for having a temperature of 39°C, 11 points for `waz=0`, 14 points for `drowsy=5`, and 15 points for `reffort=2`. Assuming that `bul.conv=ausc=feeding=0`, that infant gets 68 total points. This corresponds to $X\hat{\beta} = -0.68$ and a probability of 0.34.

3

14.11 Validating the Model

For the full CR model that was fitted using penalized maximum likelihood estimation (PMLE), we used 200 bootstrap replications to estimate and then to correct for optimism in various statistical indexes: D_{xy} , generalized R^2 , intercept and slope of a linear re-calibration equation for $X\hat{\beta}$, the maximum calibration error for $\Pr(Y = 0)$ based on the linear-logistic re-calibration (`Emax`), and the Brier quadratic probability score `B`. PMLE is used at each of the 200 resamples. During the bootstrap simulations, we sample with

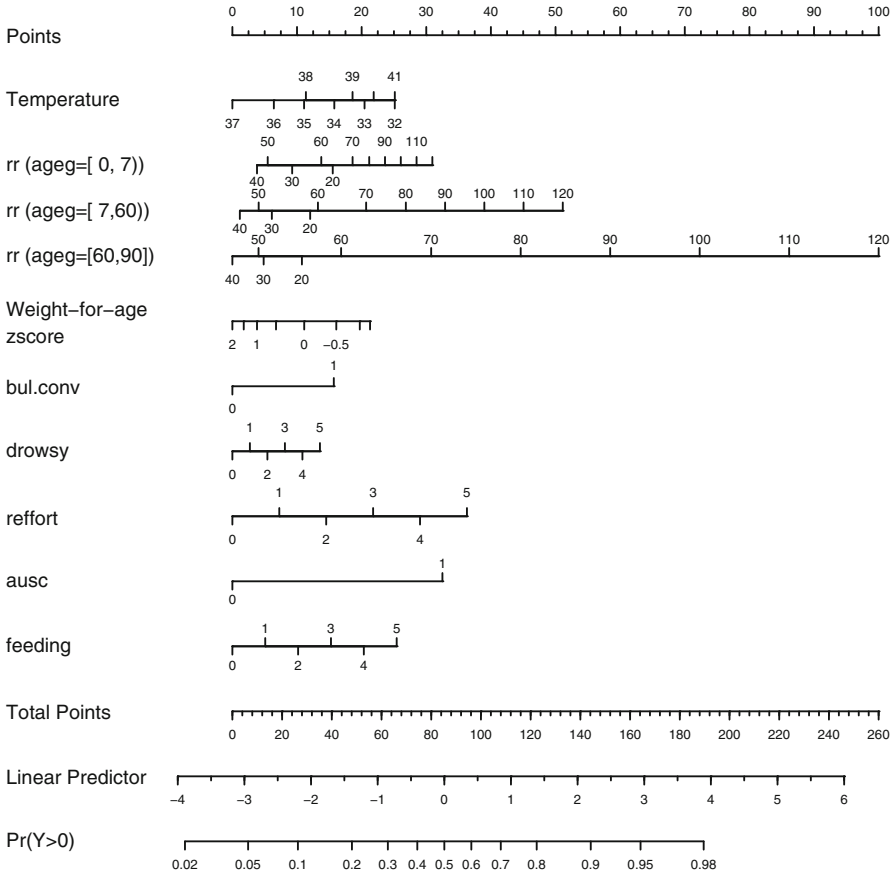


Fig. 14.10 Nomogram for predicting $\Pr(Y > 0)$ from the penalized extended CR model, using an approximate model fitted using ordinary least squares ($R^2 = 0.972$ against the full model's predicted logits).

replacement from the *patients* and not from the 5553 expanded *records*, hence the specification `cluster=u$subs`, where `u$subs` is the vector of sequential patient numbers computed from `cr.setup` above. To be able to assess predictive accuracy of a single predicted probability, the `subset` parameter is specified so that $\Pr(Y = 0)$ is being assessed even though 5553 observations are used to develop each of the 200 models.

```
set.seed(1) # so can reproduce results
v ← validate(full.pen, B=200, cluster=u$subs,
             subset=cohort=='all')
latex(v, file='', digits=2, size='smaller')
```

Index	Original Sample	Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.67	0.68	0.67	0.01	0.66	200
R^2	0.38	0.38	0.37	0.01	0.36	200
Intercept	-0.03	-0.03	0.00	-0.03	0.00	200
Slope	1.03	1.03	1.00	0.03	1.00	200
E_{\max}	0.00	0.00	0.00	0.00	0.00	200
D	0.28	0.29	0.28	0.01	0.27	200
U	0.00	0.00	0.00	0.00	0.00	200
Q	0.28	0.29	0.28	0.01	0.27	200
B	0.12	0.12	0.12	0.00	0.12	200
g	1.47	1.50	1.45	0.04	1.42	200
g_p	0.22	0.23	0.22	0.00	0.22	200

```
v ← round(v, 3)
```

We see that for the apparent $D_{xy} = 0.672$ and that the optimism from overfitting was estimated to be 0.011 for the PMLE model, so the bias-corrected estimate of predictive discrimination is 0.661. The intercept and slope needed to re-calibrate $X\hat{\beta}$ to a 45° line are very near (0, 1). The estimate of the maximum calibration error in predicting $\Pr(Y = 0)$ is 0.001, which is quite satisfactory. The corrected Brier score is 0.122.

The simple calibration statistics just listed do not address the issue of whether predicted values from the model are miscalibrated in a nonlinear way, so now we estimate an overfitting-corrected calibration curve nonparametrically.

```
cal ← calibrate(full.pen, B=200, cluster=u$subs,
               subset=cohort=='all')
err ← plot(cal) # Figure 14.11
```

```
n=5553 Mean absolute error=0.017 Mean squared error=0.00043
0.9 Quantile of absolute error=0.038
```

The results are shown in Figure 14.11. One can see a slightly nonlinear calibration function estimate, but the overfitting-corrected calibration is excellent everywhere, being only slightly worse than the apparent calibration. The estimated maximum calibration error is 0.044. The excellent validation for both predictive discrimination and calibration are a result of the large sample size, frequency distribution of Y , initial data reduction, and PMLE.

14.12 Summary

Clinically guided variable clustering and item weighting resulted in a great reduction in the number of candidate predictor degrees of freedom and hence increased the true predictive accuracy of the model. Scores summarizing clusters of clinical signs, along with temperature, respiration rate, and weight-for-age after suitable nonlinear transformation and allowance for interactions

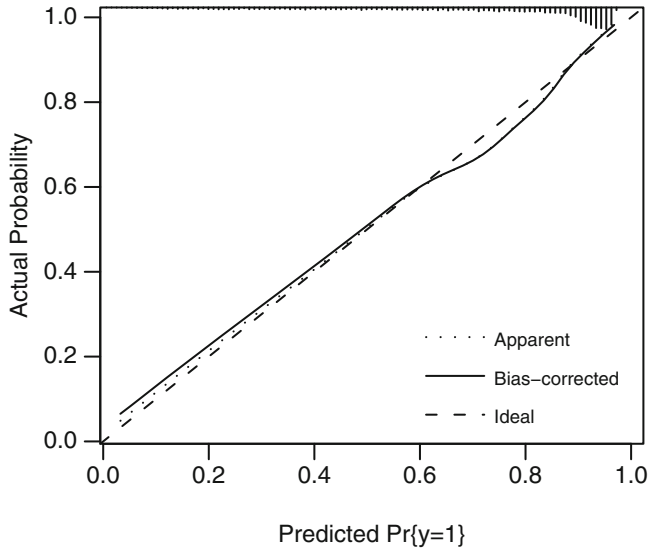


Fig. 14.11 Bootstrap calibration curve for the full penalized extended CR model. 200 bootstrap repetitions were used in conjunction with the `loess` smoother.¹¹¹ Also shown is a “rug plot” to demonstrate how effective this model is in discriminating patients into low- and high-risk groups for $\Pr(Y = 0)$ (which corresponds with the derived variable value $y = 1$ when `cohort='all'`).

with age, are powerful predictors of the ordinal response. Graphical methods are effective for detecting lack of fit in the PO and CR models and for diagramming the final model. Model approximation allowed development of parsimonious clinical prediction tools. Approximate models inherit the shrinkage from the full model. For the ordinal model developed here, substantial shrinkage of the full model was needed.

14.13 Further Reading

- ① See Moons et al.⁴⁶² for another case study in penalized maximum likelihood estimation.
- ② The *lasso* method of Tibshirani^{608,609} also incorporates shrinkage into variable selection.
- ③ To see how this compares with predictions using the full model, the extra clinical signs in that model that are not in the approximate model were predicted individually on the basis of $X\hat{\beta}$ from the reduced model along with the signs that are in that model, using ordinary linear regression. The signs not specified when evaluating the approximate model were then set to predicted values based on the values given for the 6-day-old infant above. The resulting $X\hat{\beta}$ for the full model is -0.81 and the predicted probability is 0.31 , as compared with -0.68 and 0.34 quoted above.

14.14 Problems

Develop a proportional odds ordinal logistic model predicting the severity of functional disability (`sfdm2`) in `SUPPORT`. The highest level of this variable corresponds to patients dying before the two-month follow-up interviews. Consider this level as the most severe outcome. Consider the following predictors: `age`, `sex`, `dzgroup`, `num.co`, `scoma`, `race` (use all levels), `meanbp`, `hrt`, `temp`, `pafi`, `alb`, `adlsc`. The last variable is the baseline level of functional disability from the “activities of daily living scale.”

1. For the variables `adlsc`, `sex`, `age`, `meanbp`, and others if you like, make plots of means of predictors stratified by levels of the response, to check for ordinality. On the same plot, show estimates of means assuming the proportional odds relationship between predictors and response holds. Comment on the evidence for ordinality and for proportional odds.
2. To allow for maximum adjustment of baseline functional status, treat this predictor as nominal (after rounding it to the nearest whole number; fractional values are the result of imputation) in remaining steps, so that all dummy variables will be generated. Make a single chart showing proportions of various outcomes stratified (individually) by `adlsc`, `sex`, `age`, `meanbp`. For continuous predictors use quartiles. You can pass the following function to the `summary` (`summary.formula`) function to obtain the proportions of patients having `sfdm2` at or worse than each of its possible levels (other than the first level). An easy way to do this is to use the `cumcategory` function with the `Hmisc` package’s `summary.formula` function. `cumcategorysummary.formula` Print estimates to only two significant digits of precision. Manually check the calculations for the `sex` variable using `table(sex, sfdm2)`. Then plot all estimates on a single graph using `plot(object, which=1:4)`, where `object` was created by `summary` (actually `summary.formula`). Note: for printing tables you may want to convert `sfdm2` to a 0–4 variable so that column headers are short and so that later calculations are simpler. You can use for example:

```
sfdm ← as.integer(sfdm2) - 1
```

3. Use an R function such as the following to compute the *logits* of the cumulative proportions.

```
sf ← function(y)
  c('Y ≥ 1' = qlogis(mean(y ≥ 1)),
    'Y ≥ 2' = qlogis(mean(y ≥ 2)),
    'Y ≥ 3' = qlogis(mean(y ≥ 3)),
    'Y ≥ 4' = qlogis(mean(y ≥ 4)))
```

As the $Y = 3$ category is rare, it may be even better to omit the $Y \geq 4$ column above, as was done in Section 13.3.9 and Figure 13.1. For each predictor pick two rows of the `summary` table having reasonable sample sizes, and take the difference between the two rows. Comment on the

validity of the proportional odds assumption by assessing how constant the row differences are across columns. Note: constant differences in log odds (logits) mean constant ratios of odds or constant relative effects of the predictor across outcome levels.

4. Make two plots nonparametrically relating `age` to all of the cumulative proportions or their logits. You can use commands such as the following (to use the R `Hmisc` package).

```
for(i in 1:4)
  plsmo(age, sfdm ≥ i, add=i>1,
        ylim=c(.2,.8), ylab='Proportion Y ≥ j')
for(i in 1:4)
  plsmo(age, sfdm ≥ i, add=i>1, fun=qlogis,
        ylim=qlogis(c(.2,.8)), ylab='logit')
```

Comment on the linearity of the `age` effect (which of the two plots do you use?) and on the proportional odds assumption for `age`, by assessing parallelism in the second plot.

5. Impute `race` using the most frequent category and `pafi` and `alb` using “normal” values.
6. Fit a model to predict the ordinal response using all predictors. For continuous ones assume a smooth relationship but allow it to be nonlinear. Quantify the ability of the model to discriminate patients in the five outcomes. Do an overall likelihood ratio test for whether any variables are associated with the level of functional disability.
7. Compute partial tests of association for each predictor and a test of nonlinearity for continuous ones. Compute a global test of nonlinearity. Graphically display the ranking of importance of the predictors.
8. Display the shape of how each predictor relates to the log odds of exceeding any level of `sfdm2` you choose, setting other predictors to typical values (one value per predictor). By default, `Predict` will make predictions for the second response category, which is a satisfactory choice here.
9. Use resampling to validate the Somers’ D_{xy} rank correlation between predicted logit and the ordinal outcome. Also validate the generalized R^2 , and slope shrinkage coefficient, all using a single R command. Comment on the quality (potential “export-ability”) of the model.

Chapter 15

Regression Models for Continuous Y and Case Study in Ordinal Regression

This chapter concerns univariate continuous Y . There are many multivariable models for predicting such response variables, such as

- linear models with assumed normal residuals, fitted with ordinary least squares
- generalized linear models and other parametric models based on special distributions such as the gamma
- generalized additive models (GAMs)²⁷⁷
- generalization of GAMs to also nonparametrically transform Y (see Chapter 16)
- quantile regression (see Section 15.2)
- other robust regression models that, like quantile regression, use an objective different from minimizing the sum of squared errors⁶³⁵
- semiparametric models based on the ranks of Y , such as the Cox proportional hazards model (Chapter 20) and the proportional odds ordinal logistic model (Chapters 13 and 14)
- cumulative probability models (often called *cumulative link models*) which are semiparametric models from a wider class of families than the logistic.

Semiparametric models that treat Y as ordinal but not interval-scaled have many advantages including robustness and freedom from all distributional assumptions for Y conditional on any given set of predictors. Advantages are demonstrated in a case study of a cumulative probability ordinal model. Some of the results are compared to quantile regression and OLS. Many of the methods used in the case study also apply to ordinary linear models.

15.1 The Linear Model

The most popular multivariable model for analyzing a univariate continuous Y is the linear model

$$E(Y|X) = X\beta, \quad (15.1)$$

where β is estimated using ordinary least squares, that is, by solving for $\hat{\beta}$ to minimize $\sum(Y_i - X\hat{\beta})^2$.

To compute P -values and confidence limits using parametric methods we would have to assume that $Y|X$ is normal with mean $X\beta$ and constant variance σ^{2a} . One could estimate conditional means of Y without any distributional assumptions, but least squares estimators are not robust to outliers or high-leverage points, and the model would be inaccurate in estimating conditional quantiles of $Y|X$ or $\text{Prob}[Y \geq c|X]$ unless normality of residuals holds. To be accurate in estimating all quantities, the linear model assumes that the Gaussian distribution of $Y|X_1$ is a simple shift from the distribution of $Y|X_2$.

15.2 Quantile Regression

Quantile regression^{355,357} is a different approach to modeling Y . It makes no distributional assumptions other than continuity of Y , while having all the usual right hand side assumptions. Quantile regression provides essentially the same estimates as sample quantiles if there is only an intercept or a categorical predictor in the model. Quantile regression is transformation invariant — pre-transforming Y is not important.

Quantile regression is a natural generalization of sample quantiles. Let $\rho_\tau(y) = y(\tau - [y < 0])$. The τ^{th} sample quantile is the minimizer q of $\sum_{i=1}^n \rho_\tau(y_i - q)$. For a conditional τ^{th} quantile of $Y|X$ the corresponding quantile regression estimator $\hat{\beta}_\tau$ minimizes $\sum_{i=1}^n \rho_\tau(Y_i - X\beta)$.

In non-large samples, quantile regression is not as efficient at estimating quantiles as is ordinary least squares at estimating the mean, if the latter's assumptions hold.

Koenker's `quantreg` package in \mathbb{R} ³⁵⁶ implements quantile regression, and the `rms` package's `Rq` function provides a front-end that gives rise to various graphics and inference tools.

Using quantile regression, we directly model the median as a function of covariates so that only the $X\beta$ structure need be correct. Other quantiles (e.g., 90th percentile) can be modeled but standard errors will be much larger as it is more difficult to precisely estimate outer quantiles.

^a The latter assumption may be dispensed with if we use a robust Huber–White or bootstrap covariance matrix estimate. Normality may sometimes be dispensed with by using bootstrap confidence intervals.

15.3 Ordinal Regression Models for Continuous Y

A different robust semiparametric regression approach than quantile regression is the cumulative probability ordinal model. Semiparametric models have several advantages over parametric models such as OLS. While quantile regression has no restriction in the parameters when modeling one quantile versus another^b, ordinal cumulative probability models assume a connection between distributions of Y for different X . Ordinal regression even makes one less assumption than quantile regression about the distribution of Y for a specific X : the distribution need not be continuous.

Applying an increasing 1–1 transformation to Y results in no change to regression coefficient estimates with ordinal regression^c. Regression coefficient estimates are completely robust to extreme Y values^d. Estimates of quantiles of Y from ordinal regression are exactly transformation-preserving, e.g., the estimate of the median of $\log Y$ is exactly the log of the estimate of the median Y .

For a general continuous distribution function $F(y)$, an ordinal regression model based on cumulative probabilities may be stated as follows^e. Let the ordered unique values of Y be denoted by y_1, y_2, \dots, y_k and let the intercepts associated with y_1, \dots, y_k be $\alpha_1, \alpha_2, \dots, \alpha_k$, where $\alpha_1 = \infty$ because $\text{Prob}[Y \geq y_1] = 1$. Let $\alpha_y = \alpha_i, i : y_i = y$. Then

$$\text{Prob}[Y \geq y_i|X] = F(\alpha_i + X\beta) = F(\alpha_{y_i} + X\beta) \quad (15.2)$$

For the OLS fully parametric case, the model may be restated

$$\text{Prob}[Y \geq y|X] = \text{Prob}\left[\frac{Y - X\beta}{\sigma} \geq \frac{y - X\beta}{\sigma}\right] \quad (15.3)$$

$$= 1 - \Phi\left(\frac{y - X\beta}{\sigma}\right) = \Phi\left(\frac{-y}{\sigma} + \frac{X\beta}{\sigma}\right) \quad (15.4)$$

^b Quantile regression allows the estimated value of the 0.5 quantile to be higher than the estimated value of the 0.6 quantile for some values of X . Composite quantile regression⁶⁹⁰ removes this possibility by forcing all the X coefficients to be the same across multiple quantiles, a restriction not unlike what cumulative probability ordinal models make.

^c For symmetric distributions applying a decreasing transformation will negate the coefficients. For asymmetric distributions (e.g., Gumbel), reversing the order of Y will do more than change signs.

^d Only an estimate of mean Y from these $\hat{\beta}$ s is non-robust.

^e It is more traditional to state the model in terms of $\text{Prob}[Y \leq y|X]$ but we use $\text{Prob}[Y \geq y|X]$ so that higher predicted values are associated with higher Y .

Table 15.1 Distribution families used in ordinal cumulative probability models. Φ denotes the Gaussian cumulative distribution function. For the Connection column, $P_1 = \text{Prob}[Y \geq y|X_1], P_2 = \text{Prob}[Y \geq y|X_2], \Delta = (X_2 - X_1)\beta$. The connection specifies the only distributional assumption if the model is fitted semiparametrically, i.e. contains an intercept for every unique Y value less one. For parametric models, P_1 must be specified absolutely instead of just requiring a relationship between P_1 and P_2 . For example, the traditional Gaussian parametric model specifies that $\text{Prob}[Y \geq y|X] = 1 - \Phi(\frac{y-X\beta}{\sigma}) = \Phi(\frac{-y+X\beta}{\sigma})$.

Distribution	F	Inverse (Link Function)	Link Name	Connection
Logistic	$[1 + \exp(-y)]^{-1}$	$\log(\frac{y}{1-y})$	logit	$\frac{P_2}{1-P_2} = \frac{P_1}{1-P_1} \exp(\Delta)$
Gaussian	$\Phi(y)$	$\Phi^{-1}(y)$	probit	$P_2 = \Phi(\Phi^{-1}(P_1) + \Delta)$
Gumbel maximum value	$\exp(-\exp(-y))$	$\log(-\log(y))$	log - log	$P_2 = P_1^{\exp(\Delta)}$
Gumbel minimum value	$1 - \exp(-\exp(y))$	$\log(-\log(1-y))$	complementary log - log	$1 - P_2 = (1 - P_1)^{\exp(\Delta)}$
Cauchy	$\frac{1}{\pi} \tan^{-1}(y) + \frac{1}{2}$	$\tan[\pi(y - \frac{1}{2})]$	cauchit	

so that to within an additive constant^f $\alpha_y = \frac{-y}{\sigma}$ (intercepts α are linear in y whereas they are arbitrarily descending in the ordinal model), and σ is absorbed in β to put the OLS model into the new notation.

The general ordinal regression model assumes that for fixed X_1, X_2 ,

$$F^{-1}(\text{Prob}[Y \geq y|X_2]) - F^{-1}(\text{Prob}[Y \geq y|X_1]) \tag{15.5}$$

$$= (X_2 - X_1)\beta \tag{15.6}$$

independent of the α s (parallelism assumption). If $F = [1 + \exp(-y)]^{-1}$, this is the proportional odds assumption.

Common choices of F , implemented in the R `rms orm` function, are shown in Table 15.1. The Gumbel maximum value distribution is also called the extreme value type I distribution. This distribution (log - log link) also represents a continuous time proportional hazards model. The hazard ratio when X changes from X_1 to X_2 is $\exp(-(X_2 - X_1)\beta)$.

The mean of $Y|X$ is easily estimated from a fitted cumulative probability ordinal model by computing

$$\sum_{i=1}^n y_i \widehat{\text{Prob}}[Y = y_i|X] \tag{15.7}$$

and the q^{th} quantile of $Y|X$ is y such that $F^{-1}(1 - q) - X\hat{\beta} = \hat{\alpha}_y$.^g

^f $\hat{\alpha}_y$ are unchanged if a constant is added to all y .

^g The intercepts have to be shifted to the left one position in solving this equation because the quantile is such that $\text{Prob}[Y \leq y] = q$ whereas the model is stated in terms of $\text{Prob}[Y \geq y]$.

The `orm` function in the `rms` package takes advantage of the information matrix being of a sparse tri-band diagonal form for the intercept parameters. This makes the computations efficient even for hundreds of intercepts (i.e., unique values of Y). `orm` is made to handle continuous Y .

Ordinal regression has nice properties in addition to those listed above, allowing for

- estimation of quantiles as efficiently as quantile regression if the parallel slopes assumptions hold
- efficient estimation of mean Y
- direct estimation of $\text{Prob}[Y \geq y|X]$
- arbitrary clumping of values of Y , while still estimating β and mean Y efficiently^h
- solutions for $\hat{\beta}$ using ordinary Newton-Raphson or other popular optimization techniques
- being based on a standard likelihood function, penalized estimation can be straightforward
- Wald, score, and likelihood ratio χ^2 tests that are more powerful than tests from quantile regression.

On the last point, if there is a single predictor in the model and it is binary, the score test from the proportional odds model is essentially the Wilcoxon test, and the score test from the Gumbel log-log cumulative probability model is essentially the log-rank test.

15.3.1 Minimum Sample Size Requirement

When Y is continuous and the purpose of an ordinal model includes semi-parametric estimation of probabilities or quantiles, the accuracy of estimates is limited even more by the accuracy of estimating the empirical cumulative distribution of Y than by estimating β . When $\beta = 0$, intercept estimates are transformations of the empirical distribution step function. As described in Section 20.3, the sample size must be 184 to estimate the entire distribution of Y with a global margin of error not exceeding 0.1. For estimating the mean of Y , smaller sample sizes may be needed.

^h But it is not sensible to estimate quantiles of Y when there are heavy ties in Y in the area containing the quantile.

15.4 Comparison of Assumptions of Various Models

Quantile regression makes the fewest left-hand-side model assumptions except for the assumption that Y be continuous, but can have less estimator precision than other models and has lower power. To summarize how assumptions of parametric models compare to assumptions of semiparametric ordinal models, consider the ordinary linear model or its special case the equal variance two-sample t -test, vs. the probit or logit (proportional odds) ordinal model or their special cases the Van der Waerden (normal-scores) two-sample rank test or the Wilcoxon two-sample test. All the assumptions of the linear model other than independence of residuals are captured in the following, using the more standard $Y \leq y$ notation:

$$F(y|X) = \text{Prob}[Y \leq y|X] = \Phi\left(\frac{y - X\beta}{\sigma}\right) \tag{15.8}$$

$$\Phi^{-1}(F(y|X)) = \frac{y - X\beta}{\sigma} \tag{15.9}$$

On the other hand, ordinal models assume the following:

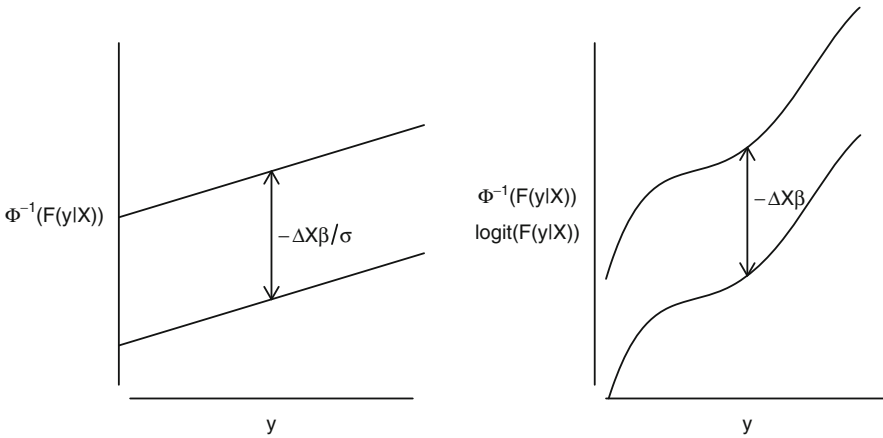


Fig. 15.1 Assumptions of the linear model (left panel) and semiparametric ordinal probit or logit (proportional odds) models (right panel). Ordinal models do not assume any shape for the distribution of Y for a given X ; they only assume parallelism. The linear model can relax the parallelism assumption if σ is allowed to vary, but in practice it is difficult to know how to vary it except for the unequal variance two-sample t -test.

$$\text{Prob}[Y \leq y|X] = F(g(y) - X\beta), \tag{15.10}$$

where g is unknown and may be discontinuous. This translates to the parallelism assumption in the right panel of Figure 15.1, whereas the linear model

makes the additional strong assumption of linearity of normal inverse cumulative distribution function, which arises from the Gaussian distribution assumption.

15.5 Dataset and Descriptive Statistics

Diabetes Mellitus (DM) type II (adult onset diabetes) is strongly associated with obesity. The currently best laboratory test for diabetes measures glycosylated hemoglobin (HbA_{1c}), also called glycated hemoglobin, glycohemoglobin, or hemoglobin A_{1c} . HbA_{1c} reflects average blood glucose for the preceding 60 to 90 days. $\text{HbA}_{1c} > 7.0$ is sometimes taken as a positive diagnosis of diabetes even though there are no data to support the use of a threshold.

The goals of this analyses are to better understand effects of body size measurements on risk of DM and to enhance screening for DM. The best way to develop a model for DM screening is **not** to fit a binary logistic model with $\text{HbA}_{1c} > 7$ as the response variable. There are at least two reasons for this. First, when the relationship between a measurement and its ultimate clinical impact is smooth, all cutpoints are arbitrary. There is no justification for any putative cut on HbA_{1c} . Second, such an analysis loses information by treating $\text{HbA}_{1c}=2$ the same as $\text{HbA}_{1c}=6.9$, and by treating $\text{HbA}_{1c}=7.1$ as equal to $\text{HbA}_{1c}=10$. Failure to use all available information results in larger standard errors of $\hat{\beta}$, lower power, and wider confidence bands. It is better to predict continuous HbA_{1c} using a continuous response model, then use that model to estimate the probability that HbA_{1c} exceeds any cutoff, or estimate the 0.9 quantile of HbA_{1c} .

The data used here are from the National Health and Nutrition Examination Survey (NHANES) 2009–2010 from the U.S. National Center for Health Statistics/Centers for Disease Control. The original data may be obtained from <http://www.cdc.gov/nchs/nhanes.htm>⁹⁴; the analysis file used here, called `nhgh`, may be obtained from the `DataSets` wiki page, along with R code used to download and create the file. Note that CDC coded age ≥ 80 as 80. We use the subset of subjects with age ≥ 21 who have neither been diagnosed nor treated for DM. Descriptive statistics are shown below.

```
require(rms)
```

```
getHdata(nhgh)
w ← subset(nhgh, age ≥ 21 & dx==0 & tx==0, select=-c(dx,tx))
latex(describe(w), file='')
```

18 Variables ^W
4629 Observations

```

seqn : Respondent sequence number
      n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4629   0 4629   1 56902 52136 52633 54284 56930 59495 61079 61641

lowest : 51624 51629 51630 51645 51647
highest: 62152 62153 62155 62157 62158

```

```

sex
  n missing unique
4629   0     2

male (2259, 49%), female (2370, 51%)

```

```

age : Age [years]
      n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4629   0   703   1 48.57 23.33 26.08 33.92 46.83 61.83 74.83 80.00

lowest : 21.00 21.08 21.17 21.25 21.33
highest: 79.67 79.75 79.83 79.92 80.00

```

```

re : Race/Ethnicity
     n missing unique
4629   0     5

Mexican American (832, 18%), Other Hispanic (474, 10%)
Non-Hispanic White (2318, 50%), Non-Hispanic Black (756, 16%)
Other Race Including Multi-Racial (249, 5%)

```

```

income : Family Income
        n missing unique
4389    240    14

[0,5000) (162, 4%), [5000,10000) (216, 5%), [10000,15000) (371, 8%)
[15000,20000) (300, 7%), [20000,25000) (374, 9%)
[25000,35000) (535, 12%), [35000,45000) (421, 10%)
[45000,55000) (346, 8%), [55000,65000) (257, 6%), [65000,75000) (188, 4%)
> 20000 (149, 3%), < 20000 (52, 1%), [75000,100000) (399, 9%)
>= 100000 (619, 14%)

```

```

wt : Weight [kg]
      n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4629   0   890   1 80.49 52.44 57.18 66.10 77.70 91.40 106.52 118.00

lowest : 33.2 36.1 37.9 38.5 38.7
highest: 184.3 186.9 195.3 196.6 203.0


```

```

ht : Standing Height [cm]
     n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4629   0   512   1 167.5 151.1 154.4 160.1 167.2 175.0 181.0 184.8


lowest : 123.3 135.4 137.5 139.4 139.8
highest: 199.2 199.3 199.6 201.7 202.7

```

bmi : Body Mass Index [kg/m²] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4629	0	1994	1	28.59	20.02	21.35	24.12	27.60	31.88	36.75	40.68

lowest : 13.18 14.59 15.02 15.40 15.49
highest: 61.20 62.81 65.62 71.30 84.87

leg : Upper Leg Length [cm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4474	155	216	1	38.39	32.0	33.5	36.0	38.4	41.0	43.3	44.6

lowest : 20.4 24.9 25.0 25.1 26.4, highest: 49.0 49.5 49.8 50.0 50.3

arml : Upper Arm Length [cm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4502	127	156	1	37.01	32.6	33.5	35.0	37.0	39.0	40.6	41.7

lowest : 24.8 27.0 27.5 29.2 29.5, highest: 45.2 45.5 45.6 46.0 47.0

armc : Arm Circumference [cm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4499	130	290	1	32.87	25.4	26.9	29.5	32.5	35.8	39.1	41.4

lowest : 17.9 19.0 19.3 19.5 19.9, highest: 54.2 54.9 55.3 56.0 61.0

waist : Waist Circumference [cm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4465	164	716	1	97.62	74.8	78.6	86.9	96.3	107.0	117.8	125.0

lowest : 59.7 60.0 61.5 62.0 62.4
highest: 160.0 160.6 162.2 162.7 168.7

tri : Triceps Skinfold [mm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4295	334	342	1	18.94	7.2	8.8	12.0	18.0	25.2	31.0	33.8

lowest : 2.6 3.1 3.2 3.3 3.4, highest: 39.6 39.8 40.0 40.2 40.6

sub : Subscapular Skinfold [mm] 


n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
3974	655	329	1	20.8	8.60	10.30	14.40	20.30	26.58	32.00	35.00

lowest : 3.8 4.2 4.6 4.8 4.9, highest: 40.0 40.1 40.2 40.3 40.4

gh : Glycohemoglobin [%] 

n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4629	0	63	0.99	5.533	4.8	5.0	5.2	5.5	5.8	6.0	6.3

lowest : 4.0 4.1 4.2 4.3 4.4, highest: 11.9 12.0 12.1 12.3 14.5

albumin : Albumin [g/dL] 

n	missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
4576	53	26	0.99	4.261	3.7	3.9	4.1	4.3	4.5	4.7	4.8

lowest : 2.6 2.7 3.0 3.1 3.2, highest: 4.9 5.0 5.1 5.2 5.3

```

bun : Blood urea nitrogen [mg/dL]
  n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4576      53      50 0.99 13.03 7 8 10 12 15 19 22
lowest : 1 2 3 4 5, highest: 49 53 55 56 63

```

```

SCr : Creatinine [mg/dL]
  n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
4576      53      167 1 0.8887 0.58 0.62 0.72 0.84 0.99 1.14 1.25
lowest : 0.34 0.38 0.39 0.40 0.41
highest: 5.98 6.34 9.13 10.98 15.66

```

```
dd ← datadist(w); options(datadist='dd')
```

15.5.1 Checking Assumptions of OLS and Other Models

First let's see if `gh` would make a Gaussian residuals model fit. Use ordinary regression on four key variables to collapse these into one variable (predicted mean from the OLS model). Stratify the predicted means into six quantile groups. Apply the normal inverse cumulative distribution function Φ^{-1} to the empirical cumulative distribution functions (ECDF) of `gh` using these strata, and check for normality and constant σ^2 . The ECDF estimates $\text{Prob}[Y \leq y|X]$ but for ordinal modeling we want to state models in terms of $\text{Prob}[Y \geq y|X]$ so take one minus the ECDF before inverse transforming.

```

f ← ols(gh ~ rcs(age,5) + sex + re + rcs(bmi, 3), data=w)
pgh ← fitted(f)

p ← function(fun, row, col) {
  f ← substitute(fun); g ← function(F) eval(f)
  z ← Ecdf(~ gh, groups=cut2(pgh, g=6),
           fun=function(F) g(1 - F),
           ylab=as.expression(f), xlim=c(4.5, 7.75), data=w,
           label.curve=FALSE)
  print(z, split=c(col, row, 2, 2), more=row < 2 | col < 2)
}
p(log(F/(1-F)), 1, 1)
p(qnorm(F), 1, 2)
p(-log(-log(F)), 2, 1)
p(log(-log(1-F)), 2, 2)
# Get slopes of pgh for some cutoffs of Y
# Use glm complementary log-log link on Prob(Y < cutoff) to
# get log-log link on Prob(Y ≥ cutoff)
r ← NULL
for(link in c('logit', 'probit', 'cloglog'))
  for(k in c(5, 5.5, 6)) {
    co ← coef(glm(gh < k ~ pgh, data=w, family=binomial(link)))

```

```
r ← rbind(r, data.frame(link=link, cutoff=k,
                        slope=round(co[2],2)))
}
print(r, row.names=FALSE)
```

link	cutoff	slope
logit	5.0	-3.39
logit	5.5	-4.33
logit	6.0	-5.62
probit	5.0	-1.69
probit	5.5	-2.61
probit	6.0	-3.07
cloglog	5.0	-3.18
cloglog	5.5	-2.97
cloglog	6.0	-2.51

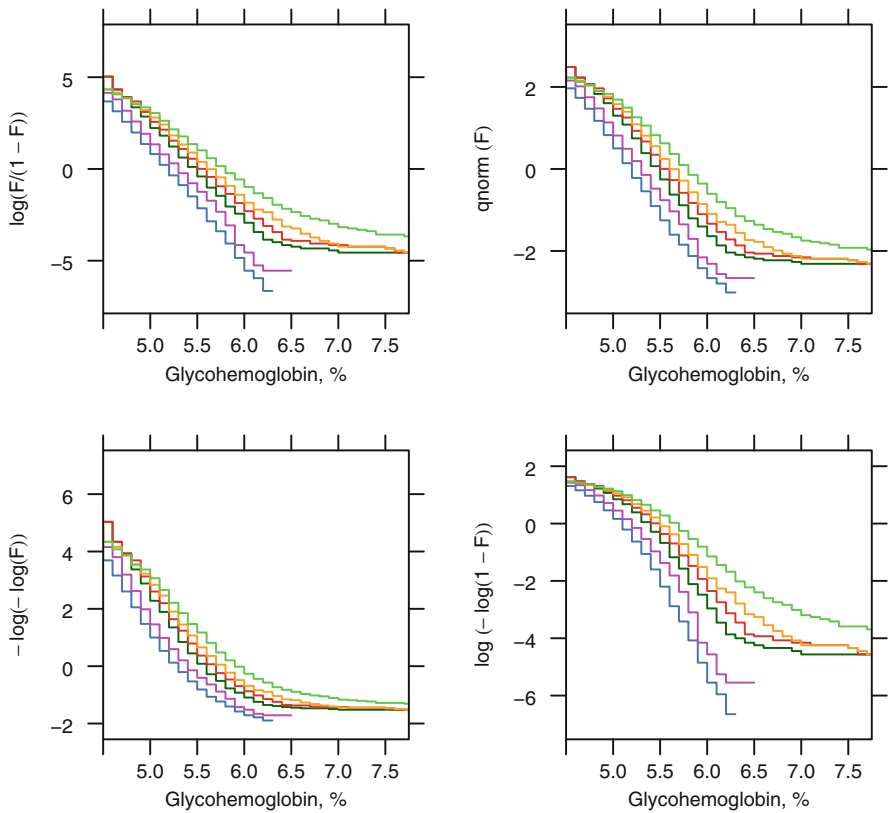


Fig. 15.2 Examination of normality and constant variance assumption, and assumptions for various ordinal models

The upper right curves in Figure 15.2 are not linear, implying that a normal conditional distribution cannot work for gh^i . There is non-parallelism for the logit model. The other graphs will be used to guide selection of an ordinal model below.

15.6 Ordinal Regression Applied to HbA_{1c}

In the upper left panel of Figure 15.2, logit inverse curves are not parallel so the proportional odds assumption does not hold when predicting HbA_{1c} . The log-log link yields the highest degree of parallelism and most constant regression coefficients across cutoffs of gh , so we use this link in an ordinal regression model (linearity of the curves is not required).

15.6.1 Checking Fit for Various Models Using Age

Another way to examine model fit is to flexibly fit the single most important predictor (age) using a variety of methods, and compare predictions to sample quantiles and means based on subsets on age. We use overlapping subsets to gain resolution, with each subset composed of those subjects having age within five years of the point being predicted by the models. Here we predict the 0.5, 0.75, and 0.9 quantiles and the mean. For quantiles we can compare to quantile regression (discussed below) and for means we compare to OLS.

```

ag ← 25:75
lag ← length(ag)
q2 ← q3 ← p90 ← means ← numeric(lag)
for(i in 1:lag) {
  s ← which(abs(w$age - ag[i]) < 5)
  y ← w$gh[s]
  a ← quantile(y, probs=c(.5, .75, .9))
  q2[i] ← a[1]
  q3[i] ← a[2]
  p90[i] ← a[3]
  means[i] ← mean(y)
}
fams ← c('logistic', 'probit', 'loglog', 'cloglog')
fe ← function(pred, target) mean(abs(pred$yhat - target))
mod ← gh ~ rcs(age,6)
P ← Er ← list()
for(est in c('q2', 'q3', 'p90', 'mean')) {
  meth ← if(est == 'mean') 'ols' else 'QR'
  p ← list()
  er ← rep(NA, 5)
  names(er) ← c(fams, meth)
  for(family in fams) {
    h ← orm(mod, family=family, data=w)
    fun ← if(est == 'mean') Mean(h)
    else {
      qu ← Quantile(h)

```

ⁱ They are not parallel either.

```

      switch(est, q2 = function(x) qu(.5, x),
                q3 = function(x) qu(.75, x),
                p90 = function(x) qu(.9, x))
    }
    p[[family]] ← z ← Predict(h, age=ag, fun=fun, conf.int=FALSE)
    er[[family]] ← fe(z, switch(est, mean=means, q2=q2, q3=q3, p90=p90))
  }
  h ← switch(est,
             mean= ols(mod, data=w),
             q2 = Rq (mod, data=w),
             q3 = Rq (mod, tau=0.75, data=w),
             p90 = Rq (mod, tau=0.90, data=w))
  p[[meth]] ← z ← Predict(h, age=ag, conf.int=FALSE)
  er[[meth]] ← fe(z, switch(est, mean=means, q2=q2, q3=q3, p90=p90))

  Er[[est]] ← er
  pr ← do.call('rbind', p)
  pr$est ← est
  P ← rbind.data.frame(P, pr)
}

xyplot(yhat ~ age | est, groups=.set., data=P, type='l', # Figure 15.3
       auto.key=list(x=.75, y=.2, points=FALSE, lines=TRUE),
       panel=function(..., subscripts) {
         panel.xyplot(..., subscripts=subscripts)
         est ← P$est[subscripts[1]]
         lpoints(ag, switch(est, mean=means, q2=q2, q3=q3, p90=p90),
                col=gray(.7))
         er ← format(round(Er[[est]],3), nsmall=3)
         ltext(26, 6.15, paste(names(er), collapse='\n'),
              cex=.7, adj=0)
         ltext(40, 6.15, paste(er, collapse='\n'),
              cex=.7, adj=1)}

```

It can be seen in Figure 15.3 that models dedicated to a specific task (quantile regression for quantiles and OLS for means) were best for those tasks. Although the log-log ordinal cumulative probability model did not estimate the median as accurately as some other methods, it does well for the 0.75 and 0.9 quantiles and is the best compromise overall because of its ability to also directly predict the mean as well as quantities such as $\text{Prob}[\text{HbA}_{1c} > 7|X]$.

From here on we focus on the log-log ordinal model. Returning to the bottom left of Figure 15.2, let's look at quantile groups of predicted HbA_{1c} by OLS and plot predicted distributions of actual HbA_{1c} against empirical distributions.

```

w$pghg ← cut2(pgh, g=6)
f ← orm(gh ~ pghg, data=w)
lp ← predict(f, newdata=data.frame(pghg=levels(w$pghg)))
ep ← ExProb(f) # Exceedance prob. functn. generator in rms
z ← ep(lp)
j ← order(w$pghg) # puts in order of lp (levels of pghg)
plot(z, xlim=c(4, 7.5), data=w[j,c('pghg', 'gh')]) # Fig. 15.4

```

Agreement between predicted and observed exceedance probability distributions is excellent in Figure 15.4.

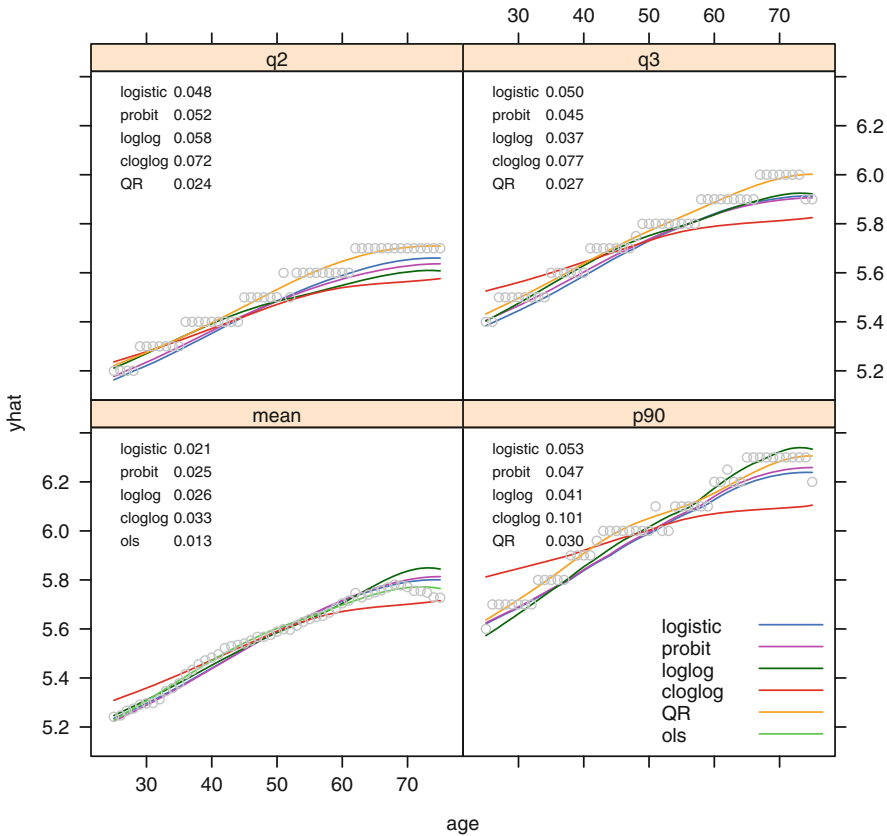


Fig. 15.3 Three estimated quantiles and estimated mean using 6 methods, compared against caliper-matched sample quantiles/means (circles). Numbers are mean absolute differences between predicted and sample quantities using overlapping intervals of age and caliper matching. QR:quantile regression.

To return to the initial look at a linear model with assumed Gaussian residuals, fit a probit ordinal model and compare the estimated intercepts to the linear relationship with gh that is assumed by the normal distribution.

```
f ← orm(gh ~ rcs(age,6), family=probit, data=w)
g ← ols(gh ~ rcs(age,6), data=w)
s ← g$stats['Sigma']
yu ← f$yunique[-1]
r ← quantile(w$gh, c(.005, .995))
alphas ← coef(f)[1:num.intercepts(f)]
plot(-yu / s, alphas, type='l', xlim=rev(- r / s), # Fig. 15.5
      xlab=expression(-y/hat(sigma)), ylab=expression(alpha[y]))
```

Figure 15.5 depicts a significant departure from the linear form implied by Gaussian residuals (Eq. 15.4).

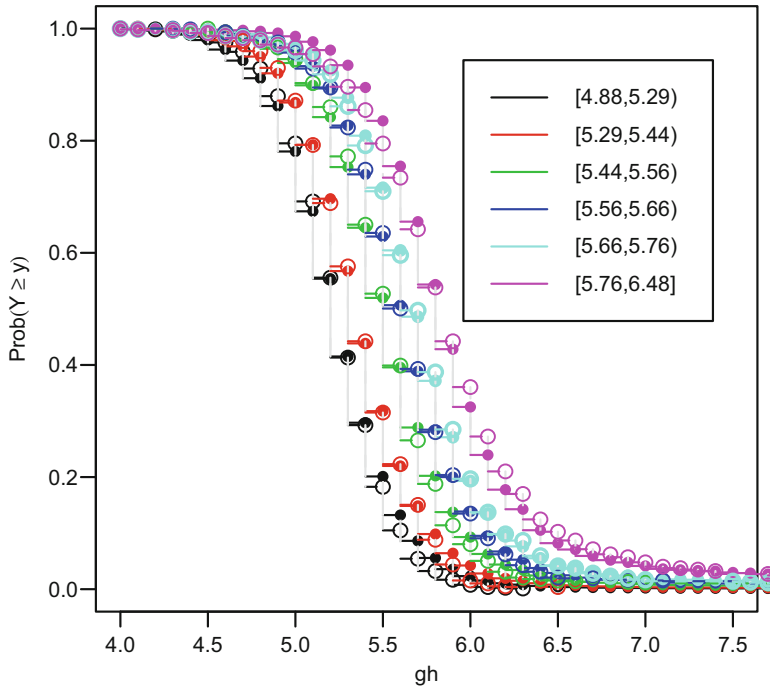


Fig. 15.4 Observed (dashed lines, open circles) and predicted (solid lines, closed circles) exceedance probability distributions from a model using 6-tiles of OLS-predicted HbA_{1c}. Key shows quantile group intervals of predicted mean HbA_{1c}.

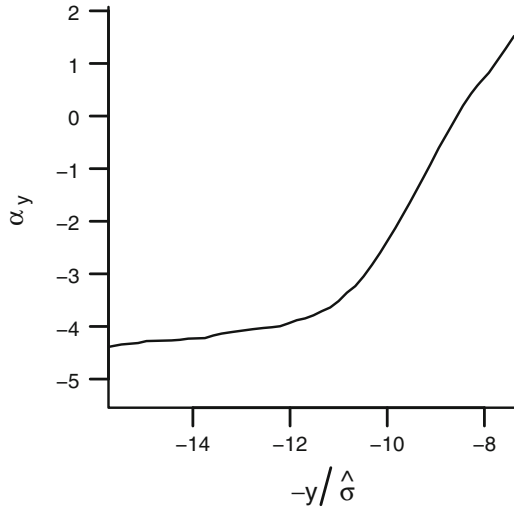


Fig. 15.5 Estimated intercepts from probit model. Linearity would have indicated Gaussian residuals.

15.6.2 Examination of BMI

Body mass index (BMI, weight divided by height²) is commonly used as an obesity measure because it is well correlated with abdominal visceral fat. But it is not obvious that BMI is the correct summary of height and weight for predicting pre-clinical diabetes, and it may be the case that body size measures other than height and weight are better predictors.

Use the log-log ordinal model to check the adequacy of BMI, adjusting for age (without assuming linearity). This can be done by examining the ratio of coefficients of log height and log weight, and also by using AIC to judge whether BMI is an adequate summary of height and weight when compared to nonlinear functions of the logs, and to a tensor spline interaction surface.

```
f ← orm(gh ~ rcs(age,5) + log(ht) + log(wt),
        family=loglog, data=w)
print(f, latex=TRUE)
```

-log-log Ordinal Regression Model

```
orm(formula = gh ~ rcs(age, 5) + log(ht) + log(wt), data = w,
    family = loglog)
```

	Model Likelihood	Discrimination	Rank Discrim.
	Ratio Test	Indexes	Indexes
Obs 4629	LR χ^2 1126.94	R^2 0.217	ρ 0.486
Unique Y 63	d.f. 6	g 0.627	
$Y_{0.5}$ 5.5	$\Pr(> \chi^2) < 0.0001$	g_r 1.872	
$\max \frac{\partial \log L}{\partial \beta} $	Score χ^2 1262.81	$ \Pr(Y \geq Y_{0.5}) - \frac{1}{2} $ 0.153	
1×10^{-6}	$\Pr(> \chi^2) < 0.0001$		

	Coef	S.E.	Wald Z	$\Pr(> Z)$
age	0.0398	0.0055	7.29	< 0.0001
age'	-0.0158	0.0275	-0.57	0.5657
age''	-0.0072	0.0866	-0.08	0.9333
age'''	0.0309	0.1135	0.27	0.7853
ht	-3.0680	0.2789	-11.00	< 0.0001
wt	1.2748	0.0704	18.10	< 0.0001

```
aic ← NULL
for(mod in list(gh ~ rcs(age,5) + rcs(log(bmi),5),
               gh ~ rcs(age,5) + rcs(log(ht),5) + rcs(log(wt),5),
               gh ~ rcs(age,5) + rcs(log(ht),4) * rcs(log(wt),4)))
  aic ← c(aic, AIC(orm(mod, family=loglog, data=w)))
print(aic)
```

[1] 25910.77 25910.17 25906.03

The ratio of the coefficient of log height to the coefficient of log weight is -2.4, which is between what BMI uses and the more dimensionally reasonable weight / height³. By AIC, a spline interaction surface between height and weight does slightly better than BMI in predicting HbA_{1c}, but a nonlinear function of BMI is barely worse. It will require other body size measures to displace BMI as a predictor.

As an aside, compare this model fit to that from the Cox proportional hazards model. The Cox model uses a conditioning argument to obtain a partial likelihood free of the intercepts α (and requires a second step to estimate these log discrete hazard components) whereas we are using a full marginal likelihood of the ranks of Y ³³⁰.

```
print(cph(Surv(gh) ~ rcs(age,5) + log(ht) + log(wt), data=w),
      latex=TRUE)
```

Cox Proportional Hazards Model

```
cph(formula = Surv(gh) ~ rcs(age, 5) + log(ht)
    + log(wt), data = w)
```

		Model Tests		Discrimination Indexes	
Obs	4629	LR χ^2	1120.20	R^2	0.215
Events	4629	d.f.	6	D_{xy}	0.359
Center	8.3792	Pr(> χ^2)	0.0000	g	0.622
		Score χ^2	1258.07	g_r	1.863
		Pr(> χ^2)	0.0000		

	Coef	S.E.	Wald Z	Pr(> Z)
age	-0.0392	0.0054	-7.24	< 0.0001
age'	0.0148	0.0274	0.54	0.5888
age''	0.0093	0.0862	0.11	0.9144
age'''	-0.0321	0.1131	-0.28	0.7767
ht	3.0477	0.2779	10.97	< 0.0001
wt	-1.2653	0.0701	-18.04	< 0.0001

Close agreement of the two is seen, as expected.

15.6.3 Consideration of All Body Size Measurements

Next we examine all body size measures, and check their redundancies.

```
v <- varclus(~ wt + ht + bmi + leg + arml + armc + waist +
             tri + sub + age + sex + re, data=w)
plot(v)
```

```
# Omit wt so it won't be removed before bmi
redun(~ ht + bmi + leg + arml + armc + waist + tri + sub,
      data=w, r2=.75)
```

```
Redundancy Analysis

redun(formula = ~ht + bmi + leg + arml + armc + waist + tri +
      sub, data = w, r2 = 0.75)

n: 3853          p: 8          nk: 3

Number of NAs:  776
Frequencies of Missing Values Due to Each Variable
  ht   bmi   leg  arml  armc  waist  tri   sub
  0     0   155  127   130   164   334   655

Transformation of target variables forced to be linear

R2 cutoff: 0.75          Type: ordinary

R2 with which each variable can be predicted from all other variables:

  ht   bmi   leg  arml  armc  waist  tri   sub
0.829 0.924 0.682 0.748 0.843 0.864 0.531 0.594

Rendundant variables:

bmi ht

Predicted from variables:

leg arml armc waist tri sub

Variable Deleted  R2 R2 after later deletions
1          bmi 0.924                      0.909
2          ht 0.792
```

Six size measures adequately capture the entire set. Height and BMI are removed (Figure 15.6). An advantage of removing height is that it is age-dependent due to vertebral compression in the elderly:

```
f <- orm(ht ~ rcs(age,4)*sex, data=w) # Prop. odds model
qu <- Quantile(f); med <- function(x) qu(.5, x)
ggplot(Predict(f, age, sex, fun=med, conf.int=FALSE),
       ylab='Predicted Median Height, cm')
```

However, upper leg length has the same declining trend, implying a survival bias or birth year effect.

In preparing to create a multivariable model, degrees of freedom are allocated according to the generalized Spearman ρ^2 (Figure 15.7)^j.

```
s <- spearman2(gh ~ age + sex + re + wt + leg + arml + armc +
              waist + tri + sub, data=w, p=2)
plot(s)
```

Parameters will be allocated in descending order of ρ^2 . But note that subscapular skinfold has a large number of NAs and other predictors also have NAs. Suboptimal casewise deletion will be used until the final model is fitted (Figure 15.8).

^j Competition between collinear size measures hurts interpretation of partial tests of association in a saturated additive model.

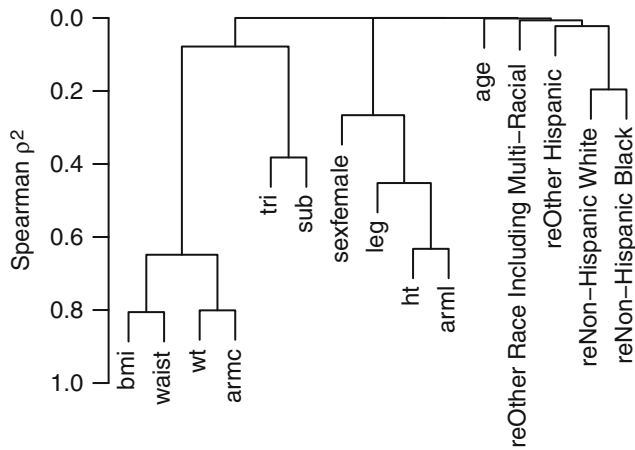


Fig. 15.6 Variable clustering for all potential predictors

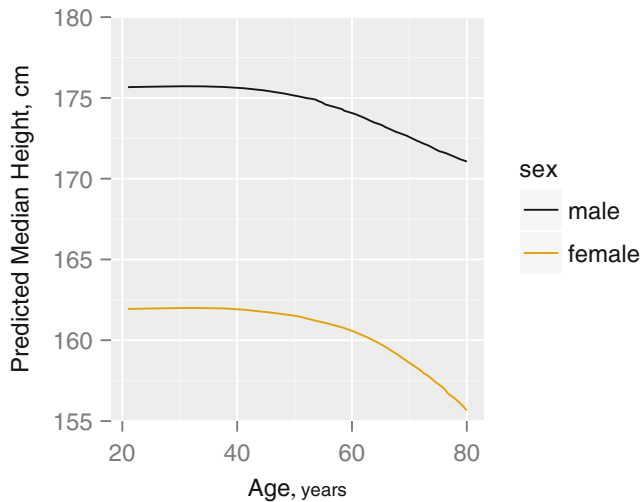


Fig. 15.7 Estimated median height as a smooth function of age, allowing age to interact with sex, from a proportional odds model

Because there are many competing body measures, we use backwards step-down to arrive at a set of predictors. The bootstrap will be used to penalize predictive ability for variable selection. First the full model is fit using casewise deletion, then we do a composite test to assess whether any of the frequently-missing predictors is important.

```
f ← orm(gh ~ rcs(age,5) + sex + re + rcs(wt,3) + rcs(leg,3) + arml +
        rcs(armc,3) + rcs(waist,4) + tri + rcs(sub,3),
        family='loglog', data=w, x=TRUE, y=TRUE)
print(f, latex=TRUE, coefs=FALSE)
```

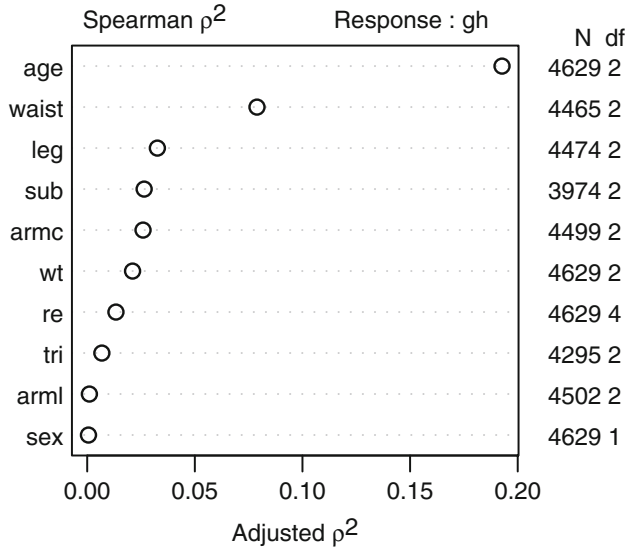
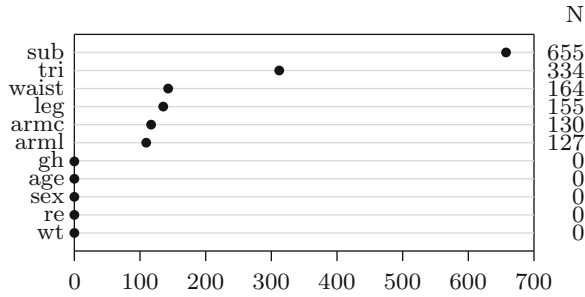


Fig. 15.8 Generalized squared rank correlations

-log-log Ordinal Regression Model

```
orm(formula = gh ~ rcs(age, 5) + sex + re + rcs(wt, 3)
    + rcs(leg, 3) + arml + rcs(armc, 3) + rcs(waist, 4)
    + tri + rcs(sub, 3), data = w, x = TRUE, y = TRUE,
    family = "loglog")
```

Frequencies of Missing Values Due to Each Variable



		Model Likelihood	Discrimination		Rank Discrim.		
		Ratio Test	Indexes		Indexes		
Obs	3853	LR χ^2	1180.13	R^2	0.265	ρ	0.520
Unique Y	60	d.f.	22	g	0.732		
$Y_{0.5}$	5.5	$\Pr(> \chi^2) < 0.0001$		g_r	2.080		
$\max \frac{\partial \log L}{\partial \beta} $		Score χ^2	1298.88	$ \Pr(Y \geq Y_{0.5}) - \frac{1}{2} $		0.172	
3×10^{-5}		$\Pr(> \chi^2) < 0.0001$					

```
## Composite test:
lan <- function(a) latex(a, table.env=FALSE, file='')
lan(anova(f, leg, arml, armc, waist, tri, sub))
```

	χ^2	d.f.	P
leg	8.30	2	0.0158
<i>Nonlinear</i>	3.32	1	0.0685
arml	0.16	1	0.6924
armc	6.66	2	0.0358
<i>Nonlinear</i>	3.29	1	0.0695
waist	29.40	3	< 0.0001
<i>Nonlinear</i>	4.29	2	0.1171
tri	16.62	1	< 0.0001
sub	40.75	2	< 0.0001
<i>Nonlinear</i>	4.50	1	0.0340
TOTAL NONLINEAR	14.95	5	0.0106
TOTAL	128.29	11	< 0.0001

The model achieves Spearman $\rho = 0.52$, the rank correlation between predicted and observed HbA_{1c}.

We show the predicted mean and median HbA_{1c} as a function of age, adjusting other variables to their median or mode (Figure 15.9). Compare the estimate of the median and 90th percentile with that from quantile regression.

```
M <- Mean(f)
qu <- Quantile(f)
med <- function(x) qu(.5, x)
p90 <- function(x) qu(.9, x)
fq <- Rq(formula(f), data=w)
fq90 <- Rq(formula(f), data=w, tau=.9)
```

```
pmean <- Predict(f, age, fun=M, conf.int=FALSE)
pmed <- Predict(f, age, fun=med, conf.int=FALSE)
p90 <- Predict(f, age, fun=p90, conf.int=FALSE)
pmedqr <- Predict(fq, age, conf.int=FALSE)
p90qr <- Predict(fq90, age, conf.int=FALSE)
z <- rbind('orm mean'=pmean, 'orm median'=pmed, 'orm P90'=p90,
          'QR median'=pmedqr, 'QR P90'=p90qr)
ggplot(z, groups='.set.',
       adj.subtitle=FALSE, legend.label=FALSE)
```



```
print(fastbw(f, rule='p'), estimates=FALSE)
```

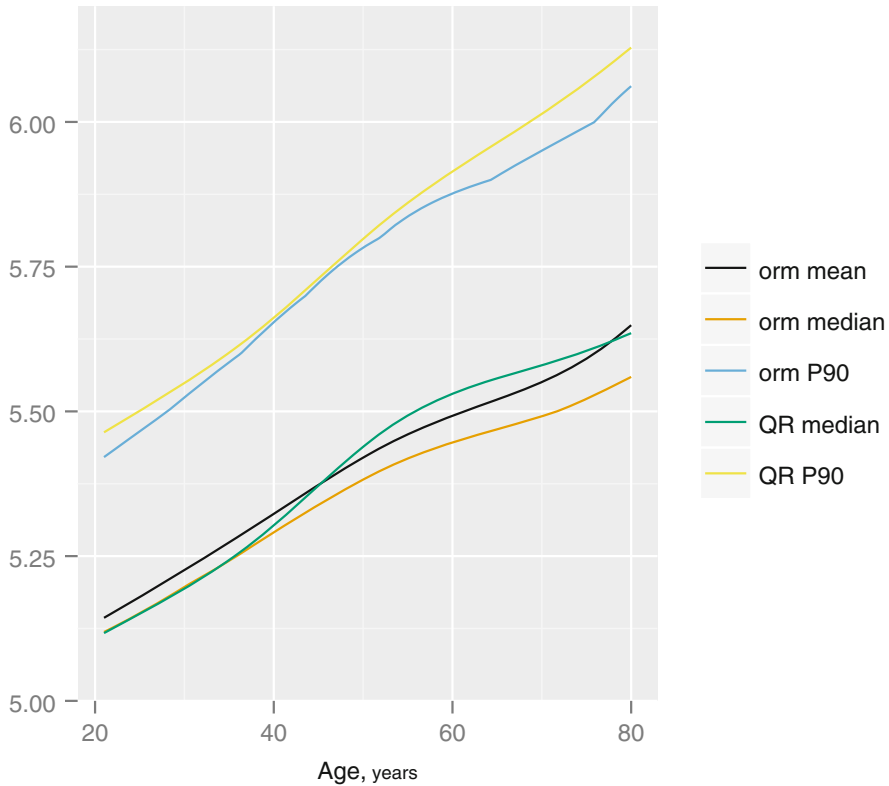


Fig. 15.9 Estimated mean and 0.5 and 0.9 quantiles from the log-log ordinal model using casewise deletion, along with predictions of 0.5 and 0.9 quantiles from quantile regression (QR). Age is varied and other predictors are held constant to medians/-modes.

```
Deleted Chi-Sq d.f. P      Residual d.f. P      AIC
arml  0.16  1  0.6924 0.16  1  0.6924 -1.84
sex   0.45  1  0.5019 0.61  2  0.7381 -3.39
wt    5.72  2  0.0572 6.33  4  0.1759 -1.67
armc  3.32  2  0.1897 9.65  6  0.1400 -2.35

Factors in Final Model

[1] age  re  leg  waist tri  sub
```

```
set.seed(13) # so can reproduce results
v <- validate(f, B=100, bw=TRUE, estimates=FALSE, rule='p')
```


Frequencies of Numbers of Factors Retained

5	6	7	8	9	10
1	19	29	46	4	1

Next we fit the reduced model, using multiple imputation to impute missing predictors (Figure 15.10).

```
a ← aregImpute(~ gh + wt + ht + bmi + leg + arml + armc + waist +
               tri + sub + age + re, data=w, n.impute=5, pr=FALSE)
g ← fit.mult.impute(gh ~ rcs(age,5) + re + rcs(leg,3) +
                   rcs(waist,4) + tri + rcs(sub,4),
                   orm, a, family=loglog, data=w, pr=FALSE)
```

```
print(g, latex=TRUE, needspace='1.5in')
```

-log-log Ordinal Regression Model

```
fit.mult.impute(formula = gh ~ rcs(age, 5) + re + rcs(leg, 3)
                + rcs(waist, 4) + tri + rcs(sub, 4), fitter = orm,
                xtrans = a, data = w, pr = FALSE, family = loglog)
```

		Model Likelihood Ratio Test		Discrimination Indexes		Rank Discrim. Indexes	
Obs	4629	LR χ^2	1448.42	R^2	0.269	ρ	0.513
Unique Y	63	d.f.	17	g	0.743		
$Y_{0.5}$	5.5	$\Pr(> \chi^2) < 0.0001$		g_r	2.102		
$\max \left \frac{\partial \log L}{\partial \beta} \right $	1×10^{-5}	Score χ^2	1569.21	$\left \Pr(Y \geq Y_{0.5}) - \frac{1}{2} \right $	0.173		
		$\Pr(> \chi^2) < 0.0001$					

	Coef	S.E.	Wald Z	$\Pr(> Z)$
age	0.0404	0.0055	7.29	< 0.0001
age'	-0.0228	0.0279	-0.82	0.4137
age''	0.0126	0.0876	0.14	0.8857
age'''	0.0424	0.1148	0.37	0.7116
re=Other Hispanic	-0.0766	0.0597	-1.28	0.1992
re=Non-Hispanic White	-0.4121	0.0449	-9.17	< 0.0001
re=Non-Hispanic Black	0.0645	0.0566	1.14	0.2543
re=Other Race Including Multi-Racial	-0.0555	0.0750	-0.74	0.4593
leg	-0.0339	0.0091	-3.73	0.0002
leg'	0.0153	0.0105	1.46	0.1434
waist	0.0073	0.0050	1.47	0.1428
waist'	0.0304	0.0158	1.93	0.0536
waist''	-0.0910	0.0508	-1.79	0.0732
tri	-0.0163	0.0026	-6.28	< 0.0001
sub	-0.0027	0.0097	-0.28	0.7817
sub'	0.0674	0.0289	2.33	0.0198
sub''	-0.1895	0.0922	-2.06	0.0398

```
an <- anova(g)
lan(an)
```

	χ^2	d.f.	<i>P</i>
age	692.50	4	< 0.0001
<i>Nonlinear</i>	28.47	3	< 0.0001
re	168.91	4	< 0.0001
leg	24.37	2	< 0.0001
<i>Nonlinear</i>	2.14	1	0.1434
waist	128.31	3	< 0.0001
<i>Nonlinear</i>	4.05	2	0.1318
tri	39.44	1	< 0.0001
sub	39.30	3	< 0.0001
<i>Nonlinear</i>	6.63	2	0.0363
TOTAL NONLINEAR	46.80	8	< 0.0001
TOTAL	1464.24	17	< 0.0001

```
b <- anova(g, leg, waist, tri, sub)
# Add new lines to the plot with combined effect of 4 size var.
s <- rbind(an, size=b['TOTAL', ])
class(s) <- 'anova.rms'
plot(s)
```

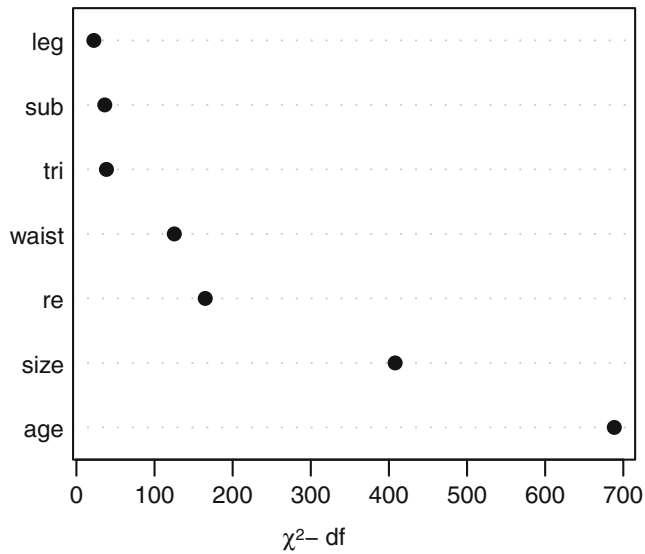


Fig. 15.10 ANOVA for reduced model, after multiple imputation, with addition of a combined effect for four size variables

```
ggplot(Predict(g), abbrev=TRUE, ylab=NULL) # Figure 15.11
```

Compare the estimated age partial effects and confidence intervals with those from a model using casewise deletion, and with bootstrap nonparametric confidence intervals (also with casewise deletion).

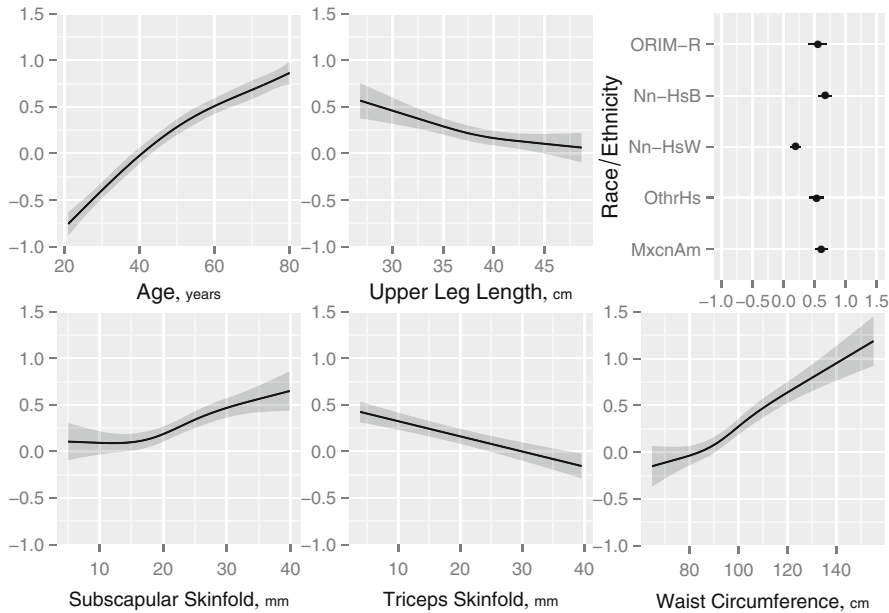


Fig. 15.11 Partial effects (log hazard or log-log cumulative probability scale) of all predictors in reduced model, after multiple imputation

```
gc ← orm(gh ~ rcs(age,5) + re + rcs(leg,3) +
         rcs(waist,4) + tri + rcs(sub,4),
         family=loglog, data=w, x=TRUE, y=TRUE)
gb ← bootcov(gc, B=300)
```

```
bootclb ← Predict(gb, age, boot.type='basic')
bootclp ← Predict(gb, age, boot.type='percentile')
multimp ← Predict(g, age)
plot(Predict(gc, age), addpanel=function(...) {
  with(bootclb, {llines(age, lower, col='blue')
                llines(age, upper, col='blue')})
  with(bootclp, {llines(age, lower, col='blue', lty=2)
                llines(age, upper, col='blue', lty=2)})
  with(multimp, {llines(age, lower, col='red')
                llines(age, upper, col='red')
                llines(age, yhat, col='red')}) },
     col.fill=gray(.9), adj.subtitle=FALSE) # Figure 15.12
```

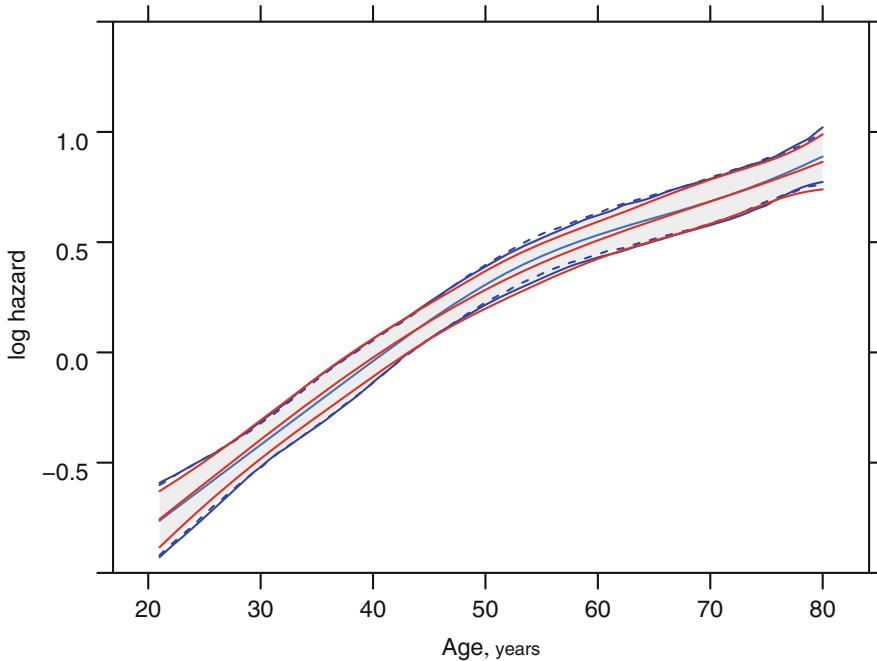


Fig. 15.12 Partial effect for age from multiple imputation (center red line) and casewise deletion (center blue line) with symmetric Wald 0.95 confidence bands using casewise deletion (gray shaded area), basic bootstrap confidence bands using casewise deletion (blue lines), percentile bootstrap confidence bands using casewise deletion (dashed blue lines), and symmetric Wald confidence bands accounting for multiple imputation (red lines).

Figure 15.13 depicts the relationship between various predicted quantities, demonstrating that the ordinal model makes fewer model assumptions that dictate their connections. A Gaussian or log-Gaussian model would have a straight-line relationship between the predicted mean and median.

```
M ← Mean(g)
qu ← Quantile(g)
med ← function(lp) qu(.5, lp)
q90 ← function(lp) qu(.9, lp)
lp ← predict(g)
lpr ← quantile(predict(g), c(.002, .998), na.rm=TRUE)
lps ← seq(lpr[1], lpr[2], length=200)
pmn ← M(lps)
pme ← med(lps)
p90 ← q90(lps)
plot(pmn, pme, # Figure 15.13
      xlab=expression(paste('Predicted Mean ', HbA["1c"])),
      ylab='Median and 0.9 Quantile', type='l',
      xlim=c(4.75, 8.0), ylim=c(4.75, 8.0), bty='n')
box(col=gray(.8))
```

```

lines(pmn, p90, col='blue')
abline(a=0, b=1, col=gray(.8))
text(6.5, 5.5, 'Median')
text(5.5, 6.3, '0.9', col='blue')
nint ← 350
scat1d(M(lp), nint=nint)
scat1d(med(lp), side=2, nint=nint)
scat1d(q90(lp), side=4, col='blue', nint=nint)

```

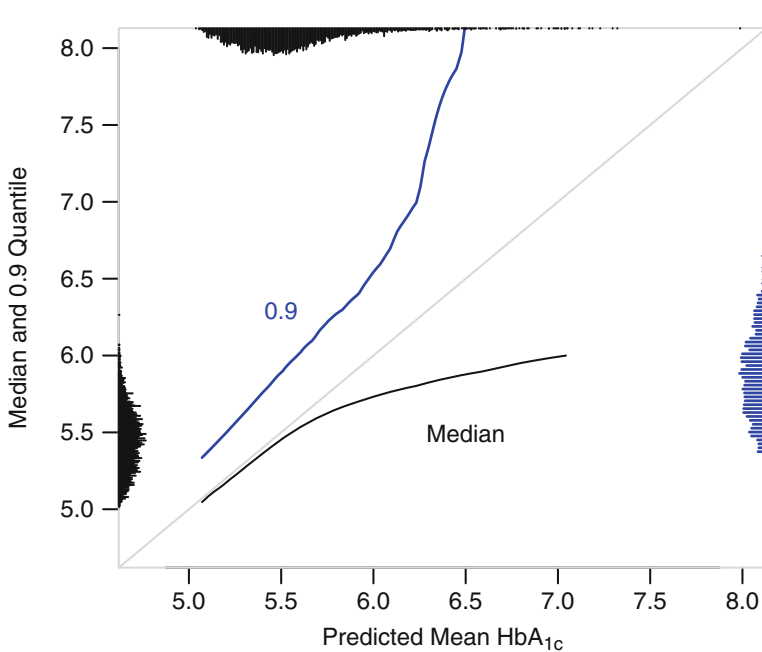


Fig. 15.13 Predicted mean HbA_{1c} vs. predicted median and 0.9 quantile along with their marginal distributions

Finally, let us draw a nomogram that shows the full power of ordinal models, by predicting five quantities of interest.

```

g ← Newlevels(g, list(re=abbreviate(levels(w$re))))
exprob ← ExProb(g)
nom ←
  nomogram(g, fun=list(Mean=M,
    'Median Glycohemoglobin' = med,
    '0.9 Quantile' = q90,
    'Prob(HbA1c ≥ 6.5)' =
      function(x) exprob(x, y=6.5),
    'Prob(HbA1c ≥ 7.0)' =
      function(x) exprob(x, y=7),
    'Prob(HbA1c ≥ 7.5)' =

```

```

function(x) exprob(x, y=7.5)),
fun.at=list(seq(5, 8, by=.5),
c(5,5.25,5.5,5.75,6,6.25),
c(5.5,6,6.5,7,8,10,12,14),
c(.01,.05,.1,.2,.3,.4),
c(.01,.05,.1,.2,.3,.4),
c(.01,.05,.1,.2,.3,.4)))
plot(nom, lmgp=.28) # Figure 15.14

```

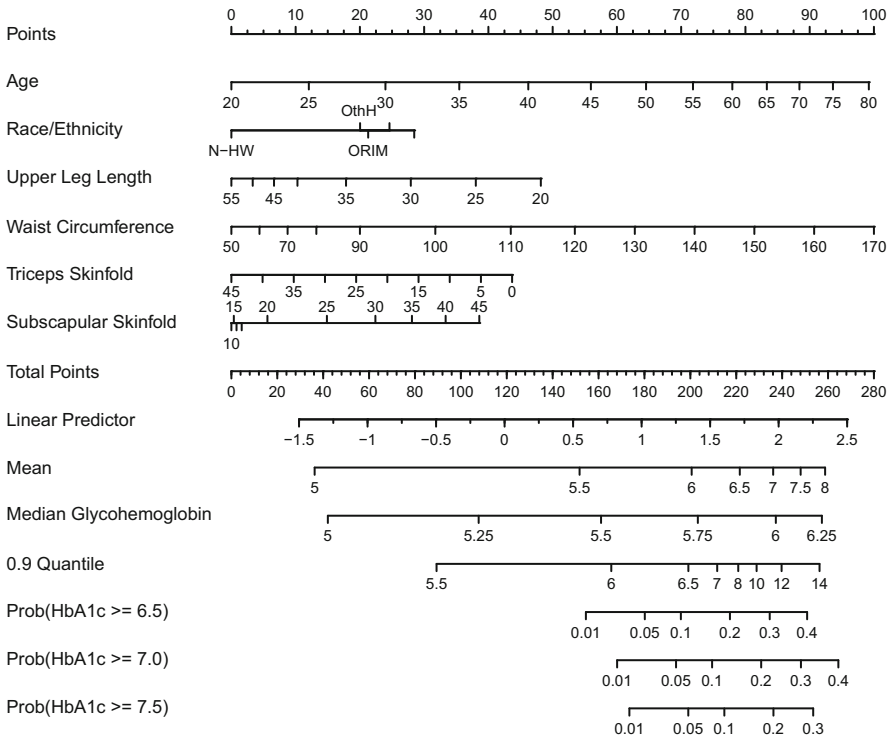


Fig. 15.14 Nomogram for predicting median, mean, and 0.9 quantile of glycohemoglobin, along with the estimated probability that HbA_{1c} ≥ 6.5, 7, or 7.5, all from the log-log ordinal model

Chapter 16

Transform-Both-Sides Regression

16.1 Background

Fitting multiple regression models by the method of least squares is one of the most commonly used methods in statistics. There are a number of challenges to the use of least squares, even when it is only used for estimation and not inference, including the following.

1. How should continuous predictors be transformed so as to get a good fit?
2. Is it better to transform the response variable? How does one find a good transformation that simplifies the right-hand side of the equation?
3. What if Y needs to be transformed non-monotonically (e.g., $|Y - 100|$) before it will have any correlation with X ?

When one is trying to draw an inference about population effects using confidence limits or hypothesis tests, the most common approach is to assume that the residuals have a normal distribution. This is equivalent to assuming that the conditional distribution of the response Y given the set of predictors X is normal with mean depending on X and variance that is (one hopes) a constant independent of X . The need for a distributional assumption to enable us to draw inferences creates a number of other challenges such as the following.

1. If for the untransformed original scale of the response Y the distribution of the residuals is not normal with constant spread, ordinary methods will not yield correct inferences (e.g., confidence intervals will not have the desired coverage probability and the intervals will need to be asymmetric).
2. Quite often there is a transformation of Y that will yield well-behaving residuals. How do you find this transformation? Can you find a transformation for the X s at the same time?

3. All classical statistical inferential methods assume that the full model was pre-specified, that is, the model was not modified after examining the data. How does one correct confidence limits, for example, for data-based model and transformation selection?

16.2 Generalized Additive Models

Hastie and Tibshirani²⁷⁵ have developed *generalized additive models* (GAMs) for a variety of distributions for Y . There are semiparametric GAMs, but most GAMs for continuous Y assume that the conditional distribution of Y is from a specific distribution family. GAMs nicely estimate the transformation each continuous X requires so as to optimize a fitting criterion such as sum of squared errors or log likelihood, subject to the degrees of freedom the analyst desires to spend on each predictor. However, GAMs assume that Y has already been transformed to fit the specified distribution family.

There is excellent software available for fitting a wide variety of GAMs, such as the R packages `gam`, `mgcv`, and `robustgam`.

16.3 Nonparametric Estimation of Y -Transformation

When the model's left-hand side also needs transformation, either to improve R^2 or to achieve constant variance of the residuals (which increases the chances of satisfying a normality assumption), there are a few approaches available. One approach is Breiman and Friedman's *alternating conditional expectation* (ACE) method.⁶⁸ ACE simultaneously transforms both Y and each of the X s so as to maximize the multiple R^2 between the transformed Y and the transformed X s. The model is given by

$$g(Y) = f_1(X_1) + f_2(X_2) + \dots + f_p(X_p). \quad (16.1)$$

ACE allows the analyst to impose restrictions on the transformations such as monotonicity. It allows for categorical predictors, whose categories will automatically be given numeric scores. The transformation for Y is allowed to be non-monotonic. One feature of ACE is its ability to estimate the *maximal correlation* between an X and the response Y . Unlike the ordinary correlation coefficient (which assumes linearity) or Spearman's rank correlation (which assumes monotonicity), the maximal correlation has the property that it is zero if and only if X and Y are statistically independent. This property holds because ACE allows for non-monotonic transformations of all variables. The "super smoother" (see the S `supsmu` function) is the basis for the nonparametric estimation of transformations for continuous X s.

Tibshirani developed a different algorithm for nonparametric additive regression based on least squares, *additivity and variance stabilization* (AVAS).⁶⁰⁷ Unlike ACE, AVAS forces $g(Y)$ to be monotonic. AVAS's fitting criterion is to maximize R^2 while forcing the transformation for Y to result in nearly constant variance of residuals. The model specification is the same as for ACE (Equation 16.3).

ACE and AVAS are powerful fitting algorithms, but they can result in overfitting (R^2 can be greatly inflated when one fits many predictors), and they provide no statistical inferential measures. As discussed earlier, the process of estimating transformations (especially those for Y) can result in significant variance under-estimation, especially for small sample sizes. The bootstrap can be used to correct the apparent R^2 (R_{app}^2) for overfitting. As before, it estimates the optimism (bias) in R_{app}^2 and subtracts this optimism from R_{app}^2 to get a more trustworthy estimate. The bootstrap can also be used to compute confidence limits for all estimated transformations, and confidence limits for estimated predictor effects that take fully into account the uncertainty associated with the transformations. To do this, all steps involved in fitting the additive models must be repeated fresh for each re-sample.

Limited testing has shown that the sample size needs to exceed 100 for ACE and AVAS to provide stable estimates. In small sample sizes the bootstrap bias-corrected estimate of R^2 will be zero because the sample information did not support simultaneous estimation of all transformations.

16.4 Obtaining Estimates on the Original Scale

A common practice in least squares fitting is to attempt to rectify lack of fit by taking parametric transformations of Y before fitting; the logarithm is the most common transformation.^a If after transformation the model's residuals have a population median of zero, the inverse transformation of a predicted transformed value estimates the population median of Y given X . This is because unlike means, quantiles are transformation-preserving. Many analysts make the mistake of not reporting which population parameter is being estimated when inverse transforming $X\hat{\beta}$, and sometimes they even report that the mean is being estimated.

How would one go about estimating the population mean or other parameter on the untransformed scale? If the residuals are assumed to be normally distributed and if $\log(Y)$ is the transformation, the mean of the log-normal distribution, a function of both the mean and the variance of the residuals, can be used to derive the desired quantity. However, if the residuals are not normally distributed, this procedure will not result in the correct estimator.

^a A disadvantage of transform-both-sides regression is this difficulty of interpreting estimates on the original scale. Sometimes the use of a special generalized linear model can allow for a good fit without transforming Y .

Duan¹⁶⁵ developed a “smearing” estimator for more nonparametrically obtaining estimates of parameters on the original scale. In the simple one-sample case without predictors in which one has computed $\hat{\theta} = \sum_{i=1}^n \log(Y_i)/n$, the residuals from this fitted value are given by $e_i = \log(Y_i) - \hat{\theta}$. The smearing estimator of the population mean is $\sum \exp[\hat{\theta} + e_i]/n$. In this simple case the result is the ordinary sample mean \bar{Y} .

The worth of Duan’s smearing estimator is in regression modeling. Suppose that the regression was run on $g(Y)$ from which estimated values $\hat{g}(Y_i) = X_i\hat{\beta}$ and residuals on the transformed scale $e_i = \hat{g}(Y_i) - X_i\hat{\beta}$ were obtained. Instead of restricting ourselves to estimating the population mean, let $W(y_1, y_2, \dots, y_n)$ denote any function of a vector of untransformed response values. To estimate the population mean in the homogeneous one-sample case, W is the simple average of all of its arguments. To estimate the population 0.25 quantile, W is the sample 0.25 quantile of y_1, \dots, y_n . Then the smearing estimator of the population parameter estimated by W given X is $W(g^{-1}(a + e_1), g^{-1}(a + e_2), \dots, g^{-1}(a + e_n))$, where g^{-1} is the inverse of the g transformation and $a = X\hat{\beta}$.

When using the AVAS algorithm, the monotonic transformation g is estimated from the data, and the predicted value of $\hat{g}(Y)$ is given by Equation 16.3. So we extend the smearing estimator as $W(\hat{g}^{-1}(a + e_1), \dots, \hat{g}^{-1}(a + e_n))$, where a is the predicted transformed response given X . As \hat{g} is nonparametric (i.e., a table look-up), the `areg.boot` function described below computes \hat{g}^{-1} using reverse linear interpolation.

If residuals from $\hat{g}(Y)$ are assumed to be symmetrically distributed, their population median is zero and we can estimate the median on the untransformed scale by computing $\hat{g}^{-1}(X\hat{\beta})$. To be safe, `areg.boot` adds the median residual to $X\hat{\beta}$ when estimating the population median (the median residual can be ignored by specifying `statistic='fitted'` to functions that operate on objects created by `areg.boot`).

When quantiles of Y are of major interest, a more direct way to obtain estimates is through the use of quantile regression³⁵⁷. An excellent case study including comparisons with other methods such as Cox regression can be found in Austin et al.³⁸.

16.5 R Functions

The R `acepack` package’s `ace` function implements all the features of the ACE algorithm, and its `avas` function does likewise for AVAS. The bootstrap and smearing capabilities mentioned above are offered for these estimation functions by the `areg.boot` (“additive regression using the bootstrap”) function in the `Hmisc` package. Unlike the `ace` and `avas` functions, `areg.boot` uses the R modeling language, making it easier for the analyst to specify the predic-

tor variables and what is assumed about their relationships with the transformed Y . `areg.boot` also implements a parametric transform-both-sides approach using restricted cubic splines and canonical variates, and offers various estimation options with and without smearing. It can estimate the effect of changing one predictor, holding others constant, using the ordinary bootstrap to estimate the standard deviation of difference in two possibly transformed estimates (for two values of X), assuming normality of such differences. Normality is assumed to avoid generating a large number of bootstrap replications of time-consuming model fits. It would not be very difficult to add non-parametric bootstrap confidence limit capabilities to the software. `areg.boot` re-samples every aspect of the modeling process it uses, just as Faraway¹⁸⁶ did for parametric least squares modeling.

`areg.boot` implements a variety of methods as shown in the simple example below. The `monotone` function restricts a variable's transformation to be monotonic, while the `I` function restricts it to be linear.

```
f ← areg.boot(Y ~ monotone(age) +
              sex + weight + I(blood.pressure))

plot(f)          #show transformations, CLs
Function(f)     #generate S functions
                #defining transformations
predict(f)      #get predictions, smearing estimates
summary(f)      #compute CLs on effects of each X
smearingEst()  #generalized smearing estimators
Mean(f)         #derive S function to
                #compute smearing mean Y
Quantile(f)     #derive function to compute smearing quantile
```

The methods are best described in a case study.

16.6 Case Study

Consider simulated data where the conditional distribution of Y is log-normal given X , but where transform-both-sides regression methods use unlogged Y . Predictor X_1 is linearly related to $\log Y$, X_2 is related by $|X_2 - \frac{1}{2}|$, and categorical X_3 has reference group a effect of zero, group b effect of 0.3, and group c effect of 0.5.

```
require(rms)

set.seed(7)
n ← 400
x1 ← runif(n)
x2 ← runif(n)
x3 ← factor(sample(c('a','b','c'), n, TRUE))
y ← exp(x1 + 2*abs(x2 - .5) + .3*(x3=='b') + .5*(x3=='c') +
        .5*rnorm(n))
```

```
# For reference fit appropriate OLS model
print(ols(log(y) ~ x1 + rcs(x2, 5) + x3), coefs=FALSE,
      latex=TRUE)
```

Linear Regression Model

```
ols(formula = log(y) ~ x1 + rcs(x2, 5) + x3)
```

	Model Likelihood Ratio Test		Discrimination Indexes	
Obs 400	LR χ^2	236.87	R^2	0.447
σ 0.4722	d.f.	7	R^2_{adj}	0.437
d.f. 392	Pr(> χ^2)	0.0000	g	0.482

Residuals
 Min 1Q Median 3Q Max
 -1.346 -0.3075 -0.0134 0.327 1.527

Now fit the `avas` model. We use 300 bootstrap repetitions but only plot the first 20 estimates to see clearly how the bootstrap re-estimates of transformations vary. Had we wanted to restrict transformations to be linear, we would have specified the identity function, for example, `I(x1)`.

```
f ← areg.boot(y ~ x1 + x2 + x3, method='avas', B=300)
```

```
f
```

```
avas Additive Regression Model

areg.boot(x = y ~ x1 + x2 + x3, B = 300, method = "avas")

Predictor Types

  type
x1    s
x2    s
x3    c

y type: s

n= 400   p= 3

Apparent R2 on transformed Y scale: 0.444
Bootstrap validated R2                : 0.42

Coefficients of standardized transformations:

      Intercept           x1           x2           x3
-3.443111e-16  9.702960e-01  1.224320e+00  9.881150e-01

Residuals on transformed scale:
```

	Min	1Q	Median	3Q	Max
	-1.877152e+00	-5.252194e-01	-3.732200e-02	5.339122e-01	2.172680e+00
	Mean	S.D.			
	8.673617e-19	7.420788e-01			

Note that the coefficients above do not mean very much as the scale of the transformations is arbitrary. We see that the model was very slightly overfitted (R^2 dropped from 0.44 to 0.42), and the R^2 are in agreement with the OLS model fit above.

Next we plot the transformations, 0.95 confidence bands, and a sample of the bootstrap estimates.

```
plot(f, boot=20) # Figure 16.1
```

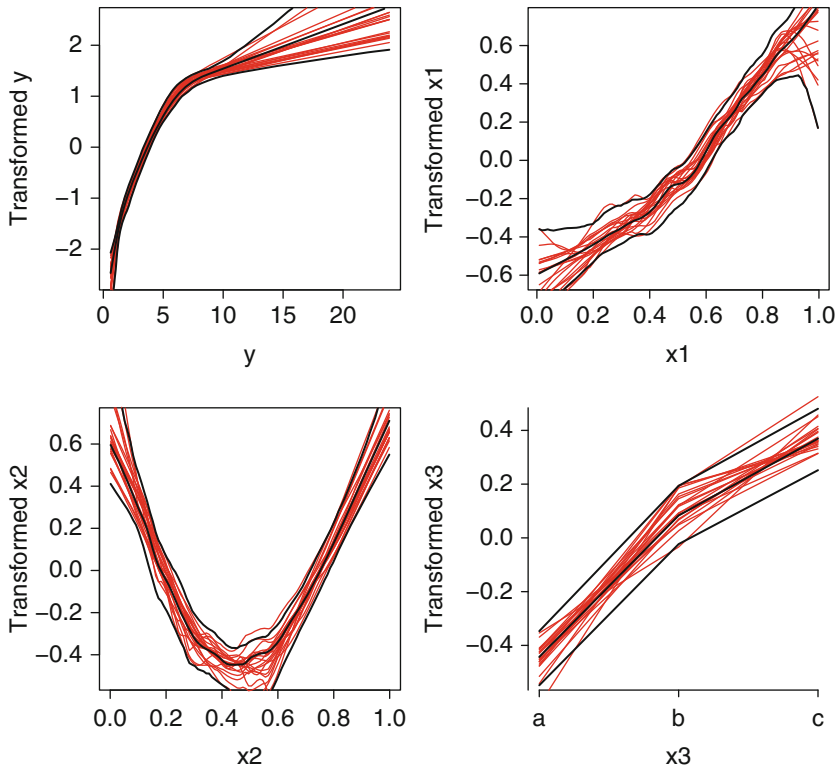


Fig. 16.1 *avas* transformations: overall estimates, pointwise 0.95 confidence bands, and 20 bootstrap estimates (red lines).

The plot is shown in Figure 16.1. The nonparametrically estimated transformation of x_1 is almost linear, and the transformation of x_2 is close to $|x_2 - 0.5|$. We know that the true transformation of y is $\log(y)$, so variance stabilization and normality of residuals will be achieved if the estimated y -transformation is close to $\log(y)$.

```

ys <- seq(.8, 20, length=200)
ytrans <- Function(f)$y # Function outputs all transforms
plot(log(ys), ytrans(ys), type='l') # Figure 16.2
abline(lm(ytrans(ys) ~ log(ys)), col=gray(.8))

```

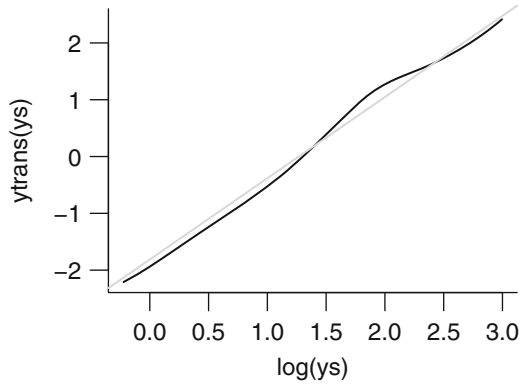


Fig. 16.2 Checking estimated against optimal transformation

Approximate linearity indicates that the estimated transformation is very log-like.^b

Now let us obtain approximate tests of effects of each predictor. `summary` does this by setting all other predictors to reference values (e.g., medians), and comparing predicted responses for a given level of the predictor X with predictions for the lowest setting of X . The default predicted response for `summary` is the median, which is used here. Therefore tests are for differences in medians.

```
summary(f, values=list(x1=c(.2, .8), x2=c(.1, .5)))
```

```
summary.areg.boot(object = f, values = list(x1 = c(0.2, 0.8),
      x2 = c(0.1, 0.5)))
```

Estimates based on 300 resamples

Values to which predictors are set when estimating effects of other predictors:

	y	x1	x2	x3
	3.728843	0.500000	0.300000	2.000000

^b Beware that use of a data-derived transformation in an ordinary model, as this will result in standard errors that are too small. This is because model selection is not taken into account.¹⁸⁶

Estimates of differences of effects on Median Y (from first X value), and bootstrap standard errors of these differences. Settings for X are shown as row headings.

Predictor: x1

x	Differences	S.E	Lower	0.95	Upper	0.95	Z	Pr(Z)
0.2	0.000000	NA	NA	NA	NA	NA	NA	NA
0.8	1.546992	0.2099959	1.135408	1.958577	7.366773	1.747491e-13		

Predictor: x2

x	Differences	S.E	Lower	0.95	Upper	0.95	Z	Pr(Z)
0.1	0.000000	NA	NA	NA	NA	NA	NA	NA
0.5	-1.658961	0.3163361	-2.278968	-1.038953	-5.244298	1.568786e-07		

Predictor: x3

x	Differences	S.E	Lower	0.95	Upper	0.95	Z	Pr(Z)
a	0.000000	NA	NA	NA	NA	NA	NA	NA
b	0.8447422	0.1768244	0.4981728	1.191312	4.777295	1.776692e-06		
c	1.3526151	0.2206395	0.9201697	1.785061	6.130431	8.764127e-10		

For example, when x_1 increases from 0.2 to 0.8 we predict an increase in median y by 1.55 with bootstrap standard error 0.21, when all other predictors are held to constants. Setting them to other constants will yield different estimates of the x_1 effect, as the transformation of y is nonlinear.

Next depict the fitted model by plotting predicted values, with x_2 varying on the x -axis, and three curves corresponding to three values of x_3 . x_1 is set to 0.5. Figure 16.3 shows estimates of both the median and the mean y .

```
newdat <- expand.grid(x2=seq(.05, .95, length=200),
                    x3=c('a', 'b', 'c'), x1=.5,
                    statistic=c('median', 'mean'))
yhat <- c(predict(f, subset(newdat, statistic=='median'),
                  statistic='median'),
         predict(f, subset(newdat, statistic=='mean'),
                  statistic='mean'))
newdat <-
  upData(newdat,
        lp = x1 + 2*abs(x2 - .5) + .3*(x3=='b') +
            .5*(x3=='c'),
        ytrue = ifelse(statistic=='median', exp(lp),
                       exp(lp + 0.5*(0.5^2))), pr=FALSE)
```

```
Input object size: 45472 bytes; 4 variables
Added variable lp
Added variable ytrue
Added variable pr
```

```
New object size: 69800 bytes; 7 variables
```

```
# Use Hmisc function xYplot to produce Figure 16.3
xYplot(yhat ~ x2 | statistic, groups=x3,
       data=newdat, type='l', col=1,
       ylab=expression(hat(y)),
       panel=function(...) {
```

```

panel.xyplot(...)
dat ← subset(newdat,
  statistic==c('median','mean')[current.column()])
for(w in c('a','b','c'))
  with(subset(dat, x3==w),
    llines(x2, ytrue, col=gray(.7), lwd=1.5))
  }
)

```

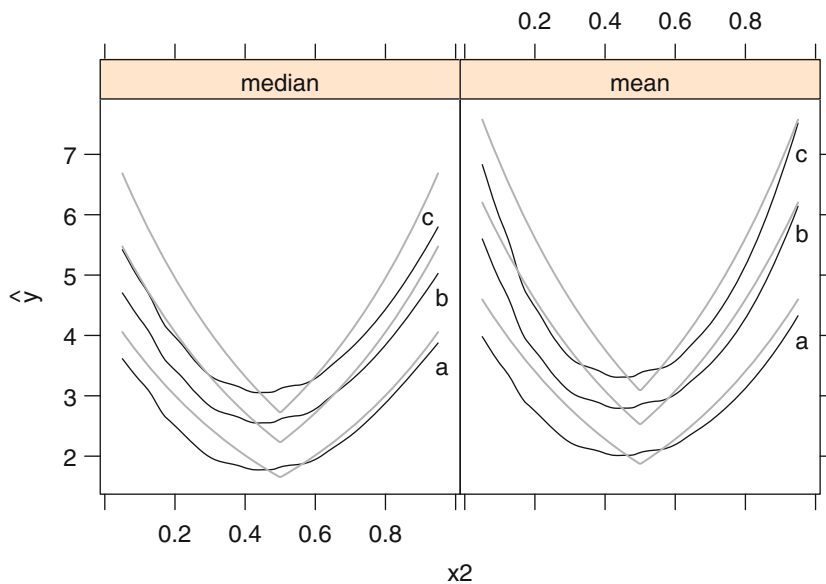


Fig. 16.3 Predicted median (left panel) and mean (right panel) y as a function of x_2 and x_3 . True population values are shown in gray.

Chapter 17

Introduction to Survival Analysis

17.1 Background

Suppose that one wished to study the occurrence of some event in a population of subjects. If the time until the occurrence of the event were unimportant, the event could be analyzed as a binary outcome using the logistic regression model. For example, in analyzing mortality associated with open heart surgery, it may not matter whether a patient dies during the procedure or he dies after being in a coma for two months. For other outcomes, especially those concerned with chronic conditions, the time until the event is important. In a study of emphysema, death at eight years after onset of symptoms is different from death at six months. An analysis that simply counted the number of deaths would be discarding valuable information and sacrificing statistical power.

Survival analysis is used to analyze data in which the time until the event is of interest. The response variable is the time until that event and is often called a *failure time*, *survival time*, or *event time*. Examples of responses of interest include the time until cardiovascular death, time until death or myocardial infarction, time until failure of a light bulb, time until pregnancy, or time until occurrence of an ECG abnormality during exercise. Bull and Spiegelhalter⁸³ have an excellent overview of survival analysis.

The response, event time, is usually continuous, but survival analysis allows the response to be incompletely determined for some subjects. For example, suppose that after a five-year follow-up study of survival after myocardial infarction a patient is still alive. That patient's survival time is *censored* on the right at five years; that is, her survival time is known only to exceed five years. The response value to be used in the analysis is 5+. Censoring can also occur when a subject is lost to follow-up.

If no responses are censored, standard regression models for continuous responses could be used to analyze the failure times by writing the expected failure time as a function of one or more predictors, assuming that

the distribution of failure time is properly specified. However, there are still several reasons for studying failure time using the specialized methods of survival analysis.

1. Time to failure can have an unusual distribution. Failure time is restricted to be positive so it has a skewed distribution and will never be normally distributed.
2. The probability of surviving past a certain time is often more relevant than the expected survival time (and expected survival time may be difficult to estimate if the amount of censoring is large).
3. A function used in survival analysis, the hazard function, helps one to understand the mechanism of failure.³⁰⁸

Survival analysis is used often in industrial life-testing experiments, and it is heavily used in clinical and epidemiologic follow-up studies. Examples include a randomized trial comparing a new drug with placebo for its ability to maintain remission in patients with leukemia, and an observational study of prognostic factors in coronary heart disease. In the latter example subjects may well be followed for varying lengths of time, as they may enter the study over a period of many years.

When regression models are used for survival analysis, all the advantages of these models can be brought to bear in analyzing failure times. Multiple, independent prognostic factors can be analyzed simultaneously and treatment differences can be assessed while adjusting for heterogeneity and imbalances in baseline characteristics. Also, patterns in outcome over time can be predicted for individual subjects.

Even in a simple well-designed experiment, survival modeling can allow one to do the following in addition to making simple comparisons.

1. Test for and describe interactions with treatment. Subgroup analyses can easily generate spurious results and they do not consider interacting factors in a dose-response manner. Once interactions are modeled, relative treatment benefits can be estimated (e.g., hazard ratios), and analyses can be done to determine if some patients are too sick or too well to have even a relative benefit.
2. Understand prognostic factors (strength and shape).
3. Model absolute effect of treatment. First, a model for the probability of surviving past time t is developed. Then differences in survival probabilities for patients on treatments A and B can be estimated. The differences will be due primarily to sickness (overall risk) of the patient and to treatment interactions.
4. Understand time course of treatment effect. The period of maximum effect or period of any substantial effect can be estimated from a plot of relative effects of treatment over time.
5. Gain power for testing treatment effects.
6. Adjust for imbalances in treatment allocation in non-randomized studies.

17.2 Censoring, Delayed Entry, and Truncation

Responses may be left-censored and interval-censored besides being right-censored. *Interval-censoring* is present, for example, when a measuring device functions only for a certain range of the response; measurements outside that range are censored at an end of the scale of the device. Interval-censoring also occurs when the presence of a medical condition is assessed during periodic exams. When the condition is present, the time until the condition developed is only known to be between the current and the previous exam. *Left-censoring* means that an event is known to have occurred before a certain time. In addition, *left-truncation* and *delayed entry* are common. Nomenclature is confusing as many authors refer to delayed entry as left-truncation. Left-truncation really means that an unknown subset of subjects failed before a certain time and the subjects didn't get into the study. For example, one might study the survival patterns of patients who were admitted to a tertiary care hospital. Patients who didn't survive long enough to be referred to the hospital compose the left-truncated group, and interesting questions such as the optimum timing of admission to the hospital cannot be answered from the data set.

Delayed entry occurs in follow-up studies when subjects are exposed to the risk of interest only after varying periods of survival. For example, in a study of occupational exposure to a toxic compound, researchers may be interested in comparing life length of employees with life expectancy in the general population. A subject must live until the beginning of employment before exposure is possible; that is, death cannot be observed before employment. The start of follow-up is delayed until the start of employment and it may be right-censored when follow-up ends. In some studies, a researcher may want to assume that for the purpose of modeling the shape of the hazard function, time zero is the day of diagnosis of disease, while patients enter the study at various times since diagnosis. Delayed entry occurs for patients who don't enter the study until some time after their diagnosis. Patients who die before study entry are left-truncated. Note that the choice of time origin is very important.^{53, 83, 112, 133}

Heart transplant studies have been analyzed by considering time zero to be the time of enrollment in the study. Pre-transplant survival is right-censored at the time of transplant. Transplant survival experience is based on delayed entry into the "risk set" to recognize that a transplant patient is not at risk of dying from transplant failure until after a donor heart is found. In other words, survival experience is not credited to transplant surgery until the day of transplant. Comparisons of transplant experience with medical treatment suffer from "waiting time bias" if transplant survival begins on the day of transplant instead of using delayed entry.^{209, 438, 570}

There are several planned mechanisms by which a response is right-censored. *Fixed type I* censoring occurs when a study is planned to end after two years of follow-up, or when a measuring device will only measure responses up to a certain limit. There the responses are observed only if they

fall below a fixed value C . In *type II censoring*, a study ends when there is a pre-specified number of events. If, for example, 100 mice are followed until 50 die, the censoring time is not known in advance.

We are concerned primarily with *random type I right-censoring* in which each subject's event time is observed only if the event occurs before a certain time, but the censoring time can vary between subjects. Whatever the cause of censoring, we assume that the censoring is *non-informative* about the event; that is, the censoring is caused by something that is independent of the impending failure. Censoring is non-informative when it is caused by planned termination of follow-up or by a subject moving out of town for reasons unrelated to the risk of the event. If subjects are removed from follow-up because of a worsening condition, the *informative censoring* will result in biased estimates and inaccurate statistical inference about the survival experience. For example, if a patient's response is censored because of an adverse effect of a drug or noncompliance to the drug, a serious bias can result if patients with adverse experiences or noncompliance are also at higher risk of suffering the outcome. In such studies, efficacy can only be assessed fairly using the *intention to treat principle*: all events should be attributed to the treatment assigned even if the subject is later removed from that treatment.

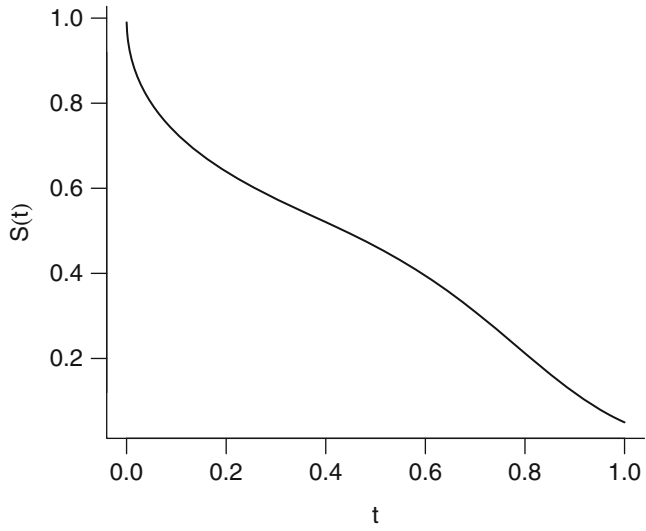
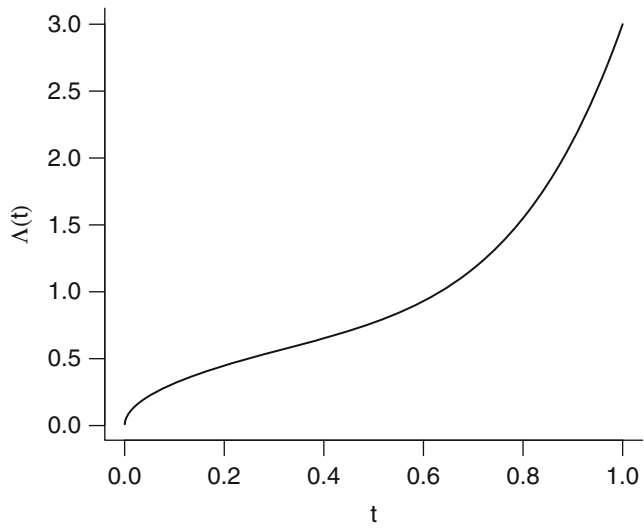
4

17.3 Notation, Survival, and Hazard Functions

In survival analysis we use T to denote the response variable, as the response is usually the time until an event. Instead of defining the statistical model for the response T in terms of the expected failure time, it is advantageous to define it in terms of the *survival function*, $S(t)$, given by

$$S(t) = \text{Prob}\{T > t\} = 1 - F(t), \quad (17.1)$$

where $F(t)$ is the cumulative distribution function for T . If the event is death, $S(t)$ is the probability that death occurs after time t , that is, the probability that the subject will survive at least until time t . $S(t)$ is always 1 at $t = 0$; all subjects survive at least to time zero. The survival function must be non-increasing as t increases. An example of a survival function is shown in Figure 17.1. In that example subjects are at very high risk of the event in the early period so that the $S(t)$ drops sharply. The risk is low for $0.1 \leq t \leq 0.6$, so $S(t)$ is somewhat flat. After $t = .6$ the risk again increases, so $S(t)$ drops more quickly. Figure 17.2 depicts the *cumulative hazard function* corresponding to the survival function in Figure 17.1. This function is denoted by $\Lambda(t)$. It describes the accumulated risk up until time t , and as is shown later, is the negative of the log of the survival function. $\Lambda(t)$ is non-decreasing as t increases; that is, the accumulated risk increases or remains the same. Another important function is the *hazard function*, $\lambda(t)$, also called the *force*

**Fig. 17.1** Survival function**Fig. 17.2** Cumulative hazard function

of mortality, or instantaneous event (death, failure) rate. The hazard at time t is related to the probability that the event will occur in a small interval around t , given that the event has not occurred before time t . By studying the event rate at a given time conditional on the event not having occurred by

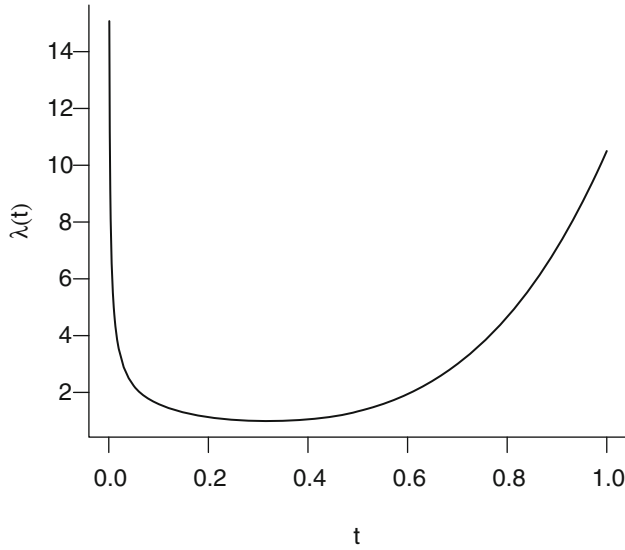


Fig. 17.3 Hazard function

that time, one can learn about the mechanisms and forces of risk over time. Figure 17.3 depicts the hazard function corresponding to $S(t)$ in Figure 17.1 and to $A(t)$ in Figure 17.2. Notice that the hazard function allows one to more easily determine the phases of increased risk than looking for sudden drops in $S(t)$ or $A(t)$.

The hazard function is defined formally by

$$\lambda(t) = \lim_{u \rightarrow 0} \frac{\text{Prob}\{t < T \leq t + u | T > t\}}{u}, \quad (17.2)$$

which using the law of conditional probability becomes

$$\begin{aligned} \lambda(t) &= \lim_{u \rightarrow 0} \frac{\text{Prob}\{t < T \leq t + u\} / \text{Prob}\{T > t\}}{u} \\ &= \lim_{u \rightarrow 0} \frac{[F(t + u) - F(t)] / u}{S(t)} \\ &= \frac{\partial F(t) / \partial t}{S(t)} \\ &= \frac{f(t)}{S(t)}, \end{aligned} \quad (17.3)$$

where $f(t)$ is the probability density function of T evaluated at t , the derivative or slope of the cumulative distribution function $1 - S(t)$. Since

$$\frac{\partial \log S(t)}{\partial t} = \frac{\partial S(t)/\partial t}{S(t)} = -\frac{f(t)}{S(t)}, \quad (17.4)$$

the hazard function can also be expressed as

$$\lambda(t) = -\frac{\partial \log S(t)}{\partial t}, \quad (17.5)$$

the negative of the slope of the log of the survival function. Working backwards, the integral of $\lambda(t)$ is:

$$\int_0^t \lambda(v)dv = -\log S(t). \quad (17.6)$$

The integral or area under $\lambda(t)$ is defined to be $A(t)$, the cumulative hazard function. Therefore

$$A(t) = -\log S(t), \quad (17.7)$$

or

$$S(t) = \exp[-A(t)]. \quad (17.8)$$

So knowing any one of the functions $S(t)$, $A(t)$, or $\lambda(t)$ allows one to derive the other two functions. The three functions are different ways of describing the same distribution.

One property of $A(t)$ is that the expected value of $A(T)$ is unity, since if $T \sim S(t)$, the density of T is $\lambda(t)S(t)$ and

$$\begin{aligned} E[A(T)] &= \int_0^\infty A(t)\lambda(t)\exp(-A(t))dt \\ &= \int_0^\infty u\exp(-u)du \\ &= 1. \end{aligned} \quad (17.9)$$

Now consider properties of the distribution of T . The population q th quantile (100 q th percentile), T_q , is the time by which a fraction q of the subjects will fail. It is the value t such that $S(t) = 1 - q$; that is

$$T_q = S^{-1}(1 - q). \quad (17.10)$$

The median life length is the time by which half the subjects will fail, obtained by setting $S(t) = 0.5$:

$$T_{0.5} = S^{-1}(0.5). \quad (17.11)$$

The q th quantile of T can also be computed by setting $\exp[-A(t)] = 1 - q$, giving

$$\begin{aligned} T_q &= \Lambda^{-1}[-\log(1 - q)] \text{ and as a special case,} \\ T_{.5} &= \Lambda^{-1}(\log 2). \end{aligned} \quad (17.12)$$

The mean or expected value of T (the expected failure time) is the area under the survival function for t ranging from 0 to ∞ :

$$\mu = \int_0^{\infty} S(v)dv. \quad (17.13)$$

Irwin has defined *mean restricted life* (see [334, 335]), which is the area under $S(t)$ up to a fixed time (usually chosen to be a point at which there is still adequate follow-up information).

The random variable T denotes a random failure time from the survival distribution $S(t)$. We need additional notation for the response and censoring information for the i th subject. Let T_i denote the response for the i th subject. This response is the time until the event of interest, and it may be censored if the subject is not followed long enough for the event to be observed. Let C_i denote the censoring time for the i th subject, and define the event indicator as

$$\begin{aligned} e_i &= 1 \text{ if the event was observed } (T_i \leq C_i), \\ &= 0 \text{ if the response was censored } (T_i > C_i). \end{aligned} \quad (17.14)$$

The observed response is

$$Y_i = \min(T_i, C_i), \quad (17.15)$$

which is the time that occurred first: the failure time or the censoring time. The pair of values (Y_i, e_i) contains all the response information for most purposes (i.e., the potential censoring time C_i is not usually of interest if the event occurred before C_i).

Figure 17.4 demonstrates this notation. The line segments start at study entry (survival time $t = 0$).

A useful property of the cumulative hazard function can be derived as follows. Let z be any cutoff time and consider the expected value of Λ evaluated at the earlier of the cutoff time or the actual failure time.

$$\begin{aligned} E[\Lambda(\min(T, z))] &= E[\Lambda(T)[T \leq z] + \Lambda(z)[T > z]] \\ &= E[\Lambda(T)[T \leq z]] + \Lambda(z)S(z). \end{aligned} \quad (17.16)$$

The first term in the right-hand side is

$$\begin{aligned} &\int_0^{\infty} \Lambda(t)[t \leq z]\lambda(t) \exp(-\Lambda(t))dt \\ &= \int_0^z \Lambda(t)\lambda(t) \exp(-\Lambda(t))dt \end{aligned} \quad (17.17)$$

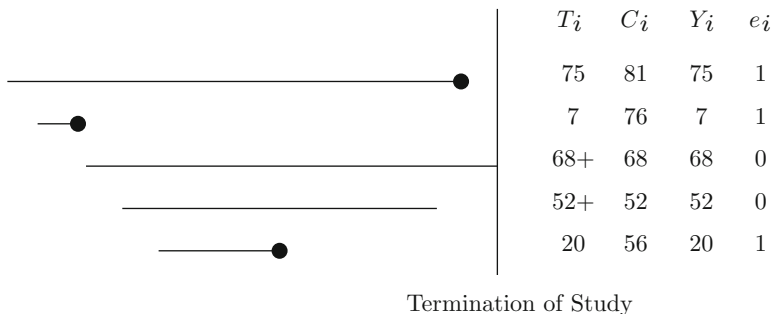


Fig. 17.4 Some censored data. Circles denote events.

$$\begin{aligned}
 &= -[u \exp(-u) + \exp(-u)]_0^{A(z)} \\
 &= 1 - S(z)[A(z) + 1].
 \end{aligned}$$

Adding $A(z)S(z)$ results in

$$E[A(\min(T, z))] = 1 - S(z) = F(z). \tag{17.18}$$

It follows that $\sum_{i=1}^n A(\min(T_i, z))$ estimates the expected number of failures occurring before time z among the n subjects.

5

17.4 Homogeneous Failure Time Distributions

In this section we assume that each subject in the sample has the same distribution of the random variable T that represents the time until the event. In particular, there are no covariables that describe differences between subjects in the distribution of T . As before we use $S(t)$, $\lambda(t)$, and $A(t)$ to denote, respectively, the survival, hazard, and cumulative hazard functions.

The form of the true population survival distribution function $S(t)$ is almost always unknown, and many distributional forms have been used for describing failure time data. We consider first the two most popular parametric survival distributions: the exponential and Weibull distributions. The exponential distribution is a very simple one in which the hazard function is constant; that is, $\lambda(t) = \lambda$. The cumulative hazard and survival functions are then

$$\begin{aligned}
 A(t) &= \lambda t \quad \text{and} \\
 S(t) &= \exp(-A(t)) = \exp(-\lambda t).
 \end{aligned} \tag{17.19}$$

The median life length is $A^{-1}(\log 2)$ or

$$T_{0.5} = \log(2)/\lambda. \quad (17.20)$$

The time by which 1/2 of the subjects will have failed is then proportional to the reciprocal of the constant hazard rate λ . This is true also of the expected or mean life length, which is $1/\lambda$.

The exponential distribution is one of the few distributions for which a closed-form solution exists for the estimator of its parameter when censoring is present. This estimator is a function of the number of events and the total *person-years* of exposure. Methods based on person-years in fact implicitly assume an exponential distribution. The exponential distribution is often used to model events that occur “at random in time.”³²³ It has the property that the future lifetime of a subject is the same, no matter how “old” it is, or

$$\text{Prob}\{T > t_0 + t | T > t_0\} = \text{Prob}\{T > t\}. \quad (17.21)$$

This “ageless” property also makes the exponential distribution a poor choice for modeling human survival except over short time periods.

The Weibull distribution is a generalization of the exponential distribution. Its hazard, cumulative hazard, and survival functions are given by

$$\begin{aligned} \lambda(t) &= \alpha\gamma t^{\gamma-1} \\ A(t) &= \alpha t^\gamma \\ S(t) &= \exp(-\alpha t^\gamma). \end{aligned} \quad (17.22)$$

The Weibull distribution with $\gamma = 1$ is an exponential distribution (with constant hazard). When $\gamma > 1$, its hazard is increasing with t , and when $\gamma < 1$ its hazard is decreasing. Figure 17.5 depicts some of the shapes of the hazard function that are possible. If T has a Weibull distribution, the median of T is

$$T_{0.5} = [(\log 2)/\alpha]^{1/\gamma}. \quad (17.23)$$

There are many other traditional parametric survival distributions, some of which have hazards that are “bathtub shaped” as in Figure 17.3.^{243,323} The restricted cubic spline function described in Section 2.4.5 is an alternative basis for $\lambda(t)$.^{286,287} This function family allows for any shape of smooth $\lambda(t)$ since the number of knots can be increased as needed, subject to the number of events in the sample. Nonlinear terms in the spline function can be tested to assess linearity of hazard (Rayleigh-ness) or constant hazard (exponentiality).

6

The restricted cubic spline hazard model with k knots is

$$\lambda_k(t) = a + bt + \sum_{j=1}^{k-2} \gamma_j w_j(t), \quad (17.24)$$

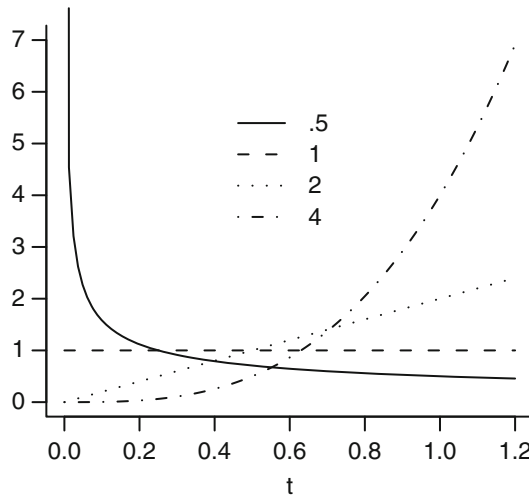


Fig. 17.5 Some Weibull hazard functions with $\alpha = 1$ and various values of γ .

where the $w_j(t)$ are the restricted cubic spline terms of Equation 2.25. There terms are cubic terms in t . A set of knots v_1, \dots, v_k is selected from the quantiles of the uncensored failure times (see Section 2.4.5 and [286]).

The cumulative hazard function for this model is

$$\Lambda(t) = at + \frac{1}{2}t^2 + \frac{1}{4} \times \text{quartic terms in } t. \tag{17.25}$$

Standard maximum likelihood theory is used to obtain estimates of the k unknown parameters to derive, for example, smooth estimates of $\lambda(t)$ with confidence bands. The flexible estimates of $S(t)$ using this method are as efficient as Kaplan–Meier estimates, but they are smooth and can be used as a basis for modeling predictor variables. The spline hazard model is particularly useful for fitting steeply falling and gently rising hazard functions that are characteristic of high-risk medical procedures.

17.5 Nonparametric Estimation of S and Λ

17.5.1 Kaplan–Meier Estimator

As the true form of the survival distribution is seldom known, it is useful to estimate the distribution without making any assumptions. For many analyses, this may be the last step, while in others this step helps one select a statistical model for more in-depth analyses. When no event times are censored, a nonparametric estimator of $S(t)$ is $1 - F_n(t)$ where $F_n(t)$ is the usual

Table 17.1 Kaplan–Meier computations

Day No.	Subjects At Risk	Deaths	Censored	Cumulative Survival
12	100	1	0	$99/100 = .99$
30	99	2	1	$97/99 \times 99/100 = .97$
60	96	0	3	$96/96 \times .97 = .97$
72	93	3	0	$90/93 \times .97 = .94$
.
.

empirical cumulative distribution function based on the observed failure times T_1, \dots, T_n . Let $S_n(t)$ denote this empirical survival function. $S_n(t)$ is given by the fraction of observed failure times that exceed t :

$$S_n(t) = [\text{number of } T_i > t]/n. \quad (17.26)$$

When censoring is present, $S(t)$ can be estimated (at least for t up until the end of follow-up) by the Kaplan–Meier³³³ product-limit estimator. This method is based on conditional probabilities. For example, suppose that every subject has been followed for 39 days or has died within 39 days so that the proportion of subjects surviving at least 39 days can be computed. After 39 days, some subjects may be lost to follow-up besides those removed from follow-up because of death within 39 days. The proportion of those still followed 39 days who survive day 40 is computed. The probability of surviving 40 days from study entry equals the probability of surviving day 40 after living 39 days, multiplied by the chance of surviving 39 days.

The life table in Table 17.1 demonstrates the method in more detail. We suppose that 100 subjects enter the study and none die or are lost before day 12.

Times in a life table should be measured as precisely as possible. If the event being analyzed is death, the failure time should usually be specified to the nearest day. We assume that deaths occur on the day indicated and that being censored on a certain day implies the subject survived through the end of that day. The data used in computing Kaplan–Meier estimates consist of $(Y_i, e_i), i = 1, 2, \dots, n$ using notation defined previously. Primary data collected to derive (Y_i, e_i) usually consist of entry date, event date (if subject failed), and censoring date (if subject did not fail). Instead, the entry date, date of event/censoring, and event/censoring indicator e_i may be specified.

The Kaplan–Meier estimator is called the product-limit estimator because it is the limiting case of actuarial survival estimates as the time periods shrink so that an entry is made for each failure time. An entry need not be in the table for censoring times (when no failures occur at that time) as long as the number of subjects censored is subtracted from the next number

Table 17.2 Summaries used in Kaplan–Meier computations

i	t_i	n_i	d_i	$(n_i - d_i)/n_i$
1	1	7	1	6/7
2	3	6	2	4/6
3	9	2	1	1/2

at risk. Kaplan–Meier estimates are preferred to actuarial estimates because they provide more resolution and make fewer assumptions. In constructing a yearly actuarial life table, for example, it is traditionally assumed that subjects censored between two years were followed 0.5 years.

The product-limit estimator is a nonparametric maximum likelihood estimator [331, pp. 10–13]. The formula for the Kaplan–Meier product-limit estimator of $S(t)$ is as follows. Let k denote the number of failures in the sample and let t_1, t_2, \dots, t_k denote the unique event times (ordered for ease of calculation). Let d_i denote the number of failures at t_i and n_i be the number of subjects *at risk* at time t_i ; that is, $n_i =$ number of failure/censoring times $\geq t_i$. The estimator is then

$$S_{\text{KM}}(t) = \prod_{i:t_i \leq t} (1 - d_i/n_i). \quad (17.27)$$

The Kaplan–Meier estimator of $A(t)$ is $A_{\text{KM}}(t) = -\log S_{\text{KM}}(t)$. An estimate of quantile q of failure time is $S_{\text{KM}}^{-1}(1 - q)$, if follow-up is long enough so that $S_{\text{KM}}(t)$ drops as low as $1 - q$. If the last subject followed failed so that $S_{\text{KM}}(t)$ drops to zero, the expected failure time can be estimated by computing the area under the Kaplan–Meier curve.

To demonstrate computation of $S_{\text{KM}}(t)$, imagine a sample of failure times given by

$$1 \quad 3 \quad 3 \quad 6^+ \quad 8^+ \quad 9 \quad 10^+,$$

where $+$ denotes a censored time. The quantities needed to compute S_{KM} are in Table 17.2. Thus

$$\begin{aligned} S_{\text{KM}}(t) &= 1, & 0 \leq t < 1 \\ &= 6/7 = .85, & 1 \leq t < 3 \\ &= (6/7)(4/6) = .57, & 3 \leq t < 9 \\ &= (6/7)(4/6)(1/2) = .29, & 9 \leq t < 10. \end{aligned} \quad (17.28)$$

Note that the estimate of $S(t)$ is undefined for $t > 10$ since not all subjects have failed by $t = 10$ but no follow-up extends beyond $t = 10$. A graph of the Kaplan–Meier estimate is found in Figure 17.6.

```
require(rms)

tt ← c(1,3,3,6,8,9,10)
stat ← c(1,1,1,0,0,1,0)
S ← Surv(tt, stat)
survplot(npsurv(S ~ 1), conf="bands", n.risk=TRUE,
         xlab=expression(t))
survplot(npsurv(S ~ 1, type="fleming-harrington",
               conf.int=FALSE), add=TRUE, lty=3)
```

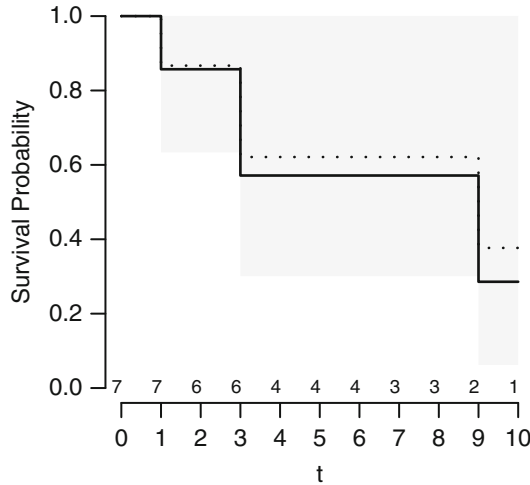


Fig. 17.6 Kaplan–Meier product–limit estimator with 0.95 confidence bands. The Altschuler–Nelson–Aalen–Fleming–Harrington estimator is depicted with the dotted lines.

The variance of $S_{KM}(t)$ can be estimated using Greenwood’s formula [331, p. 14], and using normality of $S_{KM}(t)$ in large samples this variance can be used to derive a confidence interval for $S(t)$. A better method is to derive an asymmetric confidence interval for $S(t)$ based on a symmetric interval for $\log \Lambda(t)$. This latter method ensures that a confidence limit does not exceed one or fall below zero, and is more accurate since $\log \Lambda_{KM}(t)$ is more normally distributed than $S_{KM}(t)$. Once a confidence interval, say $[a, b]$ is determined for $\log \Lambda(t)$, the confidence interval for $S(t)$ is computed by $[\exp\{-\exp(b)\}, \exp\{-\exp(a)\}]$. The formula for an estimate of the variance of interest is [331, p. 15]:

$$\text{Var}\{\log \Lambda_{KM}(t)\} = \frac{\sum_{i:t_i \leq t} d_i / [n_i(n_i - d_i)]}{\{\sum_{i:t_i \leq t} \log[(n_i - d_i)/n_i]\}^2}. \quad (17.29)$$

Letting s denote the square root of this variance estimate, an approximate $1 - \alpha$ confidence interval for $\log \Lambda(t)$ is given by $\log \Lambda_{\text{KM}}(t) \pm zs$, where z is the $1 - \alpha/2$ standard normal critical value. After simplification, the confidence interval for $S(t)$ becomes

$$S_{\text{KM}}(t)^{\exp(\pm zs)}. \quad (17.30)$$

Even though the $\log \Lambda$ basis for confidence limits has theoretical advantages, on the $\log - \log$ scale the estimate of $S(t)$ has the greatest instability where much information is available: when $S(t)$ falls just below 1.0. For that reason, the recommended default confidence limits are on the $\Lambda(t)$ scale using

$$\text{Var}\{\Lambda_{\text{KM}}(t)\} = \sum_{i:t_i \leq t} \frac{d_i}{[n_i(n_i - d_i)]}. \quad (17.31)$$

Letting s denote its square root, an approximate $1 - \alpha$ confidence interval for $S(t)$ is given by

$$\exp(\pm zs)S_{\text{KM}}(t), \quad (17.32)$$

truncated to $[0, 1]$. 7

17.5.2 Altschuler–Nelson Estimator

Altschuler¹⁹, Nelson⁴⁷², Aalen¹ and Fleming and Harrington¹⁹⁶ proposed estimators of $\Lambda(t)$ or of $S(t)$ based on an estimator of $\Lambda(t)$:

$$\begin{aligned} \hat{\Lambda}(t) &= \sum_{i:t_i \leq t} \frac{d_i}{n_i} \\ S_{\Lambda}(t) &= \exp(-\hat{\Lambda}(t)). \end{aligned} \quad (17.33)$$

$S_{\Lambda}(t)$ has advantages over $S_{\text{KM}}(t)$. First, $\sum_{i=1}^n \hat{\Lambda}(Y_i) = \sum_{i=1}^n e_i$ [605, Appendix 3]. In other words, the estimator gives the correct expected number of events. Second, there is a wealth of asymptotic theory based on the Altschuler–Nelson estimator.¹⁹⁶

See Figure 17.6 for an example of the $S_{\Lambda}(t)$ estimator. This estimator has the same variance as $S_{\text{KM}}(t)$ for large enough samples. 8

17.6 Analysis of Multiple Endpoints

Clinical studies frequently assess multiple endpoints. A cancer clinical trial may, for example, involve recurrence of disease and death, whereas a cardiovascular trial may involve nonfatal myocardial infarction and death. Endpoints may be combined, and the new event (e.g., time until infarction or

death) may be analyzed with any of the tools of survival analysis because only the usual censoring mechanism is used. Sometimes the various endpoints may need separate study, however, because they may have different risk factors.

When the multiple endpoints represent multiple causes of a terminating event (e.g., death), Prentice et al. have developed standard methods for analyzing cause-specific hazards⁵¹³ [331, pp. 163–178]. Their methods allow each cause of failure to be analyzed separately, censoring on the other causes. They do not assume any mechanism for cause removal nor make any assumptions regarding the interrelation among causes of failure. However, analyses of competing events using data where some causes of failure are removed in a different way from the original dataset will give rise to different inferences.

When the multiple endpoints represent a mixture of fatal and nonfatal outcomes, the analysis may be more complex. The same is true when one wishes to jointly study an event-time endpoint and a repeated measurement.

9

17.6.1 *Competing Risks*

When events are independent, each event may also be analyzed separately by censoring on all other events as well as censoring on loss to follow-up. This will yield an unbiased estimate of an easily interpreted cause-specific $\lambda(t)$ or $S(t)$ because censoring is non-informative [331, pp. 168–169]. One minus $S_{KM}(t)$ computed in this manner will correctly estimate the probability of failing from the event in the absence of other events. Even when the competing events are not independent, the cause-specific hazard model may lead to valid results, but the resulting model does not allow one to estimate risks conditional on removal of one or more causes of the event. See Kay³⁴⁰ for a nice example of competing risks analysis when a treatment reduces the risk of death from one cause but increases the risk of death from another cause.

10

Larson and Dinse³⁷⁶ have an interesting approach that jointly models the time until (any) failure and the failure type. For r failure types, they use an r -category polytomous logistic model to predict the probability of failing from each cause. They assume that censoring is unrelated to cause of event.

17.6.2 *Competing Dependent Risks*

In many medical and epidemiologic studies one is interested in analyzing multiple causes of death. If the goal is to estimate cause-specific failure probabilities, treating subjects dying from extraneous causes as censored and then computing the ordinary Kaplan–Meier estimate results in biased (high) survival estimates^{212, 225}. If cause m is of interest, the cause-specific hazard

function is defined as

$$\lambda_m(t) = \lim_{u \rightarrow 0} \frac{\Pr\{\text{fail from cause } m \text{ in } [t, t+u) | \text{alive at } t\}}{u}. \quad (17.34)$$

The *cumulative incidence function* or probability of failure from cause m by time t is given by

$$F_m(t) = \int_0^t \lambda_m(u) S(u) du, \quad (17.35)$$

where $S(u)$ is the probability of surviving (ignoring cause of death), which equals $\exp[-\int_0^u (\sum \lambda_m(x)) dx]$ [212]; [444, Chapter 10]; [102, 408]. As previously mentioned, $1 - F_m(t) = \exp[-\int_0^t \lambda_m(u) du]$ only if failures due to other causes are eliminated and if the cause-specific hazard of interest remains unchanged in doing so.²¹²

Again letting t_1, t_2, \dots, t_k denote the unique ordered failure times, a non-parametric estimate of $F_m(t)$ is given by

$$\hat{F}_m(t) = \sum_{i:t_i \leq t} \frac{d_{mi}}{n_i} S_{\text{KM}}(t_{i-1}), \quad (17.36)$$

where d_{mi} is the number of failures of type m at time t_i and n_i is the number of subjects at risk of failure at t_i .

Pepe and others^{494, 496, 497} showed how to use a combination of Kaplan–Meier estimators to derive an estimator of the probability of being free of event 1 by time t given event 2 has not occurred by time t (see also [349]). Let T_1 and T_2 denote, respectively, the times until events 1 and 2. Let $S_1(t)$ and $S_2(t)$ denote, respectively, the two survival functions. Let us suppose that event 1 is not a terminating event (e.g., is not death) and that even after event 1 subjects are followed to ascertain occurrences of event 2. The probability that $T_1 > t$ given $T_2 > t$ is

$$\begin{aligned} \text{Prob}\{T_1 > t | T_2 > t\} &= \frac{\text{Prob}\{T_1 > t \text{ and } T_2 > t\}}{\text{Prob}\{T_2 > t\}} \\ &= \frac{S_{12}(t)}{S_2(t)}, \end{aligned} \quad (17.37)$$

where $S_{12}(t)$ is the survival function for $\min(T_1, T_2)$, the earlier of the two events. Since $S_{12}(t)$ does not involve any informative censoring (assuming as always that loss to follow-up is non-informative), S_{12} may be estimated by the Kaplan–Meier estimator $S_{\text{KM}_{12}}$ (or by S_A). For the type of event 1 we have discussed above, S_2 can also be estimated without bias by S_{KM_2} . Thus we estimate, for example, the probability that a subject still alive at time t will be free of myocardial infarction as of time t by $S_{\text{KM}_{12}}/S_{\text{KM}_2}$.

Another quantity that can easily be computed from ordinary survival estimates is $S_2(t) - S_{12}(t) = [1 - S_{12}(t)] - [1 - S_2(t)]$, which is the probability that event 1 occurs by time t and that event 2 has not occurred by time t .

The ratio estimate above is used to estimate the survival function for one event given that another has not occurred. Another function of interest is the *crude survival function* which is a *marginal* distribution; that is, it is the probability that $T_1 > t$ whether or not event 2 occurs:³⁶²

$$\begin{aligned} S_c(t) &= 1 - F_1(t) \\ F_1(t) &= \text{Prob}\{T_1 \leq t\}, \end{aligned} \tag{17.38}$$

where $F_1(t)$ is the *crude incidence function* defined previously. Note that the $T_1 \leq t$ implies that the occurrence of event 1 is part of the probability being computed. If event 2 is a terminating event so that some subjects can never suffer event 1, the crude survival function for T_1 will never drop to zero. The crude survival function can be interpreted as the survival distribution of W where $W = T_1$ if $T_1 < T_2$ and $W = \infty$ otherwise.³⁶²

11

17.6.3 State Transitions and Multiple Types of Nonfatal Events

In many studies there is one final, absorbing state (death, all causes) and multiple live states. The live states may represent different health states or phases of a disease. For example, subjects may be completely free of cancer, have an isolated tumor, metastasize to a distant organ, and die. Unlike this example, the live states need not have a definite ordering. One may be interested in estimating *transition probabilities*, for example, the probability $\pi_{ij}(t_1, t_2)$ that an individual in state i at time t_1 is in state j after an additional time t_2 . Strauss and Shavelle⁵⁹⁶ have developed an extended Kaplan–Meier estimator for this situation. Let $S_{KM}^i(t|t_1)$ denote the ordinary Kaplan–Meier estimate of the probability of not dying before time t (ignoring distinctions between multiple live states) for a cohort of subjects beginning follow-up at time t_1 in state i . This is an estimate of the probability of surviving an additional t time units (in any live state) given that the subject was alive and in state i at time t_1 . Strauss and Shavelle’s estimator is given by

$$\pi_{ij}(t_1, t_2) = \frac{n_{ij}(t_1, t_2)}{n_i(t_1, t_2)} S_{KM}^i(t_2|t_1), \tag{17.39}$$

where $n_i(t_1, t_2)$ is the number of subjects in live state i at time t_1 who are alive and uncensored t_2 time units later, and $n_{ij}(t_1, t_2)$ is the number of such subjects in state j t_2 time units beyond t_1 .

12

17.6.4 Joint Analysis of Time and Severity of an Event

In some studies, an endpoint is given more weight if it occurs earlier or if it is more severe clinically, or both. For example, the event of interest may be myocardial infarction, which may be of any severity from minimal damage to the left ventricle to a fatal infarction. Berridge and Whitehead⁵² have provided a promising model for the analysis of such endpoints. Their method assumes that the severity of endpoints which do occur is measured on an ordinal categorical scale and that severity is assessed at the time of the event. Berridge and Whitehead's example was time until first headache, with severity of headaches graded on an ordinal scale. They proposed a joint hazard of an individual who responds with ordered category j :

$$\lambda_j(t) = \lambda(t)\pi_j(t), \quad (17.40)$$

where $\lambda(t)$ is the hazard for the failure time and $\pi_j(t)$ is the probability of an individual having event severity j given she fails at time t . Note that a shift in the distribution of response severity is allowed as the time until the event increases.

13

17.6.5 Analysis of Multiple Events

It is common to choose as an endpoint in a clinical trial an event that can recur. Examples include myocardial infarction, gastric ulcer, pregnancy, and infection. Using only the time until the first event can result in a loss of statistical information and power.^a There are specialized multivariate survival models (whose assumptions are extremely difficult to verify) for handling this setup, but in many cases a simpler approach will be efficient.

The simpler approach involves modeling the marginal distribution of the time until each event.^{407, 495} Here one forms one record per subject per event, and the survival time is the time to the first event for the first record, or is the time from the previous event to the next event for all later records. This approach yields consistent estimates of distribution parameters as long as the marginal distributions are correctly specified.⁶⁵⁵ One can allow the number of previous events to influence the hazard function of another event by modeling this count as a covariable.

The multiple events within subject are not independent, so variance estimates must be corrected for intracluster correlation. The clustered sandwich covariance matrix estimator described in Section 9.5 and in [407] will provide

^a An exception to this is the case in which once an event occurs for the first time, that event is likely to recur multiple times for any patient. Then the latter occurrences are redundant.

consistent estimates of variances and covariances even if the events are dependent. Lin⁴⁰⁷ also discussed how this method can easily be used to model multiple events of differing types.

14

17.7 R Functions

The `event.chart` function of Lee et al.³⁹⁴ will draw a variety of charts for displaying raw survival time data, for both single and multiple events per subject. Relationships with covariables can also be displayed. The `event.history` function of Dubin et al.¹⁶⁶ draws an event history graph for right-censored survival data, including time-dependent covariate status. These functions are in the `Hmisc` package.

The analyses described in this chapter can be viewed as special cases of the Cox proportional hazards model.¹³² The programs for Cox model analyses described in Section 20.13 can be used to obtain the results described here, as long as there is at least one stratification factor in the model. There are, however, several R functions that are pertinent to the homogeneous or stratified case. The R function `survfit`, and its particular renditions of the `print`, `plot`, `lines`, and `points` generic functions (all part of the `survival` package written by Terry Therneau), will compute, print, and plot Kaplan–Meier and Nelson survival estimates. Confidence intervals for $S(t)$ may be based on S , A , or $\log A$. The `rms` package’s front-end to the `survival` package’s `survfit` function is `npsurv` for “nonparametric survival”. It and other functions described in later chapters use Therneau’s `Surv` function to combine the response variable and event indicator into a single R “survival time” object. In its simplest form, use `Surv(y, event)`, where `y` is the failure/right-censoring time and `event` is the event/censoring indicator, usually coded T/F, 0 = censored 1 = event or 1 = censored 2 = event. If the event status variable has other coding (e.g., 3 means death), use `Surv(y, s==3)`. To handle interval time-dependent covariables, or to use Andersen and Gill’s *counting process* formulation of the Cox model,²³ use the notation `Surv(tstart, tstop, status)`. The counting process notation allows subjects to enter and leave risk sets at random. For each time interval for each subject, the interval is made up of `tstart < t ≤ tstop`. For time-dependent stratification, there is an optional `origin` argument to `Surv` that indicates the hazard shape time origin at the time of crossover to a new stratum. A `type` argument is used to handle left- and interval-censoring, especially for parametric survival models. Possible values of `type` are "right", "left", "interval", "counting", "interval2", "mstate".

The `Surv` expression will usually be used inside another function, but it is fine to save the result of `Surv` in another object and to use this object in the particular fitting function.

`npsurv` is invoked by the following, with default parameter settings indicated.

```
require(rms)
units(y) ← "Month"
# Default is "Day" - used for axis labels, etc.
npsurv(Surv(y, event) ~ svar1 + svar2 + ... , data, subset,
       type=c("kaplan-meier", "fleming-harrington", "fh2"),
       error=c("greenwood", "tsiatis"), se.fit=TRUE,
       conf.int=.95,
       conf.type=c("log", "log-log", "plain", "none"), ...)
```

If there are no stratification variables (`svar1, ...`), omit them. To print a table of estimates, use

```
f ← npsurv(...)
print(f) # print brief summary of f
summary(f, times, censored=FALSE) # in survival
```

For failure times stored in days, use

```
f ← npsurv(Surv(futime, event) ~ sex)
summary(f, seq(30, 180, by=30))
```

to print monthly estimates.

There is a plot method To plot the object returned by `survfit` and `npsurv`. This invokes `plot.survfit`.

Objects created by `npsurv` can be passed to the more comprehensive plotting function `survplot` (here, actually `survplot.npsurv`) for other options that include automatic curve labeling and showing the number of subjects at risk at selected times. See Figure 17.6 for an example. Stratified estimates, with four treatments distinguished by line type and curve labels, could be drawn by

```
units(y) ← "Year"
f ← npsurv(Surv(y, stat) ~ treatment)
survplot(f, ylab="Fraction Pain-Free")
```

The `groupkm` in `rms` computes and optionally plots $S_{KM}(u)$ or $\log A_{KM}(u)$ (if `loglog=TRUE`) for fixed u with automatic stratification on a continuous predictor x . As in `cut2` (Section 6.2) you can specify the number of subjects per interval (default is `m=50`), the number of quantile groups (`g`), or the actual cut-points (`cuts`). `groupkm` plots the survival or log–log survival estimate against mean x in each x interval.

The `bootkm` function in the `Hmisc` package bootstraps Kaplan–Meier survival estimates or Kaplan–Meier estimates of quantiles of the survival time distribution. It is easy to use `bootkm` to compute, for example, a nonparametric confidence interval for the ratio of median survival times for two groups.

See the Web site for a list of functions from other users for nonparametric estimation of $S(t)$ with left-, right-, and interval-censored data. The adaptive linear spline log-hazard fitting function `heft`³⁶¹ is freely available.

17.8 Further Reading

- [1] Some excellent general references for survival analysis are [57, 83, 114, 133, 154, 197, 282, 308, 331, 350, 382, 392, 444, 484, 574, 604]. Govindarajulu et al.²²⁹ have a nice review of frailty models in survival analysis, for handling clustered time-to-event data.
- [2] See Goldman,²²⁰ Bull and Spiegelhalter,⁸³ Lee et al.³⁹⁴, and Dubin et al.¹⁶⁶ for ways to construct descriptive graphs depicting right-censored data.
- [3] Some useful references for left-truncation are [83, 112, 244, 524]. Mandel⁴³⁵ carefully described the difference between censoring and truncation.
- [4] See [384, p. 164] for some ideas for detecting informative censoring. Bilker and Wang⁵⁴ discuss *right-truncation* and contrast it with right-censoring.
- [5] Arjas²⁹ has applications based on properties of the cumulative hazard function.
- [6] Kooperberg et al.^{361, 594} have an adaptive method for fitting hazard functions using linear splines in the log hazard. Binquet et al.⁵⁶ studied a related approach using quadratic splines. Mudholkar et al.⁴⁶⁶ presented a generalized Weibull model allowing for a variety of hazard shapes.
- [7] Hollander et al.²⁹⁹ provide a nonparametric *simultaneous* confidence band for $S(t)$, surprisingly using likelihood ratio methods. Miller⁴⁵⁹ showed that if the parametric form of $S(t)$ is known to be Weibull with known shape parameter (an unlikely scenario), the Kaplan–Meier estimator is very inefficient (i.e., has high variance) when compared with the parametric maximum likelihood estimator. See [666] for a discussion of how the efficiency of Kaplan–Meier estimators can be improved by interpolation as opposed to piecewise flat step functions. That paper also discusses a variety of other estimators, some of which are significantly more efficient than Kaplan–Meier.
- [8] See [112, 244, 438, 570, 614, 619] for methods of estimating S or A in the presence of left-truncation. See Turnbull⁶¹⁶ for nonparametric estimation of $S(t)$ with left-, right-, and interval-censoring, and Kooperberg and Clarkson³⁶⁰ for a flexible parametric approach to modeling that allows for interval-censoring. Lindsey and Ryan⁴¹³ have a nice tutorial on the analysis of interval-censored data.
- [9] Hogan and Laird^{297, 298} developed methods for dealing with mixtures of fatal and nonfatal outcomes, including some ideas for handling outcome-related dropouts on the repeated measurements. See also Finkelstein and Schoenfeld.¹⁹³ The 30 April 1997 issue of *Statistics in Medicine* (Vol. 16) is devoted to methods for analyzing multiple endpoints as well as designing multiple endpoint studies. The papers in that issue are invaluable, as is Therneau and Hamilton⁶⁰⁶ and Therneau and Grambsch.⁶⁰⁴ Huang and Wang³¹¹ presented a joint model for recurrent events and a terminating event, addressing such issues as the frequency of recurrent events by the time of the terminating event.
- [10] See Lunn and McNeil⁴²⁹ and Marubini and Valsecchi [444, Chapter 10] for practical approaches to analyzing competing risks using ordinary Cox proportional hazards models. A nice overview of competing risks with comparisons of various approaches is found in Tai et al.⁵⁹⁹, Geskus²¹⁴, and Koller et al.³⁵⁸. Bryant and Dignam⁷⁸ developed a semiparametric procedure in which competing risks are adjusted for nonparametrically while a parametric cumulative incidence function is used for the event of interest, to gain precision. Fine and Gray¹⁹² developed methods for analyzing competing risks by estimating sub-distribution functions. Nishikawa et al.⁴⁷⁸ developed some novel approaches to competing risk analysis involving time to adverse drug events competing with time to withdrawal from therapy. They also dealt with different severities of events in an interesting way. Putter et al.⁵¹⁷ has a nice tutorial on competing risks, multi-state models, and associated R software. Fiocco et al.¹⁹⁴ developed

an approach to avoid the problems caused by having to estimate a large number of regression coefficients in multi-state models. Ambrogi et al.²² provide clinically useful estimates from competing risks analyses.

- [11] Jiang, Chappell, and Fine³²² present methods for estimating the distribution of event times of nonfatal events in the presence of terminating events such as death.
- [12] Shen and Thall⁵⁶⁸ have developed a flexible parametric approach to multi-state survival analysis.
- [13] Lancar et al.³⁷² developed a method for analyzing repeated events of varying severities.
- [14] Lawless and Nadeau³⁸⁴ have a very good description of models dealing with recurrent events. They use the notion of the *cumulative mean function*, which is the expected number of events experienced by a subject by a certain time. Lawless³⁸³ contrasts this approach with other approaches. See Aalen et al.³ for a nice example in which multivariate failure times (time to failure of fillings in multiple teeth per subject) are analyzed. Francis and Fuller²⁰⁴ developed a graphical device for depicting complex event history data. Therneau and Hamilton⁶⁰⁶ have very informative comparisons of various methods for modeling multiple events, showing the importance of whether the analyst starts the clock over after each event. Kelly and Lim³⁴³ have another very useful paper comparing various methods for analyzing recurrent events. Wang and Chang⁶⁵⁰ demonstrated the difficulty of using Kaplan–Meier estimates for recurrence time data.

17.9 Problems

1. Make a rough drawing of a hazard function from birth for a man who develops significant coronary artery disease at age 50 and undergoes coronary artery bypass surgery at age 55.
2. Define in words the relationship between the hazard function and the survival function.
3. In a study of the life expectancy of light bulbs as a function of the bulb's wattage, 100 bulbs of various wattage ratings were tested until each had failed. What is wrong with using the product-moment linear correlation test to test whether wattage is associated with life length concerning (a) distributional assumptions and (b) other assumptions?
4. A placebo-controlled study is undertaken to ascertain whether a new drug decreases mortality. During the study, some subjects are withdrawn because of moderate to severe side effects. Assessment of side effects and withdrawal of patients is done on a blinded basis. What statistical technique can be used to obtain an unbiased treatment comparison of survival times? State at least one efficacy endpoint that can be analyzed unbiasedly.
5. Consider long-term follow-up of patients in the **support** dataset. What proportion of the patients have censored survival times? Does this imply that one cannot make accurate estimates of chances of survival? Make a histogram or empirical distribution function estimate of the *censored* follow-up times. What is the typical follow-up duration for a patient in the study

who has survived so far? What is the typical survival time for patients who have died? Taking censoring into account, what is the median survival time from the Kaplan–Meier estimate of the overall survival function? Estimate the median graphically or using any other sensible method.

6. Plot Kaplan–Meier survival function estimates stratified by `dzclass`. Estimate the median survival time and the first quartile of time until death for each of the four disease classes.
7. Repeat Problem 6 except for tertiles of `meanbp`.
8. The commonly used log-rank test for comparing survival times between groups of patients is a special case of the test of association between the grouping variable and survival time in a Cox proportional hazards regression model. Depending on how one handles tied failure times, the log-rank χ^2 statistic exactly equals the score χ^2 statistic from the Cox model, and the likelihood ratio and Wald χ^2 test statistics are also appropriate. To obtain global score or LR χ^2 tests and P -values you can use a statement as the following, where `cph` is in the `rms` package. It is similar to the `survival` package's `coxph` function.

```
cph(Survobject ~ predictor)
```

Here `Survobject` is a survival time object created by the `Surv` function. Obtain the log-rank (score) χ^2 statistic, degrees of freedom, and P -value for testing for differences in survival time between levels of `dzclass`. Interpret this test, referring to the graph you produced in Problem 6 if needed.

9. Do preliminary analyses of survival time using the Mayo Clinic primary biliary cirrhosis dataset described in Section 8.9. Make graphs of Altshuler–Nelson or Kaplan–Meier survival estimates stratified separately by a few categorical predictors and by categorized versions of one or two continuous predictors. Estimate median failure time for the various strata. You may want to suppress confidence bands when showing multiple strata on one graph. See [361] for parametric fits to the survival and hazard function for this dataset.

Chapter 18

Parametric Survival Models

18.1 Homogeneous Models (No Predictors)

The nonparametric estimator of $S(t)$ is a very good descriptive statistic for displaying survival data. For many purposes, however, one may want to make more assumptions to allow the data to be modeled in more detail. By specifying a functional form for $S(t)$ and estimating any unknown parameters in this function, one can

1. easily compute selected quantiles of the survival distribution;
2. estimate (usually by extrapolation) the expected failure time;
3. derive a concise equation and smooth function for estimating $S(t)$, $\Lambda(t)$, and $\lambda(t)$; and
4. estimate $S(t)$ more precisely than $S_{KM}(t)$ or $S_{\Lambda}(t)$ if the parametric form is correctly specified.

18.1.1 Specific Models

Parametric modeling requires choosing one or more distributions. The Weibull and exponential distributions were discussed in Chapter 18. Other commonly used survival distributions are obtained by transforming T and using a standard distribution. The log transformation is most commonly employed. The *log-normal* distribution specifies that $\log(T)$ has a normal distribution with mean μ and variance σ^2 . Stated another way, $\log(T) \sim \mu + \sigma\epsilon$, where ϵ has a standard normal distribution. Then $S(t) = 1 - \Phi((\log(t) - \mu)/\sigma)$, where Φ is the standard normal cumulative distribution function. The *log-logistic* distribution is given by $S(t) = [1 + \exp(-(\log(t) - \mu)/\sigma)]^{-1}$. Here $\log(T) \sim \mu + \sigma\epsilon$ where ϵ follows a logistic distribution $[1 + \exp(-u)]^{-1}$. The *log*

extreme value distribution is given by $S(t) = \exp[-\exp((\log(t) - \mu)/\sigma)]$, and $\log(T) \sim \mu + \sigma\epsilon$, where $\epsilon \sim 1 - \exp[-\exp(u)]$.

The generalized gamma and generalized F distributions provide a richer variety of distribution and hazard functions^{127,128}. Spline hazard models^{286,287,361} are other excellent alternatives.

18.1.2 Estimation

Maximum likelihood (ML) estimation is used to estimate the unknown parameters of $S(t)$. The general method presented in Chapter 9 must be augmented, however, to allow for censored failure times. The basic idea is as follows. Again let T be a random variable representing time until the event, T_i be the (possibly censored) failure time for the i th observation, and Y_i denote the observed failure or censoring time $\min(T_i, C_i)$, where C_i is the censoring time. If Y_i is uncensored, observation i contributes a factor to the likelihood equal to the density function for T evaluated at Y_i , $f(Y_i)$. If Y_i instead represents a censored time so that $T_i = Y_i^+$, it is only known that T_i exceeds Y_i . The contribution to the likelihood function is the probability that $T_i > C_i$ (equal to $\text{Prob}\{T_i > Y_i\}$). This probability is $S(Y_i)$. The joint likelihood over all observations $i = 1, 2, \dots, n$ is

$$L = \prod_{i:Y_i \text{ uncensored}}^n f(Y_i) \prod_{i:Y_i \text{ censored}}^n S(Y_i). \quad (18.1)$$

There is one more component to L : the distribution of censoring times if these are not fixed in advance. Recall that we assume that censoring is non-informative, that is, it is independent of the risk of the event. This independence implies that the likelihood component of the censoring distribution simply multiplies L and that the censoring distribution contains little information about the survival distribution. In addition, the censoring distribution may be very difficult to specify. For these reasons we can maximize L separately to estimate parameters of $S(t)$ and ignore the censoring distribution.

Recalling that $f(t) = \lambda(t)S(t)$ and $\Lambda(t) = -\log S(t)$, the log likelihood can be written as

$$\log L = \sum_{i:Y_i \text{ uncensored}}^n \log \lambda(Y_i) - \sum_{i=1}^n \Lambda(Y_i). \quad (18.2)$$

All observations then contribute an amount to the log likelihood equal to the negative of the cumulative hazard evaluated at the failure/censoring time. In addition, uncensored observations contribute an amount equal to the log of the hazard function evaluated at the time of failure. Once L or $\log L$ is specified, the general ML methods outlined earlier can be used without

change in most situations. The principal difference is that censored observations contribute less information to the statistical inference than uncensored observations. For distributions such as the log-normal that are written only in terms of $S(t)$, it may be easier to write the likelihood in terms of $S(t)$ and $f(t)$.

As an example, we turn to the exponential distribution, for which $\log L$ has a simple form that can be maximized explicitly. Recall that for this distribution $\lambda(t) = \lambda$ and $A(t) = \lambda t$. Therefore,

$$\log L = \sum_{i: Y_i \text{ uncensored}}^n \log \lambda - \sum_{i=1}^n \lambda Y_i. \quad (18.3)$$

Letting n_u denote the number of uncensored event times,

$$\log L = n_u \log \lambda - \sum_{i=1}^n \lambda Y_i. \quad (18.4)$$

Letting w denote the sum of all failure/censoring times (“person years of exposure”):

$$w = \sum_{i=1}^n Y_i, \quad (18.5)$$

the derivatives of $\log L$ are given by

$$\begin{aligned} \frac{\partial \log L}{\partial \lambda} &= n_u / \lambda - w \\ \frac{\partial^2 \log L}{\partial \lambda^2} &= -n_u / \lambda^2. \end{aligned} \quad (18.6)$$

Equating the derivative of $\log L$ to zero implies that the MLE of λ is

$$\hat{\lambda} = n_u / w \quad (18.7)$$

or the number of failures per person-years of exposure. By inserting the MLE of λ into the formula for the second derivative we obtain the observed estimated information, w^2/n_u . The estimated variance of $\hat{\lambda}$ is thus n_u/w^2 and the standard error is $n_u^{1/2}/w$. The precision of the estimate depends primarily on n_u .

Recall that the expected life length μ is $1/\lambda$ for the exponential distribution. The MLE of μ is w/n_u and its estimated variance is w^2/n_u^3 . The MLE of $S(t)$, $\hat{S}(t)$, is $\exp(-\hat{\lambda}t)$, and the estimated variance of $\log(\hat{A}(t))$ is simply $1/n_u$.

As an example, consider the sample listed previously,

$$1 \quad 3 \quad 3 \quad 6^+ \quad 8^+ \quad 9 \quad 10^+.$$

Here $n_u = 4$ and $w = 40$, so the MLE of λ is 0.1 failure per person-period. The estimated standard error is $2/40 = 0.05$. Estimated expected life length is 10 units with a standard error of 5 units. Estimated median failure time is $\log(2)/0.1 = 6.931$. The estimated survival function is $\exp(-0.1t)$, which at $t = 1, 3, 9, 10$ yields 0.90, 0.74, 0.41, and 0.37, which can be compared to the product limit estimates listed earlier (0.85, 0.57, 0.29, 0.29).

Now consider the Weibull distribution. The log likelihood function is

$$\log L = \sum_{i:Y_i \text{ uncensored}}^n \log[\alpha\gamma Y_i^{\gamma-1}] - \sum_{i=1}^n \alpha Y_i^\gamma. \quad (18.8)$$

Although $\log L$ can be simplified somewhat, it cannot be solved explicitly for α and γ . An iterative method such as the Newton–Raphson method is used to compute the MLEs of α and γ . Once these estimates are obtained, the estimated variance–covariance matrix and other derived quantities such as $\hat{S}(t)$ can be obtained in the usual manner.

For the dataset used in the exponential fit, the Weibull fit follows.

$$\begin{aligned} \hat{\alpha} &= 0.0728 \\ \hat{\gamma} &= 1.164 \\ \hat{S}(t) &= \exp(-0.0728t^{1.164}) \\ \hat{S}^{-1}(0.5) &= [(\log 2)/\hat{\alpha}]^{1/\hat{\gamma}} = 6.935 \text{ (estimated median)}. \end{aligned} \quad (18.9)$$

This fit is very close to the exponential fit since $\hat{\gamma}$ is near 1.0. Note that the two medians are almost equal. The predicted survival probabilities for the Weibull model for $t = 1, 3, 9, 10$ are, respectively, 0.93, 0.77, 0.39, 0.35.

Sometimes a formal test can be made to assess the fit of the proposed parametric survival distribution. For the data just analyzed, a formal test of exponentiality versus a Weibull alternative is obtained by testing $H_0 : \gamma = 1$ in the Weibull model. A score test yielded $\chi^2 = 0.14$ with 1 d.f., $p = 0.7$, showing little evidence for non-exponentiality (note that the sample size is too small for this test to have any power).

18.1.3 Assessment of Model Fit

The fit of the hypothesized survival distribution can often be checked easily using graphical methods. Nonparametric estimates of $S(t)$ and $A(t)$ are primary tools for this purpose. For example, the Weibull distribution $S(t) = \exp(-\alpha t^\gamma)$ can be rewritten by taking logarithms twice:

$$\log[-\log S(t)] = \log A(t) = \log \alpha + \gamma(\log t). \quad (18.10)$$

The fit of a Weibull model can be assessed by plotting $\log \hat{\Lambda}(t)$ versus $\log t$ and checking whether the curve is approximately linear. Also, the plotted curve provides approximate estimates of α (the antilog of the intercept) and γ (the slope). Since an exponential distribution is a special case of a Weibull distribution when $\gamma = 1$, exponentially distributed data will tend to have a graph that is linear with a slope of 1.

For any assumed distribution $S(t)$, a graphical assessment of goodness of fit can be made by plotting $S^{-1}[S_A(t)]$ or $S^{-1}[S_{\text{KM}}(t)]$ against t and checking for linearity. For log distributions, S specifies the distribution of $\log(T)$, so we plot against $\log t$. For a log-normal distribution we thus plot $\Phi^{-1}[S_A(t)]$ against $\log t$, where Φ^{-1} is the inverse of the standard normal cumulative distribution function. For a log-logistic distribution we plot $\text{logit}[S_A(t)]$ versus $\log t$. For an extreme value distribution we use log–log plots as with the Weibull distribution. Parametric model fits can also be checked by plotting the fitted $\hat{S}(t)$ and $S_A(t)$ against t on the same graph.

18.2 Parametric Proportional Hazards Models

In this section we present one way to generalize the survival model to a survival regression model. In other words, we allow the sample to be heterogeneous by adding predictor variables $X = \{X_1, X_2, \dots, X_k\}$. As with other regression models, X can represent a mixture of binary, polytomous, continuous, spline-expanded, and even ordinal predictors (if the categories are scored to satisfy the linearity assumption). Before discussing ways in which the regression part of a survival model might be specified, first recall how regression effects have been modeled in other settings. In multiple linear regression, the regression effect $X\beta = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_k X_k$ can be thought of as an increment in the expected value of the response Y . In binary logistic regression, $X\beta$ specifies the log odds that $Y = 1$, or $\exp(X\beta)$ multiplies the odds that $Y = 1$.

18.2.1 Model

The most widely used survival regression specification is to allow the hazard function $\lambda(t)$ to be multiplied by $\exp(X\beta)$. The survival model is thus generalized from a hazard function $\lambda(t)$ for the failure time T to a hazard function $\lambda(t)\exp(X\beta)$ for the failure time given the predictors X :

$$\lambda(t|X) = \lambda(t)\exp(X\beta). \quad (18.11)$$

This regression formulation is called the *proportional hazards (PH)* model. The $\lambda(t)$ part of $\lambda(t|X)$ is sometimes called an *underlying hazard function* or a *hazard function for a standard subject*, which is a subject with $X\beta = 0$. Any parametric hazard function can be used for $\lambda(t)$, and as we show later, $\lambda(t)$ can be left completely unspecified without sacrificing the ability to estimate β , by the use of Cox's semi-parametric PH model.¹³² Depending on whether the underlying hazard function $\lambda(t)$ has a constant scale parameter, $X\beta$ may or may not include an intercept β_0 . The term $\exp(X\beta)$ can be called a *relative hazard function* and in many cases it is the function of primary interest as it describes the (relative) effects of the predictors.

The PH model can also be written in terms of the cumulative hazard and survival functions:

$$\begin{aligned} \Lambda(t|X) &= \Lambda(t) \exp(X\beta) \\ S(t|X) &= \exp[-\Lambda(t) \exp(X\beta)] = \exp[-\Lambda(t)]^{\exp(X\beta)}. \end{aligned} \quad (18.12)$$

$\Lambda(t)$ is an “underlying” cumulative hazard function. $S(t|X)$, the probability of surviving past time t given the values of the predictors X , can also be written as

$$S(t|X) = S(t)^{\exp(X\beta)}, \quad (18.13)$$

where $S(t)$ is the “underlying” survival distribution, $\exp(-\Lambda(t))$. The effect of the predictors is to multiply the hazard and cumulative hazard functions by a factor $\exp(X\beta)$, or equivalently to raise the survival function to a power equal to $\exp(X\beta)$.

18.2.2 Model Assumptions and Interpretation of Parameters

In the general regression notation of Section 2.2, the log hazard or log cumulative hazard can be used as the property of the response T evaluated at time t that allows distributional and regression parts to be isolated and checked. The PH model can be linearized with respect to $X\beta$ using the following identities.

$$\begin{aligned} \log \lambda(t|X) &= \log \lambda(t) + X\beta \\ \log \Lambda(t|X) &= \log \Lambda(t) + X\beta. \end{aligned} \quad (18.14)$$

No matter which of the three model statements are used, there are certain assumptions in a parametric PH survival model. These assumptions are listed below.

1. The true form of the underlying functions (λ , Λ , and S) should be specified correctly.

2. The relationship between the predictors and log hazard or log cumulative hazard should be linear in its simplest form. In the absence of interaction terms, the predictors should also operate additively.
3. The way in which the predictors affect the distribution of the response should be by multiplying the hazard or cumulative hazard by $\exp(X\beta)$ or equivalently by adding $X\beta$ to the log hazard or log cumulative hazard at each t . The effect of the predictors is assumed to be the same at all values of t since $\log \lambda(t)$ can be separated from $X\beta$. In other words, the PH assumption implies no t by predictor interaction.

The regression coefficient for X_j , β_j , is the increase in log hazard or log cumulative hazard at any fixed point in time if X_j is increased by one unit and all other predictors are held constant. This can be written formally as

$$\beta_j = \log \lambda(t|X_1, X_2, \dots, X_j + 1, X_{j+1}, \dots, X_k) - \log \lambda(t|X_1, \dots, X_j, \dots, X_k), \quad (18.15)$$

which is equivalent to the log of the ratio of the hazards at time t . The regression coefficient can just as easily be written in terms of a ratio of hazards at time t . The ratio of hazards at $X_j + d$ versus X_j , all other factors held constant, is $\exp(\beta_j d)$. Thus the effect of increasing X_j by d is to increase the hazard of the event by a factor of $\exp(\beta_j d)$ at all points in time, assuming X_j is linearly related to $\log \lambda(t)$. In general, the ratio of hazards for an individual with predictor variable values X^* compared to an individual with predictors X is

$$\begin{aligned} X^* : X \text{ hazard ratio} &= [\lambda(t) \exp(X^* \beta)] / [\lambda(t) \exp(X \beta)] \\ &= \exp(X^* \beta) / \exp(X \beta) = \exp[(X^* - X) \beta]. \end{aligned} \quad (18.16)$$

If there is only one predictor X_1 and that predictor is binary, the PH model can be written

$$\begin{aligned} \lambda(t|X_1 = 0) &= \lambda(t) \\ \lambda(t|X_1 = 1) &= \lambda(t) \exp(\beta_1). \end{aligned} \quad (18.17)$$

Here $\exp(\beta_1)$ is the $X_1 = 1 : X_1 = 0$ hazard ratio. This simple case has no regression assumption but assumes PH and a form for $\lambda(t)$. If the single predictor X_1 is continuous, the model becomes

$$\lambda(t|X_1) = \lambda(t) \exp(\beta_1 X). \quad (18.18)$$

Without further modification (such as taking a transformation of the predictor), the model assumes a straight line in the log hazard or that for all t , an increase in X by one unit increases the hazard by a factor of $\exp(\beta_1)$.

As in logistic regression, much more general regression specifications can be made, including interaction effects. Unlike logistic regression, however, a model containing, say age, sex, and age \times sex interaction is not equivalent to

fitting two separate models. This is because even though males and females are allowed to have unequal age slopes, both sexes are assumed to have the

Table 18.1 Mortality differences and ratios when hazard ratio is 0.5

Subject	5-Year Survival		Difference	Mortality Ratio (T/C)
	C	T		
	1	0.98		
2	0.80	0.89	0.09	$0.11/0.2 = 0.55$
3	0.25	0.50	0.25	$0.5/0.75 = 0.67$

underlying hazard function proportional to $\lambda(t)$ (i.e., the PH assumption holds for sex in addition to age).

18.2.3 Hazard Ratio, Risk Ratio, and Risk Difference

Other ways of modeling predictors can also be specified besides a multiplicative effect on the hazard. For example, one could postulate that the effect of a predictor is to add to the hazard of failure instead of to multiply it by a factor. The effect of a predictor could also be described in terms of a mortality ratio (relative risk), risk difference, odds ratio, or increase in expected failure time. However, just as an odds ratio is a natural way to describe an effect on a binary response, a hazard ratio is often a natural way to describe an effect on survival time. One reason is that a hazard ratio *can* be constant.

Table 18.1 provides treated (T) to control (C) survival (mortality) differences and mortality ratios for three hypothetical types of subjects. We suppose that subjects 1, 2, and 3 have increasingly worse prognostic factors. For example, the age at baseline of the subjects might be 30, 50, and 70 years, respectively. We assume that the treatment affects the hazard by a constant multiple of 0.5 (i.e., PH is in effect and the constant hazard ratio is 0.5). Note that $S_T = S_C^{0.5}$. Notice that the mortality difference and ratio depend on the survival of the control subject. A control subject having “good” predictor values will leave little room for an improved prognosis from the treatment.

The hazard ratio is a basis for describing the mechanism of an effect. In the above example, it is reasonable that the treatment affect each subject by lowering her hazard of death by a factor of 2, even though less sick subjects have a low mortality difference. Hazard ratios also lead to good statistical tests for differences in survival patterns and to predictive models. Once the model is developed, however, survival differences may better capture the impact of a risk factor. Absolute survival differences rather than relative differences

(hazard ratios) also relate more closely to statistical power. For example, even if the effect of a treatment is to halve the hazard rate, a population where the control survival is 0.99 will require a much larger sample than will a population where the control survival is 0.3.

Figure 18.1 depicts the relationship between survival $S(t)$ of a control subject at any time t , relative reduction in hazard (h), and difference in survival $S(t) - S(t)^h$. This figure demonstrates that absolute clinical benefit

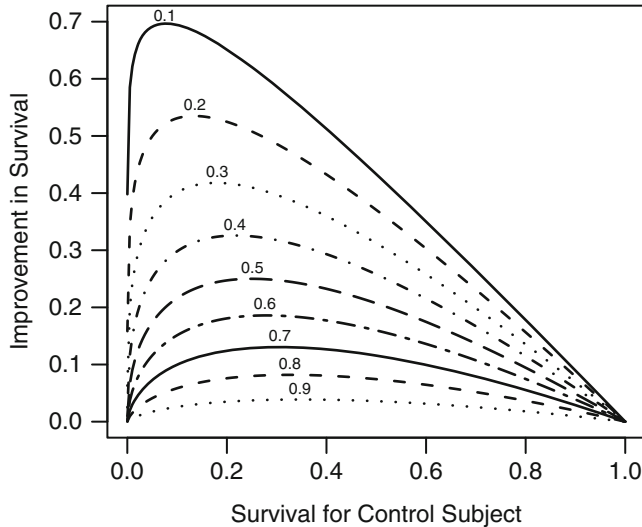


Fig. 18.1 Absolute clinical benefit as a function of survival in a control subject and the relative benefit (hazard ratio). The hazard ratios are given for each curve.

is primarily a function of the baseline risk of a subject. Clinical benefit will also be a function of factors that interact with treatment, that is, factors that modify the relative benefit of treatment. Once a model is developed for estimating $S(t|X)$, this model can be used to estimate absolute benefit as a function of baseline risk factors as well as factors that interact with a treatment. Let X_1 be a binary treatment indicator and let $A = \{X_2, \dots, X_p\}$ be the other factors (which for convenience we assume do not interact with X_1). Then the estimate of $S(t|X_1 = 0, A) - S(t|X_1 = 1, A)$ can be plotted against $S(t|X_1 = 0)$ or against levels of variables in A to display absolute benefit versus overall risk or specific subject characteristics.

1

18.2.4 Specific Models

Let $X\beta$ denote the linear combination of predictors excluding an intercept term. Using the PH formulation, an exponential survival regression model²¹⁸ can be stated as

$$\begin{aligned}\lambda(t|X) &= \lambda \exp(X\beta) \\ S(t|X) &= \exp[-\lambda t \exp(X\beta)] = \exp(-\lambda t)^{\exp(X\beta)}.\end{aligned}\quad (18.19)$$

The parameter λ can be thought of as the antilog of an intercept term since the model could be written $\lambda(t|X) = \exp[(\log \lambda) + X\beta]$. The effect of X on the expected or median failure time is as follows.

$$\begin{aligned}E\{T|X\} &= 1/[\lambda \exp(X\beta)] \\ T_{0.5}|X &= (\log 2)/[\lambda \exp(X\beta)].\end{aligned}\quad (18.20)$$

The exponential regression model can be written in another form that is more numerically stable by replacing the λ parameter with an intercept term in $X\beta$, specifically $\lambda = \exp(\beta_0)$. After redefining $X\beta$ to include β_0 , λ can be dropped in all the above formulas.

The Weibull regression model is defined by one of the following functions (assuming that $X\beta$ does not contain an intercept).

$$\begin{aligned}\lambda(t|X) &= \alpha \gamma t^{\gamma-1} \exp(X\beta) \\ \Lambda(t|X) &= \alpha t^\gamma \exp(X\beta) \\ S(t|X) &= \exp[-\alpha t^\gamma \exp(X\beta)] \\ &= [\exp(-\alpha t^\gamma)]^{\exp(X\beta)}.\end{aligned}\quad (18.21)$$

Note that the parameter α in the homogeneous Weibull model has been replaced with $\alpha \exp(X\beta)$. The median survival time is given by

$$T_{0.5}|X = \{\log 2/[\alpha \exp(X\beta)]\}^{1/\gamma}.\quad (18.22)$$

As with the exponential model, the parameter α could be dropped (and replaced with $\exp(\beta_0)$) if an intercept β_0 is added to $X\beta$.

For numerical reasons it is sometimes advantageous to write the Weibull PH model as

$$S(t|X) = \exp(-\Lambda(t|X)),\quad (18.23)$$

where

$$\Lambda(t|X) = \exp(\gamma \log t + X\beta).\quad (18.24)$$

18.2.5 Estimation

The parameters in λ and β are estimated by maximizing a log likelihood function constructed in the same manner as described in Section 18.1. The only difference is the insertion of $\exp(X_i\beta)$ in the likelihood function:

$$\log L = \sum_{i: Y_i \text{ uncensored}}^n \log[\lambda(Y_i) \exp(X_i \beta)] - \sum_{i=1}^n \Lambda(Y_i) \exp(X_i \beta). \quad (18.25)$$

Once $\hat{\beta}$, the MLE of β , is computed along with the large-sample standard error estimates, hazard ratio estimates and their confidence intervals can readily be computed. Letting s denote the estimated standard error of $\hat{\beta}_j$, a $1 - \alpha$ confidence interval for the $X_j + 1 : X_j$ hazard ratio is given by $\exp[\hat{\beta}_j \pm zs]$, where z is the $1 - \alpha/2$ critical value for the standard normal distribution.

Once the parameters of the underlying hazard function are estimated, the MLE of $\lambda(t)$, $\hat{\lambda}(t)$, can be derived. The MLE of $\lambda(t|X)$, the hazard as a function of t and X , is given by

$$\hat{\lambda}(t|X) = \hat{\lambda}(t) \exp(X \hat{\beta}). \quad (18.26)$$

The MLE of $\Lambda(t)$, $\hat{\Lambda}(t)$, can be derived from the integral of $\hat{\lambda}(t)$ with respect to t . Then the MLE of $S(t|X)$ can be derived:

$$\hat{S}(t|X) = \exp[-\hat{\Lambda}(t) \exp(X \hat{\beta})]. \quad (18.27)$$

For the Weibull model, we denote the MLEs of the hazard parameters α and γ by $\hat{\alpha}$ and $\hat{\gamma}$. The MLE of $\lambda(t|X)$, $\Lambda(t|X)$, and $S(t|X)$ for this model are

$$\begin{aligned} \hat{\lambda}(t|X) &= \hat{\alpha} \hat{\gamma} t^{\hat{\gamma}-1} \exp(X \hat{\beta}) \\ \hat{\Lambda}(t|X) &= \hat{\alpha} t^{\hat{\gamma}} \exp(X \hat{\beta}) \\ \hat{S}(t|X) &= \exp[-\hat{\Lambda}(t|X)]. \end{aligned} \quad (18.28)$$

Confidence intervals for $S(t|X)$ are best derived using general matrix notation to obtain an estimate s of the standard error of $\log[\hat{\lambda}(t|X)]$ from the estimated information matrix of all hazard and regression parameters. A confidence interval for \hat{S} will be of the form

$$\hat{S}(t|X)^{\exp(\pm zs)}. \quad (18.29)$$

The MLEs of β and of the hazard shape parameters lead directly to MLEs of the expected and median life length. For the Weibull model the MLE of the median life length given X is

$$\hat{T}_{0.5}|X = \{\log 2 / [\hat{\alpha} \exp(X \hat{\beta})]\}^{1/\hat{\gamma}}. \quad (18.30)$$

For the exponential model, the MLE of the expected life length for a subject having predictor values X is given by

$$\hat{E}(T|X) = [\hat{\lambda} \exp(X \hat{\beta})]^{-1}, \quad (18.31)$$

where $\hat{\lambda}$ is the MLE of λ .

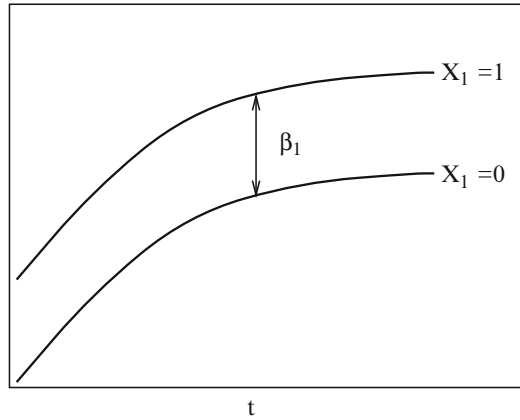


Fig. 18.2 PH model with one binary predictor. Y-axis is $\log \lambda(t)$ or $\log A(t)$. For $\log A(t)$, the curves must be non-decreasing. For $\log \lambda(t)$, they may be any shape.

18.2.6 Assessment of Model Fit

Three assumptions of the parametric PH model were listed in Section 18.2.2. We now lay out in more detail what relationships need to be satisfied. We first assume a PH model with a single binary predictor X_1 . For a general underlying hazard function $\lambda(t)$, all assumptions of the model are displayed in Figure 18.2. In this case, the assumptions are PH and a shape for $\lambda(t)$.

If $\lambda(t)$ is Weibull, the two curves will be linear if $\log t$ is plotted instead of t on the x -axis. Note also that if there is no association between X and survival ($\beta_1 = 0$), estimates of the two curves will be close and will intertwine due to random variability. In this case, PH is not an issue.

If the single predictor is continuous, the relationships in Figures 18.3 and 18.4 must hold. Here linearity is assumed (unless otherwise specified) besides PH and the form of $\lambda(t)$. In Figure 18.3, the curves must be parallel for any choices of times t_1 and t_2 as well as each individual curve being linear. Also, the difference between ordinates needs to conform to the assumed distribution. This difference is $\log[\lambda(t_2)/\lambda(t_1)]$ or $\log[A(t_2)/A(t_1)]$.

Figure 18.4 highlights the PH assumption. The relationship between the two curves must hold for any two values c and d of X_1 . The shape of the function for a given value of X_1 must conform to the assumed $\lambda(t)$. For a Weibull model, the functions should each be linear in $\log t$.

When there are multiple predictors, the PH assumption can be displayed in a way similar to Figures 18.2 and 18.4 but with the population additionally cross-classified by levels of the other predictors besides X_1 . If there is one binary predictor X_1 and one continuous predictor X_2 , the relationship in

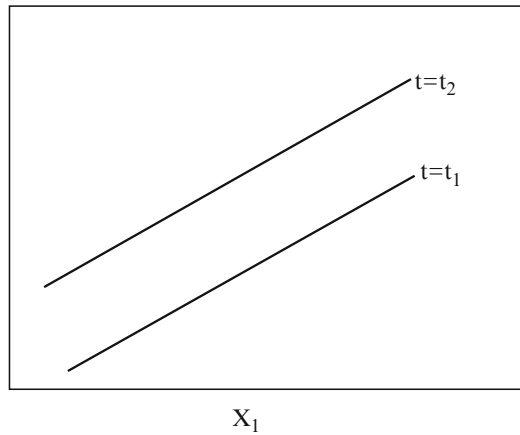


Fig. 18.3 PH model with one continuous predictor. Y -axis is $\log \lambda(t)$ or $\log \Lambda(t)$; for $\log \Lambda(t)$, drawn for $t_2 > t_1$. The slope of each line is β_1 .

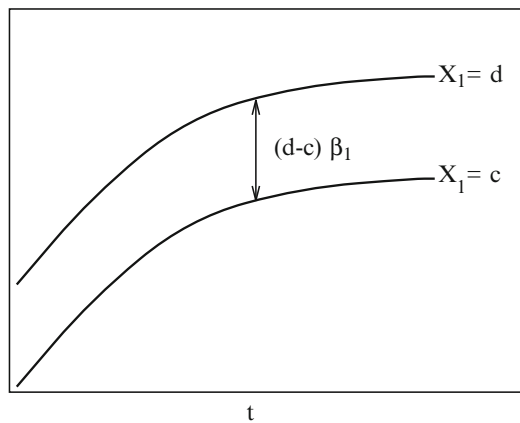


Fig. 18.4 PH model with one continuous predictor. Y -axis is $\log \lambda(t)$ or $\log \Lambda(t)$. For $\log \lambda$, the functions need not be monotonic.

Figure 18.5 must hold at each time t if linearity is assumed for X_2 and there is no interaction between X_1 and X_2 . Methods for verifying the regression assumptions (e.g., splines and residuals) and the PH assumption are covered in detail under the Cox PH model in Chapter 20.

The method for verifying the assumed shape of $S(t)$ in Section 18.1.3 is also useful when there are a limited number of categorical predictors. To validate a Weibull PH model one can stratify on X and plot $\log \Lambda_{\text{KM}}(t|X \text{ stratum})$ against $\log t$. This graph simultaneously assesses PH in addition to shape assumptions—all curves should be parallel as well as straight. Straight but nonparallel (non-PH) curves indicate that a series of Weibull models with differing γ parameters will fit.

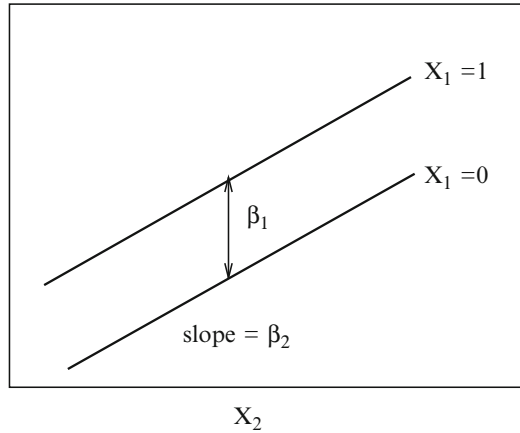


Fig. 18.5 Regression assumptions, linear additive PH or AFT model with two predictors. For PH, Y -axis is $\log \lambda(t)$ or $\log \Lambda(t)$ for a fixed t . For AFT, Y -axis is $\log(T)$.

18.3 Accelerated Failure Time Models

18.3.1 Model

Besides modeling the effect of predictors by a multiplicative effect on the hazard function, other regression effects can be specified. The *accelerated failure time (AFT) model* is commonly used; it specifies that the predictors act multiplicatively on the failure time or additively on the log failure time. The effect of a predictor is to alter the rate at which a subject proceeds along the time axis (i.e., to accelerate the time to failure [331, pp. 33–35]). The model is

$$S(t|X) = \psi((\log(t) - X\beta)/\sigma), \quad (18.32)$$

where ψ is any standardized survival distribution function. The parameter σ is called the *scale parameter*. The model can also be stated as $(\log(T) - X\beta)/\sigma \sim \psi$ or $\log(T) = X\beta + \sigma\epsilon$, where ϵ is a random variable from the distribution ψ . Sometimes the untransformed T is used in place of $\log(T)$. When the log form is used, the models are said to be log-normal, log-logistic, and so on.

The exponential and Weibull are the only two distributions that can describe either a PH or an AFT model.

18.3.2 Model Assumptions and Interpretation of Parameters

The $\log \lambda$ or $\log \Lambda$ transformation of the PH model has the following equivalent for AFT models.

$$\psi^{-1}[S(t|X)] = (\log(t) - X\beta)/\sigma. \quad (18.33)$$

Letting as before ϵ denote a random variable from the distribution S , the model is also

$$\log(T) = X\beta + \sigma\epsilon. \quad (18.34)$$

So the property of the response T of interest for regression modeling is $\log(T)$. In the absence of censoring, we could check the model by plotting an X against $\log T$ and checking that the residuals $\log(T) - X\hat{\beta}$ are distributed as ψ to within a scale factor.

The assumptions of the AFT model are thus the following.

1. The true form of ψ (the distributional family) is correctly specified.
2. In the absence of nonlinear and interaction terms, each X_j affects $\log(T)$ or $\psi^{-1}[S(t|X)]$ linearly.
3. Implicit in these assumptions is that σ is a constant independent of X .

A one-unit change in X_j is then most simply understood as a β_j change in the log of the failure time. The one-unit change in X_j increases the failure time by a factor of $\exp(\beta_j)$.

The median survival time is obtained by solving $\psi((\log(t) - X\beta)/\sigma) = 0.5$ giving

$$T_{0.5}|X = \exp[X\beta + \sigma\psi^{-1}(0.5)] \quad (18.35)$$

18.3.3 Specific Models

Common choices for the distribution function ψ in Equation 18.32 are the extreme value distribution $\psi(u) = \exp(-\exp(u))$, the logistic distribution $\psi(u) = [1 + \exp(u)]^{-1}$, and the normal distribution $\psi(u) = 1 - \Phi(u)$. The AFT model equivalent of the Weibull model is obtained by using the extreme value distribution, negating β , and replacing γ with $1/\sigma$ in Equation 18.24:

$$\begin{aligned} S(t|X) &= \exp[-\exp((\log(t) - X\beta)/\sigma)] \\ T_{0.5}|X &= [\log(2)]^\sigma \exp(X\beta). \end{aligned} \quad (18.36)$$

The exponential model is obtained by restricting $\sigma = 1$ in the extreme value distribution.

The log-normal regression model is

$$S(t|X) = 1 - \Phi((\log(t) - X\beta)/\sigma), \quad (18.37)$$

and the log-logistic model is

$$S(t|X) = [1 + \exp((\log(t) - X\beta)/\sigma)]^{-1}. \quad (18.38)$$

The t distribution allows for more flexibility by varying the degrees of freedom. Figure 18.6 depicts possible hazard functions for the log t distribution for varying σ and degrees of freedom. However, this distribution does not have a late increasing hazard phase typical of human survival.

```
require(rms)

haz <- survreg.auxinfo$t$hazard
times <- c(seq(0, .25, length=100), seq(.26, 2, length=150))
high <- c(6, 1.5, 1.5, 1.75)
low <- c(0, 0, 0, .25)
dfs <- c(1, 2, 3, 5, 7, 15, 500)
cols <- rep(1, 7)
ltys <- 1:7
i <- 0
for(scale in c(.25, .6, 1, 2)) {
  i <- i + 1
  plot(0, 0, xlim=c(0,2), ylim=c(low[i], high[i]),
       xlab=expression(t), ylab=expression(lambda(t)), type="n")
  col <- 1.09
  j <- 0
  for(df in dfs) {
    j <- j+1
    ## Divide by t to get hazard for log t distribution
    lines(times,
          haz(log(times), 0, c(log(scale), df))/times,
          col=cols[j], lty=ltys[j])
    if(i==1) text(1.7, .23 + haz(log(1.7), 0,
                                c(log(scale),df))/1.7, format(df))
  }
  title(paste("Scale:", format(scale)))
} # Figure 18.6
```

All three of these parametric survival models have median survival time $T_{0.5}|X = \exp(X\beta)$.

18.3.4 Estimation

Maximum likelihood estimation is used much the same as in Section 18.2.5. Care must be taken in the choice of initial values; iterative methods are especially prone to problems in choosing the initial $\hat{\sigma}$. Estimation works better if σ is parameterized as $\exp(\delta)$. Once β and σ ($\exp(\delta)$) are estimated, MLEs of secondary parameters such as survival probabilities and medians can readily be obtained:

$$\begin{aligned}\hat{S}(t|X) &= \psi((\log(t) - X\hat{\beta})/\hat{\sigma}) \\ \hat{T}_{0.5}|X &= \exp[X\hat{\beta} + \hat{\sigma}\psi^{-1}(0.5)].\end{aligned}\tag{18.39}$$

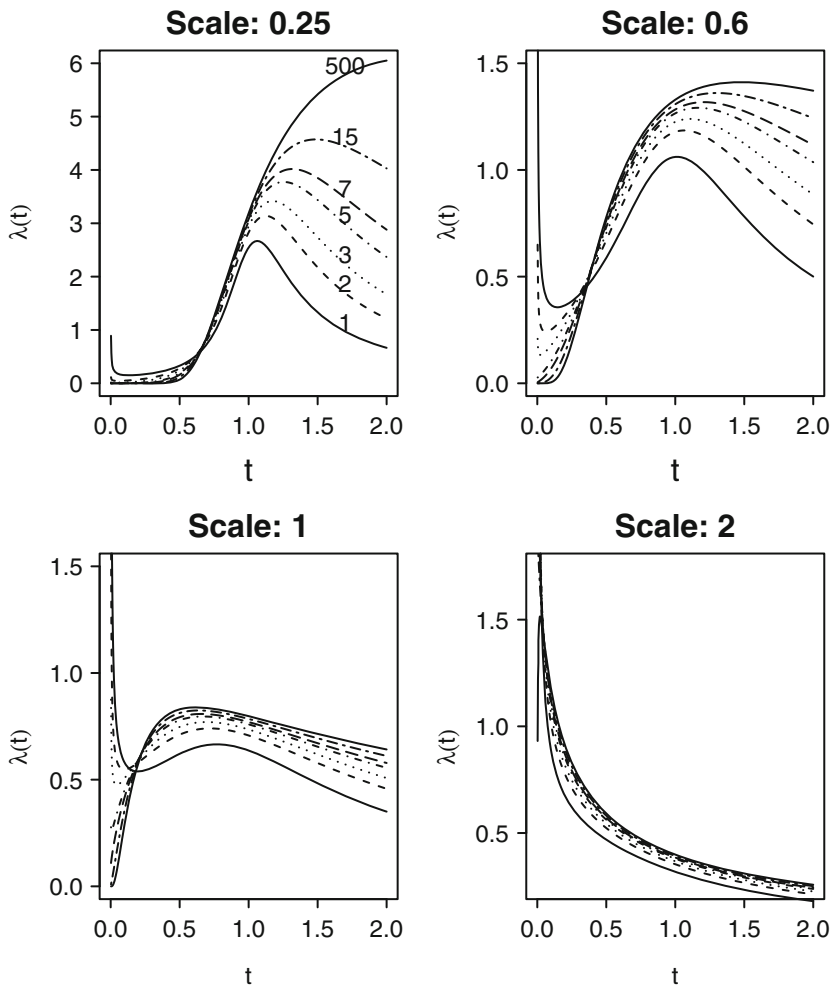


Fig. 18.6 $\log(T)$ distribution for $\sigma = 0.25, 0.6, 1, 2$ and for degrees of freedom 1, 2, 3, 5, 7, 15, 500 (almost log-normal). The top left plot has degrees of freedom written in the plot.

For normal and logistic distributions, $\hat{T}_{0.5}|X = \exp(X\hat{\beta})$. The MLE of the effect on $\log(T)$ of increasing X_j by d units is $\hat{\beta}_j d$ if X_j is linear and additive.

The delta (statistical differential) method can be used to compute an estimate of the variance of $f = [\log(t) - X\hat{\beta}]/\hat{\sigma}$. Let $(\hat{\beta}, \hat{\delta})$ denote the estimated parameters, and let \hat{V} denote the estimated covariance matrix for these parameter estimates. Let F denote the vector of derivatives of f with respect to $(\beta_0, \beta_1, \dots, \beta_p, \delta)$; that is, $F = [-1, -X_1, -X_2, \dots, -X_p, -(\log(t) - X\hat{\beta})/\hat{\sigma}]$. The variance of f is then approximately

$$\text{Var}(f) = F\hat{V}F'. \tag{18.40}$$

Letting s be the square root of the variance estimate and $z_{1-\alpha/2}$ be the normal critical value, a $1 - \alpha$ confidence limit for $S(t|X)$ is

$$\psi((\log(t) - X\hat{\beta})/\hat{\sigma} \pm z_{1-\alpha/2} \times s). \quad (18.41)$$

18.3.5 Residuals

For an AFT model, standardized residuals are simply

$$r = (\log(T) - X\hat{\beta})/\sigma. \quad (18.42)$$

4 When T is right-censored, r is right-censored. Censoring must be taken into account, for example, by displaying Kaplan–Meier estimates based on groups of residuals rather than showing individual residuals. The residuals can be used to check for lack of fit as described in the next section. Note that examining individual uncensored residuals is not appropriate, as their distribution is conditional on $T_i < C_i$, where C_i is the censoring time.

Cox and Snell¹³⁴ proposed a type of general residuals that also work for censored data. Using their method on the cumulative probability scale results in the probability integral transformation. If the probability of failure before time t given X is $S(t|X)$, $F(T|X) = 1 - S(T|X)$ has a uniform $[0, 1]$ distribution, where T is a subject's actual failure time. When T is right-censored, so is $1 - S(T|X)$. Substituting \hat{S} for S results in an approximate uniform $[0, 1]$ distribution for any value of X . One minus the Kaplan–Meier estimate of $1 - \hat{S}(T|X)$ (using combined data for all X) is compared against a 45° line to check for goodness of fit. A more stringent assessment is obtained by repeating this process while stratifying on X .

18.3.6 Assessment of Model Fit

For a single binary predictor, all assumptions of the AFT model are depicted in Figure 18.7. That figure also shows the assumptions for any two values of a single continuous predictor that behaves linearly. For a single continuous predictor, the relationships in Figure 18.8 must hold for any two follow-up times. The regression assumptions are isolated in Figure 18.5.

To verify the fit of a log-logistic model with age as the only predictor, one could stratify by quartiles of age and check for linearity and parallelism of the four logit $S_A(t)$ or $S_{KM}(t)$ curves over increasing t as in Figure 18.7, which stresses the distributional assumption (no T by X interaction and linearity vs. $\log(t)$). To stress the linear regression assumption while checking for absence of time interactions (part of the distributional assumptions), one could make

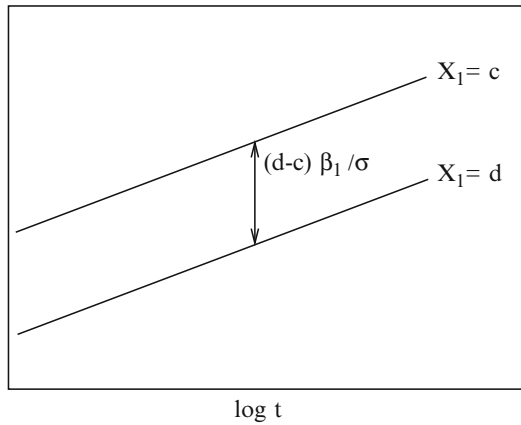


Fig. 18.7 AFT model with one predictor. Y-axis is $\psi^{-1}[S(t|X)] = (\log(t) - X\beta)/\sigma$. Drawn for $d > c$. The slope of the lines is σ^{-1} .

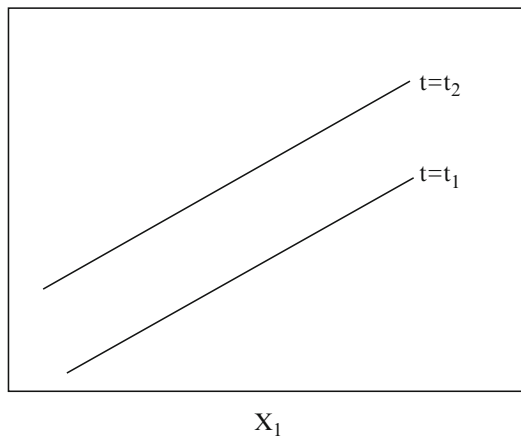


Fig. 18.8 AFT model with one continuous predictor. Y-axis is $\psi^{-1}[S(t|X)] = (\log(t) - X\beta)/\sigma$. Drawn for $t_2 > t_1$. The slope of each line is β_1/σ and the difference between the lines is $\log(t_2/t_1)/\sigma$.

a plot like Figure 18.8. For each decile of age, the logit transformation of the 1-, 3-, and 5-year survival estimates for that decile would be plotted against the mean age in the decile. This checks for linearity and constancy of the age effect over time. Regression splines will be a more effective method for checking linearity and determining transformations. This is demonstrated in Chapter 20 with the Cox model, but identical methods apply here.

As an example, consider data from Kalbfleisch and Prentice [331, pp. 1–2], who present data from Pike⁵⁰⁸ on the time from exposure to the carcinogen DMBA to mortality from vaginal cancer in rats. The rats are divided into two groups on the basis of a pre-treatment regime. Survival times in days (with censored times marked +) are found in Table 18.2.

Table 18.2 Rat vaginal cancer data from Pike⁵⁰⁸

Group 1	143	164	188	188	190	192	206	209	213	216
	220	227	230	234	246	265	304	216 ⁺	244 ⁺	
Group 2	142	156	163	198	205	232	232	233	233	233
	233	239	240	261	280	280	296	296	323	204 ⁺
										344 ⁺

```

getHdata(kprats)
kprats$group ← factor(kprats$group, 0:1, c('Group 1', 'Group 2'))
dd ← datadist(kprats); options(datadist="dd")

S ← with(kprats, Surv(t, death))
f ← npsurv(S ~ group, type="fleming", data=kprats)
survplot(f, n.risk=TRUE, conf='none', # Figure 18.9
         label.curves=list(keys='lines'), levels.only=TRUE)
title(sub="Nonparametric estimates", adj=0, cex=.7)

# Check fits of Weibull, log-logistic, log-normal
xl ← c(4.8, 5.9)
survplot(f, loglog=TRUE, logt=TRUE, conf="none", xlim=xl,
         label.curves=list(keys='lines'), levels.only=TRUE)
title(sub="Weibull (extreme value)", adj=0, cex=.7)
survplot(f, fun=function(y)log(y/(1-y)), ylab="logit S(t)",
         logt=TRUE, conf="none", xlim=xl,
         label.curves=list(keys='lines'), levels.only=TRUE)
title(sub="Log-logistic", adj=0, cex=.7)
survplot(f, fun=qnorm, ylab="Inverse Normal S(t)",
         logt=TRUE, conf="none",
         xlim=xl, cex.label=.7,
         label.curves=list(keys='lines'), levels.only=TRUE)
title(sub="Log-normal", adj=0, cex=.7)

```

The top left plot in Figure 18.9 displays nonparametric survival estimates for the two groups, with the number of rats “at risk” at each 30-day mark written above the x -axis. The remaining three plots are for checking assumptions of three models. None of the parametric models presented will completely allow for such a long period with no deaths. Neither will any allow for the early crossing of survival curves. Log-normal and log-logistic models yield very similar results due to the similarity in shapes between $\Phi(z)$ and $[1 + \exp(-z)]^{-1}$ for non-extreme z . All three transformations show good parallelism after the early crossing. The log-logistic and log-normal transformations are slightly more linear. The fitted models are:

```

fw ← psm(S ~ group, data=kprats, dist='weibull')
fl ← psm(S ~ group, data=kprats, dist='loglogistic',
         y=TRUE)
fn ← psm(S ~ group, data=kprats, dist='lognormal')
latex(fw, fi='')

```

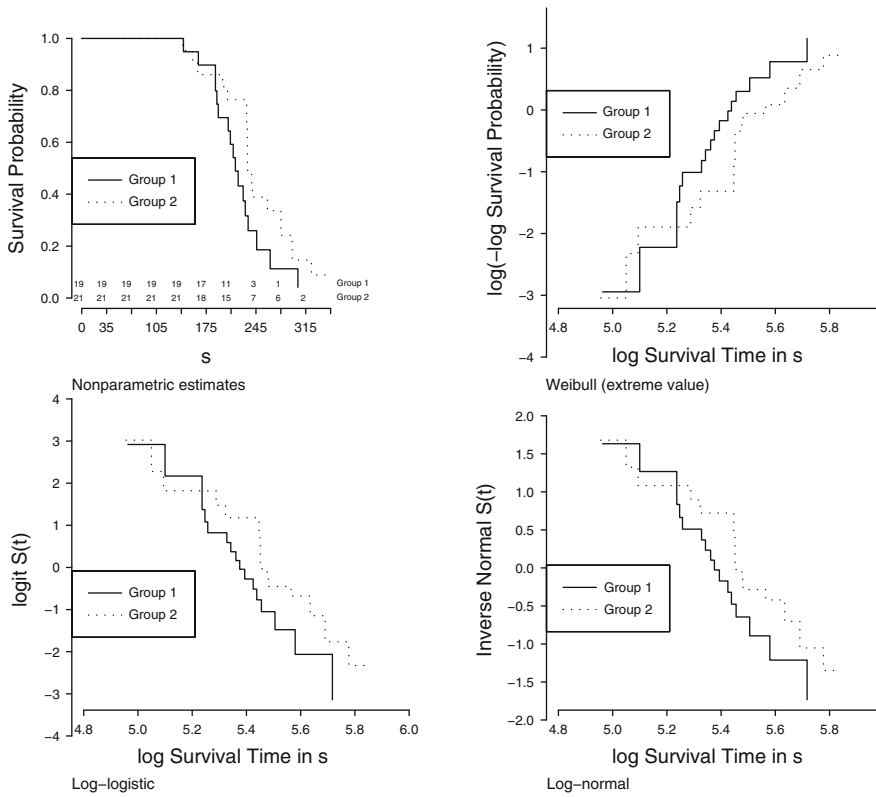


Fig. 18.9 Altschuler–Nelson–Fleming–Harrington nonparametric survival estimates for rats treated with DMBA,⁵⁰⁸ along with various transformations of the estimates for checking distributional assumptions of three parametric survival models.

$$\text{Prob}\{T \geq t\} = \exp\left[-\exp\left(\frac{\log(t) - X\hat{\beta}}{0.1832976}\right)\right] \text{ where}$$

$$X\hat{\beta} = 5.450859 + 0.131983[\text{Group } 2]$$

and $[c] = 1$ if subject is in group c , 0 otherwise.

```
latex(f1, fi='')
```

Table 18.3 Group effects from three survival models

Model	Group 2:1	Median Survival Time	
	Failure Time Ratio	Group 1	Group 2
Extreme Value (Weibull)	1.14	217	248
Log-logistic	1.11	217	241
Log-normal	1.10	217	238

$$\text{Prob}\{T \geq t\} = [1 + \exp(\frac{\log(t) - X\hat{\beta}}{0.1159753})]^{-1} \quad \text{where}$$

$$\begin{aligned} X\hat{\beta} = & \\ & 5.375675 \\ & +0.1051005[\text{Group } 2] \end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise.

```
latex(fn, fi='')
```

$$\text{Prob}\{T \geq t\} = 1 - \Phi(\frac{\log(t) - X\hat{\beta}}{0.2100184}) \quad \text{where}$$

$$\begin{aligned} X\hat{\beta} = & \\ & 5.375328 \\ & +0.0930606[\text{Group } 2] \end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise.

The estimated failure time ratios and median failure times for the two groups are given in Table 18.3. For example, the effect of going from Group 1 to Group 2 is to increase log failure time by 0.132 for the extreme value model, giving a Group 2:1 failure time ratio of $\exp(0.132) = 1.14$. This ratio is also the ratio of median survival times. We choose the log-logistic model for its simpler form. The fitted survival curves are plotted with the nonparametric estimates in Figure 18.10. Excellent agreement is seen, except for 150 to 180 days for Group 2. The standard error of the regression coefficient for group in the log-logistic model is 0.0636 giving a Wald χ^2 for group differences of $(.105/.0636)^2 = 2.73, P = 0.1$.

```
survplot(f, conf.int=FALSE, # Figure 18.10
         levels.only=TRUE, label.curves=list(keys='lines'))
survplot(f1, add=TRUE, label.curves=FALSE, conf.int=FALSE)
```

The Weibull PH form of the fitted extreme value model, using Equation 18.24, is

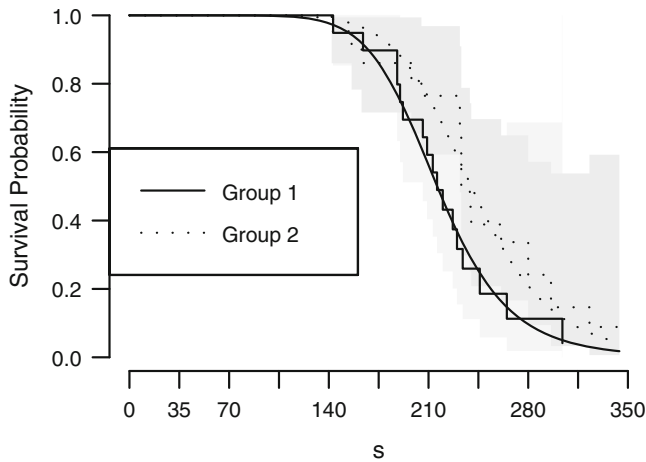


Fig. 18.10 Agreement between fitted log-logistic model and nonparametric survival estimates for rat vaginal cancer data.

$$\text{Prob}\{T \geq t\} = \exp\{-t^{5.456} \exp(X\hat{\beta})\} \quad \text{where}$$

$$X\hat{\beta} = \begin{aligned} & -29.74 \\ & -0.72[\text{Group } 2] \end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise.

A sensitive graphical verification of the distributional assumptions of the AFT model is obtained by plotting the estimated survival distribution of standardized residuals (Equation 18.3.5), censored identically to the way T is censored. This distribution is plotted along with the theoretical distribution ψ . The assessment may be made more stringent by stratifying the residuals by important subject characteristics and plotting separate survival function estimates; they should all have the same standardized distribution (e.g., same σ).

```
r <- resid(fl, 'cens')
survplot(npsurv(r ~ group, data=kprats),
         conf='none', xlab='Residual',
         label.curves=list(keys='lines'), levels.only=TRUE)
survplot(npsurv(r ~ 1), conf='none', add=TRUE, col='red')
lines(r, lwd=1, col='blue') # Figure 18.11
```

As an example, Figure 18.11 shows the Kaplan–Meier estimate of the distribution of residuals, Kaplan–Meier estimates stratified by group, and the assumed log-logistic distribution.

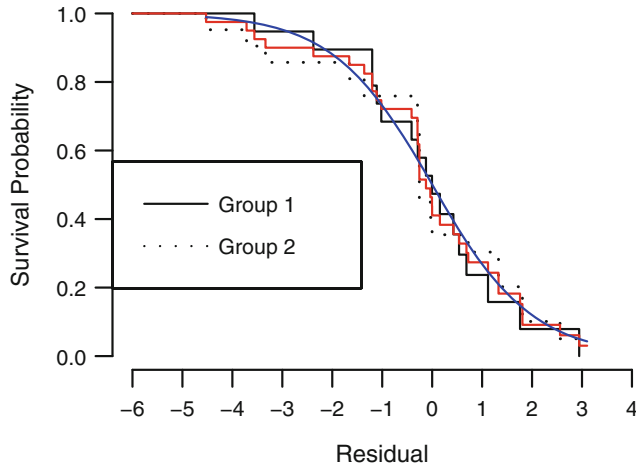


Fig. 18.11 Kaplan–Meier estimates of distribution of standardized censored residuals from the log-logistic model, along with the assumed standard log-logistic distribution (dashed curve). The step functions in red is the estimated distribution of all residuals, and the step functions in black are the estimated distributions of residuals stratified by group, as indicated. The blue curve is the assumed log-logistic distribution.

Section 19.2 has a more in-depth example of this approach.

18.3.7 Validating the Fitted Model

AFT models may be validated for both calibration and discrimination accuracy using the same methods that are presented for the Cox model in Section 20.11. The methods discussed there for checking calibration are based on choosing a single follow-up time. Checking the distributional assumptions of the parametric model is also a check of calibration accuracy in a sense. Another indirect calibration assessment may be obtained from a set of Cox–Snell residuals (Section 18.3.5) or by using ordinary residuals as just described. A higher resolution indirect calibration assessment based on plotting individual uncensored failure times is available when the theoretical censoring times for those observations are known. Let C denote a subject’s censoring time and F the cumulative distribution of a failure time T . The expected value of $F(T|X)$ is 0.5 when T is an actual failure time random variable. The expected value for an event time that is observed *because it is uncensored* is the expected value of $F(T|T \leq C, X) = 0.5F(C|X)$. A smooth plot (using, say, `loess`) of $F(T|X) - 0.5F(C|X)$ against $X\hat{\beta}$ should be a flat line through $y = 0$ if the model is well calibrated. A smooth plot of $2F(T|X)/F(C|X)$ against $X\hat{\beta}$ (or anything else) should be a flat line through $y = 1$. This method assumes that the model is calibrated well enough that we can substitute $1 - \hat{S}(C|X)$ for $F(C|X)$.

18.4 Buckley–James Regression Model

Buckley and James⁸¹ developed a method for estimating regression coefficients using least squares after imputing censored residuals. Their method does not assume a distribution for survival time or the residuals, but is aimed at estimating expected survival time or expected log survival time given predictor variables. This method has been generalized to allow for smooth non-linear effects and interactions in the Sbj function in the `rms` package, written by Stare and Harrell⁵⁸⁵.

18.5 Design Formulations

Various designs can be formulated with survival regression models just as with other regression models. By constructing the proper dummy variables, ANOVA and ANOCOVA models can easily be specified for testing differences in survival time between multiple treatments. Interactions and complex non-linear effects may also be modeled.

18.6 Test Statistics

As discussed previously, likelihood ratio, score, and Wald statistics can be derived from the maximum likelihood analysis, and the choice of test statistic depends on the circumstance and on computational convenience.

18.7 Quantifying Predictive Ability

See Section 20.10 for a generalized measure of concordance between predicted and observed survival time (or probability of survival) for right-censored data.

18.8 Time-Dependent Covariates

Time-dependent covariates (predictors) requires special likelihood functions and add significant complexity to analyses in exchange for greater versatility and enhanced predictive discrimination⁶⁰⁴. Nicolaie et al.⁴⁷⁷ and D'Agostino et al.¹⁴⁵ provide useful static covariate approaches to modeling time-dependent predictors using landmark analysis.

18.9 R Functions

Therneau's `survreg` function (part of his `survival` package) can fit regression models in the AFT family with left-, right-, or interval-censoring. The time variable can be untransformed or log-transformed (the default). Distributions supported are extreme value (Weibull and exponential), normal, logistic, and Student- t . The version of `survreg` in `rms` that fits parametric survival models in the same framework as `lrm`, `ols`, and `cph` is called `psm`. `psm` works with `print`, `coef`, `formula`, `specs`, `summary`, `anova`, `predict`, `Predict`, `fastbw`, `latex`, `nomogram`, `validate`, `calibrate`, `survest`, and `survplot` functions for obtaining and plotting predicted survival probabilities. The `dist` argument to `psm` can be "exponential", "extreme", "gaussian", "logistic", "loglogistic", "lognormal", "t", or "weibull". To fit a model with no covariables, use the command

```
psm(Surv(d.time, event) ~ 1)
```

To restate a Weibull or exponential model in PH form, use the `pphsm` function. An example of how many of the functions are used is found below.

```
units(d.time) ← "Year"
f ← psm(Surv(d.time, cdeath) ~ lsp(age, 65)*sex)
# default is Weibull
anova(f)
summary(f)           # summarize effects with delta log T
latex(f)             # typeset math. form of fitted model
survest(f, times=1) # 1y survival est. for all subjects
survest(f, expand.grid(sex="female", age=30:80), times=1:2)
# 1y, 2y survival estimates vs. age, for females
survest(f, data.frame(sex="female", age=50))
# survival curve for an individual subject
survplot(f, sex=NA, age=50, n.risk=T)
# survival curves for each sex, adjusting age to 50
f.ph ← pphsm(f)      # convert from AFT to PH
summary(f.ph)        # summarize with hazard ratios
# instead of changes in log(T)
```

Special functions work with objects created by `psm` to create S functions that contain the analytic form for predicted survival probabilities (`Survival`), hazard functions (`Hazard`), quantiles of survival time (`Quantile`), and mean or expected survival time (`Mean`). Once the S functions are constructed, they can be used in a variety of contexts. The `survplot` and `survest` functions have a special argument for `psm` fits: `what`. The default is `what="survival"` to estimate or plot survival probabilities. Specifying `what="hazard"` will plot hazard functions. `Predict` also has a special argument for `psm` fits: `time`. Specifying a single value for `time` results in survival probability for that time being plotted instead of $X\hat{\beta}$. Examples of many of the functions appear below, with the output of the `survplot` command shown in Figure 18.12.

```
med ← Quantile(f1)
meant ← Mean(f1)
```

```
haz ← Hazard(fl)
surv ← Survival(fl)
latex(surv, file='', type='Sinput')
```

```
surv ← function (times = NULL, lp = NULL,
                 parms = -2.15437773933124)
{
  1/(1 + exp((logb(times) - lp)/exp(parms)))
}
```

```
# Plot estimated hazard functions and add median
# survival times to graph
survplot(fl, group, what="hazard") # Figure 18.12
# Compute median survival time
m ← med(lp=predict(fl,
                  data.frame(group=levels(kprats$group))))
m
```

```
      1      2
216.0857 240.0328
```

```
med(lp=range(fl$linear.predictors))
```

```
[1] 216.0857 240.0328
```

```
m ← format(m, digits=3)
text(68, .02, paste("Group 1 median: ", m[1], "\n",
                  "Group 2 median: ", m[2], sep=""))
# Compute survival probability at 210 days
xbeta ← predict(fl,
                data.frame(group=c("Group 1", "Group 2")))
surv(210, xbeta)
```

```
      1      2
0.5612718 0.7599776
```

The S object called `survreg.distributions` in Therneau's `survival` package and the object `survreg.auxinfo` in the `rms` package have detailed information for extreme-value, logistic, normal, and t distributions. For each distribution, components include the deviance function, an algorithm for obtaining starting parameter estimates, a \LaTeX representation of the survival function, and S functions defining the survival, hazard, quantile functions, and basic survival inverse function (which could have been used in Figure 18.9). See Figure 18.6 for examples. `rms`'s `val.surv` function is useful for indirect external validation of parametric models using Cox–Snell residuals and other approaches of Section 18.3.7. The `plot` method for an object created by `val.surv` makes it easy to stratify all computations by a variable of interest to more stringently validate the fit with respect to that variable.

`rms`'s `bj` function fits the Buckley–James model for right-censored responses.

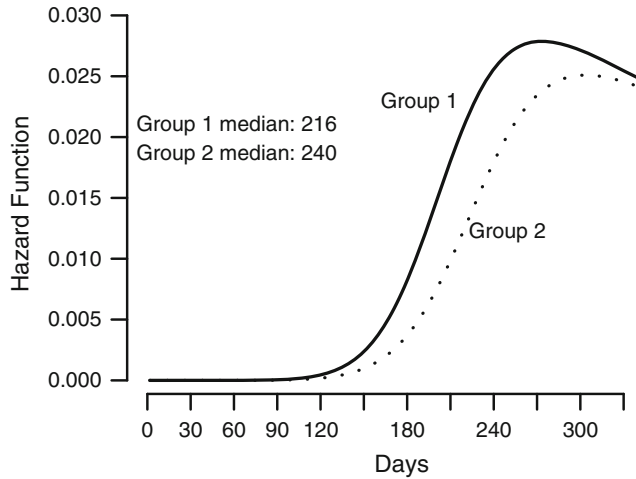


Fig. 18.12 Estimated hazard functions for log-logistic fit to rat vaginal cancer data, along with median survival times.

Kooperberg et al.'s adaptive linear spline log-hazard model^{360,361,594} has been implemented in the S function `hars`. Their procedure searches for second-order interactions involving predictors (and linear splines of them) and linear splines in follow-up time (allowing for non-proportional hazards). `hars` is also used to estimate calibration curves for parametric survival models (rms function `calibrate`) as it is for Cox models.

18.10 Further Reading

- 1 Wellek⁶⁵⁷ developed a test statistic for a specified maximum survival difference after relating this difference to a hazard ratio.
- 2 Hougaard³⁰⁸ compared accelerated failure time models with proportional hazard models.
- 3 Gore et al.²²⁶ discuss how an AFT model (the log-logistic model) gives rise to varying hazard ratios.
- 4 See Hillis²⁹³ for other types of residuals and plots that use them.
- 5 See Gore et al.²²⁶ and Lawless³⁸² for other methods of checking assumptions for AFT models. Lawless is an excellent text for in-depth discussion of parametric survival modeling. Kwong and Hutton³⁶⁹ present other methods of choosing parametric survival models, and discuss the robustness of estimates when fitting an incorrectly chosen accelerated failure time model.

18.11 Problems

1. For the failure times (in days)

$$1 \quad 3 \quad 3^+ \quad 6^+ \quad 7^+$$

- compute MLEs of the following parameters of an exponential distribution by hand: λ , μ , $T_{0.5}$, and $S(3 \text{ days})$. Compute 0.95 confidence limits for λ and $S(3)$, basing the latter on $\log[A(t)]$.
2. For the same data in Problem 1, compute MLEs of parameters of a Weibull distribution. Also compute the MLEs of $S(3)$ and $T_{0.5}$.

Chapter 19

Case Study in Parametric Survival Modeling and Model Approximation

Consider the random sample of 1000 patients from the SUPPORT study,³⁵² described in Section 3.12. In this case study we develop a parametric survival time model (accelerated failure time model) for time until death for the acute disease subset of SUPPORT (acute respiratory failure, multiple organ system failure, coma). We eliminate the chronic disease categories because the shapes of the survival curves are different between acute and chronic disease categories. To fit both acute and chronic disease classes would require a log-normal model with σ parameter that is disease-specific.

Patients had to survive until day 3 of the study to qualify. The baseline physiologic variables were measured during day 3.

19.1 Descriptive Statistics

First we create a variable `acute` to flag the categories of interest, and print univariable descriptive statistics for the data subset.

```
require(rms)

getHdata(support)      # Get data frame from web site
acute <- support$dzclass %in% c('ARF/MOSF', 'Coma')
latex(describe(support[acute,]), file='')
```


support[acute,]
35 Variables 537 Observations

age : Age

n missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
537	0	529	1	60.7	28.49	35.22	47.93	63.67	74.49	81.54 85.56
lowest : 18.04 18.41 19.76 20.30 20.31										
highest: 91.62 91.82 91.93 92.74 95.51										

death : Death at any time up to NDI date:31DEC94

n missing	unique	Info	Sum	Mean
537	0	2	0.67	356 0.6629

sex

n missing	unique
537	0 2
female (251, 47%), male (286, 53%)	

hospdead : Death in Hospital

n missing	unique	Info	Sum	Mean
537	0	2	0.7	201 0.3743

slos : Days from Study Entry to Discharge

n missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
537	0	85	1	23.44	4.0	5.0	9.0	15.0	27.0	47.4 68.2
lowest : 3 4 5 6 7, highest: 145 164 202 236 241										

d.time : Days of Follow-Up

n missing	unique	Info	Mean	.05	.10	.25	.50	.75	.90	.95
537	0	340	1	446.1	4	6	16	182	724	1421 1742
lowest : 3 4 5 6 7, highest: 1977 1979 1982 2011 2022										

dzgroup

n missing	unique
537	0 3
ARF/MOSF w/Sepsis (391, 73%), Coma (60, 11%), MOSF w/Malig (86, 16%)	

dzclass

n missing	unique
537	0 2
ARF/MOSF (477, 89%), Coma (60, 11%)	

num.co : number of comorbidities

n missing	unique	Info	Mean				
537	0	7	0.93 1.525				
Frequency	0	1	2	3	4	5	6
	111	196	133	51	31	10	5
%	21	36	25	9	6	2	1

edu : Years of Education

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 411 126 22 0.96 12.03 7 8 10 12 14 16 17
 lowest : 0 1 2 3 4, highest: 17 18 19 20 22

income

n missing unique
 335 202 4
 under \$11k (158, 47%), \$11-\$25k (79, 24%), \$25-\$50k (63, 19%)
 >\$50k (35, 10%)

scoma : SUPPORT Coma Score based on Glasgow D3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 11 0.82 19.24 0 0 0 0 37 55 100

	0	9	26	37	41	44	55	61	89	94	100
Frequency	301	50	44	19	17	43	11	6	8	6	32
%	56	9	8	4	3	8	2	1	1	1	6

charges : Hospital Charges

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 517 20 516 1 86652 11075 15180 27389 51079 100904 205562 283411
 lowest : 3448 4432 4574 5555 5849
 highest: 504660 538323 543761 706577 740010

totcst : Total RCC cost

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 471 66 471 1 46360 6359 8449 15412 29308 57028 108927 141569
 lowest : 0 2071 2522 3191 3325
 highest: 269057 269131 338955 357919 390460

totmctst : Total micro-cost

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 331 206 328 1 39022 6131 8283 14415 26323 54102 87495 111920
 lowest : 0 1562 2478 2626 3421
 highest: 144234 154709 198047 234876 271467

avtisst : Average TISS, Days 3-25

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 536 1 205 1 29.83 12.46 14.50 19.62 28.00 39.00 47.17 50.37
 lowest : 4.000 5.667 8.000 9.000 9.500
 highest: 58.500 59.000 60.000 61.000 64.000

race

n missing unique
 535 2 5

	white	black	asian	other	hispanic
Frequency	417	84	4	8	22
%	78	16	1	1	4

meanbp : Mean Arterial Blood Pressure Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 109 1 83.28 41.8 49.0 59.0 73.0 111.0 124.4 135.0
 lowest : 0 20 27 30 32, highest: 155 158 161 162 180

wbhc : White Blood Cell Count Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 532 5 241 1 14.1 0.8999 4.5000 7.9749 12.3984 18.1992 25.1891 30.1873
 lowest : 0.05000 0.06999 0.09999 0.14999 0.19998
 highest: 51.39844 58.19531 61.19531 79.39062 100.00000

hrt : Heart Rate Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 111 1 105 51 60 75 111 126 140 155
 lowest : 0 11 30 36 40, highest: 189 193 199 232 300

resp : Respiration Rate Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 45 1 23.72 8 10 12 24 32 39 40
 lowest : 0 4 6 7 8, highest: 48 49 52 60 64

temp : Temperature (celcius) Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 61 1 37.52 35.50 35.80 36.40 37.80 38.50 39.09 39.50
 lowest : 32.50 34.00 34.09 34.90 35.00
 highest: 40.20 40.59 40.90 41.00 41.20

pafi : PaO2/(.01*FiO2) Day 3

n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 500 37 357 1 227.2 86.99 105.08 137.88 202.56 290.00 390.49 433.31
 lowest : 45.00 48.00 53.33 54.00 55.00
 highest: 574.00 595.12 640.00 680.00 869.38

alb : Serum Albumin Day 3


n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 346 191 34 1 2.668 1.700 1.900 2.225 2.600 3.100 3.400 3.800
 lowest : 1.100 1.200 1.300 1.400 1.500
 highest: 4.100 4.199 4.500 4.699 4.800


bili : Bilirubin Day 3

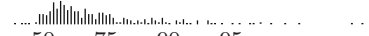
n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 386 151 88 1 2.678 0.3000 0.4000 0.6000 0.8999 2.0000 6.5996 13.1743
 lowest : 0.09999 0.19998 0.29999 0.39996 0.50000
 highest: 22.59766 30.00000 31.50000 35.00000 39.29688

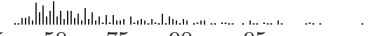
crea : Serum creatinine Day 3

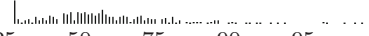
n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 84 1 2.232 0.6000 0.7000 0.8999 1.3999 2.5996 5.2395 7.3197
 lowest : 0.3 0.4 0.5 0.6 0.7, highest: 10.4 10.6 11.2 11.6 11.8

sod : Serum sodium Day 3 
 n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 537 0 38 1 138.1 129 131 134 137 142 147 150
 lowest : 118 120 121 126 127, highest: 156 157 158 168 175

ph : Serum pH (arterial) Day 3 
 n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 500 37 49 1 7.416 7.270 7.319 7.380 7.420 7.470 7.510 7.529
 lowest : 6.960 6.989 7.069 7.119 7.130
 highest: 7.560 7.569 7.590 7.600 7.659

glucose : Glucose Day 3 
 n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 297 240 179 1 167.7 76.0 89.0 106.0 141.0 200.0 292.4 347.2
 lowest : 30 42 52 55 68, highest: 446 468 492 576 598

bun : BUN Day 3 
 n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 304 233 100 1 38.91 8.00 11.00 16.75 30.00 56.00 79.70 100.70
 lowest : 1 3 4 5 6, highest: 123 124 125 128 146

urine : Urine Output Day 3 
 n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
 303 234 262 1 2095 20.3 364.0 1156.5 1870.0 2795.0 4008.6 4817.5
 lowest : 0 5 8 15 20, highest: 6865 6920 7360 7560 7750

adlp : ADL Patient Day 3
 n missing unique Info Mean
 104 433 8 0.87 1.577

	0	1	2	3	4	5	6	7
Frequency	51	19	7	6	4	7	8	2
%	49	18	7	6	4	7	8	2

adls : ADL Surrogate Day 3
 n missing unique Info Mean
 392 145 8 0.89 1.86

	0	1	2	3	4	5	6	7
Frequency	185	68	22	18	17	20	39	23
%	47	17	6	5	4	5	10	6

sfdm2
 n missing unique
 468 69 5

no(M2 and SIP pres) (134, 29%), adl>=4 (>=5 if sur) (78, 17%)
 SIP>=30 (30, 6%), Coma or Intub (5, 1%), <2 mo. follow-up (221, 47%)

```

adlsc : Imputed ADL Calibrated to Surrogate | . . . . .
      n missing unique Info Mean .05 .10 .25 .50 .75 .90 .95
      537      0    144 0.96 2.119 0.000 0.000 0.000 1.839 3.375 6.000 6.000

lowest : 0.0000 0.4948 0.4948 1.0000 1.1667
highest: 5.7832 6.0000 6.3398 6.4658 7.0000

```

Next, patterns of missing data are displayed.

```
plot(naclus(support[acute,])) # Figure 19.1
```

The `hmisc::varclus` function is used to quantify and depict associations between predictors, allowing for general nonmonotonic relationships. This is done by using Hoeffding's D as a similarity measure for all possible pairs of predictors instead of the default similarity, Spearman's ρ .

```

ac <- support[acute,]
ac$dzgroup <- ac$dzgroup[drop=TRUE] # Remove unused levels
label(ac$dzgroup) <- 'Disease Group'
attach(ac)
vc <- varclus(~ age + sex + dzgroup + num.co + edu + income +
              scoma + race + meanbp + wblc + hrt + resp +
              temp + pafi + alb + bili + crea + sod + ph +
              glucose + bun + urine + adlsc, sim='hoeffding')
plot(vc) # Figure 19.2

```

19.2 Checking Adequacy of Log-Normal Accelerated Failure Time Model

Let us check whether a parametric survival time model will fit the data, with respect to the key prognostic factors. First, Kaplan–Meier estimates stratified by disease group are computed, and plotted after inverse normal transformation, against $\log t$. Parallelism and linearity indicate goodness of fit to the log normal distribution for disease group. Then a more stringent assessment is made by fitting an initial model and computing right-censored residuals. These residuals, after dividing by $\hat{\sigma}$, should all have a normal distribution if the model holds. We compute Kaplan–Meier estimates of the distribution of the residuals and overlay the estimated survival distribution with the theoretical Gaussian one. This is done overall, and then to get more stringent assessments of fit, residuals are stratified by key predictors and plots are produced that contain multiple Kaplan–Meier curves along with a single theoretical normal curve. All curves should hover about the normal distribution. To gauge the natural variability of stratified residual distribution estimates, the residuals are also stratified by a random number that has no bearing on the goodness of fit.

```

dd <- datadist(ac)
# describe distributions of variables to rms

```

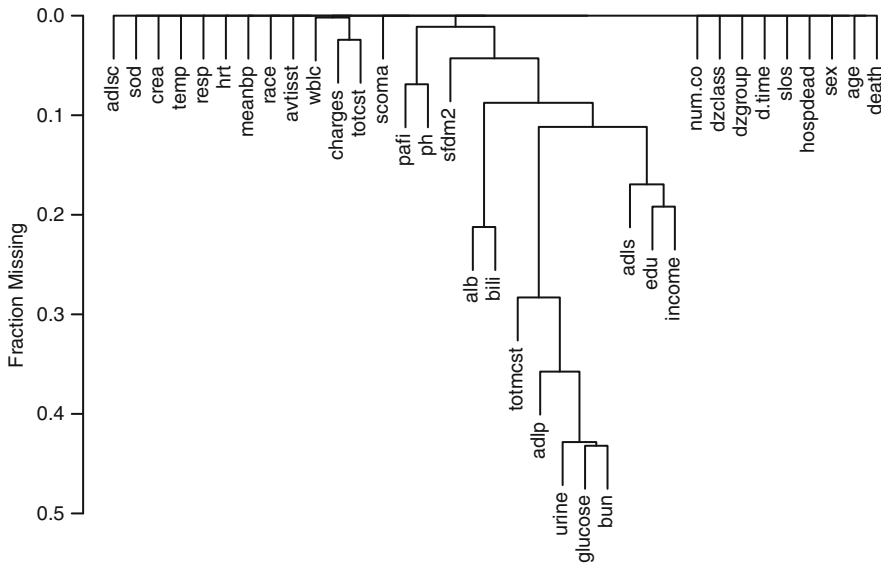


Fig. 19.1 Cluster analysis showing which predictors tend to be missing on the same patients

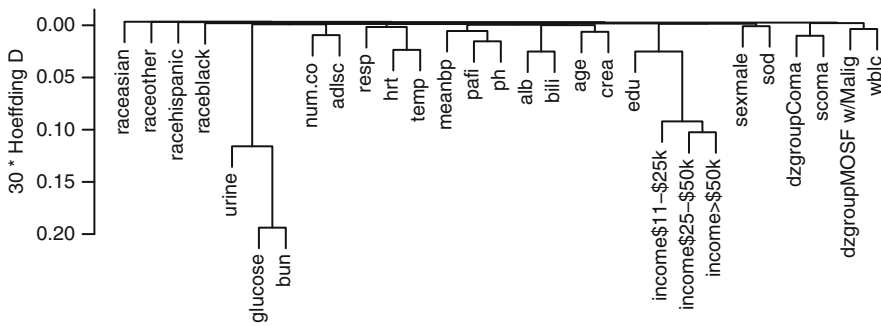


Fig. 19.2 Hierarchical clustering of potential predictors using Hoeffding D as a similarity measure. Categorical predictors are automatically expanded into dummy variables.

```
options(datadist='dd')

# Generate right-censored survival time variable
years <- d.time/365.25
units(years) <- 'Year'
S <- Surv(years, death)

# Show normal inverse Kaplan-Meier estimates
# stratified by dzgroup
survplot(npsurv(S ~ dzgroup), conf='none',
         fun=qnorm, logt=TRUE) # Figure 19.3
```

```
f ← psm(S ~ dzgroup + rcs(age,5) + rcs(meanbp,5),
        dist='lognormal', y=TRUE)
r ← resid(f)

survplot(r, dzgroup, label.curve=FALSE)
survplot(r, age, label.curve=FALSE)
survplot(r, meanbp, label.curve=FALSE)
random ← runif(length(age)); label(random) ← 'Random Number'
survplot(r, random, label.curve=FALSE) # Fig. 19.4
```

Now remove from consideration predictors that are missing in more than 0.2 of patients. Many of these were collected only for the second half of SUPPORT. Of those variables to be included in the model, find which ones have enough potential predictive power to justify allowing for nonlinear relationships or multiple categories, which spend more d.f. For each variable compute Spearman ρ^2 based on multiple linear regression of $\text{rank}(x)$, $\text{rank}(x)^2$, and the

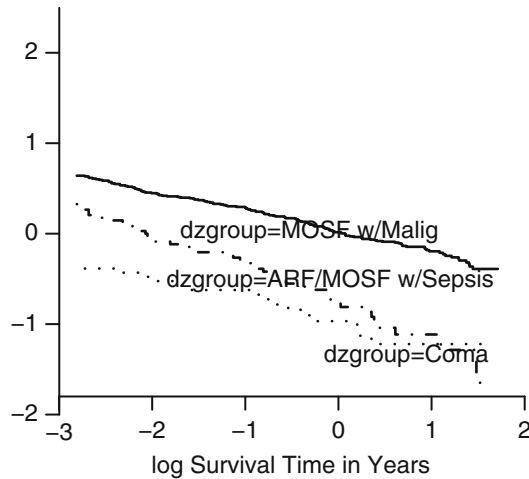


Fig. 19.3 $\Phi^{-1}(S_{KM}(t))$ stratified by `dzgroup`. Linearity and semi-parallelism indicate a reasonable fit to the log-normal accelerated failure time model with respect to one predictor.

survival time, truncating survival time at the shortest follow-up for survivors (356 days; see Section 4.1).

```
shortest.follow.up ← min(d.time[death==0], na.rm=TRUE)
d.timet ← pmin(d.time, shortest.follow.up)

w ← spearman2(d.timet ~ age + num.co + scoma + meanbp +
              hrt + resp + temp + crea + sod + adlsc +
              wblc + pafi + ph + dzgroup + race, p=2)
plot(w, main='') # Figure 19.5
```

A better approach is to use the complete information in the failure and censoring times by computing Somers' D_{xy} rank correlation allowing for censoring.

```
w ← rcorrcens(S ~ age + num.co + scoma + meanbp + hrt + resp +
             temp + crea + sod + adlsc + wblc + pafi + ph +
             dzgroup + race)
plot(w, main='') # Figure 19.6
```

Remaining missing values are imputed using the “most normal” values, a procedure found to work adequately for this particular study. Race is imputed using the modal category.

```
# Compute number of missing values per variable
sapply(1:11, function(i) sum(is.na(x[[i]])))
```

age	num.co	scoma	meanbp	hrt	resp	temp	crea	sod	adlsc
0	0	0	0	0	0	0	0	0	0
wblc	pafi	ph							
5	37	37							

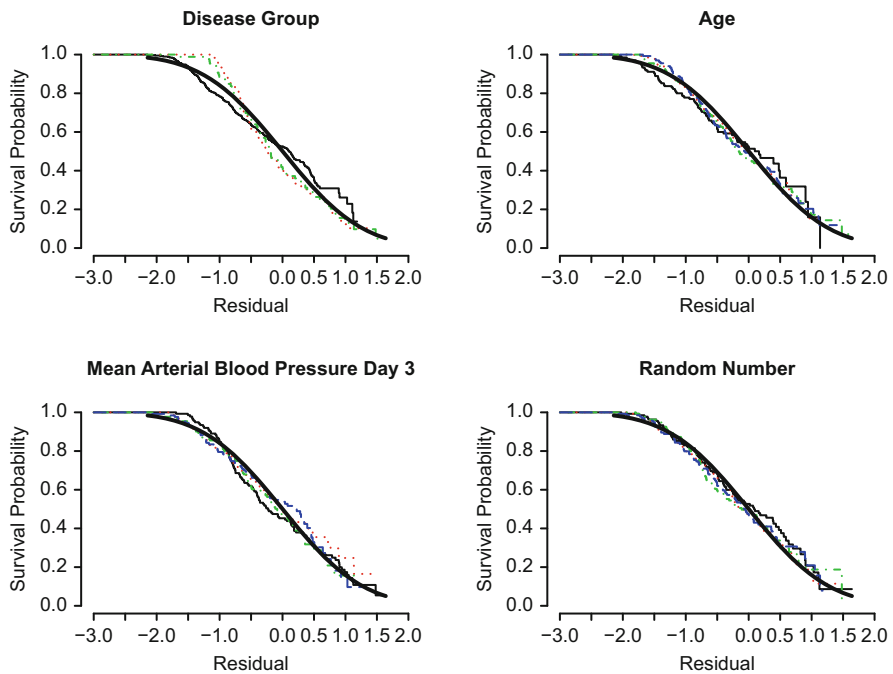


Fig. 19.4 Kaplan-Meier estimates of distributions of normalized, right-censored residuals from the fitted log-normal survival model. Residuals are stratified by important variables in the model (by quartiles of continuous variables), plus a random variable to depict the natural variability (in the lower right plot). Theoretical standard Gaussian distributions of residuals are shown with a thick solid line.

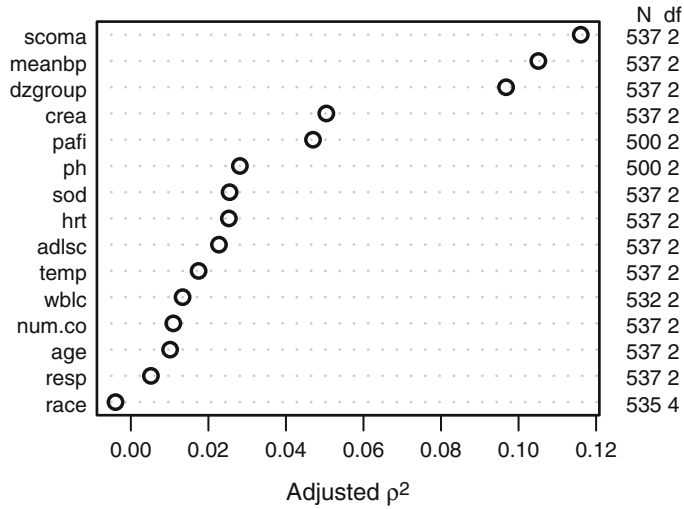


Fig. 19.5 Generalized Spearman ρ^2 rank correlation between predictors and truncated survival time

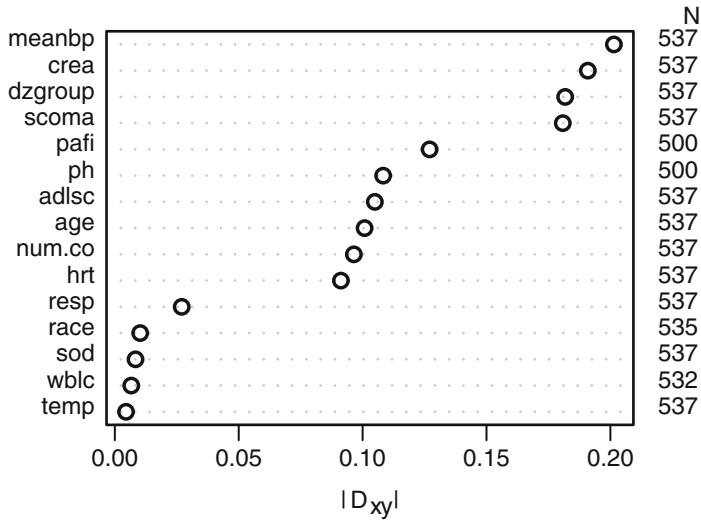


Fig. 19.6 Somers' D_{xy} rank correlation between predictors and original survival time. For dzgroup or race, the correlation coefficient is the maximum correlation from using a dummy variable to represent the most frequent or one to represent the second most frequent category. scap=Somers' D_{xy} rank correlation between predictors and original survival time

```
# Can also do naplot(naclus(support[acute,]))
# Can also use the Hmisc naclus and naplot functions
# Impute missing values with normal or modal values
wblc.i ← impute(wblc, 9)
pafi.i ← impute(pafi, 333.3)
ph.i ← impute(ph, 7.4)
race2 ← race
levels(race2) ← list(white='white', other=levels(race)[-1])
race2[is.na(race2)] ← 'white'
dd ← datadist(dd, wblc.i, pafi.i, ph.i, race2)
```

Now that missing values have been imputed, a formal multivariable redundancy analysis can be undertaken. The `Hmisc` package's `redun` function goes farther than the `varclus` pairwise correlation approach and allows for non-monotonic transformations in predicting each predictor from all the others.

```
redun(~ crea + age + sex + dzgroup + num.co + scoma + adlsc +
      race2 + meanbp + hrt + resp + temp + sod + wblc.i +
      pafi.i + ph.i, nk=4)
```

Redundancy Analysis

```
redun(formula = ~crea + age + sex + dzgroup + num.co + scoma +
      adlsc + race2 + meanbp + hrt + resp + temp + sod + wblc.i +
      pafi.i + ph.i, nk = 4)
```

```
n: 537  p: 16  nk: 4
```

```
Number of NAs: 0
```

```
Transformation of target variables forced to be linear
```

```
R2 cutoff: 0.9  Type: ordinary
```

```
R2 with which each variable can be predicted from all other variables:
```

crea	age	sex	dzgroup	num.co	scoma	adlsc	race2	meanbp
0.133	0.246	0.132	0.451	0.147	0.418	0.153	0.151	0.178
hrt	resp	temp	sod	wblc.i	pafi.i	ph.i		
0.258	0.131	0.197	0.135	0.093	0.143	0.171		

```
No redundant variables
```

Now turn to a more efficient approach for gauging the potential of each predictor, one that makes maximal use of failure time and censored data is to all continuous variables to have a maximum number of knots in a log-normal survival model. This approach must use imputation to have an adequate sample size. A semi-saturated main effects additive log-normal model is fitted. It is necessary to limit restricted cubic splines to 4 knots, force `scoma` to be linear, and to omit `ph.i` in order to avoid a singular covariance matrix in the fit.

```
k ← 4
f ← psm(S ~ rcs(age, k) + sex + dzgroup + pol(num.co, 2) + scoma +
      pol(adlsc, 2) + race + rcs(meanbp, k) + rcs(hrt, k) +
```

```

rcs(resp,k)+rcs(temp,k)+rcs(crea,3)+rcs(sod,k)+
rcs(wblc.i,k)+rcs(pafi.i,k), dist='lognormal')
plot(anova(f)) # Figure 19.7

```

Figure 19.7 properly blinds the analyst to the form of effects (tests of linearity). Next fit a log-normal survival model with number of parameters corresponding to nonlinear effects determined from the partial χ^2 tests in Figure 19.7. For the most promising predictors, five knots can be allocated, as there are fewer singularity problems once less promising predictors are simplified.

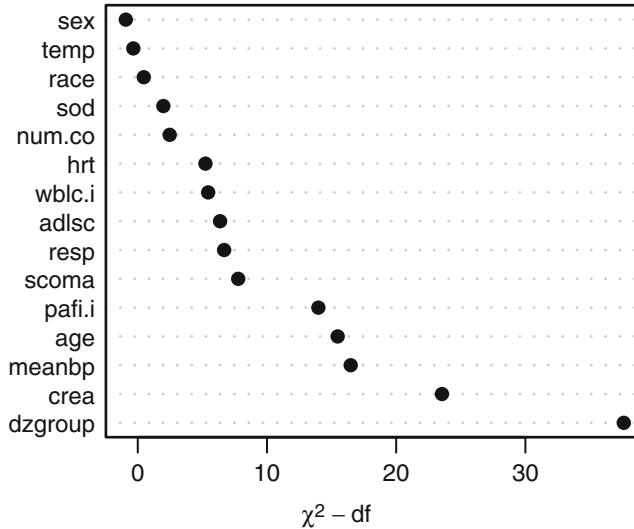


Fig. 19.7 Partial χ^2 statistics for association of each predictor with response from saturated main effects model, penalized for d.f.

```

f ← psm(S ~ rcs(age,5)+sex+dzgroup+num.co+
scoma+pol(adlsc,2)+race2+rcs(meanbp,5)+
rcs(hrt,3)+rcs(resp,3)+temp+
rcs(crea,4)+sod+rcs(wblc.i,3)+rcs(pafi.i,4),
dist='lognormal')
print(f, latex=TRUE, coefs=FALSE)

```

Parametric Survival Model: Log Normal Distribution

```

psm(formula = S ~ rcs(age, 5) + sex + dzgroup + num.co + scoma +
pol(adlsc, 2) + race2 + rcs(meanbp, 5) + rcs(hrt, 3) + rcs(resp,
3) + temp + rcs(crea, 4) + sod + rcs(wblc.i, 3) + rcs(pafi.i,
4), dist = "lognormal")

```

		Model Likelihood Ratio Test	Discrimination Indexes
Obs	537	LR χ^2 236.83	R^2 0.594
Events	356	d.f. 30	D_{xy} 0.485
σ	2.230782	$\Pr(> \chi^2) < 0.0001$	g 0.033
			g_r 1.959

```
a ← anova(f)
```

Table 19.1 Wald Statistics for S

	χ^2	d.f.	P
age	15.99	4	0.0030
<i>Nonlinear</i>	0.23	3	0.9722
sex	0.11	1	0.7354
dzgroup	45.69	2	< 0.0001
num.co	4.99	1	0.0255
scoma	10.58	1	0.0011
adlsc	8.28	2	0.0159
<i>Nonlinear</i>	3.31	1	0.0691
race2	1.26	1	0.2624
meanbp	27.62	4	< 0.0001
<i>Nonlinear</i>	10.51	3	0.0147
hrt	11.83	2	0.0027
<i>Nonlinear</i>	1.04	1	0.3090
resp	11.10	2	0.0039
<i>Nonlinear</i>	8.56	1	0.0034
temp	0.39	1	0.5308
crea	33.63	3	< 0.0001
<i>Nonlinear</i>	21.27	2	< 0.0001
sod	0.08	1	0.7792
wblc.i	5.47	2	0.0649
<i>Nonlinear</i>	5.46	1	0.0195
pafi.i	15.37	3	0.0015
<i>Nonlinear</i>	6.97	2	0.0307
TOTAL NONLINEAR	60.48	14	< 0.0001
TOTAL	261.47	30	< 0.0001

19.3 Summarizing the Fitted Model

First let's plot the shape of the effect of each predictor on log survival time. All effects are centered so that they can be placed on a common scale. This allows the relative strength of various predictors to be judged. Then Wald χ^2 statistics, penalized for d.f., are plotted in descending order. Next, relative effects of varying predictors over reasonable ranges (survival time ratios varying continuous predictors from the first to the third quartile) are charted.

```
ggplot(Predict(f, ref.zero=TRUE), vnames='names',
       sepdiscrte='vertical', anova=a) # Figure 19.8
```

```
latex(a, file='', label='tab:support-anovat') # Table 19.1
```

```
plot(a) # Figure 19.9
```

```
options(digits=3)
plot(summary(f), log=TRUE, main='') # Figure 19.10
```

19.4 Internal Validation of the Fitted Model Using the Bootstrap

Let us decide whether there was significant overfitting during the development of this model, using the bootstrap.

```
# First add data to model fit so bootstrap can re-sample
# from the data
g ← update(f, x=TRUE, y=TRUE)
set.seed(717)
latex(validate(g, B=300), digits=2, size='Ssize')
```

Index	Original Sample	Training Sample	Test Sample	Optimism	Corrected Index	<i>n</i>
D_{xy}	0.49	0.51	0.46	0.05	0.43	300
R^2	0.59	0.66	0.54	0.12	0.47	300
Intercept	0.00	0.00	-0.05	0.05	-0.05	300
Slope	1.00	1.00	0.90	0.10	0.90	300
D	0.48	0.55	0.42	0.13	0.35	300
U	0.00	0.00	-0.01	0.01	-0.01	300
Q	0.48	0.56	0.43	0.12	0.36	300
g	1.96	2.05	1.87	0.19	1.77	300

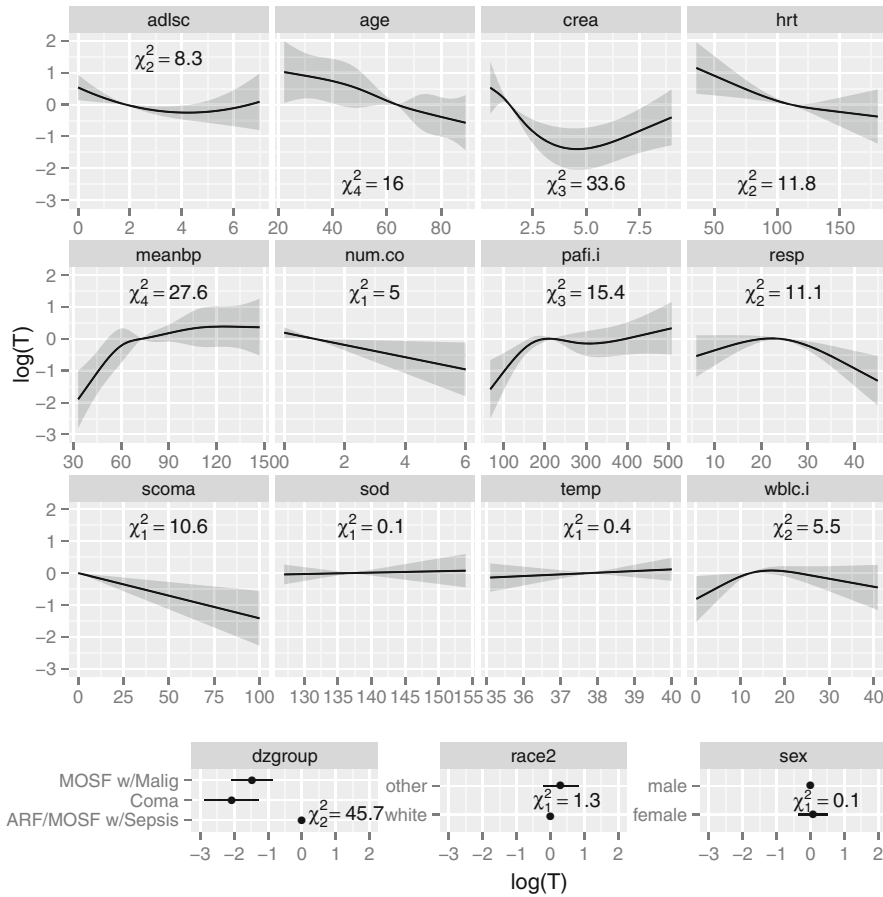


Fig. 19.8 Effect of each predictor on log survival time. Predicted values have been centered so that predictions at predictor reference values are zero. Pointwise 0.95 confidence bands are also shown. As all y -axes have the same scale, it is easy to see which predictors are strongest.

Judging from D_{xy} and R^2 there is a moderate amount of overfitting. The slope shrinkage factor (0.9) is not troublesome, however. An almost unbiased estimate of future predictive discrimination on similar patients is given by the corrected D_{xy} of 0.43. This index equals the difference between the probability of concordance and the probability of discordance of pairs of predicted survival times and pairs of observed survival times, accounting for censoring.

Next, a bootstrap overfitting-corrected calibration curve is estimated. Patients are stratified by the predicted probability of surviving one year, such that there are at least 60 patients in each group.

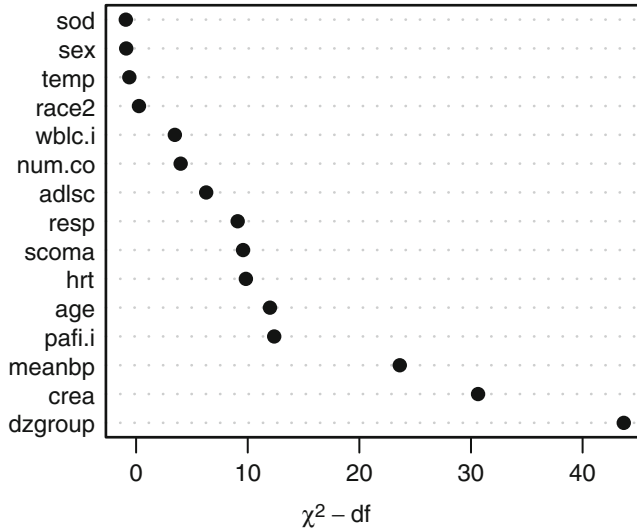


Fig. 19.9 Contribution of variables in predicting survival time in log-normal model

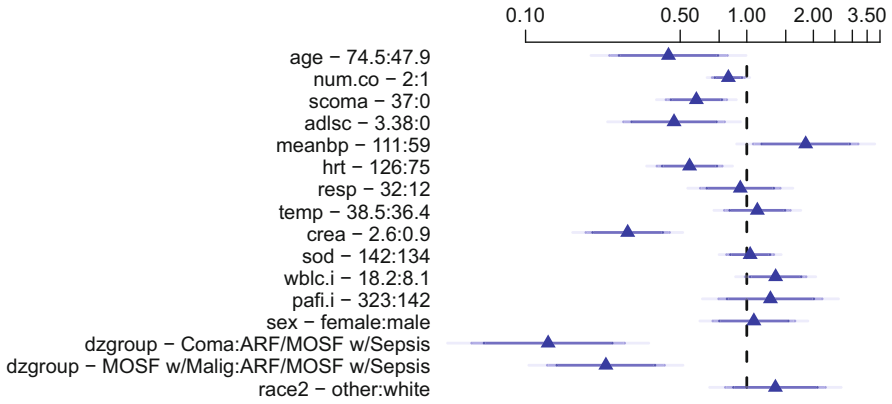


Fig. 19.10 Estimated survival time ratios for default settings of predictors. For example, when age changes from its lower quartile to the upper quartile (47.9y to 74.5y), median survival time decreases by more than half. Different shaded areas of bars indicate different confidence levels (.9, 0.95, 0.99).

```
set.seed(717)
cal ← calibrate(g, u=1, B=300)
plot(cal, subtitles=FALSE)
cal ← calibrate(g, cmethod='KM', u=1, m=60, B=120, pr=FALSE)
plot(cal, add=TRUE) # Figure 19.11
```

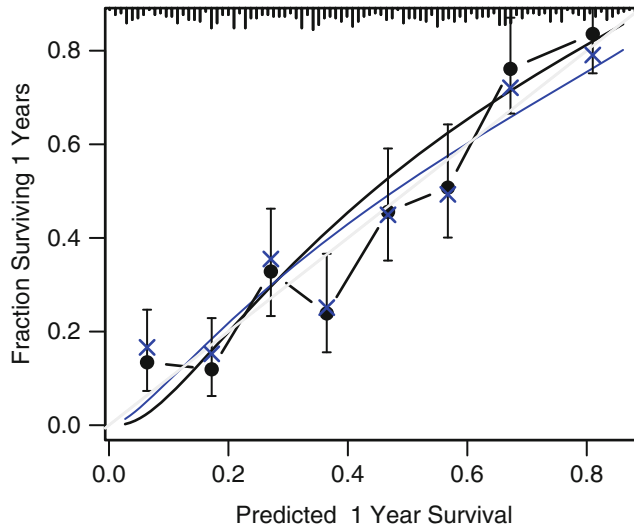


Fig. 19.11 Bootstrap validation of calibration curve. Dots represent apparent calibration accuracy; \times are bootstrap estimates corrected for overfitting, based on binning predicted survival probabilities and computing Kaplan-Meier estimates. Black curve is the estimated observed relationship using `harc` and the blue curve is the overfitting-corrected `harc` estimate. The gray-scale line depicts the ideal relationship.

19.5 Approximating the Full Model

The fitted log-normal model is perhaps too complex for routine use and for routine data collection. Let us develop a simplified model that can predict the predicted values of the full model with high accuracy ($R^2 = 0.967$). The simplification is done using a fast backward step-down against the full model predicted values.

```
Z <- predict(f)      # X*beta hat
a <- ols(Z ~ rcs(age,5)+sex+dzgroup+num.co+
          scoma+pol(adlsc,2)+race2+
          rcs(meanbp,5)+rcs(hrt,3)+rcs(resp,3)+
          temp+rcs(crea,4)+sod+rcs(wblc.i,3)+
          rcs(pafi.i,4), sigma=1)
# sigma=1 is used to prevent sigma hat from being zero when
# R2=1.0 since we start out by approximating Z with all
# component variables
fastbw(a, aics=10000) # fast backward stepdown
```

Deleted	Chi-Sq	d.f.	P	Residual	d.f.	P	AIC	R2
sod	0.43	1	0.512	0.43	1	0.5117	-1.57	1.000
sex	0.57	1	0.451	1.00	2	0.6073	-3.00	0.999
temp	2.20	1	0.138	3.20	3	0.3621	-2.80	0.998
race2	6.81	1	0.009	10.01	4	0.0402	2.01	0.994
wblc.i	29.52	2	0.000	39.53	6	0.0000	27.53	0.976


```

num.co  30.84 1    0.000  70.36  7    0.0000  56.36  0.957
resp    54.18 2    0.000  124.55  9    0.0000  106.55  0.924
adlsc   52.46 2    0.000  177.00 11    0.0000  155.00  0.892
pafi.i  66.78 3    0.000  243.79 14    0.0000  215.79  0.851
scoma   78.07 1    0.000  321.86 15    0.0000  291.86  0.803
hrt     83.17 2    0.000  405.02 17    0.0000  371.02  0.752
age     68.08 4    0.000  473.10 21    0.0000  431.10  0.710
crea    314.47 3    0.000  787.57 24    0.0000  739.57  0.517
meanbp  403.04 4    0.000 1190.61 28    0.0000 1134.61  0.270
dzgroup 441.28 2    0.000 1631.89 30    0.0000 1571.89  0.000

```

Approximate Estimates after Deleting Factors

```

          Coef      S.E. Wald Z P
[1,] -0.5928 0.04315 -13.74 0

```

Factors in Final Model

None

```

f.approx ← ols(Z ~ dzgroup + rcs(meanbp,5) + rcs(crea,4) +
              rcs(age,5) + rcs(hrt,3) + scoma +
              rcs(pafi.i,4) + pol(adlsc,2)+
              rcs(resp,3), x=TRUE)
f.approx$stats

```

n	Model L.R.	d.f.	R2	g	Sigma
537.000	1688.225	23.000	0.957	1.915	0.370

We can estimate the variance–covariance matrix of the coefficients of the reduced model using Equation 5.2 in Section 5.5.2. The computations below result in a covariance matrix that does not include elements related to the scale parameter. In the code x is the matrix T in Section 5.5.2.

```

V ← vcov(f, regcoef.only=TRUE) # var(full model)
X ← cbind(Intercept=1, g$x) # full model design
x ← cbind(Intercept=1, f.approx$x) # approx. model design
w ← solve(t(x) %*% x, t(x)) %*% X # contrast matrix
v ← w %*% V %*% t(w)

```

Let's compare the variance estimates (diagonals of v) with variance estimates from a reduced model that is fitted against the actual outcomes.

```

f.sub ← psm(S ~ dzgroup + rcs(meanbp,5) + rcs(crea,4) +
            rcs(age,5) + rcs(hrt,3) + scoma + rcs(pafi.i,4) +
            pol(adlsc,2)+ rcs(resp,3), dist='lognormal')
diag(v)/diag(vcov(f.sub, regcoef.only=TRUE))

```

Intercept	dzgroup=Coma	dzgroup=MOSF	w/Malig
0.981	0.979		0.979
meanbp	meanbp'		meanbp''
0.977	0.979		0.979
meanbp'''	crea		crea'
0.979	0.979		0.979
crea''	age		age'
0.979	0.982		0.981
age'''	age'''		hrt
0.981	0.980		0.978

hrt'	scoma	pafi.i
0.976	0.979	0.980
pafi.i'	pafi.i''	adlsc
0.980	0.980	0.981
adlsc^2	resp	resp'
0.981	0.978	0.977

```
r ← diag(v)/diag(vcov(f.sub, regcoef.only=TRUE))
r[c(which.min(r), which.max(r))]
```

hrt'	age
0.976	0.982

The estimated variances from the reduced model are actually slightly smaller than those that would have been obtained from stepwise variable selection in this case, had variable selection used a stopping rule that resulted in the same set of variables being selected. Now let us compute Wald statistics for the reduced model.

```
f.approx$var ← v
latex(anova(f.approx, test='Chisq', ss=FALSE), file='',
      label='tab:support-anovaa')
```

The results are shown in Table 19.2. Note the similarity of the statistics to those found in the table for the full model. This would not be the case had deleted variables been very collinear with retained variables.

The equation for the simplified model follows. The model is also depicted graphically in Figure 19.12. The nomogram allows one to calculate mean and median survival time. Survival probabilities could have easily been added as additional axes.

```
# Typeset mathematical form of approximate model
latex(f.approx, file='')
```

$$E(Z) = X\beta, \text{ where}$$

$$\begin{aligned}
 X\hat{\beta} = & \\
 & -2.51 \\
 & -1.94[\text{Coma}] - 1.75[\text{MOSF w/Malig}] \\
 & +0.068\text{meanbp} - 3.08 \times 10^{-5}(\text{meanbp} - 41.8)_+^3 + 7.9 \times 10^{-5}(\text{meanbp} - 61)_+^3 \\
 & -4.91 \times 10^{-5}(\text{meanbp} - 73)_+^3 + 2.61 \times 10^{-6}(\text{meanbp} - 109)_+^3 - 1.7 \times 10^{-6}(\text{meanbp} - 135)_+^3 \\
 & -0.553\text{crea} - 0.229(\text{crea} - 0.6)_+^3 + 0.45(\text{crea} - 1.1)_+^3 - 0.233(\text{crea} - 1.94)_+^3 \\
 & +0.0131(\text{crea} - 7.32)_+^3 \\
 & -0.0165\text{age} - 1.13 \times 10^{-5}(\text{age} - 28.5)_+^3 + 4.05 \times 10^{-5}(\text{age} - 49.5)_+^3 \\
 & -2.15 \times 10^{-5}(\text{age} - 63.7)_+^3 - 2.68 \times 10^{-5}(\text{age} - 72.7)_+^3 + 1.9 \times 10^{-5}(\text{age} - 85.6)_+^3 \\
 & -0.0136\text{hrt} + 6.09 \times 10^{-7}(\text{hrt} - 60)_+^3 - 1.68 \times 10^{-6}(\text{hrt} - 111)_+^3 + 1.07 \times 10^{-6}(\text{hrt} - 140)_+^3 \\
 & -0.0135\text{scoma} \\
 & +0.0161\text{pafi.i} - 4.77 \times 10^{-7}(\text{pafi.i} - 88)_+^3 + 9.11 \times 10^{-7}(\text{pafi.i} - 167)_+^3
 \end{aligned}$$

Table 19.2 Wald Statistics for Z

	χ^2	d.f.	P
dzgroup	55.94	2	< 0.0001
meanbp	29.87	4	< 0.0001
<i>Nonlinear</i>	9.84	3	0.0200
crea	39.04	3	< 0.0001
<i>Nonlinear</i>	24.37	2	< 0.0001
age	18.12	4	0.0012
<i>Nonlinear</i>	0.34	3	0.9517
hrt	9.87	2	0.0072
<i>Nonlinear</i>	0.40	1	0.5289
scoma	9.85	1	0.0017
pafi.i	14.01	3	0.0029
<i>Nonlinear</i>	6.66	2	0.0357
adlsc	9.71	2	0.0078
<i>Nonlinear</i>	2.87	1	0.0904
resp	9.65	2	0.0080
<i>Nonlinear</i>	7.13	1	0.0076
TOTAL NONLINEAR	58.08	13	< 0.0001
TOTAL	252.32	23	< 0.0001

$$\begin{aligned}
& -5.02 \times 10^{-7} (\text{pafi.i} - 276)_+^3 + 6.76 \times 10^{-8} (\text{pafi.i} - 426)_+^3 - 0.369 \text{ adlsc} + 0.0409 \text{ adlsc}^2 \\
& + 0.0394 \text{ resp} - 9.11 \times 10^{-5} (\text{resp} - 10)_+^3 + 0.000176 (\text{resp} - 24)_+^3 - 8.5 \times 10^{-5} (\text{resp} - 39)_+^3
\end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise; $(x)_+ = x$ if $x > 0$, 0 otherwise.

```

# Derive S functions that express mean and quantiles
# of survival time for specific linear predictors
# analytically
expected.surv ← Mean(f)
quantile.surv ← Quantile(f)
latex(expected.surv, file='', type='Sinput')

```

```

expected.surv ← function (lp = NULL,
                          parms = 0.802352037606488)
{
  names(parms) ← NULL
  exp(lp + exp(2 * parms)/2)
}

```

```

latex(quantile.surv, file='', type='Sinput')

```

```

quantile.surv ← function (q = 0.5, lp = NULL,
                          parms = 0.802352037606488)

```

```
{
  names(parms) ← NULL
  f ← function(lp, q, parms) lp + exp(parms) * qnorm(q)
  names(q) ← format(q)
  drop(exp(outer(lp, q, FUN = f, parms = parms)))
}
```

```
median.surv ← function(x) quantile.surv(lp=x)
```

```
# Improve variable labels for the nomogram
f.approx ← Newlabels(f.approx, c('Disease Group',
  'Mean Arterial BP', 'Creatinine', 'Age', 'Heart Rate',
  'SUPPORT Coma Score', 'PaO2/(.01*FiO2)', 'ADL',
  'Resp. Rate'))
nom ←
  nomogram(f.approx,
    pafi.i=c(0, 50, 100, 200, 300, 500, 600, 700, 800,
      900),
    fun=list('Median Survival Time'=median.surv,
      'Mean Survival Time' =expected.surv),
    fun.at=c(.1, .25, .5, 1, 2, 5, 10, 20, 40))
plot(nom, cex.var=1, cex.axis=.75, lmgp=.25)
# Figure 19.12
```

19.6 Problems

Analyze the Mayo Clinic PBC dataset.

1. Graphically assess whether Weibull (extreme value), exponential, log-logistic, or log-normal distributions will fit the data, using a few apparently important stratification factors.
2. For the best fitting parametric model from among the four examined, fit a model containing several sensible covariables, both categorical and continuous. Do a Wald test for whether each factor in the model has an association with survival time, and a likelihood ratio test for the simultaneous contribution of all predictors. For classification factors having more than two levels, be sure that the Wald test has the appropriate degrees of freedom. For continuous factors, verify or relax linearity assumptions. If using a Weibull model, test whether a simpler exponential model would be appropriate. Interpret all estimated coefficients in the model. Write the full survival model in mathematical form. Generate a predicted survival curve for a patient with a given set of characteristics.

See [361] for an analysis of this dataset using linear splines in time and in the covariables.

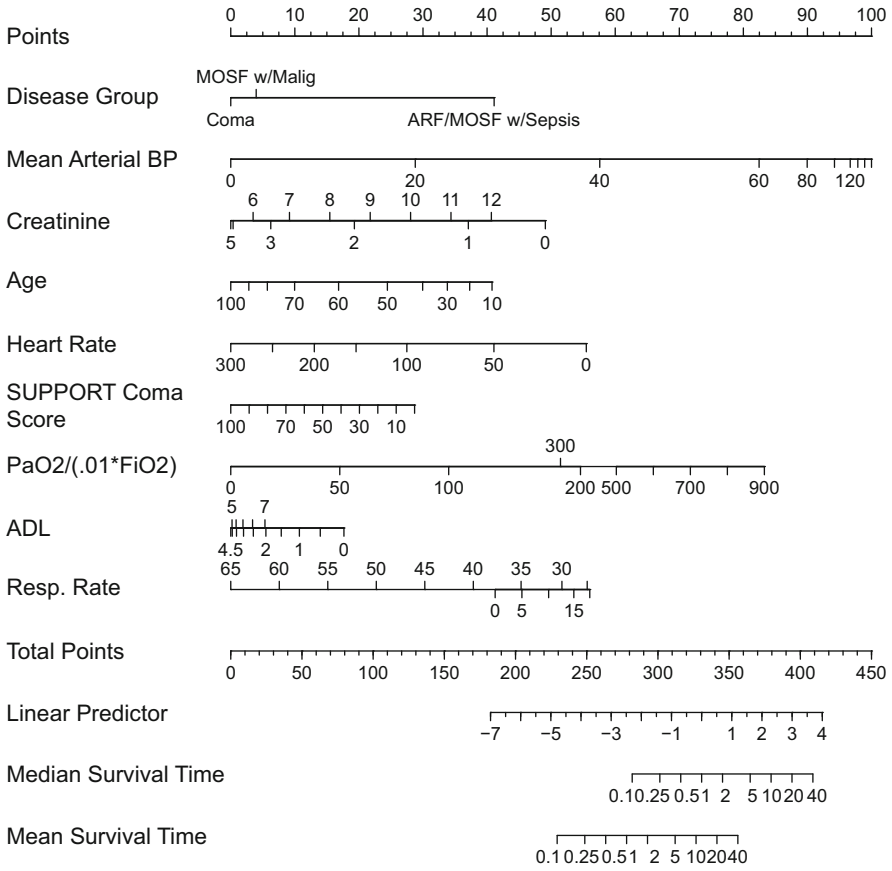


Fig. 19.12 Nomogram for predicting median and mean survival time, based on approximation of full model

Chapter 20

Cox Proportional Hazards Regression Model

20.1 Model

20.1.1 Preliminaries

The Cox proportional hazards model¹³² is the most popular model for the analysis of survival data. It is a semiparametric model; it makes a parametric assumption concerning the effect of the predictors on the hazard function, but makes no assumption regarding the nature of the hazard function $\lambda(t)$ itself. The Cox PH model assumes that predictors act multiplicatively on the hazard function but does not assume that the hazard function is constant (i.e., exponential model), Weibull, or any other particular form. The regression portion of the model is fully parametric; that is, the regressors are linearly related to log hazard or log cumulative hazard. In many situations, either the form of the true hazard function is unknown or it is complex, so the Cox model has definite advantages. Also, one is usually more interested in the effects of the predictors than in the shape of $\lambda(t)$, and the Cox approach allows the analyst to essentially ignore $\lambda(t)$, which is often not of primary interest.

1

The Cox PH model uses only the rank ordering of the failure and censoring times and thus is less affected by outliers in the failure times than fully parametric methods. The model contains as a special case the popular log-rank test for comparing survival of two groups. For estimating and testing regression coefficients, the Cox model is as efficient as parametric models (e.g., Weibull model with PH) even when all assumptions of the parametric model are satisfied.¹⁷¹

When a parametric model's assumptions are not true (e.g., when a Weibull model is used and the population is not from a Weibull survival distribution so that the choice of model is incorrect), the Cox analysis is more efficient

than the parametric analysis. As shown below, diagnostics for checking Cox model assumptions are very well developed.

20.1.2 Model Definition

The Cox PH model is most often stated in terms of the hazard function:

$$\lambda(t|X) = \lambda(t) \exp(X\beta). \quad (20.1)$$

We do not include an intercept parameter in $X\beta$ here. Note that this is identical to the parametric PH model stated earlier. There is an important difference, however, in that now we do not assume any specific shape for $\lambda(t)$. For the moment, we are not even interested in estimating $\lambda(t)$. The reason for this departure from the fully parametric approach is due to an ingenious conditional argument by Cox.¹³² Cox argued that when the PH model holds, information about $\lambda(t)$ is not very useful in estimating the parameters of primary interest, β . By special conditioning in formulating the log likelihood function, Cox showed how to derive a valid estimate of β that does not require estimation of $\lambda(t)$ as $\lambda(t)$ dropped out of the new likelihood function. Cox's derivation focuses on using the information in the data that relates to the relative hazard function $\exp(X\beta)$.

20.1.3 Estimation of β

Cox's derivation of an estimator of β can be loosely described as follows. Let $t_1 < t_2 < \dots < t_k$ represent the unique ordered failure times in the sample of n subjects; assume for now that there are no tied failure times (tied censoring times are allowed) so that $k = n$. Consider the set of individuals at risk of failing an instant before failure time t_i . This set of individuals is called the *risk set* at time t_i , and we use R_i to denote this risk set. R_i is the set of subjects j such that the subject had not failed or been censored by time t_i ; that is, the risk set R_i includes subjects with failure/censoring time $Y_j \geq t_i$.

The conditional probability that individual i is the one that failed at t_i , given that the subjects in the set R_i are at risk of failing, and given further that exactly one failure occurs at t_i , is

$$\begin{aligned} \text{Prob}\{\text{subject } i \text{ fails at } t_i | R_i \text{ and one failure at } t_i\} &= \\ &= \frac{\text{Prob}\{\text{subject } i \text{ fails at } t_i | R_i\}}{\text{Prob}\{\text{one failure at } t_i | R_i\}} \end{aligned} \quad (20.2)$$

using the rules of conditional probability. This conditional probability equals

$$\frac{\lambda(t_i) \exp(X_i\beta)}{\sum_{j \in R_i} \lambda(t_i) \exp(X_j\beta)} = \frac{\exp(X_i\beta)}{\sum_{j \in R_i} \exp(X_j\beta)} = \frac{\exp(X_i\beta)}{\sum_{Y_j \geq t_i} \exp(X_j\beta)} \quad (20.3)$$

independent of $\lambda(t)$. To understand this likelihood, consider a special case where the predictors have no effect; that is, $\beta = 0$ [93, pp. 48–49]. Then $\exp(X_i\beta) = \exp(X_j\beta) = 1$ and $\text{Prob}\{\text{subject } i \text{ is the subject that failed at } t_i | R_i \text{ and one failure occurred at } t_i\}$ is $1/n_i$ where n_i is the number of subjects at risk at time t_i .

By arguing that these conditional probabilities are themselves conditionally independent across the different failure times, a total likelihood can be computed by multiplying these individual likelihoods over all failure times. Cox termed this a *partial likelihood* for β :

$$L(\beta) = \prod_{Y_i \text{ uncensored}} \frac{\exp(X_i\beta)}{\sum_{Y_j \geq Y_i} \exp(X_j\beta)}. \quad (20.4)$$

The log partial likelihood is

$$\log L(\beta) = \sum_{Y_i \text{ uncensored}} \{X_i\beta - \log[\sum_{Y_j \geq Y_i} \exp(X_j\beta)]\}. \quad (20.5)$$

Cox and others have shown that this partial log likelihood can be treated as an ordinary log likelihood to derive valid (partial) MLEs of β . Note that this log likelihood is unaffected by the addition of a constant to any or all of the X s. This is consistent with the fact that an intercept term is unnecessary and cannot be estimated since the Cox model is a model for the relative hazard and does not directly estimate the underlying hazard $\lambda(t)$.

When there are tied failure times in the sample, the true partial log likelihood function involves permutations so it can be time-consuming to compute. When the number of ties is not large, Breslow⁷⁰ has derived a satisfactory approximate log likelihood function. The formula given above, when applied without modification to samples containing ties, actually uses Breslow's approximation. If there are ties so that $k < n$ and t_1, \dots, t_k denote the unique failure times as we originally intended, Breslow's approximation is written as

$$\log L(\beta) = \sum_{i=1}^k \{S_i\beta - d_i \log[\sum_{Y_j \geq t_i} \exp(X_j\beta)]\}, \quad (20.6)$$

where $S_i = \sum_{j \in D_i} X_j$, D_i is the set of indexes j for subjects failing at time t_i , and d_i is the number of failures at t_i .

Efron¹⁷¹ derived another approximation to the true likelihood that is significantly more accurate than the Breslow approximation and often yields estimates that are very close to those from the more cumbersome permutation likelihood.²⁸⁸

$$\begin{aligned} \log L(\beta) = & \sum_{i=1}^k \{S_i\beta - \sum_{j=1}^{d_i} \log[\sum_{Y_j \geq t_i} \exp(X_j\beta) \\ & - \frac{j-1}{d_i} \sum_{l \in D_i} \exp(X_l\beta)]\}. \end{aligned} \quad (20.7)$$

In the special case when all tied failure times are from subjects with identical $X_i\beta$, the Efron approximation yields the exact (permutation) marginal likelihood (Therneau, personal communication, 1993).

Kalbfleisch and Prentice³³⁰ showed that Cox's partial likelihood, in the absence of predictors that are functions of time, is a marginal distribution of the *ranks* of the failure/censoring times.

See Therneau and Grambsch⁶⁰⁴ and Huang and Harrington³¹⁰ for descriptions of penalized partial likelihood estimation methods for improving mean squared error of estimates of β in a similar fashion to what was discussed in Section 9.10.

20.1.4 Model Assumptions and Interpretation of Parameters

The Cox PH regression model has the same assumptions as the parametric PH model except that no assumption is made regarding the shape of the underlying hazard or survival functions $\lambda(t)$ and $S(t)$. The Cox PH model assumes, in its most basic form, linearity and additivity of the predictors with respect to log hazard or log cumulative hazard. It also assumes the PH assumption of no time by predictor interactions; that is, the predictors have the same effect on the hazard function at all values of t . The relative hazard function $\exp(X\beta)$ is constant through time and the survival functions for subjects with different values of X are powers of each other. If, for example, the hazard of death at time t for treated patients is half that of control patients at time t , this same hazard ratio is in effect at any other time point. In other words, treated patients have a consistently better hazard of death over all follow-up time.

The regression parameters are interpreted the same as in the parametric PH model. The only difference is the absence of hazard shape parameters in the model, since the hazard shape is not estimated in the Cox partial likelihood procedure.

20.1.5 Example

Consider again the rat vaginal cancer data from Section 18.3.6. Figure 20.1 displays the nonparametric survival estimates for the two groups along with estimates derived from the Cox model (by a method discussed later).

```
require(rms)
```

```

group ← c(rep('Group 1',19),rep('Group 2',21))
group ← factor(group)
dd ← datadist(group); options(datadist='dd')
days ←
  c(143,164,188,188,190,192,206,209,213,216,220,227,230,
    234,246,265,304,216,244,142,156,163,198,205,232,232,
    233,233,233,233,239,240,261,280,280,296,296,323,204,344)
death ← rep(1,40)
death[c(18,19,39,40)] ← 0
units(days) ← 'Day'
df ← data.frame(days, death, group)
S ← Surv(days, death)

f ← npsurv(S ~ group, type='fleming')
for(meth in c('exact', 'breslow', 'efron')) {
  g ← cph(S ~ group, method=meth, surv=TRUE, x=TRUE, y=TRUE)
  # print(g) to see results
}
f.exp ← psm(S ~ group, dist='exponential')
fw ← psm(S ~ group, dist='weibull')
phform ← pphsm(fw)

```

```

co ← gray(c(0, .8))
survplot(f, lty=c(1, 1), lwd=c(1, 3), col=co,
  label.curves=FALSE, conf='none')
survplot(g, lty=c(3, 3), lwd=c(1, 3), col=co, # Efron approx.
  add=TRUE, label.curves=FALSE, conf.type='none')
legend(c(2, 160), c(.38, .54),
  c('Nonparametric Estimates', 'Cox-Breslow Estimates'),
  lty=c(1, 3), cex=.8, bty='n')
legend(c(2, 160), c(.18, .34), cex=.8,
  c('Group 1', 'Group 2'), lwd=c(1,3), col=co, bty='n')

```

The predicted survival curves from the fitted Cox model are in good agreement with the nonparametric estimates, again verifying the PH assumption for these data. The estimates of the group effect from a Cox model (using the exact likelihood since there are ties, along with both Efron's and Breslow's approximations) as well as from a Weibull model and an exponential model are shown in Table 20.1. The exponential model, with its constant hazard, cannot accommodate the long early period with no failures. The group predictor was coded as $X_1 = 0$ and $X_1 = 1$ for Groups 1 and 2, respectively. For this example, the Breslow likelihood approximation resulted in $\hat{\beta}$ closer to that from maximizing the exact likelihood. Note how the group effect (47% reduction in hazard of death by the exact Cox model) is underestimated by the exponential model (9% reduction in hazard). The hazard ratio from the Weibull fit agrees with the Cox fit.

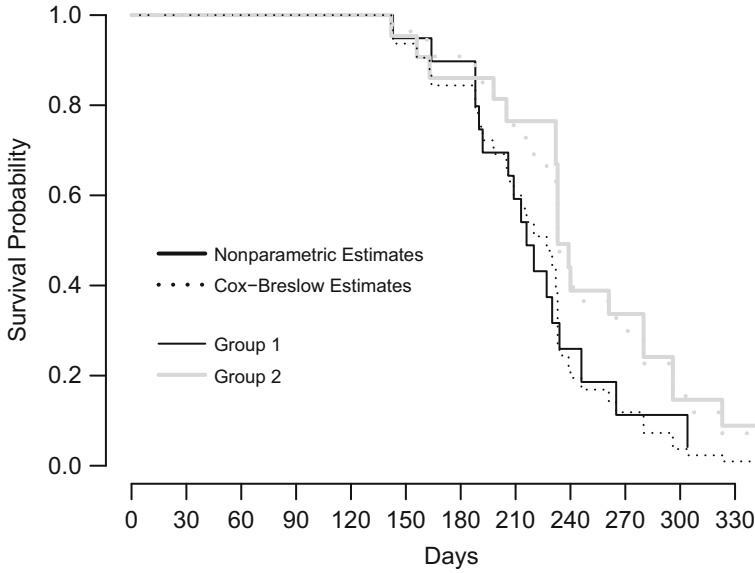


Fig. 20.1 Altschuler–Nelson–Fleming–Harrington nonparametric survival estimates and Cox-Breslow estimates for rat data⁵⁰⁸

Table 20.1 Group effects using three versions of the partial likelihood and three parametric models

Model	Group Regression Coefficient	S.E.	Wald <i>P</i> -Value	Group 2:1 Hazard Ratio
Cox (Exact)	−0.629	0.361	0.08	0.533
Cox (Efron)	−0.569	0.347	0.10	0.566
Cox (Breslow)	−0.596	0.348	0.09	0.551
Exponential	−0.093	0.334	0.78	0.911
Weibull (AFT)	0.132	0.061	0.03	–
Weibull (PH)	−0.721	–	–	0.486

20.1.6 Design Formulations

Designs are no different for the Cox PH model than for other models except for one minor distinction. Since the Cox model does not have an intercept parameter, the group omitted from X in an ANOVA model will go into the underlying hazard function. As an example, consider a three-group model for treatments A, B, and C. We use the two dummy variables

$$\begin{aligned}
 X_1 &= 1 \quad \text{if treatment is A, } 0 \text{ otherwise, and} \\
 X_2 &= 1 \quad \text{if treatment is B, } 0 \text{ otherwise.}
 \end{aligned}$$

The parameter β_1 is the A : C log hazard ratio or difference in hazards at any time t between treatment A and treatment C. β_2 is the B : C log hazard ratio ($\exp(\beta_2)$ is the B : C hazard ratio, etc.). Since there is no intercept parameter, there is no direct estimate of the hazard function for treatment C or any other treatment; only relative hazards are modeled.

As with all regression models, a Wald, score, or likelihood ratio test for differences between any treatments is conducted by testing $H_0 : \beta_1 = \beta_2 = 0$ with 2 d.f.

20.1.7 Extending the Model by Stratification

A unique feature of the Cox PH model is its ability to adjust for factors that are not modeled. Such factors usually take the form of polytomous stratification factors that either are too difficult to model or do not satisfy the PH assumption. For example, a subject's occupation or clinical study site may take on dozens of levels and the sample size may not be large enough to model this nominal variable with dozens of dummy variables. Also, one may know that a certain predictor (either a polytomous one or a continuous one that is grouped) may not satisfy PH and it may be too complex to model the hazard ratio for that predictor as a function of time.

The idea behind the *stratified* Cox PH model is to allow the form of the underlying hazard function to vary across levels of the stratification factors. A stratified Cox analysis ranks the failure times separately within strata. Suppose that there are b strata indexed by $j = 1, 2, \dots, b$. Let C denote the stratum identification. For example, $C = 1$ or 2 may stand for the female and male strata, respectively. The stratified PH model is

$$\begin{aligned}\lambda(t|X, C = j) &= \lambda_j(t) \exp(X\beta), \quad \text{or} \\ S(t|X, C = j) &= S_j(t)^{\exp(X\beta)}.\end{aligned}\tag{20.8}$$

Here $\lambda_j(t)$ and $S_j(t)$ are, respectively, the underlying hazard and survival functions for the j th stratum. The model does not assume any connection between the shapes of these functions for different strata.

In this stratified analysis, the data are stratified by C but, by default, a common vector of regression coefficients is fitted across strata. These common regression coefficients can be thought of as “pooled” estimates. For example, a Cox model with age as a (modeled) predictor and sex as a stratification variable essentially estimates the common slope of age by pooling information about the age effect over the two sexes. The effect of age is adjusted by sex differences, but no assumption is made about how sex affects survival. There is no PH assumption for sex. Levels of the stratification factor C can represent multiple stratification factors that are cross-classified. Since these factors are not modeled, no assumption is made regarding interactions among them.

At first glance it appears that stratification causes a loss of efficiency. However, in most cases the loss is small as long as the number of strata is not too large with regard to the total number of events. A stratum that contains no events contributes no information to the analysis, so such a situation should be avoided if possible.

The stratified or “pooled” Cox model is fitted by formulating a separate log likelihood function for each stratum, but with each log likelihood having a common β vector. If different strata are made up of independent subjects, the strata are independent and the likelihood functions are multiplied together to form a joint likelihood over strata. Log likelihood functions are thus added over strata. This total log likelihood function is maximized once to derive a pooled or stratified estimate of β and to make an inference about β . No inference can be made about the stratification factors. They are merely “adjusted for.”

Stratification is useful for checking the PH and linearity assumptions for one or more predictors. Predicted Cox survival curves (Section 20.2) can be derived by modeling the predictors in the usual way, and then stratified survival curves can be estimated by using those predictors as stratification factors. Other factors for which PH is assumed can be modeled in both instances. By comparing the modeled versus stratified survival estimates, a graphical check of the assumptions can be made. Figure 20.1 demonstrates this method although there are no other factors being adjusted for and stratified Cox estimates are KM estimates. The stratified survival estimates are derived by stratifying the dataset to obtain a separate underlying survival curve for each stratum, while pooling information across strata to estimate coefficients of factors that are modeled.

Besides allowing a factor to be adjusted for without modeling its effect, a stratified Cox PH model can also allow a modeled factor to interact with strata.^{143, 180, 603} For the age–sex example, consider the following model with X_1 denoting age and $C = 1, 2$ denoting females and males, respectively.

$$\begin{aligned}\lambda(t|X_1, C = 1) &= \lambda_1(t) \exp(\beta_1 X_1) \\ \lambda(t|X_1, C = 2) &= \lambda_2(t) \exp(\beta_1 X_1 + \beta_2 X_1).\end{aligned}\tag{20.9}$$

This model can be simplified to

$$\lambda(t|X_1, C = j) = \lambda_j(t) \exp(\beta_1 X_1 + \beta_2 X_2)\tag{20.10}$$

if X_2 is a product interaction term equal to 0 for females and X_1 for males. The β_2 parameter quantifies the interaction between age and sex: it is the difference in the age slope between males and females. Thus the interaction between age and sex can be quantified and tested, even though the effect of sex is not modeled!

The stratified Cox model is commonly used to adjust for hospital differences in a multicenter randomized trial. With this method, one can allow

for differences in outcome between q hospitals without estimating $q - 1$ parameters. Treatment \times hospital interactions can be tested efficiently without computational problems by estimating only the treatment main effect, after stratifying on hospital. The score statistic (with $q - 1$ d.f.) for testing $q - 1$ treatment \times hospital interaction terms is then computed (“residual χ^2 ” in a stepwise procedure with treatment \times hospital terms as candidate predictors).

The stratified Cox model turns out to be a generalization of the conditional logistic model for analyzing matched set (e.g., case-control) data.⁷¹ Each stratum represents a set, and the number of “failures” in the set is the number of “cases” in that set. For $r : 1$ matching (r may vary across sets), the Breslow⁷⁰ likelihood may be used to fit the conditional logistic model exactly. For $r : m$ matching, an exact Cox likelihood must be computed.

20.2 Estimation of Survival Probability and Secondary Parameters

As discussed above, once a partial log likelihood function is derived, it is used as if it were an ordinary log likelihood function to estimate β , estimate standard errors of β , obtain confidence limits, and make statistical tests. Point and interval estimates of hazard ratios are obtained in the same fashion as with parametric PH models discussed earlier.

The Cox model and parametric survival models differ markedly in how one estimates $S(t|X)$. Since the Cox model does not depend on a choice of the underlying survival function $S(t)$, fitting a Cox model does not result directly in an estimate of $S(t|X)$. However, several authors have derived secondary estimates of $S(t|X)$. One method is the *discrete hazard model* of Kalbfleisch and Prentice [331, pp. 36–37, 84–87]. Their estimator has two advantages: it is an extension of the Kaplan–Meier estimator and is identical to S_{KM} if the estimated value of β happened to be zero or there are no covariables being modeled; and it is not affected by the choice of what constitutes a “standard” subject having the underlying survival function $S(t)$. In other words, it would not matter whether the standard subject is one having age equal to the mean age in the sample or the median age in the sample; the estimate of $S(t|X)$ as a function of $X = \text{age}$ would be the same (this is also true of another estimator which follows).

Let t_1, t_2, \dots, t_k denote the unique failure times in the sample. The discrete hazard model assumes that the probability of failure is greater than zero only at observed failure times. The probability of failure at time t_j given that the subject has not failed before that time is also the hazard of failure at time t_j since the model is discrete. The hazard at t_j for the standard subject is written λ_j . Letting $\alpha_j = 1 - \lambda_j$, the underlying survival function can be written

$$S(t_i) = \prod_{j=0}^{i-1} \alpha_j, i = 1, 2, \dots, k \quad (\alpha_0 = 1). \quad (20.11)$$

A separate equation can be solved using the Newton–Raphson method to estimate each α_j . If there is only one failure at time t_i , there is a closed-form solution for the maximum likelihood estimate of α_i , $\hat{\alpha}_i$, letting j denote the subject who failed at t_i . $\hat{\beta}$ denotes the partial MLE of β .

$$\hat{\alpha}_i = [1 - \exp(X_j \hat{\beta}) / \sum_{Y_m \geq Y_j} \exp(X_m \hat{\beta})] \exp(-X_j \hat{\beta}). \quad (20.12)$$

If $\hat{\beta} = 0$, this formula reduces to a conditional probability component of the product-limit estimator, $1 - (1/\text{number at risk})$.

The estimator of the underlying survival function is

$$\hat{S}(t) = \prod_{j:t_j \leq t} \hat{\alpha}_j, \quad (20.13)$$

and the estimate of the probability of survival past time t for a subject with predictor values X is

$$\hat{S}(t|X) = \hat{S}(t)^{\exp(X \hat{\beta})}. \quad (20.14)$$

When the model is stratified, estimation of the α_j and S is carried out separately within each stratum once $\hat{\beta}$ is obtained by pooling over strata. The stratified survival function estimates can be thought of as stratified Kaplan–Meier estimates adjusted for X , with the adjustment made by assuming PH and linearity. As mentioned previously, these stratified adjusted survival estimates are useful for checking model assumptions and for providing a simple way to incorporate factors that violate PH.

The stratified estimates are also useful in themselves as descriptive statistics without making assumptions about a major factor. For example, in a study from Califf et al.⁸⁸ to compare medical therapy with coronary artery bypass grafting (CABG), the model was stratified by treatment but adjusted for a variety of baseline characteristics by modeling. These adjusted survival estimates do not assume a form for the effect of surgery. Figure 20.2 displays unadjusted (Kaplan–Meier) and adjusted survival curves, with baseline predictors adjusted to their mean levels in the combined sample. Notice that valid adjusted survival estimates are obtained even though the curves cross (i.e., PH is violated for the treatment variable). These curves are essentially product limit estimates with respect to treatment and Cox PH estimates with respect to the baseline descriptor variables.

The Kalbfleisch–Prentice discrete underlying hazard model estimates of the α_j are one minus estimates of the hazard function at the discrete failure times. However, these estimated hazard functions are usually too “noisy” to be useful unless the sample size is very large or the failure times have been grouped (say by rounding).

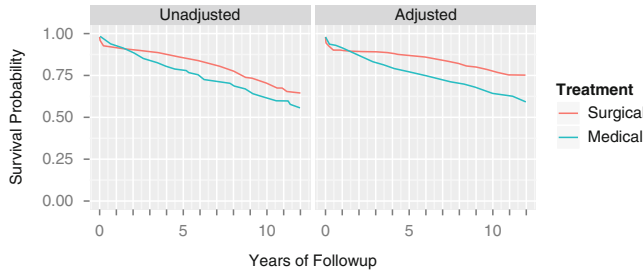


Fig. 20.2 Unadjusted (Kaplan–Meier) and adjusted (Cox–Kalbfleisch–Prentice) estimates of survival. Left, Kaplan–Meier estimates for patients treated medically and surgically at Duke University Medical Center from November 1969 through December 1984. These survival curves are not adjusted for baseline prognostic factors. Right, survival curves for patients treated medically or surgically after adjusting for all known important baseline prognostic characteristics.⁸⁸

Just as Kalbfleisch and Prentice have generalized the Kaplan–Meier estimator to allow for covariables, Breslow⁷⁰ has generalized the Altschuler–Nelson–Aalen–Fleming–Harrington estimator to allow for covariables. Using the notation in Section 20.1.3, Breslow’s estimate is derived through an estimate of the cumulative hazard function:

$$\hat{\Lambda}(t) = \sum_{i:t_i < t} \frac{d_i}{\sum_{Y_i \geq t_i} \exp(X_i \hat{\beta})}. \tag{20.15}$$

For any X , the estimates of Λ and S are

$$\begin{aligned} \hat{\Lambda}(t|X) &= \hat{\Lambda}(t) \exp(X \hat{\beta}) \\ \hat{S}(t|X) &= \exp[-\hat{\Lambda}(t) \exp(X \hat{\beta})]. \end{aligned} \tag{20.16}$$

More asymptotic theory has been derived from the Breslow estimator than for the Kalbfleisch–Prentice estimator. Another advantage of the Breslow estimator is that it does not require iterative computations for $d_i > 1$. Lawless [382, p. 362] states that the two survival function estimators differ little except in the right-hand tail when all d_i s are unity. Like the Kalbfleisch–Prentice estimator, the Breslow estimator is invariant under different choices of “standard subjects” for the underlying survival $S(t)$.

Somewhat complex formulas are available for computing confidence limits of $\hat{S}(t|X)$.⁶¹⁵

2

3

20.3 Sample Size Considerations

One way of estimating the minimum sample size for a Cox model analysis aimed at estimating survival probabilities is to consider the simplest case where there are no covariates. Thus the problem reduces to using the Kaplan-Meier estimate to estimate $S(t)$. Let's further simplify things to assume there is no censoring. Then the Kaplan-Meier estimate is just one minus the empirical cumulative distribution function. By the Dvoretzky-Kiefer-Wolfowitz inequality, the maximum absolute error in an empirical distribution function estimate of the true continuous distribution function is less than or equal to ϵ with probability of at least $1 - 2e^{-2n\epsilon^2}$. For the probability to be at least 0.95, $n = 184$. Thus in the case of no censoring, one needs 184 subjects to estimate the survival curve to within a margin of error of 0.1 everywhere. To estimate the subject-specific survival curves ($S(t|X)$) will require greater sample sizes, as will having censored data. It is a fair approximation to think of 184 as the needed number of subjects suffering the event or being censored "late."

Turning to estimation of a hazard ratio for a single binary predictor X that has equal numbers of $X = 0$ and $X = 1$, if the total sample size is n and the number of events in the two categories are respectively e_0 and e_1 , the variance of the log hazard ratio is approximately $v = \frac{1}{e_0} + \frac{1}{e_1}$. Letting z denote the $1 - \alpha/2$ standard normal critical value, the multiplicative margin of error (MMOE) with confidence $1 - \alpha$ is given by $\exp(z\sqrt{v})$. To achieve a MMOE of 1.2 in estimating $e^{\hat{\beta}}$ with equal numbers of events in the two groups and $\alpha = 0.05$ requires a total of 462 events.

20.4 Test Statistics

Wald, score, and likelihood ratio statistics are useful and valid for drawing inferences about β in the Cox model. The score test deserves special mention here. If there is a single binary predictor in the model that describes two groups, the score test for assessing the importance of the binary predictor is virtually identical to the Mantel-Haenszel log-rank test for comparing the two groups. If the analysis is stratified for other (nonmodeled) factors, the score test from a stratified Cox model is equivalent to the corresponding stratified log-rank test. Of course, the likelihood ratio or Wald tests could also be used in this situation, and in fact the likelihood ratio test may be better than the score test (i.e., type I errors by treating the likelihood ratio test statistic as having a χ^2 distribution may be more accurate than using the log-rank statistic).

The Cox model can be thought of as a generalization of the log-rank procedure since it allows one to test continuous predictors, perform simultaneous

tests of various predictors, and adjust for other continuous factors without grouping them. Although a stratified log-rank test does not make assumptions regarding the effect of the adjustment (stratifying) factors, it makes the same assumption (i.e., PH) as the Cox model regarding the treatment effect for the statistical test of no difference in survival between groups.

20.5 Residuals

Therneau et al.⁶⁰⁵ discussed four types of residuals from the Cox model: martingale, score, Schoenfeld, and deviance. The first three have been proven to be very useful, as indicated in Table 20.2.

4

Table 20.2 Types of residuals for the Cox model

Residual	Purposes
Martingale	Assessing adequacy of a hypothesized predictor transformation. Graphing an estimate of a predictor transformation (Section 20.6.1).
Score	Detecting overly influential observations (Section 20.9). Robust estimate of covariance matrix of $\hat{\beta}$ (Section 9.5). ⁴¹⁰
Schoenfeld	Testing PH assumption (Section 20.6.2). Graphing estimate of hazard ratio function (Section 20.6.2).

20.6 Assessment of Model Fit

As stated before, the Cox model makes the same assumptions as the parametric PH model except that it does not assume a given shape for $\lambda(t)$ or $S(t)$. Because the Cox PH model is so widely used, methods of assessing its fit are dealt with in more detail than was done with the parametric PH models.

20.6.1 Regression Assumptions

Regression assumptions (linearity, additivity) for the PH model are displayed in Figures 18.3 and 18.5. As mentioned earlier, the regression assumptions can be verified by stratifying by X and examining $\log \hat{A}(t|X)$ or $\log[\Lambda_{KM}(t|X)]$ estimates as a function of X at fixed time t . However, as was pointed out

in logistic regression, the stratification method is prone to problems of high variability of estimates. The sample size must be moderately large before estimates are precise enough to observe trends through the “noise.” If one wished to divide the sample by quintiles of age and 15 events were thought to be needed in each stratum to derive a reliable estimate of $\log[A_{KM}(2 \text{ years})]$, there would need to be 75 events in the entire sample. If the Kaplan–Meier estimates were needed to be adjusted for another factor that was binary, twice as many events would be needed to allow the sample to be stratified by that factor.

Figure 20.3 displays Kaplan–Meier three-year log cumulative hazard estimates stratified by sex and decile of age. The simulated sample consists of 2000 hypothetical subjects (389 of whom had events), with 1174 males (146 deaths) and 826 females (243 deaths). The sample was drawn from a population with a known survival distribution that is exponential with hazard function

$$\lambda(t|X_1, X_2) = .02 \exp[.8X_1 + .04(X_2 - 50)], \quad (20.17)$$

where X_1 represents the sex group (0 = male, 1 = female) and X_2 age in years, and censoring is uniform. Thus for this population PH, linearity, and additivity hold. Notice the amount of variability and wide confidence limits in the stratified nonparametric survival estimates.

```
n ← 2000
set.seed(3)
age ← 50 + 12 * rnorm(n)
label(age) ← 'Age'
sex ← factor(1 + (runif(n) ≤ .4), 1:2, c('Male', 'Female'))
cens ← 15 * runif(n)
h ← .02 * exp(.04 * (age - 50) + .8 * (sex == 'Female'))
ft ← -log(runif(n)) / h
e ← ifelse(ft ≤ cens, 1, 0)
print(table(e))
```

```
e
  0   1
1611 389
```

```
ft ← pmin(ft, cens)
units(ft) ← 'Year'
Srv ← Surv(ft, e)
age.dec ← cut2(age, g=10, levels.mean=TRUE)
label(age.dec) ← 'Age'
dd ← datadist(age, sex, age.dec); options(datadist='dd')
f.np ← cph(Srv ~ strat(age.dec) + strat(sex), surv=TRUE)
# surv=TRUE speeds up computations, and confidence limits when
# there are no covariables are still accurate.
p ← Predict(f.np, age.dec, sex, time=3, loglog=TRUE)
# Treat age.dec as a numeric variable (means within deciles)
p$age.dec ← as.numeric(as.character(p$age.dec))
ggplot(p, ylim=c(-5, -.5))
```

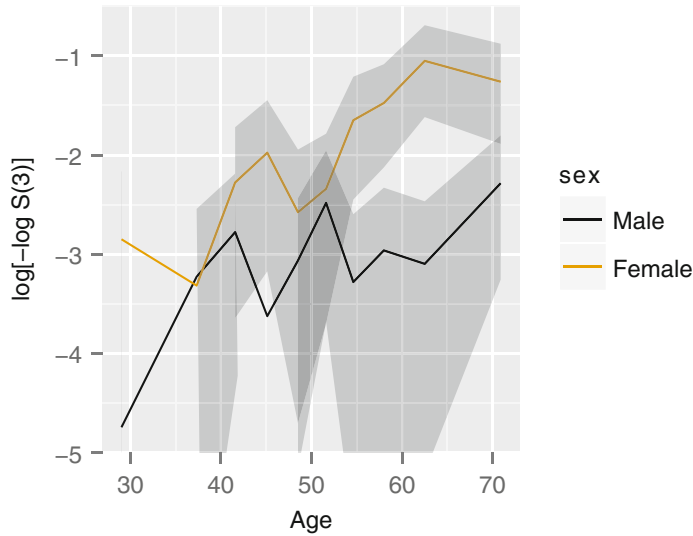


Fig. 20.3 Kaplan–Meier log λ estimates by sex and deciles of age, with 0.95 confidence limits. Solid line is for males, dashed line for females.

As with the logistic model and other regression models, the restricted cubic spline function is an excellent tool for modeling the regression relationship with very few assumptions. A four-knot spline Cox PH model in two variables (X_1, X_2) that assumes linearity in X_1 and no $X_1 \times X_2$ interaction is given by

$$\begin{aligned} \lambda(t|X) &= \lambda(t) \exp(\beta_1 X_1 + \beta_2 X_2 + \beta_3 X'_2 + \beta_4 X''_2), \\ &= \lambda(t) \exp(\beta_1 X_1 + f(X_2)), \end{aligned} \tag{20.18}$$

where X'_2 and X''_2 are spline component variables as described earlier and $f(X_2)$ is the spline function or spline transformation of X_2 given by

$$f(X_2) = \beta_2 X_2 + \beta_3 X'_2 + \beta_4 X''_2. \tag{20.19}$$

In linear form the Cox model without assuming linearity in X_2 is

$$\log \lambda(t|X) = \log \lambda(t) + \beta_1 X_1 + f(X_2). \tag{20.20}$$

By computing partial MLEs of $\beta_2, \beta_3,$ and $\beta_4,$ one obtains the estimated transformation of X_2 that yields linearity in log hazard or log cumulative hazard.

A similar model that does not assume PH in X_1 is the Cox model stratified on X_1 . Letting the stratification factor be $C = X_1,$ this model is

$$\begin{aligned}\log \lambda(t|X_2, C = j) &= \log \lambda_j(t) + \beta_1 X_2 + \beta_2 X_2' + \beta_3 X_2'' \\ &= \log \lambda_j(t) + f(X_2).\end{aligned}\quad (20.21)$$

This model does assume no $X_1 \times X_2$ interaction.

Figure 20.4 displays the estimated spline function relating age and sex to $\log[\Lambda(3)]$ in the simulated dataset, using the additive model stratified on sex.

```
f.noia <- cph(Srv ~ rcs(age,4) + strat(sex), x=TRUE, y=TRUE)
# Get accurate C.L. for any age by specifying x=TRUE y=TRUE
# Note: for evaluating shape of regression, we would not
# ordinarily bother to get 3-year survival probabilities -
# would just use X * beta
# We do so here to use same scale as nonparametric estimates
w <- latex(f.noia, inline=TRUE, digits=3)
latex(anova(f.noia), table.env=FALSE, file='')
```

	χ^2	d.f.	P
age	72.33	3	< 0.0001
Nonlinear	0.69	2	0.7067
TOTAL	72.33	3	< 0.0001

```
p <- Predict(f.noia, age, sex, time=3, loglog=TRUE)
ggplot(p, ylim=c(-5, -.5))
```

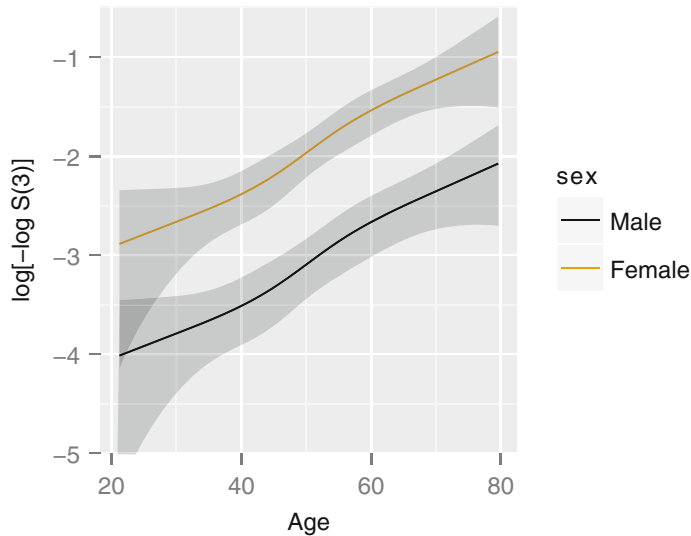


Fig. 20.4 Cox PH model stratified on sex, using spline function for age, no interaction. 0.95 confidence limits also shown. Solid line is for males, dashed line is for females.

A formal test of the linearity assumption of the Cox PH model in the above example is obtained by testing $H_0 : \beta_2 = \beta_3 = 0$. The χ^2 statistic with 2 d.f. is 0.69, $P = 0.7$. The fitted equation, after simplifying the restricted cubic spline to simpler (unrestricted) form, is $X\hat{\beta} = -1.46 + 0.0255\text{age} + 2.59 \times 10^{-5}(\text{age} - 30.3)^3 - 0.000101(\text{age} - 45.1)_+^3 + 9.73 \times 10^{-5}(\text{age} - 54.6)_+^3 - 2.22 \times 10^{-5}(\text{age} - 69.6)_+^3$. Notice that the spline estimates are closer to the true linear relationships than were the Kaplan–Meier estimates, and the confidence limits are much tighter. The spline estimates impose a smoothness on the relationship and also use more information from the data by treating age as a continuous ordered variable. Also, unlike the stratified Kaplan–Meier estimates, the modeled estimates can make the assumption of no age \times sex interaction. When this assumption is true, modeling effectively boosts the sample size in estimating a common function for age across both sex groups. Of course, this assumption can be tested and interactions can be modeled if necessary.

A Cox model that still does not assume PH for $X_1 = C$ but which allows for an $X_1 \times X_2$ interaction is

$$\begin{aligned} \log \lambda(t|X_2, C = j) = & \log \lambda_j(t) + \beta_1 X_2 + \beta_2 X_2' + \beta_3 X_2'' \\ & + \beta_4 X_1 X_2 + \beta_5 X_1 X_2' \\ & + \beta_6 X_1 X_2''. \end{aligned} \quad (20.22)$$

This model allows the relationship between X_2 and log hazard to be a smooth nonlinear function and the shape of the X_2 effect to be completely different for each level of X_1 if X_1 is dichotomous. Figure 20.5 displays a fit of this model at $t = 3$ years for the simulated dataset.

```
f.ia <- cph(Srv ~ rcs(age,4) * strat(sex), x=TRUE, y=TRUE,
           surv=TRUE)
w <- latex(f.ia, inline=TRUE, digits=3)
latex(anova(f.ia), table.env=FALSE, file='')
```

	χ^2	d.f.	P
age (Factor+Higher Order Factors)	72.82	6	< 0.0001
<i>All Interactions</i>	1.05	3	0.7886
<i>Nonlinear (Factor+Higher Order Factors)</i>	1.80	4	0.7728
age \times sex (Factor+Higher Order Factors)	1.05	3	0.7886
<i>Nonlinear</i>	1.05	2	0.5911
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	1.05	2	0.5911
TOTAL NONLINEAR	1.80	4	0.7728
TOTAL NONLINEAR + INTERACTION	1.80	5	0.8763
TOTAL	72.82	6	< 0.0001

```
p <- Predict(f.ia, age, sex, time=3, loglog=TRUE)
ggplot(p, ylim=c(-5, -.5))
```

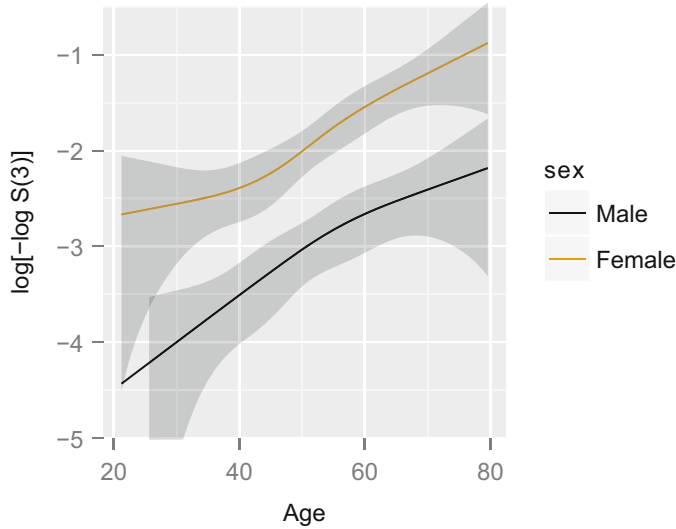


Fig. 20.5 Cox PH model stratified on sex, with interaction between age spline and sex. 0.95 confidence limits are also shown. Solid line is for males, dashed line for females.

The fitted equation is $X\hat{\beta} = -1.8 + 0.0493\text{age} - 2.15 \times 10^{-6}(\text{age} - 30.3)_+^3 - 2.82 \times 10^{-5}(\text{age} - 45.1)_+^3 + 5.18 \times 10^{-5}(\text{age} - 54.6)_+^3 - 2.15 \times 10^{-5}(\text{age} - 69.6)_+^3 + [\text{Female}][-0.0366\text{age} + 4.29 \times 10^{-5}(\text{age} - 30.3)_+^3 - 0.00011(\text{age} - 45.1)_+^3 + 6.74 \times 10^{-5}(\text{age} - 54.6)_+^3 - 2.32 \times 10^{-7}(\text{age} - 69.6)_+^3]$. The test for interaction yielded $\chi^2 = 1.05$ with 3 d.f., $P = 0.8$. The simultaneous test for linearity and additivity yielded $\chi^2 = 1.8$ with 5 d.f., $P = 0.9$. Note that allowing the model to be very flexible (not assuming linearity in age, additivity between age and sex, and PH for sex) still resulted in estimated regression functions that are very close to the true functions. However, confidence limits in this unrestricted model are much wider.

Figure 20.6 displays the estimated relationship between left ventricular ejection fraction (LVEF) and log hazard ratio for cardiovascular death in a sample of patients with significant coronary artery disease. The relationship is estimated using three knots placed at quantiles 0.05, 0.5, and 0.95 of LVEF. Here there is significant nonlinearity (Wald $\chi^2 = 9.6$ with 1 d.f.). The graphs leads to a transformation of LVEF that better satisfies the linearity assumption: $\min(\text{LVEF}, 0.5)$. This transformation has the best log likelihood “for the money” as judged by the Akaike information criterion ($\text{AIC} = -2 \log \text{L.R.} - 2 \times \text{no. parameters} = 127$). The AICs for 3, 4, 5, and 6-knot spline fits were, respectively, 126, 124, 122, and 120.

Had the suggested transformation been more complicated than a truncation, a tentative transformation could have been checked for adequacy by expanding the new transformed variable into a new spline function and testing it for linearity.

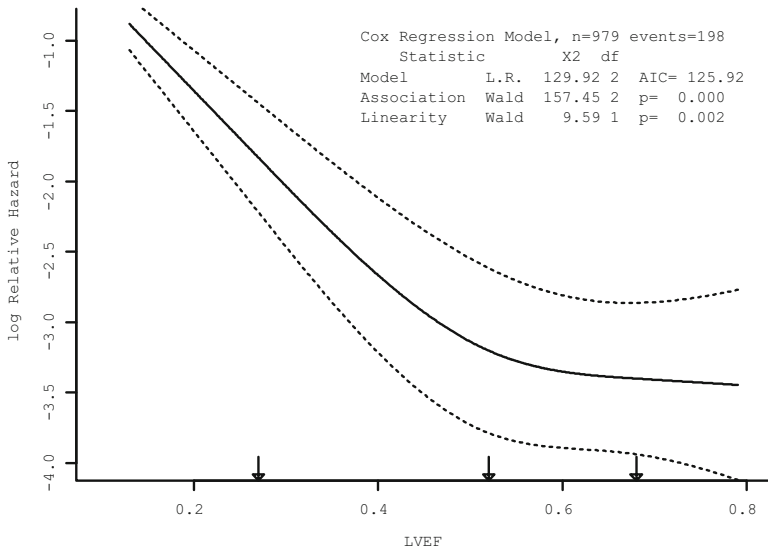


Fig. 20.6 Restricted cubic spline estimate of relationship between LVEF and relative log hazard from a sample of 979 patients and 198 cardiovascular deaths. Data from the Duke Cardiovascular Disease Databank.

Other methods based on smoothed residual plots are also valuable tools for selecting predictor transformations. Therneau et al.⁶⁰⁵ describe residuals based on martingale theory that can estimate transformations of any number of predictors omitted from a Cox model fit, after adjusting for other variables included in the fit. Figure 20.7 used various smoothing methods on the points (LVEF, residual). First, the R `loess` function⁹⁶ was used to obtain a smoothed scatterplot fit and approximate 0.95 confidence bars. Second, an ordinary least squares model, representing LVEF as a restricted cubic spline with five default knots, was fitted. Ideally, both fits should have used weighted regression as the residuals do not have equal variance. Predicted values from this fit along with 0.95 confidence limits are shown. The `loess` and spline-linear regression agree extremely well. Third, Cleveland's `lowess` scatterplot smoother¹¹¹ was used on the martingale residuals against LVEF. The suggested transformation from all three is very similar to that of Figure 20.6. For smaller sample sizes, the raw residuals should also be displayed. There is one vector of martingale residuals that is plotted against all of the predictors. When correlations among predictors are mild, plots of estimated predictor transformations without adjustment for other predictors (i.e., marginal transformations) may be useful. Martingale residuals may be obtained quickly by fixing $\hat{\beta} = 0$ for all predictors. Then smoothed plots of predictor against residual may be made for all predictors. Table 20.3 summarizes some of the

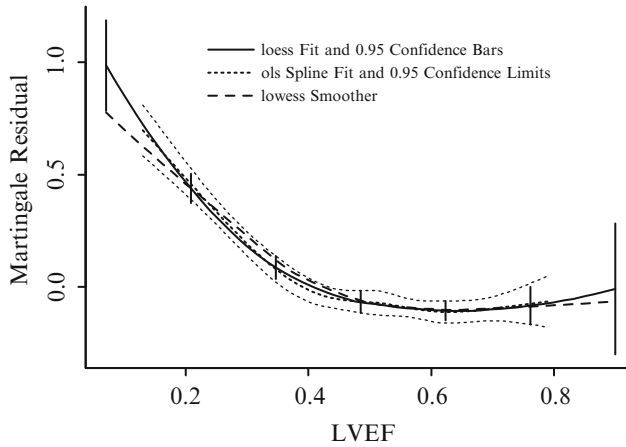


Fig. 20.7 Three smoothed estimates relating martingale residuals⁶⁰⁵ to LVEF.

Table 20.3 Uses of martingale residuals for estimating predictor transformations

Purpose	Method
Estimate transformation for a single variable	Force $\hat{\beta}_1 = 0$ and compute residuals from the null regression
Check linearity assumption for a single variable	Compute $\hat{\beta}_1$ and compute residuals from the linear regression
Estimate marginal transformations for p variables	Force $\hat{\beta}_1, \dots, \hat{\beta}_p = 0$ and compute residuals from the global null model
Estimate transformation for variable i adjusted for other $p - 1$ variables	Estimate $p - 1$ β s, forcing $\hat{\beta}_i = 0$ Compute residuals from mixed global/null model

6 ways martingale residuals may be used. See section 10.5 for more information on checking the regression assumptions. The methods for examining interaction surfaces described there apply without modification to the Cox model (except that the nonparametric regression surface does not apply because of censoring).

20.6.2 Proportional Hazards Assumption

Even though assessment of fit of the regression part of the Cox PH model corresponds with other regression models such as the logistic model, the Cox model has its own distributional assumption in need of validation. Here, of course, the distributional assumption is not as stringent as with other survival

models, but we do need to validate how the survival or hazard functions for various subjects are connected. There are many graphical and analytical methods of verifying the PH assumption. Two of the methods have already been discussed: a graphical examination of parallelism of $\log \Lambda$ plots, and a comparison of stratified with unstratified models (as in Figure 20.1). Muenz⁴⁶⁷ suggested a simple modification that will make nonproportional hazards more apparent: plot $\Lambda_{KM_1}(t)/\Lambda_{KM_2}(t)$ against t and check for flatness. The points on this curve can be passed through a smoother. One can also plot differences in $\log(-\log S(t))$ against t .¹⁴³ Arjas²⁹ developed a graphical method based on plotting the estimated cumulative hazard versus the cumulative number of events in a stratum as t progresses.

There are other methods for assessing whether PH holds that may be more direct. Gore et al.,²²⁶ Harrell and Lee,²⁶⁶ and Kay³⁴⁰ (see also Anderson and Senthilselvan²⁷) describe a method for allowing the log hazard ratio (Cox regression coefficient) for a predictor to be a function of time by fitting specially stratified Cox models. Their method assumes that the predictor being examined for PH already satisfies the linear regression assumption. Follow-up time is stratified into intervals and a separate model is fitted to compute the regression coefficient within each interval, assuming that the effect of the predictor is constant only within that small interval. It is recommended that intervals be constructed so that there is roughly an equal number of events in each. The number of intervals should allow at least 10 or 20 events per interval.

The interval-specific log hazard ratio is estimated by excluding all subjects with event/censoring time before the start of the interval and censoring all events that occur after the end of the interval. This process is repeated for all desired time intervals. By plotting the log hazard ratio and its confidence limits versus the interval, one can assess the importance of a predictor as a function of follow-up time and learn how to model non-PH using more complicated models containing predictor by time interactions. If the hazard ratio is approximately constant within broad time intervals, the time stratification method can be used for fitting and testing the predictor \times time interaction [266, p. 827]; [98].

Consider as an example the rat vaginal cancer data used in Figures 18.9, 18.10, and 20.1. Recall that the PH assumption appeared to be satisfied for the two groups although Figure 18.9 demonstrated some non-Weibullness. Figure 20.8 contains a Λ ratio plot.⁴⁶⁷

```
f <- cph(S ~ strat(group), surv=TRUE)
# For both strata, eval. S(t) at combined set of death times
times <- sort(unique(days[death == 1]))
est <- survest(f, data.frame(group=levels(group)),
              times=times, conf.type="none")$surv
cumhaz <- -log(est)
plot(times, cumhaz[2,] / cumhaz[1,], xlab="Days",
      ylab="Cumulative Hazard Ratio", type="s")
abline(h=1, col=gray(.80))
```

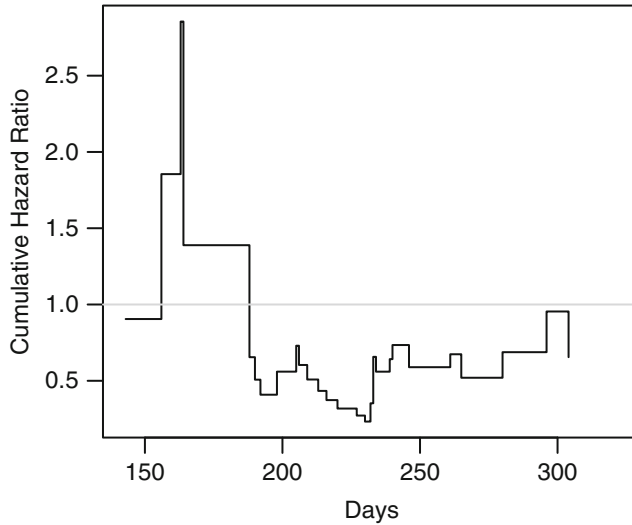


Fig. 20.8 Estimate of Λ_2/Λ_1 based on $-\log$ of Altschuler–Nelson–Fleming–Harrington nonparametric survival estimates.

Table 20.4 Interval-specific group effects from rat data by artificial censoring

Time Interval	Observations	Deaths	Log Hazard Ratio	Standard Error
[0, 209)	40	12	-0.47	0.59
[209, 234)	27	12	-0.72	0.58
234 +	14	12	-0.50	0.64

```
hazard.ratio.plot(g$x, g$y, e=12, pr=TRUE)
```

The number of observations is declining over time because computations in each interval were based on animals followed at least to the start of that interval. The overall Cox regression coefficient was -0.57 with a standard error of 0.35 . There does not appear to be any trend in the hazard ratio over time, indicating a constant hazard ratio or proportional hazards (Table 20.4).

Now consider the Veterans Administration Lung Cancer dataset [331, pp. 60, 223–4]. Log Λ plots indicated that the four cell types did not satisfy PH. To simplify the problem, omit patients with “large” cell type and let the binary predictor be 1 if the cell type is “squamous” and 0 if it is “small” or “adeno.” We are assessing whether survival patterns for the two groups “squamous” versus “small” or “adeno” have PH. Interval-specific estimates of the squamous : small,adeno log hazard ratios (using Efron’s likelihood) are found in Table 20.5. Times are in days.

Table 20.5 Interval-specific effects of squamous cell cancer in VA lung cancer data

Time Interval	Observations	Deaths	Log Hazard Ratio	Standard Error
[0, 21)	110	26	-0.46	0.47
[21, 52)	84	26	-0.90	0.50
[52, 118)	59	26	-1.35	0.50
118 +	28	26	-1.04	0.45

Table 20.6 Interval-specific effects of performance status in VA lung cancer data

Time Interval	Observations	Deaths	Log Hazard Ratio	Standard Error
[0, 19]	137	27	-0.053	0.010
[19, 49)	112	26	-0.047	0.009
[49, 99)	85	27	-0.036	0.012
99 +	28	26	-0.012	0.014

```

getHdata(valung)
with(valung, {
  hazard.ratio.plot(1 * (cell == 'Squamous'), Surv(t, dead),
                    e=25, subset=cell != 'Large',
                    pr=TRUE, pl=FALSE)
  hazard.ratio.plot(1 * kps, Surv(t, dead), e=25,
                    pr=TRUE, pl=FALSE) })

```

There is evidence of a trend of a decreasing hazard ratio over time which is consistent with the observation that squamous cell patients had equal or worse survival in the early period but decidedly better survival in the late phase.

From the same dataset now examine the PH assumption for Karnofsky performance status using data from all subjects, if the linearity assumption is satisfied. Interval-specific regression coefficients for this predictor are given in Table 20.6. There is good evidence that the importance of performance status is decreasing over time and that it is not a prognostic factor after roughly 99 days. In other words, once a patient survives 99 days, the performance status does not contain much information concerning whether the patient will survive 120 days. This non-PH would be more difficult to detect from Kaplan-Meier plots stratified on performance status unless performance status was stratified carefully.

Figure 20.9 displays a log hazard ratio plot for a larger dataset in which more time strata can be formed. In 3299 patients with coronary artery disease, 827 suffered cardiovascular death or nonfatal myocardial infarction. Time

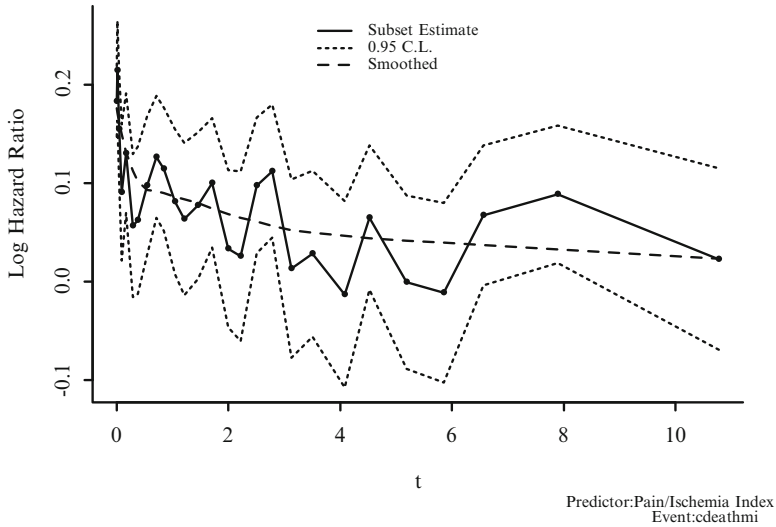


Fig. 20.9 Stratified hazard ratios for pain/ischemia index over time. Data from the Duke Cardiovascular Disease Databank.

was stratified into intervals containing approximately 30 events, and within each interval the Cox regression coefficient for an index of anginal pain and ischemia was estimated. The pain/ischemia index, one component of which is unstable angina, is seen to have a strong effect for only six months. After that, survivors have stabilized and knowledge of the angina status in the previous six months is not informative.

Another method for graphically assessing the log hazard ratio over time is based on Schoenfeld's *partial residuals*^{503, 557} with respect to each predictor in the fitted model. The residual is the contribution of the first derivative of the log likelihood function with respect to the predictor's regression coefficient, computed separately at each risk set or unique failure time. In Figure 20.10 the "loess-smoothed"⁹⁶ (with approximate 0.95 confidence bars) and "super-smoothed"²⁰⁷ relationship between the residual and unique failure time is shown for the same data as Figure 20.9. For smaller n , the raw residuals should also be displayed to convey the proper sense of variability. The agreement with the pattern in Figure 20.9 is evident.

Pettitt and Bin Daud⁵⁰³ suggest scaling the partial residuals by the information matrix components. They also propose a score test for PH based on the Schoenfeld residuals. Grambsch and Therneau²³³ found that the Pettitt–Bin Daud standardization is sometimes misleading in that non-PH in one variable may cause the residual plot for another variable to display non-PH. The Grambsch–Therneau weighted residual solves this problem and also yields a residual that is on the same scale as the log relative hazard ratio. Their residual is

$$\hat{\beta} + dR\hat{V}, \quad (20.23)$$

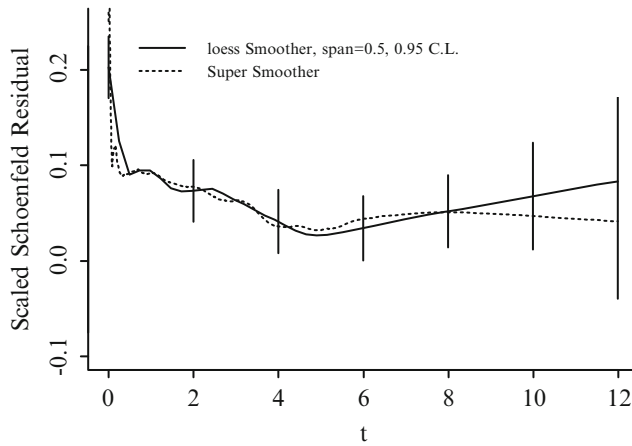


Fig. 20.10 Smoothed weighted²³³ Schoenfeld⁵⁵⁷ residuals for the same data in Figure 20.9. Test for PH based on the correlation (ρ) between the individual weighted Schoenfeld residuals and the rank of failure time yielded $\rho = -0.23$, $z = -6.73$, $P = 2 \times 10^{-11}$.

where d is the total number of events, R is the $n \times p$ matrix of Schoenfeld residuals, and \hat{V} is the estimated covariance matrix for $\hat{\beta}$. This new residual can also be the basis for tests for PH, by correlating a user-specified function of unique failure times with the weighted residuals.

8

The residual plot is computationally very attractive since the score residual components are byproducts of Cox maximum likelihood estimation. Another attractive feature is the lack of need to categorize the time axis. Unless approximate confidence intervals are derived from smoothing techniques, a lack of confidence intervals from most software is one disadvantage of the method.

9

Formal tests for PH can be based on time-stratified Cox regression estimates.^{27, 266} Alternatively, more complex (and probably more efficient) formal tests for PH can be derived by specifying a form for the time by predictor interaction (using what is called a time-dependent covariable in the Cox model) and testing coefficients of such interactions for significance. The obsolete Version 5 SAS PHGLM procedure used a computationally fast procedure based on an approximate score statistic that tests for linear correlation between the rank order of the failure times in the sample and Schoenfeld’s partial residuals.^{258, 266} This test is available in R (for both weighted and unweighted residuals) using Therneau’s `cox.zph` function in the `survival` package. For the results in Figure 20.10, the test for PH is highly significant (correlation coefficient = -0.23 , normal deviate $z = -6.73$). Since there is only one regression parameter, the weighted residuals are a constant multiple of the unweighted ones, and have the same correlation coefficient.

10

11

Table 20.7 Time-specific hazard ratio estimates of squamous cell cancer effect in VA lung cancer data, by fitting two Weibull distributions with unequal shape parameters

t	log Hazard Ratio
10	-0.36
36	-0.64
83.5	-0.83
200	-1.02

Another method for checking the PH assumption which is especially applicable to a polytomous predictor involves taking ratios of parametrically estimated hazard functions estimated separately for each level of the predictor. For example, suppose that a risk factor X is either present ($X = 1$) or absent ($X = 0$), and suppose that separate Weibull distributions adequately fit the survival pattern of each group. If there are no other predictors to adjust for, define the hazard function for $X = 0$ as $\alpha\gamma t^{\gamma-1}$ and the hazard for $X = 1$ as $\delta\theta t^{\theta-1}$. The $X = 1 : X = 0$ hazard ratio is

$$\frac{\alpha\gamma t^{\gamma-1}}{\delta\theta t^{\theta-1}} = \frac{\alpha\gamma}{\delta\theta} t^{\gamma-\theta}. \quad (20.24)$$

The hazard ratio is constant if the two Weibull shape parameters (γ and θ) are equal. These Weibull parameters can be estimated separately and a Wald test statistic of $H_0 : \gamma = \theta$ can be computed by dividing the square of their difference by the sum of the squares of their estimated standard errors, or better by a likelihood ratio test. A plot of the estimate of the hazard ratio above as a function of t may also be informative.

12

In the VA lung cancer data, the MLEs of the Weibull shape parameters for squamous cell cancer is 0.77 and for the combined small + adeno is 0.99. Estimates of the reciprocals of these parameters, provided by some software packages, are 1.293 and 1.012 with respective standard errors of 0.183 and 0.0912. A Wald test for differences in these reciprocals provides a rough test for a difference in the shape estimates. The Wald χ^2 is 1.89 with 1 d.f. indicating slight evidence for non-PH.

The fitted Weibull hazard function for squamous cell cancer is $.0167t^{0.23}$ and for adeno + small is $0.0144t^{-0.01}$. The estimated hazard ratio is then $1.16t^{-0.22}$ and the log hazard ratio is $0.148 - 0.22 \log t$. By evaluating this Weibull log hazard ratio at interval midpoints (arbitrarily using $t = 200$ for the last (open) interval) we obtain log hazard ratios that are in good agreement with those obtained by time-stratifying the Cox model (Table 20.5) as shown in Table 20.7.

There are many methods of assessing PH using time-dependent covariables in the Cox model.^{226, 583} Gray^{237, 238} mentions a flexible and efficient method of estimating the hazard ratio function using time-dependent covariables that are $X \times$ spline term interactions. Gray's method uses B-splines and

requires one to maximize a *penalized* log-likelihood function. Verweij and van Houwelingen⁶⁴¹ developed a more nonparametric version of this approach. Hess²⁸⁹ uses simple restricted cubic splines to model the time-dependent covariable effects (see also [4, 287, 398, 498]). Suppose that $k = 4$ knots are used and that a covariable X is already transformed correctly. The model is

$$\log \lambda(t|X) = \log \lambda(t) + \beta_1 X + \beta_2 X t + \beta_3 X t' + \beta_4 X t'', \quad (20.25)$$

where t', t'' are constructed spline variables (Equation 2.25). The $X + 1 : X$ log hazard ratio function is estimated by

$$\hat{\beta}_1 + \hat{\beta}_2 t + \hat{\beta}_3 t' + \hat{\beta}_4 t''. \quad (20.26)$$

This method can be generalized to allow for simultaneous estimation of the shape of the X effect and $X \times t$ interaction using spline surfaces in (X, t) instead of (X_1, X_2) (Section 2.7.2).

Table 20.8 summarizes many facets of verifying assumptions for PH models. The trade-offs of the various methods for assessing proportional hazards are given in Table 20.9.

13

14

20.7 What to Do When PH Fails

When a factor violates the PH assumption and a test of association is not needed, the factor can be adjusted for through stratification as mentioned earlier. This is especially attractive if the factor is categorical. For continuous predictors, one may want to stratify into quantile groups. The continuous version of the predictor can still be adjusted for as a covariable to account for any residual linearity within strata.

When a test of significance is needed and the P -value is impressive, the “principle of conservatism” could be invoked, as the P -value would likely have been more impressive had the factor been modeled correctly. Predicted survival probabilities using this approach will be erroneous in certain time intervals.

An efficient test of association can be done using time-dependent covariables [444, pp. 208–217]. For example, in the model

$$\lambda(t|X) = \lambda_0(t) \exp(\beta_1 X + \beta_2 X \times \log(t + 1)) \quad (20.27)$$

one tests $H_0 : \beta_1 = \beta_2 = 0$ with 2 d.f. This is similar to the approach used by [72]. Stratification on time intervals can also be used:^{27, 226, 266}

$$\lambda(t|X) = \lambda_0(t) \exp(\beta_1 X + \beta_2 X \times [t > c]). \quad (20.28)$$

Table 20.8 Assumptions of the Proportional Hazards Model

Variables	Assumptions	Verification
Response Variable T Time Until Event	Shape of $\lambda(t X)$ for fixed X as $t \uparrow$ Cox: none Weibull: t^θ	Shape of $S_{KM}(t)$
Interaction Between X and T	Proportional hazards—effect of X does not depend on T (e.g., treatment effect is constant over time)	<ul style="list-style-type: none"> • Categorical X: check parallelism of stratified $\log[-\log S(t)]$ plots as $t \uparrow$ • Muenz⁴⁶⁷ cum. hazard ratio plots • Arjas²⁹ cum. hazard plots • Check agreement of stratified and modeled estimates • Hazard ratio plots • Smoothed Schoenfeld residual plots and correlation test (time vs. residual) • Test time-dependent covariable such as $X \times \log(t + 1)$ • Ratio of parametrically estimated $\lambda(t)$
Individual Predictors X	Shape of $\lambda(t X)$ for fixed t as $X \uparrow$ Linear: $\log \lambda(t X) = \log \lambda(t) + \beta X$ Nonlinear: $\log \lambda(t X) = \log \lambda(t) + f(X)$	<ul style="list-style-type: none"> • k-level ordinal X : linear term + $k - 2$ dummy variables • Continuous X: polynomials, spline functions, smoothed martingale residual plots
Interaction Between X_1 and X_2	Additive effects: effect of X_1 on $\log \lambda$ is independent of X_2 and vice versa	Test nonadditive terms (e.g., products)

If this step-function model holds, and if a sufficient number of subjects have late follow-up, you can also fit a model for early outcomes and a separate one for late outcomes using interval-specific censoring as discussed in Section 20.6.2. The dual model approach provides easy to interpret models, assuming that proportional hazards is satisfied within each interval.

Kronborg and Aaby³⁶⁷ and Dabrowska et al.¹⁴³ provide tests for differences in $\Lambda(t)$ at specific t based on stratified PH models. These can also be used to test for treatment effects when PH is violated for treatment but not for

adjustment variables. Differences in mean restricted life length (differences in areas under survival curves up to a fixed finite time) can also be useful for comparing therapies when PH fails.³³⁵

Table 20.9 Comparison of methods for checking the proportional hazards assumption and for allowing for non-proportional hazards

Method	Requires Grouping X	Requires Grouping t	Computational Efficiency	Yields Formal Test	Yields Estimate of $\lambda_2(t)/\lambda_1(t)$	Requires Fitting 2 Models	Must Choose Smoothing Parameter
log[-log], Muenz, Arjas plots	x		x			x	
Dabrowska log \hat{A} difference plots	x		x	x		x	
Stratified vs. Modeled Estimates	x		x			x	
Hazard ratio plot		x		?	x	x	?
Schoenfeld residual plot			x		x		x
Schoenfeld residual correlation test			x	x			
Fit time-dependent covariables				x	x		
Ratio of parametric estimates of $\lambda(t)$	x		x	x	x	x	

Parametric models that assume an effect other than PH, for example, the log-logistic model,²²⁶ can be used to allow a predictor to have a constantly increasing or decreasing effect over time. If one predictor satisfies PH but another does not, this approach will not work.

15

20.8 Collinearity

See Section 4.6 for the general approach using variance inflation factors.

20.9 Overly Influential Observations

Therneau et al.⁶⁰⁵ describe the use of *score residuals* for assessing influence in Cox and related regression models. They show that the *infinitesimal jackknife* estimate of the influence of observation i on β equals Vs' , where V is the estimated variance–covariance matrix of the p regression estimates b and $s = (s_{i1}, s_{i2}, \dots, s_{ip})$ is the vector of score residuals for the p regression coefficients for the i th observation. Let $S_{n \times p}$ denote the matrix of score residuals over all observations. Then an approximation to the unstandardized change in b (DFBETA) is SV . Standardizing by the standard errors of b found from the diagonals of V , $e = (V_{11}, V_{22}, \dots, V_{pp})^{1/2}$, yields

$$\text{DFBETAS} = SV \text{Diag}(e)^{-1}, \quad (20.29)$$

where $\text{Diag}(e)$ is a diagonal matrix containing the estimated standard errors.

As discussed in Section 20.13, identification of overly influential observations is facilitated by printing, for each predictor, the list of observations containing $\text{DFBETAS} > u$ for any parameter associated with that predictor. The choice of cutoff u depends on the sample size among other things. A typical choice might be $u = 0.2$ indicating a change in a regression coefficient of 0.2 standard errors.

20.10 Quantifying Predictive Ability

To obtain a unitless measure of predictive ability for a Cox PH model we can use the R index described in Section 9.8.3, which is the square root of the fraction of log likelihood explained by the model of the log likelihood that could be explained by a perfect model, penalized for the complexity of the model. The lowest (best) possible $-2 \log$ likelihood for the Cox model is zero, which occurs when the predictors can perfectly rank order the survival times. Therefore, as was the case with the logistic model, the quantity L^* from Section 9.8.3 is zero and an R index that is penalized for the number of parameters in the model is given by

$$R^2 = (\text{LR} - 2p)/L^0, \quad (20.30)$$

where p is the number of parameters estimated and L^0 is the $-2 \log$ likelihood when β is restricted to be zero (i.e., there are no predictors in the model). R will be near one for a perfectly predictive model and near zero for a model that does not discriminate between short and long survival times. The R index does not take into account any stratification factors. If stratification factors are present, R will be near one if survival times can be perfectly ranked within strata even though there is overlap between strata.

Schemper⁵⁴⁶ and Korn and Simon³⁶⁵ have reported that R^2 is too sensitive to the distribution of censoring times and have suggested alternatives based on the distance between estimated Cox survival probabilities (using predictors) and Kaplan–Meier estimates (ignoring predictors). Kent and O’Quigley³⁴⁵ also report problems with R^2 and suggest a more complex measure. Schemper⁵⁴⁸ investigated the Maddala–Magee^{431, 432} index R_{LR}^2 described in Section 9.8.3, applied to Cox regression:

$$\begin{aligned} R_{LR}^2 &= 1 - \exp(-LR/n) \\ &= 1 - \omega^{2/n}, \end{aligned} \tag{20.31}$$

where ω is the null model likelihood divided by the fitted model likelihood.

For many situations, R_{LR}^2 performed as well as Schemper’s more complex measure^{546, 549} and hence it is preferred because of its ease of calculation (assuming that PH holds). Ironically, Schemper⁵⁴⁸ demonstrated that the n in the formula for this index is the total number of observations, not the number of events (but see O’Quigley, Xu, and Stare⁴⁸¹). To make the R^2 index have a maximum value of 1.0, we use the Nagelkerke⁴⁷¹ R_N^2 discussed in Section 9.8.3.

An easily interpretable index of discrimination for survival models is derived from Kendall’s τ and Somers’ D_{xy} rank correlation,⁵⁷⁹ the Gehan–Wilcoxon statistic for comparing two samples for survival differences, and the Brown–Hollander–Korwar nonparametric test of association for censored data.^{76, 170, 262, 268} This index, c , is a generalization of the area under the ROC curve discussed under the logistic model, in that it applies to a continuous response variable that can be censored. The c index is the proportion of all pairs of subjects whose survival time can be ordered such that the subject with the higher predicted survival is the one who survived longer. Two subjects’ survival times cannot be ordered if both subjects are censored or if one has failed and the follow-up time of the other is less than the failure time of the first. The c index is a probability of concordance between predicted and observed survival, with $c = 0.5$ for random predictions and $c = 1$ for a perfectly discriminating model. The c index is mildly affected by the amount of censoring. D_{xy} is obtained from $2(c - 0.5)$. While c (and D_{xy}) is a good measure of pure discrimination ability of a single model, it is not sensitive enough to allow multiple models to be compared⁴⁴⁷.

Since high hazard means short survival time, when the linear predictor $X\hat{\beta}$ from a Cox model is compared with observed survival time, D_{xy} will be negative. Some analysts may want to negate reported values of D_{xy} .

16

17

20.11 Validating the Fitted Model

Separate bootstrap or cross-validation assessments can be made for calibration and discrimination of Cox model survival and log relative hazard estimates.

18

20.11.1 Validation of Model Calibration

One approach to validation of the calibration of predictions is to obtain unbiased estimates of the difference between Cox predicted and Kaplan–Meier survival estimates at a fixed time u . Here is one sequence of steps.

1. Obtain cutpoints (e.g., deciles) of predicted survival at time u so as to have a given number of subjects (e.g., 50) in each interval of predicted survival. These cutpoints are based on the distribution of $\hat{S}(u|X)$ in the whole sample for the “final” model (for data-splitting, instead use the model developed in the training sample). Let k denote the number of intervals used.
2. Compute the average $\hat{S}(u|X)$ in each interval.
3. Compare this with the Kaplan–Meier survival estimates at time u , stratified by intervals of $\hat{S}(u|X)$. Let the differences be denoted by $d = (d_1, \dots, d_k)$.
4. Use bootstrapping or cross-validation to estimate the overoptimism in d and then to correct d to get a more fair assessment of these differences. For each repetition, repeat any stepwise variable selection or stagewise significance testing using the same stopping rules as were used to derive the “final” model. No more than $B = 200$ replications are needed to obtain accurate estimates.
5. If desired, the bias-corrected d can be added to the original stratified Kaplan–Meier estimates to obtain a bias-corrected calibration curve.

However, any statistical method that uses binning of continuous variables (here, the predicted risk), is arbitrary and has lower precision than smooth estimates that allow for interpolation. A far better approach to estimating calibration curves for survival models is to use the flexible adaptive hazard regression approach of Kooperberg et al.³⁶¹ as discussed on P. 450. Their method does not assume linearity or proportional hazards. Hazard regression can be used to estimate the relationship between (suitably transformed) predicted survival probabilities and observed outcomes, i.e., to derive a calibration curve. The bootstrap is used to de-bias the estimates to correct for overfitting, allowing estimation of the likely future calibration performance of the fitted model.

As an example, consider a dataset of 20 random uniformly distributed predictors for a sample of size 200. Let the failure time be another random

uniform variable that is independent of *all* the predictors, and censor half of the failure times at random. Due to fitting 20 predictors to 100 events, there will apparently be fair agreement between predicted and observed survival over all strata (smooth black curve from hazard regression in Figure 20.11). However, the bias-corrected calibration (blue curve from hazard regression) gives a more truthful answer: examining the X s across levels of predicted survival demonstrate that predicted and observed survival are weakly related, in more agreement with how the data were generated. For the more arbitrary Kaplan-Meier approach, we divide the observations into quintiles of predicted 0.5-year survival, so that there are 40 observations per stratum.

```
n ← 200
p ← 20
set.seed(6)
xx ← matrix(rnorm(n * p), nrow=n, ncol=p)
y ← runif(n)
units(y) ← "Year"
e ← c(rep(0, n / 2), rep(1, n / 2))
f ← cph(Surv(y, e) ~ xx, x=TRUE, y=TRUE,
        time.inc=.5, surv=TRUE)
cal ← calibrate(f, u=.5, B=200)
```

Using Cox survival estimates at 0.5 Years

```
plot(cal, ylim=c(.4, 1), subtitles=FALSE)
calkm ← calibrate(f, u=.5, m=40, cmethod='KM', B=200)
```

Using Cox survival estimates at 0.5 Years

```
plot(calkm, add=TRUE) # Figure 20.11
```

20.11.2 Validation of Discrimination and Other Statistical Indexes

Here bootstrapping and cross-validation are used as for logistic models (Section 10.9). We can obtain bootstrap bias-corrected estimates of c or equivalently D_{xy} . To instead obtain a measure of relative calibration or slope shrinkage, we can bootstrap the apparent estimate of $\gamma = 1$ in the model

$$\lambda(t|X) = \lambda(t) \exp(\gamma Xb). \quad (20.32)$$

Besides being a measure of calibration in itself, the bootstrap estimate of γ also leads to an unreliability index U which measures how far the model maximum log likelihood (which allows for an overall slope correction) is from the log likelihood evaluated at “frozen” regression coefficients ($\gamma = 1$) (see [267] and Section 10.9).

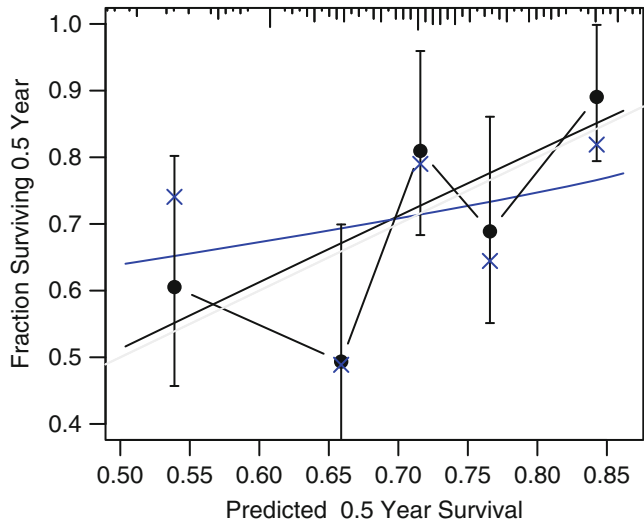


Fig. 20.11 Calibration of random predictions using Efron’s bootstrap with $B = 200$ resamples. Dataset has $n = 200$, 100 uncensored observations, 20 random predictors, model $\chi^2_{20} = 19$. The smooth black line is the apparent calibration estimated by adaptive linear spline hazard regression³⁶¹, and the blue line is the bootstrap bias– (overfitting–) corrected calibration curve estimated also by hazard regression. The gray scale line is the line of identity representing perfect calibration. Black dots represent apparent calibration accuracy obtained by stratifying into intervals of predicted 0.5y survival containing 40 events per interval and plotting the mean predicted value within the interval against the stratum’s Kaplan-Meier estimate. The blue \times represent bootstrap bias-corrected Kaplan-Meier estimates.

$$U = \frac{\text{LR}(\hat{\gamma}Xb) - \text{LR}(Xb)}{L^0}, \tag{20.33}$$

where L^0 is the $-2 \log$ likelihood for the null model (Section 9.8.3). Similarly, a discrimination index D ²⁶⁷ can be derived from the $-2 \log$ likelihood at the shrunken linear predictor, penalized for estimating one parameter (γ) (see also [633, p. 1318] and [123]):

$$D = \frac{\text{LR}(\hat{\gamma}Xb) - 1}{L^0}. \tag{20.34}$$

D is the same as R^2 discussed above when $p = 1$ (indicating only one reestimated parameter, γ), the penalized proportion of explainable log likelihood that was explained by the model. Because of the remark of Schemper,⁵⁴⁶ all of these indexes may unfortunately be functions of the censoring pattern.

An index of overall quality that penalizes discrimination for unreliability is

$$Q = D - U = \frac{\text{LR}(Xb) - 1}{L^0}. \tag{20.35}$$

Q is a normalized and penalized -2 log likelihood that is evaluated at the uncorrected linear predictor.

For the random predictions used in Figure 20.11, the bootstrap estimates with $B = 200$ resamples are found in Table 20.10.

```
latex(validate(f, B=200), digits=3, file='', caption='',
        table.env=TRUE, label='tab:cox-val-random')
```

Table 20.10 Bootstrap validation of a Cox model with random predictors

Index	Original Sample	Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.213	0.335	0.147	0.188	0.025	200
R^2	0.092	0.191	0.042	0.150	-0.058	200
Slope	1.000	1.000	0.389	0.611	0.389	200
D	0.021	0.048	0.009	0.039	-0.019	200
U	-0.002	-0.002	0.028	-0.031	0.028	200
Q	0.023	0.050	-0.020	0.070	-0.047	200
g	0.516	0.878	0.339	0.539	-0.023	200

It can be seen that the apparent correlation ($D_{xy} = -0.21$) does not hold up after correcting for overfitting ($D_{xy} = -0.02$). Also, the slope shrinkage (0.39) indicates extreme overfitting.

See [633, Section 6] and [640] and Section 18.3.7 for still more useful methods for validating the Cox model.

20.12 Describing the Fitted Model

As with logistic modeling, once a Cox PH model has been fitted and all its assumptions verified, the final model needs to be presented and interpreted. The fastest way to describe the model is to interpret each effect in it. For each predictor the change in log hazard per desired units of change in the predictor value may be computed, or the antilog of this quantity, $\exp(\beta_j \times \text{change in } X_j)$, may be used to estimate the hazard ratio holding all other factors constant. When X_j is a nonlinear factor, changes in predicted $X\beta$ for sensible values of X_j such as quartiles can be used as described in Section 10.10. Of course for nonmodeled stratification factors, this method is of no help. Figure 20.12 depicts a way to display estimated surgical : medical hazard ratios in the presence of a significant treatment by disease severity interaction and a secular trend in the benefit of surgical therapy (treatment by year of diagnosis interaction).

Often, the use of predicted survival probabilities may make the model more interpretable. If the effect of only one factor is being displayed and

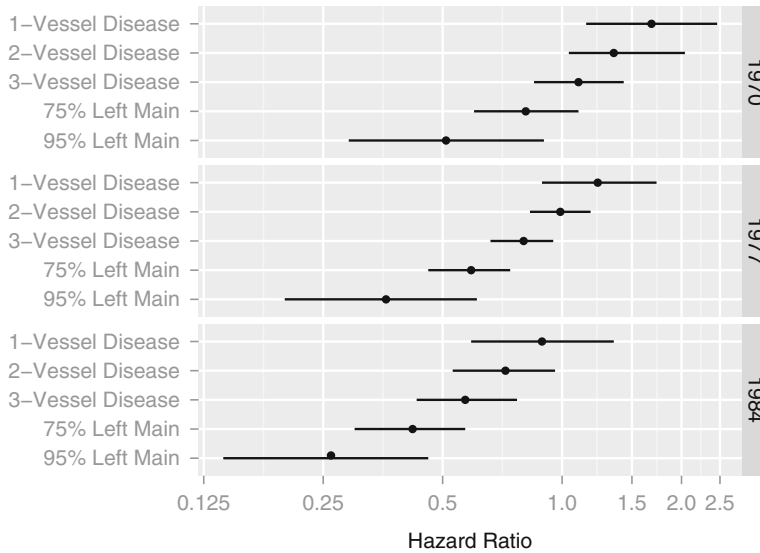


Fig. 20.12 A display of an interaction between treatment and extent of disease, and between treatment and calendar year of start of treatment. Comparison of medical and surgical average hazard ratios for patients treated in 1970, 1977, and 1984 according to coronary disease severity. Circles represent point estimates; bars represent 0.95 confidence limits of hazard ratios. Ratios less than 1.0 indicate that coronary bypass surgery is more effective.⁸⁸

that factor is polytomous or predictions are made for specific levels, survival curves (with or without adjustment for other factors not shown) can be drawn for each level of the predictor of interest, with follow-up time on the x -axis. Figure 20.2 demonstrated this for a factor which was a stratification factor. Figure 20.13 extends this by displaying survival estimates stratified by treatment but adjusted to various levels of two modeled factors, one of which, year of diagnosis, interacted with treatment.

When a continuous predictor is of interest, it is usually more informative to display that factor on the x -axis with estimated survival at one or more time points on the y -axis. When the model contains only one predictor, even if that predictor is represented by multiple terms such as a spline expansion, one may simply plot that factor against the predicted survival. Figure 20.14 depicts the relationship between treadmill exercise score, which is a weighted linear combination of several predictors in a Cox model, and the probability of surviving five years.

When displaying the effect of a single factor after adjusting for multiple predictors which are not displayed, care only need be taken for the values to which the predictors are adjusted (e.g., grand means). When instead the desire is to display the effect of multiple predictors simultaneously, an important continuous predictor can be displayed on the x -axis while separate

curves or graphs are made for levels of other factors. Figure 20.15, which corresponds to the log Λ plots in Figure 20.5, displays the joint effects of age and sex on the three-year survival probability. Age is modeled with a cubic spline function, and the model includes terms for an age \times sex interaction.

```
p ← Predict(f.ia, age, sex, time=3)
ggplot(p)
```

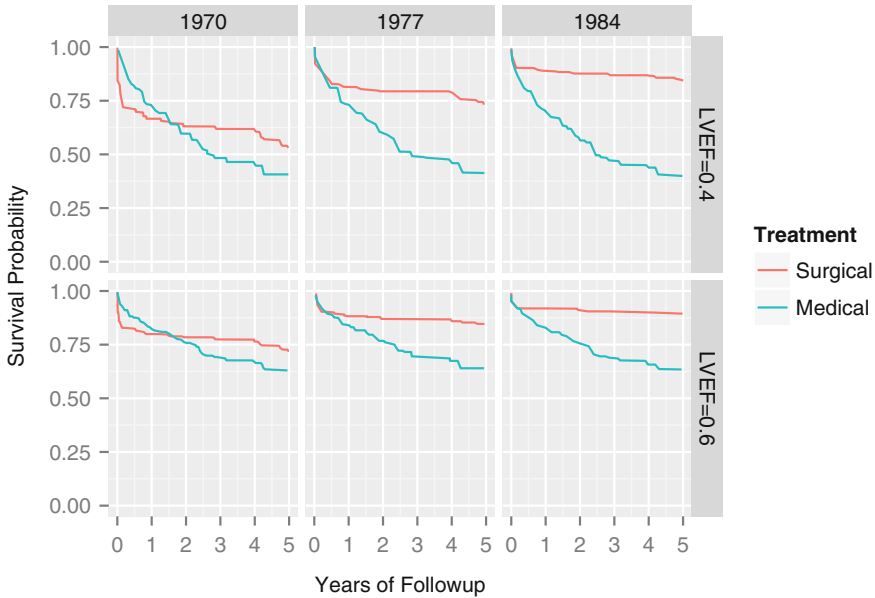


Fig. 20.13 Cox–Kalbfleisch–Prentice survival estimates stratifying on treatment and adjusting for several predictors, showing a secular trend in the efficacy of coronary artery bypass surgery. Estimates are for patients with left main disease and normal (LVEF=0.6) or impaired (LVEF=0.4) ventricular function.⁵¹⁶

Besides making graphs of survival probabilities estimated for given levels of the predictors, nomograms have some utility in specifying a fitted Cox model. A nomogram can be used to compute $X\hat{\beta}$, the estimated log hazard for a subject with a set of predictor values X relative to the “standard” subject. The central line in the nomogram will be on this linear scale unlike the logistic model nomograms given in Section 10.10 which further transformed $X\hat{\beta}$ into $[1 + \exp(-X\hat{\beta})]^{-1}$. Alternatively, the central line could be on the nonlinear $\exp(X\hat{\beta})$ hazard ratio scale or survival at fixed t .

A graph of the estimated underlying survival function $\hat{S}(t)$ as a function of t can be coupled with the nomogram used to compute $X\hat{\beta}$. The survival for a specific subject, $\hat{S}(t|X)$ is obtained from $\hat{S}(t)^{\exp(X\hat{\beta})}$. Alternatively, one could graph $\hat{S}(t)^{\exp(X\hat{\beta})}$ for various values of $X\hat{\beta}$ (e.g., $X\hat{\beta} = -2, -1, 0, 1, 2$)

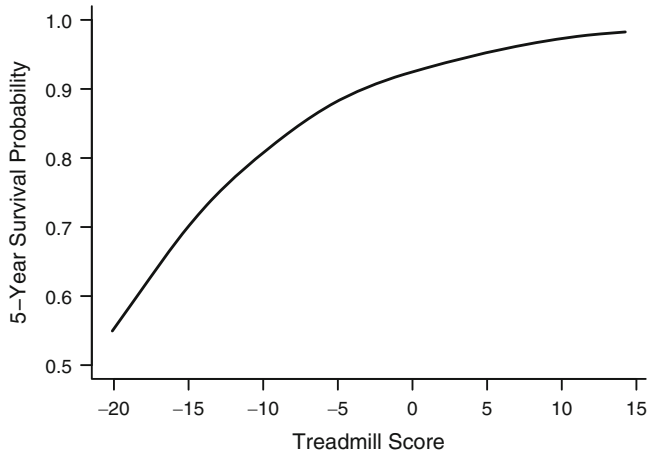


Fig. 20.14 Cox model predictions with respect to a continuous variable. X -axis shows the range of the treadmill score seen in clinical practice and Y -axis shows the corresponding five-year survival probability predicted by the Cox regression model for the 2842 study patients.⁴⁴⁰

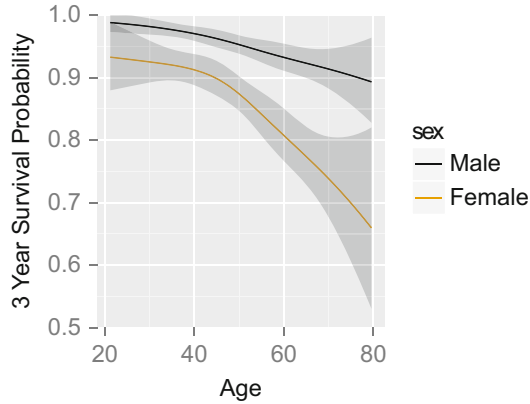


Fig. 20.15 Survival estimates for model stratified on sex, with interaction.

so that the desired survival curve could be read directly, at least to the nearest tabulated $X\hat{\beta}$. For estimating survival at a fixed time, say two years, one only need to provide the constant $\hat{S}(t)$. The nomogram could even be adapted to include a nonlinear scale $\hat{S}(2)^{\exp(X\hat{\beta})}$ to allow direct computation of two-year survival.

20.13 R Functions

Harrell's `cpower`, `spower`, and `ciapower` (in the `Hmisc` package) perform power calculations for Cox tests in follow-up studies. `cpower` computes power for a two-sample Cox (log-rank) test with random patient entry over a fixed duration and a given length of minimum follow-up. The expected number of events in each group is estimated by assuming exponential survival. `cpower` uses a slight modification of the method of Schoenfeld⁵⁵⁸ (see [501]). Separate specification of noncompliance in the active treatment arm and “drop-in” from the control arm into the active arm are allowed, using the method of Lachin and Foulkes.³⁷⁰ The `ciapower` function computes power of the Cox interaction test in a 2×2 setup using the method of Peterson and George.⁵⁰¹ It does not take noncompliance into account. The `spower` function simulates power for two-sample tests (the log-rank test by default) allowing for very complex conditions such as continuously varying treatment effect and noncompliance probabilities.

The `rms` package `cph` function is a slight modification of the `coxph` function written by Terry Therneau (in his `survival` package to work in the `rms` framework. `cph` computes MLEs of Cox and stratified Cox PH models, overall score and likelihood ratio χ^2 statistics for the model, martingale residuals, the linear predictor ($X\hat{\beta}$ centered to have mean 0), and collinearity diagnostics. Efron, Breslow, and exact partial likelihoods are supported (although the exact likelihood is very computationally intensive if ties are frequent). The function also fits the Andersen–Gill²³ generalization of the Cox PH model. This model allows for predictor values to change over time in the form of step functions as well as allowing time-dependent stratification (subjects can jump to different hazard function shapes). The Andersen–Gill formulation allows multiple events per subject and permits subjects to move in and out of risk at any desired time points. The latter feature allows time zero to have a more general definition. (See Section 9.5 for methods of adjusting the variance–covariance matrix of $\hat{\beta}$ for dependence in the events per subject.) The printing function corresponding to `cph` prints the Nagelkerke index R_N^2 described in Section 20.10, and has a `latex` option for better output. `cph` works in conjunction with the generic functions such as `specs`, `predict`, `summary`, `anova`, `fastbw`, `which.influence`, `latex`, `residuals`, `coef`, `nomogram`, and `Predict` described in Section 20.13, the same as the logistic regression function `lrm` does. For the purpose of plotting predicted survival at a single time, `Predict` has an additional argument `time` for plotting `cph` fits. It also has an argument `loglog` which if `TRUE` causes instead log-log survival to be plotted on the y -axis. `cph` has all the arguments described in Section 20.13 and some that are specific to it.

Similar to functions for `psm`, there are `Survival`, `Quantile`, and `Mean` functions which create other R functions to evaluate survival probabilities and perform other calculations, based on a `cph` fit with `surv=TRUE`. These functions, unlike all the others, allow polygon (linear interpolation) estimation of survival

probabilities, quantiles, and mean survival time as an option. `Quantile` is the only automatic way for obtaining survival quantiles with `cph`. Quantile estimates will be missing when the survival curve does not extend long enough. Likewise, survival estimates will be missing for $t >$ maximum follow-up time, when the last event time is censored. `Mean` computes the mean survival time if the last failure time in each stratum is uncensored. Otherwise, `Mean` may be used to compute restricted mean lifetime using a user-specified truncation point.³³⁴ `Quantile` and `Mean` are especially useful with `plot` and `nomogram`. `Survival` is useful with `nomogram`.

The R program below demonstrates how several `cph`-related functions work well with the `nomogram` function. Here predicted three-year survival probabilities and median survival time (when defined) are displayed against age and sex from the previously simulated dataset. The fact that a nonlinear effect interacts with a stratified factor is taken into account.

```

surv      <- Survival(f.ia)
surv.f    <- function(lp) surv(3, lp, stratum='sex=Female')
surv.m    <- function(lp) surv(3, lp, stratum='sex=Male')
quant     <- Quantile(f.ia)
med.f     <- function(lp) quant(.5, lp, stratum='sex=Female')
med.m     <- function(lp) quant(.5, lp, stratum='sex=Male')
at.surv   <- c(.01, .05, seq(.1,.9,by=.1), .95, .98, .99, .999)
at.med    <- c(0, .5, 1, 1.5, seq(2, 14, by=2))
n <- nomogram(f.ia, fun=list(surv.m, surv.f, med.m, med.f),
             funlabel=c('S(3 | Male)', 'S(3 | Female)',
                       'Median (Male)', 'Median (Female)'),
             fun.at=list(c(.8,.9,.95,.98,.99),
                        c(.1,.3,.5,.7,.8,.9,.95,.98),
                        c(8,10,12),c(1,2,4,8,12)))
plot(n, col.grid=FALSE, lmgp=.2)
latex(f.ia, file='', digits=3)

```

$$\text{Prob}\{T \geq t \mid \text{sex} = i\} = S_i(t)^{e^{X\beta}}, \quad \text{where}$$

$$\begin{aligned}
X\hat{\beta} = & \\
& -1.8 \\
& +0.0493\text{age} - 2.15 \times 10^{-6}(\text{age} - 30.3)_+^3 - 2.82 \times 10^{-5}(\text{age} - 45.1)_+^3 \\
& +5.18 \times 10^{-5}(\text{age} - 54.6)_+^3 - 2.15 \times 10^{-5}(\text{age} - 69.6)_+^3 \\
& +[\text{Female}][-0.0366\text{age} + 4.29 \times 10^{-5}(\text{age} - 30.3)_+^3 - 0.00011(\text{age} - 45.1)_+^3 \\
& +6.74 \times 10^{-5}(\text{age} - 54.6)_+^3 - 2.32 \times 10^{-7}(\text{age} - 69.6)_+^3]
\end{aligned}$$

and $[c] = 1$ if subject is in group c , 0 otherwise; $(x)_+ = x$ if $x > 0$, 0 otherwise.

t	$S_{Male}(t)$	$S_{Female}(t)$
0	1.000	1.000
1	0.993	0.902
2	0.984	0.825
3	0.975	0.725
4	0.967	0.648
5	0.956	0.576
6	0.947	0.520
7	0.938	0.481
8	0.928	0.432
9	0.920	0.395
10	0.909	0.358
11	0.904	0.314
12	0.892	0.268
13	0.886	0.223
14	0.877	0.203

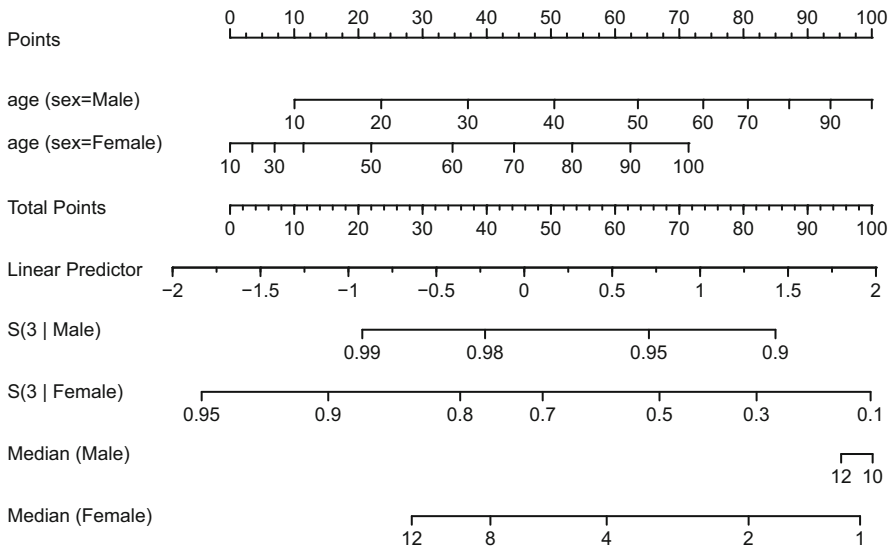


Fig. 20.16 Nomogram from a fitted stratified Cox model that allowed for interaction between age and sex, and nonlinearity in age. The axis for median survival time is truncated on the left where the median is beyond the last follow-up time.

```
rcspline.plot(lvef, d.time, event=cdeath, nk=3)
```

The corresponding smoothed martingale residual plot for LVEF in Figure 20.7 was created with

```

cox ← cph(Surv(d.time, cdeath) ~ lvef, iter.max=0)
res ← resid(cox)
g ~ loess(res ~ lvef)
plot(g, coverage=0.95, confidence=7, xlab="LVEF",
      ylab="Martingale Residual")
g ← ols(res ~ rcs(lvef,5))
plot(g, lvef=NA, add=T, lty=2)
lines(lowess(lvef, res, iter=0), lty=3)
legend(.3, 1.15, c("loess Fit and 0.95 Confidence Bars",
                  "ols Spline Fit and 0.95 Confidence Limits",
                  "lowess Smoother"), lty=1:3, bty="n")

```

Because we desired residuals with respect to the omitted predictor LVEF, the parameter `iter.max=0` had to be given to make `cph` stop the estimation process at the starting parameter estimates (default of zero). The effect of this is to ignore the predictors when computing the residuals; that is, to compute residuals from a flat line rather than the usual residuals from a fitted straight line.

The `residuals` function is a slight modification of Therneau's `residuals.coxph` function to obtain martingale, Schoenfeld, score, deviance residuals, or approximate DFBETA or DFBETAS. Since martingale residuals are always stored by `cph` (assuming there are covariables present), `residuals` merely has to pick them off the fit object and reinsert rows that were deleted due to missing values. For other residuals, you must have stored the design matrix and `Surv` object with the fit by using `..., x=TRUE, y=TRUE`. Storing the design matrix with `x=TRUE` ensures that the same transformation parameters (e.g., knots) are used in evaluating the model as were used in fitting it. To use `residuals` you can use the abbreviation `resid`. See the help file for `residuals` for an example of how martingale residuals may be used to quickly plot univariable (unadjusted) relationships for several predictors.

Figure 20.10, which used smoothed scaled Schoenfeld partial residuals⁵⁵⁷ to estimate the form of a predictor's log hazard ratio over time, was made with

```

Srv ← Surv(dm.time, cdeathmi)
cox ← cph(Srv ~ pi, x=T, y=T)
cox.zph(cox, "rank") # Test for PH for each column of X
res ← resid(cox, "scaledsch")
time ← as.numeric(names(res))
# Use dimnames(res)[[1]] if more than one predictor
f ← loess(res ~ time, span=0.50)
plot(f, coverage=0.95, confidence=7, xlab="t",
      ylab="Scaled Schoenfeld Residual", ylim=c(-.1, .25))
lines(supsmu(time, res), lty=2)
legend(1.1, .21, c("loess Smoother with span=0.50 and 0.95 C.L.",
                  "Super Smoother"), lty=1:2, bty="n")

```

The computation and plotting of scaled Schoenfeld residuals could have been done automatically in this case by using the single command `plot(cox.zph(cox))`, although `cox.zph` defaults to plotting against the Kaplan–Meier transformation of follow-up time.

The `hazard.ratio.plot` function in `rms` repeatedly estimates Cox regression coefficients and confidence limits within time intervals. The log hazard ratios are plotted against the mean failure/censoring time within the interval. Figure 20.9 was created with

```
hazard.ratio.plot(pi, S) # S was Surv(dm.time, ...)
```

If you have multiple degree of freedom factors, you may want to score them into linear predictors before using `hazard.ratio.plot`. The `predict` function with argument `type="terms"` will produce a matrix with one column per factor to do this (Section 20.13).

Therneau's `cox.zph` function implements Harrell's Schoenfeld residual correlation test for PH. This function also stores results that can easily be passed to a plotting method for `cox.zph` to automatically plot smoothed residuals that estimate the effect of each predictor over time.

Therneau has also written an R function `survdiff` that compares two or more survival curves using the $G - \rho$ family of rank tests (Harrington and Fleming²⁷³).

The `rcorr.cens` function in the `Hmisc` library computes the c index and the corresponding generalization of Somers' D_{xy} rank correlation for a censored response variable. `rcorr.cens` also works for uncensored and binary responses (see ROC area in Section 10.8), although its use of all possible pairings makes it slow for this purpose. The `survival` package's `survConcordance` has an extremely fast algorithm for the c index and a fairly accurate estimator of its standard error.

The `calibrate` function for `cph` constructs a bootstrap or cross-validation optimism-corrected calibration curve for a single time point by resampling the differences between average Cox predicted survival and Kaplan–Meier estimates (see Section 20.11.1). But more precise is `calibrate`'s default method based on adaptive semiparametric regression discussed in the same section. Figure 20.11 is an example.

The `validate` function for `cph` fits validates several statistics describing Cox model fits—slope shrinkage, R_N^2 , D , U , Q , and D_{xy} . The `val.surv` function can also be of use in externally validating a Cox model using the methods presented in Section 18.3.7.

20

20.14 Further Reading

- [1] Good general texts for the Cox PH model include Cox and Oakes,¹³³ Kalbfleisch and Prentice,³³¹ Lawless,³⁸² Collett,¹¹⁴ Marubini and Valsecchi,⁴⁴⁴ and Klein and Moeschberger.³⁵⁰ Therneau and Grambsch⁶⁰⁴ describe the many ways the standard Cox model may be extended.
- [2] Cupples et al.¹⁴¹ and Marubini and Valsecchi [444, pp. 201–206] present good description of various methods of computing “adjusted survival curves.”
- [3] See Altman and Andersen¹⁵ for simpler approximate formulas. Cheng et al.¹⁰³ derived methods for obtaining pointwise and simultaneous confidence bands for

- $S(t)$ for future subjects, and Henderson²⁸² has a comprehensive discussion of the use of Cox models to estimate survival time for individual subjects.
- [4] Aalen² and Valsecchi et al.⁶²⁵ discuss other residuals useful in graphically checking survival model assumptions. León and Tsai⁴⁰⁰ derived residuals for estimating covariate transformations that are different from martingale residuals.
 - [5] [411] has other methods for generating confidence intervals for martingale residual plots.
 - [6] Lin et al.⁴¹¹ describe other methods of checking transformations using *cumulative* martingale residuals.
 - [7] A parametric analysis of the VA dataset using linear splines and incorporating $X \times t$ interactions is found in [361].
 - [8] Winnett and Sasieni⁶⁷¹ show how to use scaled Schoenfeld residuals in an iterative fashion to actually model effects that are not in proportional hazards.
 - [9] See [233, 503] for some methods for obtaining confidence bands for Schoenfeld residual plots. Winnett and Sasieni⁶⁷⁰ discuss conditions in which the Grambsch–Therneau scaling of the Schoenfeld residuals does not perform adequately for estimating $\beta(t)$.
 - [10] [475, 519] compared the power of the test for PH based on the correlation between failure time and Schoenfeld residuals with the power of several other tests.
 - [11] See Lin et al.⁴¹¹ for another approach to deriving a formal test of PH using residuals. Other graphical methods for examining the PH assumption are due to Gray,²³⁶ who used hazard smoothing to estimate hazard ratios as a function of time, and Thaler,⁶⁰² who developed a nonparametric estimator of the hazard ratio over time for time-dependent covariables. See Valsecchi et al.⁶²⁵ for other useful graphical assessments of PH.
 - [12] A related test of constancy of hazard ratios may be found in [519]. Also, see Schemper⁵⁴⁷ for related methods.
 - [13] See [547] for a variation of the standard Cox likelihood to allow for non-PH.
 - [14] An excellent review of graphical methods for assessing PH may be found in Hess.²⁹⁰ Sahoo and Sengupta⁵³⁷ provide some new graphical methods for assessing PH irrespective of satisfaction of the other model assumptions.
 - [15] Schemper⁵⁴⁷ provides a way to determine the effect of falsely assuming PH by comparing the Cox regression coefficient with a well-described average log hazard ratio. Zucker⁶⁹¹ shows how dependent a weighted log-rank test is on the true hazard ratio function, when the weights are derived from a hypothesized hazard ratio function. Valsecchi et al.⁶²⁵ proposed a method that is robust to non-PH that occurs in the late follow-up period. Their method uses down-weighting of certain types of “outliers.” See Herndon and Harrell²⁸⁷ for a flexible parametric PH model with time-dependent covariables, which uses the restricted cubic spline function to specify $\lambda(t)$. Putter et al.⁵¹⁸ and Muggeo and Tagliavia⁴⁶⁸ have nice approaches that use time-dependent covariates to model time interactions to allow non-proportional hazards. Perperoglou et al.^{498, 499} developed a systematic approach that allows one to continuously vary the amount of non PH allowed, through the use of a structure matrix that connects predictors with functions of time. Schuabel et al.⁵⁴³ have a good exposition of internal time-dependent covariates.
 - [16] See van Houwelingen and le Cessie [633, Eq. 61] and Verweij and van Houwelingen⁶⁴⁰ for an interesting index of cross-validated predictive accuracy. Schemper and Henderson⁵⁵² relate explained variation to predictive accuracy in Cox models. Hielscher et al.²⁹¹ compares and illustrates several measures of explained variation as does Choodari-Oskooei et al.¹⁰⁶. Choodari-Oskooei et al.¹⁰⁵ studied explained randomness and predictive accuracy measures.
 - [17] See similar indexes in Schemper⁵⁴⁴ and a related idea in [633, Eq. 63]. Mandel, Galai, and Simchen⁴³⁶ presented a time-varying c index. See Korn and

Simon,³⁶⁵ Schemper and Stare,⁵⁵⁴ and Henderson²⁸² for nice comparisons of various measures. Pencina and D'Agostino⁴⁸⁹ provide more details about the c index and derived new interval estimates. They also discussed the relationship between c and a version of Kendall's τ . Pencina et al.⁴⁹¹ found advantages of c . Uno et al.⁶¹⁸ described exactly how c depends on the amount of censoring and proposed a new index, requiring one to choose a time cutoff, that is invariant to the amount of censoring. Henderson et al.²⁸³ discussed the benefits of using the probability of a serious prognostication error (e.g., being off by a factor of 2.0 or worse on the time scale) as an accuracy measure. Schemper⁵⁵⁰ shows that models with very important predictors can have very low absolute prediction ability, and he discusses measures of predictive accuracy from a general standpoint. Lawless and Yuan³⁸⁶ present prediction error estimators and confidence limits, focusing on such measures as error in predicted median or mean survival time. Schmid and Potapov⁵⁵⁵ studied the bias of several variations on the c index under non-proportional hazards and/or nonrandom censoring. Gönen and Heller²²³ developed a c -index that is censoring-independent.

- [18] Altman and Royston¹⁸ have a good discussion of validation of prognostic models and present several examples of validation using a simple discrimination index. Thomas Gerds has an R package `pec` that provides many validation methods and accuracy indexes.
- [19] Kattan et al.³³⁸ describe how to make nomograms for deriving predicted survival probabilities when there are competing risks.
- [20] Hielscher et al.²⁹¹ provides an overview of software for computing accuracy indexes with censored data.

Chapter 21

Case Study in Cox Regression

21.1 Choosing the Number of Parameters and Fitting the Model

Consider the randomized trial of estrogen for treatment of prostate cancer⁸⁷ described in Chapter 8. Let us now develop a model for time until death (of any cause). There are 354 deaths among the 502 patients. To be able to efficiently estimate treatment benefit, to test for differential treatment effect, or to estimate prognosis or absolute treatment benefit for individual patients, we need a multivariable survival model. In this case study we do not make use of data reductions obtained in Chapter 8 but show simpler (partial) approaches to data reduction. We do use the `transcan` results for imputation.

First let's assess the wisdom of fitting a full additive model that does not assume linearity of effect for any predictor. Categorical predictors are expanded using dummy variables. For `pf` we could lump the last two categories as before since the last category has only two patients. Likewise, we could combine the last two levels of `ekg`. Continuous predictors are expanded by fitting four-knot restricted cubic spline functions, which contain two nonlinear terms and thus have a total of three d.f. Table 21.1 defines the candidate predictors and lists their d.f. The variable `stage` is not listed as it can be predicted with high accuracy from `sz,sg,ap,bm` (`stage` could have been used as a predictor for imputing missing values on `sz, sg`). There are a total of 36 candidate d.f. that should not be artificially reduced by “univariable screening” or graphical assessments of association with death. This is about 1/10 as many predictor d.f. as there are deaths, so there is some hope that a fitted model may validate. Let us also examine this issue by estimating the amount of shrinkage using Equation 4.3. We first use `transcan` impute missing data.

```
require(rms)
```

Table 21.1 Initial allocation of degrees of freedom

Predictor	Name	d.f.	Original Levels
Dose of estrogen	rx	3	placebo, 0.2, 1.0, 5.0 mg estrogen
Age in years	age	3	
Weight index: $wt(kg) - ht(cm) + 200$	wt	3	
Performance rating	pf	2	normal, in bed < 50% of time, in bed > 50%, in bed always
History of cardiovascular disease	hx	1	present/absent
Systolic blood pressure/10	sbp	3	
Diastolic blood pressure/10	dbp	3	
Electrocardiogram code	ekg	5	normal, benign, rhythm disturb., block, strain, old myocardial infarction, new MI
Serum hemoglobin (g/100ml)	hg	3	
Tumor size (cm ²)	sz	3	
Stage/histologic grade combination	sg	3	
Serum prostatic acid phosphatase	ap	3	
Bone metastasis	bm	1	present/absent

```

getHdata(prostate)
levels(prostate$ekg)[levels(prostate$ekg) %in%
  c('old MI','recent MI')] ← 'MI'
# combines last 2 levels and uses a new name, MI

prostate$pf.coded ← as.integer(prostate$pf)
# save original pf, re-code to 1-4
levels(prostate$pf) ← c(levels(prostate$pf)[1:3],
  levels(prostate$pf)[3])
# combine last 2 levels

w ← transcan(~ sz + sg + ap + sbp + dbp + age +
  wt + hg + ekg + pf + bm + hx, imputed=TRUE,
  data=prostate, pl=FALSE, pr=FALSE)

attach(prostate)
sz ← impute(w, sz, data=prostate)
sg ← impute(w, sg, data=prostate)
age ← impute(w, age, data=prostate)
wt ← impute(w, wt, data=prostate)
ekg ← impute(w, ekg, data=prostate)

dd ← datadist(prostate); options(datadist='dd')

```

```

units(dtime) ← 'Month'
S ← Surv(dtime, status != 'alive')

f ← cph(S ~ rx + rcs(age,4) + rcs(wt,4) + pf + hx +
        rcs(sbp,4) + rcs(dbp,4) + ekg + rcs(hg,4) +
        rcs(sg,4) + rcs(sz,4) + rcs(log(ap),4) + bm)

print(f, latex=TRUE, coefs=FALSE)

```

Cox Proportional Hazards Model

```

cph(formula = S ~ rx + rcs(age, 4) + rcs(wt, 4) + pf + hx
    + rcs(sbp, 4) + rcs(dbp, 4) + ekg + rcs(hg, 4)
    + rcs(sg, 4) + rcs(sz, 4) + rcs(log(ap), 4) + bm)

```

		Model Tests		Discrimination Indexes	
Obs	502	LR χ^2	136.22	R^2	0.238
Events	354	d.f.	36	D_{xy}	0.333
Center	-2.9933	$\Pr(> \chi^2)$	0.0000	g	0.787
		Score χ^2	143.62	g_r	2.196
		$\Pr(> \chi^2)$	0.0000		

The likelihood ratio χ^2 statistic is 136.2 with 36 d.f. This test is highly significant so some modeling is warranted. The AIC value (on the χ^2 scale) is $136.2 - 2 \times 36 = 64.2$. The rough shrinkage estimate is 0.74 ($100.2/136.2$) so we estimate that 0.26 of the model fitting will be noise, especially with regard to calibration accuracy. The approach of Spiegelhalter⁵⁸² is to fit this full model and to shrink predicted values. We instead try to do data reduction (blinded to individual χ^2 statistics from the above model fit) to see if a reliable model can be obtained without shrinkage. A good approach at this point might be to do a variable clustering analysis followed by single degree of freedom scoring for individual predictors or for clusters of predictors. Instead we do an informal data reduction. The strategy is described in Table 21.2. For `ap`, more exploration is desired to be able to model the shape of effect with such a highly skewed distribution. Since we expect the tumor variables to be strong prognostic factors we retain them as separate variables. No assumption is made for the dose-response shape for estrogen, as there is reason to expect a non-monotonic effect due to competing risks for cardiovascular death.

```

heart ← hx + ekg %nin% c('normal', 'benign')
label(heart) ← 'Heart Disease Code'
map ← (2*dbp + sbp)/3
label(map) ← 'Mean Arterial Pressure/10'
dd ← datadist(dd, heart, map)

f ← cph(S ~ rx + rcs(age,4) + rcs(wt,3) + pf.coded +

```

Table 21.2 Final allocation of degrees of freedom

Variables	Reductions	d.f. Saved
wt	Assume variable not important enough for 4 knots; use 3 knots	1
pf	Assume linearity	1
hx,ekg	Make new 0,1,2 variable and assume linearity: 2 = hx and ekg not normal or benign, 1 = either, 0 = none	5
sbp,dbp	Combine into mean arterial bp and use 3 knots: map = (2 dbp + sbp)/3	4
sg	Use 3 knots	1
sz	Use 3 knots	1
ap	Look at shape of effect of ap in detail, and take log before expanding as spline to achieve numerical stability: add 1 knots	-1

```
heart + rcs(map,3) + rcs(hg,4) +
rcs(sg,3) + rcs(sz,3) + rcs(log(ap),5) + bm,
x=TRUE, y=TRUE, surv=TRUE, time.inc=5*12)
print(f, latex=TRUE, coefs=3)
```

Cox Proportional Hazards Model

```
cph(formula = S ~ rx + rcs(age, 4) + rcs(wt, 3) + pf.coded +
heart + rcs(map, 3) + rcs(hg, 4) + rcs(sg, 3) +
rcs(sz, 3) + rcs(log(ap), 5) + bm, x = TRUE, y = TRUE,
surv = TRUE, time.inc = 5 * 12)
```

		Model Tests	Discrimination Indexes
Obs	502	LR χ^2 118.37	R^2 0.210
Events	354	d.f. 24	D_{xy} 0.321
Center	-2.4307	$\Pr(> \chi^2)$ 0.0000	g 0.717
		Score χ^2 125.58	g_r 2.049
		$\Pr(> \chi^2)$ 0.0000	

	Coef	S.E.	Wald Z	$\Pr(> Z)$
rx=0.2 mg estrogen	-0.0002	0.1493	0.00	0.9987
rx=1.0 mg estrogen	-0.4160	0.1657	-2.51	0.0121
rx=5.0 mg estrogen	-0.1107	0.1571	-0.70	0.4812
...				

Table 21.3 Wald Statistics for S

	χ^2	d.f.	<i>P</i>
rx	8.01	3	0.0459
age	13.84	3	0.0031
<i>Nonlinear</i>	9.06	2	0.0108
wt	8.21	2	0.0165
<i>Nonlinear</i>	2.54	1	0.1110
pf.coded	3.79	1	0.0517
heart	23.51	1	< 0.0001
map	0.04	2	0.9779
<i>Nonlinear</i>	0.04	1	0.8345
hg	12.52	3	0.0058
<i>Nonlinear</i>	8.25	2	0.0162
sg	1.64	2	0.4406
<i>Nonlinear</i>	0.05	1	0.8304
sz	12.73	2	0.0017
<i>Nonlinear</i>	0.06	1	0.7990
ap	6.51	4	0.1639
<i>Nonlinear</i>	6.22	3	0.1012
bm	0.03	1	0.8670
TOTAL NONLINEAR	23.81	11	0.0136
TOTAL	119.09	24	< 0.0001

```
# x, y for predict, validate, calibrate;
# surv, time.inc for calibrate
latex(anova(f), file='', label='tab:coxcase-anova1')# Table 21.3
```

The total savings is thus 12 d.f. The likelihood ratio χ^2 is 118 with 24 d.f., with a slightly improved AIC of 70. The rough shrinkage estimate is slightly better at 0.80, but still worrisome. A further data reduction could be done, such as using the `transcan` transformations determined from self-consistency of predictors, but we stop here and use this model.

From Table 21.3 there are 11 parameters associated with nonlinear effects, and the overall test of linearity indicates the strong presence of nonlinearity for at least one of the variables `age`, `wt`, `map`, `hg`, `sz`, `sg`, `ap`. There is no strong evidence for a difference in survival time between doses of estrogen.

21.2 Checking Proportional Hazards

Now that we have a tentative model, let us examine the model's distributional assumptions using smoothed scaled Schoenfeld residuals. A messy detail is how to handle multiple regression coefficients per predictor. Here we do an

approximate analysis in which each predictor is scored by adding up all that predictor's terms in the model, to transform that predictor to optimally relate to the log hazard (at least if the *shape* of the effect does not change with time). In doing this we are temporarily ignoring the fact that the individual regression coefficients were estimated from the data. For dose of estrogen, for example, we code the effect as 0 (placebo), -0.00025 (0.2 mg), -0.416 (1.0 mg), and -0.111 (5.0 mg), and `age` is transformed using its fitted spline function. In the `rms` package the `predict` function easily summarizes multiple terms and produces a matrix (here, `z`) containing the total effects for each predictor. Matrix factors can easily be included in model formulas.

```
z <- predict(f, type='terms')
# required x=T above to store design matrix
f.short <- cph(S ~ z, x=TRUE, y=TRUE)
# store raw x, y so can get residuals
```

The fit `f.short` based on the matrix of single d.f. predictors `z` has the same LR χ^2 of 118 as the fit `f`, but with a falsely low 11 d.f. All regression coefficients are unity.

Now we compute scaled Schoenfeld residuals separately for each predictor and test the PH assumption using the “correlation with time” test. Also plot smoothed trends in the residuals. The `plot` method for `cox.zph` objects uses cubic splines to smooth the relationship.

```
phptest <- cox.zph(f.short, transform='identity')
phptest
```

	rho	chisq	p
rx	0.10232	4.00823	0.0453
age	-0.05483	1.05850	0.3036
wt	0.01838	0.11632	0.7331
pf.coded	-0.03429	0.41884	0.5175
heart	0.02650	0.30052	0.5836
map	0.02055	0.14135	0.7069
hg	-0.00362	0.00511	0.9430
sg	-0.05137	0.94589	0.3308
sz	-0.01554	0.08330	0.7729
ap	0.01720	0.11858	0.7306
bm	0.04957	0.95354	0.3288
GLOBAL	NA	7.18985	0.7835

```
plot(phptest, var='rx') # Figure 21.1
```

Perhaps only the drug effect significantly changes over time ($P = 0.05$ for testing the correlation `rho` between the scaled Schoenfeld residual and time), but when a global test of PH is done penalizing for 11 d.f., the P value is 0.78. A graphical examination of the trends doesn't find anything interesting for the last 10 variables. A residual plot is drawn for `rx` alone and is shown in Figure 21.1. We ignore the possible increase in effect of estrogen over time. If this non-PH is real, a more accurate model might be obtained by stratifying on `rx` or by using a `time × rx` interaction as a time-dependent covariable.

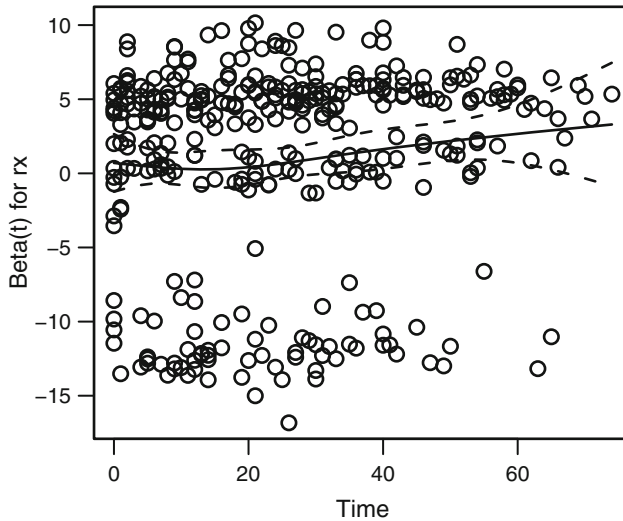


Fig. 21.1 Raw and spline-smoothed scaled Schoenfeld residuals for dose of estrogen, nonlinearly coded from the Cox model fit, with ± 2 standard errors.

21.3 Testing Interactions

Note that the model has several insignificant predictors. These are not deleted, as that would not improve predictive accuracy and it would make accurate confidence intervals hard to obtain. At this point it would be reasonable to test prespecified interactions. Here we test all interactions with dose. Since the multiple terms for many of the predictors (and for `rx`) make for a great number of d.f. for testing interaction (and a loss of power), we do approximate tests on the data-driven coding of predictors. P -values for these tests are likely to be somewhat anti-conservative.

```
z.dose ← z[,"rx"] # same as saying z[,1] - get first column
z.other ← z[,-1] # all but the first column of z
f.ia ← cph(S ~ z.dose * z.other) # Figure 21.4:
latex(anova(f.ia), file='', label='tab:coxcase-anova2')
```

The global test of additivity in Table 21.4 has $P = 0.27$, so we ignore the interactions (and also forget to penalize for having looked for them below!).

21.4 Describing Predictor Effects

Let us plot how each predictor is related to the log hazard of death, including 0.95 confidence bands. Note in Figure 21.2 that due to a peculiarity of the Cox model the standard error of the predicted $X\hat{\beta}$ is zero at the reference values (medians here, for continuous predictors).

Table 21.4 Wald Statistics for \mathbf{S}

	χ^2	d.f.	P
z.dose (Factor+Higher Order Factors)	18.74	11	0.0660
<i>All Interactions</i>	12.17	10	0.2738
z.other (Factor+Higher Order Factors)	125.89	20	< 0.0001
<i>All Interactions</i>	12.17	10	0.2738
z.dose \times z.other (Factor+Higher Order Factors)	12.17	10	0.2738
TOTAL	129.10	21	< 0.0001

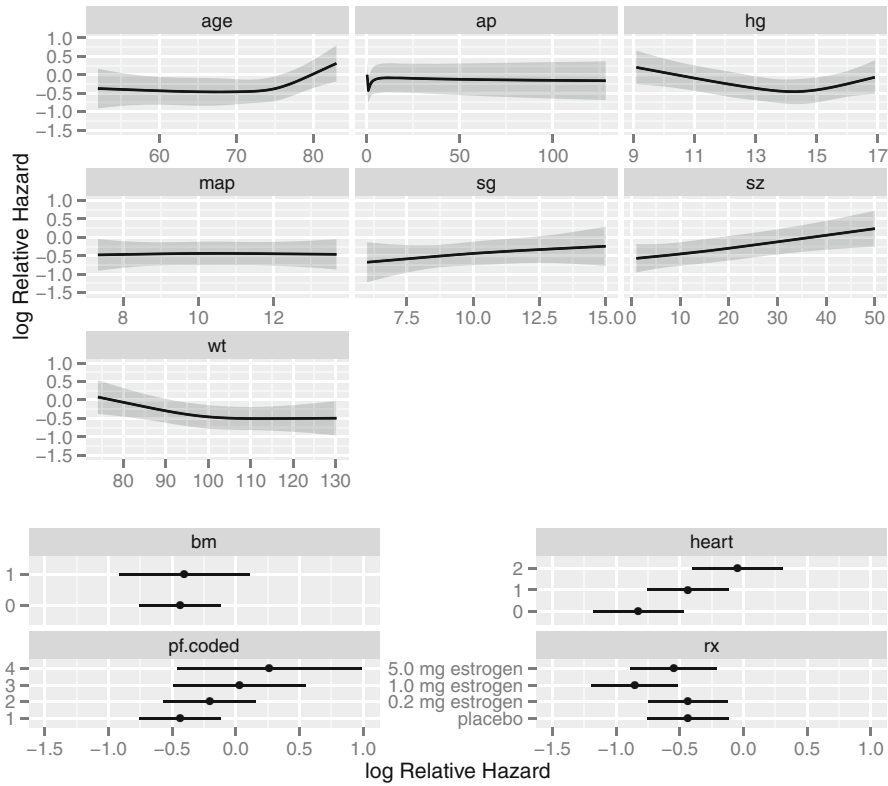


Fig. 21.2 Shape of each predictor on log hazard of death. Y -axis shows $X\hat{\beta}$, but the predictors not plotted are set to reference values. Note the highly non-monotonic relationship with *ap*, and the increased slope after age 70 which occurs in outcome models for various diseases.

```
ggplot(Predict(f), sepdiscrete='vertical', nlevels=4,
       vnames='names') # Figure 21.2
```

21.5 Validating the Model

We first validate this model for Somers' D_{xy} rank correlation between predicted log hazard and observed survival time, and for slope shrinkage. The bootstrap is used (with 300 resamples) to penalize for possible overfitting, as discussed in Section 5.3.

```
set.seed(1) # so can reproduce results
v <- validate(f, B=300)
```

Divergence or singularity in 83 samples

```
latex(v, file='')
```

	Original Sample	Training Sample	Test Sample	Optimism	Corrected Index	n
D_{xy}	0.3208	0.3454	0.2954	0.0500	0.2708	217
R^2	0.2101	0.2439	0.1754	0.0685	0.1417	217
Slope	1.0000	1.0000	0.7941	0.2059	0.7941	217
D	0.0292	0.0348	0.0238	0.0110	0.0182	217
U	-0.0005	-0.0005	0.0023	-0.0028	0.0023	217
Q	0.0297	0.0353	0.0216	0.0138	0.0159	217
g	0.7174	0.7918	0.6273	0.1645	0.5529	217

Here “training” refers to accuracy when evaluated on the bootstrap sample used to fit the model, and “test” refers to the accuracy when this model is applied without modification to the original sample. The apparent D_{xy} is 0.32, but a better estimate of how well the model will discriminate prognoses in the future is $D_{xy} = 0.27$. The bootstrap estimate of slope shrinkage is 0.79, close to the simple heuristic estimate. The shrinkage coefficient could easily be used to shrink predictions to yield better calibration.

Finally, we validate the model (without using the shrinkage coefficient) for calibration accuracy in predicting the probability of surviving five years. The bootstrap is used to estimate the optimism in how well predicted five-year survival from the final Cox model tracks flexible smooth estimates, without any binning of predicted survival probabilities or assuming proportional hazards.

```
cal ← calibrate(f, B=300, u=5*12, maxdim=4)
```

```
Using Cox survival estimates at 60 Months
```

```
plot(cal, subtitles=FALSE) # Figure 21.3
```

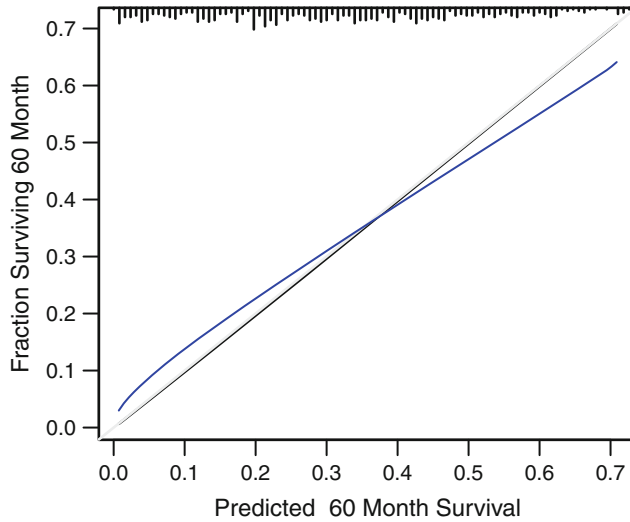


Fig. 21.3 Bootstrap estimate of calibration accuracy for 5-year estimates from the final Cox model, using adaptive linear spline hazard regression³⁶¹. The line nearer the ideal line corresponds to apparent predictive accuracy. The blue curve corresponds to bootstrap-corrected estimates.

The estimated calibration curves are shown in Figure 21.3, similar to what was done in Figure 19.11. Bootstrap calibration demonstrates some overfitting, consistent with regression to the mean. The absolute error is appreciable for 5-year survival predicted to be very low or high.

21.6 Presenting the Model

To present point and interval estimates of predictor effects we draw a hazard ratio chart (Figure 21.4), and to make a final presentation of the model we draw a nomogram having multiple “predicted value” axes. Since the ap relationship is so non-monotonic, use a 20 : 1 hazard ratio for this variable.

```
plot(summary(f, ap=c(1,20)), log=TRUE, main='') # Figure 21.4
```

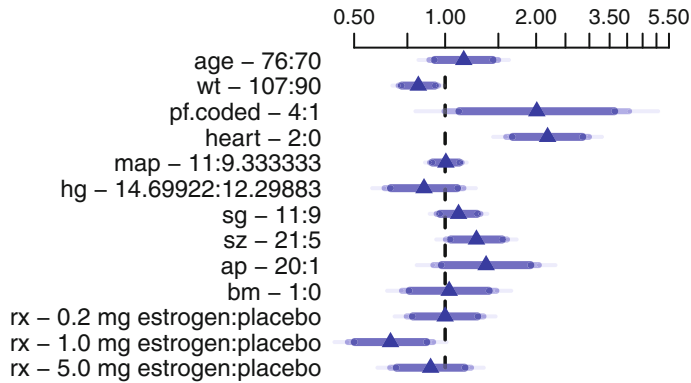


Fig. 21.4 Hazard ratios and multi-level confidence bars for effects of predictors in model, using default ranges except for ap

The ultimate graphical display for this model will be a nomogram relating the predictors to $X\hat{\beta}$, estimated three- and five-year survival probabilities and median survival time. It is easy to add as many “output” axes as desired to a nomogram.

```

surv  <- Survival(f)
surv3 <- function(x) surv(3*12, lp=x)
surv5 <- function(x) surv(5*12, lp=x)
quan  <- Quantile(f)
med   <- function(x) quan(lp=x)/12
ss    <- c(.05, .1, .2, .3, .4, .5, .6, .7, .8, .9, .95)

nom  <- nomogram(f, ap=c(.1, .5, 1, 2, 3, 4, 5, 10, 20, 30, 40),
                fun=list(surv3, surv5, med),
                funlabel=c('3-year Survival ', '5-year Survival ',
                           'Median Survival Time (years)'),
                fun.at=list(ss, ss, c(.5, 1:6)))
plot(nom, xfrac=.65, lmgp=.35) # Figure 21.5

```

21.7 Problems

Perform Cox regression analyses of survival time using the Mayo Clinic PBC dataset described in Section 8.9. Provide model descriptions, parameter estimates, and conclusions.

1. Assess the nature of the association of several predictors of your choice. For polytomous predictors, perform a log-rank-type score test (or k -sample ANOVA extension if there are more than two levels). For continuous predictors, plot a smooth curve that estimates the relationship between the predictor and the log hazard or log-log survival. Use both parametric and nonparametric (using martingale residuals) approaches. Make a test of H_0 : predictor is not associated with outcome versus H_a : predictor

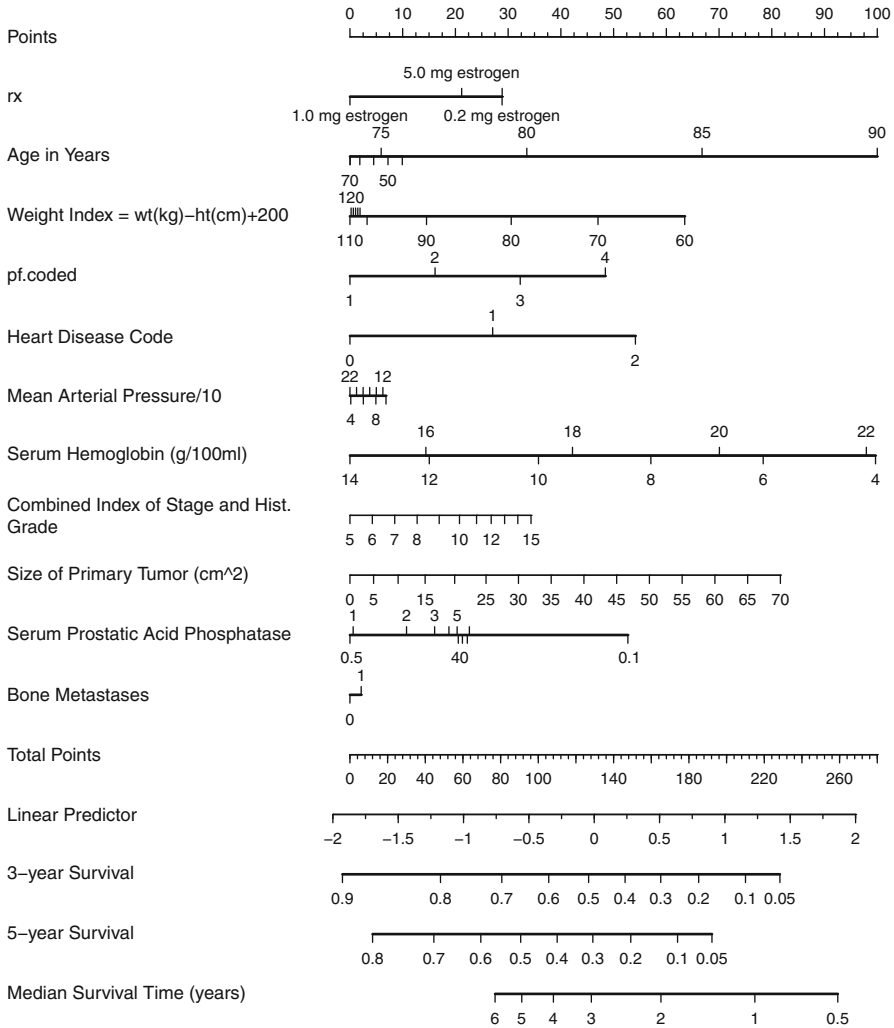


Fig. 21.5 Nomogram for predicting death in prostate cancer trial

is associated (by a smooth function). The test should have more than 1 d.f. If there is no evidence that the predictor is associated with outcome. Make a formal test of linearity of each remaining continuous predictor. Use restricted cubic spline functions with four knots. If you feel that you can't narrow down the number of candidate predictors without examining the outcomes, and the number is too great to be able to derive a reliable model, use a data reduction technique and combine many of the variables into a summary index.

2. For factors that remain, assess the PH assumption using at least two methods, after ensuring that continuous predictors are transformed to be as linear as possible. In addition, for polytomous predictors, derive log cumulative hazard estimates adjusted for continuous predictors that do not assume anything about the relationship between the polytomous factor and survival.
3. Derive a final Cox PH model. Stratify on polytomous factors that do not satisfy the PH assumption. Decide whether to categorize and stratify on continuous factors that may strongly violate PH. Remember that in this case you can still model the continuous factor to account for any residual regression after adjusting for strata intervals. Include an interaction between two predictors of your choosing. Interpret the parameters in the final model. Also interpret the final model by providing some predicted survival curves in which an important continuous predictor is on the x -axis, predicted survival is on the y -axis, separate curves are drawn for levels of another factor, and any other factors in the model are adjusted to specified constants or to the grand mean. The estimated survival probabilities should be computed at $t = 730$ days.
4. Verify, in an unbiased fashion, your “final” model, for either calibration or discrimination. Validate intermediate steps, not just the final parameter estimates.

Appendix A

Datasets, R Packages, and Internet Resources

Central Web Site and Datasets

The web site for information related to this book is biostat.mc.vanderbilt.edu/rms, and a related web site for a full-semester course based on the book is <http://biostat.mc.vanderbilt.edu/CourseBios330>. The main site contains links to several other web sites and a link to the dataset repository that holds most of the datasets mentioned in the text for downloading. These datasets are in fully annotated R `save` (`.sav` suffixes) files^a; some of these are also available in other formats. The datasets were selected because of the variety of types of response and predictor variables, sample size, and numbers of missing values. In R they may be read using the `load` function, `load(url())` to read directly from the Web, or by using the `Hmisc` package's `getHdata` function to do the same (as is done in code in the case studies). From the web site there are links to other useful dataset sources. Links to presentations and technical reports related to the text are also found on this site, as is information for instructors for obtaining quizzes and answer sheets, extra problems, and solutions to these and to many of the problems in the text. Details about short courses based on the text are also found there. The main site also has Chapter 7 from the first edition, which is a case study in ordinary least squares modeling.

R Packages

The `rms` package written by the author maintains detailed information about a model's design matrix so that many analyses using the model fit are automated. `rms` is a large package of R functions. Most of the functions in `rms` analyze model fits, validate them, or make presentation graphics from them,

^a By convention these should have had `.rda` suffixes.

but the packages also contain special model-fitting functions for binary and ordinal logistic regression (optionally using penalized maximum likelihood), unpenalized ordinal regression with a variety of link functions, penalized and unpenalized least squares, and parametric and semiparametric survival models. In addition, `rms` handles quantile regression and longitudinal analysis using generalized least squares. The `rms` package pays special attention to computing predicted values in that design matrix attributes (e.g., knots for splines, categories for categorical predictors) are “remembered” so that predictors are properly transformed while predictions are being generated. The functions makes extensive use of a wealth of survival analysis software written by Terry Therneau of the Mayo Foundation. This `survival` package is a standard part of R.

The author’s `Hmisc` package contains other miscellaneous functions used in the text. These are functions that do not operate on model fits that used the enhanced design attributes stored by the `rms` package. Functions in `Hmisc` include facilities for data reduction, imputation, power and sample size calculation, advanced table making, recoding variables, translating SAS datasets into R data frames while preserving all data attributes (including variable and value labels and special missing values), drawing and annotating plots, and converting certain R objects to \LaTeX ³⁷¹ typeset form. The latter capability, provided by a family of `latex` functions, completes the conversion to \LaTeX of many of the objects created by `rms`. The packages contain several \LaTeX methods that create \LaTeX code for typesetting model fits in algebraic notation, for printing ANOVA and regression effect (e.g., odds ratio) tables, and other applications. The \LaTeX methods were used extensively in the text, especially for writing restricted cubic spline function fits in simplest notation.

The latest version of the `rms` package is available from CRAN (see below). It is necessary to install the `Hmisc` package in order to use `rms` package. The Web site also contains more in-depth overviews of the packages, which run on UNIX, Linux, Mac, and Microsoft Windows systems. The packages may be automatically downloaded and installed using R’s `install.packages` function or using menus under R graphical user interfaces.

R-help, CRAN, and Discussion Boards

To subscribe to the highly informative and helpful `R-help` e-mail group, see the Web site. `R-help` is appropriate for asking general questions about R including those about finding or writing functions to do specific analyses (for questions specific to a package, contact the author of that package). Another resource is the CRAN repository at www.r-project.org. Another excellent resource for asking questions about R is stackoverflow.com/questions/tagged/r. There is a Google group `regmod` devoted to the book and courses.

Multiple Imputation

The `Impute` E-mail list maintained by Juned Siddique of Northwestern University is an invaluable source of information regarding missing data problems. To subscribe to this list, see the Web site. Other excellent sources of on-line information are Joseph Schafer's "Multiple Imputation Frequently Asked Questions" site and Stef van Buuren and Karin Oudshoorn's "Multiple Imputation Online" site, for which links exist on the main Web site.

Bibliography

An extensive annotated bibliography containing all the references in this text as well as other references concerning predictive methods, survival analysis, logistic regression, prognosis, diagnosis, modeling strategies, model validation, practical Bayesian methods, clinical trials, graphical methods, papers for teaching statistical methods, the bootstrap, and many other areas may be found at <http://www.citeulike.org/user/harrelfe>.

SAS

SAS macros for fitting restricted cubic splines and for other basic operations are freely available from the main Web site. The Web site also has notes on SAS usage for some of the methods presented in the text.

References

Numbers following \diamond are page numbers of citations.

1. O. O. Aalen. Nonparametric inference in connection with multiple decrement models. *Scan J Stat*, 3:15–27, 1976. \diamond 413
2. O. O. Aalen. Further results on the non-parametric linear regression model in survival analysis. *Stat Med*, 12:1569–1588, 1993. \diamond 518
3. O. O. Aalen, E. Bjertness, and T. Sønju. Analysis of dependent survival data applied to lifetimes of amalgam fillings. *Stat Med*, 14:1819–1829, 1995. \diamond 421
4. M. Abrahamowicz, T. MacKenzie, and J. M. Esdaile. Time-dependent hazard ratio: Modeling and hypothesis testing with applications in lupus nephritis. *JAMA*, 91:1432–1439, 1996. \diamond 501
5. A. Agresti. A survey of models for repeated ordered categorical response data. *Stat Med*, 8:1209–1224, 1989. \diamond 324
6. A. Agresti. *Categorical data analysis*. Wiley, Hoboken, NJ, second edition, 2002. \diamond 271
7. H. Ahn and W. Loh. Tree-structured proportional hazards regression modeling. *Biometrics*, 50:471–485, 1994. \diamond 41, 178
8. J. Aitchison and S. D. Silvey. The generalization of probit analysis to the case of multiple responses. *Biometrika*, 44:131–140, 1957. \diamond 324
9. K. Akazawa, T. Nakamura, and Y. Palesch. Power of logrank test and Cox regression model in clinical trials with heterogeneous samples. *Stat Med*, 16:583–597, 1997. \diamond 4
10. O. O. Al-Radi, F. E. Harrell, C. A. Caldarone, B. W. McCrindle, J. P. Jacobs, M. G. Williams, G. S. Van Arsdell, and W. G. Williams. Case complexity scores in congenital heart surgery: A comparative study of the Aristotal Basic Complexity score and the Risk Adjustment in Congenital Heart Surg (RACHS-1) system. *J Thorac Cardiovasc Surg*, 133:865–874, 2007. \diamond 215
11. J. M. Alho. On the computation of likelihood ratio and score test based confidence intervals in generalized linear models. *Stat Med*, 11:923–930, 1992. \diamond 214
12. P. D. Allison. *Missing Data*. Sage University Papers Series on Quantitative Applications in the Social Sciences, 07-136. Sage, Thousand Oaks CA, 2001. \diamond 49, 58

13. D. G. Altman. Categorising continuous covariates (letter to the editor). *Brit J Cancer*, 64:975, 1991. ◊11, 19
14. D. G. Altman. Suboptimal analysis using ‘optimal’ cutpoints. *Brit J Cancer*, 78:556–557, 1998. ◊19
15. D. G. Altman and P. K. Andersen. A note on the uncertainty of a survival probability estimated from Cox’s regression model. *Biometrika*, 73:722–724, 1986. ◊11, 517
16. D. G. Altman and P. K. Andersen. Bootstrap investigation of the stability of a Cox regression model. *Stat Med*, 8:771–783, 1989. ◊68, 70, 341
17. D. G. Altman, B. Lausen, W. Sauerbrei, and M. Schumacher. Dangers of using ‘optimal’ cutpoints in the evaluation of prognostic factors. *J Nat Cancer Inst*, 86:829–835, 1994. ◊11, 19, 20
18. D. G. Altman and P. Royston. What do we mean by validating a prognostic model? *Stat Med*, 19:453–473, 2000. ◊6, 122, 519
19. B. Altschuler. Theory for the measurement of competing risks in animal experiments. *Math Biosci*, 6:1–11, 1970. ◊413
20. C. F. Alzola and F. E. Harrell. An Introduction to S and the Hmisc and Design Libraries, 2006. Electronic book, 310 pages. ◊129
21. G. Ambler, A. R. Brady, and P. Royston. Simplifying a prognostic model: a simulation study based on clinical data. *Stat Med*, 21(24):3803–3822, Dec. 2002. ◊121
22. F. Ambrogi, E. Biganzoli, and P. Boracchi. Estimates of clinically useful measures in competing risks survival analysis. *Stat Med*, 27:6407–6425, 2008. ◊421
23. P. K. Andersen and R. D. Gill. Cox’s regression model for counting processes: A large sample study. *Ann Stat*, 10:1100–1120, 1982. ◊418, 513
24. G. L. Anderson and T. R. Fleming. Model misspecification in proportional hazards regression. *Biometrika*, 82:527–541, 1995. ◊4
25. J. A. Anderson. Regression and ordered categorical variables. *J Roy Stat Soc B*, 46:1–30, 1984. ◊324
26. J. A. Anderson and P. R. Philips. Regression, discrimination and measurement models for ordered categorical variables. *Appl Stat*, 30:22–31, 1981. ◊324
27. J. A. Anderson and A. Senthilselvan. A two-step regression model for hazard functions. *Appl Stat*, 31:44–51, 1982. ◊495, 499, 501
28. D. F. Andrews and A. M. Herzberg. *Data*. Springer-Verlag, New York, 1985. ◊161
29. E. Arjas. A graphical method for assessing goodness of fit in Cox’s proportional hazards model. *J Am Stat Assoc*, 83:204–212, 1988. ◊420, 495, 502
30. H. R. Arkes, N. V. Dawson, T. Speroff, F. E. Harrell, C. Alzola, R. Phillips, N. Desbiens, R. K. Oye, W. Knaus, A. F. Connors, and T. Investigators. The covariance decomposition of the probability score and its use in evaluating prognostic estimates. *Med Decis Mak*, 15:120–131, 1995. ◊257
31. B. G. Armstrong and M. Sloan. Ordinal regression models for epidemiologic data. *Am J Epi*, 129:191–204, 1989. See letter to editor by Peterson. ◊319, 320, 321, 324
32. D. Ashby, C. R. West, and D. Ames. The ordered logistic regression model in psychiatry: Rising prevalence of dementia in old people’s homes. *Stat Med*, 8:1317–1326, 1989. ◊324
33. A. C. Atkinson. A note on the generalized information criterion for choice of a model. *Biometrika*, 67:413–418, 1980. ◊69, 204
34. P. C. Austin. A comparison of regression trees, logistic regression, generalized additive models, and multivariate adaptive regression splines for predicting AMI mortality. *Stat Med*, 26:2937–2957, 2007. ◊41

35. P. C. Austin. Bootstrap model selection had similar performance for selecting authentic and noise variables compared to backward variable elimination: a simulation study. *J Clin Epi*, 61:1009–1017, 2008. ◊70
36. P. C. Austin and E. W. Steyerberg. Events per variable (EPV) and the relative performance of different strategies for estimating the out-of-sample validity of logistic regression models. *Statistical methods in medical research*, Nov. 2014. ◊112
37. P. C. Austin and E. W. Steyerberg. Graphical assessment of internal and external calibration of logistic regression models by using loess smoothers. *Stat Med*, 33(3):517–535, Feb. 2014. ◊105
38. P. C. Austin, J. V. Tu, P. A. Daly, and D. A. Alter. Tutorial in Biostatistics: The use of quantile regression in health care research: a case study examining gender differences in the timeliness of thrombolytic therapy. *Stat Med*, 24:791–816, 2005. ◊392
39. D. Bamber. The area above the ordinal dominance graph and the area below the receiver operating characteristic graph. *J Mathe Psych*, 12:387–415, 1975. ◊257
40. J. Banks. Nomograms. In S. Kotz and N. L. Johnson, editors, *Encyclopedia of Stat Scis*, volume 6. Wiley, New York, 1985. ◊104, 267
41. J. Barnard and D. B. Rubin. Small-sample degrees of freedom with multiple imputation. *Biometrika*, 86:948–955, 1999. ◊58
42. S. A. Barnes, S. R. Lindborg, and J. W. Seaman. Multiple imputation techniques in small sample clinical trials. *Stat Med*, 25:233–245, 2006. ◊47, 58
43. F. Barzi and M. Woodward. Imputations of missing values in practice: Results from imputations of serum cholesterol in 28 cohort studies. *Am J Epi*, 160:34–45, 2004. ◊50, 58
44. R. A. Becker, J. M. Chambers, and A. R. Wilks. *The New S Language*. Wadsworth and Brooks/Cole, Pacific Grove, CA, 1988. ◊127
45. H. Belcher. The concept of residual confounding in regression models and some applications. *Stat Med*, 11:1747–1758, 1992. ◊11, 19
46. D. A. Belsley. *Conditioning Diagnostics: Collinearity and Weak Data in Regression*. Wiley, New York, 1991. ◊101
47. D. A. Belsley, E. Kuh, and R. E. Welsch. *Regression Diagnostics: Identifying Influential Data and Sources of Collinearity*. Wiley, New York, 1980. ◊91
48. R. Bender and A. Benner. Calculating ordinal regression models in SAS and S-Plus. *Biometrical J*, 42:677–699, 2000. ◊324
49. J. K. Benedetti, P. Liu, H. N. Sather, J. Seinfeld, and M. A. Epton. Effective sample size for tests of censored survival data. *Biometrika*, 69:343–349, 1982. ◊73
50. K. Berhane, M. Hauptmann, and B. Langholz. Using tensor product splines in modeling exposure–time–response relationships: Application to the Colorado Plateau Uranium Miners cohort. *Stat Med*, 27:5484–5496, 2008. ◊37
51. K. N. Berk and D. E. Booth. Seeing a curve in multiple regression. *Technometrics*, 37:385–398, 1995. ◊272
52. D. M. Berridge and J. Whitehead. Analysis of failure time data with ordinal categories of response. *Stat Med*, 10:1703–1710, 1991. ◊319, 320, 324, 417
53. C. Berzuini and D. Clayton. Bayesian analysis of survival on multiple time scales. *Stat Med*, 13:823–838, 1994. ◊401
54. W. B. Bilker and M. Wang. A semiparametric extension of the Mann-Whitney test for randomly truncated data. *Biometrics*, 52:10–20, 1996. ◊420
55. D. A. Binder. Fitting Cox’s proportional hazards models from survey data. *Biometrika*, 79:139–147, 1992. ◊213, 215
56. C. Binquet, M. Abrahamowicz, A. Mahboubi, V. Jooste, J. Faivre, C. Bonithon-Kopp, and C. Quantin. Empirical study of the dependence of the results of multivariable flexible survival analyses on model selection strategy. *Stat Med*, 27:6470–6488, 2008. ◊420

57. E. H. Blackstone. Analysis of death (survival analysis) and other time-related events. In F. J. Macartney, editor, *Current Status of Clinical Cardiology*, pages 55–101. MTP Press Limited, Lancaster, UK, 1986. ◊420
58. S. E. Bleeker, H. A. Moll, E. W. Steyerberg, A. R. T. Donders, G. Derkson-Lubsen, D. E. Grobbee, and K. G. M. Moons. External validation is necessary in prediction research: A clinical example. *J Clin Epi*, 56:826–832, 2003. ◊122
59. M. Blettner and W. Sauerbrei. Influence of model-building strategies on the results of a case-control study. *Stat Med*, 12:1325–1338, 1993. ◊123
60. D. D. Boos. On generalized score tests. *Ann Math Stat*, 46:327–333, 1992. ◊123
61. J. G. Booth and S. Sarkar. Monte Carlo approximation of bootstrap variances. *Am Statistician*, 52:354–357, 1998. ◊122
62. R. Bordley. Statistical decisionmaking without math. *Chance*, 20(3):39–44, 2007. ◊5
63. R. Brant. Assessing proportionality in the proportional odds model for ordinal logistic regression. *Biometrics*, 46:1171–1178, 1990. ◊324
64. S. R. Brazer, F. S. Pancotto, T. T. Long III, F. E. Harrell, K. L. Lee, M. P. Tyor, and D. B. Pryor. Using ordinal logistic regression to estimate the likelihood of colorectal neoplasia. *J Clin Epi*, 44:1263–1270, 1991. ◊324
65. A. R. Brazzale and A. C. Davison. Accurate parametric inference for small samples. *Statistical Sci*, 23(4):465–484, 2008. ◊214
66. L. Breiman. The little bootstrap and other methods for dimensionality selection in regression: X-fixed prediction error. *J Am Stat Assoc*, 87:738–754, 1992. ◊69, 100, 112, 114, 123, 204
67. L. Breiman. Statistical modeling: The two cultures (with discussion). *Statistical Sci*, 16:199–231, 2001. ◊11
68. L. Breiman and J. H. Friedman. Estimating optimal transformations for multiple regression and correlation (with discussion). *J Am Stat Assoc*, 80:580–619, 1985. ◊82, 176, 390
69. L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone. *Classification and Regression Trees*. Wadsworth and Brooks/Cole, Pacific Grove, CA, 1984. ◊30, 41, 142
70. N. E. Breslow. Covariance analysis of censored survival data. *Biometrics*, 30:89–99, 1974. ◊477, 483, 485
71. N. E. Breslow, N. E. Day, K. T. Halvorsen, R. L. Prentice, and C. Sabai. Estimation of multiple relative risk functions in matched case-control studies. *Am J Epi*, 108:299–307, 1978. ◊483
72. N. E. Breslow, L. Edler, and J. Berger. A two-sample censored-data rank test for acceleration. *Biometrics*, 40:1049–1062, 1984. ◊501
73. G. W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Rev*, 78:1–3, 1950. ◊257
74. W. M. Briggs and R. Zaretzki. The skill plot: A graphical technique for evaluating continuous diagnostic tests (with discussion). *Biometrics*, 64:250–261, 2008. ◊5, 11
75. G. Bron. The loss of the “Titanic”. *The Sphere*, 49:103, May 1912. The results analysed and shown in a special “Sphere” diagram drawn from the official figures given in the House of Commons. ◊291
76. B. W. Brown, M. Hollander, and R. M. Korwar. Nonparametric tests of independence for censored data, with applications to heart transplant studies. In F. Proschan and R. J. Serfling, editors, *Reliability and Biometry*, pages 327–354. SIAM, Philadelphia, 1974. ◊505
77. D. Brownstone. Regression strategies. In *Proceedings of the 20th Symposium on the Interface between Computer Science and Statistics*, pages 74–79, Washington, DC, 1988. American Statistical Association. ◊116
78. J. Bryant and J. J. Dignam. Semiparametric models for cumulative incidence functions. *Biometrics*, 69:182–190, 2004. ◊420

79. S. F. Buck. A method of estimation of missing values in multivariate data suitable for use with an electronic computer. *J Roy Stat Soc B*, 22:302–307, 1960. ◊52
80. S. T. Buckland, K. P. Burnham, and N. H. Augustin. Model selection: An integral part of inference. *Biometrics*, 53:603–618, 1997. ◊10, 11, 214
81. J. Buckley and I. James. Linear regression with censored data. *Biometrika*, 66:429–36, 1979. ◊447
82. P. Buettner, C. Garbe, and I. Guggenmoos-Holzmann. Problems in defining cutoff points of continuous prognostic factors: Example of tumor thickness in primary cutaneous melanoma. *J Clin Epi*, 50:1201–1210, 1997. ◊11, 19
83. K. Bull and D. Spiegelhalter. Survival analysis in observational studies. *Stat Med*, 16:1041–1074, 1997. ◊399, 401, 420
84. K. P. Burnham and D. R. Anderson. *Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach*. Springer, 2nd edition, Dec. 2003. ◊69
85. S. Buuren. *Flexible imputation of missing data*. Chapman & Hall/CRC, Boca Raton, FL, 2012. ◊54, 55, 58, 304
86. M. Buyse. R^2 : A useful measure of model performance when predicting a dichotomous outcome. *Stat Med*, 19:271–274, 2000. Letter to the Editor regarding *Stat Med* 18:375–384; 1999. ◊272
87. D. P. Byar and S. B. Green. The choice of treatment for cancer patients based on covariate information: Application to prostate cancer. *Bulletin Cancer, Paris*, 67:477–488, 1980. ◊161, 275, 521
88. R. M. Califf, F. E. Harrell, K. L. Lee, J. S. Rankin, and Others. The evolution of medical and surgical therapy for coronary artery disease. *JAMA*, 261:2077–2086, 1989. ◊484, 485, 510
89. R. M. Califf, H. R. Phillips, and Others. Prognostic value of a coronary artery jeopardy score. *J Am College Cardiol*, 5:1055–1063, 1985. ◊207
90. R. M. Califf, L. H. Woodlief, F. E. Harrell, K. L. Lee, H. D. White, A. Guerci, G. I. Barbash, R. Simes, W. Weaver, M. L. Simoons, E. J. Topol, and T. Investigators. Selection of thrombolytic therapy for individual patients: Development of a clinical model. *Am Heart J*, 133:630–639, 1997. ◊4
91. A. J. Canty, A. C. Davison, D. V. Hinkley, and V. Venture. Bootstrap diagnostics and remedies. *Can J Stat*, 34:5–27, 2006. ◊122
92. J. Carpenter and J. Bithell. Bootstrap confidence intervals: when, which, what? A practical guide for medical statisticians. *Stat Med*, 19:1141–1164, 2000. ◊122, 214
93. W. H. Carter, G. L. Wampler, and D. M. Stablein. *Regression Analysis of Survival Data in Cancer Chemotherapy*. Marcel Dekker, New York, 1983. ◊477
94. Centers for Disease Control and Prevention CDC. National Center for Health Statistics NCHS. National Health and Nutrition Examination Survey, 2010. ◊365
95. M. S. Cepeda, R. Boston, J. T. Farrar, and B. L. Strom. Comparison of logistic regression versus propensity score when the number of events is low and there are multiple confounders. *Am J Epi*, 158:280–287, 2003. ◊272
96. J. M. Chambers and T. J. Hastie, editors. *Statistical Models in S*. Wadsworth and Brooks/Cole, Pacific Grove, CA, 1992. ◊x, 29, 41, 128, 142, 245, 269, 493, 498
97. L. E. Chambless and K. E. Boyle. Maximum likelihood methods for complex sample data: Logistic regression and discrete proportional hazards models. *Comm Stat A*, 14:1377–1392, 1985. ◊215
98. R. Chappell. A note on linear rank tests and Gill and Schumacher’s tests of proportionality. *Biometrika*, 79:199–201, 1992. ◊495
99. C. Chatfield. Avoiding statistical pitfalls (with discussion). *Statistical Sci*, 6:240–268, 1991. ◊91

100. C. Chatfield. Model uncertainty, data mining and statistical inference (with discussion). *J Roy Stat Soc A*, 158:419–466, 1995. ◊vii, 9, 10, 11, 68, 100, 123, 204
101. S. Chatterjee and A. S. Hadi. *Regression Analysis by Example*. Wiley, New York, fifth edition, 2012. ◊78, 101
102. S. C. Cheng, J. P. Fine, and L. J. Wei. Prediction of cumulative incidence function under the proportional hazards model. *Biometrics*, 54:219–228, 1998. ◊415
103. S. C. Cheng, L. J. Wei, and Z. Ying. Predicting Survival Probabilities with Semiparametric Transformation Models. *JASA*, 92(437):227–235, Mar. 1997. ◊517
104. F. Chiaromonte, R. D. Cook, and B. Li. Sufficient dimension reduction in regressions with categorical predictors. *Appl Stat*, 30:475–497, 2002. ◊101
105. B. Choodari-Oskooei, P. Royston, and M. K. B. Parmar. A simulation study of predictive ability measures in a survival model II: explained randomness and predictive accuracy. *Stat Med*, 31(23):2644–2659, 2012. ◊518
106. B. Choodari-Oskooei, P. Royston, and M. K. B. Parmar. A simulation study of predictive ability measures in a survival model I: Explained variation measures. *Stat Med*, 31(23):2627–2643, 2012. ◊518
107. A. Ciampi, A. Negassa, and Z. Lou. Tree-structured prediction for censored survival data and the Cox model. *J Clin Epi*, 48:675–689, 1995. ◊41
108. A. Ciampi, J. Thiffault, J. P. Nakache, and B. Asselain. Stratification by stepwise regression, correspondence analysis and recursive partition. *Comp Stat Data Analysis*, 1986:185–204, 1986. ◊41, 81
109. L. A. Clark and D. Pregibon. Tree-Based Models. In J. M. Chambers and T. J. Hastie, editors, *Statistical Models in S*, chapter 9, pages 377–419. Wadsworth and Brooks/Cole, Pacific Grove, CA, 1992. ◊41
110. T. G. Clark and D. G. Altman. Developing a prognostic model in the presence of missing data: an ovarian cancer case study. *J Clin Epi*, 56:28–37, 2003. ◊57
111. W. S. Cleveland. Robust locally weighted regression and smoothing scatterplots. *J Am Stat Assoc*, 74:829–836, 1979. ◊29, 141, 238, 315, 356, 493
112. A. Cnaan and L. Ryan. Survival analysis in natural history studies of disease. *Stat Med*, 8:1255–1268, 1989. ◊401, 420
113. T. J. Cole, C. J. Morley, A. J. Thornton, M. A. Fowler, and P. H. Hewson. A scoring system to quantify illness in babies under 6 months of age. *J Roy Stat Soc A*, 154:287–304, 1991. ◊324
114. D. Collett. *Modelling Survival Data in Medical Research*. Chapman and Hall, London, 1994. ◊420, 517
115. D. Collett. *Modelling Binary Data*. Chapman and Hall, London, second edition, 2002. ◊213, 272, 315
116. A. F. Connors, T. Speroff, N. V. Dawson, C. Thomas, F. E. Harrell, D. Wagner, N. Desbiens, L. Goldman, A. W. Wu, R. M. Califf, W. J. Fulkerson, H. Vidaillet, S. Broste, P. Bellamy, J. Lynn, W. A. Knaus, and T. S. Investigators. The effectiveness of right heart catheterization in the initial care of critically ill patients. *JAMA*, 276:889–897, 1996. ◊3
117. E. F. Cook and L. Goldman. Asymmetric stratification: An outline for an efficient method for controlling confounding in cohort studies. *Am J Epi*, 127:626–639, 1988. ◊31, 231
118. N. R. Cook. Use and misuses of the receiver operating characteristic curve in risk prediction. *Circulation*, 115:928–935, 2007. ◊93, 101, 273
119. R. D. Cook. Fisher Lecture: Dimension reduction in regression. *Statistical Sci*, 22:1–26, 2007. ◊101
120. R. D. Cook and L. Forzani. Principal fitted components for dimension reduction in regression. *Statistical Sci*, 23(4):485–501, 2008. ◊101

121. J. Copas. The effectiveness of risk scores: The logit rank plot. *Appl Stat*, 48:165–183, 1999. ◊273
122. J. B. Copas. Regression, prediction and shrinkage (with discussion). *J Roy Stat Soc B*, 45:311–354, 1983. ◊100, 101
123. J. B. Copas. Cross-validation shrinkage of regression predictors. *J Roy Stat Soc B*, 49:175–183, 1987. ◊115, 123, 273, 508
124. J. B. Copas. Unweighted sum of squares tests for proportions. *Appl Stat*, 38:71–80, 1989. ◊236
125. J. B. Copas and T. Long. Estimating the residual variance in orthogonal regression with variable selection. *The Statistician*, 40:51–59, 1991. ◊68
126. C. Cox. Location-scale cumulative odds models for ordinal data: A generalized non-linear model approach. *Stat Med*, 14:1191–1203, 1995. ◊324
127. C. Cox. The generalized f distribution: An umbrella for parametric survival analysis. *Stat Med*, 27:4301–4313, 2008. ◊424
128. C. Cox, H. Chu, M. F. Schneider, and A. Muñoz. Parametric survival analysis and taxonomy of hazard functions for the generalized gamma distribution. *Stat Med*, 26:4352–4374, 2007. ◊424
129. D. R. Cox. The regression analysis of binary sequences (with discussion). *J Roy Stat Soc B*, 20:215–242, 1958. ◊14, 220
130. D. R. Cox. Two further applications of a model for binary regression. *Biometrika*, 45(3/4):562–565, 1958. ◊259
131. D. R. Cox. Further results on tests of separate families of hypotheses. *J Roy Stat Soc B*, 24:406–424, 1962. ◊205
132. D. R. Cox. Regression models and life-tables (with discussion). *J Roy Stat Soc B*, 34:187–220, 1972. ◊39, 41, 172, 207, 213, 314, 418, 428, 475, 476
133. D. R. Cox and D. Oakes. *Analysis of Survival Data*. Chapman and Hall, London, 1984. ◊401, 420, 517
134. D. R. Cox and E. J. Snell. A general definition of residuals (with discussion). *J Roy Stat Soc B*, 30:248–275, 1968. ◊440
135. D. R. Cox and E. J. Snell. *The Analysis of Binary Data*. Chapman and Hall, London, second edition, 1989. ◊206
136. D. R. Cox and N. Wermuth. A comment on the coefficient of determination for binary responses. *Am Statistician*, 46:1–4, 1992. ◊206, 256
137. J. G. Cragg and R. Uhler. The demand for automobiles. *Canadian Journal of Economics*, 3:386–406, 1970. ◊206, 256
138. S. L. Crawford, S. L. Tennstedt, and J. B. McKinlay. A comparison of analytic methods for non-random missingness of outcome data. *J Clin Epi*, 48:209–219, 1995. ◊58
139. N. J. Crichton and J. P. Hinde. Correspondence analysis as a screening method for indicants for clinical diagnosis. *Stat Med*, 8:1351–1362, 1989. ◊81
140. N. J. Crichton, J. P. Hinde, and J. Marchini. Models for diagnosing chest pain: Is CART useful? *Stat Med*, 16:717–727, 1997. ◊41
141. L. A. Cupples, D. R. Gagnon, R. Ramaswamy, and R. B. D’Agostino. Age-adjusted survival curves with application in the Framingham Study. *Stat Med*, 14:1731–1744, 1995. ◊517
142. E. E. Cureton and R. B. D’Agostino. *Factor Analysis, An Applied Approach*. Erlbaum, Hillsdale, NJ, 1983. ◊81, 87, 101
143. D. M. Dabrowska, K. A. Doksum, N. J. Feduska, R. Husing, and P. Neville. Methods for comparing cumulative hazard functions in a semi-proportional hazard model. *Stat Med*, 11:1465–1476, 1992. ◊482, 495, 502
144. R. B. D’Agostino, A. J. Belanger, E. W. Markson, M. Kelly-Hayes, and P. A. Wolf. Development of health risk appraisal functions in the presence of multiple indicators: The Framingham Study nursing home institutionalization model. *Stat Med*, 14:1757–1770, 1995. ◊81, 101

145. R. B. D'Agostino, M. L. Lee, A. J. Belanger, and L. A. Cupples. Relation of pooled logistic regression to time dependent Cox regression analysis: The Framingham Heart Study. *Stat Med*, 9:1501–1515, 1990. ◊447
146. D'Agostino, Jr and D. B. Rubin. Estimating and using propensity scores with partially missing data. *J Am Stat Assoc*, 95:749–759, 2000. ◊58
147. C. E. Davis, J. E. Hyde, S. I. Bangdiwala, and J. J. Nelson. An example of dependencies among variables in a conditional logistic regression. In S. H. Moolgavkar and R. L. Prentice, editors, *Modern Statistical Methods in Chronic Disease Epi*, pages 140–147. Wiley, New York, 1986. ◊79, 138, 255
148. C. S. Davis. *Statistical Methods for the Analysis of Repeated Measurements*. Springer, New York, 2002. ◊143, 149
149. R. B. Davis and J. R. Anderson. Exponential survival trees. *Stat Med*, 8:947–961, 1989. ◊41
150. A. C. Davison and D. V. Hinkley. *Bootstrap Methods and Their Application*. Cambridge University Press, Cambridge, 1997. ◊70, 106, 109, 122
151. R. J. M. Dawson. The 'Unusual Episode' data revisited. *J Stat Edu*, 3(3), 1995. Online journal at www.amstat.org/publications/jse/v3n3/datasets.-dawson.html. ◊291
152. C. de Boor. *A Practical Guide to Splines*. Springer-Verlag, New York, revised edition, 2001. ◊23, 40
153. J. de Leeuw and P. Mair. Gifi methods for optimal scaling in r: The package `homals`. *J Stat Software*, 31(4):1–21, Aug. 2009. ◊101
154. E. R. DeLong, C. L. Nelson, J. B. Wong, D. B. Pryor, E. D. Peterson, K. L. Lee, D. B. Mark, R. M. Califf, and S. G. Pauker. Using observational data to estimate prognosis: an example using a coronary artery disease registry. *Stat Med*, 20:2505–2532, 2001. ◊420
155. S. Derksen and H. J. Keselman. Backward, forward and stepwise automated subset selection algorithms: Frequency of obtaining authentic and noise variables. *British J Math Stat Psych*, 45:265–282, 1992. ◊68
156. T. F. Devlin and B. J. Weeks. Spline functions for logistic regression modeling. In *Proceedings of the Eleventh Annual SAS Users Group International Conference*, pages 646–651, Cary, NC, 1986. SAS Institute, Inc. ◊21, 24
157. T. DiCiccio and B. Efron. More accurate confidence intervals in exponential families. *Biometrika*, 79:231–245, 1992. ◊214
158. E. R. Dickson, P. M. Grambsch, T. R. Fleming, L. D. Fisher, and A. Langworthy. Prognosis in primary biliary cirrhosis: Model for decision making. *Hepatology*, 10:1–7, 1989. ◊178
159. P. J. Diggle, P. Heagerty, K.-Y. Liang, and S. L. Zeger. *Analysis of Longitudinal Data*. Oxford University Press, Oxford UK, second edition, 2002. ◊143, 147
160. N. Doganaksoy and J. Schmee. Comparisons of approximate confidence intervals for distributions used in life-data analysis. *Technometrics*, 35:175–184, 1993. ◊198, 214
161. Donders, G. J. M. G. van der Heijden, T. Stijnen, and K. G. M. Moons. Review: A gentle introduction to imputation of missing values. *J Clin Epi*, 59:1087–1091, 2006. ◊49, 58
162. A. Donner. The relative effectiveness of procedures commonly used in multiple regression analysis for dealing with missing values. *Am Statistician*, 36:378–381, 1982. ◊48, 52
163. D. Draper. Assessment and propagation of model uncertainty (with discussion). *J Roy Stat Soc B*, 57:45–97, 1995. ◊10, 11
164. M. Drum and P. McCullagh. Comment on regression models for discrete longitudinal responses by G. M. Fitzmaurice, N. M. Laird, and A. G. Rotnitzky. *Stat Sci*, 8:300–301, 1993. ◊197
165. N. Duan. Smearing estimate: A nonparametric retransformation method. *J Am Stat Assoc*, 78:605–610, 1983. ◊392

166. J. A. Dubin, H. Müller, and J. Wang. Event history graphs for censored data. *Stat Med*, 20:2951–2964, 2001. ◊418, 420
167. R. Dudley, F. E. Harrell, L. Smith, D. B. Mark, R. M. Califf, D. B. Pryor, D. Glower, J. Lipscomb, and M. Hlatky. Comparison of analytic models for estimating the effect of clinical factors on the cost of coronary artery bypass graft surgery. *J Clin Epi*, 46:261–271, 1993. ◊x
168. S. Durrleman and R. Simon. Flexible regression models with cubic splines. *Stat Med*, 8:551–561, 1989. ◊40
169. J. P. Eaton and C. A. Haas. *Titanic: Triumph and Tragedy*. W. W. Norton, New York, second edition, 1995. ◊291
170. B. Efron. The two sample problem with censored data. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, volume 4, pages 831–853. 1967. ◊505
171. B. Efron. The efficiency of Cox’s likelihood function for censored data. *J Am Stat Assoc*, 72:557–565, 1977. ◊475, 477
172. B. Efron. Estimating the error rate of a prediction rule: Improvement on cross-validation. *J Am Stat Assoc*, 78:316–331, 1983. ◊70, 113, 114, 115, 116, 123, 259
173. B. Efron. How biased is the apparent error rate of a prediction rule? *J Am Stat Assoc*, 81:461–470, 1986. ◊101, 114
174. B. Efron. Missing data, imputation, and the bootstrap (with discussion). *J Am Stat Assoc*, 89:463–479, 1994. ◊52, 54
175. B. Efron and G. Gong. A leisurely look at the bootstrap, the jackknife, and cross-validation. *Am Statistician*, 37:36–48, 1983. ◊114
176. B. Efron and C. Morris. Stein’s paradox in statistics. *Sci Am*, 236(5):119–127, 1977. ◊77
177. B. Efron and R. Tibshirani. Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy. *Statistical Sci*, 1:54–77, 1986. ◊70, 106, 114, 197
178. B. Efron and R. Tibshirani. *An Introduction to the Bootstrap*. Chapman and Hall, New York, 1993. ◊70, 106, 114, 115, 122, 197, 199
179. B. Efron and R. Tibshirani. Improvements on cross-validation: The .632+ bootstrap method. *J Am Stat Assoc*, 92:548–560, 1997. ◊123, 124
180. G. E. Eide, E. Omenaas, and A. Gulsvik. The semi-proportional hazards model revisited: Practical reparameterizations. *Stat Med*, 15:1771–1777, 1996. ◊482
181. C. Faes, G. Molenberghs, M. Aerts, G. Verbeke, and M. G. Kenward. The effective sample size and an alternative small-sample degrees-of-freedom method. *Am Statistician*, 63(4):389–399, 2009. ◊148
182. M. W. Fagerland and D. W. Hosmer. A goodness-of-fit test for the proportional odds regression model. *Stat Med*, 32(13):2235–2249, 2013. ◊317
183. J. Fan and R. A. Levine. To amnio or not to amnio: That is the decision for Bayes. *Chance*, 20(3):26–32, 2007. ◊5
184. D. Faraggi, M. LeBlanc, and J. Crowley. Understanding neural networks using regression trees: an application to multiple myeloma survival data. *Stat Med*, 20:2965–2976, 2001. ◊120
185. D. Faraggi and R. Simon. A simulation study of cross-validation for selecting an optimal cutpoint in univariate survival analysis. *Stat Med*, 15:2203–2213, 1996. ◊11, 19
186. J. J. Faraway. The cost of data analysis. *J Comp Graph Stat*, 1:213–229, 1992. ◊10, 11, 97, 100, 115, 116, 322, 393, 396
187. V. Fedorov, F. Mannino, and R. Zhang. Consequences of dichotomization. *Pharm Stat*, 8:50–61, 2009. ◊5, 19
188. Z. Feng, D. McLerran, and J. Grizzle. A comparison of statistical methods for clustered data analysis with Gaussian error. *Stat Med*, 15:1793–1806, 1996. ◊197, 213

189. L. Ferré. Determining the dimension in sliced inverse regression and related methods. *J Am Stat Assoc*, 93:132–149, 1998. ◊101
190. S. E. Fienberg. *The Analysis of Cross-Classified Categorical Data*. Springer, New York, second edition, 2007. ◊311, 319
191. P. Filzmoser, H. Fritz, and K. Kalcher. *pcaPP: Robust PCA by Projection Pursuit*, 2012. R package version 1.9–48. ◊175
192. J. P. Fine and R. J. Gray. A proportional hazards model for the subdistribution of a competing risk. *J Am Stat Assoc*, 94:496–509, 1999. ◊420
193. D. M. Finkelstein and D. A. Schoenfeld. Combining mortality and longitudinal measures in clinical trials. *Stat Med*, 18:1341–1354, 1999. ◊420
194. M. Fiocco, H. Putter, and H. C. van Houwelingen. Reduced-rank proportional hazards regression and simulation-based prediction for multi-state models. *Stat Med*, 27:4340–4358, 2008. ◊420
195. G. M. Fitzmaurice. A caveat concerning independence estimating equations with multivariate binary data. *Biometrics*, 51:309–317, 1995. ◊214
196. T. R. Fleming and D. P. Harrington. Nonparametric estimation of the survival distribution in censored data. *Comm Stat Th Meth*, 13(20):2469–2486, 1984. ◊413
197. T. R. Fleming and D. P. Harrington. *Counting Processes & Survival Analysis*. Wiley, New York, 1991. ◊178, 420
198. I. Ford, J. Norrie, and S. Ahmadi. Model inconsistency, illustrated by the Cox proportional hazards model. *Stat Med*, 14:735–746, 1995. ◊4
199. E. B. Fowlkes. Some diagnostics for binary logistic regression via smoothing. *Biometrika*, 74:503–515, 1987. ◊272
200. J. Fox. *Applied Regression Analysis, Linear Models, and Related Methods*. SAGE Publications, Thousand Oaks, CA, 1997. ◊viii
201. J. Fox. *An R and S-PLUS Companion to Applied Regression*. SAGE Publications, Thousand Oaks, CA, 2002. ◊viii
202. J. Fox. *Applied Regression Analysis and Generalized Linear Models*. SAGE Publications, Thousand Oaks, CA, second edition, 2008. ◊121
203. Fox, John. Bootstrapping Regression Models: An Appendix to An R and S-PLUS Companion to Applied Regression, 2002. ◊202
204. B. Francis and M. Fuller. Visualization of event histories. *J Roy Stat Soc A*, 159:301–308, 1996. ◊421
205. D. Freedman, W. Navidi, and S. Peters. *On the Impact of Variable Selection in Fitting Regression Equations*, pages 1–16. Lecture Notes in Economics and Mathematical Systems. Springer-Verlag, New York, 1988. ◊115
206. D. A. Freedman. On the so-called “Huber sandwich estimator” and “robust standard errors”. *Am Statistician*, 60:299–302, 2006. ◊213
207. J. H. Friedman. A variable span smoother. Technical Report 5, Laboratory for Computational Statistics, Department of Statistics, Stanford University, 1984. ◊29, 82, 141, 210, 273, 498
208. L. Friedman and M. Wall. Graphical views of suppression and multicollinearity in multiple linear regression. *Am Statistician*, 59:127–136, 2005. ◊101
209. M. H. Gail. Does cardiac transplantation prolong life? A reassessment. *Ann Int Med*, 76:815–817, 1972. ◊401
210. M. H. Gail and R. M. Pfeiffer. On criteria for evaluating models of absolute risk. *Biostatistics*, 6(2):227–239, 2005. ◊5
211. J. C. Gardiner, Z. Luo, and L. A. Roman. Fixed effects, random effects and GEE: What are the differences? *Stat Med*, 28:221–239, 2009. ◊160
212. J. J. Gaynor, E. J. Feuer, C. C. Tan, D. H. Wu, C. R. Little, D. J. Straus, D. D. Clarkson, and M. F. Brennan. On the use of cause-specific failure and conditional failure probabilities: Examples from clinical oncology data. *J Am Stat Assoc*, 88:400–409, 1993. ◊414, 415

213. A. Gelman. Scaling regression inputs by dividing by two standard deviations. *Stat Med*, 27:2865–2873, 2008. ◊121
214. R. B. Geskus. Cause-specific cumulative incidence estimation and the Fine and Gray model under both left truncation and right censoring. *Biometrics*, 67(1):39–49, 2011. ◊420
215. A. Giannoni, R. Baruah, T. Leong, M. B. Rehman, L. E. Pastormerlo, F. E. Harrell, A. J. Coats, and D. P. Francis. Do optimal prognostic thresholds in continuous physiological variables really exist? Analysis of origin of apparent thresholds, with systematic review for peak oxygen consumption, ejection fraction and BNP. *PLoS ONE*, 9(1), 2014. ◊19, 20
216. J. H. Giudice, J. R. Fieberg, and M. S. Lenarz. Spending degrees of freedom in a poor economy: A case study of building a sightability model for moose in northeastern minnesota. *J Wildlife Manage*, 2011. ◊100
217. S. A. Glantz and B. K. Slinker. *Primer of Applied Regression and Analysis of Variance*. McGraw-Hill, New York, 1990. ◊78
218. M. Glasser. Exponential survival with covariance. *J Am Stat Assoc*, 62:561–568, 1967. ◊431
219. T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. *J Am Stat Assoc*, 102:359–378, 2007. ◊4, 5, 273
220. A. I. Goldman. EVENTCHARTS: Visualizing survival and other timed-events data. *Am Statistician*, 46:13–18, 1992. ◊420
221. H. Goldstein. Restricted unbiased iterative generalized least-squares estimation. *Biometrika*, 76(3):622–623, 1989. ◊146, 147
222. R. Goldstein. The comparison of models in discrimination cases. *Jurimetrics J*, 34:215–234, 1994. ◊215
223. M. Gönen and G. Heller. Concordance probability and discriminatory power in proportional hazards regression. *Biometrika*, 92(4):965–970, Dec. 2005. ◊122, 519
224. G. Gong. Cross-validation, the jackknife, and the bootstrap: Excess error estimation in forward logistic regression. *J Am Stat Assoc*, 81:108–113, 1986. ◊ 114
225. T. A. Gooley, W. Leisenring, J. Crowley, and B. E. Storer. Estimation of failure probabilities in the presence of competing risks: New representations of old estimators. *Stat Med*, 18:695–706, 1999. ◊414
226. S. M. Gore, S. J. Pocock, and G. R. Kerr. Regression models and non-proportional hazards in the analysis of breast cancer survival. *Appl Stat*, 33:176–195, 1984. ◊450, 495, 500, 501, 503
227. H. H. H. Göring, J. D. Terwilliger, and J. Blangero. Large upward bias in estimation of locus-specific effects from genomewide scans. *Am J Hum Gen*, 69:1357–1369, 2001. ◊100
228. W. Gould. Confidence intervals in logit and probit models. *Stata Tech Bull*, STB-14:26–28, July 1993. <http://www.stata.com/products/stb/journals/stb14.pdf>. ◊186
229. U. S. Govindarajulu, H. Lin, K. L. Lunetta, and R. B. D’Agostino. Frailty models: Applications to biomedical and genetic studies. *Stat Med*, 30(22):2754–2764, 2011. ◊420
230. U. S. Govindarajulu, D. Spiegelman, S. W. Thurston, B. Ganguli, and E. A. Eisen. Comparing smoothing techniques in Cox models for exposure-response relationships. *Stat Med*, 26:3735–3752, 2007. ◊40
231. I. M. Graham and E. Clavel. Communicating risk — coronary risk scores. *J Roy Stat Soc A*, 166:217–223, 2003. ◊122
232. J. W. Graham, A. E. Olchowski, and T. D. Gilreath. How many imputations are really needed? Some practical clarifications of multiple imputation theory. *Prev Sci*, 8:206–213, 2007. ◊54

233. P. Grambsch and T. Therneau. Proportional hazards tests and diagnostics based on weighted residuals. *Biometrika*, 81:515–526, 1994. Amendment and corrections in 82: 668 (1995). ◊314, 498, 499, 518
234. P. M. Grambsch and P. C. O’Brien. The effects of transformations and preliminary tests for non-linearity in regression. *Stat Med*, 10:697–709, 1991. ◊32, 36, 68
235. B. I. Graubard and E. L. Korn. Regression analysis with clustered data. *Stat Med*, 13:509–522, 1994. ◊214
236. R. J. Gray. Some diagnostic methods for Cox regression models through hazard smoothing. *Biometrics*, 46:93–102, 1990. ◊518
237. R. J. Gray. Flexible methods for analyzing survival data using splines, with applications to breast cancer prognosis. *J Am Stat Assoc*, 87:942–951, 1992. ◊30, 41, 77, 209, 210, 211, 345, 346, 500
238. R. J. Gray. Spline-based tests in survival analysis. *Biometrics*, 50:640–652, 1994. ◊30, 41, 500
239. M. J. Greenacre. Correspondence analysis of multivariate categorical data by weighted least-squares. *Biometrika*, 75:457–467, 1988. ◊81
240. S. Greenland. Alternative models for ordinal logistic regression. *Stat Med*, 13:1665–1677, 1994. ◊324
241. S. Greenland. When should epidemiologic regressions use random coefficients? *Biometrics*, 56:915–921, 2000. ◊68, 100, 215
242. S. Greenland and W. D. Finkle. A critical look at methods for handling missing covariates in epidemiologic regression analyses. *Am J Epi*, 142:1255–1264, 1995. ◊46, 59
243. A. J. Gross and V. A. Clark. *Survival Distributions: Reliability Applications in the Biomedical Sciences*. Wiley, New York, 1975. ◊408
244. S. T. Gross and T. L. Lai. Nonparametric estimation and regression analysis with left-truncated and right-censored data. *J Am Stat Assoc*, 91:1166–1180, 1996. ◊420
245. A. Guisan and F. E. Harrell. Ordinal response regression models in ecology. *J Veg Sci*, 11:617–626, 2000. ◊324
246. J. Guo, G. James, E. Levina, G. Michailidis, and J. Zhu. Principal component analysis with sparse fused loadings. *J Comp Graph Stat*, 19(4):930–946, 2011. ◊101
247. M. J. Gurka, L. J. Edwards, and K. E. Muller. Avoiding bias in mixed model inference for fixed effects. *Stat Med*, 30(22):2696–2707, 2011. ◊160
248. P. Gustafson. Bayesian regression modeling with interactions and smooth effects. *J Am Stat Assoc*, 95:795–806, 2000. ◊41
249. P. Hall and H. Miller. Using generalized correlation to effect variable selection in very high dimensional problems. *J Comp Graph Stat*, 18(3):533–550, 2009. ◊100
250. P. Hall and H. Miller. Using the bootstrap to quantify the authority of an empirical ranking. *Ann Stat*, 37(6B):3929–3959, 2009. ◊117
251. M. Halperin, W. C. Blackwelder, and J. I. Verter. Estimation of the multivariate logistic risk function: A comparison of the discriminant function and maximum likelihood approaches. *J Chron Dis*, 24:125–158, 1971. ◊272
252. D. Hand and M. Crowder. *Practical Longitudinal Data Analysis*. Chapman & Hall, London, 1996. ◊143
253. D. J. Hand. *Construction and Assessment of Classification Rules*. Wiley, Chichester, 1997. ◊273
254. T. L. Hankins. Blood, dirt, and nomograms. *Chance*, 13(1):26–37, 2000. ◊104, 122, 267
255. J. A. Hanley and B. J. McNeil. The meaning and use of the area under a receiver operating characteristic (ROC) curve. *Radiology*, 143:29–36, 1982. ◊257

256. O. Harel and X. Zhou. Multiple imputation: Review of theory, implementation and software. *Stat Med*, 26:3057–3077, 2007. ◊46, 50, 58
257. F. E. Harrell. The LOGIST Procedure. In *SUGI Supplemental Library Users Guide*, pages 269–293. SAS Institute, Inc., Cary, NC, Version 5 edition, 1986. ◊69
258. F. E. Harrell. The PHGLM Procedure. In *SUGI Supplemental Library Users Guide*, pages 437–466. SAS Institute, Inc., Cary, NC, Version 5 edition, 1986. ◊499
259. F. E. Harrell. Comparison of strategies for validating binary logistic regression models. Unpublished manuscript, 1991. ◊115, 259
260. F. E. Harrell. Semiparametric modeling of health care cost and resource utilization. Available from hesweb1.med.virginia.edu/biostat/presentations, 1999. ◊x
261. F. E. Harrell. *rms*: R functions for biostatistical/epidemiologic modeling, testing, estimation, validation, graphics, prediction, and typesetting by storing enhanced model design attributes in the fit, 2013. Implements methods in *Regression Modeling Strategies*, New York:Springer, 2001. ◊127
262. F. E. Harrell, R. M. Califf, D. B. Pryor, K. L. Lee, and R. A. Rosati. Evaluating the yield of medical tests. *JAMA*, 247:2543–2546, 1982. ◊505
263. F. E. Harrell and R. Goldstein. A survey of microcomputer survival analysis software: The need for an integrated framework. *Am Statistician*, 51:360–373, 1997. ◊142
264. F. E. Harrell and K. L. Lee. A comparison of the *discrimination* of discriminant analysis and logistic regression under multivariate normality. In P. K. Sen, editor, *Biostatistics: Statistics in Biomedical, Public Health, and Environmental Sciences. The Bernard G. Greenberg Volume*, pages 333–343. North-Holland, Amsterdam, 1985. ◊205, 207, 258, 272
265. F. E. Harrell and K. L. Lee. The practical value of logistic regression. In *Proceedings of the Tenth Annual SAS Users Group International Conference*, pages 1031–1036, 1985. ◊237
266. F. E. Harrell and K. L. Lee. Verifying assumptions of the Cox proportional hazards model. In *Proceedings of the Eleventh Annual SAS Users Group International Conference*, pages 823–828, Cary, NC, 1986. SAS Institute, Inc. ◊495, 499, 501
267. F. E. Harrell and K. L. Lee. Using logistic model calibration to assess the quality of probability predictions. Unpublished manuscript, 1987. ◊259, 269, 507, 508
268. F. E. Harrell, K. L. Lee, R. M. Califf, D. B. Pryor, and R. A. Rosati. Regression modeling strategies for improved prognostic prediction. *Stat Med*, 3:143–152, 1984. ◊72, 101, 332, 505
269. F. E. Harrell, K. L. Lee, and D. B. Mark. Multivariable prognostic models: Issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat Med*, 15:361–387, 1996. ◊xi, 100
270. F. E. Harrell, K. L. Lee, D. B. Matchar, and T. A. Reichert. Regression models for prognostic prediction: Advantages, problems, and suggested solutions. *Ca Trt Rep*, 69:1071–1077, 1985. ◊41, 72
271. F. E. Harrell, K. L. Lee, and B. G. Pollock. Regression models in clinical studies: Determining relationships between predictors and response. *J Nat Cancer Inst*, 80:1198–1202, 1988. ◊30, 40
272. F. E. Harrell, P. A. Margolis, S. Gove, K. E. Mason, E. K. Mulholland, D. Lehmann, L. Muhe, S. Gatchalian, and H. F. Eichenwald. Development of a clinical prediction model for an ordinal outcome: The World Health Organization ARI Multicentre Study of clinical signs and etiologic agents of pneumonia, sepsis, and meningitis in young infants. *Stat Med*, 17:909–944, 1998. ◊xi, 77, 96, 327

273. D. P. Harrington and T. R. Fleming. A class of rank test procedures for censored survival data. *Biometrika*, 69:553–566, 1982. ◊517
274. T. Hastie. Discussion of “The use of polynomial splines and their tensor products in multivariate function estimation” by C. J. Stone. *Appl Stat*, 22:177–179, 1994. ◊37
275. T. Hastie and R. Tibshirani. *Generalized Additive Models*. Chapman and Hall, London, 1990. ◊29, 41, 142, 390
276. T. J. Hastie, J. L. Botha, and C. M. Schnitzler. Regression with an ordered categorical response. *Stat Med*, 8:785–794, 1989. ◊324
277. T. J. Hastie and R. J. Tibshirani. *Generalized Additive Models*. Chapman & Hall/CRC, Boca Raton, FL, 1990. ISBN 9780412343902. ◊90, 359
278. W. W. Hauck and A. Donner. Wald’s test as applied to hypotheses in logit analysis. *J Am Stat Assoc*, 72:851–863, 1977. ◊193, 234
279. X. He and L. Shen. Linear regression after spline transformation. *Biometrika*, 84:474–481, 1997. ◊82
280. Y. He and A. M. Zaslavsky. Diagnosing imputation models by applying target analyses to posterior replicates of completed data. *Stat Med*, 31(1):1–18, 2012. ◊59
281. G. Heinze and M. Schemper. A solution to the problem of separation in logistic regression. *Stat Med*, 21(16):2409–2419, 2002. ◊203
282. R. Henderson. Problems and prediction in survival-data analysis. *Stat Med*, 14:161–184, 1995. ◊420, 518, 519
283. R. Henderson, M. Jones, and J. Stare. Accuracy of point predictions in survival analysis. *Stat Med*, 20:3083–3096, 2001. ◊519
284. A. V. Hernández, M. J. Eijkemans, and E. W. Steyerberg. Randomized controlled trials with time-to-event outcomes: how much does prespecified covariate adjustment increase power? *Annals of epidemiology*, 16(1):41–48, Jan. 2006. ◊231
285. A. V. Hernández, E. W. Steyerberg, and J. D. F. Habbema. Covariate adjustment in randomized controlled trials with dichotomous outcomes increases statistical power and reduces sample size requirements. *J Clin Epi*, 57:454–460, 2004. ◊231
286. J. E. Herndon and F. E. Harrell. The restricted cubic spline hazard model. *Comm Stat Th Meth*, 19:639–663, 1990. ◊408, 409, 424
287. J. E. Herndon and F. E. Harrell. The restricted cubic spline as baseline hazard in the proportional hazards model with step function time-dependent covariables. *Stat Med*, 14:2119–2129, 1995. ◊408, 424, 501, 518
288. I. Hertz-Picciotto and B. Rockhill. Validity and efficiency of approximation methods for tied survival times in Cox regression. *Biometrics*, 53:1151–1156, 1997. ◊477
289. K. R. Hess. Assessing time-by-covariate interactions in proportional hazards regression models using cubic spline functions. *Stat Med*, 13:1045–1062, 1994. ◊501
290. K. R. Hess. Graphical methods for assessing violations of the proportional hazards assumption in Cox regression. *Stat Med*, 14:1707–1723, 1995. ◊518
291. T. Hielscher, M. Zucknick, W. Werft, and A. Benner. On the prognostic value of survival models with application to gene expression signatures. *Stat Med*, 29:818–829, 2010. ◊518, 519
292. J. Hilden and T. A. Gerds. A note on the evaluation of novel biomarkers: do not rely on integrated discrimination improvement and net reclassification index. *Statist. Med.*, 33(19):3405–3414, Aug. 2014. ◊101
293. S. L. Hillis. Residual plots for the censored data linear regression model. *Stat Med*, 14:2023–2036, 1995. ◊450
294. S. G. Hilsenbeck and G. M. Clark. Practical p -value adjustment for optimally selected cutpoints. *Stat Med*, 15:103–112, 1996. ◊11, 19

295. W. Hoeffding. A non-parametric test of independence. *Ann Math Stat*, 19:546–557, 1948. ◊81, 166
296. H. Hofmann. Simpson on board the Titanic? Interactive methods for dealing with multivariate categorical data. *Stat Comp Graphics News ASA*, 9(2):16–19, 1999. <http://stat-computing.org/newsletter/issues/scgn-09-2.pdf>. ◊291
297. J. W. Hogan and N. M. Laird. Mixture models for the joint distribution of repeated measures and event times. *Stat Med*, 16:239–257, 1997. ◊420
298. J. W. Hogan and N. M. Laird. Model-based approaches to analysing incomplete longitudinal and failure time data. *Stat Med*, 16:259–272, 1997. ◊420
299. M. Hollander, I. W. McKeague, and J. Yang. Likelihood ratio-based confidence bands for survival functions. *J Am Stat Assoc*, 92:215–226, 1997. ◊420
300. N. Holländer, W. Sauerbrei, and M. Schumacher. Confidence intervals for the effect of a prognostic factor after selection of an ‘optimal’ cutpoint. *Stat Med*, 23:1701–1713, 2004. ◊19, 20
301. N. J. Horton and K. P. Kleinman. Much ado about nothing: A comparison of missing data methods and software to fit incomplete data regression models. *Am Statistician*, 61(1):79–90, 2007. ◊59
302. N. J. Horton and S. R. Lipsitz. Multiple imputation in practice: Comparison of software packages for regression models with missing variables. *Am Statistician*, 55:244–254, 2001. ◊54
303. D. W. Hosmer, T. Hosmer, S. le Cessie, and S. Lemeshow. A comparison of goodness-of-fit tests for the logistic regression model. *Stat Med*, 16:965–980, 1997. ◊236
304. D. W. Hosmer and S. Lemeshow. Goodness-of-fit tests for the multiple logistic regression model. *Comm Stat Th Meth*, 9:1043–1069, 1980. ◊236
305. D. W. Hosmer and S. Lemeshow. *Applied Logistic Regression*. Wiley, New York, 1989. ◊255, 272
306. D. W. Hosmer and S. Lemeshow. Confidence interval estimates of an index of quality performance based on logistic regression models. *Stat Med*, 14:2161–2172, 1995. See letter to editor 16:1301-3,1997. ◊272
307. T. Hothorn, F. Bretz, and P. Westfall. Simultaneous inference in general parametric models. *Biometrical J*, 50(3):346–363, 2008. ◊xii, 199, 202
308. P. Hougaard. Fundamentals of survival data. *Biometrics*, 55:13–22, 1999. ◊400, 420, 450
309. B. Hu, M. Palta, and J. Shao. Properties of R^2 statistics for logistic regression. *Stat Med*, 25:1383–1395, 2006. ◊272
310. J. Huang and D. Harrington. Penalized partial likelihood regression for right-censored data with bootstrap selection of the penalty parameter. *Biometrics*, 58:781–791, 2002. ◊215, 478
311. Y. Huang and M. Wang. Frequency of recurrent events at failure times: Modeling and inference. *J Am Stat Assoc*, 98:663–670, 2003. ◊420
312. P. J. Huber. The behavior of maximum likelihood estimates under nonstandard conditions. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, volume 1: Statistics, pages 221–233. University of California Press, Berkeley, CA, 1967. ◊196
313. S. Hunsberger, D. Murray, C. Davis, and R. R. Fabsitz. Imputation strategies for missing data in a school-based multi-center study: the Pathways study. *Stat Med*, 20:305–316, 2001. ◊59
314. C. M. Hurvich and C. Tsai. Regression and time series model selection in small samples. *Biometrika*, 76:297–307, 1989. ◊214, 215
315. C. M. Hurvich and C. Tsai. Model selection for extended quasi-likelihood models in small samples. *Biometrics*, 51:1077–1084, 1995. ◊214
316. C. M. Hurvich and C. L. Tsai. The impact of model selection on inference in linear regression. *Am Statistician*, 44:214–217, 1990. ◊100

317. L. I. Iezzoni. Dimensions of Risk. In L. I. Iezzoni, editor, *Risk Adjustment for Measuring Health Outcomes*, chapter 2, pages 29–118. Foundation of the American College of Healthcare Executives, Ann Arbor, MI, 1994. ◊7
318. R. Ihaka and R. Gentleman. R: A language for data analysis and graphics. *J Comp Graph Stat*, 5:299–314, 1996. ◊127
319. K. Imai, G. King, and O. Lau. Towards a common framework for statistical analysis and development. *J Comp Graph Stat*, 17(4):892–913, 2008. ◊142
320. J. E. Jackson. *A User's Guide to Principal Components*. Wiley, New York, 1991. ◊101
321. K. J. Janssen, A. R. Donders, F. E. Harrell, Y. Vergouwe, Q. Chen, D. E. Grobbee, and K. G. Moons. Missing covariate data in medical research: To impute is better than to ignore. *J Clin Epi*, 63:721–727, 2010. ◊54
322. H. Jiang, R. Chapell, and J. P. Fine. Estimating the distribution of nonterminal event time in the presence of mortality or informative dropout. *Controlled Clin Trials*, 24:135–146, 2003. ◊421
323. N. L. Johnson, S. Kotz, and N. Balakrishnan. *Distributions in Statistics: Continuous Univariate Distributions*, volume 1. Wiley-Interscience, New York, second edition, 1994. ◊408
324. I. T. Jolliffe. Discarding variables in a principal component analysis. I. Artificial data. *Appl Stat*, 21:160–173, 1972. ◊101
325. I. T. Jolliffe. *Principal Component Analysis*. Springer-Verlag, New York, second edition, 2010. ◊101, 172
326. M. P. Jones. Indicator and stratification methods for missing explanatory variables in multiple linear regression. *J Am Stat Assoc*, 91:222–230, 1996. ◊49, 58
327. L. Joseph, P. Belisle, H. Tamim, and J. S. Sampalis. Selection bias found in interpreting analyses with missing data for the prehospital index for trauma. *J Clin Epi*, 57:147–153, 2004. ◊58
328. M. Julien and J. A. Hanley. Profile-specific survival estimates: Making reports of clinical trials more patient-relevant. *CT*, 5:107–115, 2008. ◊122
329. A. C. Justice, K. E. Covinsky, and J. A. Berlin. Assessing the generalizability of prognostic information. *Ann Int Med*, 130:515–524, 1999. ◊122
330. J. D. Kalbfleisch and R. L. Prentice. Marginal likelihood based on Cox's regression and life model. *Biometrika*, 60:267–278, 1973. ◊375, 478
331. J. D. Kalbfleisch and R. L. Prentice. *The Statistical Analysis of Failure Time Data*. Wiley, New York, 1980. ◊411, 412, 414, 420, 436, 441, 483, 496, 517
332. G. Kalton and D. Kasprzyk. The treatment of missing survey data. *Surv Meth*, 12:1–16, 1986. ◊58
333. E. L. Kaplan and P. Meier. Nonparametric estimation from incomplete observations. *J Am Stat Assoc*, 53:457–481, 1958. ◊410
334. T. Karrison. Restricted mean life with adjustment for covariates. *J Am Stat Assoc*, 82:1169–1176, 1987. ◊406, 514
335. T. G. Karrison. Use of Irwin's restricted mean as an index for comparing survival in different treatment groups—Interpretation and power considerations. *Controlled Clin Trials*, 18:151–167, 1997. ◊406, 503
336. J. Karvanen and F. E. Harrell. Visualizing covariates in proportional hazards model. *Stat Med*, 28:1957–1966, 2009. ◊104
337. R. E. Kass and A. E. Raftery. Bayes factors. *J Am Stat Assoc*, 90:773–795, 1995. ◊71, 214
338. M. W. Kattan, G. Heller, and M. F. Brennan. A competing-risks nomogram for sarcoma-specific death following local recurrence. *Stat Med*, 22:3515–3525, 2003. ◊519
339. M. W. Kattan and J. Marasco. What is a real nomogram? *Sem Onc*, 37(1): 23–26, Feb. 2010. ◊104, 122

340. R. Kay. Treatment effects in competing-risks analysis of prostate cancer data. *Biometrics*, 42:203–211, 1986. ◊276, 414, 495
341. R. Kay and S. Little. Assessing the fit of the logistic model: A case study of children with the haemolytic uraemic syndrome. *Appl Stat*, 35:16–30, 1986. ◊272
342. S. Keleş and M. R. Segal. Residual-based tree-structured survival analysis. *Stat Med*, 21:313–326, 2002. ◊41
343. P. J. Kelly and L. Lim. Survival analysis for recurrent event data: An application to childhood infectious diseases. *Stat Med*, 19:13–33, 2000. ◊421
344. D. M. Kent and R. Hayward. Limitations of applying summary results of clinical trials to individual patients. *JAMA*, 298:1209–1212, 2007. ◊4
345. J. T. Kent and J. O’Quigley. Measures of dependence for censored survival data. *Biometrika*, 75:525–534, 1988. ◊505
346. M. G. Kenward, I. R. White, and J. R. Carpenter. Should baseline be a covariate or dependent variable in analyses of change from baseline in clinical trials? (letter to the editor). *Stat Med*, 29:1455–1456, 2010. ◊160
347. H. J. Keselman, J. Algina, R. K. Kowalchuk, and R. D. Wolfinger. A comparison of two approaches for selecting covariance structures in the analysis of repeated measurements. *Comm Stat - Sim Comp*, 27:591–604, 1998. ◊69, 160
348. V. Kipnis. Relevancy criterion for discriminating among alternative model specifications. In K. Berk and L. Malone, editors, *Proceedings of the 21st Symposium on the Interface between Computer Science and Statistics*, pages 376–381, Alexandria, VA, 1989. American Statistical Association. ◊123
349. J. P. Klein, N. Keiding, and E. A. Copelan. Plotting summary predictions in multistate survival models: Probabilities of relapse and death in remission for bone marrow transplantation patients. *Stat Med*, 12:2314–2332, 1993. ◊415
350. J. P. Klein and M. L. Moeschberger. *Survival Analysis: Techniques for Censored and Truncated Data*. Springer, New York, 1997. ◊420, 517
351. W. A. Knaus, F. E. Harrell, C. J. Fisher, D. P. Wagner, S. M. Opan, J. C. Sadoff, E. A. Draper, C. A. Walawander, K. Conboy, and T. H. Grasela. The clinical evaluation of new drugs for sepsis: A prospective study design based on survival analysis. *JAMA*, 270:1233–1241, 1993. ◊4
352. W. A. Knaus, F. E. Harrell, J. Lynn, L. Goldman, R. S. Phillips, A. F. Connors, N. V. Dawson, W. J. Fulkerson, R. M. Califf, N. Desbiens, P. Layde, R. K. Oye, P. E. Bellamy, R. B. Hakim, and D. P. Wagner. The SUPPORT prognostic model: Objective estimates of survival for seriously ill hospitalized adults. *Ann Int Med*, 122:191–203, 1995. ◊59, 84, 86, 453
353. M. J. Knol, K. J. M. Janssen, R. T. Donders, A. C. G. Egberts, E. R. Heerding, D. E. Grobbee, K. G. M. Moons, and M. I. Geerlings. Unpredictable bias when using the missing indicator method or complete case analysis for missing confounder values: an empirical example. *J Clin Epi*, 63:728–736, 2010. ◊47, 49
354. G. G. Koch, I. A. Amara, and J. M. Singer. A two-stage procedure for the analysis of ordinal categorical data. In P. K. Sen, editor, *BIOSTATISTICS: Statistics in Biomedical, Public Health and Environmental Sciences*. Elsevier Science Publishers B. V. (North-Holland), Amsterdam, 1985. ◊324
355. R. Koenker. *Quantile Regression*. Cambridge University Press, New York, 2005. ISBN-10: 0-521-60827-9; ISBN-13: 978-0-521-60827-5. ◊360
356. R. Koenker. *quantreg: Quantile Regression*, 2009. R package version 4.38. ◊131, 360
357. R. Koenker and G. Bassett. Regression quantiles. *Econometrica*, 46:33–50, 1978. ◊131, 360, 392
358. M. T. Koller, H. Raatz, E. W. Steyerberg, and M. Wolbers. Competing risks and the clinical community: irrelevance or ignorance? *Stat Med*, 31(11–12):1089–1097, 2012. ◊420

359. S. Konishi and G. Kitagawa. *Information Criteria and Statistical Modeling*. Springer, New York, 2008. ISBN 978-0-387-71886-6. ◊204
360. C. Kooperberg and D. B. Clarkson. Hazard regression with interval-censored data. *Biometrics*, 53:1485–1494, 1997. ◊420, 450
361. C. Kooperberg, C. J. Stone, and Y. K. Truong. Hazard regression. *J Am Stat Assoc*, 90:78–94, 1995. ◊178, 419, 420, 422, 424, 450, 473, 506, 508, 518, 530
362. E. L. Korn and F. J. Dorey. Applications of crude incidence curves. *Stat Med*, 11:813–829, 1992. ◊416
363. E. L. Korn and B. I. Graubard. Analysis of large health surveys: Accounting for the sampling design. *J Roy Stat Soc A*, 158:263–295, 1995. ◊208
364. E. L. Korn and B. I. Graubard. Examples of differing weighted and unweighted estimates from a sample survey. *Am Statistician*, 49:291–295, 1995. ◊208
365. E. L. Korn and R. Simon. Measures of explained variation for survival data. *Stat Med*, 9:487–503, 1990. ◊206, 215, 505, 519
366. E. L. Korn and R. Simon. Explained residual variation, explained risk, and goodness of fit. *Am Statistician*, 45:201–206, 1991. ◊206, 215, 273
367. D. Kronborg and P. Aaby. Piecewise comparison of survival functions in stratified proportional hazards models. *Biometrics*, 46:375–380, 1990. ◊502
368. W. F. Kuhfeld. The PRINQUAL procedure. In *SAS/STAT 9.2 User's Guide*. SAS Publishing, Cary, NC, second edition, 2009. ◊82, 167
369. G. P. S. Kwong and J. L. Hutton. Choice of parametric models in survival analysis: applications to monotherapy for epilepsy and cerebral palsy. *Appl Stat*, 52:153–168, 2003. ◊450
370. J. M. Lachin and M. A. Foulkes. Evaluation of sample size and power for analyses of survival with allowance for nonuniform patient entry, losses to follow-up, noncompliance, and stratification. *Biometrics*, 42:507–519, 1986. ◊513
371. L. Lamport. *L^AT_EX: A Document Preparation System*. Addison-Wesley, Reading, MA, second edition, 1994. ◊536
372. R. Lancar, A. Kramar, and C. Haie-Meder. Non-parametric methods for analysing recurrent complications of varying severity. *Stat Med*, 14:2701–2712, 1995. ◊421
373. J. M. Landwehr, D. Pregibon, and A. C. Shoemaker. Graphical methods for assessing logistic regression models (with discussion). *J Am Stat Assoc*, 79:61–83, 1984. ◊272, 315
374. T. P. Lane and W. H. DuMouchel. Simultaneous confidence intervals in multiple regression. *Am Statistician*, 48:315–321, 1994. ◊199
375. K. Larsen and J. Merlo. Appropriate assessment of neighborhood effects on individual health: integrating random and fixed effects in multilevel logistic regression. *American journal of epidemiology*, 161(1):81–88, Jan. 2005. ◊122
376. M. G. Larson and G. E. Dinse. A mixture model for the regression analysis of competing risks data. *Appl Stat*, 34:201–211, 1985. ◊276, 414
377. P. W. Laud and J. G. Ibrahim. Predictive model selection. *J Roy Stat Soc B*, 57:247–262, 1995. ◊214
378. A. Laupacis, N. Sekar, and I. G. Stiell. Clinical prediction rules: A review and suggested modifications of methodological standards. *JAMA*, 277:488–494, 1997. ◊x, 6
379. B. Lausen and M. Schumacher. Evaluating the effect of optimized cutoff values in the assessment of prognostic factors. *Comp Stat Data Analysis*, 21(3):307–326, 1996. ◊11, 19
380. P. W. Lavori, R. Dawson, and T. B. Mueller. Causal estimation of time-varying treatment effects in observational studies: Application to depressive disorder. *Stat Med*, 13:1089–1100, 1994. ◊231
381. P. W. Lavori, R. Dawson, and D. Shera. A multiple imputation strategy for clinical trials with truncation of patient data. *Stat Med*, 14:1913–1925, 1995. ◊47

382. J. F. Lawless. *Statistical Models and Methods for Lifetime Data*. Wiley, New York, 1982. ◊420, 450, 485, 517
383. J. F. Lawless. The analysis of recurrent events for multiple subjects. *Appl Stat*, 44:487–498, 1995. ◊421
384. J. F. Lawless and C. Nadeau. Some simple robust methods for the analysis of recurrent events. *Technometrics*, 37:158–168, 1995. ◊420, 421
385. J. F. Lawless and K. Singhal. Efficient screening of nonnormal regression models. *Biometrics*, 34:318–327, 1978. ◊70, 137
386. J. F. Lawless and Y. Yuan. Estimation of prediction error for survival models. *Stat Med*, 29:262–274, 2010. ◊519
387. S. le Cessie and J. C. van Houwelingen. A goodness-of-fit test for binary regression models, based on smoothing methods. *Biometrics*, 47:1267–1282, 1991. ◊236
388. S. le Cessie and J. C. van Houwelingen. Ridge estimators in logistic regression. *Appl Stat*, 41:191–201, 1992. ◊77, 209
389. M. LeBlanc and J. Crowley. Survival trees by goodness of fit. *J Am Stat Assoc*, 88:457–467, 1993. ◊41
390. M. LeBlanc and R. Tibshirani. Adaptive principal surfaces. *J Am Stat Assoc*, 89:53–64, 1994. ◊101
391. A. Leclerc, D. Luce, F. Lert, J. F. Chastang, and P. Logeay. Correspondence analysis and logistic modelling: Complementary use in the analysis of a health survey among nurses. *Stat Med*, 7:983–995, 1988. ◊81
392. E. T. Lee. *Statistical Methods for Survival Data Analysis*. Lifetime Learning Publications, Belmont, CA, second edition, 1980. ◊420
393. E. W. Lee, L. J. Wei, and D. A. Amato. Cox-type regression analysis for large numbers of small groups of correlated failure time observations. In J. P. Klein and P. K. Goel, editors, *Survival Analysis: State of the Art*, NATO ASI, pages 237–247. Kluwer Academic, Boston, 1992. ◊197
394. J. J. Lee, K. R. Hess, and J. A. Dubin. Extensions and applications of event charts. *Am Statistician*, 54:63–70, 2000. ◊418, 420
395. K. L. Lee, D. B. Pryor, F. E. Harrell, R. M. Califf, V. S. Behar, W. L. Floyd, J. J. Morris, R. A. Waugh, R. E. Whalen, and R. A. Rosati. Predicting outcome in coronary disease: Statistical models versus expert clinicians. *Am J Med*, 80:553–560, 1986. ◊205
396. S. Lee, J. Z. Huang, and J. Hu. Sparse logistic principal components analysis for binary data. *Ann Appl Stat*, 4(3):1579–1601, 2010. ◊101
397. E. L. Lehmann. Model specification: The views of Fisher and Neyman and later developments. *Statistical Sci*, 5:160–168, 1990. ◊8, 10
398. S. Lehr and M. Schemper. Parsimonious analysis of time-dependent effects in the Cox model. *Stat Med*, 26:2686–2698, 2007. ◊501
399. F. Leisch. Sweave: Dynamic Generation of Statistical Reports Using Literate Data Analysis. In W. Härdle and B. Rönz, editors, *Compstat 2002 — Proceedings in Computational Statistics*, pages 575–580. Physica Verlag, Heidelberg, 2002. ISBN 3-7908-1517-9. ◊138
400. L. F. León and C. Tsai. Functional form diagnostics for Cox’s proportional hazards model. *Biometrics*, 60:75–84, 2004. ◊518
401. M. A. H. Levine, A. I. El-Nahas, and B. Asa. Relative risk and odds ratio data are still portrayed with inappropriate scales in the medical literature. *J Clin Epi*, 63:1045–1047, 2010. ◊122
402. C. Li and B. E. Shepherd. A new residual for ordinal outcomes. *Biometrika*, 99(2):473–480, 2012. ◊315
403. K. Li, J. Wang, and C. Chen. Dimension reduction for censored regression data. *Ann Stat*, 27:1–23, 1999. ◊101
404. K. C. Li. Sliced inverse regression for dimension reduction. *J Am Stat Assoc*, 86:316–327, 1991. ◊101

405. K.-Y. Liang and S. L. Zeger. Longitudinal data analysis of continuous and discrete responses for pre-post designs. *Sankhyā*, 62:134–148, 2000. ◊160
406. J. G. Liao and D. McGee. Adjusted coefficients of determination for logistic regression. *Am Statistician*, 57:161–165, 2003. ◊273
407. D. Y. Lin. Cox regression analysis of multivariate failure time data: The marginal approach. *Stat Med*, 13:2233–2247, 1994. ◊197, 213, 417, 418
408. D. Y. Lin. Non-parametric inference for cumulative incidence functions in competing risks studies. *Stat Med*, 16:901–910, 1997. ◊415
409. D. Y. Lin. On fitting Cox’s proportional hazards models to survey data. *Biometrika*, 87:37–47, 2000. ◊215
410. D. Y. Lin and L. J. Wei. The robust inference for the Cox proportional hazards model. *J Am Stat Assoc*, 84:1074–1078, 1989. ◊197, 213, 487
411. D. Y. Lin, L. J. Wei, and Z. Ying. Checking the Cox model with cumulative sums of martingale-based residuals. *Biometrika*, 80:557–572, 1993. ◊518
412. D. Y. Lin and Z. Ying. Semiparametric regression analysis of longitudinal data with informative drop-outs. *Biostatistics*, 4:385–398, 2003. ◊47
413. J. C. Lindsey and L. M. Ryan. Tutorial in biostatistics: Methods for interval-censored data. *Stat Med*, 17:219–238, 1998. ◊420
414. J. K. Lindsey. *Models for Repeated Measurements*. Clarendon Press, 1997. ◊143
415. J. K. Lindsey and B. Jones. Choosing among generalized linear models applied to medical data. *Stat Med*, 17:59–68, 1998. ◊11
416. K. Linnet. Assessing diagnostic tests by a strictly proper scoring rule. *Stat Med*, 8:609–618, 1989. ◊114, 123, 257, 258
417. S. R. Lipsitz, L. P. Zhao, and G. Molenberghs. A semiparametric method of multiple imputation. *J Roy Stat Soc B*, 60:127–144, 1998. ◊54
418. R. Little and H. An. Robust likelihood-based analysis of multivariate data with missing values. *Statistica Sinica*, 14:949–968, 2004. ◊57, 59
419. R. J. Little. Missing Data. In *Ency of Biostatistics*, pages 2622–2635. Wiley, New York, 1998. ◊59
420. R. J. A. Little. Missing-data adjustments in large surveys. *J Bus Econ Stat*, 6:287–296, 1988. ◊51
421. R. J. A. Little. Regression with missing X ’s: A review. *J Am Stat Assoc*, 87:1227–1237, 1992. ◊50, 51, 54
422. R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data*. Wiley, New York, second edition, 2002. ◊48, 52, 54, 59
423. G. F. Liu, K. Lu, R. Mogg, M. Mallick, and D. V. Mehrotra. Should baseline be a covariate or dependent variable in analyses of change from baseline in clinical trials? *Stat Med*, 28:2509–2530, 2009. ◊160
424. K. Liu and A. R. Dyer. A rank statistic for assessing the amount of variation explained by risk factors in epidemiologic studies. *Am J Epi*, 109:597–606, 1979. ◊206, 256
425. R. Lockhart, J. Taylor, R. J. Tibshirani, and R. Tibshirani. A significance test for the lasso. Technical report, arXiv, 2013. ◊68
426. J. S. Long and L. H. Ervin. Using heteroscedasticity consistent standard errors in the linear regression model. *Am Statistician*, 54:217–224, 2000. ◊213
427. J. Lubsen, J. Pool, and E. van der Does. A practical device for the application of a diagnostic or prognostic function. *Meth Info Med*, 17:127–129, 1978. ◊104
428. D. J. Lunn, J. Wakefield, and A. Racine-Poon. Cumulative logit models for ordinal data: a case study involving allergic rhinitis severity scores. *Stat Med*, 20:2261–2285, 2001. ◊324
429. M. Lunn and D. McNeil. Applying Cox regression to competing risks. *Biometrics*, 51:524–532, 1995. ◊420
430. X. Luo, L. A. Stfanski, and D. D. Boos. Tuning variable selection procedures by adding noise. *Technometrics*, 48:165–175, 2006. ◊11, 100

431. G. S. Maddala. *Limited-Dependent and Qualitative Variables in Econometrics*. Cambridge University Press, Cambridge, UK, 1983. ◊206, 256, 505
432. L. Magee. R^2 measures based on Wald and likelihood ratio joint significance tests. *Am Statistician*, 44:250–253, 1990. ◊206, 256, 505
433. L. Magee. Nonlocal behavior in polynomial regressions. *Am Statistician*, 52:20–22, 1998. ◊21
434. C. Mallows. The zeroth problem. *Am Statistician*, 52:1–9, 1998. ◊11
435. M. Mandel. Censoring and truncation—Highlighting the differences. *Am Statistician*, 61(4):321–324, 2007. ◊420
436. M. Mandel, N. Galae, and E. Simchen. Evaluating survival model performance: a graphical approach. *Stat Med*, 24:1933–1945, 2005. ◊518
437. N. Mantel. Why stepdown procedures in variable selection. *Technometrics*, 12:621–625, 1970. ◊70
438. N. Mantel and D. P. Byar. Evaluation of response-time data involving transient states: An illustration using heart-transplant data. *J Am Stat Assoc*, 69:81–86, 1974. ◊401, 420
439. P. Margolis, E. K. Mulholland, F. E. Harrell, S. Gove, and the WHO Young Infants Study Group. Clinical prediction of serious bacterial infections in young infants in developing countries. *Pediatr Infect Dis J*, 18S:S23–S31, 1999. ◊327
440. D. B. Mark, M. A. Hlatky, F. E. Harrell, K. L. Lee, R. M. Califf, and D. B. Pryor. Exercise treadmill score for predicting prognosis in coronary artery disease. *Ann Int Med*, 106:793–800, 1987. ◊512
441. G. Marshall, F. L. Grover, W. G. Henderson, and K. E. Hammermeister. Assessment of predictive models for binary outcomes: An empirical approach using operative death from cardiac surgery. *Stat Med*, 13:1501–1511, 1994. ◊101
442. G. Marshall, B. Warner, S. MaWhinney, and K. Hammermeister. Prospective prediction in the presence of missing data. *Stat Med*, 21:561–570, 2002. ◊57
443. R. J. Marshall. The use of classification and regression trees in clinical epidemiology. *J Clin Epi*, 54:603–609, 2001. ◊41
444. E. Marubini and M. G. Valsecchi. *Analyzing Survival Data from Clinical Trials and Observational Studies*. Wiley, Chichester, 1995. ◊213, 214, 415, 420, 501, 517
445. J. M. Massaro. Battery Reduction. 2005. ◊87
446. S. E. Maxwell and H. D. Delaney. Bivariate median splits and spurious statistical significance. *Psych Bull*, 113:181–190, 1993. ◊19
447. M. May, P. Royston, M. Egger, A. C. Justice, and J. A. C. Sterne. Development and validation of a prognostic model for survival time data: application to prognosis of HIV positive patients treated with antiretroviral therapy. *Stat Med*, 23:2375–2398, 2004. ◊505
448. G. P. McCabe. Principal variables. *Technometrics*, 26:137–144, 1984. ◊101
449. P. McCullagh. Regression models for ordinal data. *J Roy Stat Soc B*, 42:109–142, 1980. ◊313, 324
450. P. McCullagh and J. A. Nelder. *Generalized Linear Models*. Chapman and Hall/CRC, second edition, Aug. 1989. ◊viii
451. D. R. McNeil, J. Trussell, and J. C. Turner. Spline interpolation of demographic data. *Demography*, 14:245–252, 1977. ◊40
452. W. Q. Meeker and L. A. Escobar. Teaching about approximate confidence regions based on maximum likelihood estimation. *Am Statistician*, 49:48–53, 1995. ◊214
453. N. Meinshausen. Hierarchical testing of variable importance. *Biometrika*, 95(2):265–278, 2008. ◊101
454. S. Menard. Coefficients of determination for multiple logistic regression analysis. *Am Statistician*, 54:17–24, 2000. ◊215, 272
455. X. Meng. Multiple-imputation inferences with uncongenial sources of input. *Stat Sci*, 9:538–558, 1994. ◊58

456. G. Michailidis and J. de Leeuw. The Gifi system of descriptive multivariate analysis. *Statistical Sci*, 13:307–336, 1998. ◊81
457. M. E. Miller, S. L. Hui, and W. M. Tierney. Validation techniques for logistic regression models. *Stat Med*, 10:1213–1226, 1991. ◊259
458. M. E. Miller, T. M. Morgan, M. A. Espeland, and S. S. Emerson. Group comparisons involving missing data in clinical trials: a comparison of estimates and power (size) for some simple approaches. *Stat Med*, 20:2383–2397, 2001. ◊58
459. R. G. Miller. What price Kaplan–Meier? *Biometrics*, 39:1077–1081, 1983. ◊420
460. S. Minkin. Profile-likelihood-based confidence intervals. *Appl Stat*, 39:125–126, 1990. ◊214
461. M. Mittlböck and M. Schemper. Explained variation for logistic regression. *Stat Med*, 15:1987–1997, 1996. ◊215, 273
462. K. G. M. Moons, Donders, E. W. Steyerberg, and F. E. Harrell. Penalized maximum likelihood estimation to directly adjust diagnostic and prognostic prediction models for overoptimism: a clinical example. *J Clin Epi*, 57:1262–1270, 2004. ◊215, 273, 356
463. K. G. M. Moons, R. A. R. T. Donders, T. Stijnen, and F. E. Harrell. Using the outcome for imputation of missing predictor values was preferred. *J Clin Epi*, 59:1092–1101, 2006. ◊54, 55, 59
464. B. J. T. Morgan, K. J. Palmer, and M. S. Ridout. Negative score test statistic (with discussion). *Am Statistician*, 61(4):285–295, 2007. ◊213
465. B. K. Moser and L. P. Coombs. Odds ratios for a continuous outcome variable without dichotomizing. *Stat Med*, 23:1843–1860, 2004. ◊19
466. G. S. Mudholkar, D. K. Srivastava, and G. D. Kollia. A generalization of the Weibull distribution with application to the analysis of survival data. *J Am Stat Assoc*, 91:1575–1583, 1996. ◊420
467. L. R. Muenz. Comparing survival distributions: A review for nonstatisticians. II. *Ca Invest*, 1:537–545, 1983. ◊495, 502
468. V. M. R. Muggego and M. Tagliavia. A flexible approach to the crossing hazards problem. *Stat Med*, 29:1947–1957, 2010. ◊518
469. H. Murad, A. Fleischman, S. Sadetzki, O. Geyer, and L. S. Freedman. Small samples and ordered logistic regression: Does it help to collapse categories of outcome? *Am Statistician*, 57:155–160, 2003. ◊324
470. R. H. Myers. *Classical and Modern Regression with Applications*. PWS-Kent, Boston, 1990. ◊78
471. N. J. D. Nagelkerke. A note on a general definition of the coefficient of determination. *Biometrika*, 78:691–692, 1991. ◊206, 256, 505
472. W. B. Nelson. Theory and applications of hazard plotting for censored failure data. *Technometrics*, 14:945–965, 1972. ◊413
473. R. Newson. Parameters behind “nonparametric” statistics: Kendall’s tau, Somers’ D and median differences. *Stata Journal*, 2(1), 2002. <http://www.stata-journal.com/article.html?article=st0007>. ◊273
474. R. Newson. Confidence intervals for rank statistics: Somers’ D and extensions. *Stata J*, 6(3):309–334, 2006. ◊273
475. N. H. Ng’andu. An empirical comparison of statistical tests for assessing the proportional hazards assumption of Cox’s model. *Stat Med*, 16:611–626, 1997. ◊518
476. T. G. Nick and J. M. Hardin. Regression modeling strategies: An illustrative case study from medical rehabilitation outcomes research. *Am J Occ Ther*, 53:459–470, 1999. ◊viii, 100
477. M. A. Nicolaie, H. C. van Houwelingen, T. M. de Witte, and H. Putter. Dynamic prediction by landmarking in competing risks. *Stat Med*, 32(12):2031–2047, 2013. ◊447
478. M. Nishikawa, T. Tango, and M. Ogawa. Non-parametric inference of adverse events under informative censoring. *Stat Med*, 25:3981–4003, 2006. ◊420

479. P. C. O'Brien. Comparing two samples: Extensions of the t , rank-sum, and log-rank test. *J Am Stat Assoc*, 83:52–61, 1988. ◊231
480. P. C. O'Brien, D. Zhang, and K. R. Bailey. Semi-parametric and non-parametric methods for clinical trials with incomplete data. *Stat Med*, 24:341–358, 2005. ◊47
481. J. O'Quigley, R. Xu, and J. Stare. Explained randomness in proportional hazards models. *Stat Med*, 24(3):479–489, 2005. ◊505
482. W. Original. *survival: Survival analysis, including penalised likelihood*, 2009. R package version 2.37-7. ◊131
483. M. Y. Park and T. Hastie. Penalized logistic regression for detecting gene interactions. *Biostat*, 9(1):30–50, 2008. ◊215
484. M. K. B. Parmar and D. Machin. *Survival Analysis: A Practical Approach*. Wiley, Chichester, 1995. ◊420
485. D. Paul, E. Bair, T. Hastie, and R. Tibshirani. “Preconditioning” for feature selection and regression in high-dimensional problems. *Ann Stat*, 36(4):1595–1619, 2008. ◊121
486. P. Peduzzi, J. Concato, A. R. Feinstein, and T. R. Holford. Importance of events per independent variable in proportional hazards regression analysis. II. Accuracy and precision of regression estimates. *J Clin Epi*, 48:1503–1510, 1995. ◊100
487. P. Peduzzi, J. Concato, E. Kemper, T. R. Holford, and A. R. Feinstein. A simulation study of the number of events per variable in logistic regression analysis. *J Clin Epi*, 49:1373–1379, 1996. ◊73, 100
488. N. Peek, D. G. T. Arts, R. J. Bosman, P. H. J. van der Voort, and N. F. de Keizer. External validation of prognostic models for critically ill patients required substantial sample sizes. *J Clin Epi*, 60:491–501, 2007. ◊93
489. M. J. Pencina and R. B. D'Agostino. Overall C as a measure of discrimination in survival analysis: model specific population value and confidence interval estimation. *Stat Med*, 23:2109–2123, 2004. ◊519
490. M. J. Pencina, R. B. D'Agostino, and O. V. Demler. Novel metrics for evaluating improvement in discrimination: net reclassification and integrated discrimination improvement for normal variables and nested models. *Stat Med*, 31(2):101–113, 2012. ◊101, 142, 273
491. M. J. Pencina, R. B. D'Agostino, and L. Song. Quantifying discrimination of Framingham risk functions with different survival C statistics. *Stat Med*, 31(15):1543–1553, 2012. ◊519
492. M. J. Pencina, R. B. D'Agostino, and E. W. Steyerberg. Extensions of net reclassification improvement calculations to measure usefulness of new biomarkers. *Stat Med*, 30:11–21, 2011. ◊101, 142
493. M. J. Pencina, R. B. D'Agostino Sr, R. B. D'Agostino Jr, and R. S. Vasan. Evaluating the added predictive ability of a new marker: From area under the ROC curve to reclassification and beyond. *Stat Med*, 27:157–172, 2008. ◊93, 101, 142, 273
494. M. S. Pepe. Inference for events with dependent risks in multiple endpoint studies. *J Am Stat Assoc*, 86:770–778, 1991. ◊415
495. M. S. Pepe and J. Cai. Some graphical displays and marginal regression analyses for recurrent failure times and time dependent covariates. *J Am Stat Assoc*, 88:811–820, 1993. ◊417
496. M. S. Pepe, G. Longton, and M. Thornquist. A qualifier Q for the survival function to describe the prevalence of a transient condition. *Stat Med*, 10: 413–421, 1991. ◊415
497. M. S. Pepe and M. Mori. Kaplan–Meier, marginal or conditional probability curves in summarizing competing risks failure time data? *Stat Med*, 12: 737–751, 1993. ◊415

498. A. Perperoglou, A. Keramopoulos, and H. C. van Houwelingen. Approaches in modelling long-term survival: An application to breast cancer. *Stat Med*, 26:2666–2685, 2007. ◊501, 518
499. A. Perperoglou, S. le Cessie, and H. C. van Houwelingen. Reduced-rank hazard regression for modelling non-proportional hazards. *Stat Med*, 25:2831–2845, 2006. ◊518
500. S. A. Peters, M. L. Bots, H. M. den Ruijter, M. K. Palmer, D. E. Grobbee, J. R. Crouse, D. H. O’Leary, G. W. Evans, J. S. Raichlen, K. G. Moons, H. Koffijberg, and METEOR study group. Multiple imputation of missing repeated outcome measurements did not add to linear mixed-effects models. *J Clin Epi*, 65(6):686–695, 2012. ◊160
501. B. Peterson and S. L. George. Sample size requirements and length of study for testing interaction in a $1 \times k$ factorial design when time-to-failure is the outcome. *Controlled Clin Trials*, 14:511–522, 1993. ◊513
502. B. Peterson and F. E. Harrell. Partial proportional odds models for ordinal response variables. *Appl Stat*, 39:205–217, 1990. ◊315, 321, 324
503. A. N. Pettitt and I. Bin Daud. Investigating time dependence in Cox’s proportional hazards model. *Appl Stat*, 39:313–329, 1990. ◊498, 518
504. A. N. Phillips, S. G. Thompson, and S. J. Pocock. Prognostic scores for detecting a high risk group: Estimating the sensitivity when applied to new data. *Stat Med*, 9:1189–1198, 1990. ◊100, 101
505. R. R. Picard and K. N. Berk. Data splitting. *Am Statistician*, 44:140–147, 1990. ◊122
506. R. R. Picard and R. D. Cook. Cross-validation of regression models. *J Am Stat Assoc*, 79:575–583, 1984. ◊123
507. L. W. Pickle. Maximum likelihood estimation in the new computing environment. *Stat Comp Graphics News ASA*, 2(2):6–15, Nov. 1991. ◊213
508. M. C. Pike. A method of analysis of certain class of experiments in carcinogenesis. *Biometrics*, 22:142–161, 1966. ◊441, 442, 443, 480
509. J. C. Pinheiro and D. M. Bates. *Mixed-Effects Models in S and S-PLUS*. Springer, New York, 2000. ◊131, 143, 146, 147, 148
510. R. F. Potthoff and S. N. Roy. A generalized multivariate analysis of variance model useful especially for growth curve problems. *Biometrika*, 51:313–326, 1964. ◊146
511. D. Pregibon. Logistic regression diagnostics. *Ann Stat*, 9:705–724, 1981. ◊255
512. D. Pregibon. Resistant fits for some commonly used logistic models with medical applications. *Biometrics*, 38:485–498, 1982. ◊272
513. R. L. Prentice, J. D. Kalbfleisch, A. V. Peterson, N. Flournoy, V. T. Farewell, and N. E. Breslow. The analysis of failure times in the presence of competing risks. *Biometrics*, 34:541–554, 1978. ◊414
514. S. J. Press and S. Wilson. Choosing between logistic regression and discriminant analysis. *J Am Stat Assoc*, 73:699–705, 1978. ◊272
515. D. B. Pryor, F. E. Harrell, K. L. Lee, R. M. Califf, and R. A. Rosati. Estimating the likelihood of significant coronary artery disease. *Am J Med*, 75:771–780, 1983. ◊273
516. D. B. Pryor, F. E. Harrell, J. S. Rankin, K. L. Lee, L. H. Muhlbaier, H. N. Oldham, M. A. Hlatky, D. B. Mark, J. G. Reves, and R. M. Califf. The changing survival benefits of coronary revascularization over time. *Circulation (Supplement V)*, 76:13–21, 1987. ◊511
517. H. Putter, M. Fiocco, and R. B. Geskus. Tutorial in biostatistics: Competing risks and multi-state models. *Stat Med*, 26:2389–2430, 2007. ◊420
518. H. Putter, M. Sasako, H. H. Hartgrink, C. J. H. van de Velde, and J. C. van Houwelingen. Long-term survival with non-proportional hazards: results from the Dutch Gastric Cancer Trial. *Stat Med*, 24:2807–2821, 2005. ◊518

519. C. Quantin, T. Moreau, B. Asselain, J. Maccaria, and J. Lellouch. A regression survival model for testing the proportional hazards assumption. *Biometrics*, 52:874–885, 1996. ◊518
520. R Development Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2013. ◊127
521. D. R. Ragland. Dichotomizing continuous outcome variables: Dependence of the magnitude of association and statistical power on the cutpoint. *Epi*, 3:434–440, 1992. See letters to editor May 1993 P. 274-, Vol 4 No. 3. ◊11, 19
522. B. M. Reilly and A. T. Evans. Translating clinical research into clinical practice: Impact of using prediction rules to make decisions. *Ann Int Med*, 144:201–209, 2006. ◊6
523. M. Reilly and M. Pepe. The relationship between hot-deck multiple imputation and weighted likelihood. *Stat Med*, 16:5–19, 1997. ◊59
524. B. D. Ripley and P. J. Solomon. Statistical models for prevalent cohort data. *Biometrics*, 51:373–374, 1995. ◊420
525. J. S. Roberts and G. M. Capalbo. A SAS macro for estimating missing values in multivariate data. In *Proceedings of the Twelfth Annual SAS Users Group International Conference*, pages 939–941, Cary, NC, 1987. SAS Institute, Inc. ◊52
526. J. M. Robins, S. D. Mark, and W. K. Newey. Estimating exposure effects by modeling the expectation of exposure conditional on confounders. *Biometrics*, 48:479–495, 1992. ◊231
527. L. D. Robinson and N. P. Jewell. Some surprising results about covariate adjustment in logistic regression models. *Int Stat Rev*, 59:227–240, 1991. ◊231
528. E. B. Roecker. Prediction error and its estimation for subset-selected models. *Technometrics*, 33:459–468, 1991. ◊100, 112
529. W. H. Rogers. Regression standard errors in clustered samples. *Stata Tech Bull*, STB-13:19–23, May 1993. <http://www.stata.com/products/stb/journals/stb13.pdf>. ◊197
530. P. R. Rosenbaum and D. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70:41–55, 1983. ◊3, 231
531. P. R. Rosenbaum and D. B. Rubin. Assessing sensitivity to an unobserved binary covariate in an observational study with binary outcome. *J Roy Stat Soc B*, 45:212–218, 1983. ◊231
532. P. Royston and D. G. Altman. Regression using fractional polynomials of continuous covariates: Parsimonious parametric modelling. *ApplStat*, 43:429–453, 1994. Discussion pp. 453–467. ◊40
533. P. Royston, D. G. Altman, and W. Sauerbrei. Dichotomizing continuous predictors in multiple regression: a bad idea. *Stat Med*, 25:127–141, 2006. ◊19
534. P. Royston and S. G. Thompson. Comparing non-nested regression models. *Biometrics*, 51:114–127, 1995. ◊215
535. D. Rubin and N. Schenker. Multiple imputation in health-care data bases: An overview and some applications. *Stat Med*, 10:585–598, 1991. ◊46, 50, 59
536. D. B. Rubin. *Multiple Imputation for Nonresponse in Surveys*. Wiley, New York, 1987. ◊54, 59
537. S. Sahoo and D. Sengupta. Some diagnostic plots and corrective adjustments for the proportional hazards regression model. *J Comp Graph Stat*, 20(2):375–394, 2011. ◊518
538. S. Sardy. On the practice of rescaling covariates. *Int Stat Rev*, 76:285–297, 2008. ◊215
539. W. Sarle. The VARCLUS procedure. In *SAS/STAT User's Guide*, volume 2, chapter 43, pages 1641–1659. SAS Institute, Inc., Cary, NC, fourth edition, 1990. ◊79, 81, 101

540. SAS Institute, Inc. *SAS/STAT User's Guide*, volume 2. SAS Institute, Inc., Cary NC, fourth edition, 1990. ◊315
541. W. Sauerbrei and M. Schumacher. A bootstrap resampling procedure for model building: Application to the Cox regression model. *Stat Med*, 11:2093–2109, 1992. ◊70, 113, 177
542. J. L. Schafer and J. W. Graham. Missing data: Our view of the state of the art. *Psych Meth*, 7:147–177, 2002. ◊58
543. D. E. Schaebel, R. A. Wolfe, and R. M. Merion. Estimating the effect of a time-dependent treatment by levels of an internal time-dependent covariate: Application to the contrast between liver wait-list and posttransplant mortality. *J Am Stat Assoc*, 104(485):49–59, 2009. ◊518
544. M. Schemper. Analyses of associations with censored data by generalized Mantel and Breslow tests and generalized Kendall correlation. *Biometrical J*, 26:309–318, 1984. ◊518
545. M. Schemper. Non-parametric analysis of treatment-covariate interaction in the presence of censoring. *Stat Med*, 7:1257–1266, 1988. ◊41
546. M. Schemper. The explained variation in proportional hazards regression (correction in 81:631, 1994). *Biometrika*, 77:216–218, 1990. ◊505, 508
547. M. Schemper. Cox analysis of survival data with non-proportional hazard functions. *The Statistician*, 41:445–455, 1992. ◊518
548. M. Schemper. Further results on the explained variation in proportional hazards regression. *Biometrika*, 79:202–204, 1992. ◊505
549. M. Schemper. The relative importance of prognostic factors in studies of survival. *Stat Med*, 12:2377–2382, 1993. ◊215, 505
550. M. Schemper. Predictive accuracy and explained variation. *Stat Med*, 22:2299–2308, 2003. ◊519
551. M. Schemper and G. Heinze. Probability imputation revisited for prognostic factor studies. *Stat Med*, 16:73–80, 1997. ◊52, 177
552. M. Schemper and R. Henderson. Predictive accuracy and explained variation in Cox regression. *Biometrics*, 56:249–255, 2000. ◊518
553. M. Schemper and T. L. Smith. Efficient evaluation of treatment effects in the presence of missing covariate values. *Stat Med*, 9:777–784, 1990. ◊52
554. M. Schemper and J. Stare. Explained variation in survival analysis. *Stat Med*, 15:1999–2012, 1996. ◊215, 519
555. M. Schmid and S. Potapov. A comparison of estimators to evaluate the discriminatory power of time-to-event models. *Stat Med*, 31(23):2588–2609, 2012. ◊519
556. C. Schmoor, K. Ulm, and M. Schumacher. Comparison of the Cox model and the regression tree procedure in analysing a randomized clinical trial. *Stat Med*, 12:2351–2366, 1993. ◊41
557. D. Schoenfeld. Partial residuals for the proportional hazards regression model. *Biometrika*, 69:239–241, 1982. ◊314, 498, 499, 516
558. D. A. Schoenfeld. Sample size formulae for the proportional hazards regression model. *Biometrics*, 39:499–503, 1983. ◊513
559. G. Schulgen, B. Lausen, J. Olsen, and M. Schumacher. Outcome-oriented cut-points in quantitative exposure. *Am J Epi*, 120:172–184, 1994. ◊19, 20
560. G. Schwarz. Estimating the dimension of a model. *Ann Stat*, 6:461–464, 1978. ◊214
561. S. C. Scott, M. S. Goldberg, and N. E. Mayo. Statistical assessment of ordinal outcomes in comparative studies. *J Clin Epi*, 50:45–55, 1997. ◊324
562. M. R. Segal. Regression trees for censored data. *Biometrics*, 44:35–47, 1988. ◊41
563. S. Senn. Change from baseline and analysis of covariance revisited. *Stat Med*, 25:4334–4344, 2006. ◊159, 160

564. S. Senn and S. Julious. Measurement in clinical trials: A neglected issue for statisticians? (with discussion). *Stat Med*, 28:3189–3225, 2009. ◊313
565. J. Shao. Linear model selection by cross-validation. *J Am Stat Assoc*, 88:486–494, 1993. ◊100, 113, 122
566. J. Shao and R. R. Sitter. Bootstrap for imputed survey data. *J Am Stat Assoc*, 91:1278–1288, 1996. ◊54
567. X. Shen, H. Huang, and J. Ye. Inference after model selection. *J Am Stat Assoc*, 99:751–762, 2004. ◊102
568. Y. Shen and P. F. Thall. Parametric likelihoods for multiple non-fatal competing risks and death. *Stat Med*, 17:999–1015, 1998. ◊421
569. J. Siddique. Multiple imputation using an iterative hot-deck with distance-based donor selection. *Stat Med*, 27:83–102, 2008. ◊58
570. R. Simon and R. W. Makuch. A non-parametric graphical representation of the relationship between survival and the occurrence of an event: Application to responder versus non-responder bias. *Stat Med*, 3:35–44, 1984. ◊401, 420
571. J. S. Simonoff. The “Unusual Episode” and a second statistics course. *J Stat Edu*, 5(1), 1997. Online journal at www.amstat.org/publications/jse/v5n1/~simonoff.html. ◊291
572. S. L. Simpson, L. J. Edwards, K. E. Muller, P. K. Sen, and M. A. Styner. A linear exponent AR(1) family of correlation structures. *Stat Med*, 29:1825–1838, 2010. ◊148
573. J. C. Sinclair and M. B. Bracken. Clinically useful measures of effect in binary analyses of randomized trials. *J Clin Epi*, 47:881–889, 1994. ◊272
574. J. D. Singer and J. B. Willett. Modeling the days of our lives: Using survival analysis when designing and analyzing longitudinal studies of duration and the timing of events. *Psych Bull*, 110:268–290, 1991. ◊420
575. L. A. Sleeper and D. P. Harrington. Regression splines in the Cox model with application to covariate effects in liver disease. *J Am Stat Assoc*, 85:941–949, 1990. ◊23, 40
576. A. F. M. Smith and D. J. Spiegelhalter. Bayes factors and choice criteria for linear models. *J Roy Stat Soc B*, 42:213–220, 1980. ◊214
577. L. R. Smith, F. E. Harrell, and L. H. Muhlbaier. Problems and potentials in modeling survival. In M. L. Grady and H. A. Schwartz, editors, *Medical Effectiveness Research Data Methods (Summary Report)*, AHCPR Pub. No. 92-0056, pages 151–159. US Dept. of Health and Human Services, Agency for Health Care Policy and Research, Rockville, MD, 1992. ◊72
578. P. L. Smith. Splines as a useful and convenient statistical tool. *Am Statistician*, 33:57–62, 1979. ◊40
579. R. H. Somers. A new asymmetric measure of association for ordinal variables. *Am Soc Rev*, 27:799–811, 1962. ◊257, 505
580. A. Spanos, F. E. Harrell, and D. T. Durack. Differential diagnosis of acute meningitis: An analysis of the predictive value of initial observations. *JAMA*, 262:2700–2707, 1989. ◊266, 267, 268
581. I. Spence and R. F. Garrison. A remarkable scatterplot. *Am Statistician*, 47:12–19, 1993. ◊91
582. D. J. Spiegelhalter. Probabilistic prediction in patient management and clinical trials. *Stat Med*, 5:421–433, 1986. ◊97, 101, 115, 116, 523
583. D. M. Stablein, W. H. Carter, and J. W. Novak. Analysis of survival data with nonproportional hazard functions. *Controlled Clin Trials*, 2:149–159, 1981. ◊500
584. N. Stallard. Simple tests for the external validation of mortality prediction scores. *Stat Med*, 28:377–388, 2009. ◊237
585. J. Stare, F. E. Harrell, and H. Heinzl. BJ: An S-PLUS program to fit linear regression models to censored data using the Buckley and James method. *Comp Meth Prog Biomed*, 64:45–52, 2001. ◊447

586. E. W. Steyerberg. *Clinical Prediction Models*. Springer, New York, 2009. ◊viii
587. E. W. Steyerberg, S. E. Bleeker, H. A. Moll, D. E. Grobbee, and K. G. M. Moons. Internal and external validation of predictive models: A simulation study of bias and precision in small samples. *Journal of Clinical Epi*, 56(5):441–447, May 2003. ◊123
588. E. W. Steyerberg, P. M. M. Bossuyt, and K. L. Lee. Clinical trials in acute myocardial infarction: Should we adjust for baseline characteristics? *Am Heart J*, 139:745–751, 2000. Editorial, pp. 761–763. ◊4, 231
589. E. W. Steyerberg, M. J. C. Eijkemans, F. E. Harrell, and J. D. F. Habbema. Prognostic modelling with logistic regression analysis: A comparison of selection and estimation methods in small data sets. *Stat Med*, 19:1059–1079, 2000. ◊69, 100, 286
590. E. W. Steyerberg, M. J. C. Eijkemans, F. E. Harrell, and J. D. F. Habbema. Prognostic modeling with logistic regression analysis: In search of a sensible strategy in small data sets. *Med Decis Mak*, 21:45–56, 2001. ◊100, 271
591. E. W. Steyerberg, F. E. Harrell, G. J. J. M. Borsboom, M. J. C. Eijkemans, Y. Vergouwe, and J. D. F. Habbema. Internal validation of predictive models: Efficiency of some procedures for logistic regression analysis. *J Clin Epi*, 54:774–781, 2001. ◊115
592. E. W. Steyerberg, A. J. Vickers, N. R. Cook, T. Gerds, M. Gonen, N. Obuchowski, M. J. Pencina, and M. W. Kattan. Assessing the performance of prediction models: a framework for traditional and novel measures. *Epi (Cambridge, Mass.)*, 21(1):128–138, Jan. 2010. ◊101
593. C. J. Stone. Comment: Generalized additive models. *Statistical Sci*, 1:312–314, 1986. ◊26, 28
594. C. J. Stone, M. H. Hansen, C. Kooperberg, and Y. K. Truong. Polynomial splines and their tensor products in extended linear modeling (with discussion). *Ann Stat*, 25:1371–1470, 1997. ◊420, 450
595. C. J. Stone and C. Y. Koo. Additive splines in statistics. In *Proceedings of the Statistical Computing Section ASA*, pages 45–48, Washington, DC, 1985. ◊24, 28, 41
596. D. Strauss and R. Shavelle. An extended Kaplan–Meier estimator and its applications. *Stat Med*, 17:971–982, 1998. ◊416
597. S. Suissa and L. Blais. Binary regression with continuous outcomes. *Stat Med*, 14:247–255, 1995. ◊11, 19
598. G. Sun, T. L. Shook, and G. L. Kay. Inappropriate use of bivariable analysis to screen risk factors for use in multivariable analysis. *J Clin Epi*, 49:907–916, 1996. ◊72
599. B. Tai, D. Machin, I. White, and V. GebSKI. Competing risks analysis of patients with osteosarcoma: a comparison of four different approaches. *Stat Med*, 20:661–684, 2001. ◊420
600. J. M. G. Taylor, A. L. Siqueira, and R. E. Weiss. The cost of adding parameters to a model. *J Roy Stat Soc B*, 58:593–607, 1996. ◊101
601. R. D. C. Team. *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria, 2015. ISBN 3-900051-07-0. ◊127
602. H. T. Thaler. Nonparametric estimation of the hazard ratio. *J Am Stat Assoc*, 79:290–293, 1984. ◊518
603. P. F. Thall and J. M. Lachin. Assessment of stratum-covariate interactions in Cox’s proportional hazards regression model. *Stat Med*, 5:73–83, 1986. ◊482
604. T. Therneau and P. Grambsch. *Modeling Survival Data: Extending the Cox Model*. Springer-Verlag, New York, 2000. ◊420, 447, 478, 517
605. T. M. Therneau, P. M. Grambsch, and T. R. Fleming. Martingale-based residuals for survival models. *Biometrika*, 77:216–218, 1990. ◊197, 413, 487, 493, 494, 504

606. T. M. Therneau and S. A. Hamilton. rhDNase as an example of recurrent event analysis. *Stat Med*, 16:2029–2047, 1997. ◊420, 421
607. R. Tibshirani. Estimating transformations for regression via additivity and variance stabilization. *J Am Stat Assoc*, 83:394–405, 1988. ◊391
608. R. Tibshirani. Regression shrinkage and selection via the lasso. *J Roy Stat Soc B*, 58:267–288, 1996. ◊71, 215, 356
609. R. Tibshirani. The lasso method for variable selection in the Cox model. *Stat Med*, 16:385–395, 1997. ◊71, 356
610. R. Tibshirani and K. Knight. Model search and inference by bootstrap “bumping”. Technical report, Department of Statistics, University of Toronto, 1997. <http://www-stat.stanford.edu/tibs>. Presented at the Joint Statistical Meetings, Chicago, August 1996. ◊xii, 214
611. R. Tibshirani and K. Knight. The covariance inflation criterion for adaptive model selection. *J Roy Stat Soc B*, 61:529–546, 1999. ◊11, 123
612. N. H. Timm. The estimation of variance-covariance and correlation matrices from incomplete data. *Psychometrika*, 35:417–437, 1970. ◊52
613. T. Tjur. Coefficients of determination in logistic regression models—A new proposal: The coefficient of discrimination. *Am Statistician*, 63(4):366–372, 2009. ◊257, 272
614. W. Y. Tsai, N. P. Jewell, and M. C. Wang. A note on the product limit estimator under right censoring and left truncation. *Biometrika*, 74:883–886, 1987. ◊420
615. A. A. Tsiatis. A large sample study of Cox’s regression model. *Ann Stat*, 9:93–108, 1981. ◊485
616. B. W. Turnbull. Nonparametric estimation of a survivorship function with doubly censored data. *J Am Stat Assoc*, 69:169–173, 1974. ◊420
617. J. Twisk, M. de Boer, W. de Vente, and M. Heymans. Multiple imputation of missing values was not necessary before performing a longitudinal mixed-model analysis. *J Clin Epi*, 66(9):1022–1028, 2013. ◊58
618. H. Uno, T. Cai, M. J. Pencina, R. B. D’Agostino, and L. J. Wei. On the C -statistics for evaluating overall adequacy of risk prediction procedures with censored survival data. *Stat Med*, 30:1105–1117, 2011. ◊519
619. Ü. Uzunoğullari and J.-L. Wang. A comparison of hazard rate estimators for left truncated and right censored data. *Biometrika*, 79:297–310, 1992. ◊420
620. W. Vach. *Logistic Regression with Missing Values in the Covariates*, volume 86 of *Lecture Notes in Statistics*. Springer-Verlag, New York, 1994. ◊59
621. W. Vach. Some issues in estimating the effect of prognostic factors from incomplete covariate data. *Stat Med*, 16:57–72, 1997. ◊52, 59
622. W. Vach and M. Blettner. Logistic regression with incompletely observed categorical covariates—Investigating the sensitivity against violation of the missing at random assumption. *Stat Med*, 14:1315–1329, 1995. ◊59
623. W. Vach and M. Blettner. Missing Data in Epidemiologic Studies. In *Ency of Biostatistics*, pages 2641–2654. Wiley, New York, 1998. ◊52, 58, 59
624. W. Vach and M. Schumacher. Logistic regression with incompletely observed categorical covariates: A comparison of three approaches. *Biometrika*, 80:353–362, 1993. ◊59
625. M. G. Valsecchi, D. Silvestri, and P. Sasieni. Evaluation of long-term survival: Use of diagnostics and robust estimators with Cox’s proportional hazards model. *Stat Med*, 15:2763–2780, 1996. ◊518
626. S. van Buuren, H. C. Boshuizen, and D. L. Knook. Multiple imputation of missing blood pressure covariates in survival analysis. *Stat Med*, 18:681–694, 1999. ◊58
627. S. van Buuren, J. P. L. Brand, C. G. M. Groothuis-Oudshoorn, and D. B. Rubin. Fully conditional specification in multivariate imputation. *J Stat Computation Sim*, 76(12):1049–1064, 2006. ◊55

628. G. J. M. G. van der Heijden, Donders, T. Stijnen, and K. G. M. Moons. Imputation of missing values is superior to complete case analysis and the missing-indicator method in multivariable diagnostic research: A clinical example. *J Clin Epi*, 59:1102–1109, 2006. ◊48, 49
629. T. van der Ploeg, P. C. Austin, and E. W. Steyerberg. Modern modelling techniques are data hungry: a simulation study for predicting dichotomous endpoints. *BMC Medical Research Methodology*, 14(1):137+, Dec. 2014. ◊41, 100
630. M. J. van Gorp, E. W. Steyerberg, M. Kallewaard, and Y. van der Graaf. Clinical prediction rule for 30-day mortality in Björk-Shiley convexo-concave valve replacement. *J Clin Epi*, 56:1006–1012, 2003. ◊122
631. H. C. van Houwelingen and J. Thorogood. Construction, validation and updating of a prognostic model for kidney graft survival. *Stat Med*, 14:1999–2008, 1995. ◊100, 101, 123, 215
632. J. C. van Houwelingen and S. le Cessie. Logistic regression, a review. *Statistica Neerlandica*, 42:215–232, 1988. ◊271
633. J. C. van Houwelingen and S. le Cessie. Predictive value of statistical models. *Stat Med*, 9:1303–1325, 1990. ◊77, 101, 113, 115, 123, 204, 214, 215, 258, 259, 273, 508, 509, 518
634. W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S-Plus*. Springer-Verlag, New York, third edition, 1999. ◊101
635. W. N. Venables and B. D. Ripley. *Modern Applied Statistics with S*. Springer-Verlag, New York, fourth edition, 2003. ◊xi, 127, 129, 143, 359
636. D. J. Venzon and S. H. Moolgavkar. A method for computing profile-likelihood-based confidence intervals. *Appl Stat*, 37:87–94, 1988. ◊214
637. G. Verbeke and G. Molenberghs. *Linear Mixed Models for Longitudinal Data*. Springer, New York, 2000. ◊143
638. Y. Vergouwe, E. W. Steyerberg, M. J. C. Eijkemans, and J. D. F. Habbema. Substantial effective sample sizes were required for external validation studies of predictive logistic regression models. *J Clin Epi*, 58:475–483, 2005. ◊122
639. P. Verweij and H. C. van Houwelingen. Penalized likelihood in Cox regression. *Stat Med*, 13:2427–2436, 1994. ◊77, 209, 210, 211, 215
640. P. J. M. Verweij and H. C. van Houwelingen. Cross-validation in survival analysis. *Stat Med*, 12:2305–2314, 1993. ◊100, 123, 207, 215, 509, 518
641. P. J. M. Verweij and H. C. van Houwelingen. Time-dependent effects of fixed covariates in Cox regression. *Biometrics*, 51:1550–1556, 1995. ◊209, 211, 501
642. A. J. Vickers. Decision analysis for the evaluation of diagnostic tests, prediction models, and molecular markers. *Am Statistician*, 62(4):314–320, 2008. ◊5
643. S. K. Vines. Simple principal components. *Appl Stat*, 49:441–451, 2000. ◊101
644. E. Vittinghoff and C. E. McCulloch. Relaxing the rule of ten events per variable in logistic and Cox regression. *Am J Epi*, 165:710–718, 2006. ◊100
645. P. T. von Hippel. Regression with missing ys: An improved strategy for analyzing multiple imputed data. *Soc Meth*, 37(1):83–117, 2007. ◊47
646. H. Wainer. Finding what is not there through the unfortunate binning of results: The Mendel effect. *Chance*, 19(1):49–56, 2006. ◊19, 20
647. S. H. Walker and D. B. Duncan. Estimation of the probability of an event as a function of several independent variables. *Biometrika*, 54:167–178, 1967. ◊14, 220, 311, 313
648. A. R. Walter, A. R. Feinstein, and C. K. Wells. Coding ordinal independent variables in multiple regression analyses. *Am J Epi*, 125:319–323, 1987. ◊39
649. A. Wang and E. A. Gehan. Gene selection for microarray data analysis using principal component analysis. *Stat Med*, 24:2069–2087, 2005. ◊101
650. M. Wang and S. Chang. Nonparametric estimation of a recurrent survival function. *J Am Stat Assoc*, 94:146–153, 1999. ◊421
651. R. Wang, J. Sedransk, and J. H. Jinn. Secondary data analysis when there are missing observations. *J Am Stat Assoc*, 87:952–961, 1992. ◊53

652. Y. Wang and J. M. G. Taylor. Inference for smooth curves in longitudinal data with application to an AIDS clinical trial. *Stat Med*, 14:1205–1218, 1995. ◊215
653. Y. Wang, G. Wahba, C. Gu, R. Klein, and B. Klein. Using smoothing spline ANOVA to examine the relation of risk factors to the incidence and progression of diabetic retinopathy. *Stat Med*, 16:1357–1376, 1997. ◊41
654. Y. Wax. Collinearity diagnosis for a relative risk regression analysis: An application to assessment of diet-cancer relationship in epidemiological studies. *Stat Med*, 11:1273–1287, 1992. ◊79, 138, 255
655. L. J. Wei, D. Y. Lin, and L. Weissfeld. Regression analysis of multivariate incomplete failure time data by modeling marginal distributions. *J Am Stat Assoc*, 84:1065–1073, 1989. ◊417
656. R. E. Weiss. The influence of variable selection: A Bayesian diagnostic perspective. *J Am Stat Assoc*, 90:619–625, 1995. ◊100
657. S. Wellek. A log-rank test for equivalence of two survivor functions. *Biometrics*, 49:877–881, 1993. ◊450
658. T. L. Wenger, F. E. Harrell, K. K. Brown, S. Lederman, and H. C. Strauss. Ventricular fibrillation following canine coronary reperfusion: Different outcomes with pentobarbital and α -chloralose. *Can J Phys Pharm*, 62:224–228, 1984. ◊266
659. H. White. A heteroskedasticity-consistent covariance matrix estimator and a direct test for heteroskedasticity. *Econometrica*, 48:817–838, 1980. ◊196
660. I. R. White and J. B. Carlin. Bias and efficiency of multiple imputation compared with complete-case analysis for missing covariate values. *Stat Med*, 29:2920–2931, 2010. ◊59
661. I. R. White and P. Royston. Imputing missing covariate values for the Cox model. *Stat Med*, 28:1982–1998, 2009. ◊54
662. I. R. White, P. Royston, and A. M. Wood. Multiple imputation using chained equations: Issues and guidance for practice. *Stat Med*, 30(4):377–399, 2011. ◊53, 54, 58
663. A. Whitehead, R. Z. Omar, J. P. T. Higgins, E. Savaluny, R. M. Turner, and S. G. Thompson. Meta-analysis of ordinal outcomes using individual patient data. *Stat Med*, 20:2243–2260, 2001. ◊324
664. J. Whitehead. Sample size calculations for ordered categorical data. *Stat Med*, 12:2257–2271, 1993. See letter to editor SM 15:1065-6 for binary case; see errata in SM 13:871 1994; see kol95com, jul96sam. ◊2, 73, 313, 324
665. J. Whittaker. Model interpretation from the additive elements of the likelihood function. *Appl Stat*, 33:52–64, 1984. ◊205, 207
666. A. S. Whittemore and J. B. Keller. Survival estimation using splines. *Biometrics*, 42:495–506, 1986. ◊420
667. H. Wickham. *ggplot2: elegant graphics for data analysis*. Springer, New York, 2009. ◊xi
668. R. E. Wiegand. Performance of using multiple stepwise algorithms for variable selection. *Stat Med*, 29:1647–1659, 2010. ◊100
669. A. R. Willan, W. Ross, and T. A. MacKenzie. Comparing in-patient classification systems: A problem of non-nested regression models. *Stat Med*, 11:1321–1331, 1992. ◊205, 215
670. A. Winnett and P. Sasieni. A note on scaled Schoenfeld residuals for the proportional hazards model. *Biometrika*, 88:565–571, 2001. ◊518
671. A. Winnett and P. Sasieni. Iterated residuals and time-varying covariate effects in Cox regression. *J Roy Stat Soc B*, 65:473–488, 2003. ◊518
672. D. M. Witten and R. Tibshirani. Testing significance of features by lassoed principal components. *Ann Appl Stat*, 2(3):986–1012, 2008. ◊175
673. A. M. Wood, I. R. White, and S. G. Thompson. Are missing outcome data adequately handled? A review of published randomized controlled trials in major medical journals. *Clin Trials*, 1:368–376, 2004. ◊58

674. S. N. Wood. *Generalized Additive Models: An Introduction with R*. Chapman & Hall/CRC, Boca Raton, FL, 2006. ISBN 9781584884743. ◊90
675. C. F. J. Wu. Jackknife, bootstrap and other resampling methods in regression analysis. *Ann Stat*, 14(4):1261–1350, 1986. ◊113
676. Y. Xiao and M. Abrahamowicz. Bootstrap-based methods for estimating standard errors in Cox’s regression analyses of clustered event times. *Stat Med*, 29:915–923, 2010. ◊213
677. Y. Xie. *knitr: A general-purpose package for dynamic report generation in R*, 2013. R package version 1.5. ◊xi, 138
678. J. Ye. On measuring and correcting the effects of data mining and model selection. *J Am Stat Assoc*, 93:120–131, 1998. ◊10
679. T. W. Yee and C. J. Wild. Vector generalized additive models. *J Roy Stat Soc B*, 58:481–493, 1996. ◊324
680. F. W. Young, Y. Takane, and J. de Leeuw. The principal components of mixed measurement level multivariate data: An alternating least squares method with optimal scaling features. *Psychometrika*, 43:279–281, 1978. ◊81
681. R. M. Yucel and A. M. Zaslavsky. Using calibration to improve rounding in imputation. *Am Statistician*, 62(2):125–129, 2008. ◊56
682. H. Zhang. Classification trees for multiple binary responses. *J Am Stat Assoc*, 93:180–193, 1998. ◊41
683. H. Zhang, T. Holford, and M. B. Bracken. A tree-based method of analysis for prospective studies. *Stat Med*, 15:37–49, 1996. ◊41
684. B. Zheng and A. Agresti. Summarizing the predictive power of a generalized linear model. *Stat Med*, 19:1771–1781, 2000. ◊215, 273
685. X. Zheng and W. Loh. Consistent variable selection in linear models. *J Am Stat Assoc*, 90:151–156, 1995. ◊214
686. H. Zhou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *J Comp Graph Stat*, 15:265–286, 2006. ◊101
687. X. Zhou. Effect of verification bias on positive and negative predictive values. *Stat Med*, 13:1737–1745, 1994. ◊328
688. X. Zhou, G. J. Eckert, and W. M. Tierney. Multiple imputation in public health research. *Stat Med*, 20:1541–1549, 2001. ◊59
689. H. Zou, T. Hastie, and R. Tibshirani. On the “degrees of freedom” of the lasso. *Ann Stat*, 35:2173–2192, 2007. ◊11
690. H. Zou and M. Yuan. Composite quantile regression and the oracle model selection theory. *Ann Stat*, 36(3):1108–1126, 2008. ◊361
691. D. M. Zucker. The efficiency of a weighted log-rank test under a percent error misspecification model for the log hazard ratio. *Biometrics*, 48:893–899, 1992. ◊518

Index

Entries in **this font** are names of software components. Page numbers in **bold** denote the most comprehensive treatment of the topic.

Symbols

D_{xy} , 105, 142, **257**, 257–259, 269, 284, 318, 461, 505, 529

censored data, 505, 517

R^2 , 110, 111, 206, 272, 390, 391

adjusted, 74, 77, 105

generalized, 207

significant difference in, 215

c index, 93, 100, 105, 142, **257**, 257, 259, 318, 505, 517

censored data, 505

generalized, 318, 505

HbA_{1c}, 365

15:1 rule, 72, 100

A

Aalen survival function estimator,
see survival function

abs.error.pred, 102

accelerated failure time, *see*
model

accuracy, 104, 111, 113, 114, 210,
354, 446

g -index, 105

absolute, 93, 102

apparent, 114, 269, 529

approximation, 119, 275,
287, 348, 469

bias-corrected, 100, 109,
114, 115, 141, 391, 529

calibration, **72–78**,

88, 92, 93, 105, 111, 115, 141,
236, 237, 259, 260,

264, 269, 271, 284, 301, 322,
446, 467, 506

discrimination, 72, 92, 93,

105, **111**, 111, 257, 259,

269, 284, 287, 318, 331, 346,
467, 505, 506, 508

future, 211

index, 122, 123, 141

ACE, 82, 176, 179, **390**, 391, 392

ace, 176, 392

acepack package, 176, 392

actuarial survival, 410

adequacy index, 207

AIC, 28, 69, 78, 88, 172, **204**, 204,
210, 211, 214, 215,

240, 241, 269, 275, 277, 332,
374, 375

- AIC, 134, 135, 277
- Akaike information criterion, *see* AIC
- analysis of covariance, *see* ANOCOVA
- ANOCOVA, 16, 223, 230, 447
- ANOVA, 13, 32, 75, 230, 235, 317, 447, 480, 531
- anova*, 65, 127, 133, 134, 136, 149, 155, 278, 302, 306, 336, 342, 346, 464, 466
- anova.gls*, 149
- areg.boot*, 392–394
- aregImpute*, 51, 53–56, 59, 304, 305
- Arjas plot, 495
- asis*, 132, 133
- assumptions
 - accelerated failure time, 436, 437, 458
 - additivity, 37, 248
 - continuation ratio, **320**, 321, 338
 - correlation pattern, 148, 153
 - distributional, 39, 97, 148, 317, 446, 525
 - linearity, 21–26
 - ordinality, **312**, 319, 333, 340
 - proportional hazards, 429, **494–503**
 - proportional odds, **313**, 315, 317, 336, 362
- AVAS, 390–392
 - case study, 393–398
- avas*, 392, 394, 395
- B**
- B-spline, *see* spline function
- battery reduction, 87
- Bayesian modeling, 71, 209, 215
- BIC, 211, 214, 269
- binary response, *see* response
- bj*, 131, 135, 447, 449
- bootcov*, 134–136, 198–202, 319
- bootkm*, 419
- bootstrap, 106–109, 114–116
 - .632, 115, 123
 - adjusting for imputation, 53
 - approximate Bayesian, 50
 - basic, 202, 203
 - BCa, 202, 203
 - cluster, 135, 197, 199, 213
 - conditional, 115, 122, 197
 - confidence intervals, *see* confidence intervals, 199
 - covariance matrix, 135, 198
 - density, 107, 136
 - distribution, 201
 - estimating shrinkage, 77, 115
 - model uncertainty, 11, 113, 304
 - overfitting correction, 112, 114, 115, 257, 391
 - ranks, 117
 - variable selection, 70, 97, 113, 177, 260, 275, 282, 286
- bpplot*, 134
- Breslow survival function
 - estimator, *see* survival function
- Brier score, 142, 237, 257–259, 271, 318
- C**
- CABG, 484
- calibrate*, 135, **141**, 269, 271, 284, 300, 319, 323, 355, 450, 467, 517
- calibration, *see* accuracy
- caliper matching, 372
- cancor*, 141
- canonical correlation, 141
- canonical variate, 82, 83, 129, 141, 167, 169, 393
- CART, *see* recursive partitioning
- casewise deletion, *see* missing data
- categorical predictor, *see* predictor
- categorization of continuous variable, 8, **18–21**

- catg, 132, 133
 - causal inference, 103
 - cause removal, 414
 - censoring, 401–402, 406, 424
 - informative, 402, 414, 415, 420
 - interval, 401, 418, 420
 - left, 401
 - right, 402, 418
 - type I, 401
 - type II, 402
 - ciapower, 513
 - classification, 4, 6
 - classifier, 4, 6
 - clustered data, 197, 417
 - clustering
 - hierarchical, 129, 166, 330
 - variable, 81, 101, 175, 355
 - ClustOfVar, 101
 - coef, 134
 - coefficient of discrimination, *see*
 - accuracy
 - collinearity, 78–79
 - competing risks, 414, 420
 - concordance probability, *see c*
 - index
 - conditional logistic model, *see*
 - logistic model
 - conditional probability, 320, 404, 476, 484
 - confidence intervals, 10, 30, 35, 64, 66, 96, 136, 185, 198, 273, 282, 391
 - bootstrap, 107, 109, 119, 122, 135, 149, 199, 201–203, 214, 217
 - coverage, 35, 198, 199, 389
 - simultaneous, 136, 199, 202, 214, 420, 517
 - confounding, 31, 103, 231
 - confplot, 214
 - contingency table, 195, 228, 230, 235
 - contrast, *see* hypothesis test
 - contrast, 134, 136, 192, 193, 198, 199
 - convergence, 193, 264
 - coronary artery disease, 48, 207, 240, 245, 252, 492, 497
 - correlation structures, 147, 148
 - correspondence analysis, 81, 129
 - cost-effectiveness, 4
 - Cox model, 362, 375, 392, 475–517
 - case study, 521–531
 - data reduction example, 172
 - multiple imputation, 54
 - cox.zph, 499, 516, 517, 526
 - coxph, 131, 422, 513
 - cph, 131, 133, 135, 172, 422, 448, 513, 513, 514, 516, 517
 - cpower, 513
 - cr.setup, 323, 340, 354
 - cross-validation, *see* validation of
 - model
 - cubic spline, *see* spline function
 - cumcategory, 357
 - cumulative hazard function, *see*
 - hazard function
 - cumulative probability model, 359, 361–363, 370, 371
 - cut2, 129, 133, 334, 419
 - cutpoint, 21
- ## D
- data reduction, 79–88, 275
 - case study 1, 161–177
 - case study 2, 277
 - case study 3, 329–333
 - data-splitting, *see* validation of
 - model
 - data.frame, 309
 - datadist, 130, 130, 138, 292, 463
 - datasets, 535
 - cdystonia, 149
 - cervical dystonia, 149
 - diabetes, 317
 - meningitis, 266, 267
 - NHANES, 365
 - prostate, 161, 275, 521
 - SUPPORT, 59, 453

Titanic, 291
 degrees of freedom, 193
 effective, 30, 41, 77, 96, 136, 210, 269
 generalized, 10
 phantom, 35, 111
 delayed entry, 401
 delta method, 439
 describe, 129, 291, 453
 deviance, 236, 449, 487, 516
 DFBETA, 91
 DFBETAS, 91
 DFFIT, 91
 DFFITS, 91
 diabetes, *see* datasets, 365
 difference in predictions, 192, 201
 dimensionality, 88
 discriminant analysis, 220, 230, 272
 discrimination, *see* accuracy, *see* accuracy
 distribution, 317
 t, 186
 binomial, 73, 181, 194, 235
 Cauchy, 362
 exponential, 142, 407, 408, 425, 427, 451
 extreme value, 362, 363, 427, 437
 Gumbel, 362, 363
 log-logistic, 9, 423, 427, 440, 442, 503
 log-normal, 9, 106, 391, 423, 427, 442, 463, 464
 normal, 187
 Weibull, 39, 408, 408, 420, 426, 432–437, 444, 448
 dose-response, 523
 doubly nonlinear, 131
 drop-in, 513
 dropouts, 143
 dummy variable, 1, *see* indicator variable, 75, 129, 130, 209, 210

E

economists, 71
 effective.df, 134, 136, 345, 346
 Emax, 353
 epidemiology, 38
 estimation, 2, 98, 104
 estimator
 Buckley–James, 447, 449
 maximum likelihood, 181
 mean, 362
 penalized, *see* maximum likelihood, 175
 quantile, 362
 self-consistent, 525
 smearing, 392, 393
 explained variation, 273
 exponential distribution, *see* distribution
 ExProb, 135
 external validation, *see* validation of model

F

failure time, 399
 fastbw, 133, 134, 137, 280, 286, 351, 469
 feature selection, 94
 financial data, 3
 fit.mult.impute, 54, 306
 Fleming–Harrington survival function estimator, *see* survival function
 formula, 134
 fractional polynomial, 40
 Function, 134, 135, 138, 149, 310, 395
 functions, generating R code, 395

G

GAM, *see* generalized additive model, *see* generalized additive model
 gam package, 390
 GDF, *see* degrees of freedom
 GEE, 147

- Gehan–Wilcoxon test, *see* hypothesis test
- gendata, 134, 136
- generalized additive model, 29, 41, 138, 142, 390
 case study, 393–398
- getHdata, 59, 178, 535
- ggplot, 134
- ggplot2 package, xi, 134, 294
- gIndex, 105
- glht, 199
- Glm, 131, 135, 271
- glm, 131, 141, 271
- Gls, 131, 135, 149
- gls, 131, 149
- goodness of fit, 236, 269, 427, 440, 458
- Greenwood’s formula, *see* survival function
- groupkm, 419
- H**
- hare, 450
- hat matrix, 91
- Hazard, 135, 448
- hazard function, 135, 362, 375, 400, 402, 405, 409, 427, 475, 476
 bathtub, 408
 cause-specific, 414, 415
 cumulative, 402–409
- hazard ratio, 429–431, 433, 478, 479, 481
 interval-specific, 495–497, 502
- hazard.ratio.plot, 517
- hclust, 129
- heft, 419
- heterogeneity, unexplained, 4, 231, 400
- histSpikeg, 294
- Hmisc package, xi, 129, 133, 137, 167, 176, 273, 277, 294, 304, 319, 357, 392, 418, 458, 463, 513, 536
- hoeffd, 129
- Hoeffding D , 129, 166, 458
- Hosmer–Lemeshow test, 236, 237
- Hotelling test, *see* hypothesis test
- Huber–White estimator, 196
- hypothesis test, 1, 18, 32, 99
 additivity, 37, 248
 association, 2, 18, 32, 43, 66, 129, 235, 338, 486
 contrast, 157, 192, 193, 198
 equal slopes, 315, 321, 322, 338, 339, 458, 460, 495
 exponentiality, 408, 426
 Gehan–Wilcoxon, 505
 global, 69, 97, 189, 205, 230, 232, 342, 526
 Hotelling, 230
 independence, 129, 166
 Kruskal–Wallis, 2, 66, 129
 linearity, 18, 32, 35, 36, 39, 42, 66, 91, 238
 log-rank, 41, 363, 422, 475, 486, 513, 518
 Mantel–Haenszel, 486
 normal scores, 364
 partial, 190
 Pearson χ^2 , 195, 235
 robust, 9, 81, 311
 Van der Waerden, 364
 Wilcoxon, 1, 73, 129, 230, 257, 311, 313, 325, 363, 364
- I**
- ignorable nonresponse, *see* missing data
- imbalances, baseline, 400
- improveProb, 142
- imputation, 47–57, 83
 chained equations, 55, 304
 model for, 49, 50, 50–52, 59, 84, 129
 multiple, 47, 53, 54, 54–56, 95, 129, 304, 382, 537
 censored data, 54

- predictive mean matching, [51](#), [52](#), [55](#)
 - single, [52](#), [56](#), [57](#), [138](#), [171](#), [275](#), [276](#), [334](#)
 - `impute`, [129](#), [135](#), [138](#), [171](#), [276](#), [277](#), [334](#), [461](#)
 - incidence
 - crude, [416](#)
 - cumulative, [415](#)
 - incomplete principal component regression, [170](#), [275](#)
 - indicator variable, [16](#), [17](#), [38](#), [39](#)
 - infinite regression coefficient, [234](#)
 - influential observations, [90–92](#), [116](#), [255](#), [256](#), [269](#), [504](#)
 - information function, [182](#), [183](#)
 - information matrix, [79](#), [188](#), [189](#), [191](#), [196](#), [208](#), [211](#), [232](#), [346](#)
 - informative missing, *see* missing data
 - interaction, [16](#), [36](#), [375](#)
 - interquartile-range effect, [104](#), [136](#)
 - intracluster correlation, [135](#), [141](#), [197](#), [417](#)
 - isotropic correlation structure, *see* correlation structures
- J**
- jackknife, [113](#), [504](#)
- K**
- Kalbfleisch–Prentice estimator, *see* survival function
 - Kaplan–Meier estimator, *see* survival function
 - knots, [22](#)
 - Kullback–Leibler information, [215](#)
- L**
- landmark survival time analysis, [447](#)
 - lasso, [71](#), [100](#), [121](#), [175](#), [356](#)
 - L^AT_EX, [129](#), [536](#)
 - `latex`, [129](#), [134](#), [135](#), [137](#), [138](#), [149](#), [246](#), [282](#), [292](#), [336](#), [342](#), [346](#), [453](#), [466](#), [470](#), [536](#)
 - `lattice` package, [134](#)
 - least squares
 - censored, [447](#)
 - leave-out-one, *see* validation of model
 - left truncation, [401](#), [420](#)
 - life expectancy, [4](#), [408](#), [472](#)
 - lift curve, [5](#)
 - likelihood function, [182](#), [187](#), [188](#), [190](#), [194](#), [195](#), [424](#), [425](#), [476](#)
 - partial, [477](#)
 - likelihood ratio test, [185–186](#), [189–191](#), [193–195](#), [198](#), [204](#), [205](#), [207](#), [228](#), [240](#)
 - linear model, [73](#), [74](#), [143](#), [311](#), [359](#), [361](#), [362](#), [364](#), [368](#), [370](#), [372](#)
 - case study, [143](#)
 - linear spline, *see* spline function
 - link function, [15](#)
 - Cauchy, [362](#)
 - complementary log-log, [362](#)
 - log-log, [362](#)
 - probit, [362](#)
 - `lm`, [131](#)
 - `lme`, [149](#)
 - local regression, *see* nonparametric
 - loess, *see* nonparametric
 - `loess`, [29](#), [142](#), [493](#)
 - log-rank, *see* hypothesis test
 - LOGISTIC, [315](#)
 - logistic model
 - binary, [219–231](#)
 - case study 1, [275–288](#)
 - case study 2, [291–310](#)
 - conditional, [483](#)
 - continuation ratio, [319–323](#)
 - case study, [338–340](#)
 - extended continuation ratio, [321–322](#)
 - case study, [340–355](#)

- ordinal, 311
- proportional odds, 73, 311, 312, **313–319**, 333, 362, 364
 - case study, 333–338
- logLik, 134, 135
- longitudinal data, 143
- lowess, *see* nonparametric
- lowess, 141, 294
- lrm, 65, 131, 134, 135, 201, **269**, 269, 273, 277, 278, 296, 297, 302, 306, 319, 323, 335, 337, 339, 341, 342, 448, 513
- lrtest, 134, 135
- lsp, 133
- M**
- Mallows' C_p , 69
- Mantel–Haenszel test, *see* hypothesis test
- marginal distribution, 26, 417, 478
- marginal estimates, *see* unconditioning
- martingale residual, 487, 493, 494, 515, 516
- matrix, 133
- matrx, 133
- maximal correlation, 390
- maximum generalized variance, 82, 83
- maximum likelihood, 147
 - estimation, **181**, 231, 424, 425, 477
 - penalized, 11, 77, 78, 115, 136, **209–212**, 269, 327, 328, 353
 - case study, 342–355
 - weighted, 208
- maximum total variance, 81
- Mean, 135, 319, 448, 472, 513, 514
- meningitis, *see* datasets
- mgcv package, 390
- MGV, *see* maximum generalized variance
- MICE, 54, 55, 59
- missing data, 143, 302
 - casewise deletion, 47, 48, 81, 296, 307, 384
 - describing patterns, *see* naclus, naplot
 - imputation, *see* imputation
 - informative, 46, 424
 - random, 46
- MLE, *see* maximum likelihood
- model
 - accelerated failure time, **436–446**, 453
 - case study, 453–473
 - Andersen–Gill, 513
 - approximate, **119–123**, 275, 287, 349, 352–354, 356
 - Buckley–James, 447, 449
 - comparing more than one, 92
 - Cox, *see* Cox model
 - cumulative link, *see* cumulative probability model
 - cumulative probability, *see* cumulative probability model
 - extended linear, 146
 - generalized additive, *see* generalized additive model, 359
 - generalized linear, 146, 359
 - growth curve, 146
 - linear, *see* linear model, 117, 199, 287, 317, 389
 - log-logistic, 437
 - log-normal, 437, 453
 - logistic, *see* logistic model
 - longitudinal, 143
 - ols, 146
 - ordinal, *see* ordinal model
 - parametric proportional hazards, 427
 - quantile regression, *see* quantile regression
 - semiparametric, *see* semiparametric model

validation, *see* validation of model
 model approximation, *see* model
 model uncertainty, 170, 304
 model validation, *see* validation of model
 modeling strategy, *see* strategy
 monotone, 393
 monotonicity, 66, 83, 84, 95, 129, 166, 389, 390, 393, 458
 MTV, *see* maximum total variance
 multcomp package, 199, 202
 multi-state model, 420
 multiple events, 417

N

na.action, 131
 na.delete, 131, 132
 na.detail.response, 131
 na.fail, 132
 na.fun.response, 131
 na.omit, 132
 naclus, 47, 142, 302, 458, 461
 naplot, 47, 302, 461
 naprint, 135
 naresid, 132, 135
 natural spline, *see* restricted cubic spline
 nearest neighbor, 51
 Nelson estimator, *see* survival function, 422
 NewLabels, 473
 Newton–Raphson algorithm, 193, 195, 196, 209, 231, 426
 NHANES, 365
 nlme package, 131, 148, 149
 noise, 34, 68, 69, 72, 209, 488, 523
 nomogram, 104, 268, 310, 318, 353, 514, 531
 nomogram, 135, 138, 149, 282, 319, 353, 473, 514
 non-proportional hazards, 73, 450, 506

noncompliance, 402, 513
 nonignorable nonresponse, *see* missing data
 nonparametric
 correlation, 66
 censored data, 517
 generalized Spearman correlation, 66, 376
 independence test, 129, 166
 regression, 29, 41, 105, 142, 245, 285
 test, 2, 66, 129
 nonproportional hazards, 495
 npsurv, 418, 419
 ns, 132, 133
 nuisance parameter, 190, 191

O

object-oriented program, x, 127, 133
 observational study, 3, 58, 230, 400
 odds ratio, 222, 224, 318
 OLS, *see* linear model
 ols, 131, 135, 137, 350, 351, 448, 469, 470
 optimism, 109, 111, 114, 391
 ordered, 133
 ordinal model, 311, 359, 361–363, 370, 371
 case study, 327–356, 359–387
 probit, 364
 ordinal response, *see* response
 ordinality, *see* assumptions
 orm, 131, 135, 319, 362, 363
 outlier, 116, 294
 overadjustment, 2
 overfitting, 72, 109–110

P

parsimony, 87, 97, 119
 partial effect plot, 104, 318
 partial residual, *see* residual
 partial test, *see* hypothesis test
 PC, *see* principal component, 170, 172, 175, 275

- pcaPP package, 175
 pec package, 519
 penalized maximum likelihood,
 see maximum likelihood
 pentrace, 134, 136, 269, 323, 342,
 344
 person-years, 408, 425
 plclust, 129
 plot.lrm.partial, 339
 plot.xmean.ordinaly, 319, 323, 333
 plsmo, 358
 Poisson model, 271
 pol, 133
 poly, 132, 133
 polynomial, 21
 popower, 319
 posamsize, 319
 power calculation, *see* cpower,
 spower, ciapower, popower
 pphsm, 448
 prcomp, 141
 preconditioning, 118, 123
 predab.resample, 141, 269, 323
 Predict, 130, 134, 136, 149,
 198, 199, 202, 278, 299, 307,
 319, 448, 466
 predict, 127, 132, 136, 140, 309,
 319, 469, 517, 526
 predictor
 continuous, 21, 40
 nominal, 16, 210
 ordinal, 38
 principal component, 81, 87,
 101, 275
 sparse, 101, 175
 princomp, 141, 171
 PRINQUAL, 82, 83
 product-limit estimator, *see*
 survival function
 propensity score, 3, 58, 231
 proportional hazards model, *see*
 Cox model
 proportional odds model, *see*
 logistic model
 prostate, *see* datasets
 psm, 131, 135, 448, 448,
 460, 464, 513
Q
 Q-R decomposition, 23
 Q-Q plot, 148
 qr, 192
 Quantile, 135, 448, 472, 513, 514
 quantile regression, 359, 360, 364,
 370, 379, 392
 composite, 361
 quantreg, 131, 360
R
 random forests, 100
 rank correlation, *see*
 nonparametric
 Rao score test, 186–187,
 191, 193–195, 198
 rcorr, 166
 rcorr.cens, 142, 461, 517
 rcorr.cens, 461
 rcorr.p.cens, 142
 rcs, 133, 296, 297
 rcspline.eval, 129
 rcspline.plot, 273
 rcspline.restate, 129
 receiver operating characteristic
 curve, 6, 11
 area, 92, 93, 111, 257, 346
 area, generalized, 318, 505
 recursive partitioning, 10, 30, 31,
 41, 46, 47, 51, 52, 83, 87,
 100, 120, 142, 302, 349
 redun, 80, 463
 redundancy analysis, 80, 175
 regression to the mean, 75, 530
 resampling, 105, 112
 resid, 134, 336, 337, 460, 516
 residual
 logistic score, 314, 336
 martingale, 487, 493, 494,
 515, 516
 partial, 34, 272, 315, 321, 337

Schoenfeld score, 314, **487**,
498, 499, 516, 517, 525, 526
 residuals, 132, 134, 269, 336, 337,
460, 516
 residuals.coxph, 516
 response
 binary, 219–221
 censored or truncated, 401
 continuous, **389–398**
 ordinal, 311, 327, 359
 restricted cubic spline, *see* spline
 function
 ridge regression, 77, 115, 209, 210
 risk difference, 224, 430
 risk ratio, 224, 430
 rms package, xi, 129, **130–141**,
149, 192, 193, 198, 199, 211,
214, 319, 362, 363, 418,
422, 535
 robcov, 134, 135, 198, 202
 robust covariance estimator, *see*
 variance–covariance matrix
 robustgam package, 390
 ROC, *see* receiver operating
 characteristic curve, 105
 rpart, 142, 302, 303
 Rq, 131, 135, 360
 rq, 131
 runif, 460

S

sample size, 73, 74, 148,
233, 363, 486
 sample survey, 135, 197, 208, 417
 sas.get, 129
 sascode, 138
 scientific quantity, 20
 score function, 182, 183, 186
 score test, *see* Rao score test,
235, 363
 score.binary, 86
 scored, 132, 133
 scoring, hierarchical, 86
 scree plot, 172

semiparametric model, 311, 359,
361–363, 370, 371, 475
 sensuc, 134
 shrinkage, 75–78, 87, 88,
209–212, 342–348
 similarity measure, 81, 330, 458
 smearing estimator, *see* estimator
 smoother, 390
 Somers' rank correlation, *see* D_{xy}
 somers2, 346
 spca package, 175
 sPCAgrid, 175, 179
 Spearman rank correlation, *see*
 nonparametric
 spearman2, 129, 460
 specs, 134, 135
 spline function, 22, 30,
167, 192, 393
 B-spline, 23, 41, 132, 500
 cubic, 23
 linear, 22, 133
 normalization, 26
 restricted cubic, 24–28
 tensor, 37, 247, 374, 375
 spower, 513
 standardized regression
 coefficient, 103
 state transition, 416, 420
 step, 134
 step halving, 196
 strat, 133
 strata, 133
 strategy, 63
 comparing models, 92
 data reduction, 79
 describing model, 103, 318
 developing imputations, 49
 developing model for effect
 estimation, 98
 developing models for
 hypothesis testing, 99
 developing predictive model, 95
 global, 94
 in a nutshell, ix, 95
 influential observations, 90

- maximum number of
 - parameters, 72
 - model approximation, 118, 275, 287
 - multiple imputation, 53
 - prespecification of complexity, 64
 - shrinkage, 77
 - validation, 109, 110
 - variable selection, 63, 67
 - stratification, 225, 237, 238, **254**, 418, 419, **481–483**, 488
 - subgroup estimates, 34, 241, 400
 - summary, 127, 130, 134, 136, 149, 167, 198, 199, 201, 278, 292, 466
 - summary.formula, 302, 319, 357
 - summary.gls, 149
 - super smoother, 29
 - SUPPORT study, *see* datasets
 - suppression, 101
 - supsmu, 141, 273, 390
 - Surv, 172, 418, 422, 458, 516
 - survConcordance, 517
 - survdiff, 517
 - survest, 135, 448
 - survfit, 135, 418, 419
 - Survival, 135, 448, 513, 514
 - survival function
 - Aalen estimator, 412, 413
 - Breslow estimator, 485
 - crude, 416
 - Fleming–Harrington estimator, 412, 413, 485
 - Kalbfleisch–Prentice estimator, 484, 485
 - Kaplan–Meier estimator, **409–413**, 414–416, 420
 - multiple state estimator, 416, 420
 - Nelson estimator, 412, 413, 418, 485
 - standard error, 412
 - survival package, 131, 418, 422, 499, 513, 517, 536
 - survplot, 135, 419, 448, 458, 460
 - survreg, 131, 448
 - survreg.auxinfo, 449
 - survreg.distributions, 449
- T**
- test of linearity, *see* hypothesis test
 - test statistic, *see* hypothesis test
 - time to event, 399
 - and severity of event, 417
 - time-dependent covariable, 322, 418, 447, 499–503, 513, 518, 526
 - Titanic, *see* datasets
 - training sample, 111–113, 122
 - transace, 176, 177
 - transcan, 51, 55, 80, **83**, 83–85, 129, 135, 138, 167, 170–172, 175–177, 276, 277, 330, 334, 335, 521, 525
 - transform both sides regression, 176, 389, 392
 - transformation, 389, 393, 395
 - post, 133
 - pre, 179
 - tree model, *see* recursive partitioning
 - truncation, 401
- U**
- unconditioning, 119
 - uniqueness analysis, 94
 - univariable screening, 72
 - univarLR, 134, 135
 - unsupervised learning, 79
- V**
- val.prob, 109, 135, 271
 - val.surv, 109, 449, 517
 - validate, 135, **141**, 142, 260, 269, 271, 282, 286, 300, 301, 319, 323, 354, 466, 517

- validation of model, **109–116**,
259, 299, 318, 322, 353, 446,
466, 506, 529
 - bootstrap, 114–116
 - cross, 113, 115, 116, 210
 - data-splitting, **111**, 112, 271
 - external, 109, 110, 237,
271, 449, 517
 - leave-out-one, 113, 122,
215, 255
 - quantities to validate, 110
 - randomization, 113
 - `varclus`, 79, 129, 167, 330, 458,
463
 - variable selection, **67–72**, 171
 - step-down, 70, 137,
275, 280, 282, 286, 377
 - variance inflation factors, 79, 135,
138, 255
 - variance stabilization, 390
 - variance–covariance matrix,
51, 54, 120, 129, 189,
191, 193, 196–198, 208,
211, 215
 - cluster sandwich, 197, 202
 - Huber–White estimator, 147
 - sandwich, 147, 211, 217
 - variogram, 148, 153
 - `vcov`, 134, 135
 - `vif`, 135, 138
- W**
- waiting time, 401
 - Wald statistic, **186**, 189, 191, 192,
194, 196, 198, 206, **244**, 278
 - weighted analysis, *see* maximum
likelihood
 - `which.influence`, 134, 137, 269
 - working independence model, 197