# Springer Series in Statistics

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

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

*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

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                                            Frank E. Harrell, Jr.
July 2015

# Contents

# Typographical Conventions

Boxed numbers in the margins such as $\boxed{1}$ 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.