# Applied Regression Including Computing and Graphics

R. DENNIS COOK

SANFORD WEISBERG
The University of Minnesota
St. Paul, Minnesota

This Page Intentionally Left Blank

Applied Regression
Including Computing and Graphics

# Applied Regression Including Computing and Graphics

R. DENNIS COOK

SANFORD WEISBERG
The University of Minnesota
St. Paul, Minnesota

For ordering and customer service, call 1-800-CALL WILEY.

This book is dedicated to the students
who turn words and graphs into ideas and discoveries.

This Page Intentionally Left Blank

# Contents

This Page Intentionally Left Blank

# Preface

This textbook is about regression, the study of how a response variable depends on one or more predictors. Regression is one of the fundamental tools of data analysis. This book is intended for anyone interested in applying regression to data. The mathematical level required is intentionally low; prerequisite is only one semester of basic statistical methods.

The main features of the book are as follows:

- Emphasis is placed on *seeing* results through graphs. We present many easily used graphical methods, most based on simple two-dimensional scatterplots, that provide analysts with more insight into their data than would have been possible otherwise, including a deeper appreciation for interpretation.
- We provide user-friendly computer software called *Arc* that lets the reader immediately apply the ideas we present, both to the examples in the book and to their own data. Many of the methods we use are impossible or at best difficult to implement in standard statistical software packages.
- This book is a *complete* textbook on applied regression suitable for a one semester course. We start with very basic ideas, progressing from standard ideas for linear models integrated with newer graphical approaches, to regression graphics and more complex models like generalized linear models.
- The book includes over 300 figures. The reader following along on the computer can reproduce almost all of them.
- Most of the examples and homework problems are based on real data. All of the data sets used in the book are included with *Arc*.
- A companion Internet site includes the software, more problems, help for the reader, additional statistical material, extensions to the software, and much more.

We have used drafts of this textbook for several years as a basis for two different courses in applied regression analysis: One is intended primarily for

advanced undergraduate and beginning graduate students in fields other than
statistics, while the other is intended for first-year statistics graduate students.
Drafts have also been used by others for students in social science, business,
and other disciplines. After completing these courses, our students can analyze
regression problems with greater ease and more depth than could our students
before we began to use this book.

## *Arc*

Integrated into the text is discussion of a computer package called *Arc*. This
is a user-friendly program designed specifically for studying this material, as
well as for applying the ideas learned to other data sets. The program permits
the user to do the analyses discussed in the book easily and quickly. It can
be down-loaded for free from the Internet site for this book at the address
given in the Appendix; versions are available for Windows, Macintosh, and
Unix.

For those readers who prefer to use other statistical packages, we include
on our Internet site descriptions of how a few of the major packages can be
used for some of the calculations.

## ORGANIZATION AND STYLE

When writing the book, we envisioned the reader sitting at a computer and
reworking the examples in the text. Indeed, some of the text can be read as
if we were next to the reader, suggesting what to do next. To maintain this
low-key style, we have tried to avoid heavy algebra. References and technical
comments are collected in the complements section at the end of each chapter.

## PATHS

The book is divided into four parts:

*I: Introduction.* The first part consists of five chapters that present many
basic ideas of regression to set the stage for later developments, including
the one-sample problem, using and interpreting histograms and scatterplots,
smoothing including density estimates and scatterplot smoothers, and bivariate
distributions like the normal. These chapters weave together standard and non-
standard topics that provide a basis for understanding regression. For example,
smoothing is presented before simple linear regression.

*II: Multiple Linear Regression.* The second part of the book weaves stan-
dard linear model ideas, starting with simple regression, with the basics of
graphics, including 2D and 3D scatterplots and scatterplot matrices. All these

graph types are used repeatedly throughout the rest of the book. Emphasis is split between basic results assuming the multiple linear regression model holds, graphical ideas, the use of transformations, and graphs as the basis for diagnostic methods.

*III: Graphics.* The third part of the book, which is unique to this work, shows how graphs can be used to better understand regression problems in which no model is available, or else an appropriate model is in doubt. These methods allow the analyst to see appropriate answers, and consequently have increased faith that the data sustain any models that are developed.

*IV: Other Models.* The last part of the book gives the fundamentals of fitting generalized linear models. Most of the graphical methods included here are also unique to this book.

For a one-quarter (30-lecture) course, we recommend a very quick tour of Part I (3–5 lectures), followed by about 18 lectures on Part II, with the remainder of the course spent in Part III. For a semester course, the introduction can be expanded, and presentation of Parts II and III can be slowed down as well. We have used Part IV of the book in the second quarter of a two-quarter course; completing the whole book in one semester would require a fairly rapid pace. Some teachers may prefer substituting Part IV for Part III; we don't recommend this because we believe that the methodologies described in Part III are too useful to be skipped.

Some of the material relating to the construction of graphical displays is suitable for presentation in a laboratory or recitation session by a teaching assistant. For example, we have covered Section 7.1 and most of Chapters 5 and 8 in this way. Further suggestions on teaching from this book are available in a teacher's manual; see the Internet site for more information.

## OTHER BOOKS

The first book devoted entirely to regression was probably by the American economist Mordechai Ezekiel (1924, revised 1930 and 1941 and finally as Ezekiel and Fox, 1959). Books on regression have proliferated since the advent of computers in universities. We have contributed to this literature. Weisberg (1985) provided an introduction to applied regression, particularly to multiple linear regression, but at a modestly higher mathematical level. This book includes virtually all the material in Weisberg (1985); even some of the examples, but none of the prose, are common between the two books. Cook and Weisberg (1994b) provided an introduction to graphics and regression. Nearly all of the material in that book, though again little of the prose, has found its way into this book. Chapter 15 contains a low-level introduction to the material in the research monograph on residual and influence analysis by Cook and Weisberg (1982). Finally, Cook (1998b) provides a mathematical

and rigorous treatment of regression through graphics that is the core of this book.

On the Internet site for this book, we provide additional references for reading for those who would like a more theoretical approach to this area.

## ACKNOWLEDGMENTS

<div align="right">

R. DENNIS COOK
SANFORD WEISBERG

</div>

*St. Paul, Minnesota*
*May, 1999*

Applied Regression
Including Computing and Graphics

This Page Intentionally Left Blank