# LINEAR MODELS IN STATISTICS

Each generation has its unique needs and aspirations. When Charles Wiley first opened his small printing shop in lower Manhattan in 1807, it was a generation of boundless potential searching for an identity. And we were there, helping to define a new American literary tradition. Over half a century later, in the midst of the Second Industrial Revolution, it was a generation focused on building the future. Once again, we were there, supplying the critical scientific, technical, and engineering knowledge that helped frame the world. Throughout the 20th Century, and into the new millennium, nations began to reach out beyond their own borders and a new international community was born. Wiley was there, expanding its operations around the world to enable a global exchange of ideas, opinions, and know-how.

For 200 years, Wiley has been an integral part of each generation's journey, enabling the flow of information and understanding necessary to meet their needs and fulfill their aspirations. Today, bold new technologies are changing the way we live and learn. Wiley will be there, providing you the must-have knowledge you need to imagine new worlds, new possibilities, and new opportunities.

Generations come and go, but you can always count on Wiley to provide you the knowledge you need, when and where you need it!

**WILLIAM J. PESCE**
PRESIDENT AND CHIEF EXECUTIVE OFFICER

**PETER BOOTH WILEY**
CHAIRMAN OF THE BOARD

# LINEAR MODELS IN STATISTICS

**Second Edition**

**Alvin C. Rencher and G. Bruce Schaalje**
*Department of Statistics, Brigham Young University, Provo, Utah*

# CONTENTS

# PREFACE

In the second edition, we have added chapters on Bayesian inference in linear models (Chapter 11) and linear mixed models (Chapter 17), and have upgraded the material in all other chapters. Our continuing objective has been to introduce the theory of linear models in a clear but rigorous format.

In spite of the availability of highly innovative tools in statistics, the main tool of the applied statistician remains the linear model. The linear model involves the simplest and seemingly most restrictive statistical properties: independence, normality, constancy of variance, and linearity. However, the model and the statistical methods associated with it are surprisingly versatile and robust. More importantly, mastery of the linear model is a prerequisite to work with advanced statistical tools because most advanced tools are generalizations of the linear model. The linear model is thus central to the training of any statistician, applied or theoretical.

This book develops the basic theory of linear models for regression, analysis-of-variance, analysis–of–covariance, and linear mixed models. Chapter 18 briefly introduces logistic regression, generalized linear models, and nonlinear models. Applications are illustrated by examples and problems using real data. This combination of theory and applications will prepare the reader to further explore the literature and to more correctly interpret the output from a linear models computer package.

This introductory linear models book is designed primarily for a one-semester course for advanced undergraduates or MS students. It includes more material than can be covered in one semester so as to give an instructor a choice of topics and to serve as a reference book for researchers who wish to gain a better understanding of regression and analysis-of-variance. The book would also serve well as a text for PhD classes in which the instructor is looking for a one-semester introduction, and it would be a good supplementary text or reference for a more advanced PhD class for which the students need to review the basics on their own.

Our overriding objective in the preparation of this book has been *clarity of exposition*. We hope that students, instructors, researchers, and practitioners will find this linear models text more comfortable than most. In the final stages of development, we asked students for written comments as they read each day's assignment. They made many suggestions that led to improvements in readability of the book. We are grateful to readers who have notified us of errors and other suggestions for improvements of the text, and we will continue to be very grateful to readers who take the time to do so for this second edition.

Another objective of the book is to tie up loose ends. There are many approaches to teaching regression, for example. Some books present estimation of regression coefficients for fixed $x$'s only, other books use random $x$'s, some use centered models, and others define estimated regression coefficients in terms of variances and covariances or in terms of correlations. Theory for linear models has been presented using both an algebraic and a geometric approach. Many books present classical (frequentist) inference for linear models, while increasingly the Bayesian approach is presented. We have tried to cover all these approaches carefully and to show how they relate to each other. We have attempted to do something similar for various approaches to analysis-of-variance. We believe that this will make the book useful as a reference as well as a textbook. An instructor can choose the approach he or she prefers, and a student or researcher has access to other methods as well.

The book includes a large number of theoretical problems and a smaller number of applied problems using real datasets. The problems, along with the extensive set of answers in Appendix A, extend the book in two significant ways: (1) the theoretical problems and answers fill in nearly all gaps in derivations and proofs and also extend the coverage of material in the text, and (2) the applied problems and answers become additional examples illustrating the theory. As instructors, we find that having answers available for the students saves a great deal of class time and enables us to cover more material and cover it better. The answers would be especially useful to a reader who is engaging this material outside the formal classroom setting.

The mathematical prerequisites for this book are multivariable calculus and matrix algebra. The review of matrix algebra in Chapter 2 is intended to be sufficiently complete so that the reader with no previous experience can master matrix manipulation up to the level required in this book. Statistical prerequisites include some exposure to statistical theory, with coverage of topics such as distributions of random variables, expected values, moment generating functions, and an introduction to estimation and testing hypotheses. These topics are briefly reviewed as each is introduced. One or two statistical methods courses would also be helpful, with coverage of topics such as $t$ tests, regression, and analysis-of-variance.

We have made considerable effort to maintain consistency of notation throughout the book. We have also attempted to employ standard notation as far as possible and to avoid exotic characters that cannot be readily reproduced on the chalkboard. With a few exceptions, we have refrained from the use of abbreviations and mnemonic devices. We often find these annoying in a book or journal article.

Equations are numbered sequentially throughout each chapter; for example, (3.29) indicates the twenty-ninth numbered equation in Chapter 3. Tables and figures are also numbered sequentially throughout each chapter in the form "Table 7.4" or "Figure 3.2." On the other hand, examples and theorems are numbered sequentially within a section, for example, Theorems 2.2a and 2.2b.

The solution of most of the problems with real datasets requires the use of the computer. We have not discussed command files or output of any particular program, because there are so many good packages available. Computations for the numerical examples and numerical problems were done with SAS. The datasets and SAS

command files for all the numerical examples and problems in the text are available on the Internet; see Appendix B.

The references list is not intended to be an exhaustive survey of the literature. We have provided original references for some of the basic results in linear models and have also referred the reader to many up-to-date texts and reference books useful for further reading. When citing references in the text, we have used the standard format involving the year of publication. For journal articles, the year alone suffices, for example, Fisher (1921). But for a specific reference in a book, we have included a page number or section, as in Hocking (1996, p. 216).

Our selection of topics is intended to prepare the reader for a better understanding of applications and for further reading in topics such as mixed models, generalized linear models, and Bayesian models. Following a brief introduction in Chapter 1, Chapter 2 contains a careful review of all aspects of matrix algebra needed to read the book. Chapters 3, 4, and 5 cover properties of random vectors, matrices, and quadratic forms. Chapters 6, 7, and 8 cover simple and multiple linear regression, including estimation and testing hypotheses and consequences of misspecification of the model. Chapter 9 provides diagnostics for model validation and detection of influential observations. Chapter 10 treats multiple regression with random $x$'s. Chapter 11 covers Bayesian multiple linear regression models along with Bayesian inferences based on those models. Chapter 12 covers the basic theory of analysis-of-variance models, including estimability and testability for the overparameterized model, reparameterization, and the imposition of side conditions. Chapters 13 and 14 cover balanced one-way and two-way analysis-of-variance models using an over-parameterized model. Chapter 15 covers unbalanced analysis-of-variance models using a cell means model, including a section on dealing with empty cells in two-way analysis-of-variance. Chapter 16 covers analysis of covariance models. Chapter 17 covers the basic theory of linear mixed models, including residual maximum likelihood estimation of variance components, approximate small-sample inferences for fixed effects, best linear unbiased prediction of random effects, and residual analysis. Chapter 18 introduces additional topics such as nonlinear regression, logistic regression, loglinear models, Poisson regression, and generalized linear models.

In our class for first-year master's-level students, we cover most of the material in Chapters 2–5, 7–8, 10–12, and 17. Many other sequences are possible. For example, a thorough one-semester regression and analysis-of-variance course could cover Chapters 1–10, and 12–15.

Al's introduction to linear models came in classes taught by Dale Richards and Rolf Bargmann. He also learned much from the books by Graybill, Scheffé, and Rao. Al expresses thanks to the following for reading the first edition manuscript and making many valuable suggestions: David Turner, John Walker, Joel Reynolds, and Gale Rex Bryce. Al thanks the following students at Brigham Young University (BYU) who helped with computations, graphics, and typing of the first edition: David Fillmore, Candace Baker, Scott Curtis, Douglas Burton, David Dahl, Brenda Price, Eric Hintze, James Liechty, and Joy Willbur. The students

in Al's Linear Models class went through the manuscript carefully and spotted many typographical errors and passages that needed additional clarification.

Bruce's education in linear models came in classes taught by Mel Carter, Del Scott, Doug Martin, Peter Bloomfield, and Francis Giesbrecht, and influential short courses taught by John Nelder and Russ Wolfinger.

We thank Bruce's Linear Models classes of 2006 and 2007 for going through the book and new chapters. They made valuable suggestions for improvement of the text. We thank Paul Martin and James Hattaway for invaluable help with LaTex. The Department of Statistics, Brigham Young University provided financial support and encouragement throughout the project.

## Second Edition

For the second edition we added Chapter 11 on Bayesian inference in linear models (including Gibbs sampling) and Chapter 17 on linear mixed models.

We also added a section in Chapter 2 on vector and matrix calculus, adding several new theorems and covering the Lagrange multiplier method. In Chapter 4, we presented a new proof of the conditional distribution of a subvector of a multivariate normal vector. In Chapter 5, we provided proofs of the moment generating function and variance of a quadratic form of a multivariate normal vector. The section on the geometry of least squares was completely rewritten in Chapter 7, and a section on the geometry of least squares in the overparameterized linear model was added to Chapter 12. Chapter 8 was revised to provide more motivation for hypothesis testing and simultaneous inference. A new section was added to Chapter 15 dealing with two-way analysis-of-variance when there are empty cells. This material is not available in any other textbook that we are aware of.

This book would not have been possible without the patience, support, and encouragement of Al's wife LaRue and Bruce's wife Lois. Both have helped and supported us in more ways than they know. This book is dedicated to them.

ALVIN C. RENCHER AND G. BRUCE SCHAALJE

*Department of Statistics*
*Brigham Young University*
*Provo, Utah*