# Introduction to
# Linear Regression Analysis

## FIFTH EDITION

DOUGLAS C. MONTGOMERY

ELIZABETH A. PECK

G. GEOFFREY VINING

WILEY

ftp://
SITE AVAILABLE

# INTRODUCTION TO LINEAR REGRESSION ANALYSIS

Fifth Edition

**DOUGLAS C. MONTGOMERY**

Arizona State University
School of Computing, Informatics, and Decision Systems Engineering
Tempe, AZ

**ELIZABETH A. PECK**

The Coca-Cola Company (retired)
Atlanta, GA

**G. GEOFFREY VINING**

Virginia Tech
Department of Statistics
Blacksburg, VA

# INTRODUCTION TO LINEAR REGRESSION ANALYSIS

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

A complete list of the titles in this series appears at the end of this volume.

# PREFACE

Regression analysis is one of the most widely used techniques for analyzing multi-factor data. Its broad appeal and usefulness result from the conceptually logical process of using an equation to express the relationship between a variable of interest (the response) and a set of related predictor variables. Regression analysis is also interesting theoretically because of elegant underlying mathematics and a well-developed statistical theory. Successful use of regression requires an appreciation of both the theory and the practical problems that typically arise when the technique is employed with real-world data.

This book is intended as a text for a basic course in regression analysis. It contains the standard topics for such courses and many of the newer ones as well. It blends both theory and application so that the reader will gain an understanding of the basic principles necessary to apply regression model-building techniques in a wide variety of application environments. The book began as an outgrowth of notes for a course in regression analysis taken by seniors and first-year graduate students in various fields of engineering, the chemical and physical sciences, statistics, mathematics, and management. We have also used the material in many seminars and industrial short courses for professional audiences. We assume that the reader has taken a first course in statistics and has familiarity with hypothesis tests and confidence intervals and the normal, $t$, $\chi^2$, and $F$ distributions. Some knowledge of matrix algebra is also necessary.

The computer plays a significant role in the modern application of regression. Today even spreadsheet software has the capability to fit regression equations by least squares. Consequently, we have integrated many aspects of computer usage into the text, including displays of both tabular and graphical output, and general discussions of capabilities of some software packages. We use Minitab®, JMP®, SAS®, and R for various problems and examples in the text. We selected these packages because they are widely used both in practice and in teaching regression and they have good regression. Many of the homework problems require software

for their solution. All data sets in the book are available in electronic form from the publisher. The ftp site ftp://ftp.wiley.com/public/sci_tech_med/introduction_linear_regression hosts the data, problem solutions, PowerPoint files, and other material related to the book.

## CHANGES IN THE FIFTH EDITION

We have made extensive changes in this edition of the book. This includes the reorganization of text material, new examples, new exercises, a new chapter on time series regression, and new material on designed experiments for regression models. Our objective was to make the book more useful as both a text and a reference and to update our treatment of certain topics.

Chapter 1 is a general introduction to regression modeling and describes some typical applications of regression. Chapters 2 and 3 provide the standard results for least-squares model fitting in simple and multiple regression, along with basic inference procedures (tests of hypotheses, confidence and prediction intervals). Chapter 4 discusses some introductory aspects of model adequacy checking, including residual analysis and a strong emphasis on residual plots, detection and treatment of outliers, the PRESS statistic, and testing for lack of fit. Chapter 5 discusses how transformations and weighted least squares can be used to resolve problems of model inadequacy or to deal with violations of the basic regression assumptions. Both the Box–Cox and Box–Tidwell techniques for analytically specifying the form of a transformation are introduced. Influence diagnostics are presented in Chapter 6, along with an introductory discussion of how to deal with influential observations. Polynomial regression models and their variations are discussed in Chapter 7. Topics include the basic procedures for fitting and inference for polynomials and discussion of centering in polynomials, hierarchy, piecewise polynomials, models with both polynomial and trigonometric terms, orthogonal polynomials, an overview of response surfaces, and an introduction to nonparametric and smoothing regression techniques. Chapter 8 introduces indicator variables and also makes the connection between regression and analysis-of-variance models. Chapter 9 focuses on the multicollinearity problem. Included are discussions of the sources of multicollinearity, its harmful effects, diagnostics, and various remedial measures. We introduce biased estimation, including ridge regression and some of its variations and principal-component regression. Variable selection and model-building techniques are developed in Chapter 10, including stepwise procedures and all-possible-regressions. We also discuss and illustrate several criteria for the evaluation of subset regression models. Chapter 11 presents a collection of techniques useful for regression model validation.

The first 11 chapters are the nucleus of the book. Many of the concepts and examples flow across these chapters. The remaining four chapters cover a variety of topics that are important to the practitioner of regression, and they can be read independently. Chapter 12 in introduces nonlinear regression, and Chapter 13 is a basic treatment of generalized linear models. While these are perhaps not standard topics for a linear regression textbook, they are so important to students and professionals in engineering and the sciences that we would have been seriously remiss without giving an introduction to them. Chapter 14 covers regression

models for time series data. Chapter 15 includes a survey of several important topics, including robust regression, the effect of measurement errors in the regressors, the inverse estimation or calibration problem, bootstrapping regression estimates, classification and regression trees, neural networks, and designed experiments for regression.

In addition to the text material, Appendix C contains brief presentations of some additional topics of a more technical or theoretical nature. Some of these topics will be of interest to specialists in regression or to instructors teaching a more advanced course from the book. Computing plays an important role in many regression courses. Mintab, JMP, SAS, and R are widely used in regression courses. Outputs from all of these packages are provided in the text. Appendix D is an introduction to using SAS for regression problems. Appendix E is an introduction to R.

## USING THE BOOK AS A TEXT

Because of the broad scope of topics, this book has great flexibility as a text. For a first course in regression, we would recommend covering Chapters 1 through 10 in detail and then selecting topics that are of specific interest to the audience. For example, one of the authors (D.C.M.) regularly teaches a course in regression to an engineering audience. Topics for that audience include nonlinear regression (because mechanistic models that are almost always nonlinear occur often in engineering), a discussion of neural networks, and regression model validation. Other topics that we would recommend for consideration are multicollinearity (because the problem occurs so often) and an introduction to generalized linear models focusing mostly on logistic regression. G.G.V. has taught a regression course for graduate students in statistics that makes extensive use of the Appendix C material.

We believe the computer should be directly integrated into the course. In recent years, we have taken a notebook computer and computer projector to most classes and illustrated the techniques as they are introduced in the lecture. We have found that this greatly facilitates student understanding and appreciation of the techniques. We also require that the students use regression software for solving the homework problems. In most cases, the problems use real data or are based on real-world settings that represent typical applications of regression.

There is an instructor's manual that contains solutions to all exercises, electronic versions of all data sets, and questions/problems that might be suitable for use on examinations.

## ACKNOWLEDGMENTS

often in the form of penetrating questions, that led to rewriting or expansion of material in the book. We are also indebted to John Wiley & Sons, the American Statistical Association, and the Biometrika Trustees for permission to use copyrighted material.

DOUGLAS C. MONTGOMERY
ELIZABETH A. PECK
G. GEOFFREY VINING

# CONTENTS