Springer Texts in Statistics

Advisors: George Casella Stephen Fienberg Ingram Olkin

Springer

New York Berlin Heidelberg Barcelona Hong Kong London Milan Paris Singapore Tokyo

Springer Texts in Statistics

Alfred: Elements of Statistics for the Life and Social Sciences Berger: An Introduction to Probability and Stochastic Processes Bilodeau and Brenner: Theory of Multivariate Statistics Blom: Probability and Statistics: Theory and Applications Brockwell and Davis: An Introduction to Times Series and Forecasting Chow and Teicher: Probability Theory: Independence, Interchangeability, Martingales, Third Edition Christensen: Plane Answers to Complex Questions: The Theory of Linear Models, Second Edition Christensen: Linear Models for Multivariate, Time Series, and Spatial Data Christensen: Log-Linear Models and Logistic Regression, Second Edition Creighton: A First Course in Probability Models and Statistical Inference Dean and Voss: Design and Analysis of Experiments du Toit, Stevn, and Stumpf: Graphical Exploratory Data Analysis Durrett: Essentials of Stochastic Processes Edwards: Introduction to Graphical Modelling, Second Edition Finkelstein and Levin: Statistics for Lawyers Flury: A First Course in Multivariate Statistics Jobson: Applied Multivariate Data Analysis, Volume I: Regression and Experimental Design Jobson: Applied Multivariate Data Analysis, Volume II: Categorical and Multivariate Methods Kalbfleisch: Probability and Statistical Inference, Volume I: Probability, Second Edition Kalbfleisch: Probability and Statistical Inference, Volume II: Statistical Inference, Second Edition Karr: Probability Keyfitz: Applied Mathematical Demography, Second Edition Kiefer: Introduction to Statistical Inference Kokoska and Nevison: Statistical Tables and Formulae Kulkarni: Modeling, Analysis, Design, and Control of Stochastic Systems Lehmann: Elements of Large-Sample Theory Lehmann: Testing Statistical Hypotheses, Second Edition Lehmann and Casella: Theory of Point Estimation, Second Edition Lindman: Analysis of Variance in Experimental Design Lindsey: Applying Generalized Linear Models Madansky: Prescriptions for Working Statisticians McPherson: Applying and Interpreting Statistics: A Comprehensive Guide, Second Edition Mueller: Basic Principles of Structural Equation Modeling: An Introduction to LISREL and EOS

(continued after index)

John O. Rawlings Sastry G. Pantula David A. Dickey

Applied Regression Analysis

A Research Tool

Second Edition

With 78 Figures



John O. Rawlings Sastry G. Pantula David A. Dickey Department of Statistics North Carolina State University Raleigh, NC 27695 USA

Editorial Board

George Casella Biometrics Unit Cornell University Ithaca, NY 14853-7801 USA Stephen Fienberg Department of Statistics Carnegie Mellon University Pittsburgh, PA 15213-3890 USA Ingram Olkin Department of Statistics Stanford University Stanford, CA 94305 USA

Library of Congress Cataloging-in-Publication Data
Rawlings, John O., 1932–
Applied regression analysis: a research tool. — 2nd ed. / John
O. Rawlings, Sastry G. Pentula, David A. Dickey.
p. cm. — (Springer texts in statistics)
Includes bibliographical references and indexes.
ISBN 0-387-98454-2 (hardcover: alk. paper)
1. regression analysis. I. Pentula, Sastry G. II. Dickey, David
A. III. Title. IV. Series.
QA278.2.R38 1998
519.5'36—dc21
97-48858

Printed on acid-free paper.

© 1989 Wadsworth, Inc.

© 1998 Springer-Verlag New York, Inc.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer-Verlag New York, Inc., 175 Fifth Avenue, New York, NY 10010, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use of general descriptive names, trade names, trademarks, etc., in this publication, even if the former are not especially identified, is not to be taken as a sign that such names, as understood by the Trade Marks and Merchandise Marks Act, may accordingly be used freely by anyone.

 $9\ 8\ 7\ 6\ 5\ 4\ 3\ 2\ 1$

ISBN 0-387-98454-2 Springer-Verlag New York Berlin Heidelberg SPIN 10660129

То

Our Families

PREFACE

This text is a new and improved edition of Rawlings (1988). It is the outgrowth of several years of teaching an applied regression course to graduate students in the sciences. Most of the students in these classes had taken a two-semester introduction to statistical methods that included experimental design and multiple regression at the level provided in texts such as Steel, Torrie, and Dickey (1997) and Snedecor and Cochran (1989). For most, the multiple regression had been presented in matrix notation.

The basic purpose of the course and this text is to develop an understanding of least squares and related statistical methods without becoming excessively mathematical. The emphasis is on regression concepts, rather than on mathematical proofs. Proofs are given only to develop facility with matrix algebra and comprehension of mathematical relationships. Good students, even though they may not have strong mathematical backgrounds, quickly grasp the essential concepts and appreciate the enhanced understanding. The learning process is reinforced with continuous use of numerical examples throughout the text and with several case studies. Some numerical and mathematical exercises are included to whet the appetite of graduate students.

The first four chapters of the book provide a review of simple regression in algebraic notation (Chapter 1), an introduction to key matrix operations and the geometry of vectors (Chapter 2), and a review of ordinary least squares in matrix notation (Chapters 3 and 4). Chapter 4 also provides a foundation for the testing of hypotheses and the properties of sums of squares used in analysis of variance. Chapter 5 is a case study giving a complete multiple regression analysis using the methods reviewed in the first four chapters. Then Chapter 6 gives a brief geometric interpretation of least squares illustrating the relationships among the data vectors, the link between the analysis of variance and the lengths of the vectors, and the role of degrees of freedom. Chapter 7 discusses the methods and criteria for determining which independent variables should be included in the models. The next two chapters include special classes of multiple regression models. Chapter 8 introduces polynomial and trigonometric regression models. This chapter also discusses response curve models that are linear in the parameters. Class variables and the analysis of variance of designed experiments (models of less than full rank) are introduced in Chapter 9.

Chapters 10 through 14 address some of the problems that might be encountered in regression. A general introduction to the various kinds of problems is given in Chapter 10. This is followed by discussions of regression diagnostic techniques (Chapter 11), and scaling or transforming variables to rectify some of the problems (Chapter 12). Analysis of the correlational structure of the data and biased regression are discussed as techniques for dealing with the collinearity problem common in observational data (Chapter 13). Chapter 14 is a case study illustrating the analysis of data in the presence of collinearity.

Models that are nonlinear in the parameters are presented in Chapter 15. Chapter 16 is another case study using polynomial response models, nonlinear modeling, transformations to linearize, and analysis of residuals. Chapter 17 addresses the analysis of unbalanced data. Chapter 18 (new to this edition) introduces linear models that have more than one random effect. The ordinary least squares approach to such models is given. This is followed by the definition of the variance–covariance matrix for such models and a brief introduction to mixed effects and random coefficient models. The use of iterative maximum likelihood estimation of both the variance components and the fixed effects is discussed. The final chapter, Chapter 19, is a case study of the analysis of unbalanced data.

We are grateful for the assistance of many in the development of this book. Of particular importance have been the dedicated editing of the earlier edition by Gwen Briggs, daughter of John Rawlings, and her many suggestions for improvement. It is uncertain when the book would have been finished without her support. A special thanks goes to our former student, Virginia Lesser, for her many contributions in reading parts of the manuscript, in data analysis, and in the enlistment of many data sets from her graduate student friends in the biological sciences. We are indebted to our friends, both faculty and students, at North Carolina State University for bringing us many interesting consulting problems over the years that have stimulated the teaching of this material. We are particularly indebted to those (acknowledged in the text) who have generously allowed the use of their data. In this regard, Rick Linthurst warrants special mention for his stimulating discussions as well as the use of his data. We acknowledge the encouragement and valuable discussions of colleagues in the Department of Statistics at NCSU, and we thank Matthew Sommerville for checking answers to the exercises. We wish to thank Sharon Sullivan and Dawn Haines for their help with IATEX. Finally, we want to express appreciation for the critical reviews and many suggestions provided for the first edition by the Wadsworth Brooks/Cole reviewers: Mark Conaway, University of Iowa; Franklin Graybill, Colorado State University; Jason Hsu, Ohio State University; Kenneth Koehler, Iowa State University; B. Lindsay, The Pennsylvania State University; Michael Meridith, Cornell University; M. B. Rajarshi, University of Poona (India); Muni Srivastava, University of Toronto; and Patricia Wahl, University of Washington; and for the second edition by the Springer-Verlag reviewers.

Acknowledgment is given for the use of material in the appendix tables. Appendix Table A.7 is reproduced in part from Tables 4 and 6 of Durbin and Watson (1951) with permission of the Biometrika Trustees. Appendix Table A.8 is reproduced with permission from Shapiro and Francia (1972), *Journal of the American Statistical Association*. The remaining appendix tables have been computer generated by one of the authors. We gratefully acknowledge permission of other authors and publishers for use of material from their publications as noted in the text.

Note to the Reader

Most research is aimed at quantifing relationships among variables that either measure the end result of some process or are likely to affect the process. The process in question may be any biological, chemical, or physical process of interest to the scientist. The quantification of the process may be as simple as determining the degree of association between two variables or as complicated as estimating the many parameters of a very detailed nonlinear mathematical model of the system.

Regardless of the degree of sophistication of the model, the most commonly used statistical method for estimating the parameters of interest is the method of **least squares**. The criterion applied in least squares estimation is simple and has great intuitive appeal. The researcher chooses the model that is believed to be most appropriate for the project at hand. The parameters for the model are then estimated such that the predictions from the model and the observed data are in as good agreement as possible as measured by the **least squares criterion**, minimization of the sum of squared differences between the predicted and the observed points.

Least squares estimation is a powerful research tool. Few assumptions are required and the estimators obtained have several desirable properties. Inference from research data to the true behavior of a process, however, can be a difficult and dangerous step due to unrecognized inadequacies in the data, misspecification of the model, or inappropriate inferences of

x PREFACE

causality. As with any research tool it is important that the least squares method be thoroughly understood in order to eliminate as much misuse or misinterpretation of the results as possible. There is a distinct difference between understanding and pure memorization. Memorization can make a good technician, but it takes understanding to produce a master. A discussion of the **geometric interpretation** of least squares is given to enhance your understanding. You may find your first exposure to the geometry of least squares somewhat traumatic but the visual perception of least squares is worth the effort. We encourage you to tackle the topic in the spirit in which it is included.

The general topic of least squares has been broadened to include statistical techniques associated with **model development and testing**. The backbone of least squares is the classical multiple regression analysis using the linear model to relate several independent variables to a response or dependent variable. Initially, this classical model is assumed to be appropriate. Then methods for detecting inadequacies in this model and possible remedies are discussed.

The connection between the analysis of variance for designed experiments and multiple regression is developed to build the foundation for the analysis of **unbalanced data**. (This also emphasizes the generality of the least squares method.) Interpretation of unbalanced data is difficult. It is important that the application of least squares to the analysis of such data be understood if the results from computer programs designed for the analysis of unbalanced data are to be used correctly.

The objective of a research project determines the amount of effort to be devoted to the development of realistic models. If the intent is one of prediction only, the degree to which the model might be considered realistic is immaterial. The only requirement is that the predictions be adequately precise in the region of interest. On the other hand, realism is of primary importance if the goal is a thorough understanding of the system. The simple linear additive model can seldom be regarded as a realistic model. It is at best an approximation of the true model. Almost without exception, models developed from the basic principles of a process will be nonlinear in the parameters. The least squares estimation principle is still applicable but the mathematical methods become much more difficult. You are introduced to **nonlinear least squares regression methods** and some of the more common nonlinear models.

Least squares estimation is controlled by the correlational structure observed among the independent and dependent variables in the data set. Observational data, data collected by observing the state of nature according to some sampling plan, will frequently cause special problems for least squares estimation because of strong correlations or, more generally, near-linear dependencies among the independent variables. The seriousness of the problems will depend on the use to be made of the analyses. Understanding the correlational structure of the data is most helpful in interpreting regression results and deciding what inferences might be made. Principal component analysis is introduced as an aid in characterizing the correlational structure of the data. A graphical procedure, Gabriel's biplot, is introduced to help visualize the correlational structure. Principal component analysis also serves as an introduction to **biased regression methods**. Biased regression methods are designed to alleviate the deleterious effects of near-linear dependencies (among the independent variables) on ordinary least squares estimation.

Least squares estimation is a powerful research tool and, with modern low cost computers, is readily available. This ease of access, however, also facilitates misuse. Proper use of least squares requires an understanding of the basic method and assumptions on which it is built, and an awareness of the possible problems and their remedies. In some cases, alternative methods to least squares estimation might be more appropriate. It is the intent of this text to convey the basic understanding that will allow you to use least squares as an effective research tool.

The data sets used in this text are available on the internet at

http://www.stat.ncsu.edu/publications/rawlings/applied_least_squares or through a link at the Springer-Verlag page. The "readme" file explains the contents of each data set.

Raleigh, North Carolina March 4, 1998 John O. Rawlings Sastry G. Pantula David A. Dickey

CONTENTS

PREFACE			vii
1	REV	VIEW OF SIMPLE REGRESSION	1 2 3
	1.1	The Linear Model and Assumptions	2
	1.2	Least Squares Estimation	3
	1.3	Predicted Values and Residuals	6
	1.4	Analysis of Variation in the Dependent Variable	7
	1.5	Precision of Estimates	11
	1.6	Tests of Significance and Confidence Intervals	16
	1.7	Regression Through the Origin	21
	1.8	Models with Several Independent Variables	27
	1.9	Violation of Assumptions	28
	1.10	Summary	29
	1.11	Exercises	30
2	INT	RODUCTION TO MATRICES	37
	2.1	Basic Definitions	37
	2.2	Special Types of Matrices	39
	2.3	Matrix Operations	40
	2.4	Geometric Interpretations of Vectors	46
	2.5	Linear Equations and Solutions	50
	2.6	Orthogonal Transformations and Projections	54
	2.7	Eigenvalues and Eigenvectors	57
	2.8	Singular Value Decomposition	60

	2.9	Summary	68
	2.10	Exercises	68
3	MU	LTIPLE REGRESSION IN MATRIX NOTATION	75
	3.1	The Model	75
	3.2	The Normal Equations and Their Solution $\ldots \ldots \ldots$	78
	3.3	The Y and Residuals vectors	80
	3.4	Properties of Linear Functions of Random Vectors	82
	3.5	Properties of Regression Estimates	87
	3.6	Summary of Matrix Formulae	92
	3.7	Exercises	93
4	AN	ALYSIS OF VARIANCE	
	AN	D QUADRATIC FORMS	101
	4.1	Introduction to Quadratic Forms	102
	4.2	Analysis of Variance	107
	4.3	Expectations of Quadratic Forms	113
	4.4	Distribution of Quadratic Forms	115
	4.5	General Form for Hypothesis Testing	119
		4.5.1 The General Linear Hypothesis	119
		4.5.2 Special Cases of the General Form	121
		4.5.3 A Numerical Example	122
		4.5.4 Computing Q from Differences in Sums of Squares .	126
		4.5.5 The <i>R</i> -Notation to Label Sums of Squares	129
		4.5.6 Example: Sequential and Partial Sums of Squares	133
	4.6	Univariate and Joint Confidence Regions	135
		4.6.1 Univariate Confidence Intervals	135
		4.6.2 Simultaneous Confidence Statements	137
		4.6.3 Joint Confidence Regions	139
	4.7	Estimation of Pure Error	143
	4.8	Exercises	149
5	CAS	SE STUDY: FIVE INDEPENDENT VARIABLES	161
	5.1	Spartina Biomass Production in the Cape Fear Estuary	161
	5.2	Regression Analysis for the Full Model	162
		5.2.1 The Correlation Matrix	164
		5.2.2 Multiple Regression Results: Full Model	165
	5.3	Simplifying the Model	167
	5.4	Results of the Final Model	170
	5.5	General Comments	177
	5.6	Exercises	179
6	GEO	OMETRY OF LEAST SQUARES	183
3	6.1	Linear Model and Solution	184
	6.2	Sums of Squares and Degrees of Freedom	189
		1 0	

	6.3	Reparameterization	192
	6.4	Sequential Regressions	196
	6.5	The Collinearity Problem	197
	6.6	Summary	201
	6.7	Exercises	201
7	MC	DEL DEVELOPMENT: VARIABLE SELECTION	205
	7.1	Uses of the Regression Equation	206
	7.2	Effects of Variable Selection on Least Squares	208
	7.3	All Possible Regressions	210
	7.4	Stepwise Regression Methods	213
	7.5	Criteria for Choice of Subset Size	220
		7.5.1 Coefficient of Determination	220
		7.5.2 Residual Mean Square	222
		7.5.3 Adjusted Coefficient of Determination	222
		7.5.4 Mallows' C_p Statistic	223
		7.5.5 Information Criteria: AIC and SBC	225
		7.5.6 "Significance Levels" for Choice of Subset Size	226
	7.6	Model Validation	228
	7.7	Exercises	231
8	РО	LYNOMIAL REGRESSION	235
	8.1	Polynomials in One Variable	236
	8.2	Trigonometric Regression Models	245
	8.3	Response Curve Modeling	249
		8.3.1 Considerations in Specifying the Functional Form	249
		8.3.2 Polynomial Response Models	250
	8.4	Exercises	262
9	CL	ASS VARIABLES IN REGRESSION	269
	9.1	Description of Class Variables	270
	9.2	The Model for One-Way Structured Data	271
	9.3	Reparameterizing to Remove Singularities	273
		9.3.1 Reparameterizing with the Means Model	274
		9.3.2 Reparameterization Motivated by $\sum \tau_i = 0$	277
		9.3.3 Reparameterization Motivated by $\overline{\tau_t} = 0$	279
		9.3.4 Reparameterization: A Numerical Example	280
	9.4	Generalized Inverse Approach	282
	9.5	The Model for Two-Way Classified Data	284
	9.6	Class Variables To Test Homogeneity of Regressions	288
	9.7	Analysis of Covariance	294
	9.8	Numerical Examples	300
		9.8.1 Analysis of Variance	301
		9.8.2 Test of Homogeneity of Regression Coefficients	306
		9.8.3 Analysis of Covariance	307

	9.9	Exercises	316
10	PRO	OBLEM AREAS IN LEAST SQUARES	325
	10.1	Nonnormality	326
	10.2	Heterogeneous Variances	328
	10.3	Correlated Errors	329
	10.4	Influential Data Points and Outliers	330
	10.5	Model Inadequacies	332
	10.6	The Collinearity Problem	333
	10.7	Errors in the Independent Variables	334
	10.8	Summary	339
	10.9	Exercises	339
11	REC	GRESSION DIAGNOSTICS	341
	11.1	Residuals Analysis	342
		11.1.1 Plot of \hat{e} Versus \hat{Y}	346
		11.1.2 Plots of e Versus X_i	350
		11.1.3 Plots of e Versus Time	351
		11.1.4 Plots of e_i Versus e_{i-1}	354
		11.1.5 Normal Probability Plots	356
		11.1.6 Partial Regression Leverage Plots	359
	11.2	Influence Statistics	361
		11.2.1 Cook's D	362
		11.2.2 DFFITS	363
		11.2.3 DFBETAS	364
		11.2.4 COVRATIO	364
		11.2.5 Summary of Influence Measures	367
	11.3	Collinearity Diagnostics	369
		11.3.1 Condition Number and Condition Index	371
		11.3.2 Variance Inflation Factor	372
		11.3.3 Variance Decomposition Proportions	373
		11.3.4 Summary of Collinearity Diagnostics	377
	11.4	Regression Diagnostics on the Linthurst Data	377
		11.4.1 Plots of Residuals	378
		11.4.2 Influence Statistics	388
		11.4.3 Collinearity Diagnostics	391
	11.5	Exercises	392
12	TRA	ANSFORMATION OF VARIABLES	397
	12.1	Reasons for Making Transformations	397
	12.2	Transformations to Simplify Relationships	399
	12.3	Transformations to Stabilize Variances	407
	12.4	Transformations to Improve Normality	409
	12.5	Generalized Least Squares	411
		12.5.1 Weighted Least Squares	414

12.5.2 Generalized Least Squares	417
12.6 Summary	426
12.7 Exercises	. 427
13 COLLINEARITY	433
13.1 Understanding the Structure of the X-Space	435
13.2 Biased Regression	443
13.2.1 Explanation \ldots	443
13.2.2 Principal Component Regression	446
13.3 General Comments on Collinearity	457
13.4 Summary	. 459
13.5 Exercises	. 459
14 CASE STUDY: COLLINEARITY PROBLEMS	463
14.1 The Problem	. 463
14.2 Multiple Regression: Ordinary Least Squares	. 467
14.3 Analysis of the Correlational Structure	. 471
14.4 Principal Component Regression	. 479
14.5 Summary	. 482
14.6 Exercises	. 483
15 MODELS NONLINEAR IN THE PARAMETERS	485
15.1 Examples of Nonlinear Models	. 486
15.2 Fitting Models Nonlinear in the Parameters \ldots	. 494
15.3 Inference in Nonlinear Models	. 498
15.4 Violation of Assumptions	. 507
15.4.1 Heteroscedastic Errors	. 507
15.4.2 Correlated Errors	. 509
15.5 Logistic Regression	. 509
15.6 Exercises	. 511
16 CASE STUDY: RESPONSE CURVE MODELING	515
16.1 The Ozone–Sulfur Dioxide Response Surface (1981)	. 517
16.1.1 Polynomial Response Model	. 520
16.1.2 Nonlinear Weibull Response Model	. 524
16.2 Analysis of the Combined Soybean Data	. 530
16.3 Exercises	. 543
17 ANALYSIS OF UNBALANCED DATA	545
17.1 Sources Of Imbalance	. 546
17.2 Effects Of Imbalance	. 547
17.3 Analysis of Cell Means	. 549
17.4 Linear Models for Unbalanced Data	. 553
17.4.1 Estimable Functions with Balanced Data \ldots	. 554
17.4.2 Estimable Functions with Unbalanced Data	558

xviii CONTENTS

	17.4.3 Least Squares Means	564	
	17.5 Exercises	568	
18	MIXED EFFECTS MODELS	573	
	18.1 Random Effects Models	574	
	18.2 Fixed and Random Effects	579	
	18.3 Random Coefficient Regression Models	584	
	18.4 General Mixed Linear Models	586	
	18.5 Exercises	589	
10	CASE STUDY. ANALYSIS OF UNDALANCED DATA	502	
19	10.1 The Analysis Of Variance	506	
		090	
	19.2 Mean Square Expectations and Choice of Errors	607	
	19.3 Least Squares Means and Standard Errors	610	
	19.4 Mixed Model Analysis	615	
	19.5 Exercises	618	
\mathbf{A}	APPENDIX TABLES	621	
RF	FERENCES	635	
AU	THOR INDEX	647	
su	SUBJECT INDEX		