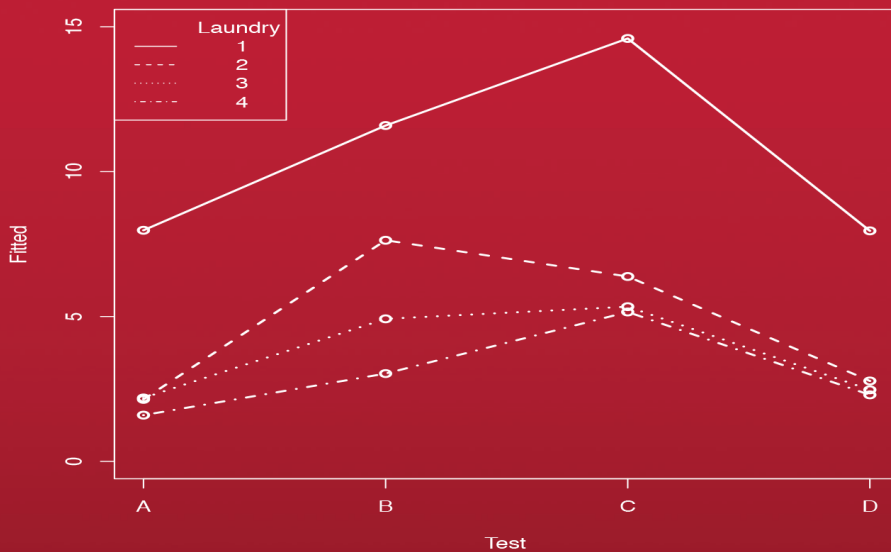


Texts in Statistical Science

Analysis of Variance, Design, and Regression

Linear Modeling for
Unbalanced Data

Second Edition



Ronald Christensen



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Analysis of Variance, Design, and Regression

**Linear Modeling for
Unbalanced Data**

Second Edition

CHAPMAN & HALL/CRC
Texts in Statistical Science Series

Series Editors

Francesca Dominici, *Harvard School of Public Health, USA*

Julian J. Faraway, *University of Bath, UK*

Martin Tanner, *Northwestern University, USA*

Jim Zidek, *University of British Columbia, Canada*

Statistical Theory: A Concise Introduction

F. Abramovich and Y. Ritov

Practical Multivariate Analysis, Fifth Edition

A. Afifi, S. May, and V.A. Clark

Practical Statistics for Medical Research

D.G. Altman

**Interpreting Data: A First Course
in Statistics**

A.J.B. Anderson

Introduction to Probability with R

K. Baclawski

**Linear Algebra and Matrix Analysis for
Statistics**

S. Banerjee and A. Roy

**Mathematical Statistics: Basic Ideas and
Selected Topics, Volume I, Second Edition**

P.J. Bickel and K. A. Doksum

**Mathematical Statistics: Basic Ideas and
Selected Topics, Volume II**

P.J. Bickel and K. A. Doksum

Analysis of Categorical Data with R

C. R. Bilder and T. M. Loughin

Statistical Methods for SPC and TQM

D. Bissell

Introduction to Probability

J. K. Blitzstein and J. Hwang

**Bayesian Methods for Data Analysis,
Third Edition**

B.P. Carlin and T.A. Louis

Second Edition

R. Caulcutt

**The Analysis of Time Series: An Introduction,
Sixth Edition**

C. Chatfield

Introduction to Multivariate Analysis

C. Chatfield and A.J. Collins

**Problem Solving: A Statistician's Guide,
Second Edition**

C. Chatfield

**Statistics for Technology: A Course in Applied
Statistics, Third Edition**

C. Chatfield

**Analysis of Variance, Design, and Regression :
Linear Modeling for Unbalanced Data, Second
Edition**

R. Christensen

**Bayesian Ideas and Data Analysis: An
Introduction for Scientists and Statisticians**

R. Christensen, W. Johnson, A. Branscum,
and T.E. Hanson

Modelling Binary Data, Second Edition

D. Collett

**Modelling Survival Data in Medical Research,
Third Edition**

D. Collett

**Introduction to Statistical Methods for
Clinical Trials**

T.D. Cook and D.L. DeMets

Applied Statistics: Principles and Examples

D.R. Cox and E.J. Snell

**Multivariate Survival Analysis and Competing
Risks**

M. Crowder

Statistical Analysis of Reliability Data

M.J. Crowder, A.C. Kimber,
T.J. Sweeting, and R.L. Smith

**An Introduction to Generalized
Linear Models, Third Edition**

A.J. Dobson and A.G. Barnett

**Nonlinear Time Series: Theory, Methods, and
Applications with R Examples**

R. Douc, E. Moulines, and D.S. Stoffer

**Introduction to Optimization Methods and
Their Applications in Statistics**

B.S. Everitt

**Extending the Linear Model with R:
Generalized Linear, Mixed Effects and
Nonparametric Regression Models**

J.J. Faraway

Linear Models with R, Second Edition

J.J. Faraway

A Course in Large Sample Theory

T.S. Ferguson

Multivariate Statistics: A Practical Approach

B. Flury and H. Riedwyl

Readings in Decision Analysis

S. French

Markov Chain Monte Carlo:

Stochastic Simulation for Bayesian Inference, Second Edition

D. Gamerman and H.F. Lopes

Bayesian Data Analysis, Third Edition

A. Gelman, J.B. Carlin, H.S. Stern, D.B. Dunson, A. Vehtari, and D.B. Rubin

Multivariate Analysis of Variance and Repeated Measures: A Practical Approach for Behavioural Scientists

D.J. Hand and C.C. Taylor

Practical Longitudinal Data Analysis

D.J. Hand and M. Crowder

Logistic Regression Models

J.M. Hilbe

Richly Parameterized Linear Models: Additive, Time Series, and Spatial Models Using Random Effects

J.S. Hodges

Statistics for Epidemiology

N.P. Jewell

Stochastic Processes: An Introduction, Second Edition

P.W. Jones and P. Smith

The Theory of Linear Models

B. Jørgensen

Principles of Uncertainty

J.B. Kadane

Graphics for Statistics and Data Analysis with R

K.J. Keen

Mathematical Statistics

K. Knight

Introduction to Multivariate Analysis: Linear and Nonlinear Modeling

S. Konishi

Nonparametric Methods in Statistics with SAS Applications

O. Korosteleva

Modeling and Analysis of Stochastic Systems, Second Edition

V.G. Kulkarni

Exercises and Solutions in Biostatistical Theory

L.L. Kupper, B.H. Neelon, and S.M. O'Brien

Exercises and Solutions in Statistical Theory

L.L. Kupper, B.H. Neelon, and S.M. O'Brien

Design and Analysis of Experiments with R

J. Lawson

Design and Analysis of Experiments with SAS

J. Lawson

A Course in Categorical Data Analysis

T. Leonard

Statistics for Accountants

S. Letchford

Introduction to the Theory of Statistical Inference

H. Liero and S. Zwanzig

Statistical Theory, Fourth Edition

B.W. Lindgren

Stationary Stochastic Processes: Theory and Applications

G. Lindgren

Statistics for Finance

E. Lindström, H. Madsen, and J. N. Nielsen

The BUGS Book: A Practical Introduction to Bayesian Analysis

D. Lunn, C. Jackson, N. Best, A. Thomas, and D. Spiegelhalter

Introduction to General and Generalized Linear Models

H. Madsen and P. Thyregod

Time Series Analysis

H. Madsen

Pólya Urn Models

H. Mahmoud

Randomization, Bootstrap and Monte Carlo Methods in Biology, Third Edition

B.F.J. Manly

Introduction to Randomized Controlled Clinical Trials, Second Edition

J.N.S. Matthews

Statistical Rethinking: A Bayesian Course with Examples in R and Stan

R. McElreath

Statistical Methods in Agriculture and Experimental Biology, Second Edition

R. Mead, R.N. Curnow, and A.M. Hasted

Statistics in Engineering: A Practical Approach

A.V. Metcalfe

Statistical Inference: An Integrated Approach, Second Edition

H. S. Migon, D. Gamerman, and F. Louzada

Beyond ANOVA: Basics of Applied Statistics

R.G. Miller, Jr.

A Primer on Linear Models

J.F. Monahan

Applied Stochastic Modelling, Second Edition

B.J.T. Morgan

Elements of Simulation

B.J.T. Morgan

Probability: Methods and Measurement

A. O'Hagan

Introduction to Statistical Limit Theory

A.M. Polansky

Applied Bayesian Forecasting and Time Series Analysis

A. Pole, M. West, and J. Harrison

Statistics in Research and Development, Time Series: Modeling, Computation, and Inference

R. Prado and M. West

Introduction to Statistical Process Control

P. Qiu

Sampling Methodologies with Applications

P.S.R.S. Rao

A First Course in Linear Model Theory

N. Ravishanker and D.K. Dey

Essential Statistics, Fourth Edition

D.A.G. Rees

Stochastic Modeling and Mathematical Statistics: A Text for Statisticians and Quantitative Scientists

F.J. Samaniego

Statistical Methods for Spatial Data Analysis

O. Schabenberger and C.A. Gotway

Bayesian Networks: With Examples in R

M. Scutari and J.-B. Denis

Large Sample Methods in Statistics

P.K. Sen and J. da Motta Singer

Spatio-Temporal Methods in Environmental Epidemiology

G. Shaddick and J.V. Zidek

Decision Analysis: A Bayesian Approach

J.Q. Smith

Analysis of Failure and Survival Data

P.J. Smith

Applied Statistics: Handbook of GENSTAT Analyses

E.J. Snell and H. Simpson

Applied Nonparametric Statistical Methods, Fourth Edition

P. Sprent and N.C. Smeeton

Data Driven Statistical Methods

P. Sprent

Generalized Linear Mixed Models: Modern Concepts, Methods and Applications

W. W. Stroup

Survival Analysis Using S: Analysis of Time-to-Event Data

M. Tableman and J.S. Kim

Applied Categorical and Count Data Analysis

W. Tang, H. He, and X.M. Tu

Elementary Applications of Probability Theory, Second Edition

H.C. Tuckwell

Introduction to Statistical Inference and Its Applications with R

M.W. Trosset

Understanding Advanced Statistical Methods

P.H. Westfall and K.S.S. Henning

Statistical Process Control: Theory and Practice, Third Edition

G.B. Wetherill and D.W. Brown

Generalized Additive Models:

An Introduction with R

S. Wood

Epidemiology: Study Design and Data Analysis, Third Edition

M. Woodward

Practical Data Analysis for Designed Experiments

B.S. Yandell

Texts in Statistical Science

Analysis of Variance, Design, and Regression

**Linear Modeling for
Unbalanced Data**

Second Edition

Ronald Christensen

University of New Mexico
Albuquerque, USA



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an informa business
A CHAPMAN & HALL BOOK

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2016 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20151221

International Standard Book Number-13: 978-1-4987-7405-5 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

To Mark, Karl, and John

It was great fun.

Contents

Preface	xvii
Edited Preface to First Edition	xxi
Computing	xxv
1 Introduction	1
1.1 Probability	1
1.2 Random variables and expectations	4
1.2.1 Expected values and variances	6
1.2.2 Chebyshev's inequality	9
1.2.3 Covariances and correlations	10
1.2.4 Rules for expected values and variances	12
1.3 Continuous distributions	13
1.4 The binomial distribution	17
1.4.1 Poisson sampling	21
1.5 The multinomial distribution	21
1.5.1 Independent Poissons and multinomials	23
1.6 Exercises	24
2 One Sample	27
2.1 Example and introduction	27
2.2 Parametric inference about μ	31
2.2.1 Significance tests	34
2.2.2 Confidence intervals	37
2.2.3 P values	38
2.3 Prediction intervals	39
2.4 Model testing	42
2.5 Checking normality	43
2.6 Transformations	48
2.7 Inference about σ^2	51
2.7.1 Theory	54
2.8 Exercises	55
3 General Statistical Inference	57
3.1 Model-based testing	58
3.1.1 An alternative F test	64
3.2 Inference on single parameters: assumptions	64
3.3 Parametric tests	66
3.4 Confidence intervals	70
3.5 P values	72
3.6 Validity of tests and confidence intervals	75
3.7 Theory of prediction intervals	75

3.8	Sample size determination and power	78
3.9	The shape of things to come	80
3.10	Exercises	85
4	Two Samples	87
4.1	Two correlated samples: Paired comparisons	87
4.2	Two independent samples with equal variances	90
4.2.1	Model testing	95
4.3	Two independent samples with unequal variances	96
4.4	Testing equality of the variances	101
4.5	Exercises	104
5	Contingency Tables	109
5.1	One binomial sample	109
5.1.1	The sign test	112
5.2	Two independent binomial samples	112
5.3	One multinomial sample	115
5.4	Two independent multinomial samples	117
5.5	Several independent multinomial samples	120
5.6	Lancaster–Irwin partitioning	123
5.7	Exercises	129
6	Simple Linear Regression	133
6.1	An example	133
6.1.1	Computer commands	137
6.2	The simple linear regression model	139
6.3	The analysis of variance table	141
6.4	Model-based inference	141
6.5	Parametric inferential procedures	143
6.6	An alternative model	145
6.7	Correlation	146
6.8	Two-sample problems	147
6.9	A multiple regression	148
6.10	Estimation formulae for simple linear regression	149
6.11	Exercises	154
7	Model Checking	157
7.1	Recognizing randomness: Simulated data with zero correlation	157
7.2	Checking assumptions: Residual analysis	159
7.2.1	Another example	163
7.2.2	Outliers	165
7.2.3	Effects of high leverage	166
7.3	Transformations	168
7.3.1	Circle of transformations	168
7.3.2	Box–Cox transformations	171
7.3.3	Constructed variables	174
7.4	Exercises	177

8	Lack of Fit and Nonparametric Regression	179
8.1	Polynomial regression	179
8.1.1	Picking a polynomial	181
8.1.2	Exploring the chosen model	183
8.2	Polynomial regression and leverages	185
8.3	Other basis functions	189
8.3.1	High-order models	191
8.4	Partitioning methods	191
8.4.1	Fitting the partitioned model	192
8.4.2	Output for categorical predictors*	194
8.4.3	Utts' method	196
8.5	Splines	198
8.6	Fisher's lack-of-fit test	200
8.7	Exercises	201
9	Multiple Regression: Introduction	205
9.1	Example of inferential procedures	205
9.1.1	Computing commands	209
9.1.2	General statement of the multiple regression model	210
9.2	Regression surfaces and prediction	211
9.3	Comparing regression models	213
9.3.1	General discussion	214
9.4	Sequential fitting	216
9.5	Reduced models and prediction	218
9.6	Partial correlation coefficients and added variable plots	219
9.7	Collinearity	221
9.8	More on model testing	223
9.9	Additive effects and interaction	227
9.10	Generalized additive models	229
9.11	Final comment	230
9.12	Exercises	230
10	Diagnostics and Variable Selection	235
10.1	Diagnostics	235
10.2	Best subset model selection	240
10.2.1	R^2 statistic	241
10.2.2	Adjusted R^2 statistic	243
10.2.3	Mallows's C_p statistic	244
10.2.4	A combined subset selection table	245
10.3	Stepwise model selection	246
10.3.1	Backwards elimination	246
10.3.2	Forward selection	247
10.3.3	Stepwise methods	248
10.4	Model selection and case deletion	248
10.5	Lasso regression	250
10.6	Exercises	252

11 Multiple Regression: Matrix Formulation	255
11.1 Random vectors	255
11.2 Matrix formulation of regression models	256
11.2.1 Simple linear regression in matrix form	256
11.2.2 The general linear model	258
11.3 Least squares estimation of regression parameters	262
11.4 Inferential procedures	266
11.5 Residuals, standardized residuals, and leverage	269
11.6 Principal components regression	270
11.7 Exercises	274
12 One-Way ANOVA	277
12.1 Example	277
12.1.1 Inferences on a single group mean	281
12.1.2 Inference on pairs of means	281
12.1.3 Inference on linear functions of means	283
12.1.4 Testing $\mu_1 = \mu_2 = \mu_3$	284
12.2 Theory	284
12.2.1 Analysis of variance tables	289
12.3 Regression analysis of ANOVA data	290
12.3.1 Testing a pair of means	292
12.3.2 Model testing	293
12.3.3 Another choice	296
12.4 Modeling contrasts	297
12.4.1 A hierarchical approach	298
12.4.2 Evaluating the hierarchy	299
12.4.3 Regression analysis	303
12.4.4 Relation to orthogonal contrasts	303
12.4.5 Theory: Difficulties in general unbalanced analyses	303
12.5 Polynomial regression and one-way ANOVA	304
12.5.1 Fisher's lack-of-fit test	310
12.5.2 More on R^2	313
12.6 Weighted least squares	314
12.6.1 Theory	316
12.7 Exercises	317
13 Multiple Comparison Methods	323
13.1 "Fisher's" least significant difference method	324
13.2 Bonferroni adjustments	326
13.3 Scheffé's method	328
13.4 Studentized range methods	330
13.4.1 Tukey's honest significant difference	331
13.5 Summary of multiple comparison procedures	332
13.6 Exercises	332
14 Two-Way ANOVA	335
14.1 Unbalanced two-way analysis of variance	335
14.1.1 Initial analysis	336
14.1.2 Hierarchy of models	339
14.1.3 Computing issues	340
14.1.4 Discussion of model fitting	341
14.1.5 Diagnostics	342

CONTENTS	xiii
14.1.6 Outlier deleted analysis	342
14.2 Modeling contrasts	346
14.2.1 Nonequivalence of tests	347
14.3 Regression modeling	349
14.4 Homologous factors	351
14.4.1 Symmetric additive effects	351
14.4.2 Skew symmetric additive effects	353
14.4.3 Symmetry	355
14.4.4 Hierarchy of models	357
14.5 Exercises	357
15 ACOVA and Interactions	361
15.1 One covariate example	361
15.1.1 Additive regression effects	362
15.1.2 Interaction models	364
15.1.3 Multiple covariates	369
15.2 Regression modeling	369
15.2.1 Using overparameterized models	370
15.3 ACOVA and two-way ANOVA	371
15.3.1 Additive effects	372
15.4 Near replicate lack-of-fit tests	375
15.5 Exercises	377
16 Multifactor Structures	379
16.1 Unbalanced three-factor analysis of variance	379
16.1.1 Computing	383
16.1.2 Regression fitting	385
16.2 Balanced three-factors	386
16.3 Higher-order structures	393
16.4 Exercises	393
17 Basic Experimental Designs	397
17.1 Experiments and causation	397
17.2 Technical design considerations	399
17.3 Completely randomized designs	401
17.4 Randomized complete block designs	401
17.4.1 Paired comparisons	405
17.5 Latin square designs	406
17.5.1 Latin square models	407
17.5.2 Discussion of Latin squares	407
17.6 Balanced incomplete block designs	408
17.6.1 Special cases	410
17.7 Youden squares	412
17.7.1 Balanced lattice squares	412
17.8 Analysis of covariance in designed experiments	413
17.9 Discussion of experimental design	415
17.10 Exercises	416

18 Factorial Treatments	421
18.1 Factorial treatment structures	421
18.2 Analysis	422
18.3 Modeling factorials	424
18.4 Interaction in a Latin square	425
18.5 A balanced incomplete block design	429
18.6 Extensions of Latin squares	433
18.7 Exercises	436
19 Dependent Data	439
19.1 The analysis of split-plot designs	439
19.1.1 Modeling with interaction	446
19.2 A four-factor example	450
19.2.1 Unbalanced subplot analysis	452
19.2.2 Whole-plot analysis	456
19.2.3 Fixing effect levels	459
19.2.4 Final models and estimates	460
19.3 Multivariate analysis of variance	463
19.4 Random effects models	472
19.4.1 Subsampling	473
19.4.2 Random effects	474
19.5 Exercises	477
20 Logistic Regression: Predicting Counts	481
20.1 Models for binomial data	481
20.2 Simple linear logistic regression	484
20.2.1 Goodness-of-fit tests	485
20.2.2 Assessing predictive ability	486
20.2.3 Case diagnostics	488
20.3 Model testing	489
20.4 Fitting logistic models	490
20.5 Binary data	493
20.5.1 Goodness-of-fit tests	494
20.5.2 Case diagnostics	496
20.5.3 Assessing predictive ability	496
20.6 Multiple logistic regression	497
20.7 ANOVA type logit models	505
20.8 Ordered categories	507
20.9 Exercises	510
21 Log-Linear Models: Describing Count Data	513
21.1 Models for two-factor tables	514
21.1.1 Lancaster–Irwin partitioning	514
21.2 Models for three-factor tables	515
21.2.1 Testing models	517
21.3 Estimation and odds ratios	518
21.4 Higher-dimensional tables	520
21.5 Ordered categories	522
21.6 Offsets	525
21.7 Relation to logistic models	526
21.8 Multinomial responses	528
21.9 Logistic discrimination and allocation	530

CONTENTS	xv
21.10 Exercises	535
22 Exponential and Gamma Regression: Time-to-Event Data	537
22.1 Exponential regression	538
22.1.1 Computing issues	540
22.2 Gamma regression	541
22.2.1 Computing issues	543
22.3 Exercises	543
23 Nonlinear Regression	545
23.1 Introduction and examples	545
23.2 Estimation	546
23.2.1 The Gauss–Newton algorithm	547
23.2.2 Maximum likelihood estimation	551
23.3 Statistical inference	551
23.4 Linearizable models	559
23.5 Exercises	560
Appendix A: Matrices and Vectors	563
A.1 Matrix addition and subtraction	564
A.2 Scalar multiplication	564
A.3 Matrix multiplication	564
A.4 Special matrices	566
A.5 Linear dependence and rank	567
A.6 Inverse matrices	568
A.7 A list of useful properties	570
A.8 Eigenvalues and eigenvectors	570
Appendix B: Tables	573
B.1 Tables of the t distribution	574
B.2 Tables of the χ^2 distribution	576
B.3 Tables of the W' statistic	580
B.4 Tables of the Studentized range	581
B.5 The Greek alphabet	585
B.6 Tables of the F distribution	586
References	599

Preface

Background

Big Data are the future of Statistics. The electronic revolution has increased exponentially our ability to measure things. A century ago, data were hard to come by. Statisticians put a premium on extracting every bit of information that the data contained. Now data are easy to collect; the problem is sorting through them to find meaning. To a large extent, this happens in two ways: doing a crude analysis on a massive amount of data or doing a careful analysis on the moderate amount of data that were isolated from the massive data as being meaningful. It is quite literally impossible to analyze a million data points as carefully as one can analyze a hundred data points, so “crude” is not a pejorative term but rather a fact of life.

The fundamental tools used in analyzing data have been around a long time. It is the emphases and the opportunities that have changed. With thousands of observations, we don’t need a perfect statistical analysis to detect a large effect. But with thousands of observations, we might look for subtle effects that we never bothered looking for before, and such an analysis must be done carefully—as must any analysis in which only a small part of the massive data are relevant to the problem at hand. The electronic revolution has also provided us with the opportunity to perform data analysis procedures that were not practical before, but in my experience, the new procedures (often called *machine learning*), are sophisticated applications of fundamental tools.

This book explains some of the fundamental tools and the ideas needed to adapt them to big data. It is not a book that analyzes big data. The book analyzes small data sets carefully but by using tools that 1) can easily be scaled to large data sets or 2) apply to the haphazard way in which small relevant data sets are now constructed. Personally, I believe that it is not safe to apply models to large data sets until you understand their implications for small data. There is also a major emphasis on tools that look for subtle effects (interactions, homologous effects) that are hard to identify.

The fundamental tools examined here are linear structures for modeling data; specifically, how to incorporate specific ideas about the structure of the data into the model for the data. Most of the book is devoted to adapting linear structures (regression, analysis of variance, analysis of covariance) to examine measurement (continuous) data. But the exact same methods apply to either-or (Yes/No, binomial) data, count (Poisson, multinomial) data, and time-to-event (survival analysis, reliability) data. The book also places strong emphasis on foundational issues, e.g., the meaning of significance tests and the interval estimates associated with them; the difference between prediction and causation; and the role of randomization.

The platform for this presentation is the revision of a book I published in 1996, *Analysis of Variance, Design, and Regression: Applied Statistical Methods*. Within a year, I knew that the book was not what I thought needed to be taught in the 21st century, cf., Christensen (2000). This book, *Analysis of Variance, Design, and Regression: Linear Modeling of Unbalanced Data*, shares with the earlier book lots of the title, much of the data, and even some of the text, but the book is radically different. The original book focused greatly on balanced analysis of variance. This book focuses on modeling unbalanced data. As such, it generalizes much of the work in the previous book. The more general methods presented here agree with the earlier methods for balanced data. Another advantage of taking a modeling approach to unbalanced data is that by making the effort to treat unbalanced analysis of variance, one can easily handle a wide range of models for nonnormal data, because the same fundamental methods apply. To that end, I have included new chapters on logistic regression,

log-linear models, and time-to-event data. These are placed near the end of the book, not because they are less important, but because the real subject of the book is modeling with linear structures and the methods for measurement data carry over almost immediately.

In early versions of this edition I made extensive comparisons between the methods used here and the balanced ANOVA methods used in the 1996 book. In particular, I emphasized how the newer methods continue to give the same results as the earlier methods when applied to balanced data. While I have toned that down, comparisons still exist. In such comparisons, I do not repeat the details of the balanced analysis given in the earlier book. CRC Press/Chapman & Hall have been kind enough to let me place a version of the 1996 book on my website so that readers can explore the comparisons in detail. Another good thing about having the old book up is that it contains a chapter on confounding and fractional replications in 2^n factorials. I regret having to drop that chapter, but the discussion is based on contrasts for balanced ANOVA and did not really fit the theme of the current edition.

When I was in high school, my two favorite subjects were math and history. On a whim, I made the good choice to major in Math for my BA. I mention my interest in history to apologize (primarily in the same sense that C.S. Lewis was a Christian “apologist”) for using so much old data. Unless you are trying to convince 18-year-olds that Statistics is sexy, I don’t think the age of the data should matter.

I need to thank Adam Branscum, my coauthor on Christensen et al. (2010). Adam wrote the first drafts of Chapter 7 and Appendix C of that book. Adam’s work on Chapter 7 definitely influenced this work and Adam’s work on Appendix C is what got me programming in R. This is also a good time to thank the people who have most influenced my career: Wes Johnson, Ed Bedrick, Don Berry, Frank Martin, and the late, great Seymour Geisser. My colleague Yan Lu taught out of a prepublication version of the book, and, with her students, pointed out a number of issues. Generally, the first person whose opinions and help I sought was my son Fletcher.

After the effort to complete this book, I’m feeling as unbalanced as the data being analyzed.

Specifics

I think of the book as something to use in the traditional Master’s level year-long course on regression and analysis of variance. If one needed to actually separate the material into a regression course and an ANOVA course, the regression material is in Chapters 6–11 and 20–23. Chapters 12–19 are traditionally viewed as ANOVA. But I much prefer to use both regression and ANOVA ideas when examining the generalized linear models of Chapters 20–22. Well-prepared students could begin with Chapter 3 and skip to Chapter 6. By well-prepared, I tautologically mean students who are already familiar with Chapters 1, 2, 4, and 5.

For less well-prepared students, obviously I would start at the beginning and deemphasize the more difficult topics. This is what I have done when teaching data analysis to upper division Statistics students and graduate students from other fields. I have tried to isolate more difficult material into clearly delineated (sub)sections. In the first semester of such a course, I would skip the end of Chapter 8, include the beginning of Chapter 12, and let time and student interest determine how much of Chapters 9, 10, and 13 to cover. But the book wasn’t written to be a text for such a course; it is written to address unbalanced multi-factor ANOVA.

The book requires very little pre-knowledge of math, just algebra, but does require that one not be afraid of math. It does not perform calculus, but it discusses that integrals provide areas under curves and, in an appendix, gives the integral formulae for means and variances. It largely avoids matrix algebra but presents enough of it to enable the matrix approach to linear models to be introduced. For a regression-ANOVA course, I would supplement the material after Chapter 11 with occasional matrix arguments. Any material described as a regression approach to an ANOVA problem lends itself to matrix discussion.

Although the book starts at the beginning mathematically, it is not for the intellectually unsophisticated. By Chapter 2 it discusses the impreciseness of our concepts of populations and how

the deletion of outliers must change those concepts. Chapter 2 also discusses the “murky” transformation from a probability interval to a confidence interval and the differences between significance testing, Neyman–Pearson hypothesis testing, and Bayesian methods. Because a lot of these ideas are subtle, and because people learn best from specifics to generalities rather than the other way around, Chapter 3 reiterates much of Chapter 2 but for general linear models. Most of the remainder of the book can be viewed as the application of Chapter 3 to specific data structures. Well-prepared students could start with Chapter 3 despite occasional references made to results in the first two chapters.

Chapter 4 considers two-sample data. Perhaps its most unique feature is, contrary to what seems popular in introductory Statistics these days, the argument that testing equality of means for two independent samples provides much less information when the variances are different than when they are the same.

Chapter 5 exists because I believe that if you teach one- and two-sample continuous data problems, you have a duty to present their discrete data analogs. Having gone that far, it seemed silly to avoid analogs to one-way ANOVA. I do not find the one-way ANOVA F test for equal group means to be all that useful. Contrasts contain more interesting information. The last two sections of Chapter 5 contain, respectively, discrete data analogs to one-way ANOVA and a method of extracting information similar to contrasts.

Chapters 6, 7, and 8 provide tools for exploring the relationship between a single dependent variable and a single measurement (continuous) predictor. A key aspect of the discussion is that the methods in Chapters 7 and 8 extend readily to more general linear models, i.e., those involving categorical and/or multiple predictors. The title of Chapter 8 arises from my personal research interest in testing lack of fit for linear models and the recognition of its relationship to nonparametric regression.

Chapters 9, 10, and 11 examine features associated with multiple regression. Of particular note are new sections on modeling interaction through generalized additive models and on lasso regression. I consider these important concepts for serious students of Statistics. The last of these chapters is where the book’s use of matrices is focused. The discussion of principal component regression is located here, not because the discussion uses matrices, but because the discussion requires matrix knowledge to understand.

The rest of the book involves categorical predictor variables. In particular, *the material after Chapter 13 is the primary reason for writing this edition*. The first edition focused on multifactor balanced data and looking at contrasts, not only in main effects but contrasts within two- and three-factor interactions. This edition covers the same material for unbalanced data.

Chapters 12 and 13 cover one-way analysis of variance (ANOVA) models and multiple comparisons but with an emphasis on the ideas needed when examining multiple categorical predictors. Chapter 12 involves one categorical predictor much like Chapter 6 involved one continuous predictor.

Chapter 14 examines the use of two categorical predictors, i.e., two-way ANOVA. It also introduces the concept of homologous factors. Chapter 15 looks at models with one continuous and one categorical factor, analysis of covariance. Chapter 16 considers models with three categorical predictors.

Chapters 17 and 18 introduce the main ideas of experimental design. Chapter 17 introduces a wide variety of standard designs and concepts of design. Chapter 18 introduces the key idea of defining treatments with factorial structure. The unusual aspect of these chapters is that the analyses presented apply when data are missing from the original design.

Chapter 19 introduces the analysis of dependent data. The primary emphasis is on the analysis of split-plot models. A short discussion is also given of multivariate analysis. Both of these methods require groups of observations that are independent of other groups but that are dependent within the groups. Both methods require balance within the groups but the groups themselves can be unbalanced. Subsection 19.2.1 even introduces a method for dealing with unbalance within groups.

It seems to have become popular to treat fixed and random effects models as merely two options