

Springer Series in Statistics

Advisors:

P. Diggle, S. Fienberg, K. Krickeberg,
I. Olkin, N. Wermuth

Springer Series in Statistics

- Andersen/Borgan/Gill/Keiding*: Statistical Models Based on Counting Processes.
- Andrews/Herzberg*: Data: A Collection of Problems from Many Fields for the Student and Research Worker.
- Anscombe*: Computing in Statistical Science through APL.
- Berger*: Statistical Decision Theory and Bayesian Analysis, 2nd edition.
- Bolfarine/Zacks*: Prediction Theory for Finite Populations.
- Brémaud*: Point Processes and Queues: Martingale Dynamics.
- Brockwell/Davis*: Time Series: Theory and Methods, 2nd edition.
- Daley/Vere-Jones*: An Introduction to the Theory of Point Processes.
- Dzhaparidze*: Parameter Estimation and Hypothesis Testing in Spectral Analysis of Stationary Time Series.
- Fahrmeir/Tutz*: Multivariate Statistical Modelling Based on Generalized Linear Models.
- Farrell*: Multivariate Calculation.
- Federer*: Statistical Design and Analysis for Intercropping Experiments.
- Fienberg/Hoaglin/Kruskal/Tanur (Eds.)*: A Statistical Model: Frederick Mosteller's Contributions to Statistics, Science and Public Policy.
- Fisher/Sen*: The Collected Works of Wassily Hoeffding.
- Good*: Permutation Tests: A Practical Guide to Resampling Methods for Testing Hypotheses.
- Goodman/Kruskal*: Measures of Association for Cross Classifications.
- Grandell*: Aspects of Risk Theory.
- Hall*: The Bootstrap and Edgeworth Expansion.
- Härdle*: Smoothing Techniques: With Implementation in S.
- Hartigan*: Bayes Theory.
- Heyer*: Theory of Statistical Experiments.
- Jolliffe*: Principal Component Analysis.
- Kotz/Johnson (Eds.)*: Breakthroughs in Statistics Volume I.
- Kotz/Johnson (Eds.)*: Breakthroughs in Statistics Volume II.
- Kres*: Statistical Tables for Multivariate Analysis.
- Le Cam*: Asymptotic Methods in Statistical Decision Theory.
- Le Cam/Yang*: Asymptotics in Statistics: Some Basic Concepts.
- Longford*: Models for Uncertainty in Educational Testing.
- Manoukian*: Modern Concepts and Theorems of Mathematical Statistics.
- Miller, Jr.*: Simultaneous Statistical Inference, 2nd edition.
- Mosteller/Wallace*: Applied Bayesian and Classical Inference: The Case of *The Federalist Papers*.
- Pollard*: Convergence of Stochastic Processes.
- Pratt/Gibbons*: Concepts of Nonparametric Theory.
- Read/Cressie*: Goodness-of-Fit Statistics for Discrete Multivariate Data.

(continued after index)

Nicholas T. Longford

Models for Uncertainty in Educational Testing

With 31 Figures



Springer-Verlag

New York Berlin Heidelberg London Paris
Tokyo Hong Kong Barcelona Budapest

Nicholas T. Longford
Research Division
15-T Educational Testing Service
Rosedale Road
Princeton, NJ 08541 USA

Library of Congress Cataloging-in-Publication Data
Longford, Nicholas T., 1955-

Models for uncertainty in educational testing / Nicholas T. Longford.
p. cm. — (Springer series in statistics)

Includes bibliographical references and index.

ISBN-13: 978-1-4613-8465-6 e-ISBN-13: 978-1-4613-8463-2

DOI: 10.1007/978-1-4613-8463-2

1. Educational tests and measurements. 2. Examinations—Validity.
3. Examinations—Interpretation. 4. Examinations—Design and
construction. 5. Examinations—Scoring. 6. Educational Testing
Service. I. Title. II. Series.

LB3051.L625 1995

371.2'6'013—dc20

95-8145

Printed on acid-free paper.

© 1995 Springer-Verlag New York, Inc.
Softcover reprint of the hardcover 1st edition 1995

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.

Production managed by Hal Henglein; manufacturing supervised by Jeffrey Taub.
Photocomposed pages prepared from the author's LaTeX file.

9 8 7 6 5 4 3 2 1

ISBN-13: 978-1-4613-8465-6

*To Elyse and Kalman,
The children of the 'lesser' Dad*

Preface

Empirical Bayes and James-Stein estimation. This closely related pair of ideas is the biggest theoretical success on my list and the biggest underachiever in number of practical applications. The irony here is that these methods potentially offer the biggest practical gains over their predecessors.

Bradley Efron, *The statistical century*¹

The idea for this book came about with my growing realization that the Educational Testing Service (ETS) and its databanks are a genuine treasure trove of interesting statistical problems that are also encountered in many other fields. That in itself is a rather unoriginal statement, which can be found in one form or another in many other statistical texts motivated by educational testing issues. However, recent changes in the educational testing industry, especially the growing emphasis on constructed-response items, have brought about new statistical challenges and novel areas of application. This book is intended to be a tasting, rather than a complete and integral enumeration of these challenges, which are likely to be in a state of flux for a few years to come.

The title of the book, arrived at after umpteen revisions, is meant to have no other than statistical connotation. The term ‘uncertainty’ may

¹Bradley Efron’s agreement and permission of the Royal Statistical Society to reprint this passage from Efron (1995) are acknowledged.

sound somewhat negativist, but I hope that its appearance in the title will contribute to the emancipation of educational research from decades of futile attempts to find near-deterministic descriptions of educational phenomena with the aid of imperfect measurement instruments. I do not see the problem in the lack of perfection (when human subjects and observers are perfect and consistent, they are no longer human); it is rather in the reliance on statistical methods that assume perfect measurement and on cursory attention to problems of generalizability from one context (in which the study is conducted) to another (for which inference is desired). I find the term ‘uncertainty’ more appealing than ‘variation’ or *unexplained* variation to describe our inability to find a perfect explanation of an observed process.

The book is not intended as a celebration of ETS, but in part a piece of criticism, hopefully constructive, and principally a collection of case studies from educational testing. ETS has been generous in funding (in a piecemeal fashion but with good intentions) much of the research described in the book. I wish to express my gratitude to the research management for their patience and endurance with my criticism on a number of issues, and for giving me a go-ahead for this project. In brief, I acknowledge the flexible exercise of their authority. Where views are expressed, they are my own, unless stated otherwise, and ETS management has had little input in their formation. The book does not represent any aspects of official ETS policy, and I do not expect that it will do so at any time in the future.

One aspect in which my employer is beyond reproach is that I have always been granted leaves in response to invitations from academic institutions. I have done some of the work on this book while visiting the Zentrum für Umfragen, Methoden und Analysen in Mannheim, Germany (Michael Wiedenbeck), the Department of Statistics, University of Adelaide, Australia (Patty Solomon) and the Department of Economics, University of Pompeu Fabra in Barcelona, Spain (Albert Satorra). Their hospitality and stimulating environment are acknowledged. Research on the material presented in Chapter 7 was supported by a generous grant from the National Center for Educational Statistics.

The book would not have come about without the encouragement and exemplary cooperation of Martin Gilchrist of Springer-Verlag, who arranged for several detailed, informative, and insightful reviews of earlier drafts of the manuscript. I am grateful to several colleagues at ETS for their contribution to the book, in a variety of forms. Being mean by nature, I name only Eric Bradlow, Charlie Lewis, Gary Marco, Bob Mislevy, and Neal Thomas. I want to give credit to Bill Strawderman for several discussions and the resulting inspiration.

Those who provided the datasets analysed in the book are acknowledged in the text. Mike Wagner and Joe Bezek attended to the well-being of all the files on my Sun workstation. The production staff at Springer-Verlag helped to iron out the problems in wordprocessing and copy-editing.

My mother, Irena Lefkovičová, who introduced me many years ago to the charm and excitement of numbers, has been a staunch supporter throughout.

Princeton, NJ
April 1995

Nicholas T. Longford

Contents

Preface	vii
1 Inference about variation	1
1.1 Imperfection and variation	2
1.2 Educational measurement and testing	5
1.3 Statistical context	8
1.3.1 Statistical objects	8
1.3.2 Estimation	10
1.3.3 Correlation structure and similarity	13
1.3.4 Notation	14
2 Reliability of essay rating	17
2.1 Introduction	17
2.2 Models	21
2.3 Estimation	24
2.4 Extensions	29
2.5 Diagnostic procedures	30
2.6 Examples	32
2.6.1 Advanced Placement tests	37
2.7 Standard errors	42
2.7.1 Simulations	43
2.8 Summary	45
2.9 Literature review	45

3	Adjusting subjectively rated scores	47
3.1	Introduction	47
3.2	Estimating severity	48
3.3	Examinee-specific shrinkage	54
3.3.1	Rating in a single session	56
3.3.2	Shrinking to the rater's mean	57
3.4	General scheme	58
3.4.1	Sensitivity and robustness	59
3.5	More diagnostics	60
3.6	Examples	61
3.6.1	Advanced Placement tests	63
3.7	Estimating linear combinations of true scores	69
3.7.1	Optimal linear combinations	70
3.8	Summary	73
	Appendix. Derivation of MSE for the general adjustment scheme	74
4	Rating several essays	77
4.1	Introduction	77
4.2	Models	79
4.3	Estimation	82
4.4	Application	86
4.4.1	Itemwise analyses	90
4.4.2	Simultaneous analysis	90
4.5	Choice of essay topics	93
4.5.1	Modelling choice	95
4.5.2	Simulations	96
4.6	Summary	102
5	Summarizing item-level properties	105
5.1	Introduction	105
5.2	Differential item functioning	107
5.3	DIF variance	110
5.4	Estimation	112
5.5	Examples	116
5.5.1	National Teachers' Examination	116
5.5.2	GRE Verbal test	118
5.6	Shrinkage estimation of DIF coefficients	119
5.7	Model criticism and diagnostics	122
5.8	Multiple administrations	124
5.8.1	Estimation	126
5.8.2	Examples	129
5.8.3	Other applications	130
5.9	Conclusion	131

6	Equating and equivalence of tests	133
6.1	Introduction	133
6.2	Equivalent scores	135
6.2.1	Equating test forms	137
6.2.2	Half-forms	138
6.2.3	Linear true-score equating	139
6.3	Estimation	142
6.4	Application	145
6.4.1	Data and analysis	147
6.4.2	Comparing validity	152
6.4.3	Model criticism	154
6.5	Summary	155
7	Inference from surveys with complex sampling design	157
7.1	Introduction	157
7.2	Sampling design	159
7.2.1	The realized sampling design	160
7.2.2	The ‘model’ sampling design	161
7.2.3	Sampling weights and non-response	162
7.3	Proficiency scores	164
7.3.1	Imputed values	164
7.4	Jackknife	165
7.5	Model-based method	167
7.5.1	Stratification and clustering	167
7.5.2	Sampling variance of the ratio estimator	169
7.5.3	Within-cluster variance	170
7.5.4	Between-cluster variance	172
7.5.5	Multivariate outcomes	174
7.6	Examples	175
7.6.1	Subpopulation means	179
7.6.2	How much do weights matter?	181
7.7	Estimating proportions	182
7.7.1	Percentiles	185
7.8	Regression with survey data	186
7.9	Estimating many subpopulation means	189
7.10	Jackknife and model-based estimators	192
7.11	Summary	196
8	Small-area estimation	199
8.1	Introduction	199
8.2	Shrinkage estimation	200
8.3	Regression with survey data	203
8.4	Fitting two-level regression	205
8.4.1	Restricted maximum likelihood	207
8.4.2	Sampling weights	210

8.5	Small-area mean prediction	211
8.6	Selection of covariates	214
8.7	Application	215
8.7.1	No adjustment	219
8.7.2	Adjustment for covariates	220
8.7.3	Prediction and cross-validation	224
8.7.4	Refinement	226
8.8	Summary and literature review	229
9	Cut scores for pass/fail decisions	231
9.1	Introduction	231
9.2	Models	234
9.3	Fitting logistic regression	235
9.3.1	Generalized linear models	236
9.3.2	Random coefficients	238
9.3.3	Cut score estimation	239
9.4	Examples	240
9.4.1	PPST Writing test	240
9.4.2	Physical Education	243
9.5	Summary	248
10	Incomplete longitudinal data	251
10.1	Introduction	251
10.2	Informative missingness	252
10.3	Longitudinal analysis	254
10.4	EM algorithm	256
10.5	Application	258
10.6	Estimation	259
10.6.1	Variation in growth	263
10.6.2	Covariate adjustment	264
10.6.3	Missing covariate data	268
10.6.4	Standard errors	271
10.6.5	Clustering	271
10.7	Summary	272
	References	273
	Index	279