

Statistics for Social and Behavioral Sciences

Advisors:

S.E. Fienberg

W.J. van der Linden

For other titles published in this series, go to
<http://www.springer.com/series/3463>

Jean-Paul Fox

Bayesian Item Response Modeling

Theory and Applications

 Springer

Jean-Paul Fox
Department of Research Methodology,
Measurement, and Data Analysis
Faculty of Behavioral Sciences
University of Twente
7500 AE Enschede
The Netherlands

Series Editors

Stephen E. Fienberg
Department of Statistics
Carnegie Mellon University
Pittsburgh, PA 15213-3890
USA

Wim J. van der Linden
CTB/McGraw-Hill
20 Ryan Ranch Road
Monterey, CA 93940
USA

ISBN 978-1-4419-0741-7 e-ISBN 978-1-4419-0742-4
DOI 10.1007/978-1-4419-0742-4
Springer New York Dordrecht Heidelberg London

Library of Congress Control Number: 2010927930

© Springer Science+Business Media, LLC 2010

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, 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 in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

To Jasmijn and Kae

Preface

The modeling of item response data is governed by item response theory, also referred to as modern test theory. The field of inquiry of item response theory has become very large and shows the enormous progress that has been made. The mainstream literature is focused on frequentist statistical methods for estimating model parameters and evaluating model fit. However, the Bayesian methodology has shown great potential, particularly for making further improvements in the statistical modeling process.

The Bayesian approach has two important features that make it attractive for modeling item response data. First, it enables the possibility of incorporating nondata information beyond the observed responses into the analysis. The Bayesian methodology is also very clear about how additional information can be used. Second, the Bayesian approach comes with powerful simulation-based estimation methods. These methods make it possible to handle all kinds of priors and data-generating models.

One of my motives for writing this book is to give an introduction to the Bayesian methodology for modeling and analyzing item response data. A Bayesian counterpart is presented to the many popular item response theory books (e.g., Baker and Kim 2004; De Boeck and Wilson, 2004; Hambleton and Swaminathan, 1985; van der Linden and Hambleton, 1997) that are mainly or completely focused on frequentist methods. The usefulness of the Bayesian methodology is illustrated by discussing and applying a range of Bayesian item response models.

Complex Assessments

The recognition of the complexity of the processes leading to responses on test items stimulated the development of more realistic response models. The flexibility of the Bayesian modeling framework makes it particularly useful for making proper adjustments when the response data violate common model assumptions. Such violations might appear due to flaws in the data collection procedure or the complexity of the sample design.

Individuals' performances are best measured under controlled conditions such that other sources of variation are under control. An experimenter will choose a test for measuring a construct consisting of a set of items that minimizes individual variation and emphasize differences between subjects. Variability between respondents due to factors other than the construct under consideration is not desirable. Measurement models can become quite complex when they are adjusted in such a way that all sources of variation are taken into account. Flaws in experimental setups may engender the need for more complex measurement models. The impact of context effects on assessment results may further complicate the modeling process. Context effects appear when items function differently due to factors like item positioning or other material correlated with the item. Test data that violate assumptions of common item response models require a more flexible model that accounts for the violations.

Other complexities may arise due to the fact that besides response information other kinds of information are known (e.g., individual or group characteristics, response times, item characteristics). Different sampling designs can be used to sample respondents and/or items (e.g., adaptive item sampling, multistage sampling, randomized response sampling, simple random sampling, stratified sampling). Different response formats (e.g., multiple choice, binary, polytomous) and the presence of clusters of items may further stimulate the use of a more complex measurement model.

Besides introducing the Bayesian methodology, my aim has been to write a book to introduce Bayesian modeling of response data from complex assessments. Often information beyond the observed response data is available that can be used. The Bayesian modeling framework will prove to be very flexible, allowing simultaneous estimation of model parameters, computation of complex statistics, and simultaneous hypothesis testing.

In the 1990s, Bayesian inference became feasible with the introduction of Bayesian computational methods such as computer simulation and Monte Carlo techniques. The development of powerful computational simulation techniques induced a tremendous positive change in the applicability of Bayesian methodology. This led to the development of more flexible statistical models for test theory but also different strategies with respect to parameter estimation and hypothesis testing. In this book, the Bayesian way of item response modeling combined with the development of powerful numerical simulation techniques that led to a new research area in modern test theory is outlined.

Outside the Scope of This Book

Designing tests and testing whether tests are suited for the intended purpose are very complex subjects. Various questions need to be answered with respect to the response format of the test, the purpose of the test, and the construction of test materials, among others. The tests developed should also be reliable

and valid; that is, consistently result in scores that reflect the construct level of each respondent and measure what they are supposed to measure. Good tests are discriminating in the sense that they show differences in the construct level of respondents. There are a number of sources where this information is readily available. For classical test theory, see, for example, Gulliksen (1950), and Lord and Novick (1968), and for item response theory, see, for example, Lord and Novick (1968) and Lord (1980). A manual of standards for the construction and use of tests has been prepared by a joint committee of the American Educational Research Association, American Psychological Association and National Council of Measurement in Education (2000).

Overview

Statistical computations are necessary for applying the Bayesian methodology, and some programming skills are needed. That is, some familiarity with a statistical software package like R or S+ is needed to perform Bayesian analysis. On the one hand, this book aims to serve those who just want to apply the models, and they can use the software implemented in R packages and S+ programs (see Section 1.5). On the other hand, others may want to learn via programming and/or implement codes by themselves to extend models or adjust priors. For them, the mathematical details of the estimation procedures are discussed in the book, and the computer codes are provided via a website associated with the book. To understand the material, a basic background in probability and statistics is needed, including some familiarity with matrix algebra at the undergraduate level. The contents as well as the algorithms with their implementations make this book self-contained. Hopefully, it will provide an introduction to the essential features of Bayesian item response modeling as well as a better understanding of more advanced topics. The contents, programs, and codes will hopefully help readers implement their own algorithms and build their own set of tools for Bayesian item response modeling.

The book is organized as follows. In Chapter 1, the typical structure of item response data and the common item response models are discussed. Basic elements of Bayesian response modeling are introduced together with the basic building blocks for making Bayesian statistical inferences. WinBUGS is used to illustrate the Bayesian modeling approach. Chapter 2 presents a hierarchical modeling approach that supports the pooling of information, which becomes important when typically limited information is observed about many individuals. The Bayesian hierarchical modeling approach is outlined, which has tremendous potential with the current developments in statistical computing. Before discussing various sampling-based estimation methods for Bayesian item response models, which will be discussed in Chapter 4, in Chapter 3 a more general introduction is given to sampling-based estimation methods, testing hypotheses, and methods for model selection. Chapter 5 discusses methods for testing hypotheses and for model selection for the Bayesian item response models described in Chapter 4.

In Chapters 6–9, more advanced item response models are discussed for response data from complex assessments, response and response time data, and responses from complex sampling designs. In Chapter 6, respondents are assumed to be nested in groups (e.g., schools, countries). A hierarchical population model for respondents is defined to account for the within- and between-group dependencies. In Chapter 7, models for relaxing common measurement invariance restrictions are discussed. Chapter 8 introduces the multivariate analysis of responses and response times for measuring the speed and accuracy of working. Chapter 9 introduces models for (randomized) response data that are masked before they are observed to invite respondents to answer honestly when asked sensitive questions. Several empirical examples are presented to illustrate the methods and the usefulness of the Bayesian approach.

Acknowledgments

I would like to thank numerous people for their assistance and/or input during the writing of this book. I acknowledge the input from collaborators on earlier research projects that were addressed in the book. The cooperation of Rinke Klein Entink, Cees Glas, Wim van der Linden, Martijn de Jong, and Jan-Benedict Steenkamp has been greatly appreciated. I am indebted to Jim Albert, who provided very helpful comments on earlier drafts. Cheryl Wyrick has kindly provided data for a randomized response application. I thank my colleagues at the Department of Research Methodology, Measurement, and Data Analysis at the University of Twente. My colleagues Rinke Klein Entink and Josine Verhagen read drafts of the book, and their suggestions and comments led to its substantial improvement. I also thank John Kimmel for his confidence and assistance during the preparation of the book. The VIDI grant of the Netherlands Organization for Scientific Research (NWO) supported the writing of this book, for which I am most grateful.

Finally, I thank my wife, Miranda, for her support, encouragement, and patience during the writing of this book.

University of Twente, Enschede
March 2010

Jean-Paul Fox

Contents

Preface	VII
1 Introduction to Bayesian Response Modeling	1
1.1 Introduction	1
1.1.1 Item Response Data Structures	3
1.1.2 Latent Variables	5
1.2 Traditional Item Response Models	6
1.2.1 Binary Item Response Models.....	7
1.2.2 Polytomous Item Response Models	12
1.2.3 Multidimensional Item Response Models	14
1.3 The Bayesian Approach	15
1.3.1 Bayes' Theorem	16
1.3.2 Posterior Inference	20
1.4 A Motivating Example Using WinBUGS	21
1.4.1 Modeling Examinees' Test Results	21
1.5 Computation and Software	24
1.6 Exercises	27
2 Bayesian Hierarchical Response Modeling	31
2.1 Pooling Strength	31
2.2 From Beliefs to Prior Distributions	33
2.2.1 Improper Priors	38
2.2.2 A Hierarchical Bayes Response Model.....	39
2.3 Further Reading.....	42
2.4 Exercises	43
3 Basic Elements of Bayesian Statistics	45
3.1 Bayesian Computational Methods	45
3.1.1 Markov Chain Monte Carlo Methods	46
3.2 Bayesian Hypothesis Testing	51
3.2.1 Computing the Bayes Factor	54

3.2.2	HPD Region Testing	58
3.2.3	Bayesian Model Choice	59
3.3	Discussion and Further Reading	61
3.4	Exercises	62
4	Estimation of Bayesian Item Response Models	67
4.1	Marginal Estimation and Integrals	67
4.2	MCMC Estimation	71
4.3	Exploiting Data Augmentation Techniques	73
4.3.1	Latent Variables and Latent Responses	74
4.3.2	Binary Data Augmentation	75
4.3.3	TIMMS 2007: Dutch Sixth-Graders' Math Achievement	81
4.3.4	Ordinal Data Augmentation	83
4.4	Identification of Item Response Models	86
4.4.1	Data Augmentation and Identifying Assumptions	87
4.4.2	Rescaling and Priors with Identifying Restrictions	88
4.5	Performance MCMC Schemes	89
4.5.1	Item Parameter Recovery	89
4.5.2	Hierarchical Priors and Shrinkage	92
4.6	European Social Survey: Measuring Political Interest	95
4.7	Discussion and Further Reading	98
4.8	Exercises	99
5	Assessment of Bayesian Item Response Models	107
5.1	Bayesian Model Investigation	107
5.2	Bayesian Residual Analysis	108
5.2.1	Bayesian Latent Residuals	109
5.2.2	Computation of Bayesian Latent Residuals	109
5.2.3	Detection of Outliers	110
5.2.4	Residual Analysis: Dutch Primary School Math Test	111
5.3	HPD Region Testing and Bayesian Residuals	112
5.3.1	Measuring Alcohol Dependence: Graded Response Analysis	116
5.4	Predictive Assessment	117
5.4.1	Prior Predictive Assessment	119
5.4.2	Posterior Predictive Assessment	122
5.5	Illustrations of Predictive Assessment	126
5.5.1	The Observed Score Distribution	126
5.5.2	Detecting Testlet Effects	127
5.6	Model Comparison and Information Criteria	130
5.6.1	Dutch Math Data: Model Comparison	131
5.7	Summary and Conclusions	131
5.8	Exercises	133
5.9	Appendix: CAPS Questionnaire	139

6	Multilevel Item Response Theory Models	141
6.1	Introduction: School Effectiveness Research	141
6.2	Nonlinear Mixed Effects Models	142
6.3	The Multilevel IRT Model	145
6.3.1	A Structural Multilevel Model	145
6.3.2	The Synthesis of IRT and Structural Multilevel Models	148
6.4	Estimating Level-3 Residuals: School Effects	153
6.5	Simultaneous Parameter Estimation of MLIRT	158
6.6	Applications of MLIRT Modeling	162
6.6.1	Dutch Primary School Mathematics Test	162
6.6.2	PISA 2003: Dutch Math Data	165
6.6.3	School Effects in the West Bank: Covariate Error	172
6.6.4	MMSE: Individual Trajectories of Cognitive Impairment	174
6.7	Summary and Further Reading	181
6.8	Exercises	183
6.9	Appendix: The Expected School Effect	188
6.10	Appendix: Likelihood MLIRT Model	190
7	Random Item Effects Models	193
7.1	Random Item Parameters	193
7.1.1	Measurement Invariance	194
7.1.2	Random Item Effects Prior	195
7.2	A Random Item Effects Response Model	198
7.2.1	Handling the Clustering of Respondents	203
7.2.2	Explaining Cross-national Variation	203
7.2.3	The Likelihood for the Random Item Effects Model	204
7.3	Identification: Linkage Between Countries	205
7.3.1	Identification Without (Designated) Anchor Items	206
7.3.2	Concluding Remarks	208
7.4	MCMC: Handling Order Restrictions	209
7.4.1	Sampling Threshold Values via an M-H Algorithm	209
7.4.2	Sampling Threshold Values via Gibbs Sampling	211
7.4.3	Simultaneous Estimation via MCMC	212
7.5	Tests for Invariance	214
7.6	International Comparisons of Student Achievement	216
7.7	Discussion	221
7.8	Exercises	222
8	Response Time Item Response Models	227
8.1	Mixed Multivariate Response Data	227
8.2	Measurement Models for Ability and Speed	228
8.3	Joint Modeling of Responses and Response Times	231
8.3.1	A Structural Multivariate Multilevel Model	232

8.3.2	The RTIRT Likelihood Model	234
8.4	RTIRT Model Prior Specifications	235
8.4.1	Multivariate Prior Model for the Item Parameters	235
8.4.2	Prior for Σ_P with Identifying Restrictions	236
8.5	Exploring the Multivariate Normal Structure	238
8.6	Model Selection Using the DIC	241
8.7	Model Fit via Residual Analysis	242
8.8	Simultaneous Estimation of RTIRT	243
8.9	Natural World Assessment Test	246
8.10	Discussion	248
8.11	Exercises	250
8.12	Appendix: DIC RTIRT Model	254
9	Randomized Item Response Models	255
9.1	Surveys about Sensitive Topics	255
9.2	The Randomized Response Technique	256
9.2.1	Related and Unrelated Randomized Response Designs	257
9.3	Extending Randomized Response Models	258
9.4	A Mixed Effects Randomized Item Response Model	259
9.4.1	Individual Response Probabilities	259
9.4.2	A Structural Mixed Effects Model	261
9.5	Inferences from Randomized Item Response Data	262
9.5.1	MCMC Estimation	265
9.5.2	Detecting Noncompliance Behavior	267
9.5.3	Testing for Fixed-Group Differences	268
9.5.4	Model Choice and Fit	270
9.6	Simulation Study	272
9.6.1	Different Randomized Response Sampling Designs	272
9.6.2	Varying Randomized Response Design Properties	274
9.7	Cheating Behavior and Alcohol Dependence	275
9.7.1	Cheating Behavior at a Dutch University	275
9.7.2	College Alcohol Problem Scale	279
9.8	Discussion	284
9.9	Exercises	285
	References	289
	Index	309