

*Statistics for Social and
Behavioral Sciences*

Wim J. van der Linden

**Linear Models for
Optimal Test
Design**



Springer

*Statistics for Social and
Behavioral Sciences*

Wim J. van der Linden

**Linear Models for
Optimal Test
Design**

Statistics for Social and Behavioral Sciences

Advisors:

S.E. Fienberg W.J. van der Linden

Wim J. van der Linden

Linear Models for Optimal Test Design

Foreword by Ronald K. Hambleton

With 44 Figures

 Springer

Wim J. van der Linden
Department of Measurement
and Data Analysis
Faculty of Behavioral Sciences
University of Twente
7500 AE Enschede
The Netherlands
w.j.vanderlinden@utwente.nl

Advisors:

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

Wim J. van der Linden
Department of Measurement
and Data Analysis
Faculty of Behavioral Sciences
University of Twente
7500 AE Enschede
The Netherlands

Library of Congress Control Number: 2005923810

ISBN-10: 0-387-20272-2
ISBN-13: 978-0387-20272-3

Printed on acid-free paper.

© 2005 Springer Science+Business Media, 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 Science+Business Media, Inc., 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 in the United States of America. (EB)

9 8 7 6 5 4 3 2 1

springeronline.com

Voor mijn lieve Tonneke

Foreword

Over my nearly forty years of teaching and conducting research in the field of psychometric methods, I have seen a number of major technical advances that respond to pressing educational and psychological measurement problems. The development of criterion-referenced assessment was the first, beginning in the late 1960s with the important work of Robert Glaser and Jim Popham, in response to the need for assessments that considered candidate performance in relation to a well-defined body of knowledge and skills rather than in relation to a norm group. The development of criterion-referenced testing methodology with a focus on decision-theoretic concepts and methods, content validity, standard-setting, and the recognition of the merits of both criterion-norm-referenced and criterion-referenced assessments has tremendously influenced current test theory and testing .

The second major advance was the introduction of item response-theory (IRT) and associated models and their applications to replace classical test theory (CTT) and related practices. Beginning slowly in the 1940s and 1950s with the pioneering work of Frederic Lord, Allan Birnbaum, and Georg Rasch, by the 1970s the measurement journals were full of important research studies describing new IRT models, technical advances in model parameter estimation and model fit, and research on applications of IRT models to equating, test development, the detection of potentially biased test items, and adaptive testing. The overall goal has been to improve and expand measurement practices by overcoming several shortcomings of classical test theory: dependence of test-item statistics and reliability estimates on examinee samples, dependence of examinee true score estimates on the particular choices of test items, and the limitation in CTT of modeling ex-

aminee performance at the test level rather than at the item level. The last two shortcomings are especially problematic for adaptive testing, where it is important to be able to assess ability independently of particular test items and closely link item statistics to examinee ability or proficiency for the optimal selection of test items to shorten testing time and improve measurement precision on a per item basis. Today, the teaching of item-response theory is common in graduate training programs in psychometric methods, and IRT models and applications dominate the field of assessment.

The third major advance was the transition of testing practices from the administration of tests via paper and pencil to administration via the computer. This transition, which began in the late 1970s in the United States with considerable research funding from the armed services and with the leadership of such important scholars as Frederic Lord, Mark Reckase, Howard Wainer, and David Weiss, is widespread, with hundreds of credentialing exams (e.g., the Uniform Certified Public Accountancy Exams, the nursing exams, and securities industry exams in the United States), admissions tests (e.g., the Graduate Record Exam, the Graduate Management Admissions Test, and the Test of English as a Foreign Language), and achievement tests (e.g., high-school graduation tests in Virginia) being administered to candidates via computers, with more tests being added every month. The computer has added flexibility (with many testing programs, candidates can now take tests when they feel they are ready or when they need to take the tests), immediate scoring capabilities (thus removing what can often be months of waiting time for candidates), and the capability of assessing knowledge and skills that could not be easily assessed with paper-and-pencil tests. On this latter point, higher-level thinking skills, complex problem-solving, conducting research using reference materials, and much more are now being included in assessments because of the power of the computer.

Assessing candidates at a computer is becoming routine, and now a number of very important lines of research have been initiated. Research on automated scoring of constructed responses will ensure that computer-based testing can include the free-response test-item format, and thus the construct validity of many assessments will be enhanced. Research on automated item generation represents the next stage in test-item development and should expedite item writing, expand item pools, and lower the costs of item development. Automated item generation also responds to one of the main threats to the validity of computer-based testing with flexible candidate scheduling, and that is the overexposure of test items. With more test items available, the problem of overexposure of test items will be reduced.

Perhaps the most researched aspect of computer-based testing concerns the choice of test design. Initially, the focus was on fully adaptive tests. How should the first test item be selected? How should the second and third items and so on, be selected? When should testing be discontinued? How should ability or proficiency following the administration of each item be

estimated? Other test designs have been studied, too: multistage computer-based test designs (instead of selecting one optimal item after another, a block of test items, sometimes called “testlets” or “modules” are selected in some optimal fashion), and linear on-the-fly test designs (random or adaptive selection of tests subject to a variety of content and statistical constraints). Even the conventional linear test has been popular with one of a number of parallel forms being selected at random for administration to a candidate at a computer. But when computer-based testing research was initiated in the late 1970s, aptitude testing was the focus (e.g., the Armed Services Vocational Aptitude Battery), and detailed content-validity considerations were not a central concern. As the focus shifted to the study of computer-based achievement tests and credentialing exams (i.e., criterion-referenced tests) and the use of test scores became more important (e.g., credentialing exams are used to determine who is qualified to obtain a license or certificate to practice in a profession), content considerations became absolutely central to test defensibility and validity, and balancing tests from one examinee to the next for the length of item stems, the balance of constructed and selected response items, minimizing the overuse of test items, meeting detailed content specifications, building tests to match target information functions, and more, considerably more sophisticated methods for item selection were needed. It was in this computer-based testing environment that automated test assembly was born.

I have probably known about automated test assembly since 1983 (Wendy Yen wrote about it in one of her many papers), but the first paper I recall reading that was dedicated to the topic, and it is a classic in the psychometric methods field today, was the paper by Professor Wim van der Linden and Ellen Boekkooi-Timminga published in *Psychometrika* in 1989. In this paper, the authors introduced the concepts underlying automated test assembly and provided some very useful examples. I was fascinated that just about any content and statistical criteria that a test developer might want to impose on a test could be specified by them in the form of linear (in)equalities. Also, a test developer could choose an “objective function” to serve as the goal for test development. With a goal for test development reflected in an “objective function,” such as with respect to a target test-information function (and perhaps even several goals), and both content and statistical specifications described in the form of linear constraints, the computer could find a set of test items that maximally met the needs of the test developer. What a breakthrough! I might add that initially there was concern by some test developers that they might be losing control of their tests, but later it became clear that the computer could be used to produce, when desired, first drafts of tests that could then be reviewed and revised by committees.

The 1989 van der Linden and Boekkooi-Timminga paper was the first that I recall that brought together three immensely important technologies, two that I have already highlighted as major advances in the psychometric

methods field—item-response theory and the use of the computer—and also operations research. But what impresses me today is that automated test assembly impacts or capitalizes on all of the major advances in the last 40 years of my career: criterion-referenced and norm-referenced assessments, item-response theory, computer-based testing, and new computer-based test designs, as well as emerging new assessment formats.

By 2004, I had accumulated a hundred papers (and probably more) on the topic. Most are by Professor Wim van der Linden and his colleagues in the Netherlands, but many other researchers have joined in and are producing important work and advancing the field. These papers overflow my files on item-response theory, test design, computerized adaptive testing, item selection, item-bank inventory, item-exposure controls, and many more topics. My filing system today is simply not capable of organizing and sequencing all of the contributions on the topic of automated test assembly since 1989, and I have lost track of the many lines of research, the most important advances, and so on. Perhaps if I were closely working in the field, the lines of research would be clearer to me, but like many measurement specialists, I have a number of research interests, and it is not possible today to be fully conversant with all of them. But from a distance, it was clear to me that automated test assembly, or optimal test design, or automated test construction, all terms that I have seen used in the field, was going to provide the next generation of test-design methods—interestingly whether or not a test was actually going to be administered at a computer! Now, with one book, van der Linden's *Linear Models for Optimal Test Design*, order in my world has been restored with respect to this immensely important topic, and future generations of assessment specialists and researchers will benefit from Professor Wim van der Linden's technical advances and succinct writing skills.

I believe *Linear Models for Optimal Test Design* should be required reading for anyone seriously interested in the psychometric methods field. Computers have brought about major changes in the way we think about tests, construct tests, administer tests, and report scores. Professor van der Linden has written a book that organizes, clarifies, and expands what is known about test design for the next generation of tests, and test design is the base or centerpiece for all future testing. He has done a superb job of organizing and synthesizing the topic of automated test assembly for readers, providing a step-by-step introduction to the topic, and offering lots of examples to support the relevant theory and practices. The field is much richer for Professor van der Linden's contribution, and I expect this book will both improve the practice of test development in the future and spur others to carry out additional research.

Ronald K. Hambleton
University of Massachusetts at Amherst

Preface

The publication of Spearman's paper "The proof and measurement of association between two things" in the *American Journal of Psychology* in 1904 was the very tentative start of a new field now known as test theory. This book appears almost exactly a century later. During this period, test theory has developed from a timid fledgling to a mature discipline, with numerous results that nowadays support item and test analysis and test scoring at nearly every testing organization around the world.

This preface is not an appropriate place to evaluate a hundred years of test theory. But two observations may help me to explain my motives for writing this book. The first is that test theory has developed by careful modeling of response processes on test items and by using sophisticated statistical tools for estimating model parameters and evaluating model fit. In doing so, it has reached a current level of perfection that no one ever thought possible, say, two or three decades ago. Second, in spite of its enormous progress, although test theory is omnipresent, its results are used in a peculiar way. Any outsider entering the testing industry would expect to find a spin-off in the form of a well-developed technology that enables us to engineer tests rigorously to our specifications. Instead, test theory is mainly used for post hoc quality control, to weed out unsuccessful items, sometimes after they have been pretested, but sometimes after they have already been in operational use. Apparently, our primary mode of operation is not to create good tests, but only to prevent bad tests. To draw a parallel with the natural sciences, it seems as if testing has led to the development of a new science, but the spin-off in the form of a technology for engineering the test has not yet been realized.

Part of the explanation for our lack of technology may be a deeply ingrained belief among some in the industry that test items are unique and that test development should be treated as an art rather than a technology. I certainly believe that test items are unique. In fact, I even hope they will remain so; testing would suffer from serious security problems if they ceased to be so. Also, as a friend of the arts, I am sensitive to the aesthetic dimension of human artifacts. The point is, however, that these qualities do not relieve testing professionals of their duty to develop a technology. To draw another parallel, architecture has a deep artistic quality to it, and good architects are true artists. But if they were to give up their technology, we would have no place to live or work.

The use of design principles is an essential difference between technology-based approaches and the approaches with post hoc quality control hinted at above. Another difference is the use of techniques to guarantee that products will operate according to our specifications. These principles and techniques are to be used in a process that goes through four different stages: (1) establishing a set of specifications for the new testing program, (2) designing an item pool to support the program, (3) developing the item pool, and (4) assembling tests from the pool to meet the specifications. Although it is essential that the first stage be completed before the others are, the three other stages are more continuous and are typically planned to optimize the use of the resources in the testing organization. But it is important to distinguish between them because each involves the use of different principles and techniques.

At a slightly more formal level, test design is not unique at all; some of its stages have much in common with entirely different areas, where professionals also develop products, have certain goals in mind, struggle with constraints, and want optimal results. In fact, in this book I borrow heavily from the techniques of linear programming, widely used in industry, business, and commerce to optimize processes and products. These techniques have been around for a long time, and to implement them, we can resort to commercial computer software not yet discovered by the testing industry. In a sense, this book does not offer anything new. Then, to demonstrate the techniques's applicability, we had to reconceptualize the process of test design, introduce a new language to deal with it, integrate the treatment of content and statistical requirements for tests, and formulate typical test-design goals and requirements as simple linear models. More importantly, we also had to demonstrate the power and nearly universal applicability of these models through a wide range of empirical examples dealing with several test-design problems.

Although the topic of this book is *test design*, the term is somewhat ambiguous. The only stage in the design process at which something is actually designed is the second stage, item-pool design. From that point on, the production of a test only involves its assembly to certain specifications from a given item pool. The stages of item-pool design and test assembly

can be based on the same techniques from linear programming. But these techniques are much more easily understood as tools of test assembly, and for didactic reasons, I first treat the problem of test assembly and return to the problem of item-pool design as one of the last topics in this book.

In particular, the book is organized as follows. Chapter 1 introduces the current practice of test development and explains some elementary concepts from test theory, such as reliability and validity, and item and test information. Chapter 2 introduces a standard language for formulating test specifications. In Chapter 3, I show how this language can be used to model test assembly problems as simple linear models. Chapter 4 discusses general approaches available in mathematical programming, more specifically integer or combinatorial programming, to solve these models. A variety of empirical examples of the applications of the techniques to test-assembly problems, including such problems as IRT-based and classical test assembly, assembling multiple test forms, assembling tests with item sets, multidimensional test assembly, and adaptive test assembly, are presented in Chapters 5–9. The topic of item-pool design for programs with fixed and adaptive tests is treated in Chapter 10 and 11, respectively. The book concludes with a few more reflective observations on the topic of test design.

My goal has been to write a book that will become a helpful resource on the desk of any test specialist. Therefore, I have done my utmost to keep the level of technical sophistication in this book at a minimum. Instead, I emphasize such aspects as problem analysis, nature of assumptions, and applicability of results. In principle, the mathematical knowledge required to understand this book comprises linear equalities and inequalities from high-school algebra and a familiarity with set theory notation. The few formulas from test theory used in this book are discussed in Chapter 1. In addition, a few concepts from linear programming that are required to understand our modeling approaches are reviewed in Appendix 1. Nevertheless, Chapter 4 had to be somewhat more technical because it deals with methods for solving optimization problems. Readers with no previous experience with this material may find the brief introductions to the various algorithms and heuristics in this chapter abstract. If they have no affinity for the subject, they should read this chapter only cursorily, skipping the details they do not understand. They can do so without losing anything needed to understand the rest of the book. Also, it is my experience that the subject of multidimensional test assembly in Chapter 8 and, for that matter, the extension of adaptive test assembly to a multidimensional item pool in the last sections of Chapter 9, is more difficult to understand, mainly because the generalization of the notion of information in a unidimensional test to the case of multidimensionality is not entirely intuitive. Readers with no interest in this subject can skip this portion of the book and go directly to Chapter 10, where we begin our treatment of the subject of item-pool design.

Although this book presents principles and techniques that can be used in the three stages of test specification, item-pool design, and test assembly, the stage of item-pool development is hardly touched. The steps of item pretesting and calibration executed in this stage are treated well in several other books and papers (e.g., Hambleton & Swaminathan, 1985; Lord, 1980; Lord & Novick, 1968), and it is not necessary to repeat this material here. As for the preceding step of writing items for a pool, I do go as far as to show how blueprints for items can be calculated at the level of specific item writers and offer suggestions on how to manage the item-writing process (Chapter 10). But I do not deal with the actual process of item writing. Current item-writing practices are challenged by rapid developments in techniques for algorithmic item writing (e.g., Irvine & Kyllonen, 2002). I find these developments, which are in the same spirit as the “engineering approach” to test design advocated in this book, most promising, and I hope that, before too long, the two technologies will meet and integrate. This integration would reserve the intellectually more challenging parts of test design for our test specialists and allow them to assign their more boring daily operations to computer algorithms.

Several of the themes in this book were addressed in earlier research projects at the Department of Research Methodology, Measurement, and Data Analysis at the University of Twente. Over a period of more than 15 years, I have had the privilege of supervising dissertations on problems in test assembly and item-pool design by Jos J. Adema, Ellen Timminga, Bernard P. Veldkamp, and, currently, Adelaide Ariel. Their cooperation, creativity, and technical skills have been greatly appreciated. Special mention is deserved by Wim M.M. Tielen, who as a software specialist has provided continuous support in numerous test-assembly projects.

The majority of the research projects in this book were done with financial support from the Law School Admissions Council (LSAC), Newtown, Pennsylvania. Its continuous belief in what I have been doing has been an important stimulus to me, for which I am much indebted to Peter J. Pashley, Lynda M. Reese, Stephen T. Schreiber, and Philip D. Shelton. My main contact with the test specialists at the LSAC was Stephen E. Luebke, who provided all of the information about the item pools and test specifications that I needed for the projects in this book.

This book was written while I was a Fellow of the Center for Advanced Study in the Behavioral Sciences, Stanford, California. My fellowship was supported by a grant to the Center from the Spencer Foundation, for which I am most grateful. The tranquil location of the Center, on the top of a hill just above the Stanford campus, and the possession of a study overlooking a beautiful portion of the Santa Cruz Mountains, enabled me to view things in a wide perspective. I thank Doug McAdam, Director, and Mark Turner, Associate Director, as well as their entire staff, for their outstanding support during my fellowship. I am indebted to Kathleen Much for her

editorial comments on a portion of this book as well as on several other papers I wrote while at the Center.

Seven chapters of this book were tried out in a course on advanced topics in educational measurement at Michigan State University by Mark D. Reckase. His critical comments and those of his students led to many improvements in the original text. Bernard P. Veldkamp read several earlier versions of the manuscript and checked all exercises, while Adelaide Ariel went far beyond her call of duty with her help with the preparation of the graphs in this book. I am also grateful to Krista Breithaupt, Simon Bussman, Britta Colver, Alexander Freund, Heiko Grossman, Donovan Hare, Heinz Holling and Tobias Kuhn, whose comments helped me tremendously to polish the final version of the manuscript. The last chapter was completed while I enjoyed a fellowship from the Invitational Fellowship Program for Research in Japan at the University of Tokyo. I am indebted to the Japan Society for the Promotion of Science (JSPS) for the fellowship and to Kazuo Shigemasu for having been such a charming host.

Last but not least, I would like to thank John Kimmel, Executive Editor, Statistics, at Springer for being a quick and helpful source of information during the production of this book.

Each of the people whose support I acknowledge here have made my task as an author much more pleasant than I anticipated when I began working on the book.

Wim J. van der Linden
University of Twente

Acknowledgment of Copyrights

Several of the figures and tables in this book are (slightly re-edited) versions of figures and tables in earlier journal articles by the author. He is grateful to *Applied Psychological Measurement* for the right to reproduce Figures 5.6, 7.1, and 7.2 and Tables 5.3, 6.2, 6.3, 7.2, 10.1, and 11.1, to the *Journal of Educational and Behavioral Statistics* for the right to reproduce Figures 9.1, 11.1, and 11.2, and to the *Journal of Educational Measurement* for the right to reproduce Figures 11.3 and 11.4 and Table 5.1.