

Lecture Notes in
Statistics

157

Anne Boomsma

Marijtje A.J. van Duijn

Tom A.B. Snijders (Editors)

Essays on
Item Response Theory

Springer-Verlag

Anne Boomsma, Marijtje A.J. van Duijn, Tom A.B. Snijders
Department of Statistics and Measurement Theory
University of Groningen
Grote Kruisstraat 2/1
9712 TS Groningen
The Netherlands

Preface

Over the past three decades, since the publishing of Lord and Novick's *Statistical Theories of Mental Test Scores* in 1968 and Fischer's *Einführung in die Theorie psychologischer Tests* in 1974, item response theory (IRT) has developed rapidly. This is demonstrated in the *Handbook of Modern Item Response Theory* (Van der Linden & Hambleton, 1997) with chapters on a wide range of topics in IRT. The study of individual responses to behavioral stimuli has clearly evolved into a major discipline of psychometric theory.

The Dutch statistician and psychometrician Ivo Molenaar has played an important role in the growth of IRT, not only in the Netherlands but worldwide. His contributions to item response modeling cover a broad area: item and person fit in both parametric and nonparametric models, for example, and polytomous nonparametric item response models including the development of the MSP program. With Gerhard Fischer he edited a book on Rasch models that was published in 1995, and in cooperation with Klaas Sijtsma he is now preparing a textbook on nonparametric item response theory. In the Netherlands in particular, he has encouraged researchers and doctoral students to advance into new areas of IRT. To honor such achievements we dedicate this volume on item response modeling to Ivo W. Molenaar.

Ivo's general attitude towards psychometrics is perhaps best characterized by the definition he gave at the European Meeting of the Psychometric Society in Lüneburg in 1999, when introducing his successor as president of that society: "Psychometrics is mathematical statistics in the service of substantive psychology." This statement reflects Ivo Molenaar's scientific efforts: after the completion of a Ph.D. in mathematical statistics and obtaining a chair in statistical analysis and measurement theory at the University of Groningen, The Netherlands, in 1971, he showed an increasing professional interest in psychometrics. In that sense Ivo's career resembles that of Georg Rasch (1901–1980) described by Andersen and Olsen in Chapter 1 of this volume. In his inaugural speech as a full professor (Molenaar, 1972), he explicitly referred to the work of Georg Rasch and Robert Mokken, exponents of parametric and nonparametric IRT, respectively. From then on he gradually used his skills in mathematics and statistics in the service of item response theory. For Ivo it was a challenge to study and further develop both classes of models with their intriguing statistical and practical problems. And in the end he became an aficionado of modern item response theory, one of his greatest scientific endeavors.

This book deals with several aspects of item response modeling. We arranged the contributions in two parts: parametric and nonparametric IRT. Within each part, the chapters are roughly ordered as (1) historic accounts and overviews, (2) new models, (3) new methods, and (4) applications and miscellaneous topics. It should be noted that some chapters, such as those of Junker and of Mellenbergh, cover both parametric and nonparametric topics.

This book honors Ivo Molenaar on the occasion of his retirement from the Department of Statistics and Measurement Theory at the University of Groningen as of September 2000. In that small department, which he guided, extended, and defended for almost 30 years, many students and colleagues were fortunate to benefit from his outstanding didactic qualities in teaching and consultation. His work is perhaps best characterized by his practical approach to the use of statistics in social science applications. The invited contributors to this volume are researchers who have cooperated more or less closely with Ivo in the area of IRT, including former doctoral students. During the past 30 years Ivo Molenaar has been a memorable actor on the psychometric stage.

This is a book that should be of special interest to psychometricians, empirical researchers, and graduate students. It gives an impression of the present state of the art of IRT models and methods. It may be seen as a complementary successor to the *Handbook of Modern Item Response Theory* (Van der Linden & Hambleton, 1997), and we hope its contents will set a direction for future developments in IRT.

We thank the reviewers for their expertise and more than useful comments on the chapters of this book. Finally, we appreciate the patience of the authors while trying to cope with our frequently demanding requests for major and minor revisions.

Groningen, September 2000

Anne Boomsma
Marijtje A.J. van Duijn
Tom A.B. Snijders

References

- Fischer, G. (1974). *Einführung in die Theorie psychologischer Tests: Grundlagen und Anwendungen*. Bern: Huber.
- Fischer, G.H., & Molenaar, I.W. (Eds.). (1995). *Rasch models: Foundations, recent developments, and applications*. New York: Springer-Verlag.
- Lord, F.M., & Novick, M.R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- Molenaar, W. (1972). *Dit is een uitdaging* [This is a challenge]. Oratie [Inaugural lecture], Rijksuniversiteit Groningen.
- Van der Linden, W.J., & Hambleton, R.K. (Eds.). (1997). *Handbook of modern item response theory*. New York: Springer-Verlag.

Contents

Preface	vii
I Parametric Item Response Theory	
1 The Life of Georg Rasch as a Mathematician and as a Statistician	3
Erling B. Andersen and Lina Wøhlk Olsen	
1 Introduction	3
2 Early Life (1901–1945)	4
3 Rasch’s First Analysis of an Intelligence Test (1945–1948)	8
4 The Analysis of Slow Readers (1952)	10
5 Measuring Intelligence (1952–1953)	13
6 Discovery of the Dichotomous Rasch Model (1952–1958)	15
7 Work on the Models (1953–1958)	16
8 The Conversation with Ragnar Frisch in 1959	17
9 Two Important Publications (1960 and 1961)	18
10 Last Years (1962–1980)	21
11 Epilogue	23
References	23
2 The Growing Family of Rasch Models	25
Jürgen Rost	
1 Introduction	25
2 What Is a Rasch Model?	26
3 Some Historic Tracks of Generalizing the Rasch Model	30
4 A Hierarchical Structure of Generalized Rasch Models	36
References	37
3 Gain Scores Revisited Under an IRT Perspective	43
Gerhard H. Fischer	
1 Introduction	43

2	Measuring Change on the Basis of a PCM	46
3	Some Technical Considerations	50
4	Statistical Assessment of Change Under a One-Sided H_1	52
5	Statistical Assessment of Change Under a Two-Sided H_1	55
6	An Example	59
7	Conclusion	64
	References	66
4	Modeling Learning in Short-Term Learning Tests	69
	Karl Christoph Klauer and Hubert Sydow	
1	Introduction	69
2	The Learning Model	72
3	Estimating the Model	74
4	Validating the Model	74
5	Application	76
6	Discussion	84
	References	86
5	An IRT Model for Multiple Raters	89
	Norman D. Verhelst and Huub H.F.M. Verstralen	
1	Introduction	89
2	The Model	91
3	A Multilevel Interpretation of the Model	93
4	The IRT Approach	96
5	The Consequences of Ignoring Dependencies	99
6	Conclusion	104
	References	106
6	Conditional Independence and Differential Item Functioning in the Two-Parameter Logistic Model	109
	Herbert Hoijtink	
1	Introduction	109
2	A Statistic for Violations of Conditional Independence	112
3	A Statistic for Differential Item Functioning	114
4	The Posterior Predictive Distribution of Fit Statistics	116
5	Performance of Fit Statistics: A Small Simulation Study	117
6	Example: Masculinity and Femininity	119
7	Discussion	124
8	Appendix: Implementation of the Gibbs Sampler	125
	References	127

7	Differential Item Functioning Depending on General Covariates	131
	Cees A.W. Glas	
1	Introduction	131
2	Modeling Differential Item Functioning	133
3	Evaluation of Differential Item Functioning	135
4	The MML Framework	136
5	A Power Study	138
6	A Comparison of the MH and LM Approach	142
7	Discussion	145
	References	145
8	Statistical Tests for Differential Test Functioning in Rasch's Model for Speed Tests	149
	Margo G.H. Jansen and Cees A.W. Glas	
1	Introduction	149
2	An IRT Model for Response Times	151
3	Estimation	152
4	Model Tests	154
5	A Lagrange Multiplier Test for DTF	155
6	A Simulation Study	156
7	An Empirical Example	158
8	Discussion	160
	References	161
9	Expected Response Functions	163
	Charles Lewis	
1	Introduction	163
2	Theoretical Development of ERFs	164
3	ERFs for the Rasch Model	165
	References	170
10	A Logistic IRT Model for Decreasing and Increasing Item Characteristic Curves	173
	Edwin L. Klinkenberg	
1	Introduction	173
2	Attitudes on Issues	174
3	The Signed One-Parameter Logistic Model	176
4	The PARELLA Model	179
5	The Traffic Data	180

6	Summary and Discussion	186
7	Appendix: Estimating Equations of the Signed OPLM . . .	187
	References	190
11	Using Parameter Expansion to Improve the Performance of the EM Algorithm for Multidimensional IRT Population-Survey Models	193
	Donald B. Rubin and Neal Thomas	
1	Educational Assessment Surveys Using IRT Models	193
2	EM and PX-EM	194
3	Statistical Model	196
4	Example	200
5	Summary	202
	References	203
12	Cross-Validating Item Parameter Estimation in Adaptive Testing	205
	Wim J. van der Linden and Cees A.W. Glas	
1	Introduction	205
2	Capitalization on Item Calibration Error	208
3	Cross-Validating Item Parameter Estimation	211
4	Empirical Study	213
5	Concluding Remarks	217
	References	218
13	Imputation of Missing Scale Data with Item Response Models	221
	Mark Huisman and Ivo W. Molenaar	
1	Introduction	221
2	Incomplete Testing Designs	222
3	Imputation of Missing Item Responses	224
4	Effects of Imputation: A Simulation Study	229
5	Incomplete Designs versus Imputation	237
6	Summary and Conclusions	241
	References	243

II Nonparametric Item Response Theory

14 On the Interplay Between Nonparametric and Parametric IRT, with Some Thoughts About the Future 247

Brian Junker

1	Introduction	247
2	Nonparametric IRT: Scale Construction	248
3	Parametric IRT: Modeling Dependence	252
4	Measurement Challenges Posed by Cognitive and Embedded Assessments	257
	References	267

15 Reversibility Revisited and Other Comparisons of Three Types of Polytomous IRT Models 277

Bas T. Hemker

1	Introduction	277
2	Three Types of Models for Polytomous Items	278
3	Summary of Comparison Studies	282
4	The Three Types of Models and Reversibility	284
5	Comparison on Psychological Agreement	289
6	Discussion	293
	References	294

16 Progress in NIRT Analysis of Polytomous Item Scores: Dilemmas and Practical Solutions 297

Klaas Sijtsma and L. Andries van der Ark

1	Mokken Scale Analysis for Polytomous Item Scores	297
2	Three Open Theoretical Problems in NIRT	303
3	Discussion	314
	References	315

17 Two-Level Nonparametric Scaling for Dichotomous Data 319

Tom A.B. Snijders

1	Introduction	319
2	A Two-Level Model for Nonparametric Scaling of Dichotomous Data	321
3	Scalability Coefficients	323
4	Interpretation of Scalability Coefficients	326
5	Estimation of Scalability Coefficients	327

6	Object Scores	328
7	Reliability	329
8	Examples for Simulated Data	331
9	Example: Assessment of Teachers by Pupils	333
10	Discussion	335
11	Appendix: Proof that Between-Subject Scalability Coefficients Are Not Larger Than Within-Subject Coefficients	337
	References	337
18	The Circles of Our Minds: A Nonparametric IRT Model for the Circumplex	339
	Robert J. Mokken, Wijbrandt H. van Schuur, and Ard Jan Leeferink	
1	Introduction	339
2	The Circumplex Scale	341
3	A Search Procedure for a Circumplex Scale	344
4	Parameter Estimation: A Scale Scoring Method	350
5	Model Fit	352
6	Discussion	353
	References	354
19	Using Resampling Methods to Produce an Improved DIMTEST Procedure	357
	William Stout, Amy Goodwin Froelich, and Furong Gao	
1	Introduction	357
2	Review of the DIMTEST Procedure	361
3	Correcting the Bias in T_L	365
4	New Bias Correction Method	366
5	DIMTEST Without AT2	368
6	Monte Carlo Simulation Study	370
7	Discussion and Conclusions	373
	References	374
20	Person Fit Across Subgroups: An Achievement Testing Example	377
	Rob R. Meijer and Edith M.L.A. van Krimpen-Stoop	
1	Introduction	377
2	Person Fit in IRT Models	379
3	Empirical Research with Person-Fit Statistics	382
4	An Empirical Example	383
5	Discussion	387
	References	388

21 Single-Peaked or Monotone Tracelines?	
On the Choice of an IRT Model for Scaling Data	391
Wendy J. Post, Marijtje A.J. van Duijn, and Berna van Baarsen	
1 Introduction	391
2 Choosing an IRT Model	392
3 Monotone versus Single-Peaked Tracelines	394
4 Reanalysis of the Loneliness Scale Data	400
5 Summary and Discussion	409
References	411
22 Outline of a Faceted Theory of Item Response Data	415
Gideon J. Mellenbergh	
1 Introduction	415
2 Stimulus Facets	416
3 Person Facets	419
4 Recording Type Facet	420
5 Scaling Type Facet	420
6 Nested Facets	421
7 TIR and IRT	424
8 Example	427
9 Conclusion	429
References	430
Index	433
Abbreviations	439

1

The Life of Georg Rasch as a Mathematician and as a Statistician

Erling B. Andersen and Lina Wøhlk Olsen¹

ABSTRACT In this chapter an account of the life of Georg Rasch as a mathematician and as a statistician is given, with emphasis on the years 1952 to 1960 when the Rasch model was developed and justified by requirements of proper measurement. But Rasch's life before 1952 is also described, especially how he, forced by the employment situation in Denmark, had to turn to statistics after having graduated and defended his doctoral dissertation in pure mathematics. Finally, it is discussed at some length how Rasch, especially after 1960, formulated his requirements of proper measurement into the formalized concept of specific objectivity.

¹Københavns Universitet, Økonomisk Institut, Studiestræde 6, 1455 København K., Danmark; e-mail: *erling.b.andersen@econ.ku.dk*, *lina@econ.ku.dk*

2

The Growing Family of Rasch Models

Jürgen Rost¹

ABSTRACT A family of Rasch models is defined in terms of the prominent properties of all Rasch models, that is, separability, sufficiency, specific objectivity, and latent additivity. It is argued that concepts such as item homogeneity, person homogeneity, and unidimensionality do not hold for all generalizations of the Rasch model (RM). Four directions of generalizing the model are discussed: the multidimensional, the ordinal polytomous, the linear logistic, and the mixture distribution generalization. A hierarchical system of generalized Rasch models is presented. Four out of these eight models are discussed in the literature and can be applied by means of reliable computer software.

¹Institute for Science Education, University of Kiel, Olshausenstr. 62, D-24098 Kiel, Federal Republic of Germany; e-mail: rost@ipn.uni-kiel.de

3

Gain Scores Revisited Under an IRT Perspective

Gerhard H. Fischer¹

ABSTRACT For the measurement and statistical assessment of individual gain scores based on item sets that satisfy the assumptions of the Rasch, Rating Scale, or Partial Credit Models, a conditional maximum likelihood estimator, Clopper–Pearson confidence intervals, uniformly most accurate confidence intervals, and uniformly most powerful unbiased tests of the hypothesis of no change are presented. All methods are grounded on the exact conditional distribution of the gain score, given the total score for both time points, so that no asymptotic approximations are required. Typical applications of the methods are mentioned.

¹Department of Psychology, University of Vienna, Liebiggasse 5, A-1010 Wien, Austria; e-mail: *gh.fischer@univie.ac.at*

4

Modeling Learning in Short-Term Learning Tests

Karl Christoph Klauer¹ and Hubert Sydow²

ABSTRACT A new model-based method for analyzing learning processes in so-called short-term learning tests is proposed and applied. Learning is measured by means of a two-dimensional item response theory model that incorporates two latent factors: ability level and learning ability. The model allows one to assess the impact of both variables and their correlation for given data sets in a manner that implicitly corrects for statistical artifacts that arise in conventional analyses. The model is applied to four short-term learning tests administered to $N = 434$ children aged five to six years. Model tests and tests for empirical validity of the model parameters succeed in establishing construct validity for one of these tests.

¹Psychologisches Institut, Rheinische Friedrich-Wilhelms-Universität Bonn, Römerstr. 164, 53117 Bonn, FR Germany; e-mail: *christoph.klauer@uni-bonn.de*

²Institut für Psychologie, Humboldt-Universität Berlin, Oranienburger Str. 18, 10178 Berlin, FR Germany; e-mail: *hubert.sydow@psychologie.hu-berlin.de*

5

An IRT Model for Multiple Raters

Norman D. Verhelst¹ and Huub H.F.M. Verstralen²

ABSTRACT An IRT model for multiple ratings is presented. If it is assumed that the quality of a student performance has a stochastic relationship with the latent variable of interest, it is shown that the ratings of several raters are not conditionally independent given the latent variable. The model gives a full account of this dependence. Several relationships with other models appear to exist. The proposed model is a special case of a nonlinear multilevel model with three levels, but it can also be seen as a linear logistic model with relaxed assumptions (LLRA). Moreover, a linearized version of the model turns out to be a special case of a generalizability model with two crossed measurement facets (items and raters) with a single first-order interaction term (persons and items). Using this linearized model, it is shown how the estimated standard errors of the parameters are affected if the dependence between the ratings is ignored.

¹National Institute for Educational Measurement (CITO), P.O. Box 1034, 6801 MG Arnhem, The Netherlands; Faculty of Educational Science and Technology, University of Twente, The Netherlands; e-mail: *norman.verhelst@cito.nl*

²National Institute for Educational Measurement (CITO), P.O. Box 1034, 6801 MG Arnhem, The Netherlands; e-mail: *huub.verstralen@cito.nl*

6

Conditional Independence and Differential Item Functioning in the Two-Parameter Logistic Model

Herbert Hoijtink¹

ABSTRACT The two-parameter logistic model is among other things characterized by conditional independence and the absence of differential item functioning. Two fit statistics are proposed that can be used to investigate conditional independence and differential item functioning globally and at the item level. The statistics are evaluated using their posterior predictive distribution.

¹Department of Methodology and Statistics, University of Utrecht, P.O. Box 80140, 3508 TC Utrecht, The Netherlands; e-mail: *h.hoijtink@fss.uu.nl*

7

Differential Item Functioning Depending on General Covariates

Cees A.W. Glas¹

ABSTRACT Item response theory (IRT) is a powerful tool for the detection of differential item functioning (DIF). It is shown that the class of IRT models with manifest predictors is a comprehensive framework for the detection of DIF. These models also support the investigation of the causes of DIF. In principle, the responses to every item in a test can be subject to DIF, and traditional IRT-based detection methods require one or more estimation runs for every single item. Therefore, Glas (1998) proposed an alternative procedure that can be performed using only a single estimate of the item parameters. This procedure is based on the Lagrange multiplier test or the equivalent Rao efficient score test. In this chapter, the procedure is generalized in various directions, the most important one being the possibility of conditioning on general covariates. A small simulation study is presented to give an impression of the power of the test. In an example using real data it is shown how the method can be applied to the identification of main and interaction effects in DIF.

¹Department of Educational Measurement and Data Analysis, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands; e-mail: glas@edte.utwente.nl

8

Statistical Tests for Differential Test Functioning in Rasch's Model for Speed Tests

Margo G.H. Jansen¹ and Cees A.W. Glas²

ABSTRACT We consider a latent trait model developed by Rasch for the response time on a set of pure speed tests, which is based on the assumption that the test response times are approximately gamma distributed with known index parameters and scale parameters depending on subject ability and test difficulty parameters. In this chapter, the principle of Lagrange multiplier tests is used to evaluate differential test functioning and subgroup invariance of the test parameters. Two numerical illustrations are given.

¹Department of Education, University of Groningen, Grote Rozenstraat 38, 9712 TJ Groningen, The Netherlands; e-mail: *g.g.h.jansen@ppsw.rug.nl*

²Department of Educational Measurement and Data Analysis, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands; e-mail: *glas@edte.utwente.nl*

9

Expected Response Functions

Charles Lewis¹

ABSTRACT Item response theory makes use of what are sometimes referred to as item response functions. Such a function is actually a probability density function (pdf) for the response of an individual to a test item, given the values of certain parameters, classified as item parameters and person parameters (or abilities). In testing, there is typically a calibration phase, in which item parameters are estimated and abilities are ignored. This is followed by an application phase, in which abilities are estimated while conditioning on the estimated values of the item parameters. A Bayesian alternative to treating the item parameters as known quantities involves replacing item response functions with another class of pdfs, referred to as expected response functions (ERFs). The latter take the uncertainty regarding item parameters into account for purposes of estimating abilities. This chapter provides a formal description of ERFs and briefly illustrates their application to the Rasch model for binary item responses.

¹Educational Testing Service, Rosedale Road, Mailstop 03-T, Princeton, NJ 08541, USA; e-mail: *clewis@ets.org*

10

A Logistic IRT Model for Decreasing and Increasing Item Characteristic Curves

Edwin L. Klinkenberg¹

ABSTRACT In item response theory, dominance relations are modeled by cumulative models, and proximity relations by unfolding models. Usually cumulative models are used for the measurement of latent abilities and unfolding models for the measurement of latent attitudes and preferences. The distinction between both types of measurement models is best represented by the shape of the item characteristic curve which is monotone increasing in the cumulative model and single-peaked in the unfolding model. A boundary case is a situation in which some items are monotone decreasing and others are monotone increasing. In this chapter, an extension of the one-parameter logistic model is proposed for that situation. It is shown that the parameters of the proposed model can be estimated by a simple data transformation. The model is illustrated using an empirical data set.

¹Department of Psychiatry and Neuropsychology, Maastricht University, P.O. Box 616, 6200 MD Maastricht, The Netherlands; e-mail: *e.klinkenberg@np.unimaas.nl*

11

Using Parameter Expansion to Improve the Performance of the EM Algorithm for Multidimensional IRT Population-Survey Models

Donald B. Rubin¹ and Neal Thomas²

ABSTRACT Several large-scale educational surveys use item response theory (IRT) models to summarize complex cognitive responses and relate them to educational and demographic variables. The IRT models are often multidimensional with prespecified traits, and maximum likelihood estimates (MLEs) are found using an EM algorithm, which is typically very slow to converge. We show that the slow convergence is due primarily to missing information about the correlations between latent traits relative to the information that would be present if these traits were observed (“complete data”). We show how to accelerate convergence using parameter extended EM (PX-EM), a recent modification of the EM algorithm. The PX-EM modification is simple to implement for the IRT survey model and reduces the number of required iterations by approximately one-half without adding any appreciable computation to each EM-iteration.

¹Datametrics Research, Inc. and Harvard University, 1 Oxford Street, Cambridge, MA 02138, USA; e-mail: *rubin@stat.harvard.edu*

²Datametrics Research, Inc. and Bristol-Myers Squibb, 61 Dreamlake Drive, Madison, CT 06443, USA; e-mail: *snthomas99@yahoo.com*

12

Cross-Validating Item Parameter Estimation in Adaptive Testing

Wim J. van der Linden and Cees A.W. Glas¹

ABSTRACT From a statistical point of view, adaptive testing belongs to the class of problems of selecting a best subset of random variables. As in any other instance of this problem, adaptive testing is likely to capitalize on estimation error. An empirical example of the phenomenon is discussed, and a strategy for cross-validating results from adaptive testing is proposed. Results from a simulation study on the strategy are presented.

¹Both authors are at the Department of Educational Measurement and Data Analysis, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands; e-mail: *vanderlinden@edte.utwente.nl* and *glas@edte.utwente.nl*

13

Imputation of Missing Scale Data with Item Response Models

Mark Huisman and Ivo W. Molenaar¹

ABSTRACT Confronted with incomplete data due to nonresponse, a researcher may want to impute missing values to estimate latent properties of respondents. In this chapter the results of a simulation study are presented, investigating the performance of several imputation techniques. Some imputation procedures are based on item response theory (IRT) models, which can also be used to estimate latent abilities directly from the incomplete data by using incomplete testing designs when data are missing by design. This strategy has some serious disadvantages in the case of item nonresponse, because nonresponse is assumed to be ignorable and computational problems arise in scales with many items. In a second simulation study, the performance of some imputation techniques is compared to the incomplete design strategy, in the case of item nonresponse. The latter strategy results in slightly better ability estimates, but imputation is almost as good, especially when it is based on IRT models.

¹Both authors are at the Department of Statistics and Measurement Theory, University of Groningen, Grote Kruisstraat 2/1, 9712 TS Groningen, The Netherlands; e-mail: *m.huisman@ppsw.rug.nl* and *w.molenaar@ppsw.rug.nl*

14

On the Interplay Between Nonparametric and Parametric IRT, with Some Thoughts About the Future

Brian Junker ¹

ABSTRACT This chapter reviews some of the important research in nonparametric and parametric item response theory (IRT) today, and considers some current measurement challenges in education and cognitive psychology. This leads to assessment models that do not look very much like today's IRT models, but for which the tools and conceptual framework of nonparametric and parametric IRT are still quite well suited.

¹Department of Statistics, Carnegie Mellon University, 232 Baker Hall, Pittsburgh, PA 15213, USA; e-mail *brian@stat.cmu.edu*

15

Reversibility Revisited and Other Comparisons of Three Types of Polytomous IRT Models

Bas T. Hemker¹

ABSTRACT Molenaar (1983) was the first to distinguish and to compare three different types of polytomous item response theory (IRT) models: cumulative, continuation, and partial credit models. In his research, and later in research by others, the three types of models were compared on various criteria. One of the criteria introduced in IRT modeling is the property of reversibility. In this chapter, the results of a number of comparison studies are summarized. Next, the property of reversibility is further investigated. Finally, the relationship between the three types of models and psychological reality is discussed.

¹National Institute for Educational Measurement (CITO), P.O. Box 1034, 6801 MG Arnhem, The Netherlands; e-mail: *bas.hemker@cito.nl*

16

Progress in NIRT Analysis of Polytomous Item Scores: Dilemmas and Practical Solutions

Klaas Sijtsma and L. Andries van der Ark¹

ABSTRACT This chapter discusses several open problems in nonparametric polytomous item response theory: (1) theoretically, the latent trait θ is not stochastically ordered by the observed total score X_+ ; (2) the models do not imply an invariant item ordering; and (3) the regression of an item score on the total score X_+ or on the restscore R is not a monotone non-decreasing function and, as a result, it cannot be used for investigating the monotonicity of the item step response function. Tentative solutions for these problems are discussed. The computer program MSP for nonparametric IRT analysis is based on models that neither imply the stochastic ordering property nor an invariant item ordering. Also, MSP uses item restscore regression for investigating item step response functions. It is discussed whether computer programs may be based (temporarily) on models that lack desirable properties and use methods that are not (yet) supported by sound psychometric theory.

¹Both authors are at the Department MTO, Tilburg University, P.O. Box 90153, 5000 LE Tilburg, The Netherlands; e-mail: *k.sijtsma@kub.nl* and *a.vdark@kub.nl*

17

Two-Level Nonparametric Scaling for Dichotomous Data

Tom A.B. Snijders¹

ABSTRACT This chapter considers a two-level design where the objects to be scaled are the higher-level units. Nested within each object are lower-level units, called subjects, and a set of dichotomous items is administered to each subject. The subjects are regarded as strictly parallel tests for the objects. Examples are the scaling of teachers on the basis of their pupils' responses, or of neighborhoods on the basis of responses by inhabitants. A two-level version of the nonparametric scaling method first proposed by Mokken (1971) is elaborated upon. The probabilities of positive responses to the items are assumed to be increasing functions of the value on a latent trait. The latent trait value for each subject is composed of an object-dependent value and a subject-dependent deviation from this value. The consistency of responses within, but also between objects is expressed by two-level versions of Loevinger's H -coefficients. The availability of parallel tests is used to calculate a reliability coefficient.

¹Department of Statistics and Measurement Theory, University of Groningen, Grote Kruisstraat 2/1, 9712 TS Groningen, The Netherlands; e-mail: *t.a.b.snijders@ppsw.rug.nl*

18

The Circles of Our Minds: A Nonparametric IRT Model for the Circumplex

Robert J. Mokken,¹ Wijbrandt H. van Schuur,² and
Ard Jan Leeferink³

ABSTRACT A nonparametric item response theory (IRT) model for the circumplex is introduced, based on Mokken's (1971) nonparametric IRT model for cumulative scaling, and Van Schuur's (1984) nonparametric IRT model for unfolding. Some examples of circumplex representations in the social sciences are given. Model fit is based first on an extension of Loevinger's coefficient of homogeneity using quadruples as elementary units of analysis. Diagnostics for the probabilistic circumplex are suggested. Assignment of (ordinal) scale values is based on the notion of item steps, as developed by Molenaar (1983). The model presented here is based on dichotomous *pick any/k* data. Suggestions for extension to rank *m/k* data and to polytomous data are discussed.

¹Faculty of Social and Behavioral Sciences, University of Amsterdam. Private address: Franklinstraat 18, 1171 BL Badhoevedorp, The Netherlands; e-mail: mokken@ccsom.uva.nl

²Department of Sociology, University of Groningen, Grote Rozenstraat 31, 9712 TG Groningen, The Netherlands; e-mail: h.van.schuur@ppsw.rug.nl

³Department of Policy Research of the City of Groningen, Postbus 1125, 9701 BC Groningen, The Netherlands; e-mail: j.g.y.leeferink@sozawe.groningen.nl

19

Using Resampling Methods to Produce an Improved DIMTEST Procedure

William Stout,¹ Amy Goodwin Froelich,² and Furong Gao³

ABSTRACT Following in the nonparametric item response theory tradition, DIMTEST (Stout, 1987) is an asymptotically justified nonparametric procedure that provides a test of hypothesis of unidimensionality of a test data set. This chapter introduces a new bias correction method for the DIMTEST procedure based on the statistical principle of resampling. A simulation study shows this new version of DIMTEST has a Type I error rate close to the nominal rate of $\alpha = 0.05$ in most cases and very high power to detect multidimensionality in a variety of realistic multidimensional models. The result with this new bias correction method is an improved DIMTEST procedure with much wider applicability and good statistical performance.

¹Department of Statistics, University of Illinois at Urbana-Champaign, 725 S. Wright Street, Champaign, IL 61820, USA; e-mail: stout@stat.uiuc.edu

²Department of Statistics, University of Illinois at Urbana-Champaign, 725 S. Wright Street, Champaign, IL 61820, USA; e-mail: amyf@stat.uiuc.edu

³CTB/McGraw-Hill Research Department, 20 Ryan Ranch Road, Monterey, CA 93940, USA; e-mail: FGao@ctb.com

Person Fit Across Subgroups: An Achievement Testing Example

Rob R. Meijer and
Edith M.L.A. van Krimpen-Stoop¹

ABSTRACT Item response theory (IRT) models are used to describe answering behavior on tests and examinations. Although items may fit an IRT model, some persons may produce misfitting item score patterns, for example, as a result of cheating or lack of motivation. Several statistics have been proposed to detect deviant item score patterns. Misfitting item score patterns may be related to group characteristics such as gender or race. Investigating misfitting item score patterns across different groups is strongly related to differential item functioning (DIF). In this study the usefulness of person fit to compare item score patterns for different groups was investigated. In particular, the effect of misspecification of a model due to DIF on person fit was explored. Empirical data of a math test were analyzed with respect to misfitting item score patterns and DIF for men and women and blacks and whites. Results indicated that there were small differences between subgroups with respect to the number of misfitting item score patterns. Also, the influence of DIF on the fit of a score pattern was small for both gender and ethnic groups. The results imply that person-fit analysis is not very sensitive to model misspecification on the item level.

¹Both authors are at the Department of Educational Measurement and Data Analysis, University of Twente, P.O. Box 217, 7500 AE Enschede, The Netherlands; e-mail: *meijer@edte.utwente.nl* and *krimpen@edte.utwente.nl*

21

Single-Peaked or Monotone Tracelines? On the Choice of an IRT Model for Scaling Data

Wendy J. Post,¹ Marijtje A.J. van Duijn,² and Berna van Baarsen³

ABSTRACT It is investigated how a motivated choice can be made for the analysis of an item set using either a cumulative item response model with monotone tracelines, modeling dominance relations, or a unimodal item response model with single-peaked tracelines, modeling proximity relations. The focus is on item sets consisting of positively and negatively formulated items (with respect to the latent trait to be measured), where the common practice is to reverse one type of item. The differences between the cumulative and unimodal model are studied theoretically, in terms of item location and item order, and empirically, in a reanalysis of a sample of De Jong Gierveld loneliness scale data. Item locations, and, in the case of the unimodal model, also subject locations are shown to be important determinants of the differences. For the loneliness scale data the analysis with the unimodal model is preferred over the cumulative model. An outline of a recommended strategy for an IRT analysis of scaling data is given.

¹Department of Medical Statistics, Leiden University Medical Center, P.O. Box 9604, 2300 RC Leiden, The Netherlands; e-mail: w.j.post@mta.azg.nl

²Department of Statistics and Measurement Theory, University of Groningen, Grote Kruisstraat 2/1, 9712 TS Groningen, The Netherlands; e-mail: m.a.j.van.duijn@ppsw.rug.nl

³Department of Philosophy and Medical Ethics, Free University, Van der Boechorststraat 7, 1081 BT Amsterdam, The Netherlands; email: b.vbaarsen.metamedica@med.vu.nl

22

Outline of a Faceted Theory of Item Response Data

Gideon J. Mellenbergh¹

ABSTRACT Coombs' (1964) *Theory of Data* is a classification system of behavioral data. In his system, item response data are individual-comparison or individual-stimulus differences data. In this contribution, item response data are further classified using the following facets: Intended Behavior Type, Task Type, Intermediating Variable Type, Construct Level Type, Construct Dimension Type, Recording Type, Scaling Type, Person \times Stimulus Interaction Type, Construct Relation Type, and Step Task Type. These facets cannot be completely crossed because some of the facets are nested within certain (combinations of) cells of crossed facets. The facets constitute a theory of item responses (TIR). The TIR and item response theory (IRT) complete each other: the TIR structures the item response data, while IRT explains these data. The TIR can be used to facilitate the choice of an appropriate IRT model.

¹Department of Psychology, University of Amsterdam, Roetersstraat 15, 1018 WB Amsterdam, The Netherlands; e-mail: ml.Mellenbergh@macmail.psy.uva.nl