

# Some topics in nonparametric and parametric IRT, with some thoughts about the future

Brian Junker<sup>1</sup>

Department of Statistics  
Carnegie Mellon University  
232 Baker Hall  
Pittsburgh PA 15213  
[brian@stat.cmu.edu](mailto:brian@stat.cmu.edu)

May 1, 2000

<sup>1</sup>On leave at Learning Research and Development Center, 3939 O'Hara Street, University of Pittsburgh, Pittsburgh PA 15260. Some of the work reported here was performed in the course of preparing a commissioned paper for the National Research Council Committee on Educational and Psychological Foundations of Assessment, United States National Academy of Sciences. Preparation of this paper was supported in part by NSF grants #SES-99.07447 and #DMS-97.05032, and the hospitality of the Learning Research and Development Center. Thanks to Anne Boomsma, Klaas Sijsma, Lou Mariano, Matt Johnson, and L. Andries van der Ark for their comments and suggestions.

## Abstract

Item response theory (IRT) has arguably been one of the most successful and widely used statistical modeling techniques in psychometrics, with applications in developmental, social, educational and cognitive psychology for example, as well as in medical research, demography and other social science settings. In this paper I will try to briefly summarize some of the important research in nonparametric and parametric IRT today. I will try to show that a broad understanding of IRT as an instance of mathematical statistics in the service of substantive psychology, together with an appreciation of some of the current measurement challenges in education and cognitive psychology, lead us 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.

## 1 Introduction

In introducing Susan Embretson's 1999 Presidential Address at the European Meeting of the Psychometric Society in Lüneburg Germany, Ivo Molenaar defined psychometrics as “mathematical statistics in the service of substantive psychology”: that definition cuts a pretty wide swath, and indicates just how general the psychometric enterprise can and should be. Item response theory (IRT; e.g., Fischer Molenaar, 1995; Van der Linden Hambleton, 1997) is a psychometric approach to modeling data from social surveys and educational and psychological tests, dating back at least to Lord (1952) and Rasch (1960), and to the work of Loevinger and Guttman before them. IRT enables us to study the characteristics of test or survey items across multiple respondent populations, and to study respondents' propensities to answer positively across various items. IRT has arguably been one of the most successful and widely used techniques in psychometrics, with applications in developmental, social, educational and cognitive psychology for example, as well as in medical research, demography and other social science settings.

In this paper I will try to briefly summarize some of the important research in nonparametric and parametric IRT today, emphasizing the interplay between parametric and nonparametric models that is the hallmark of the approach initiated in the Netherlands by Mokken and pursued by Molenaar, Sijtsma, and their colleagues, and re-ignited in the U.S. by Holland, Rosenbaum, and their colleagues. I will try to show that a broad understanding of IRT as an instance of “mathematical statistics in the service of substantive psychology”, together with an appreciation of some of the current measurement challenges in education and cognitive psychology, lead us 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 particularly well suited.

## 2 Nonparametric IRT: Scale Construction

To save space and preserve focus, my summary of nonparametric IRT will concentrate on the scaling theory techniques introduced by Mokken and pursued by Molenaar, Sijtsma, and their colleagues. Important related work on nonparametric essential unidimensionality (e.g., Stout 1987, 1990; Zhang Stout, 1996; and Stout, Habing, Douglas, Kim, Roussos Zhang, 1996), nonparametric regression estimates of item response

functions and test response surfaces (especially Ramsay, 1991, 1995, 1996), and related parametric and non-parametric work (e.g., Meijer, 1996; Cliff Donoghue, 1992; Drasgow, Levine, Tsien, Williams Mead, 1995; Samejima, 1997) will not be considered in detail.

## 2.1 Monotone Homogeneity / Strict Unidimensionality

The Mokken (1971) model of *monotone homogeneity* starts with very few assumptions. A collection of items  $\mathbf{X} = (X_1, \dots, X_J)$ , which may include dichotomous, ordered polytomous, or even continuous response items, satisfies this model, if  $\theta$ , the latent trait, is a real-valued random variable (unidimensionality); if each item step response function (ISRF)  $P[X_j > c|\theta]$  is non-decreasing in  $\theta$  for each item response variable  $X_j$  and each real threshold  $c$  (monotonicity); and if

$$P[X_1 > c_1, \dots, X_J > c_J | \theta] = \prod_{j=1}^J P[X_j > c_j | \theta] \quad (1)$$

for all possible cutoffs  $c_j$  (local independence). When the response is dichotomous (0/1) I refer to  $P_j(\theta) \equiv P[X_j = 1 | \theta]$  as the item response function (IRF). In that case, equation (1) reduces to the familiar form

$$P[X_1 = x_1, \dots, X_J = x_J] = \prod_{j=1}^J P_j(\theta)^{x_j} [1 - P_j(\theta)]^{1-x_j}. \quad (2)$$

The data we observe when examinees or subjects respond to test or survey items can be thought of as i.i.d. samples from the marginal discrete multivariate distribution,

$$P[X_1 > c_1, \dots, X_J > c_J] = \int P[X_1 > c_1, \dots, X_J > c_J | \theta] dF(\theta), \quad (3)$$

where the integrand comes from either equation (1) or equation (2), and  $dF(\theta)$  represents the distribution of  $\theta$  in the population of interest. The assumptions of unidimensionality, monotonicity and local independence can be relaxed in various ways (e.g., Sections 2.2, 3.2 and 4.1.3 below), but this basic model has been the foundation of much nonparametric scale construction.

### 2.1.1 Nonparametric scale construction

For dichotomous items ( $X_j = 0$  or 1 indicating incorrect or correct answer) a theory of scale construction—selecting groups of items that hang together well in the sense that the monotone homogeneity model is probably appropriate for them—has existed at least since Mokken (1971; 1997; see also Molenaar, 1991; 1997). The principal tools of that theory are adaptations of Loevinger's (1948)  $H$  coefficients, comparing the marginal covariance  $\text{Cov}(X_i, X_j)$  of each item pair with the maximum covariance  $\text{Cov}_{\max}(X_i, X_j)$  possible, preserving the margins of the observed  $X_i \times X_j$  table. The bound  $\text{Cov}_{\max}(X_i, X_j)$  is obtained by adjusting the table to remove Guttman errors (e.g., Molenaar, 1991); and indeed the original formulas for the  $H$  coefficients were expressed as ratios of Guttman errors (Mokken, 1997).

A related modeling condition for dichotomous items is *invariant item ordering* which says that for every item pair  $i$  and  $j$ , either  $P_i(\theta) \leq P_j(\theta)$  or the reverse inequality is maintained uniformly for all  $\theta$ .

Rosenbaum (1987a,b) and Sijtsma and Junker (1996, 1997) explore and extend this idea and provide scale construction examples. Mokken's *double monotonicity* model incorporates both the monotone homogeneity and invariant item ordering assumptions.

The  $H$  coefficients are directly sensitive only to high or low correlations between items, rather than to local independence given  $\theta$  as in equations (1) and (2). If the correlations are near zero, we may be unsatisfied to assume that such a  $\theta$  exists (see for example the discussion of co-monotonicity in Junker Ellis, 1997). While a perfect Guttman scale would produce  $H$  coefficients equal to one, large  $H$  coefficients provide only indirect evidence of a  $\theta$  "explaining" covariation in the item responses in the sense of local independence.

More direct attacks on the problem of establishing such a  $\theta$  from data analysis have been pursued by Stout, Ramsay and their students and colleagues. Stout (1990; and subsequent work, for example Stout, Habing, Douglas, Kim, Roussos Zhang, 1996) basically constructs a proxy for  $\theta$  from the total score on a specially-selected subset of the items and uses it to test a weakened version of monotone homogeneity, Stout's essential unidimensionality model. Molenaar and Stout (personal communication) are currently working to find common ground between these two approaches. Ramsay (1991) constructs nonparametric regression estimates of item (step) response functions using total score or rest score (see Section 2.1.2 below for definitions) as a proxy for  $\theta$ , which allows one to explore non-monotonicity. Ramsay (1995) constructs nonparametric regression estimates of the joint response surface of all items on the test, using an appropriate proximity measure to determine which response patterns are "close together", which enables the exploration of both non-monotonicity and lack of unidimensionality. Stout's and Ramsay's methods are generally more computationally complex, and seem to require larger examinee and item sample sizes, than the methods initiated by Mokken and developed by Molenaar, Sijtsma, Rosenbaum and their colleagues and students. Thus the Mokken techniques have been more widely used in smaller social survey and experimental psychology settings.

Extending Mokken scaling and invariant item ordering methods to the case of polytomous item responses has been more recent. Molenaar (1991) first observed that a direct generalization of the  $H$  coefficients based on covariances made sense, and provided an efficient computational method for obtaining  $\text{Cov}_{\max}(X_i, X_j)$  in the polytomous case. This provides the basis for Mokken-style exploration for monotone homogeneity in polytomous items, as illustrated by Hemker, Sijtsma and Molenaar (1995). Generalizing invariant item ordering to the polytomous case turns out to be somewhat delicate. For example, Molenaar's (1997) double monotonicity model requires  $P[X_j > c_1|\theta] \leq P[X_k > c_2|\theta]$ , or the opposite inequality, to hold uniformly in  $\theta$ , for each  $j$ ,  $k$ ,  $c_1$  and  $c_2$ . But this model fails to maintain inequalities like  $E[X_j|\theta] \leq E[X_k|\theta]$  uniformly in  $\theta$ , which are a natural generalization of the invariant item ordering condition given above for dichotomous items. Scheiblechner's (1995) ISOP model requires that  $P[X_j > c|\theta] \leq P[X_k > c|\theta]$  uniformly in  $c$  and  $\theta$ , but allows ISRF's for items  $j$  and  $k$  at different thresholds  $c_1$  and  $c_2$  to cross; this model does maintain the order of expected item scores across  $\theta$ . See Sijtsma and Hemker (1998) for a complete account.

Another approach to understanding scaling by the Mokken model—in dichotomous, polytomous, and more general settings—was initiated by Holland (1981), Rosenbaum (1984), and Holland and Rosenbaum (1986; see also Meredith, 1965), who more or less independently developed its three fundamental assumptions under the name "monotone unidimensional latent variable model", and continued by Junker (1993), Ellis and Junker (1997) and Junker and Ellis (1997) under the name "strict unidimensionality model". They

combine Holland and Rosenbaum's (1996) *conditional association (CA)* condition

$$\begin{aligned} \forall \text{ partitions } \mathbf{X} = (\mathbf{Y}, \mathbf{Z}), \forall f, g \text{ non-decreasing; } \forall h(\mathbf{Z}), \\ \text{Cov}(f(\mathbf{Y}), g(\mathbf{Y})|h(\mathbf{Z})) \geq 0, \end{aligned} \quad (4)$$

with a *vanishing conditional dependence* condition

$$\begin{aligned} \forall J, \text{ as } M \rightarrow \infty, (X_1, \dots, X_J) \text{ become independent,} \\ \text{given } (X_{J+1}, \dots, X_{J+M}), \end{aligned} \quad (5)$$

to obtain a complete characterization of an infinite-item-pool formulation of the basic monotone homogeneity model, in which  $\theta$  is both genuinely latent and consistently estimable, in terms of the joint distribution of observable item responses.

### 2.1.2 Stochastic ordering

A side effect of the effort to understand how to characterize and test the monotone homogeneity model has been a selection of other model testing criteria, such as Junker's (1993) "manifest monotonicity" property for dichotomous items following the monotone homogeneity model,

$$P[X_j = 1 | X_+^{(-j)} = s] \text{ is non-decreasing in } s. \quad (6)$$

This property, and examples showing that it does not hold when the "rest score"  $X_+^{(-j)} = \sum_{i \neq j} X_i$  is replaced by the total score  $X_+ = \sum_{i=1}^J X_i$ , are included in unpublished work of Molenaar and Tom Snijders (Junker, 1993; Junker Sijtsma, 2000).

Proving manifest monotonicity depends on establishing a "stochastic ordering" property for  $\theta$ , given the total score  $X_+$  [or equivalently the rest score  $X_+^{(-j)}$ ]:

$$P[\theta > c | X_+ = s] \text{ is non-decreasing in } s, \forall c. \quad (7)$$

Hemker, Sijtsma, Molenaar and Junker (1996, 1997) call this property "SOL" (Stochastic Ordering of the Latent trait by the sum score), and show that, surprisingly, this property does not generalize to "most" nonparametric ordered-polytomous response IRT models. Thus for example, rules based on cutoffs for  $X_+$  need not be most powerful for "mastery decisions" in the sense of  $\theta > c$ ; on the other hand, such cutoff rules for  $X_+$  are most powerful for mastery decisions in the nonparametric dichotomous response case (Grayson, 1988; Huynh, 1994).

In the process of developing these stochastic ordering ideas, Hemker et al. (1997) and Hemker and Sijtsma (1999) have developed a taxonomy of nonparametric and parametric item response models, that usefully complements the taxonomy of Thissen and Steinberg (1986). The Hemker taxonomy is based on the cumulative, continuation-ratio, and adjacent-category logits that are commonly used to define parametric families of polytomous IRT models. Common forms of graded response models (GRM; Samejima, 1997 for example), sequential models (SM; Tutz, 1990; Mellenbergh, 1995; Samejima, 1969, 1995), and partial credit models (PCM; Masters 1982) assume, respectively, that the logit functions  $\text{logit } P[X_j > c | \theta]$ ,  $\text{logit } P[X_j > c | X > c - 1, \theta]$ , and  $\text{logit } P[X_j = c + 1 | X_j \in \{c, c + 1\}, \theta]$  are *linear* in  $\theta$ . Hemker's

analogous nonparametric model classes, the np-GRM, np-SM and np-PCM, assume only that these logits are *non-decreasing* in  $\theta$ .

This taxonomy is a powerful way to organize ideas about model definition and model development in applications of both parametric and nonparametric IRT. It follows that the np-PCM class is nested within the np-SM class, which is nested within the np-GRM class; moreover all three linear-logit families above (GRM, SM and PCM) are in fact subsets of the np-PCM class. As Hemker and Sijtsma (1999) and Van der Ark (1999) show, this approach also highlights links between polytomous IRT and the machinery of generalized linear models (McCullagh Nelder, 1989), just as it has long been realized that parametric dichotomous IRT is basically multivariate mixed effects logistic regression (e.g., Douglas Qui, 1997; see also Lee Nelder, 1996). It is important to realize however that of all of the models studied by Hemker and his colleagues, only the parametric PCM of Masters (1982) and its special cases, have the nice stochastic ordering property SOL (see 7 above).

## 2.2 Some Interesting Questions

An ongoing question in this area is developing adequate data analysis methodology. Most of what can now be done, in the dichotomous and polytomous cases, is encoded in the computer program MSP (Molenaar Sijtsma, 1999). Sijtsma (1998) provides an excellent survey of nonparametric IRT approaches to the analysis of dichotomous item scores; Molenaar (1997) and Sijtsma and Van der Ark (this volume) survey extensions to the polytomous case. Snijders (this volume) introduces Mokken scaling tools for multilevel data as well. Ellis (1994) has re-examined Mokken's hypothesis testing framework for the  $H$  coefficients, and developed, in principle, new tests based on the theory of order-restricted inference of Robertson, Wright and Dykstra (1988). The same methods may be useful to develop tests of manifest monotonicity. In addition to Holland and Rosenbaum's (1986) applications of the Mantel-Haenzel test, Yuan and Clarke (1999) have also developed asymptotic theory for testing the conditions of Junker (1993), that should be adaptable to the conditions of Junker and Ellis (1997). A more direct application of the theory of order-restricted inference to testing CA and related conditions was given by Bartolucci and Forcina (in press).

In terms of the models themselves, Sijtsma and Van der Ark (this volume) discusses several interesting problems, and current progress, related to the lack of SOL (7) in ordered polytomous IRT models and to the sensitivity and specificity of the manifest monotonicity condition (6) for detecting (violations of) the monotone homogeneity model. One of the strengths of the nonparametric approach to dichotomous IRT is that it usually assures us, under very general circumstances, that simple summaries of the data are informative about inferences we wish to make, yet the current evidence suggests that there are no such simple summaries for inferences about ordered-polytomous data. Understanding the impact that this has on progress with theory and applications of polytomous IRT models will surely entail facing and solving the problems that Sijtsma and Van der Ark discuss.

Finally, some of the machinery developed to characterize monotone homogeneity models seems ready to apply to common modifications of this basic model. For example, conditional association (4) is basically an extreme sharpening of the well-known fact that inter-item correlations are nonnegative under monotone homogeneity. Post (1992; Post Snijders, 1993) has established a similar fact about a class of nonparametric probabilistic unfolding models: the inter-item correlation matrix has a band of positive correlations near the main diagonal, bordered by bands of negative correlations. Is there a sharpening of the Post result analogous to conditional association? Could this be combined with the vanishing conditional dependence condition

of (5) to produce a characterization of Post's models? Could such a result follow from a “folding” of the monotone homogeneity model to produce nonparametric unfolding models, along the lines of Verhelst and Verstralen (1993) or Andrich (1996)? In another direction, much of the work in Ellis and Junker (1997) and Junker and Ellis (1997) does *not* depend on the latent variable  $\theta$  being unidimensional. Junker and Ellis (1997) for example point out that their item-step “true scores”  $P[X_j > c|\tau(X)]$  should form a manifold of the same dimension as the underlying latent variable  $\theta$ . A characterization theorem may again result, if a weakening of conditional association to accommodate multidimensional  $\theta$  could be developed. Such theorems are very helpful in focusing our ideas about what the observable data from a monotone homogeneity model, a probabilistic unfolding model, or a multidimensional nonparametric IRT model, should behave like.

### 3 Parametric IRT: Modeling Dependence

Parametric IRT, as surveyed for example in the edited volumes of Fischer and Molenaar (1995) and Van der Linden and Hambleton (1997), is a well-established, wildly successful statistical modeling enterprise. IRT models have greatly extended the data analytic reach of psychometricians, social scientists, and educational measurement specialists. My summary of parametric IRT will be even less complete, relative to the vast parametric IRT literature, than my summary of nonparametric IRT.

A basic and familiar model in this area is the “two-parameter logistic”, or 2PL, model for dichotomous item response variables (e.g., Chapter 1 of Van der Linden Hambleton, 1997), given by the monotone homogeneity assumptions in Section 2 and the assumption of a logistic form for the item response functions

$$P_j(\theta_i; \alpha_j, \beta_j) \equiv P[X_{ij} = 1 | \theta_i, \alpha_j, \beta_j] = \frac{1}{1 + \exp(-\alpha_j[\theta_i - \beta_j])}, \quad (8)$$

describing the dichotomous response of examinee  $i$  to item  $j$ . The “discrimination” parameter  $\alpha_j$  controls the rate of increase of this logistic curve, and is directly related to the Fisher information for estimating  $\theta$ , and the “difficulty” parameter  $\beta_j$  is the location on the  $\theta$  scale at which the information is maximal; note also that at  $\theta_i = \beta_j$ ,  $P[X_{ij} = 1] = 1/2$ . The 3PL (three-parameter logistic) model extends the 2PL model by adding a non-zero lower asymptote to each item response function; on the other hand the Rasch or 1PL (one-parameter logistic) model is a restriction of the 2PL model obtained by setting  $\alpha_j$  identically equal to some constant, usually 1. Such parametric IRT models, extended by hierarchical mixture/Bayesian modeling and estimation strategies, make it possible in principle and in practice to incorporate covariates and other structure. Many violations of the basic local independence assumption of IRT models are in fact due to unmodeled heterogeneity of subjects and items, that can now be explicitly modeled using these methods.

The main purpose of this section is to introduce a general modeling framework and highlight a few developments in parametric IRT, some old and some new, that will be relevant to my discussion of applying IRT and related models to cognitive assessment problems in Section 4 below.

#### 3.1 Two-Way Hierarchical Structure

The estimation of group effects and the use of examinee and item covariates in estimating item parameters plays an important role in the analysis of large multi-site educational assessments such as the National

Assessment of Educational Progress (NAEP; e.g., Algina, 1992; Johnson, Mislevy and Thomas, 1994; and Zwick 1992). These efforts, which go back at least to Mislevy (1985; see also Mislevy Sheehan, 1989) can be recognized as the wedding of hierarchical linear or multi-level modeling methodology with standard dichotomous and polytomous IRT models. The general model is a two-way hierarchical structure for  $N$  individuals and  $J$  response variables, as follows

$$\left. \begin{array}{ll} \text{First level: } & X_{ij} \sim P(\theta_i; \gamma_j)^{X_{ij}} [1 - P(\theta_i; \gamma_j)]^{(1-X_{ij})}, \\ & i = 1, \dots, N; j = 1, \dots, J \\ \\ \text{Second level: } & \theta_i \sim f_i(\theta | \lambda_f), \text{ each } i \\ & \gamma_j \sim g_j(\gamma | \lambda_g), \text{ each } j \\ \\ \text{Third level: } & \lambda_f \sim \phi_f(\lambda_f) \\ & \lambda_g \sim \phi_g(\lambda_g) \end{array} \right\} \quad (9)$$

where  $\theta_i$  is the usual person parameter,  $\gamma_j$  is the vector of item parameters for item  $j$  (e.g.,  $\gamma_j = (\alpha_j, \beta_j)$  in the 2PL model above), and where independence is assumed between  $j$ 's conditional on  $\theta_i$  at the first level and between  $i$ 's at the second level. The terms  $\lambda_f$  and  $\lambda_g$  represent sets of hyperparameters needed to specify the person distributions  $f_i$  and item distributions  $g_j$ , with hyperprior distributions  $\phi_f$  and  $\phi_g$ , respectively. Terms in the first level, for example, are multiplied together to produce the usual joint likelihood for the  $N \times J$  item response matrix  $[X_{ij}]$ ; the second and third levels can be used to impose constraints on the first level parameters and latent variable, to deduce what integrations are needed for marginal likelihood approaches, etc. The model (9) is expressed for dichotomous items, for simplicity of exposition, but can easily be generalized to polytomous items, or combinations of item types (see for example Patz Junker, 1999a; 1999b). It is also usual to assume for  $f_i(\theta)$  a single latent trait distribution not depending on  $i$ , and similarly for  $g_j$ .

Recent advances relax these assumptions, allowing for much more flexible parametric modeling of item response data. We may allow the distribution of  $\theta$  to depend hierarchically on examinee covariates, that is, instead of taking  $f_i(\cdot)$  to be a single latent trait distribution in (9), we can allow it to depend on examinee covariates to model population heterogeneity, as in the multi-group IRT models of Mislevy (1985) and Bock and Zimowski (1997), or to reflect hierarchical linear structure as in Fox and Glas (1998). We may also elaborate  $g_j(\cdot)$ , for example by building linear structure into the item parameters. For example in the 2PL model, where  $\gamma_j = (\alpha_j, \beta_j)$ , we might take

$$\begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_{J-1} \\ \beta_J \end{bmatrix} = Q \begin{bmatrix} \psi_1 \\ \psi_2 \\ \vdots \\ \psi_K \end{bmatrix}, \quad (10)$$

where  $Q$  is an appropriate design matrix of full column rank, to reflect common sources ( $\psi_k$ 's) of item difficulty ( $\beta_j$ 's) across items. In the case of Rasch (1PL) IRF's, this is the linear logistic test model (LLTM; Scheiblechner, 1972; Fischer, 1973). This model and its various generalizations (e.g., Glas Verhelst, 1989;

Patz Junker, 1999b) continues to be used for psychological experiments with multiple outcomes per subject (e.g., Fischer and Molenaar, 1995 and the references therein) and for research in cognitively-motivated test design (Embretson, 1995; 1999). There is no reason to restrict attention to the  $\beta$ 's, and for example Embretson (1999) has explored a similar decomposition of the  $\alpha$ 's in a 2PL model.

As Bock and Zimowski (1997) point out, these generalizations of the basic IRT model both simplify and unify parametric approaches to many thorny test analysis questions, including differential item functioning and item parameter drift, nonequivalent groups and vertical equating, two-stage testing and matrix-sampled educational assessment survey work, etc. The computer programs ConQuest (Wu, Adams Wilson, 1997) and BILOG-MG (Zimowski, Muraki, Mislevy Bock, 1997) provide fairly general E-M based solutions when the underlying IRT model is the 1PL (ConQuest), or 2PL or 3PL (BILOG-MG). Fox Glas (1998) and Patz and Junker (1999a; 1999b) give two different Markov chain Monte Carlo (MCMC) approaches to the problem. Generalization to polytomous items, facets-style rated response models, and mixtures of item types are conceptually, and often computationally, straightforward; see for example Glas and Verhelst (1989) and Patz and Junker (1999a; 1999b).

### 3.2 Some Multidimensional Models

Research in multidimensional IRT models has concentrated on additive and conjunctive combinations of multiple traits to produce probabilities of response. Additive models, known as *compensatory* models in much of the literature, replace the unidimensional latent trait  $\theta$  with an item-specific, known (e.g., Embretson, 1991; Stegelmann, 1983; Kelderman and Rijkes, 1994; and Adams, Wilson Wang, 1997) or unknown (e.g., Reckase, 1985; Wilson, Wood Gibbons, 1983; Fraser MacDonald, 1988; Muraki Carlson, 1995) linear combination of components  $a_{j1}\theta_1 + \dots + a_{jd}\theta_d$  of a  $d$ -dimensional latent trait vector, for example in the dichotomous response case

$$P[X_j = 1 | \theta_1, \dots, \theta_d] = P(a_{j1}\theta_1 + \dots + a_{jd}\theta_d - \beta_j),$$

where  $P(t)$  might be the logistic or probit response function for example. Béguin and Glas (1998) survey the area well (see also several contributed chapters in Van der Linden Hamilton, 1997) and give an MCMC algorithm for estimating these models; Gibbons and Hedeker (1997) pursue related developments in biostatistical and psychiatric applications.

Conjunctive models are often referred to as *noncompensatory* or *componential* models in the literature. These models (e.g., Embretson, 1985, 1997) combine unidimensional models for components of response conjunctively, so that

$$P[X_j = 1 | \theta_1, \dots, \theta_d] = \prod_{\ell=1}^d P_{j\ell}(\theta_\ell)$$

where  $P_{j\ell}(\theta_\ell)$  are parametric unidimensional dichotomous response functions. The usual interpretation is that the  $P_{j\ell}(\theta_\ell)$  represent skills or subtasks all of which must be performed correctly in order to generate a correct response to the item itself. Janssen and De Boeck (1997) give a recent application.

Compensatory structures are attractive because of their conceptual similarity to factor analysis models. They have been very successful in aiding the understanding of how student responses can be sensitive to major content and skill components of items, and in aiding parallel test construction when the underlying

response behavior is multidimensional (e.g., Ackerman, 1994). Noncompensatory models are largely motivated from a desire to model cognitive aspects of item response, a topic to which we will return in Section 4. Embretson (1997) also reviews blends of these two approaches (her general component latent trait models; GLTM).

### 3.3 Models That Accommodate Extra Behavioral Features of Assessment

In addition to providing a way to model dependence of item responses on specific examinee and item covariates, the hierarchical or multi-level approach to IRT also allows us to model extra behavioral features of assessment. This is a largely unexplored area, but one which is worth further study, since current models and methods largely ignore these features, relegating them to the “error distribution” of the model.

One example of this sort of work involves recent efforts to more elaborately model the behavior of raters in rated item response data. When only one rater rates each item, it may be sufficient to treat each rating as a different, locally independent pseudo-item—so that the first level in (9) contains one factor for each rater  $\times$  item  $\times$  examinee combination—and to model the rater effect as a linear influence on the item’s difficulty parameter  $\beta_j$ . Mathematically this is equivalent to the LLTM model sketched above, but it has come to be known in this setting as the “Facets model” (Linacre, 1989; Engelhard, 1994, 1996).

For both formative and summative evaluation of raters, a number of multiple-read rating designs are now commonplace (Wilson Hoskens, 1999), including designs with as many as six raters per item (e.g., Sykes Heidorn, 1999). Thus each examinee performance is measured several correlated but fallible times. Junker and Patz (1998) showed that the usual Facets model formulation in which the likelihood is the product of LLTM-style factors for each rating of each item, accumulates information about  $\theta$  too optimistically, so that with only one item response, in the limit as the number of raters grows, the standard error for estimating or predicting  $\theta$  apparently goes to zero, contradicting the notion (e.g., Junker, 1993) that the number of items should tend to infinity in order to make the error of estimation of  $\theta$  vanishingly small. Instead, models are needed that appropriately accumulate information from multiple ratings to the single item response being rated, and then accumulate information across item responses to learn about  $\theta$  itself.

Wilson and Hoskens’ (1999) rater bundle model (RBM) attacks this problem by replacing the Facets product across raters for each item in level one of (9) with a log-linear model that models the dependence between ratings of the item, conditional on  $\theta$ . Their approach seems very useful for, e.g., modeling “table effects” and other rater dependence phenomena that follow when raters are allowed or encouraged to discuss ratings amongst themselves to increase rating quality and inter-rater reliability.

Patz, Junker and Johnson’s (1999) hierarchical rater model (HRM) provides an alternative approach that posits a “latent rating”  $\xi_{ij}$  (not unlike Maris’, 1995, notion of latent responses; also note the connection with the data augmentation method in Bayesian computation, e.g., Tanner, 1996) that follows a conventional IRT model, such as

$$P[\xi_{ij} = 1 | \theta_i, \alpha_j, \beta_j] = \frac{1}{1 + \exp(-\alpha_j[\theta_i - \beta_j])}$$

(or a polytomous variant), and a “signal detection model” for each rater  $r$  rating that item response such as

$$P[X_{ijr} = 1 | \xi_{ij}] = (1 - p_{01r})^{\xi_{ij}} p_{10r}^{1-\xi_{ij}},$$

where  $p_{11r}$  is the probability that rater  $r$  rates a response with latent rating  $\xi_{ij} = 1$  as a 0, and  $p_{10r}$  is the probability that rater  $r$  rates a response with latent rating  $\xi_{ij} = 0$  as a 1. The factors  $P[X_{ijr} = 1 | \xi_{ij}]$

now appear at level one of the hierarchy (9), and the factors  $P[\xi_{ij} = 1 | \theta_i, \alpha_j, \beta_j]$  effectively form a new level between levels one and two of (9). This produces a model that behaves, for multiple discrete ratings, much as a standard generalizability theory model would behave for multiple continuous ratings (see Verhelst Verstralen, 1999, for a similar development). In addition to providing an appropriate way to combine information from multiple raters to learn about student performance, the HRM makes possible calibration and monitoring of individual rater behavioral effects such as, in the case of polytomous responses, rater severity and rater precision.

### 3.4 Some Current Questions

Almost any assessment phenomenon—from between-examinee dependence due to institutional or socio-logical factors, to behavioral aspects of raters, to the analysis of item responses into requisite examinee skills or item features—can be expressed in the hierarchical mixture/Bayes modeling framework, because of its conceptual simplicity. Recent advances in computation, and MCMC methods in particular, have made it possible to estimate a vastly wider variety of these models than would have been imaginable even ten years ago. Questions motivating this work inevitably involve identifying phenomena that are worth detailed parametric modeling, and seeing if the computational machinery can be pushed to estimate models of these phenomena. Recent examples include multidimensional (Gibbons Hedeker, 1992) and hierarchical (Bradlow, Wainer Wang, 1999) modeling of testlets; blending IRT and hierarchical linear models, and behavioral models such as the rater models described above.

Speeding up the computations with approximations (including formal and informal applications of Laplace’s method such as Rigdon Tsutakawa, 1983 and Kass, Tierney Kadane, 1990; blends of Monte Carlo and E-M approaches as surveyed in Tanner, 1996; and variational methods, e.g., Jaakkola Jordan, 1999) continues to be an essential and fruitful avenue of research.

An avenue that has not been explored as much in IRT work is choosing models for which sufficient statistics are simple and interpretable. The most familiar “basic model” of this sort is of course the Rasch model, but as we shall see in the next section some current cognitively-motivated measurement problems may require a new “basic model”.

## 4 Measurement Challenges Posed by Cognitive and Embedded Assessments

In recent years, as cognitive theories of learning and instruction have become richer, and computational methods to support assessment have become more powerful, there has been increasing pressure to make assessments truly criterion referenced, that is, to “report” on student achievement relative to theory-driven lists of examinee skills, beliefs and other cognitive features needed to perform tasks in a particular assessment domain. For example Baxter and Glaser (1998) and Nichols and Sugrue (1999) present compelling cases that assessing examinees’ cognitive characteristics can and should be the focus of assessment design. In a similar vein, Resnick and Resnick (1992) advocate standards-referenced or criterion-referenced assessment closely tied to curriculum, as a way to inform instruction and enhance student learning.

Appropriate criterion-referenced testing can also be an effective teaching tool when embedded directly in teaching practice. Indeed there is substantial argument and evidence, as summarized for example by Bloom (1984), that part of what distinguishes higher student achievement in “mastery learning” and individualized

tutoring settings as opposed to the conventional classroom, is the use of frequent and relatively unobtrusive formative tests coupled with feedback for the students and corrective interventions by the instructor, and follow-up tests to determine how much the interventions helped. This approach continues to be advocated as part of a natural and effective apprenticeship style of human instruction (e.g., Gardner, 1992), and it is the basis of many computer-based intelligent tutoring systems (ITS's, e.g., Anderson, 1993; and more broadly Shute Psotka, 1996). Here too, a decomposition of assessment items into appropriate cognitive attributes is important: feedback and/or corrective action in a mastery class or from an ITS depends on knowing which cognitive attributes the examinee has and has not mastered.

Cognitive assessment models must generally deal with a more complex goal than linearly ordering examinees, or partially ordering them in a low-dimensional Euclidean space, which is what IRT has been designed and optimized to do. The goal of cognitive assessment can be thought of producing, for each examinee, a checklist of skills or other cognitive attributes that the examinee may or may not possess, based on the evidence of tasks performed by the examinee. The checklist of cognitive features in a cognitive assessment generally comes from an analysis of the cognitive attributes needed to successfully perform each task in a domain of interest. For a particular set of tasks, this analysis can be encoded in an incidence matrix, the *Q-matrix* (e.g., Tatsuoka, 1990, 1995), which we will write as a  $J \times K$  matrix  $Q = [Q_{jk}]$  of 0's and 1's with entries

$$Q_{jk} = \begin{cases} 1, & \text{if attribute } k \text{ is required by task } j \\ 0, & \text{if not} \end{cases} \quad (11)$$

While the *Q*-matrix does not capture all of the structure we may be interested in (*Q* treats the skills in a flat, non-time-ordered manner, and there may be both hierarchical and time-order structure in the skills as they are applied to a task), it is a useful bookkeeping device.

Many attempts (e.g., as surveyed by Roussos, 1994) to blend IRT and cognitive measurement are based on a linear decomposition of item parameters, as in the LLTM, or on a linear decomposition of the latent trait  $\theta$ , as in the multidimensional compensatory IRT models. Compensatory IRT models, like factor analysis models, can be sensitive to relatively large components of variation in examinee ability or propensity to answer items correctly, but they are generally not sensitive to the finer components of variation that are often of interest in cognitive assessment. Of course, whether cognitive assessment data actually supports models that track this finer level of variation is an empirical matter. LLTM-style models can be sensitive to these finer components of variation *among items* but are not at all sensitive to components of variation *among examinees*. Noncompensatory approaches (e.g., Embretson, 1985, 1997) are intended to be sensitive to finer variations among examinees, in situations in which several cognitive components are required simultaneously for successful task performance; but they may be attempting to estimate more (per-skill latent trait values and skill-response parameters) than assessment data of the type we will consider can be expected to support. These models are also related to the probability matrix decomposition models of Maris, De Boeck and Van Mechelen (1996).

#### 4.1 Three Approaches to Cognitive Assessment

To illustrate the differences between traditional IRT approaches and cognitively-motivated approaches to assessment, consider two published models intended to deal with essentially the same data: task performance by students learning the LISP programming language using one of the computer based intelligent tutoring

systems developed by John R. Anderson and his colleagues at Carnegie Mellon University (e.g., Anderson, Corbett, Koedinger Pelletier, 1995). The first model is the assessment model actually embedded in the tutoring software, as described by Corbett, Anderson and O'Brien (1995); the second is an IRT-based model for essentially the same data, as presented by Draney, Pirolli and Wilson (1995). Then I will present a third model to illustrate that, although the modeling traditions in IRT are not much like those in cognitive assessment modeling, the fundamental techniques and concepts from IRT modeling are quite useful in cognitive assessment.

#### 4.1.1 The Corbett/Anderson/O'Brien model

For the “knowledge tracing model” embedded in the LISP tutor software (Corbett, Anderson O'Brien, 1995) for this data, define

$X_{ij}(n)$	= 1 or 0	indicating whether or not student $i$ performed task $j$ correctly at time $n$
$Q_{jk}$	= 1 or 0	indicating whether or not task $j$ requires skill $k$
$\alpha_{ik}(n)$	= 1 or 0	indicating whether or not student $i$ possesses skill $k$ at time $n$
$\xi_{ij}(n)$	= $\prod_{k: Q_{jk}=1} \alpha_{ik}(n)$	indicating whether or not student $i$ has the skills needed for task $j$ at time $n$
$s$	= $P[X_{ij}(n) = 0   \xi_{ij}(n) = 1]$	a universal slip parameter
$g$	= $P[X_{ij}(n) = 1   \xi_{ij}(n) = 0]$	a universal guessing parameter

The model of task performance embedded in the knowledge tracing model is essentially

$$P[X_{ij}(n) = 1] = P[\xi_{ij}(n) = 1](1 - s) + (1 - P[\xi_{ij}(n) = 1])g, \quad (12)$$

where  $P[\xi_{ij}(n) = 1]$  follows a simple conjunctive model (though more complex expressions along the lines of Mislevy, 1996, could be imagined):

$$P[\xi_{ij}(n) = 1] = \prod_{k=1}^K P[\alpha_{ik}(n) \geq Q_{jk}]. \quad (13)$$

It is worth noting that the  $\xi_{ij}(n)$ 's [or the  $\alpha_{ik}(n)$ 's] can be interpreted as playing the role of Maris's (1995) latent responses (they can also be interpreted in terms of the method of data augmentation in Bayesian computation, e.g., Tanner, 1996).

It is the probabilities  $P[\alpha_{ik}(n) \geq Q_{jk}]$ , treating  $\alpha_{ik}(n)$  as the unknown or random quantity, that are actually of primary interest in assessment based on this model of task performance. Given current estimates of  $P[\alpha_{ik}(n) \geq Q_{jk}]$ , the tutor can both identify which skills need additional practice, and select items of suitable difficulty that exercise those skills, to assign to the student next.

Corbett, Anderson and O'Brien (1995) were particularly interested in how to gather evidence about  $P[\alpha_{ik}(n) \geq Q_{jk}]$  as the number of opportunities  $n$  to apply rule  $k$  increases—i.e. they are interested in

modeling learning. To account for the order in which correct and incorrect actions are observed when skill  $k$  is called for, they suggest treating  $\alpha_{ik}(n)$  as a hidden Markov chain with an absorbing state,

$$\begin{aligned} P[\alpha_{ik}(n) \mid \bar{\alpha}_{ik}(n-1)] &= T \\ P[\alpha_{ik}(n) \mid \alpha_{ik}(n-1)] &= 1 \\ P[\bar{\alpha}_{ik}(n) \mid \bar{\alpha}_{ik}(n-1)] &= 1 - T \\ P[\bar{\alpha}_{ik}(n) \mid \alpha_{ik}(n-1)] &= 0, \end{aligned} \tag{14}$$

where “ $\alpha_{ik}(n)$ ” stands for “ $\alpha_{ik}(n) = 1$ ” and “ $\bar{\alpha}_{ik}(n)$ ” stands for “ $\alpha_{ik}(n) = 0$ ”.

The tutoring system is arranged to directly observe evidence  $a_{ik}(n) = 0$  or 1 that the correct action was performed when skill  $k$  was called for, greatly simplifying the inferential task. Corbett et al. (1995) posit the following relationships between the hidden state  $\alpha_{ik}(n) = 0$  or 1 and the observable evidence  $a_{ik}(n) = 0$  or 1:

$$\begin{aligned} P[\alpha_{ik}(n) \mid a_{ik}(n)] &= P[\alpha_{ik}(n-1) \mid a_{ik}(n)] \\ &\quad + (1 - P[\alpha_{ik}(n-1) \mid a_{ik}(n)]) T \\ P[\alpha_{ik}(n) \mid \bar{a}_{ik}(n)] &= P[\alpha_{ik}(n-1) \mid \bar{a}_{ik}(n)] \\ &\quad + (1 - P[\alpha_{ik}(n-1) \mid \bar{a}_{ik}(n)]) T \\ P[\alpha_{ik}(n-1) \mid a_{ik}(n)] &= \{(1-s)P[\alpha_{ik}(n-1)]\} \\ &\quad / \{(1-s)P[\alpha_{ik}(n-1)] \\ &\quad + gP[\bar{\alpha}_{ik}(n-1)]\} \\ P[\alpha_{ik}(n-1) \mid \bar{a}_{ik}(n)] &= \{sP[\alpha_{ik}(n-1)]\} \\ &\quad / \{sP[\alpha_{ik}(n-1)] \\ &\quad + (1-g)P[\bar{\alpha}_{ik}(n-1)]\}. \end{aligned} \tag{15}$$

Given *a-priori* fixed values of  $T$ ,  $s$ ,  $g$ , and a probability  $p_0$  that each skill is in the learned state before the tutoring begins, we may substitute  $p_0$  for  $P[\alpha_{ik}(n-1)]$  when  $n = 1$ , and the above formulas give an algorithm for recursively updating  $P[\alpha_{ik}(n)]$  on the basis of the observed sequence of correct and incorrect actions in the first  $n$  opportunities to apply rule  $k$ . This is a particularly simple and fast computational method, capable of updating the tutor’s model of the student’s skills in real time as the student works with the tutor.

It is interesting to note that in order to produce a well-fitting model, Corbett et al. (1995) had to allow the probabilities  $T$ ,  $s$ ,  $g$ , and the probability  $p_0$  that each skill was already in the learned state before the tutoring began, to be perturbed differently from overall population values for each student  $i$ ; effectively, they allowed these parameters to become random effects. Thus, in addition to the individual differences in skills acquisition that the model had been designed to detect, there were substantial individual differences in starting knowledge of the students, in tendency to slip or guess, and in the rate of learning, under this model.

#### 4.1.2 The Draney/Pirolli/Wilson model

Draney, Pirolli and Wilson (1995) develop an LLTM-style model to analyze essentially the same data. The model they consider begins with an indicator  $a_{ijk}(n) = 1$  if student  $i$  performs correctly when the  $n$ th opportunity to use skill  $k$  occurs, under condition  $j$ ; and  $a_{ijk}(n) = 0$  otherwise. These  $a_{ijk}(n)$  differ from

Corbett et al.'s (1995)  $a_{ik}(n)$  only in that the task that provides a context for performing the skill is allowed to affect the difficulty of correct skill performance. In the Draney et al. (1995) model, the "skill response functions" are given by

$$P[a_{ijk}(n) = 1 \mid \theta_i, \tau_j, \delta_k, \gamma] = \frac{1}{1 + \exp(-\theta_i + \tau_j + \delta_k - \gamma \log(n))}.$$

This model essentially decomposes the difficulty parameter  $\beta$  in the Rasch model according to a  $Q$  matrix, as in equations (10) and (11), in terms of parameters for task,  $\tau_j$ , skill,  $\delta_k$ , and slope of the learning curve,  $\gamma$ . The logarithmic dependence on  $n$  is intended to follow the development of Anderson (1993, Appendix to Chapter 3). If the decomposition of tasks into skills is complete, and the skills are of a suitable granularity, Anderson's (1993) ACT-R theory predicts that skill "performances" will be approximately independent of one another, given the relevant difficulty and student parameters. This is essentially a statement of local independence, so that the "skill response functions" above may be multiplied together in the usual way to form an IRT likelihood.

It is interesting to compare the Corbett et al. and Draney et al. modeling approaches. For example, for a data set similar to that analyzed by Draney et al., the assessment model of Corbett, Anderson and O'Brien (1995, p. 32; see also Draney, Pirolli Wilson, 1995, p. 115) would employ essentially four continuous latent variables, and 33 dichotomous latent attribute indicators, *per student tested*—in addition to 132 parameters to characterize features of the skills being assessed. Using their variation on standard IRT models, Draney et al. provided equivalent or better fit to learning curves, employing *one* continuous latent variable per student tested (Draney, Pirolli Wilson, 1995, p. 109), and 36 parameters for the skills being tested (Draney, Pirolli Wilson, 1995, p. 115). Thus, if the goal is to model learning curves, clearly the more complex Corbett et al. model is not needed.

However, the uses to which the two models can be put, and the substantive interpretations of estimated parameters in the two models, are very different. The Corbett et al. model is immediately useful for diagnosing individual differences in student task performance behavior by relating it directly to specific skills in the task decomposition of their task domain, but can only indirectly assess the validity and reliability of the tasks in the assessment, through fit to learning curves or other summaries of student task performance.

The Draney et al. model is not immediately useful for student diagnosis, since its student parameter is one-dimensional. After fitting the model, Draney et al. go back and rank students based on (empirical Bayes) estimated  $\theta$ 's (in the context of learning curves analysis, the  $\theta$ 's essentially code students' initial facility in the task domain prior to tutoring, much as the random effect version of  $p_0$  does in the Corbett et al. model), and compare estimated skill performance difficulties to the students' aggregate  $\theta$  distribution; see for example their Figure 5.1, p. 112. Such displays allow us to predict which skills a "typical" student with some fixed value of  $\theta$  might be expected to perform correctly, and are very useful communication devices. However, detailed cognitive diagnosis of individual students on the basis of the  $\theta$ 's is not possible, without a post-hoc analysis of some sort. Indeed, for assessing whether individual students have learned particular skills, Draney et al. replace the IRT model with a Bayesian inference network that is focused on inferring the probability that an individual has learned one or more skills, using priors constructed from the fitted IRT model and new data from further attempts to perform the skill(s).

The utility of the Draney et al. IRT model for analyzing important task performance features, and aggregate examinee behaviors, should not be minimized however. For example, Junker, Koedinger and

Trottini (2000) are using essentially the same modeling framework to develop a semi-automatic stepwise variable construction/model selection procedure, with the goal of identifying skills that were either too narrowly or too broadly defined in a cognitive tutor (these result in stylized deviations, or “blips” from the theoretically predicted learning curves for the skills). In a similar vein, Huguenard, et al. (1997; see also Patz et al., 1996) applied a polytomous version of the LLTM to study the relationship between task features and working memory load, using the IRT  $\theta$  parameter to soak up residual between-subjects variation not modeled by experimental and working memory factors.

#### 4.1.3 An IRT-like cognitive assessment model

As the preceding example makes clear, traditional parametric IRT approaches may not be well suited to individual assessment tasks in computer based intelligent tutoring systems. The same issues can also arise in standalone assessments that are designed to assess presence or absence of specific skills—rather than to sort examinees along a linear scale—based on performance on a fixed set of tasks given as a standalone test.

To illustrate this we consider a modified version of the assessment model of Corbett, Anderson and O’Brien (1995). We omit the hidden-Markov learning model and assume that examinee behavior is only observed at the task level, not the skill level. This model can also be connected to multidimensional non-compensatory IRT models, since it is interpretable as a simplified version of Embretson’s (1985, 1997) multicomponent latent trait (MLTM) model. Define

$X_{ij}$	=	1 or 0	indicating whether or not student $i$ performed task $j$ correctly
$Q_{jk}$	=	1 or 0	indicating whether or not task $j$ requires skill $k$
$\alpha_{ik}$	=	1 or 0	indicating whether or not student $i$ possesses skill $k$
$\xi_{ij}$	=	$\prod_{k: Q_{jk}=1} \alpha_{ik}$	indicating whether or not student $i$ has the skills needed for task $j$
$s_j$	=	$P[X_{ij} = 0   \xi_{ij} = 1]$	a per-problem slip parameter
$g_j$	=	$P[X_{ij} = 1   \xi_{ij} = 0]$	a per-problem guessing parameter

Note that the  $Q_{ij}$  are the usual constant entries in the  $Q$ -matrix, and the  $\xi_{ij}$  are deterministic functions of the  $\alpha_{ik}$ ’s and  $Q$ . The goal is to try to estimate the  $\alpha_{ik}$ ’s, or more precisely  $P[\alpha_{ik} = 1]$ , from the task performance data.

Our basic response model [level one in the hierarchy (9)] is

$$\begin{aligned} P[X_{ij} = 1 | \boldsymbol{\xi}, \mathbf{s}, \mathbf{g}] &= \xi_{ij}(1 - s_j) + (1 - \xi_{ij})g_j \\ &= (1 - s_j)^{\xi_{ij}} g_j^{1-\xi_{ij}}, \end{aligned}$$

and so for the entire examinees by tasks matrix  $[X_{ij}]$  of task responses,

$$\begin{aligned} P[X_{ij} = x_{ij}, \forall i, j | \boldsymbol{\xi}, \mathbf{s}, \mathbf{g}] &= \prod_i \prod_j \left[ (1 - s_j)^{\xi_{ij}} g_j^{1-\xi_{ij}} \right]^{x_{ij}} \left[ 1 - (1 - s_j)^{\xi_{ij}} g_j^{1-\xi_{ij}} \right]^{1-x_{ij}} \\ &= \prod_i \prod_j \left[ (1 - s_j)^{x_{ij}} s_j^{1-x_{ij}} \right]^{\xi_{ij}} \left[ g_j^{x_{ij}} (1 - g_j)^{1-x_{ij}} \right]^{1-\xi_{ij}} \end{aligned} \tag{16}$$

To see how estimation of the parameters depends on the data, it is instructive to set up the first stages of an estimation algorithm for the model. In particular, let's assume that prior distributions have been chosen [for levels two and three in the hierarchy (9)], so that  $s_j \sim \pi_s(s_j)$ ,  $g_j \sim \pi_g(g_j)$ ,  $\alpha_{ik} \sim \pi_k^{\alpha_{ik}}(1 - \pi_k)^{1-\alpha_{ik}}$ , and perhaps  $\pi_k \sim \pi(\pi_k)$ . These prior distributions will be used as “placeholders” in the notation below; their particular form will not affect our conclusions in any substantial way. We will calculate the so-called “complete conditional” distributions (e.g., Gelman, Carlin, Stern Rubin, 1995) of each parameter, given the data and the rest of the parameters. Such a calculation is directly useful in setting up an MCMC algorithm for Bayesian estimation of the model parameters (e.g., Patz Junker, 1999a; 1999b), and is also useful in some versions of the E-M algorithm, such as ECME (e.g., Liu Rubin, 1998). However, even when E-M is “possible” for a model like (16), it need not be practical, due to the need to sum over all configurations of the latent  $\alpha$  vector (see equation 17 below). For cases in which relatively few  $\alpha$ 's carry most of the latent trait distribution, MCMC can be a more efficient—albeit approximate—way to estimate the model. This phenomenon has been found in other models with complex discrete latent structure as well (e.g., Seltman, 1999; Snijders Nowick, 1997; and Ter Hofstede, Steenkamp Wedel, 1999).

As usual in setting up an MCMC calculation, we note that the complete conditional distribution for each parameter or latent variable is proportional to the product of only those likelihood and prior factors in (9) depending on that parameter. Thus we obtain for example

$$\begin{aligned} p(s_j|\text{rest}) &\propto (1 - s_j)^{\sum_i x_{ij} \xi_{ij}} s_j^{\sum_i (1-x_{ij}) \xi_{ij}} \pi_s(s_j), \\ p(g_j|\text{rest}) &\propto g_j^{\sum_i x_{ij} (1-\xi_{ij})} (1 - g_j)^{\sum_i (1-x_{ij}) (1-\xi_{ij})} \pi_g(g_j), \end{aligned}$$

$j = 1, \dots, J$ , where “rest” stands for the data and the rest of the parameters in the model. From these it is easy to see the intuitively plausible facts that that the slip parameter  $s_j$  is sensitive only to successes and failures of examinees who we hypothesize (through the values of  $\xi_{ij}$  upon which we have conditioned) do have the requisite skills to perform task  $j$ , and similarly the guessing parameter  $g_j$  is sensitive only to successes and failures of examinees who we hypothesize do not have the requisite skills. These parameters might be given Beta distribution priors, to make life easy for estimation.

More central to the goal of assessing examinees, we see that the complete conditional distributions for  $\alpha_{ik}$  are of the form:

$$\begin{aligned} p(\alpha_{ik}|\text{rest}) &\propto \prod_{j: Q_{jk}=1} \left[ (1 - s_j)^{x_{ij}} s_j^{1-x_{ij}} \right]^{\alpha_{ik} \xi_{ij}^{(-k)}} \left[ g_j^{x_{ij}} (1 - g_j)^{1-x_{ij}} \right]^{1-\alpha_{ik} \xi_{ij}^{(-k)}} \\ &\quad \times \pi_k^{\alpha_{ik}} (1 - \pi_k)^{1-\alpha_{ik}}, \end{aligned}$$

where  $\xi_{ij}^{(-k)} = \prod_{\ell \neq k: Q_{j\ell}=1} \alpha_{i\ell}$ , which indicates presence of all skills needed for task  $j$ , *except* for skill  $k$ . From these complete conditionals we can see that

- When  $\xi_{ij}^{(-k)} = 1$ , the suggested model for  $\alpha_{ij}$  is some sort of Bernoulli, which makes sense.
- When there are no tasks such that both  $Q_{jk} = 1$  and  $\xi_{ij}^{(-k)} = 1$ , then  $\alpha_{ik}$  is sensitive only to its prior distribution  $\pi_k^{\alpha_{ik}} (1 - \pi_k)^{1-\alpha_{ik}}$ : no learning from data occurs.

The second observation is really a version of the credit/blame assignment problem (e.g., VanLehn Niu, in press): we cannot infer whether  $\alpha_{ik}$  was learned, if we are hypothesizing that another needed skill is still unlearned. Roughly, there must be information in the task performance data to allow us to assign credit (when a task is performed correctly) or blame (when it is performed incorrectly) to every cognitive attribute related to the task.

There are essentially two ways to avoid the credit/blame problem. In some situations, skills can be scored directly; this is possible for example within Anderson's ITS's for LISP, algebra and geometry, because students are required to successfully perform each subgoal/skill, with hints and repeated attempts if necessary, before moving on the next subgoal in the task (see also Embretson, 1985). If one cannot score the tasks subgoal by subgoal, one can try to design the assessment (e.g., by hand or using methods from traditional statistical experimental design) so that the tasks efficiently exercise all skills in the skillset in such a way that, taken together, the task performance data informs us about each skill. VanLehn, Niu, Siler, and Gertner (1998) illustrate the inferential difficulties that can result when item sets are not constructed with the goal of designing around the credit/blame problem.

Finally if we want to estimate the skill base rates  $\pi_k$  in the population (a measure of skill difficulty) we may include a fourth set of complete conditionals

$$p(\pi_k | \text{rest}) \propto \pi_k^{c_k} (1 - \pi_k)^{N - c_k} \pi(\pi_k),$$

where  $c_k = \sum_i \alpha_{ik}$  is the number of students who are presently estimated to have skill  $k$  (again, a Beta prior density for  $\pi_k$  is suggested).

## 4.2 A Role for Nonparametric IRT Methods?

Such models as (16) may seem very far removed from the setting in which nonparametric IRT methods are familiar. I want to suggest several ways in which they are not so far removed.

First, suppose that the skill variables  $\alpha_{ik}$  vary independently of one another in a population of students or examinees, as would be consistent with e.g., Anderson's (1993) ACT-R theory, and suppose that the slip and guessing parameters  $s_j$  and  $g_j$  are fixed and satisfy the plausible inequality  $(1 - s_j) \geq g_j$  (but  $s_j$  and  $g_j$  need not be known). A model for the task performance of a randomly-sampled examinee from the population would then be

$$\begin{aligned} P[\mathbf{X} = \mathbf{x}] &= \sum_{\boldsymbol{\alpha}} \prod_j \left[ (1 - s_j)^{\xi_j(\boldsymbol{\alpha})} g_j^{1 - \xi_j(\boldsymbol{\alpha})} \right]^{x_j} \left[ 1 - (1 - s_j)^{\xi_j(\boldsymbol{\alpha})} g_j^{1 - \xi_j(\boldsymbol{\alpha})} \right]^{1 - x_j} p(\boldsymbol{\alpha}) \\ &= \sum_{\boldsymbol{\alpha}} \prod_j P_j(\xi_j(\boldsymbol{\alpha}))^{x_j} [1 - P_j(\xi_j(\boldsymbol{\alpha}))]^{1 - x_j} p(\boldsymbol{\alpha}) \end{aligned} \quad (17)$$

where  $p(\boldsymbol{\alpha})$  is a product measure over the space of binary skills  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$ , and  $P_j(\xi_j(\boldsymbol{\alpha}))$  is the monotone function

$$P_j(\xi_j(\boldsymbol{\alpha})) = (1 - s_j)^{\xi_j(\boldsymbol{\alpha})} g_j^{1 - \xi_j(\boldsymbol{\alpha})}$$

of  $\xi_j(\boldsymbol{\alpha})$ ; hence  $P_j(\xi_j(\boldsymbol{\alpha}))$  is also monotone in the coordinates  $\alpha_k$  of  $\boldsymbol{\alpha}$ . Note the similarity of equation (17) to (2) and (3).

It follows immediately from Lemma 2 of Holland and Rosenbaum (1986; see also Kamae, Krengel O'Brien, 1977) that for any non-decreasing summary  $g(\mathbf{X})$  of  $\mathbf{X} = (X_1, \dots, X_J)$ ,  $E[g(\mathbf{X}) | \boldsymbol{\alpha}]$  is non-decreasing in each coordinate  $\alpha_{ik}$  of  $\boldsymbol{\alpha}$ ; this implies the SOM (Stochastic Ordering of the Manifest score  $X_+$  by the latent trait) property of Hemker et al. (1997), that  $P[X_+ > c | \boldsymbol{\alpha}]$  is non-decreasing in each coordinate  $\alpha_{ik}$  of  $\boldsymbol{\alpha}$ . Not much is known about the inverse and more useful property SOL (see 7) when the latent “trait” is multidimensional. We might conjecture for example that

$$P \left[ \alpha_k = 1 \mid \sum_{j:Q_{jk}=1} X_j = s \right] \quad (18)$$

would be non-decreasing in  $s$ . As in conventional dichotomous IRT models we hope that such a property is true, because it means that a most powerful test of mastery of skill  $k$  can be based on a simple total of correctly-performed tasks involving skill  $k$ .

It also follows from Theorem 8 of Holland and Rosenbaum (1986; see also Jogdeo, 1978) that since  $\boldsymbol{\alpha}$  is a collection of independent, or more generally, associated, random variables, the collection of response variables  $\mathbf{X}$  is also associated, that is, for any two non-decreasing summaries  $f(\mathbf{X})$  and  $g(\mathbf{X})$ , that  $\text{Cov}(f(\mathbf{X}), g(\mathbf{X})) \geq 0$ . In general, conditional association appears to be a difficult property to prove when the underlying latent variable,  $\boldsymbol{\alpha}$  in this case, is continuous and multidimensional, but it may be easier to either derive conditions for the binary coordinates  $\alpha_{ik}$  under which conditional association does hold—in which case monotone homogeneity IRT models begin to be competitors for modeling the same response data as the cognitively-motivated model—or at least to obtain the appropriate generalization of conditional association for cognitive assessment models such as this.

Next, an invariant item ordering property, such as  $P_j(\xi_j(\boldsymbol{\alpha})) \leq P_k(\xi_k(\boldsymbol{\alpha}))$  uniformly in  $\boldsymbol{\alpha}$ , appears to be difficult to obtain unless  $1 - s_j$ ,  $g_j$ , and  $\xi_j(\boldsymbol{\alpha})$ , are comonotone as  $j$  varies, for all  $\boldsymbol{\alpha}$ . This is a kind of Guttman scaling condition on the latent responses  $\xi_j$ , that says for example that easier guessing is associated with lower skill requirements, and vice-versa. Thus, our cognitive assessment model (16) provides fertile ground for thinking about invariant item ordering, without necessarily being tied to preconceptions about continuous unidimensional IRT models.

More broadly, Hoijtink and Molenaar (1997) show how nonparametric model features such as conditional association (4) and manifest monotonicity (6) can be directly relevant to assessing model fit in a parametric Bayesian setting. Given a complete and interesting set of nonparametric model features for models like (16), a similar approach to model fit may be taken here.

Finally, among many who work in the nonparametric IRT traditions of Mokken, Molenaar, Sijtsma, Holland, Rosenbaum and their colleagues, it is the source of some bemusement that we work so hard to establish that scores such as  $X_+$ , which are perfectly good in the Rasch model, the most stringent of logistic IRT models, also suffice to order examinees, assess monotonicity properties of the underlying model from observable data, make mastery decisions, etc., in general nonparametric settings. Certainly comparisons between Rasch and Mokken scaling are not new (e.g., Meijer, Sijtsma Smid, 1990), and more recently aspects of both Holland's (1990) “Dutch identity” work and Scheiblechner's (1995; see also Junker, 1998) “ISOP-model” work can be seen partly as attempts to formalize the connection between Rasch and nonparametric IRT methods.

I believe that the connection is actually rather simple, and is nearly obvious from Holland's (1990) work: The Rasch model is a very well-behaved exponential family model with immediately understandable sufficient statistics for items, given person parameters, and immediately understandable sufficient statistics for

persons, given item parameters. Much of the work on monotone homogeneity models and their cousins is directed at understanding just how generally these understandable, but no longer formally sufficient, statistics yield sensible inferences about examinees and items. The model (16) provides us with a different “basic”—albeit not exactly exponential-family—model for cognitive assessment, in which parameters depend on immediately understandable summaries of the data, as illustrated by the complete conditional calculations above. We may ask how complex the relationship between the examinee skill parameters  $\alpha$  and the task performance data  $X$  can get, and still have these summaries be informative about guessing, slips presence or absence of skills, etc. We may also take an excursion into parametric models, as I have done, and ask whether this is the “right” parametric model on which to base a nonparametric theory of cognitive assessment. For example, other possible models we might consider instead of (16) as a starting point for such a nonparametric theory include the constrained latent class model of Haertel (1989) and the hybrid model of Yamamoto and Gitomer (1993).

In addition to the intrinsic interest of this enterprise, the resulting nonparametric theory of cognitive assessment may have some practical utility. The relative ease with which the original Corbett, Anderson and O’Brien (1995) model can be estimated is a consequence of *both* its careful tailoring to the psychological theory, *and* the fact that data could be collected at the skill level (same granularity as the psychological theory), rather than at the task level, by the LISP tutor. But other assessment settings may not permit such tight binding of data collection design and psychological theory; and our discussion of the model (16) shows that even relatively minor modifications along these lines can make the inferential task more difficult. Estimation probably requires E-M, MCMC, or some other computationally intensive method, and great care must be taken so that every parameter is identifiable from the data. Models being proposed in the psychometric literature today to help unify the traditional IRT and discrete-cognitive-attributes approaches, such as the Unified Model of DiBello, Jiang and Stout (in press), also appear to require such treatment. In many cases, the estimation method, while feasible, may well be too slow for use in real-time feedback, as with computer based intelligent tutoring systems, and too complicated for teacher scoring of assessments embedded in instruction. A clear theory of which faster data summaries are relevant to the cognitive inferences we wish to make, over a wide variety of cognitive assessment models, would be an important contribution from the interface between nonparametric and parametric “IRT” models.

## References

- Ackerman, T. A. (1994). Using multidimensional item response theory to understand what items and tests are measuring. *Applied Measurement in Education*, 7, 255–278.
- Adams, R. J., Wilson, M., Wang, W.-C. (1997). The multidimensional random coefficients multinomial logit model. *Applied Psychological Measurement*, 21, 1–23.
- Algina, J. (1992). Special issue: the National Assessment of Educational Progress (Editor’s Note). *Journal of Educational Measurement*, 29, 93–94.
- Anderson, J. R. (ed.) (1993). *Rules of the mind*. Hillsdale, NJ: Lawrence Erlbaum Associates.
- Anderson, J. R., Corbett, A. T., Koedinger, K. R., Pelletier, R. (1995). Cognitive tutors: lessons learned. *The Journal of the Learning Sciences*, 4, 167–207.

- Andrich, D. (1996). A hyperbolic cosine latent trait model for unfolding polytomous responses: Reconciling Thurstone and Likert methodologies. *British Journal of Mathematical and Statistical Psychology*, 49, 347–365.
- Bartolucci, F., Forcina, A. (in press). A likelihood ratio test for MTP<sub>2</sub> within binary variables. In press, *Annals of Statistics*.
- Baxter, G. P., Glaser, R. (1998). Investigating the cognitive complexity of science assessments. *Educational Measurement: Issues and Practice*, 17, 37–45.
- Béguin, A.A., Glas, C.A.W. (1999). *MCMC estimation of multidimensional IRT models*. (Research Report 98-14). Department of Educational Measurement and Data Analysis, University of Twente, the Netherlands. [Online]. Available: <http://to-www.edte.utwente.nl/TO/omd/report98-rr9814.htm>. Accessed 28 April 2000.
- Bloom, B. S. (1984). The 2-sigma problem: the search for methods of group instruction as effective as one-to-one tutoring. *Educational Researcher*, 13, 4–16.
- Bock, R. D., Zimowski, M. F. (1997). Multi-group IRT. In W. J. Van der Linden, R. K. Hambleton (Eds.), *Handbook of modern item response theory* (pp. 433–448). New York: Springer Verlag.
- Bradlow, E. T., Wainer, H., Wang, X. (1999). A Bayesian random effects model for testlets. *Psychometrika*, 64, 153–168.
- Cliff, N., Donoghue, J. R. (1992). Ordinal test fidelity estimated by an item sampling model. *Psychometrika*, 57, 217–236.
- Corbett, A. T., Anderson, J. R., O'Brien, A. T. (1995). Student modeling in the ACT programming tutor. In P. D. Nichols, S. F. Chipman, R. L. Brennan (Eds.), *Cognitively diagnostic assessment* (pp. 19–41). Hillsdale, NJ: Lawrence Erlbaum Associates.
- DiBello, L., Jiang, H., Stout, W. F. (in press). A multidimensional IRT model for practical cognitive diagnosis. In Press, *Applied Psychological Measurement*.
- Douglas, J., Qui, P. (1997). *Generalized linear factor analysis with Markov chain Monte Carlo*. Unpublished manuscript.
- Draney, K. L., Pirolli, P., Wilson, M. (1995). A measurement model for a complex cognitive skill. In P. D. Nichols, S. F. Chipman, R. L. Brennan (Eds.), *Cognitively diagnostic assessment* (pp. 103–125). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Drasgow, F., Levine, M. V., Tsien, S., Williams, B., Mead, A. D. (1995). Fitting polytomous item response theory models to multiple-choice tests. *Applied Psychological Measurement*, 19, 143–165.
- Ellis, J. L. (1994). *Foundations of monotone latent variable models*. Nijmegen: Nijmegen Institute for Cognition and Information.

- Ellis, J. L., Junker, B. W. (1997). Tail-measurability in monotone latent variable models. *Psychometrika*, 62, 495–523.
- Embretson, S. E. (1985). Multicomponent latent trait models for tests design. In S. E. Embretson (Ed.), *Test design: developments in psychology and psychometrics* (pp. 195–218). New York: Academic Press.
- Embretson, S. E. (1991). A multidimensional latent trait model for measuring learning and change. *Psychometrika*, 56, 495–515.
- Embretson, S. E. (1995). Developments toward a cognitive design system for psychological tests. In D. Lubinski, R. V. Dawis (Eds.), *Assessing individual differences in human behavior: new concepts, methods and findings* (pp 17–48). Palo Alto CA: Davies-Black Publishing.
- Embretson, S. E. (1997). Multicomponent response models. In W. J. Van der Linden R. K. Hambleton (Eds.), *Handbook of modern item response theory* (pp. 305–321). New York: Springer Verlag.
- Embretson, S. E. (1999). Generating items during testing: psychometric issues and models. *Psychometrika*, 64, 407–433.
- Engelhard, G., Jr. (1994). Examining rater errors in the assessment of written composition with many-faceted Rasch models. *Journal of Educational Measurement*, 31, 93–112.
- Engelhard, G., Jr. (1996). Evaluating rater accuracy in performance assessment. *Journal of Educational Measurement*, 33, 56–70.
- Fischer, G. H. (1973). Linear logistic test model as an instrument in educational research. *Acta Psychologica*, 37, 359–374.
- Fischer, G. H., Molenaar, I. W. (eds.) (1995). *Rasch models: foundations, recent developments, and applications*. New York: Springer-Verlag.
- Fox, G.J.A., C.A.W. Glas (1998). *A multi-level IRT model with measurement error in the predictor variables*. (Research Report 98-16). Department of Educational Measurement and Data Analysis, University of Twente, the Netherlands. [Online]. Available: <http://to-www.edte.utwente.nl/TO/omd-report98/report98.htm>. Accessed 28 April 2000.
- Fraser, C., McDonald, R. P. (1988). NOHARM: Least squares item factor analysis. *Multivariate Behavioral Research* 23, 267–269.
- Gardner, H. (1992). Assessment in context: the alternative to educational testing. In B. R. Gifford, M. C. O'Connor (Eds.), *Changing assessments: alternative views of aptitude, achievement, and instruction* (pp. 77–119). Norwell, MA: Kluwer Academic Publishers.
- Gelman, A., Carlin, J. B., Stern, H. S., Rubin, D. B. (1995). *Bayesian data analysis*. New York: Chapman and Hall.
- Gibbons, R. D., Hedeker, D. R. (1992). Full-information item bi-factor analysis. *Psychometrika*, 57, 423–436.

- Gibbons, R. D., Hedeker, D. R. (1997). Random effects probit and logistic regression models for three-level data. *Biometrics*, 53, 1527–1537.
- Glas, C. A. W., Verhelst, N. D. (1989). Extensions of the partial credit model. *Psychometrika*, 54, 635–659.
- Grayson, D. A. (1988). Two-group classification in latent trait theory: scores with monotone likelihood ratio. *Psychometrika*, 53, 383–392.
- Haertel, E. H. (1989). Using restricted latent class models to map the skill structure of achievement items. *Journal of Educational Measurement*, 26, 301–321.
- Hemker, B. T., Sijtsma K. (1999, July). *A comparison of three general types of unidimensional IRT models for polytomous items*. Paper presented at a Symposium on Nonparametric Item Response Theory Models in Action, European Meeting of the Psychometric Society in Lüneburg, Germany.
- Hemker, B. T., Sijtsma, K., Molenaar, I. W. (1995). Selection of unidimensional scales from a multidimensional item bank in the polytomous Mokken IRT model. *Applied Psychological Measurement*, 19, 337–352.
- Hemker, B. T., Sijtsma K., Molenaar, I. W., Junker, B. W. (1996). Polytomous IRT models and monotone likelihood ratio of the total score. *Psychometrika*, 61, 679–693.
- Hemker, B. T., Sijtsma K., Molenaar, I. W., Junker, B. W. (1997). Stochastic ordering using the latent trait and the sum score in polytomous IRT models. *Psychometrika*, 62, 331–347.
- Hoijtink, H., Molenaar, I. W. (1997). A multidimensional item response model: Constrained latent class analysis using the Gibbs sampler and posterior predictive checks. *Psychometrika*, 62, 171–189.
- Holland, P. W. (1981). When are item response models consistent with observed data? *Psychometrika*, 46, 79–92.
- Holland, P. W. (1990). The Dutch identity: a new tool for the study of item response models. *Psychometrika*, 55, 5–18.
- Holland, P. W., Rosenbaum, P. R. (1986). Conditional association and unidimensionality in monotone latent trait models. *Annals of Statistics*, 14, 1523–1543.
- Huguenard, B. R., Lerch, F. J., Junker, B. W., Patz, R. J., Kass, R. E. (1997). Working memory failure in phone-based interaction. *ACM Transactions on Computer-Human Interaction*, 4, 67–102.
- Huynh, H. (1994). A new proof for monotone likelihood ratio for the sum of independent Bernoulli random variables. *Psychometrika*, 59, 77–79.
- Jaakkola, T. S., Jordan, M. I. (1999). Bayesian parameter estimation via variational methods. In press, *Statistics and Computing*. [Online]. Available: <http://www.cs.berkeley.edu/~jordan/publications.html>. Accessed 28 April 2000.

- Janssen, R., De Boeck, P. (1997). Psychometric modeling of componentially designed synonym tasks. *Applied Psychological Measurement*, 21, 37–50.
- Jogdeo, K. (1978). On a probability bound of Marshall and Olkin. *Annals of Statistics*, 6, 232–234.
- Johnson, E. G., Mislevy, R. J., Thomas, N. (1994). Theoretical background and philosophy of NAEP scaling procedures. In E. G. Johnson, J. Mazzeo, D. L. Kline (Eds.), *Technical Report of the NAEP 1992 Trial State Assessment Program in Reading* (pp. 133–146). Washington, DC: Office of Educational Research and Improvement, U.S. Department of Education.
- Junker, B. W. (1993). Conditional association, essential independence and monotone unidimensional item response models. *Annals of Statistics*, 21, 1359–1378.
- Junker, B. W. (1998). Some remarks on Scheiblechner's treatment of ISOP models. *Psychometrika*, 63, 73–85.
- Junker, B. W., Ellis, J. L. (1997). A characterization of monotone unidimensional latent variable models. *Annals of Statistics*, 25, 1327–1343.
- Junker, B. W., Koedinger, K. R., Trottini, M. (2000). *Finding improvements in student models for intelligent tutoring systems via variable selection for a linear logistic test model*. Paper presented at the Annual North American Meeting of the Psychometric Society, July 2000, Vancouver, BC, Canada.
- Junker, B. W., Patz, R. J. (June 1998). *The hierarchical rater model for rated test items*. Paper presented at the Annual North American Meeting of the Psychometric Society, Champaign-Urbana IL, USA.
- Junker, B. W., Sijtsma, K. (2000). Latent and manifest monotonicity in item response models. *Applied Psychological Measurement*, 24, 65–81.
- Kass, R. E., Tierney, L., Kadane, J. B. (1990). The validity of posterior expansions based on Laplace's method. In S. Geisser, J. S. Hodges, S. J. Press A. Zellner (Eds.), *Bayesian and likelihood methods in statistics and econometrics: Essays in honor of George A. Barnard* (pp. 473–488). New York: North-Holland.
- Kelderman, H., Rijkes, C. P. M. (1994). Loglinear multidimensional IRT models for polytomously scored items. *Psychometrika*, 59, 149–176.
- Kamae, T. Krengel, U., O'Brien, G. L. (1977). Stochastic inequalities on partially ordered spaces. *Annals of Probability*, 5, 899–912.
- Lee, Y., Nelder, J. A. (1996). Hierarchical generalized linear models. *Journal of the Royal Statistical Society, Series B*, 58, 619–678.
- Linacre, J. M. (1989). *Many-faceted Rasch Measurement*. Chicago: Mesa Press.
- Liu, C., Rubin, D. B. (1998). Maximum likelihood estimation of factor analysis using the ECME algorithm with complete and incomplete data. *Statistica Sinica*, 8, 729–747.

- Lord, F. M. (1952). A theory of test scores. *Psychometric Monographs, No. 7*.
- Loevinger, J. (1948). The technique of homogeneous tests compared with some aspects of “scale analysis” and factor analysis. *Psychological Bulletin, 45*, 507–530.
- Maris, E. (1995). Psychometric latent response models. *Psychometrika, 60*, 523–547.
- Maris, E., De Boeck, P., Van Mechelen, I. (1996). Probability matrix decomposition models. *Psychometrika, 61*, 7–29.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika, 47*, 149–174.
- Meijer, R. R., Sijtsma, K., Smid, N. G. (1990). Theoretical and empirical comparison of the Mokken and the Rasch approach to IRT. *Applied Psychological Measurement, 14*, 283–298.
- Meijer, R. R. (1996). Person-fit research: An introduction. (Guest editor’s introduction to the Special Issue: Person-fit research: Theory and applications.) *Applied Measurement in Education, 9*, 3–8.
- Meredith, W. (1965). Some results based on a general stochastic model for mental tests. *Psychometrika, 30*, 419–440.
- Mislevy, R. J. (1985). Estimation of latent group effects. *Journal of the American Statistical Association, 80*, 993–997.
- Mislevy, R. J. (1996). Test theory reconceived. *Journal of Educational Measurement, 33*, 379–416.
- Mislevy, R. J., Sheehan, K. M. (1989). The role of collateral information about examinees in item parameter estimation *Psychometrika, 54*, 661–679.
- Mokken, R. J. (1971). *A theory and procedure of scale analysis*. The Hague: Mouton.
- Mokken, R. J. (1997). Nonparametric models for dichotomous items. In W. J. Van der Linden, R. K. Hambleton (Eds.), *Handbook of modern item response theory* (pp. 351–368). New York: Springer Verlag.
- Molenaar, I. W. (1991). A weighted Loevinger H-coefficient extending Mokken scaling to multicategory items. *Kwantitatieve Methoden, 37*, 97–117.
- Molenaar, I. W. (1997). Nonparametric methods for polytomous responses. In W. J. Van der Linden, R. K. Hambleton (Eds.), *Handbook of modern psychometrics* (pp. 369–380). New York: Springer Verlag.
- Molenaar, I. W., Sijtsma, K. (1999). *MSP for Windows*. Groningen: iecProGAMMA.
- Molenaar, I. W., Stout, W. F. (January 5, 2000). Personal communication.
- Muraki, E., Carlson, J. E. (1995). Full-information Factor Analysis for Polytomous Item Responses. *Applied Psychological Measurement, 19*, 73–90.
- McCullagh, P., Nelder, J. A. (1989). *Generalized Linear Models* (2nd Edition). New York: Chapman and Hall.

- Nichols, P., Sugrue, B. (1999). The lack of fidelity between cognitively complex constructs and conventional test development practice. *Educational Measurement: Issues and Practice*, 18, 18–29.
- Patz, R. J., Junker, B. W. (1999a). A straightforward approach to Markov Chain Monte Carlo methods for item response models. *Journal of Educational and Behavioral Statistics*, 24, 146–178.
- Patz, R. J., Junker, B. W. (1999b). Applications and extensions of MCMC in IRT: Multiple item types, missing data, and rated responses. *Journal of Educational and Behavioral Statistics*, 24, xxx–xxx.
- Patz, R. J., Junker, B. W., Johnson, M. S. (1999, April). *The hierarchical rater model for rated test items and its application to large-scale educational assessment data*. Paper presented at a Symposium on Making the Most of Constructed Responses: The Information in Multiple Ratings, Annual Meeting of the American Educational Research Association, Montreal Canada.
- Patz, R. J., Junker, B. W., Lerch, F. J., Huguenard, B. R. (1996). *Analyzing small psychological experiments with item response models* (CMU Statistics Department technical report #644). [Online]. Available: <http://www.stat.cmu.edu/cmu-stats/tr/tr644/tr644.html>. Accessed 28 April 2000.
- Post, W. J. (1992). *Nonparametric unfolding models. A latent structure approach*. Leiden: DSWO Press, Leiden University, The Netherlands.
- Post, W. J., Snijders, T. A. B. (1993). Nonparametric unfolding models for dichotomous data. *Methodika*, 7, 130–156.
- Ramsay, J. O. (1991). Kernel smoothing approaches to nonparametric item characteristic curve estimation. *Psychometrika*, 56, 611–630.
- Ramsay, J. O. (1995). A similarity-based smoothing approach to nondimensional item analysis. *Psychometrika*, 60, 323–339.
- Ramsay, J. O. (1996). A geometrical approach to item response theory. *Behaviormetrika*, 23, 3–17.
- Reckase, M. D. (1985). The difficulty of test items that measure more than one ability. *Applied Psychological Measurement*, 9, 401–412.
- Resnick, L. B., Resnick, D. P. (1992). Assessing the thinking curriculum: new tools for educational reform. In B. R. Gifford, M. C. O'Connor (Eds.), *Changing assessments: alternative views of aptitude, achievement, and instruction* (pp 37–75). Norwell, MA: Kluwer Academic Publishers.
- Rigdon, S. E., Tsutakawa, R. K. (1983). Parameter estimation in latent trait models. *Psychometrika*, 48, 567–574.
- Robertson, T., Wright, F. T., Dykstra, R. L. (1988) *Order restricted statistical inference*. New York: Wiley.
- Rosenbaum, P. R. (1984). Testing the conditional independence and monotonicity assumptions of item response theory. *Psychometrika*, 49, 425–435.

- Rosenbaum, P. R. (1987a). Probability inequalities for latent scales. *British Journal of Mathematical and Statistical Psychology*, 40, 157–168.
- Rosenbaum, P. R. (1987b). Comparing item characteristic curves. *Psychometrika*, 52, 217–233.
- Roussos, L. (1994). Summary and review of cognitive diagnosis models. Unpublished manuscript.
- Samejima, F. (1997). Departure from normal assumptions: a promise for future psychometrics with substantive mathematical modeling. *Psychometrika*, 62, 471–493.
- Scheiblechner, H. (1972). Das Lernen und Lösen komplexer Denkaufgaben. [The learning and solving of complex reasoning items.] *Zeitschrift für Experimentelle und Angewandte Psychologie*, 3, 456–506.
- Scheiblechner, H. (1995). Isotonic ordinal probabilistic models (ISOP). *Psychometrika*, 60, 281–304.
- Seltman, H. (1999). *Hidden Stochastic Models for Biological Rhythm Data*. Unpublished Ph.D. Dissertation, Department of Statistics, Carnegie Mellon University, Pittsburgh PA, USA.
- Sijtsma, K. (1998). Methodology review: Nonparametric IRT approaches to the analysis of dichotomous item scores. *Applied Psychological Measurement*, 22, 3–32.
- Sijtsma, K., Van der Ark, L. A. (this volume). Progress in IRT analysis of polytomous item scores: dilemmas and practical solutions. In A. Boomsma, T. Snijders, M. Van Duijn (Eds.), *Essays in Item Response Modeling* (pp. xxx–xxx). New York: Springer-Verlag.
- Sijtsma, K., Hemker, B. T. (1998). Nonparametric polytomous IRT models for invariant item ordering, with results for parametric models. *Psychometrika*, 63, 183–200.
- Sijtsma, K., Junker, B. W. (1996). A survey of theory and methods of invariant item ordering. *British Journal of Mathematical and Statistical Psychology*, 49, 79–105.
- Sijtsma, K., Junker, B. W. (1997). Invariant item ordering of transitive reasoning tasks. In J. Rost, R. Langeheine (Eds.), *Applications of latent trait and latent class models in the social sciences* (pp. 97–107). Münster: Waxmann Verlag.
- Shute, V. J., Psotka, J. (1996). Intelligent tutoring systems: Past, Present and Future. In D. Jonassen (Ed.), *Handbook of Research on Educational Communications and Technology* (pp. 570–600). New York: Macmillan Press..
- Snijders, T.A.B. (this volume). Two-level non-parametric scaling for dichotomous data. In A. Boomsma, T. Snijders, M. Van Duijn (Eds.), *Essays in Item Response Modeling* (pp. xxx–xxx). New York: Springer-Verlag.
- Snijders, T.A.B. (1997). Estimation and prediction for stochastic blockmodels for graphs with latent block structure. *Journal of Classification*, 14, 75–100.
- Stegelmann, W. (1983). Expanding the Rasch model to a general model having more than one dimension. *Psychometrika*, 48, 259–267.

- Stout, W. F. (1987). A nonparametric approach for assessing latent trait unidimensionality. *Psychometrika*, 52, 589–617.
- Stout, W. F. (1990). A new item response theory modeling approach with applications to unidimensionality assessment and ability estimation. *Psychometrika*, 55, 293–325.
- Stout, W. F., Habing, B., Douglas, J., Kim, H. R., Roussos, L., Zhang, J. (1996). Conditional covariance-based nonparametric multidimensionality assessment. *Applied Psychological Measurement*, 20, 331–354.
- Sykes, R. C., Heidorn, M. (1999, April). *The assignment of raters to items: controlling for rater bias*. Paper presented at the annual meeting of the National Council on Measurement in Education, Montreal, Canada.
- Tanner, M. A. (1996). *Tools for statistical inference: methods for the exploration of posterior distributions and likelihood functions*. 3<sup>rd</sup> Edition. New York: Springer-Verlag.
- Tatsuoka, K. K. (1990). Toward an integration of item response theory and cognitive error diagnosis. In N. Fredriksen, R. Glaser, A. Lesgold, M. G. Shafto (Eds.), *Diagnostic monitoring of skill and knowledge acquisition* (pp. 453–488). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Tatsuoka, K. K. (1995). Architecture of knowledge structures and cognitive diagnosis: a statistical pattern recognition and classification approach. In P. D. Nichols, S. F. Chipman, R. L. Brennan (Eds.), *Cognitively diagnostic assessment* (pp. 327–359). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Ter Hofstede, F., Steenkamp, J.-B. E. M., Wedel, M. (1999). Identifying spatially contiguous international target markets. Manuscript submitted for publication.
- Thissen, D., Steinberg, L. (1986). A taxonomy of item response models. *Psychometrika*, 51, 567–577.
- Van der Ark, L. A. (1999, July). *A reference card for the relations between IRT models for polytomous items and some relevant properties*. Paper presented at a Symposium on Nonparametric Item Response Theory Models in Action, European Meeting of the Psychometric Society in Lüneburg, Germany.
- Van der Linden, W. J., Hambleton, R. K. (eds.) (1997). *Handbook of modern item response theory*. New York: Springer Verlag.
- VanLehn, K., Niu, Z. (in press). Bayesian student modeling, user interfaces and feedback: A sensitivity analysis. In Press, *International Journal of Artificial Intelligence in Education*.
- VanLehn, K., Niu, Z., Siler, S., Gertner, A. (1998). Student modeling from conventional test data: a Bayesian approach without priors. In B. P. Goetl, H. M. Halff, C. L. Redfield, V. J. Shute (Eds.), *Proceedings of the Intelligent Tutoring Systems Fourth International Conference, ITS 98* (pp. 434–443). Berlin: Springer-Verlag.
- Verhelst, N. D., Verstralen, H. H. F. M. (1993). A stochastic unfolding model derived from the partial credit model. *Kwantitatieve Methoden*, 42, 73–92.

- Verhelst, N.D., Verstralen, H.H.F.M. (this volume). IRT models for multiple raters. In A. Boomsma, T. Snijders, and M. Van Duijn (Eds.), *Essays in Item Response Modeling* (pp. xxx–xxx). New York: Springer-Verlag.
- Wilson, M. R., Hoskens, M. (1999, April). *The rater bundle model*. Paper presented at a Symposium on Making the Most of Constructed Responses: The Information in Multiple Ratings, Annual Meeting of the American Educational Research Association, Montreal Canada.
- Wilson, D., Wood, R. L., Gibbons, R. (1983). TESTFACT: test scoring and item factor analysis. [Computer program]. Chicago: Scientific Software Inc. Online description available: <http://ssicentral.com/irt/testfact.htm>. Accessed 28 April 2000.
- Wu, M. L., Adams, R. J., Wilson, M. R. (1997). *ConQuest: Generalized item response modeling software*. ACER.
- Yamamoto, K., Gitomer, D. H. (1993). Application of a HYBRID model to a test of cognitive skill representation. In N. Fredriksen, R. J. Mislevy (Eds.), *Test theory for a new generation of tests* (pp. 275–295). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Yuan, A., Clarke, B. (1999). Manifest characterization and testing for two latent traits. Manuscript submitted for publication.
- Zhang, J., Stout, W. F. (1996, April). A theoretical index of dimensionality and its estimation. Paper presented at the Annual Meeting of the National Council on Measurement in Education, New York.
- Zimowski, M. F., Muraki, E., Mislevy, R. J., Bock, R. D. (1997). BILOG-MG. [Computer program]. Chicago: Scientific Software Inc. Online description available: <http://ssicentral.com/irt/bilogmg.htm>. Accessed 28 April 2000.
- Zwick, R. (1992). Special issue on the National Assessment of Educational Progress. *Journal of Educational Measurement*, 17, 93–94.