

Title:

Summed Score Likelihood Based Indices for Testing Latent Variable Distribution Fit in Item Response Theory

Authors:

Zhen Li

Li Cai

Journal publication date:

2017

IES grant information:

Grant number R305D140046

Funded by National Center for Education Research (NCER)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20

Summed Score Likelihood Based Indices for Testing Latent Variable

Distribution Fit in Item Response Theory

Zhen Li

eMetric

Li Cai

University of California, Los Angeles

May 30, 2017

Part of this research was supported by the Institute of Education Sciences (R305D140046). The views expressed here belong to the authors and do not reflect the views or policies of the funding agencies.

Address all correspondence to: Li Cai, CRESST, 300 Charles E. Young Dr. North, GSEIS building, UCLA, Los Angeles, CA, USA 90095-1522. Email: lcai@ucla.edu. Phone: 310.794.7136. Fax: 310.206.5830.

1 Summed Score Likelihood Based Indices for Testing Latent Variable

2 Distribution Fit in Item Response Theory

3
4 Abstract

5 In standard item response theory (IRT) applications, the latent variable is typically assumed to
6 be normally distributed. If the normality assumption is violated, the item parameter estimates can
7 become biased. Summed score likelihood based statistics may be useful for testing latent variable
8 distribution fit. We develop Satorra-Bentler type (Satorra & Bentler, 1994) moment adjustments to
9 approximate the test statistics' tail-area probability. A simulation study was conducted to examine the
10 calibration and power of the unadjusted and adjusted statistics in various simulation conditions.
11 Results show that the proposed indices have tail-area probabilities that can be closely approximated by
12 central chi-squared random variables under the null hypothesis. Furthermore, the test statistics are
13 focused. They are powerful for detecting latent variable distributional assumption violations, and not
14 sensitive (correctly) to other forms of model misspecification such as multidimensionality. As a
15 comparison, the goodness-of-fit statistic M_2 (Maydeu-Olivares & Joe, 2005) has considerably lower
16 power against latent variable non-normality than the proposed indices. Empirical data from a
17 patient-reported health outcomes study is used as illustration.

Introduction

Item response theory (IRT) provides powerful methods supporting educational and psychological measurement (Thissen & Steinberg, 2009). The latent variable in IRT models is usually assumed to follow a normal distribution for the purpose of item parameter estimation (Bock & Lieberman, 1970; Bock & Aitkin, 1981). However, this assumption might be violated in some situations (Woods, 2006; Woods & Lin, 2009). Woods (2006) described several potential situations where θ may be nonnormal. For example, as severe symptoms of psychological disorders rarely exist in the general population and most people have low levels of psychopathological symptoms, latent variables reflecting these symptoms may be positively skewed. Another possible cause arises in the situation when the population is heterogeneous. For instance, when two or more subpopulations with different means and variances are grouped together, potentially multimodal population distributions may be the result. Calibrating the items with respect to the combined population renders the normality assumption suspect. When the assumption of normal latent variable distribution is violated, the item parameter estimates might be biased, leading to bias in subsequent inferences based on these item parameter estimates. Take Computer Adaptive Testing (CAT) as an example, the item parameter estimates are utilized for both item selection and test scoring. Thus, bias in the estimation of item parameters might result in significant bias in the reported test scores.

Although alternative approaches exist for estimating the latent variable distribution in standard IRT models (Bock & Aitkin, 1981; Woods & Thissen, 2006; Woods & Lin, 2009), these approaches are computationally more demanding and specialized software is necessary. **For example, in our**

1 experience, the empirical histogram representation of the latent prior distribution is often less stable
 2 numerically than the standard normal prior. Thus, it is worthwhile to test the assumption of latent
 3 variable normality before more “expensive” approaches are applied. Summed score likelihood based
 4 statistics may be useful for testing latent variable distribution fit. One problem is that the statistics do
 5 not asymptotically follow a chi-squared distribution. We propose a Satorra-Bentler type moment
 6 adjustment method (Satorra & Bentler, 1994) in this paper. The statistics’ tail-area probability can be
 7 approximated by making use of the item parameter error covariance matrix and a Jacobian. The
 8 properties of the adjusted and unadjusted statistics are examined by simulation and empirical studies.
 9 Additionally, a modified Lord-Wingsky algorithm for computing the Jacobian matrix is presented in
 10 the Appendix.

11 Item Response Theory Models

12 In standard IRT models, the conditional item response probabilities (also referred to as item
 13 tracelines or item characteristic curves) are represented as a function of latent variable θ and item
 14 parameters. For example, the 3-parameter logistic (3PL) model can be written as:

$$T_i(1|\theta) = g_i + \frac{1 - g_i}{1 + \exp[-(c_i + a_i\theta)]} \quad (1)$$

15 where $T_i(1|\theta)$ represents item i 's traceline for the 1 category (indicating correct/endorsement response
 16 in most contexts) as a function of θ . The item parameters include: g_i , which is the pseudo-guessing
 17 probability for the item (the lower asymptote parameter); a_i , which is the slope (the discrimination
 18 parameter), and c_i , which is the item intercept parameter. The classical difficulty (threshold) parameter
 19 is obtained as $-c_i/a_i$. If g_i is zero, the model reduces to a 2-parameter logistic (2PL) model, and if all

1 the item slopes are constrained to be equal to a common slope ($a_i \equiv a$), the 1-parameter logistic (1PL)
 2 model is the result. The incorrect/non-endorsement response probability is equal to $T_i(0|\theta) = 1 -$
 3 $T_i(1|\theta)$.

4 For an item with K_i ordered polytomous responses, the graded response model is often
 5 utilized. Let the response categories be coded as $k = 0, \dots, K_i - 1$. The cumulative response probability
 6 for item i in categories k and above is

$$T_i^+(k|\theta) = \frac{1}{1 + \exp[-(c_{ik} + a_i\theta)]} \quad (2)$$

7 for $k = 1, \dots, K_i - 1$. Having defined the boundary cases $T_i^+(0|\theta) = 1$ and $T_i^+(K_i|\theta) = 0$, the category
 8 response probabilities can be written as

$$T_i(k|\theta) = T_i^+(k|\theta) - T_i^+(k + 1|\theta), \quad (3)$$

9 for $k = 0, \dots, K_i - 1$. Let U_i be a random variable whose realization u_i is a response to item i .

10 Regardless of the number of categories, the probability mass function of U_i , conditional on θ , is that of
 11 a multinomial with trial size 1:

$$P(U_i = u_i|\theta) = \prod_{k=0}^{K_i-1} [T_i(u_i|\theta)]^{1_k(u_i)}, \quad (4)$$

12 where $1_k(u_i)$ is an indicator function such that

$$1_k(u_i) = \begin{cases} 1, & \text{if } k = u_i \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

13 The Latent Variable Distribution in IRT

14 Estimating the latent variable distribution along with item parameters using the empirical
 15 histogram (Bock & Aitkin, 1981; Mislevy, 1984; Zimowski, Muraki, Mislevy, & Bock, 1996) is an

1 established strategy for detecting and correcting latent variable nonnormality in IRT. Newer
2 semi-parametric density estimation procedures offer more efficient alternatives. These include the
3 Ramsay-Curve IRT (Woods & Thissen, 2006), and Davidian Curve IRT (Woods & Lin, 2009; Monroe &
4 Cai, 2014), as well as its multidimensional extension (Monroe, 2014). In practice, however, estimating
5 latent variable densities often requires specialized software. More complex latent variable distributions
6 also involve more parameters to be estimated from the data, increasing the need for larger calibration
7 sample sizes to achieve stable estimation. Finally, even as nonnormal latent densities may be modeled,
8 e.g., using a Ramsay curve IRT model (Woods & Thissen, 2006), and the relative model fit may be
9 evaluated against a baseline using likelihood ratio tests, it does not circumvent the need for absolute
10 goodness of fit indices to establish the adequacy of the least restrictive model in the class of models
11 being compared (see Maydeu-Olivares & Cai, 2006 for further explanation). It would be highly
12 desirable to establish a set of statistics that can be used to diagnose the extent to which a normal (or
13 non-normal) latent variable distribution may in fact be a reasonable characterization before more
14 “expensive” methods and software programs for semi-parametric density estimation are employed.

15 In developing such a group of test statistics for latent variable distribution fit, several desiderata
16 should be taken into account. First, the statistics should be easily computable, preferably using only
17 standard byproducts of the item calibration process. Second, the statistics should have well-grounded
18 heuristic motivation and theoretical justification. Third, the frequency calibration of the statistics under
19 the null hypothesis should be sufficiently accurate. Finally, the statistics should have adequate power
20 that is *focused* on latent variable distribution assumption violation and sufficient diagnostic specificity,

1 rather than becoming a surrogate of overall model fit tests.

2 The guiding insight has been provided elsewhere in the literature. For unidimensional IRT
 3 modeling, the observed and model-implied summed score distribution can be a basis for inferring the
 4 adequacy of the latent variable distribution specification in the IRT model (Thissen & Wainer, 2001; p.
 5 130). After model fitting, residual summed score probabilities may be used to construct chi-squared test
 6 statistics. While the idea itself is not new (see Ferrando & Lorenzo-seva, 2001; Hambleton & Traub,
 7 1973; Lord, 1953; Ross, 1966; Sinharay, Johnson, & Stern, 2006, among others), we utilize the recently
 8 developed theory of limited-information goodness-of-fit testing to formally demonstrate that the
 9 summed score likelihood based fit index proposed here belongs to the general family of multinomial
 10 limited-information tests.

11 The Multinomial Sampling Model and Maximum Likelihood Estimation

12 Let there be I items in a test. Under the conditional independence assumption, the IRT model
 13 specifies the conditional response pattern probability as the following product

$$P\left(\bigcap_{i=1}^I U_i = u_i \mid \theta\right) = \prod_{i=1}^I P(U_i = u_i \mid \theta). \quad (6)$$

14 Assuming that $g(\theta)$ is the distribution of the latent variable (also known as the prior distribution), the
 15 marginal response pattern probability is the following integral:

$$P\left(\bigcap_{i=1}^I U_i = u_i\right) = \int \prod_{i=1}^I P(U_i = u_i \mid \theta) g(\theta) d\theta = \pi_{\mathbf{u}}(\boldsymbol{\gamma}), \quad (7)$$

16 where $\mathbf{u} = (u_1, \dots, u_I)'$ is the response pattern, and $\boldsymbol{\gamma}$ is a $d \times 1$ vector that collects together the free
 17 item parameters from all I items. The parenthetical notation $\pi_{\mathbf{u}}(\boldsymbol{\gamma})$ in Equation (7) is used to

1 emphasize the fact that it *is* the model. The marginal response probability depends on the item
 2 parameters, the item-level response models, and the assumed latent variable distribution.

3 Recall that K_i is the number of categories for item i . For I items, the IRT model generates a
 4 total of $C = \prod_{i=1}^I K_i$ cross-classifications or possible item response patterns in the form of a
 5 contingency table. Based on a sample of N respondents, let the observed proportion associated with
 6 pattern \mathbf{u} be denoted as $p_{\mathbf{u}}$. The sampling model for this contingency table is a multinomial
 7 distribution with C cells and N trials. The multinomial log-likelihood for the item parameters $\boldsymbol{\gamma}$ is
 8 proportional to

$$\log L(\boldsymbol{\gamma}) \propto N \sum_{\mathbf{u}} p_{\mathbf{u}} \log \pi_{\mathbf{u}}(\boldsymbol{\gamma}), \quad (8)$$

9 where the summation is over all C response patterns. Maximization of the log-likelihood (e.g., with
 10 the EM algorithm; Bock & Aitkin, 1981) leads to the maximum marginal likelihood estimator $\hat{\boldsymbol{\gamma}}$.

11 Upon finding $\hat{\boldsymbol{\gamma}}$, the IRT model generates model-implied probabilities for each response
 12 pattern $\pi_{\mathbf{u}}(\hat{\boldsymbol{\gamma}}) = \hat{\pi}_{\mathbf{u}}$. Suppose the model-implied response pattern probabilities $\hat{\pi}_{\mathbf{u}}$ are collected into a
 13 $C \times 1$ vector $\hat{\boldsymbol{\pi}}$ of all model-implied response pattern probabilities. By analogy, let a $C \times 1$ vector $\boldsymbol{\pi}$
 14 contain the true (population) response pattern probabilities. Similarly, the observed proportions $p_{\mathbf{u}}$
 15 can be collected into a $C \times 1$ vector \mathbf{p} . For example, for 3 dichotomously scored items there are $2^3 = 8$
 16 item response patterns, and the response pattern probabilities and observed proportions are:

$$\boldsymbol{\pi} = \begin{pmatrix} \pi_{000} \\ \pi_{001} \\ \pi_{010} \\ \pi_{011} \\ \pi_{100} \\ \pi_{101} \\ \pi_{110} \\ \pi_{111} \end{pmatrix}, \quad \hat{\boldsymbol{\pi}} = \begin{pmatrix} \hat{\pi}_{000} \\ \hat{\pi}_{001} \\ \hat{\pi}_{010} \\ \hat{\pi}_{011} \\ \hat{\pi}_{100} \\ \hat{\pi}_{101} \\ \hat{\pi}_{110} \\ \hat{\pi}_{111} \end{pmatrix} = \begin{pmatrix} \pi_{000}(\hat{\boldsymbol{\gamma}}) \\ \pi_{001}(\hat{\boldsymbol{\gamma}}) \\ \pi_{010}(\hat{\boldsymbol{\gamma}}) \\ \pi_{011}(\hat{\boldsymbol{\gamma}}) \\ \pi_{100}(\hat{\boldsymbol{\gamma}}) \\ \pi_{101}(\hat{\boldsymbol{\gamma}}) \\ \pi_{110}(\hat{\boldsymbol{\gamma}}) \\ \pi_{111}(\hat{\boldsymbol{\gamma}}) \end{pmatrix}, \quad \boldsymbol{p} = \begin{pmatrix} p_{000} \\ p_{001} \\ p_{010} \\ p_{011} \\ p_{100} \\ p_{101} \\ p_{110} \\ p_{111} \end{pmatrix}. \quad (9)$$

1 From results in discrete multivariate analysis (e.g., Bishop, Fienberg, & Holland, 1975), $\hat{\boldsymbol{\gamma}}$ is
 2 consistent, asymptotically normal, and asymptotically efficient, which can be summarized as follows:

$$\sqrt{N}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}) \xrightarrow{D} \mathcal{N}_d(\mathbf{0}, \mathcal{F}^{-1}), \quad (10)$$

3 where $\mathcal{F} = \boldsymbol{\Delta}'[\text{diag}(\boldsymbol{\pi})]^{-1}\boldsymbol{\Delta}$ is the $d \times d$ Fisher information matrix, with the Jacobian matrix $\boldsymbol{\Delta}$
 4 defined as the $C \times d$ matrix of all first-order partial derivatives of the response patterns probabilities
 5 with respect to the item parameters:

$$\boldsymbol{\Delta} = \frac{\partial \boldsymbol{\pi}(\boldsymbol{\gamma})}{\partial \boldsymbol{\gamma}'}. \quad (11)$$

6

7

Distribution of Residuals under Maximum Likelihood Estimation

8

Based on Equation (10), it can be shown that the asymptotic distribution of the difference ($\boldsymbol{p} -$

9

$\boldsymbol{\pi}$) is C-variate normal:

$$\sqrt{N}(\boldsymbol{p} - \boldsymbol{\pi}) \xrightarrow{D} \mathcal{N}_C(\mathbf{0}, \boldsymbol{\Xi}), \quad (12)$$

10

where $\boldsymbol{\Xi} = \text{diag}(\boldsymbol{\pi}) - \boldsymbol{\pi}\boldsymbol{\pi}'$ is the covariance matrix associated with the multinomial. The residual

11

vector ($\boldsymbol{p} - \hat{\boldsymbol{\pi}}$) is asymptotically C-variate normal under maximum likelihood estimation:

$$\sqrt{N}(\boldsymbol{p} - \hat{\boldsymbol{\pi}}) \xrightarrow{D} \mathcal{N}_C(\mathbf{0}, \boldsymbol{\Gamma}), \quad (13)$$

12

where $\boldsymbol{\Gamma} = \boldsymbol{\Xi} - \boldsymbol{\Delta}\mathcal{F}^{-1}\boldsymbol{\Delta}'$, and the second term reflects variability due to estimation of item parameters.

1 *Lower-order Marginal Probabilities*

2 The IRT model implies marginal probabilities. Consider the 3-item example from above. There
 3 are 3 first order marginal probabilities $\hat{\pi}_i$ ($i = 1, \dots, 3$), one per item. There are also 3 second order
 4 marginal probabilities $\hat{\pi}_{ij}$ for the unique item pairs ($1 \leq j < i \leq 3$). In general, these probabilities
 5 correspond to the I univariate and $I(I - 1)/2$ bivariate margins that can be obtained from the full
 6 C -dimensional contingency table using a reduction operator matrix (see e.g., Maydeu-Olivares & Joe,
 7 2005). An example is given below:

$$\hat{\boldsymbol{\pi}}_2 = \begin{pmatrix} \hat{\pi}_1 \\ \hat{\pi}_2 \\ \hat{\pi}_3 \\ \hat{\pi}_{21} \\ \hat{\pi}_{31} \\ \hat{\pi}_{32} \end{pmatrix} = \mathbf{L}\hat{\boldsymbol{\pi}} = \begin{pmatrix} 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \hat{\pi}_{000} \\ \hat{\pi}_{001} \\ \hat{\pi}_{010} \\ \hat{\pi}_{011} \\ \hat{\pi}_{100} \\ \hat{\pi}_{101} \\ \hat{\pi}_{110} \\ \hat{\pi}_{111} \end{pmatrix}, \quad (14)$$

8 where \mathbf{L} is a fixed operator matrix of 0s and 1s that reduces the response pattern probabilities and
 9 proportions into marginal probabilities and proportions up to order 2. $\hat{\boldsymbol{\pi}}_2$ is the vector of first and
 10 second order marginal probabilities. Correspondingly $\mathbf{p}_2 = \mathbf{L}\mathbf{p}$ is the vector of first and second order
 11 observed marginal proportions.

12 More general versions of the reduction operator matrices for multiple categorical IRT models
 13 can be derived using similar logic (see e.g., Maydeu-Olivares & Joe, 2006; Cai & Hansen, 2013). Note
 14 that \mathbf{L} has full row rank. It implies that the marginal residual vector $(\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2) = \mathbf{L}(\mathbf{p} - \hat{\boldsymbol{\pi}})$ is a full
 15 rank linear transformation of the multinomial residual vector $(\mathbf{p} - \hat{\boldsymbol{\pi}})$. Therefore, the marginal residual
 16 vector $(\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2)$ is asymptotically normal:

$$\sqrt{N}(\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2) = \sqrt{N}\mathbf{L}(\mathbf{p} - \hat{\boldsymbol{\pi}}) \xrightarrow{D} \mathcal{N}_Q(\mathbf{0}, \boldsymbol{\Gamma}_2), \quad (15)$$

1 and $\boldsymbol{\Gamma}_2 = \mathbf{L}\boldsymbol{\Gamma}\mathbf{L}' = \mathbf{L}\boldsymbol{\Xi}\mathbf{L}' - \mathbf{L}\boldsymbol{\Delta}\mathcal{F}^{-1}\boldsymbol{\Delta}'\mathbf{L}' = \boldsymbol{\Xi}_2 - \boldsymbol{\Delta}_2\mathcal{F}^{-1}\boldsymbol{\Delta}'_2$, where $\boldsymbol{\Xi}_2 = \mathbf{L}\boldsymbol{\Xi}\mathbf{L}'$, and $\boldsymbol{\Delta}_2 = \mathbf{L}\boldsymbol{\Delta}$ is the Jacobian for
 2 the marginal probabilities. The dimensionality Q of the normal random variable is equal to the
 3 number of first and second order marginal residuals. For example, in the case of dichotomous items,
 4 the number is $Q = I + I(I - 1)/2 = I(I + 1)/2$.

5 *Summed Score Probabilities*

6 In addition to the response pattern and marginal probabilities, the IRT model also generates
 7 model-implied summed score probabilities. For a test with I items and $k = 0, \dots, K_i - 1$ coded
 8 categories for item i , there are a total of $S = 1 + \sum_{i=1}^I (K_i - 1)$ summed scores ranging from 0 to $S - 1$.
 9 Suppose the observed summed probabilities based on a sample of size N are equal to \bar{p}_s for $s =$
 10 $0, \dots, S - 1$. Under maximum likelihood estimation of item parameters, the corresponding IRT
 11 model-implied summed score probabilities are formally defined as

$$\bar{\pi}_s = \sum_{\mathbf{u}} 1_s(\|\mathbf{u}\|) \hat{\boldsymbol{\pi}}_{\mathbf{u}}, \quad (16)$$

12 where $\|\mathbf{u}\| = \sum_{i=1}^I u_i$ is a notational shorthand for the summed score associated with response pattern
 13 \mathbf{u} , and the indicator function takes a value of 1 if and only if $s = \|\mathbf{u}\|$:

$$1_s(\|\mathbf{u}\|) = \begin{cases} 1, & \text{if } s = \|\mathbf{u}\| \\ 0, & \text{otherwise} \end{cases}. \quad (17)$$

14 Equation (16) shows that the IRT model-implied probability for summed score s is a sum over all such
 15 response pattern probabilities leading to summed score s , in other words, it may also be obtained by a
 16 reduction operator matrix.

17 Let \mathbf{S} be a matrix of fixed 0s and 1s such that the pre-multiplication of $\boldsymbol{\pi}$ by \mathbf{S} yields the

1 summed score probabilities. Each row of \mathbf{S} can be understood as a set of binary logical relations. An
 2 element in row j of \mathbf{S} is equal to 1 if and only if the corresponding response pattern in $\boldsymbol{\pi}$ leads to
 3 summed score $j - 1$. In general, for I items, there are S rows and C columns in \mathbf{S} . In particular, \mathbf{S}
 4 has full row rank and the rows of \mathbf{S} are mutually orthogonal.

5 Returning to the 3-item example, there are 4 summed scores in this case: 0, 1, 2, and 3. The $4 \times$
 6 8 matrix \mathbf{S} (below) relates the summed score probabilities to the original multinomial probabilities:

$$\bar{\boldsymbol{\pi}} = \begin{pmatrix} \bar{\pi}_0 \\ \bar{\pi}_1 \\ \bar{\pi}_2 \\ \bar{\pi}_3 \end{pmatrix} = \mathbf{S}\boldsymbol{\pi} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \pi_{000} \\ \pi_{001} \\ \pi_{010} \\ \pi_{011} \\ \pi_{100} \\ \pi_{101} \\ \pi_{110} \\ \pi_{111} \end{pmatrix}, \quad \hat{\boldsymbol{\pi}} = \begin{pmatrix} \hat{\pi}_0 \\ \hat{\pi}_1 \\ \hat{\pi}_2 \\ \hat{\pi}_3 \end{pmatrix} = \mathbf{S}\hat{\boldsymbol{\pi}}. \quad (18)$$

7 The observed summed score proportions can be obtained in a similar way:

$$\bar{\boldsymbol{p}} = \begin{pmatrix} \bar{p}_0 \\ \bar{p}_1 \\ \bar{p}_2 \\ \bar{p}_3 \end{pmatrix} = \mathbf{S}\boldsymbol{p} = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} p_{000} \\ p_{001} \\ p_{010} \\ p_{011} \\ p_{100} \\ p_{101} \\ p_{110} \\ p_{111} \end{pmatrix}. \quad (19)$$

8 From Equation (13), under maximum likelihood estimation, the summed score residual vector

9 $\bar{\boldsymbol{p}} - \hat{\boldsymbol{\pi}}$ is asymptotically S -variate normally distributed:

$$\sqrt{N}(\bar{\boldsymbol{p}} - \hat{\boldsymbol{\pi}}) = \sqrt{N}(\mathbf{S}\boldsymbol{p} - \mathbf{S}\hat{\boldsymbol{\pi}}) = \sqrt{N}\mathbf{S}(\boldsymbol{p} - \hat{\boldsymbol{\pi}}) \xrightarrow{D} \mathcal{N}_S(\mathbf{0}, \bar{\boldsymbol{\Gamma}}), \quad (20)$$

10 and $\bar{\boldsymbol{\Gamma}} = \mathbf{S}\boldsymbol{\Gamma}\mathbf{S}' = \mathbf{S}diag(\boldsymbol{\pi})\mathbf{S}' - \mathbf{S}\boldsymbol{\pi}\boldsymbol{\pi}'\mathbf{S}' - \mathbf{S}\boldsymbol{\Delta}\mathcal{F}^{-1}\boldsymbol{\Delta}'\mathbf{S}' = diag(\bar{\boldsymbol{\pi}}) - \bar{\boldsymbol{\pi}}\bar{\boldsymbol{\pi}}' - \bar{\boldsymbol{\Delta}}\mathcal{F}^{-1}\bar{\boldsymbol{\Delta}}'$, with $\bar{\boldsymbol{\Delta}} = \mathbf{S}\boldsymbol{\Delta}$.

11 The reason for introducing the reduction operator matrix \mathbf{S} is primarily a theoretical one. It
 12 facilitates the subsequent derivations of summed score likelihood based indices for testing latent
 13 variable distribution fit. Pragmatically, the Lord-Wingersky (1984) algorithm should be used to

1 compute the model-implied summed score probabilities. If summed score to scale score conversion
 2 tables are computed (see Thissen & Wainer, 2001), the probabilities become automatic byproducts.

3 Goodness of Fit Statistics for IRT models

4 Existing overall goodness of fit indices may be used for testing latent variable distribution fit in
 5 IRT. The full-information test statistics such as likelihood ratio G^2 and Pearson's X^2 use residuals
 6 based on the full response pattern cross-classifications to test the IRT model against the general
 7 multinomial alternative. The comparison between $\hat{\pi}_{\mathbf{u}}$ and $p_{\mathbf{u}}$ (on logarithmic or linear scales) leads to
 8 well-known goodness of fit statistics such as the likelihood ratio G^2 and Pearson's X^2 :

$$G^2 = 2N \sum_{\mathbf{u}} p_{\mathbf{u}} \log \frac{p_{\mathbf{u}}}{\hat{\pi}_{\mathbf{u}}}, \quad X^2 = N \sum_{\mathbf{u}} \frac{(p_{\mathbf{u}} - \hat{\pi}_{\mathbf{u}})^2}{\hat{\pi}_{\mathbf{u}}}. \quad (21)$$

9 Under the null hypothesis that the IRT model fits exactly, these two statistics have the same asymptotic
 10 reference distribution, which is a central chi-square with degrees of freedom equal to $C - 1 - d$
 11 (Bishop et al., 1975). For subsequent development, it is instructive to rewrite Pearson's statistic as a
 12 quadratic form in multinomial residuals: $X^2 = N(\mathbf{p} - \hat{\boldsymbol{\pi}})'[\text{diag}(\hat{\boldsymbol{\pi}})]^{-1}(\mathbf{p} - \hat{\boldsymbol{\pi}})$.

13 Unfortunately, as the number of items increases, the number of response patterns increases
 14 exponentially. For more than a dozen or so dichotomous items (or perhaps a handful of polytomous
 15 items), the contingency table upon which the multinomial is defined becomes sparse for any realistic N .
 16 Consequently the asymptotic chi-square approximations for the full-information test statistics break
 17 down (see e.g., Bartholomew & Tzamourani, 1999) and the utility of the full-information overall
 18 goodness of fit indices for routine IRT applications becomes questionable.

1 Recently, limited-information overall fit statistics such as Maydeu-Olivares and Joe's (2005)
 2 M_2 have been developed. Limited-information fit statistics use residuals based on lower order (e.g.,
 3 first and second order) margins of the contingency table. These lower order margins are far better filled
 4 when compared to the sparse full contingency table. There is growing awareness that
 5 limited-information tests can maintain correct size and can be more powerful than the full-information
 6 tests (Cai, Maydeu-Olivares, Coffman, & Thissen, 2006; Joe & Maydeu-Olivares, 2010).

7 Under the assumption that the number of first and second order margins is larger than the
 8 number of free parameters ($Q > d$) and that $\mathbf{\Delta}_2$ has full column rank (local identification), M_2 can be
 9 written as

$$M_2 = N(\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2)' \tilde{\mathbf{\Delta}}_2 [\tilde{\mathbf{\Delta}}_2' \mathbf{\Xi}_2 \tilde{\mathbf{\Delta}}_2]^{-1} \tilde{\mathbf{\Delta}}_2' (\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2), \quad (22)$$

10 where $\tilde{\mathbf{\Delta}}_2$ is a $Q \times (Q - d)$ orthogonal complement of $\mathbf{\Delta}_2$ such that $\tilde{\mathbf{\Delta}}_2' \mathbf{\Delta}_2 = \mathbf{0}$. From Equation (15),
 11 $(\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2)$ is asymptotically normal with zero means and covariance matrix $\mathbf{\Xi}_2 - \mathbf{\Delta}_2 \mathcal{F}^{-1} \mathbf{\Delta}_2'$, which
 12 implies that the covariance matrix of $\tilde{\mathbf{\Delta}}_2' (\mathbf{p}_2 - \hat{\boldsymbol{\pi}}_2)$ is $\tilde{\mathbf{\Delta}}_2' \mathbf{\Xi}_2 \tilde{\mathbf{\Delta}}_2$. Thus, M_2 is asymptotically chi-square
 13 distributed with $Q - d$ degrees of freedom. In the current simulation study, M_2 will be used as a
 14 benchmark due to its numerous desirable properties identified in the literature (see e.g., Cai & Hansen,
 15 2013). Performance of the proposed latent variable distribution fit indices will be evaluated against M_2 .

16 While an overall test may be used to detect specification errors of latent variable distributions,
 17 the fact that they are also sensitive to other forms of model error (e.g., unmodeled multidimensionality)
 18 makes it difficult to pinpoint the source of misspecification. To that end, more specific diagnostic
 19 indices have been created for IRT. For example, Chen and Thissen's (1997) local dependence indices are

1 particularly sensitive to violations of the local independence assumption. Orlando and Thissen's (2000)
2 item fit diagnostics is another example where the extent to which the IRT model fits the empirical
3 operating characteristics for an item (e.g., whether monotonicity holds) can be examined. The next
4 section develops a set of indices that specifically target latent variable distribution fit for IRT models.

5 The Summed Score Likelihood Based Indices and Statistical Adjustments

6 There are two important lines of reasoning for the derivation of these model fit indices. The first
7 is a recognition based on heuristics: IRT model-implied summed score probabilities may provide useful
8 diagnostic information about the latent variable distributional assumption (Thissen & Wainer, 2001; p.
9 130). The second recognition is that the summed score likelihood based indices are formally
10 limited-information test statistics.

11 *A Heuristic Motivation*

12 When the latent variable distribution assumed in the IRT model does not represent the
13 population distribution of the respondents adequately, the model-implied summed score probabilities
14 $\hat{\pi}_s$ will depart from the observed summed score probabilities \bar{p}_s . Hence all that is needed is to find
15 appropriate test statistics that can summarize the degree to which the model-implied and observed
16 summed score probabilities diverge. It is also preferable if the indices are approximately chi-square
17 distributed test statistics. Pearson's X^2 introduced in the previous section meets this requirement.

18 Recall that the total number of summed scores is $S = 1 + \sum_{i=1}^I (K_i - 1)$. The Pearson-type \bar{X}^2
19 below yields a direct comparison between the model-implied summed score probabilities $\hat{\pi}_s$ and the
20 observed summed score probabilities \bar{p}_s :

$$\bar{X}^2 = N \sum_{s=0}^{S-1} \frac{(\bar{p}_s - \hat{\pi}_s)^2}{\hat{\pi}_s}, \quad (23)$$

1 where \bar{p}_s and $\hat{\pi}_s$ represent the observed and model-implied summed score probability for score s
 2 respectively. This test statistic is different from the full-information test statistic shown in Equation (21)
 3 because it is based on summed score probabilities as opposed to response pattern probabilities.

4 In preliminary studies (Li & Cai, 2012) we had conjectured that under a wide variety of
 5 conditions \bar{X}^2 may have similar asymptotic distributions whose tail-area probabilities can be
 6 approximated by a central chi-squared random variable with $S - 1 - 2$ degrees of freedom under the
 7 null hypothesis that the latent variable distribution $g(\theta)$ is correctly specified in the IRT model. This
 8 conjecture will be tested in the sequel with simulations.

9 The rationale behind the specific degrees of freedom is as follows. The S summed scores
 10 probabilities must sum to 1. The first minus 1 is to reflect that constraint. Had the item parameters been
 11 known, the degrees of freedom would have been exactly $S - 1$. When the item parameters are
 12 estimated (assuming with maximum marginal likelihood), an additional penalty must be introduced to
 13 reflect the effect of parameter estimation. While the location and scale of the latent variable θ are
 14 typically fixed for model identification, the model-implied summed score distribution does not have an
 15 inherent location and scale. The location and scale is determined as a result of estimating the item
 16 parameters. Hence the estimation of item parameters amounts to adding at least two more constraints
 17 for the model-implied summed score probability distribution. The details are of course more complex,
 18 and will be explained next.

1 *A More Formal Derivation*

2 While the proposed test statistics are not associated with particular marginal probabilities in the
 3 same manner as Maydeu-Olivares & Joe's (2005) M_2 , they are nevertheless related to the response
 4 pattern probabilities via the reduction operator matrix \mathbf{S} defined earlier (see Equations 18). It is the
 5 choice of this particular reduction operator that leads to more focused tests targeting latent variable
 6 distribution fit (see Joe & Maydeu-Olivares, 2010). For IRT models with constrained equal item
 7 discrimination parameters (e.g., the 1PL model), it is widely recognized that the summed scores are
 8 sufficient statistics for the latent variables in the model. Though the summed score sufficiency property
 9 does not hold for other IRT models such as the 2PL or the graded model, researchers have nevertheless
 10 found that summed score is an important source of information regarding the ordering of individuals
 11 along the latent variable continuum (e.g., van der Ark, 2005). One could even base parameter
 12 estimation on summed score groups (Chen & Thissen, 1999).

13 Using the reduction operator \mathbf{S} , the derivations above imply that the Pearson-type statistic \bar{X}^2
 14 can be rewritten as

$$\bar{X}^2 = N \sum_{s=0}^{S-1} \frac{(\bar{p}_s - \hat{\pi}_s)^2}{\hat{\pi}_s} = N(\bar{\mathbf{p}} - \hat{\boldsymbol{\pi}})' [\text{diag}(\hat{\boldsymbol{\pi}})]^{-1} (\bar{\mathbf{p}} - \hat{\boldsymbol{\pi}}), \quad (24)$$

15 where $(\bar{\mathbf{p}} - \hat{\boldsymbol{\pi}}) = \mathbf{S}(\mathbf{p} - \hat{\boldsymbol{\pi}})$ is the summed score residual vector (see Equation 20). Under the null
 16 hypothesis that the IRT model is correctly specified, one can obtain the probability limit of the weight
 17 matrix as $\text{plim}([\text{diag}(\hat{\boldsymbol{\pi}})]^{-1}) = [\text{diag}(\bar{\boldsymbol{\pi}})]^{-1}$ by the consistency of the maximum likelihood estimator
 18 (see Equation 9), the continuity of the mapping from $\boldsymbol{\gamma}$ to the summed score probabilities, and the

1 continuity of the matrix inverse. Following results on quadratic forms of random vectors (e.g., Mathai
 2 & Provost, 1992, p. 53), the asymptotic expected value of \bar{X}^2 is equal to

$$\begin{aligned} tr\{\bar{\Gamma}[diag(\bar{\boldsymbol{\pi}})]^{-1}\} &= tr\{[diag(\bar{\boldsymbol{\pi}}) - \bar{\boldsymbol{\pi}}\bar{\boldsymbol{\pi}}'] [diag(\bar{\boldsymbol{\pi}})]^{-1}\} - tr(\bar{\boldsymbol{\Delta}}\mathcal{F}^{-1}\bar{\boldsymbol{\Delta}}' [diag(\bar{\boldsymbol{\pi}})]^{-1}) \quad (25) \\ &= S - 1 - tr\{\mathcal{F}^{-1}\bar{\boldsymbol{\Delta}}' [diag(\bar{\boldsymbol{\pi}})]^{-1}\bar{\boldsymbol{\Delta}}\} = \mu_1. \end{aligned}$$

3 From Equations (24) and (25) we can see that the statistic \bar{X}^2 *cannot* be asymptotically
 4 chi-square distributed. Even though it is a quadratic form in asymptotically normally distributed
 5 random vectors, a key condition for its chi-squaredness is not met. That is, the product of the
 6 probability limit of the weight matrix $[diag(\bar{\boldsymbol{\pi}})]^{-1}$ and the covariance matrix of the normal random
 7 vector $\bar{\boldsymbol{\Gamma}}$ is not idempotent in general, i.e., $\bar{\boldsymbol{\Gamma}}[diag(\bar{\boldsymbol{\pi}})]^{-1}\bar{\boldsymbol{\Gamma}}[diag(\bar{\boldsymbol{\pi}})]^{-1} \neq \bar{\boldsymbol{\Gamma}}[diag(\bar{\boldsymbol{\pi}})]^{-1}$. On the other
 8 hand, Equation (25) shows that the asymptotic expected value of \bar{X}^2 is equal to $S - 1$ minus a
 9 constant that depends on the trace of $\mathcal{F}^{-1}\bar{\boldsymbol{\Delta}}' [diag(\bar{\boldsymbol{\pi}})]^{-1}\bar{\boldsymbol{\Delta}}$, which reflects additional uncertainty due to
 10 estimation of item parameters. With the first-order moment of \bar{X}^2 , the Satorra-Bentler type moment
 11 adjustment approaches can be applied to adjust the statistic, so that the tail area of its distribution can
 12 be better approximated by a chi-squared distribution (Satorra & Bentler, 1994; Cai et al., 2006).

13 *Adjustment of Statistics*

14 According to Satorra & Bentler's (1994) paper, test statistics that do not asymptotically follow a
 15 chi-squared distribution can be corrected, by matching the mean (or the mean & variance) to fixed
 16 degrees of freedom. Let df indicate the degrees of freedom of interest, and μ_1 indicate the asymptotic
 17 expected value of \bar{X}^2 . The moment adjusted statistic is

1 normal, unidimensional nonnormal or multidimensional multivariate normal).

2 In the null condition, response pattern data were simulated with a latent variable having
3 unidimensional normal distribution. In the alternative conditions, response pattern data were
4 simulated either with a nonnormally distributed latent variable or with a bivariate normally distributed
5 latent variable. The nonnormal θ s were generated from a distribution obtained from a 1:4 mixture of
6 two normally distributed densities ($M_1 = 1, SD_1 = 0.4; M_2 = 0, SD_2 = 1$). The multidimensional θ
7 distribution is standard bivariate normal with correlation equal to 0.9, representing substantial overlap
8 between the two dimensions. Half of the items loaded on each dimension in a pure between-item
9 multidimensional model. In other words, each item is only directly influenced by a single dimension,
10 but the dimensions are correlated.

11 There were 3 conditions for item parameters. For the “Equal Slopes and Equal Intercepts”
12 condition, all the slope parameters are fixed to 1, and all the intercept parameters are fixed to 0. For the
13 “Random Slopes and Random Intercepts” condition, parameters for 24 items were randomly generated
14 with properties mimicking standard educational and psychological assessments. Discrimination (a)
15 parameters were drawn from a log-normal distribution ($M = 0.5, SD = 0.2$), the threshold values (b)
16 were drawn from a normal distribution ($M = 0, SD = 0.75$), the intercepts (c) were calculated as $(-ab)$.
17 Parameters for the first 12 items were used for shorter tests. For the “Dispersed Slopes and Dispersed
18 Intercepts” condition, item slope parameters were designed to spread from 1 to 3 in equal increments,
19 while item thresholds spread from -2 to 2 across the 12 or 24 items.

20 The fitted models were standard unidimensional IRT models. In the null conditions, the

1 data-generating models and the fitted models were the same. In the alternative conditions, the fitted
2 models were mis-specified for ignoring either latent variable nonnormality or multidimensionality.
3 Bock and Aitkin's (1981) EM algorithm was used to obtain maximum likelihood estimates, and the
4 Lord-Wingersky (1984) algorithm was used to compute the model-implied summed score probabilities.

5 To compare the performance of the fit statistics, empirical Type I Error rates were computed in
6 the null conditions, and empirically observed power were computed in the alternative conditions at
7 three alpha levels: 0.01, 0.05, and 0.10. In addition, another model fit index, Maydeu-Olivares and Joe's
8 M_2 was employed as a benchmark.

9 Results

10 *Type I Error Rates*

11 Table 2 and Table 3 present the simulation study results for the unidimensional normal case
12 under the null hypothesis. The extent to which the tail areas of the proposed statistics' distribution are
13 well approximated is examined by comparing the observed Type I Error rates against the nominal
14 alpha levels. The results indicate that, when the slope and threshold parameters are equal across items,
15 the adjusted and unadjusted summed score likelihood based indices both work well. Empirical
16 rejection rates and their corresponding alpha levels are close to each other. However, when the item
17 parameters become dispersed, the adjusted statistic \bar{X}_c^2 performs better than the unadjusted statistic
18 \bar{X}^2 . These results hold across different numbers of items and different sample sizes.

19 Furthermore, as suggested earlier, the observed means of these indexes should be close to the
20 expected values of the approximating chi-squared distributions (the degrees of freedom). The results in

1 Table 2 and Table 3 confirm that when item parameters are equal, the means are close to the degrees of
2 freedom, and the variance is approximately twice the degrees of freedom. Notice that when the number
3 of items or the sample size increases, the results improve. For the 2PL model, when the item parameters
4 are dispersed, the moment adjusted statistic \bar{X}_C^2 improves upon \bar{X}^2 with a heuristic degrees of
5 freedom. However, for the graded model, both \bar{X}^2 and the adjusted statistic \bar{X}_C^2 perform well in the
6 null condition. In addition, Maydeu-Olivares and Joe's M_2 appears to be well calibrated for the
7 conditions we tested.

8 *Power*

9 From Table 4 and Table 5, it is clear that the summed score likelihood based indices have
10 substantially higher power than M_2 when the latent variable distribution is nonnormal. The
11 performance of the proposed statistics are heavily influenced by the number of items and dispersion of
12 item parameters. For both 2PL and Graded models, the power of the proposed indexes grow as the
13 sample size and number of items increase. This is to be expected as more data bring more information
14 about the latent variable distribution. When the item slope and threshold parameters are equal across
15 items, the unadjusted and adjusted statistics perform equally well. However, when the item parameters
16 are dispersed, the adjusted statistic \bar{X}_C^2 has higher power than the unadjusted statistic \bar{X}^2 . Finally,
17 Table 6 and Table 7 provide some evidence that the summed score likelihood based indices are not
18 sensitive to model misspecification related to multidimensionality, in contrast to M_2 . This is a desirable
19 feature of the proposed indices, which ought to be more targeted against specific forms of model
20 misspecification. M_2 on the other hand, is a more general index for global model fit assessment.

An Application to Empirical Data

We illustrate the test statistics with empirical data. 12 items related to positive consequences of nicotine (PCN, Tucker et al., 2014), as part of a questionnaire dealing with various attitudes, beliefs and behaviors related to smoking (Shadel, Edelen, & Tucker, 2011), were administered to a sample of 2717 daily cigarette smokers. Each item was rated on a 5-point ordinal scale. This study was part of the development of the National Institute of Health's Patient Reported Outcomes Measurement Information System (PROMIS) and extensive item and dimensional analysis was conducted prior to calibration of the items as unidimensional. The density plot (Figure 1) of the latent variable distribution for this subscale shows its deviation from a standard normal distribution that there are two maximum points in the middle instead of a "bell curve" shape. Table 8 presents the contents of the 12 items from the PROMIS smoking assessment.

Results show that when we use the normal unidimensional IRT model, $\bar{\chi}^2$ equals to 208.5, and $\bar{\chi}_C^2$ equals to 179.3, indicating significant lack of latent variable normality ($df = 46, p < 0.0001$). But when the empirical histogram latent density estimation is used instead for item parameter estimation, $\bar{\chi}^2$ is equal to 51.2 and $\bar{\chi}_C^2$ is equal to 49.4 ($df = 46, p > 0.1$). In sum, we came to the conclusion that the latent variable distribution of this set of items was probably nonnormal and our proposed indices were able to detect the violation of latent variable distribution assumption.

Discussion

Normality of latent variable distribution is a critical assumption in standard maximum marginal likelihood estimation for IRT models. However, in real world applications, the distribution of latent

1 variables can be nonnormal. The detection of latent variable nonnormality is important for item
2 analysis and test scoring. In this study we propose using summed score likelihood based indices for
3 testing departures from normality. We also develop a Satorra-Bentler type moment adjustment
4 approach to approximate the tail area probabilities of the indices.

5 In the simulation study, the performance of unadjusted and adjusted summed score likelihood
6 based statistics was compared to that of M_2 . Results show that the moment-adjusted index performs
7 well for both dichotomous data and polytomous data and maintains correct test size across number of
8 items, sample size, and type of IRT model considered. The unadjusted statistic does not work as well,
9 especially when item parameters are dispersed. Furthermore, the indices were particularly sensitive to
10 latent variable nonnormality, and not sensitive to other kinds of model misfit such as
11 multidimensionality.

12 An interesting finding is that the general goodness-of-fit statistic M_2 (Maydeu-Olivares & Joe,
13 2005) has almost no power against the nonnormal alternative and hence cannot be recommended for
14 testing latent variable distribution fit for IRT models (see also Hansen et al, 2016). This could be
15 explained by the observation that M_2 is based only on first and second order margins of the
16 underlying contingency table, but to detect latent variable distributional misfit, information from
17 higher order margins may be necessary.

18 This study is not without its limitations. First, the distributions of the proposed indices are not
19 exactly chi-squared. In our study, their tail-area probabilities were approximated to first order by a
20 chi-squared variable with the availability of the item parameter error covariance matrix and a Jacobian.

1 We focus on the first order correction due to its simplicity and the fact that we observed empirically
2 that the results of the second order correction did not differ substantively from that of the first order. In
3 the future, higher order moments could be considered to improve the performance of the adjusted
4 statistics for situations that we have not examined. Second, only a limited number of null conditions
5 and only two alternative population distributions were tested in the simulations. More extensive
6 simulations are needed to fully understand the performance of the test statistics. Third, we only studied
7 the properties of the statistics and the corrections under maximum likelihood estimation. In principle,
8 one could derive similar statistics under limited-information estimation (e.g., with weighted least
9 squares). Finally, this study only considered the conditions when item response data are assumed to
10 be unidimensional. Multidimensional IRT models (MIRT, Reckase, 2009) should be considered in
11 subsequent work. One particularly popular model in educational and psychological research is the
12 full-information item bifactor model (Gibbons & Hedeker, 1992; Cai, Yang, & Hansen, 2011; Reise, 2012).
13 In this model, all items load on a general dimension, and an item is permitted to load on at most one
14 specific dimension that influences non-overlapping subsets of items. This feature of bifactor models
15 implies that there exists valuable relation between an observed summed score and the distribution of the
16 latent general dimension (Cai, 2015). This relation implies an opportunity to test the underlying
17 assumption about the distribution of general latent dimension with summed score likelihood based
18 statistics.

References

- Bartholomew, D. J., & Tzamourani, P. (1999). The goodness-of-fit of latent trait models in attitude measurement. *Sociological Methods and Research, 27*, 525–546.
- Bishop, Y. M. M., Fienberg, S. E., & Holland, P. W. (1975). *Discrete multivariate analysis: Theory and practice*. Cambridge, MA: MIT Press.
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: Application of an EM algorithm. *Psychometrika, 46*, 443-459.
- Bock, R. D., & Lieberman, M. (1970). Fitting a response model for n dichotomously scored items. *Psychometrika, 35*, 179-197.
- Cai, L. (2013). flexMIRT® version 2: Flexible multilevel item factor analysis and test scoring [Computer software]. Chapel Hill, NC: Vector Psychometric Group, LLC.
- Cai, L. (2015). Lord-Wingsky algorithm version 2.0 for hierarchical item factor models with applications in test scoring, scale alignment, and model fit testing. *Psychometrika, 80*(2), 535-559.
- Cai, L., & Hansen, M. (2013). Limited-information goodness-of-fit testing of hierarchical item factor models. *British Journal of Mathematical and Statistical Psychology, 66*(2), 245-276.
- Cai, L., Maydeu-Olivares, A., Coffman, D. L., & Thissen, D. (2006). Limited-information goodness-of-fit testing of item response theory models for sparse 2^P tables. *British Journal of Mathematical and Statistical Psychology, 59*, 173–194.
- Cai, L., Yang, J., & Hansen, M. (2011). Generalized full-information item bifactor analysis.

Psychological Methods, 16(3), 221-248.

Chen, W. H., & Thissen, D. (1997). Local dependence indices for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, 22, 265-289.

Chen, W. H., & Thissen, D. (1999). Estimation of item parameters for the three-parameter logistic model using the marginal likelihood of summed scores. *British Journal of Mathematical and Statistical Psychology*, 52, 19-37.

Ferrando, P. J., & Lorenzo-seva, U. (2001). Checking the appropriateness of item response theory models by predicting the distribution of observed scores: The program EO-fit. *Educational and Psychological Measurement*, 61, 895-902.

Gibbons, R., & Hedeker, D. (1992). Full-information item bifactor analysis. *Psychometrika*, 57, 423-436.

Hambleton, R. K., & Traub, R. E. (1973). Analysis of empirical data using two logistic latent trait models. *British Journal of Mathematical and Statistical Psychology*, 26, 195-211.

Hansen, M., Cai, L., Monroe, S., & Li, Z. (2016). Limited-information goodness-of-fit testing of diagnostic classification item response models. *British Journal of Mathematical and Statistical Psychology*, 69(3), 225-252.

Joe, H., & Maydeu-Olivares, A. (2010). A general family of limited information goodness-of-fit statistics for multinomial data. *Psychometrika*, 75, 393-419.

Li, Z., & Cai, L. (2012, July). *Summed score based fit indices for testing latent variable distribution assumption in IRT*. Paper presented at the 2012 International Meeting of the Psychometric

Society, Lincoln, NE.

Lord, F. M. (1953). The relation of test score to the latent trait underlying the test. *Educational and Psychological Measurement*, 13, 517-548.

Lord, F. M., & Wingersky, M. S. (1984). Comparison of IRT true-score and equipercentile observed-score "equatings." *Applied Psychological Measurement*, 8, 453-461.

Mathai, A. M., & Provost, S. B. (1992). *Quadratic forms in random variables: Theory and applications*. New York: Marcel Dekker.

Maydeu-Olivares, A., & Cai, L. (2006). A cautionary note on using $G^2(\text{dif})$ to assess relative model fit in categorical data analysis. *Multivariate Behavioral Research*, 41, 55-64.

Maydeu-Olivares, A., & Joe, H. (2005). Limited and full information estimation and testing in 2^n contingency tables: A unified framework. *Journal of the American Statistical Association*, 100(471), 1009-1020.

Maydeu-Olivares, A., & Joe, H. (2006). Limited information goodness-of-fit testing in multidimensional contingency tables. *Psychometrika*, 71, 713-732.

Mislevy, R. (1984). Estimating latent distributions. *Psychometrika*, 49, 359-381.

Monroe, S. (2014). *Multidimensional item factor analysis with semi-nonparametric latent densities*. University of California, Los Angeles. (Unpublished Doctoral Dissertation).

Monroe, S. & Cai, L. (2014). Estimation of a Ramsay-curve item response theory model by the Metropolis-Hastings Robbins-Monro Algorithm. *Educational and Psychological Measurement*, 42(2), 343-369.

Orlando, M., & Thissen, D. (2000). Likelihood-based item-fit indices for dichotomous item response theory models. *Applied Psychological Measurement, 24*, 50-64.

Reckase, M. (2009). *Multidimensional item response theory (statistics for social and behavioral sciences)*. New York: Springer.

Reise, S. (2012). The rediscovery of bifactor measurement models. *Multivariate Behavioral Research, 47*(5), 667-696.

Ross, J. (1966). An empirical study of a logistic mental test model. *Psychometrika, 31*, 325-340.

Satorra, A., & Bentler, P. M. (1994). Corrections to test statistics and standard errors in covariance structure analysis. In A. von & C. C. Clogg (Eds.), *Latent variables analysis: Applications for developmental research* (pp. 399-419). Thousand Oaks, CA, US: Sage Publications, Inc.

Shadel W.G., Edelen, M. & Tucker, J. S. (2011). A unified framework for smoking assessment: the PROMIS Smoking Initiative. *Nicotine & Tobacco Research, 13*(5), 399-400.

Sinharay, S., Johnson, M. S. & Stern, H. S. (2006). Posterior predictive assessment of item response theory models. *Applied Psychological Measurement, 30*, 298-321.

Thissen, D. & Steinberg, L. (2009). Item response theory. In R. Millsap & A. Maydeu-Olivares (Eds.), *The Sage handbook of quantitative methods in psychology* (pp. 148-177). London: Sage Publications.

Thissen, D., & Wainer, H. (2001). *Test scoring*. Hillsdale, NJ: Lawrence Erlbaum Associates.

Tucker, J., Shadel, W. G., Stucky, B., Cerully, J., Li, Z., Hansen, M., & Cai, L. (2014).

- Development of the PROMIS positive emotional and sensory expectancies of smoking item banks. *Nicotine and Tobacco Research*, 16(Suppl 3), S212-222.
- van der Ark, L. A. (2005). Stochastic ordering of the latent trait by the sum score under various polytomous IRT models. *Psychometrika*, 70, 283-304.
- Woods, C. M. (2006). Ramsay-curve item response theory to detect and correct for non-normal latent variables. *Psychological Methods*, 11, 253-270.
- Woods, C. M., & Lin, N. (2009). Item response theory with estimation of the latent density using Davidian curves. *Applied Psychological Measurement*, 33, 102-117.
- Woods, C. M., & Thissen, D. (2006). Item response theory with estimation of the latent population distribution using spline-based densities. *Psychometrika*, 71, 281-301.
- Zimowski, M. F., Muraki, E., Mislevy, R. J., & Bock, R. D. (1996). BILOG-MG: Multiple-group IRT analysis and test maintenance for binary items [Computer software]. Lincolnwood, IL: Scientific Software International.

Appendix: A Modified Lord-Wingersky Algorithm for Jacobian Computations

Consider a test with n dichotomous items, calibrated by a 2PL IRT model. Recall that $T_i(1|\theta)$ is item i 's traseline for category 1 (Equation 1), with $T_i(0|\theta) = 1 - T_i(1|\theta)$ for category 0. Theoretically, there should be 2^n response patterns. The response patterns is indicated by $\mathbf{u} = (u_1, \dots, u_n)$. Under the assumption of items' conditional independence, the likelihood for a response pattern \mathbf{u} can be expressed as $L(\mathbf{u}|\theta) = \prod_{i=1}^n T_i(u_i|\theta)$. For n dichotomous items, the summed score s ranges from 0 to n . $S = n + 1$ is the number of all possible summed scores. Recall that $\|\mathbf{u}\| = \sum_{i=1}^n u_i$ is a notational shorthand for the summed score associated with response \mathbf{u} (see Equation 16). The likelihood for summed score $s = 0, \dots, n$ is defined as

$$L(s|\theta) = \sum_{\|\mathbf{u}\|=s} L(\mathbf{u}|\theta) = \sum_{s=\|\mathbf{u}\|} \prod_{i=1}^n T_i(u_i|\theta), \quad (27)$$

Clearly, the likelihood of a summed score s is the sum of all response pattern likelihoods for $\|\mathbf{u}\| = s$. In Lord-Wingersky algorithm, the summed score likelihoods are built up recursively, one at a time (Lord & Wingersky, 1984). Let $L_i(s|\theta)$ indicate the likelihood for summed score s after item i has been added into the computation. In the first step, two summed score likelihoods are computed based on the trelines of item 1: $L_1(0|\theta) = T_1(0|\theta)$ and $L_1(1|\theta) = T_1(1|\theta)$.

In the second step, we have three summed score likelihoods based on the likelihoods from step 1 and trelines of item 2:

$$L_2(0|\theta) = L_1(0|\theta)T_2(0|\theta), \quad (28)$$

$$L_2(1|\theta) = L_1(1|\theta)T_2(0|\theta) + L_1(0|\theta)T_2(1|\theta),$$

$$L_2(2|\theta) = L_1(1|\theta)T_2(1|\theta).$$

Suppose n items have been added. The likelihoods for summed scores $(0, \dots, n)$ are:

$$L_n(0|\theta) = L_{n-1}(0|\theta)T_n(0|\theta),$$

...

$$L_n(s|\theta) = L_{n-1}(s|\theta)T_n(0|\theta) + L_{n-1}(s-1|\theta)T_n(1|\theta), \quad (29)$$

...

$$L_n(n|\theta) = L_{n-1}(n-1|\theta)T_n(1|\theta).$$

To obtain the Jacobian matrix of summed score likelihoods with respect to item parameters, the Lord-Wingersky algorithm is adapted slightly. As previously mentioned, in the first step, there are only two summed score likelihoods based on item 1: $L_1(0|\theta)$ and $L_1(1|\theta)$.

The first-order derivatives of summed score likelihoods with respect to a generic item parameter γ_1 for item 1 are:

$$\begin{aligned} \frac{\partial L_1(0|\theta)}{\partial \gamma_1} &= \frac{\partial T_1(0|\theta)}{\partial \gamma_1}, \\ \frac{\partial L_1(1|\theta)}{\partial \gamma_1} &= \frac{\partial T_1(1|\theta)}{\partial \gamma_1}. \end{aligned} \quad (30)$$

In the second step, item 2 is added with a generic item parameter γ_2 . The first-order derivatives of summed score likelihoods with respect to γ_1 and γ_2 follows from the chain rule:

$$\begin{aligned} \frac{\partial L_2(0|\theta)}{\partial \gamma_1} &= \frac{\partial L_1(0|\theta)}{\partial \gamma_1} T_2(0|\theta), \\ \frac{\partial L_2(1|\theta)}{\partial \gamma_1} &= \frac{\partial L_1(1|\theta)}{\partial \gamma_1} T_2(0|\theta) + \frac{\partial L_1(0|\theta)}{\partial \gamma_1} T_2(1|\theta), \end{aligned} \quad (31)$$

$$\begin{aligned}
\frac{\partial L_2(2|\theta)}{\partial \gamma_1} &= \frac{\partial L_1(1|\theta)}{\partial \gamma_1} T_2(1|\theta), \\
\frac{\partial L_2(0|\theta)}{\partial \gamma_2} &= L_1(0|\theta) \frac{\partial T_2(0|\theta)}{\partial \gamma_2}, \\
\frac{\partial L_2(1|\theta)}{\partial \gamma_2} &= L_1(1|\theta) \frac{\partial T_2(0|\theta)}{\partial \gamma_2} + L_1(0|\theta) \frac{\partial T_2(1|\theta)}{\partial \gamma_2}, \\
\frac{\partial L_2(1|\theta)}{\partial \gamma_2} &= L_1(1|\theta) \frac{\partial T_2(1|\theta)}{\partial \gamma_2}.
\end{aligned}$$

Generalizing to n items, the first-order derivatives of summed score likelihood

functions with respect to the n item's parameters $(\gamma_1, \dots, \gamma_n)$ are:

$$\begin{aligned}
\frac{\partial L_n(0|\theta)}{\partial \gamma_1} &= \frac{\partial L_{n-1}(0|\theta)}{\partial \gamma_1} T_n(0|\theta), \\
&\dots \\
\frac{\partial L_n(s|\theta)}{\partial \gamma_1} &= \frac{\partial L_{n-1}(s|\theta)}{\partial \gamma_1} T_n(0|\theta) + \frac{\partial L_{n-1}(s-1|\theta)}{\partial \gamma_1} T_n(1|\theta), \\
&\dots \\
\frac{\partial L_n(n|\theta)}{\partial \gamma_1} &= \frac{\partial L_{n-1}(n-1|\theta)}{\partial \gamma_1} T_n(1|\theta), \\
&\dots \\
&\dots \\
&\dots \\
\frac{\partial L_n(0|\theta)}{\partial \gamma_n} &= L_{n-1}(0|\theta) \frac{\partial T_n(0|\theta)}{\partial \gamma_n}, \\
&\dots \\
\frac{\partial L_n(s|\theta)}{\partial \gamma_n} &= L_{n-1}(s|\theta) \frac{\partial T_n(0|\theta)}{\partial \gamma_n} + L_{n-1}(s-1|\theta) \frac{\partial T_n(1|\theta)}{\partial \gamma_n}, \\
&\dots \\
\frac{\partial L_n(n|\theta)}{\partial \gamma_n} &= L_{n-1}(n-1|\theta) \frac{\partial T_n(1|\theta)}{\partial \gamma_n}.
\end{aligned} \tag{32}$$

The process of modified Lord-Wingsky algorithm for calculating the Jacobian matrix is

illustrated with an example. Consider a simple test with three dichotomous items. The values of slope parameters are $\mathbf{a} = (1.0, 0.8, 1.2)$, and the values of intercept parameters are $\mathbf{c} = (-0.2, 0.6, -1.0)$. Recall that the marginal probability for summed scores with known $g(\theta)$ is

$$p(s) = \int L(s|\theta) g(\theta) d\theta, \quad (33)$$

The integrals in Equation (33) must be approximated by quadrature. We demonstrate the algorithm by showing the calculations over a set of quadrature points (Cai, 2015). We approximate the marginal probability using Q quadrature points:

$$p(s) = \int L(s|\theta) g(\theta) d\theta = \sum_{q=1}^Q L(s|X_q) W(X_q), \quad (34)$$

where X_q is a quadrature node and $W(X_q)$ is the corresponding quadrature weight. To obtain $W(X_q)$, a set of normalized ordinates of the prior density are applied (Cai, 2015), i.e., $W(X_q) = g(X_q) / \sum_{q=1}^Q g(X_q)$.

Table A1 shows the recursive computations for the parameters from item 3. It shows the values of summed score likelihoods, first-order derivatives of tracelines, and the first-order derivatives of summed score likelihoods at five equally spaced quadrature points ($Q = 5$): -2, -1, 0, 1, and 2. More quadrature points should be used for better precision (Cai, 2015). The first block presents the summed score likelihoods after the 1st and 2nd items are added in. The second block presents the first-order derivatives of item 3's tracelines with respect to its slope parameter. The third block presents the first-order derivatives of summed score likelihoods with respect to item 3's slope parameter.

Table A2 presents the first-order derivatives of summed score probabilities with respect to item 3 (the desired Jacobian elements). $W(\theta)$ indicates quadrature weights at each θ level. “Weighted derivatives” are found by multiplying (point to point) the first-order derivatives of summed score likelihoods with $W(\theta)$. The last column “Jacobian” indicates the first-order derivatives of summed score probabilities with respect to item 3’s slope parameter. It is the summation of the weighted derivatives over all quadrature points for each summed score likelihood.

Table 1

Manipulated Factors and Conditions for Simulation Study

Factor (Levels)	Conditions
Types of IRT Model (2)	2PL, Graded
Number of Items (2)	12, 24
Sample Size (3)	500, 1000, 1500
Values of Item Parameters (3)	Equal Slopes and Equal Intercepts
	Random Slopes and Random Intercepts
	Dispersed Slopes and Dispersed Intercepts
Latent Variable Distribution (3)	Normally Distributed Unidimensional
	Nonnormally Distributed Unidimensional
	Correlated Bivariate Normally Distributed

Notes. The factors are fully crossed in a $2 \times 2 \times 3 \times 3 \times 3$ design with 1000 attempted replications per cell.

Table 2

Selected Simulation Results under the Null Hypothesis: Normally Distributed Unidimensional Latent Variable in 2PL Models

n	N	Index	df	Equal Slopes and Intercepts					Random Slopes and Intercepts					Dispersed Slopes and Intercepts				
				M	Var	ERR*		KS	M	Var	ERR*		KS	M	Var	ERR*		KS
						.01	.05				.01	.05				.01	.05	
12	500	\bar{X}^2	10	10.0	19.7	.01	.04	.94	9.2	18.1	.01	.04	<u>.00</u>	9.1	18.0	.01	.03	<u>.00</u>
		\bar{X}_C^2	10	10.0	19.9	.01	.05	.98	10.1	21.5	.01	.06	.51	10.4	23.4	.02	.06	.08
		M_2	54	54.3	119	.02	.06	.33	54.5	112.8	.01	.06	.15	54.2	116.4	.01	.06	.25
12	1500	\bar{X}^2	10	10.0	18.8	.01	.05	.64	9.3	17.2	.01	.03	<u>.00</u>	8.9	15.9	.01	.03	<u>.00</u>
		\bar{X}_C^2	10	10.0	18.9	.01	.05	.49	10.1	20.1	.01	.05	.96	10.0	20.2	.01	.04	1.00
		M_2	54	53.8	112.9	.01	.06	.83	53.8	101.2	.01	.04	.79	53.8	107.0	.01	.05	.36
24	1500	\bar{X}^2	22	21.8	45.3	.01	.06	.15	21.5	40.2	.01	.04	<u>.03</u>	21.4	44.8	.01	.05	<u>.03</u>
		\bar{X}_C^2	22	21.9	45.5	.01	.06	.26	22.1	42.3	.01	.05	.75	22.2	48.0	.02	.06	.70
		M_2	252	252.9	483.3	.02	.05	.20	251.9	502.2	.01	.05	.81	251.6	526.7	.01	.04	.55

**Note: Empirical Rejection Rates at α levels 0.01 and 0.05.*

Table 3

Selected Simulation Results under the Null Hypothesis: Normally Distributed Unidimensional Latent Variable in Graded Models

n	N	Index	df	Equal Slopes and Intercepts					Random Slopes and Intercepts					Dispersed Slopes and Intercepts				
				M	Var	ERR*		KS	M	Var	ERR*		KS	M	Var	ERR*		KS
						.01	.05				p	.01				.05	p	
12	500	\bar{X}^2	34	34.5	69.9	.01	.06	.07	34.1	65.7	.01	.05	.67	33.8	65.6	.01	.04	.97
		\bar{X}_C^2	34	34.6	70.3	.01	.06	<u>.02</u>	34.4	66.8	.01	.05	.12	34.5	68.4	.01	.06	<u>.04</u>
		M_2	30	29.5	60.4	.01	.05	.06	29.8	59.0	.01	.05	.63	29.7	60.0	.01	.05	.17
12	1500	\bar{X}^2	34	33.8	69.0	.01	.04	.64	33.4	59.6	.01	.03	<u>.01</u>	33.5	61.8	.01	.03	.21
		\bar{X}_C^2	34	33.9	69.3	.01	.04	.75	33.8	60.9	.01	.04	.10	34.4	65.0	.01	.05	.26
		M_2	30	31.8	80.6	.02	.08	<u>.00</u>	29.8	56.6	.01	.04	.77	30.1	58.6	.01	.04	.24
24	1500	\bar{X}^2	70	70.9	147.1	.01	.06	<u>.03</u>	70.3	136.0	.01	.05	.46	69.9	138.6	.01	.05	.36
		\bar{X}_C^2	70	70.9	147.3	.01	.06	<u>.03</u>	70.5	137.1	.01	.06	.24	70.4	140.7	.01	.05	.15
		M_2	204	204.3	435.9	.01	.06	.58	202.0	369.9	.01	.03	<u>.00</u>	204.6	414.6	.01	.05	.29

**Note: Empirical Rejection Rates at α levels 0.01 and 0.05.*

Table 4

Selected Simulation Results under the Alternative Hypothesis: Nonnormally Distributed Unidimensional Latent Variable in 2PL Models

n	N	Index	df	Equal Slopes and Intercepts			Random Slopes and Intercepts			Dispersed Slopes and Intercepts					
				M	Power			M	Power			M	Power		
					.01	.05	.10		.01	.05	.10		.01	.05	.10
12	500	\bar{X}^2	10	13.1	.07	.17	.28	11.8	<u>.02</u>	.10	.20	10.7	<u>.02</u>	<u>.06</u>	<u>.13</u>
		\bar{X}_C^2	10	13.2	.07	.19	.28	13.0	.04	.17	.27	12.4	.04	.14	.24
		M_2	54	53.8	.01	.04	.10	55.2	.02	.07	.14	55.6	.02	.08	.13
12	1500	\bar{X}^2	10	19.3	.28	.50	.62	17.1	.16	.39	.54	14.3	.06	.21	.34
		\bar{X}_C^2	10	19.4	.28	.51	.62	18.7	.23	.47	.63	16.4	.13	.34	.48
		M_2	54	53.8	.01	.05	.11	55.4	.01	.06	.12	55.2	.02	.06	.12
24	1500	\bar{X}^2	22	40.9	.49	.73	.83	38.1	.39	.64	.76	37.2	.36	.62	.74
		\bar{X}_C^2	22	41.2	.50	.73	.84	39.4	.45	.69	.80	38.8	.42	.68	.79
		M_2	252	250.5	.01	.04	.10	255.9	.02	.08	.14	258.7	.02	.10	.17

Table 5

Selected Simulation Results under the Alternative Hypothesis: Nonnormally Distributed Unidimensional Latent Variable in Graded Models

n	N	Index	df	Equal Slopes and			Random Slopes and			Dispersed Slopes and					
				Intercepts			Intercepts			Intercepts					
				M	Power		M	Power		M	Power				
	.01	.05	.10		.01	.05	.10		.01	.05	.10				
12	500	\bar{X}^2	34	38.0	.04	.12	.21	40.9	.06	.19	.32	38.9	<u>.02</u>	.13	.22
		\bar{X}_C^2	34	38.2	.04	.13	.22	41.4	.07	.21	.33	39.9	.03	.15	.27
		M_2	30	29.4	.01	.03	.08	29.6	.01	.05	.10	30.0	.01	.06	.10
12	1500	\bar{X}^2	34	45.5	.15	.37	.49	50.9	.29	.55	.70	45.8	.14	.37	.51
		\bar{X}_C^2	34	45.6	.16	.37	.49	51.6	.31	.58	.71	47.2	.18	.42	.56
		M_2	30	31.6	.03	.09	.15	30.3	.01	.06	.11	30.3	.01	.05	.10
24	1500	\bar{X}^2	70	93.8	.32	.57	.70	96.8	.38	.66	.79	93.6	.29	.57	.71
		\bar{X}_C^2	70	94.0	.33	.57	.70	97.4	.40	.67	.80	94.5	.32	.60	.73
		M_2	204	202.5	.01	.04	.09	204.5	.01	.04	.09	206.8	.02	.07	.12

Table 6

Selected Simulation Results under the Alternative Hypothesis: Multidimensional Distributed Unidimensional Latent Variable in 2PL Models

n	N	Index	df	Equal Slopes and			Random Slopes and			Dispersed Slopes and					
				Intercepts			Intercepts			Intercepts					
				M	Power		M	Power		M	Power				
				.01	.05	.10		.01	.05	.10		.01	.05	.10	
12	500	\bar{X}^2	10	10.0	.02	.05	.09	9.8	.01	.04	.10	9.2	.00	.04	.07
		\bar{X}_C^2	10	10.1	.02	.05	.09	10.5	.02	.07	.13	10.3	.02	.07	.11
		M_2	54	83.9	.54	.74	.81	161.2	1.00	1.00	1.00	142.4	1.00	1.00	1.00
12	1500	\bar{X}^2	10	10.2	.01	.05	.11	10.1	.02	.06	.11	9.4	.00	.04	.07
		\bar{X}_C^2	10	10.2	.01	.05	.11	10.8	.03	.09	.14	10.5	.01	.07	.13
		M_2	54	140.8	1.00	1.00	1.00	377.3	1.00	1.00	1.00	320.5	1.00	1.00	1.00
24	1500	\bar{X}^2	22	21.8	.01	.06	.11	22.4	.01	.07	.11	22.0	.01	.05	.09
		\bar{X}_C^2	22	21.8	.01	.06	.11	22.9	.02	.08	.13	22.7	.02	.07	.12
		M_2	252	617.8	1.00	1.00	1.00	1281	1.00	1.00	1.00	1544	1.00	1.00	1.00

Table 7

Selected Simulation Results under the Alternative Hypothesis: Multidimensional Distributed Unidimensional Latent Variable in Graded Models

n	N	Index	df	Equal Slopes and			Random Slopes and			Dispersed Slopes and					
				Intercepts			Intercepts			Intercepts					
				M	Power		M	Power		M	Power				
				.01	.05	.10		.01	.05	.10		.01	.05	.10	
12	500	\bar{X}^2	34	34.7	.01	.06	.12	34.3	.01	.06	.10	34.1	.01	.05	.10
		\bar{X}_C^2	34	34.7	.01	.06	.12	34.5	.01	.06	.11	34.7	.01	.06	.13
		M_2	30	50.3	.42	.60	.68	183.4	1.00	1.00	1.00	205.8	1.00	1.00	1.00
12	1500	\bar{X}^2	34	34.6	.01	.07	.12	34.2	.01	.04	.10	33.1	.01	.03	.08
		\bar{X}_C^2	34	34.6	.01	.07	.12	34.4	.01	.05	.11	34.0	.01	.04	.10
		M_2	30	88.2	.84	.90	.92	429.7	1.00	1.00	1.00	455.2	1.00	1.00	1.00
24	1500	\bar{X}^2	70	71.3	.02	.07	.13	70.4	.01	.04	.10	70.2	.01	.06	.11
		\bar{X}_C^2	70	71.3	.02	.07	.13	70.6	.01	.05	.10	70.7	.01	.07	.12
		M_2	204	554.5	1.00	1.00	1.00	1961	1.00	1.00	1.00	2265	1.00	1.00	1.00

Table 8 Items from PROMIS Smoking Initiative

Item Wordings	
Item 1	Smoking helps me concentrate.
Item 2	Smoking helps me think more clearly.
Item 3	Smoking helps me stay focused.
Item 4	Smoking makes me feel better in social situations.
Item 5	Smoking makes me feel more self-confident with others.
Item 6	Smoking helps me feel more relaxed when I'm with other people.
Item 7	Smoking helps me deal with anxiety.
Item 8	Smoking calms me down.
Item 9	If I'm feeling irritable, a cigarette will help me relax.
Item 10	Smoking a cigarette energizes me.
Item 11	Smoking makes me feel less tired.
Item 12	Smoking perks me up.

Table A1

First-order Derivatives of Summed Score Likelihoods with Respect to Item 3's Slope Parameter at 5 Rectangular

Quadrature Points

Quadrature Points	-2	-1	0	1	2
Summed Score Likelihoods After Items 1 and 2					
$L_2(\mathbf{0} \theta)$.658	.423	.195	.061	.014
$L_2(\mathbf{1} \theta)$.315	.473	.515	.385	.213
$L_2(\mathbf{2} \theta)$.027	.104	.291	.553	.773
Derivatives of Tracelines with Respect to Item 3's Slope Parameter					
$\frac{\partial T_3(\mathbf{1} \theta)}{\partial a_3}$	-.063	-.090	.000	.248	.317
$\frac{\partial T_3(\mathbf{0} \theta)}{\partial a_3}$.063	.090	.000	-.248	-.317
First-order Derivatives of Summed Score Likelihoods					
$\frac{\partial L_3(\mathbf{0} \theta)}{\partial a_3} = L_2(\mathbf{0} \theta) \frac{\partial T_3(\mathbf{0} \theta)}{\partial a_3}$.041	.038	.000	-.015	-.004
$\frac{\partial L_3(\mathbf{1} \theta)}{\partial a_3} = L_2(\mathbf{1} \theta) \frac{\partial T_3(\mathbf{0} \theta)}{\partial a_3} + L_2(\mathbf{0} \theta) \frac{\partial T_3(\mathbf{1} \theta)}{\partial a_3}$	-.021	.005	.000	-.080	-.063
$\frac{\partial L_3(\mathbf{2} \theta)}{\partial a_3} = L_2(\mathbf{2} \theta) \frac{\partial T_3(\mathbf{0} \theta)}{\partial a_3} + L_2(\mathbf{1} \theta) \frac{\partial T_3(\mathbf{1} \theta)}{\partial a_3}$	-.018	-.033	.000	-.042	-.177
$\frac{\partial L_3(\mathbf{3} \theta)}{\partial a_3} = L_2(\mathbf{2} \theta) \frac{\partial T_3(\mathbf{1} \theta)}{\partial a_3}$	-.002	-.009	.000	.137	.245

Table A2

First-order Derivatives of Summed Score Probabilities with Respect to Item 3's Slope Parameter at 5 Rectangular

Quadrature Points

	Quadrature Points					
	-2	-1	0	1	2	
$W(\theta)$.054	.244	.403	.244	.054	
First-order Derivatives of Summed Score Likelihoods						
$\frac{\partial L_3(0 \theta)}{\partial a_3}$.041	.038	.000	-.015	-.004	
$\frac{\partial L_3(1 \theta)}{\partial a_3}$	-.021	.005	.000	-.080	-.063	
$\frac{\partial L_3(2 \theta)}{\partial a_3}$	-.018	-.033	.000	-.042	-.177	
$\frac{\partial L_3(3 \theta)}{\partial a_3}$	-.002	-.009	.000	.137	.245	
Weighted Derivatives					Jacobian	
$\frac{\partial L_3(0 \theta)}{\partial a_3} * W(\theta)$.002	.009	.000	-.004	.000	.008
$\frac{\partial L_3(1 \theta)}{\partial a_3} * W(\theta)$	-.001	.001	.000	-.020	-.003	-.023
$\frac{\partial L_3(2 \theta)}{\partial a_3} * W(\theta)$	-.001	-.008	.000	-.010	-.010	-.029
$\frac{\partial L_3(3 \theta)}{\partial a_3} * W(\theta)$.000	-.002	.000	.033	.013	.044

Figure 1: Latent Variable Distribution for Empirical Data. Estimated (using empirical histogram)

probability density of the latent variable is plotted (dotted line), when super-imposed on a standard normal density (solid line).

