

This is a “preproof” accepted article for *Psychometrika*.

This version may be subject to change during the production process.

DOI: 10.1017/psy.2024.20

## Generalized Bayesian Method for Diagnostic Classification Models

Kazuhiro Yamaguchi<sup>1</sup>, Yanlong Liu<sup>2</sup>, and Gongjun Xu<sup>2</sup>

<sup>1</sup> University of Tsukuba

<sup>2</sup> University of Michigan

### Author Note

The data analysis code is available in the Open Science Framework page:

<https://osf.io/sau6j/>. The authors declare no conflicts of interest. This work was supported by JSPS KAKENHI 20H01720, 21H00936, 22K13810, 23H00985, 23H00065, and 24K00485.

Correspondence concerning this article should be addressed to Kazuhiro Yamaguchi, Faculty of Human Science, University of Tsukuba, Institutes of Human Sciences A314, 1-1-1 Tennodai, Tsukuba, Ibaraki, 305-0006, Japan.

Email: [yamaguchi.kazuhir.ft@u.tsukuba.ac.jp](mailto:yamaguchi.kazuhir.ft@u.tsukuba.ac.jp)

This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivatives licence (<http://creativecommons.org/licenses/by-nc-nd/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is unaltered and is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use or in order to create a derivative work.

## Abstract

1  
2 This study extends the loss function-based parameter estimation method for diagnostic  
3 classification models proposed by C. Ma, de la Torre, et al. (2023, Psychometrika) to consider  
4 prior knowledge and uncertainty of sampling. To this end, we integrate the loss function-based  
5 estimation method with the generalized Bayesian method. We establish the consistency of  
6 attribute mastery patterns of the proposed generalized Bayesian method. The proposed  
7 generalized Bayesian method is compared in a simulation study and found to be superior to the  
8 previous nonparametric diagnostic classification method—a special case of the loss function-  
9 based method. Moreover, the proposed method is applied to real data and compared with  
10 previous parametric and nonparametric estimation methods. Finally, practical guidelines for the  
11 proposed method and future research directions are discussed.

12 *Keywords:* diagnostic classification models, parameter estimation, loss function-based  
13 method, generalized Bayesian method

14

# Generalized Bayesian Method for Diagnostic Classification Models

## 1. Introduction

Learning is an important aspect of human life. The current status of individual knowledge or depth of understanding must be evaluated to ensure efficient learning. Test analysis models called diagnostic classification models (DCMs; Rupp et al., 2010; von Davier & Lee, 2019) have been popularly employed to capture an individual's learning status. Notably, DCMs provide useful statistical tools to reveal individuals' current learning status based on the test's item responses. Latent knowledge or cognitive elements are called attributes and expressed as latent categorical variables in DCMs. Moreover, DCMs are known as restricted latent class models (e.g., Rupp & Templin, 2008; Xu, 2017), wherein each possible set of attributes represents a latent class. In other words, attribute mastery patterns indicate the attributes that are either mastered or not mastered. Therefore, one of the DCMs' final outputs is the estimate of the attribute mastery patterns of individuals or attribute mastery probabilities.

Various parameter estimation methods for the DCMs have been actively developed. Parametric and nonparametric estimation methods are commonly used in DCMs. Parametric estimation methods assume parametric item-response functions and structural models. Therefore, parametric estimation methods employ a likelihood function under the assumed model and include (penalized or regularized) maximum likelihood estimation (e.g., Chen et al., 2015; de la Torre, 2009; C. Ma, Ouyang, et al., 2023) and Bayesian estimation methods (e.g., Culpepper, 2015; Yamaguchi & Okada, 2020; Yamaguchi & Templin, 2022b), incorporating prior distributions for model parameters. Numerous parametric estimation methods have been developed and their properties have been studied (e.g., von Davier & Lee, 2019).

On the other hand, nonparametric methods (e.g., Chiu & Douglas, 2013; Chiu et al., 2018) do not use probabilistic item-response models; instead, they use an ideal response to define a type of discrepancy function, which will be formally defined in a later section. Such discrepancy functions are defined based on the distance between each item's ideal and actual responses. Intuitively, nonparametric methods directly estimate attribute mastery patterns, which

1 minimize the discrepancy function. Therefore, nonparametric methods do not require a  
2 probabilistic item-response function. Nonparametric methods exhibit satisfactory statistical  
3 properties, such as consistency under certain conditions (Chiu & Köhn, 2019; Wang & Douglas,  
4 2015).

5         Recently, a general parameter estimation method that can uniformly express parametric  
6 and nonparametric methods was developed by C. Ma, de la Torre, et al. (2023). The unified  
7 estimation method developed by C. Ma, de la Torre, et al. (2023) is a loss function-based  
8 estimation method for DCMs. If we select cross-entropy for a loss function, its minimization  
9 corresponds to maximizing the joint maximum likelihood (Chiu et al., 2016). The distance or  
10 discrepancy function in nonparametric methods is a well-known loss function. Additionally, by  
11 adding penalty terms to a cross-entropy loss function, we obtain the maximum a posteriori  
12 (MAP) estimates, classical Bayesian estimates, to minimize it. These examples indicate that the  
13 loss function-based estimation method is flexible and can represent various estimation methods  
14 in a unified manner. Furthermore, a unified estimation algorithm for the loss function-based  
15 estimation method was available.

16         However, loss function-based methods exhibit certain limitations. First, these methods  
17 only provide point estimates, which may be problematic because we cannot evaluate how point  
18 estimates vary due to sampling or estimation variations. Therefore, we cannot evaluate the  
19 uncertainty of attribute mastery using the loss function-based method. Furthermore, attribute  
20 mastery probabilities for each individual are not expressed in the loss function-based method.  
21 This is the same problem that occurs in DCMs' nonparametric estimation method. However,  
22 attribute mastery probabilities represent a more nuanced situation than attribute mastery pattern  
23 results, with or without mastery. Another limitation of these methods is that prior information on  
24 weight parameters in the generalized nonparametric method that defines generalized ideal  
25 responses is generally not considered. However, DCM users may have prior knowledge of the  
26 test items' conjunctive and disjunctive nature. If so, domain-specific knowledge must be included  
27 to improve parameter estimates.

1           It is not only loss function-based methods that have limitations that need to be addressed;  
2 several limitations of previous parametric and nonparametric estimation methods likewise need  
3 to be noted. First, parametric estimation methods need to specify data-generating distributions,  
4 which determine the likelihood function. The likelihood function provides a connection between  
5 data and model parameters such as attribute mastery patterns. Moreover, likelihood functions  
6 make it possible to evaluate estimation uncertainty with the asymptotical theory within the  
7 maximum likelihood framework or the posterior distribution within the Bayesian framework.  
8 However, the data-generating process is not always specified. DCMs are part of the educational  
9 measurement model family that need various constraints and limitations, making it difficult to  
10 specify the model.

11           Some of the limitations of the nonparametric methods are the same as those of the loss  
12 function-based methods. For instance, current nonparametric methods for DCMs cannot evaluate  
13 the uncertainty of attribute mastery estimates. The nonparametric methods for DCMs were  
14 developed in studies with small sample sizes (Chiu & Douglas, 2013). Ultimately, the  
15 nonparametric method for DCMs can be applied to individuals; however, the parameter  
16 estimates need to be evaluated with variability of parameter estimates. Currently, the  
17 nonparametric methods simply select attribute mastery patterns to minimize prespecified  
18 distance functions so the parameter uncertainty evaluation is not included in the framework. The  
19 parameter estimates with nonparametric methods can be changed by small differences of the loss  
20 function. One main purpose of DCMs is the diagnosis of individual knowledge. Thus, such  
21 variations in parameter estimates due to small differences in the distance functions may be a  
22 fundamental problem for application.

23           To overcome these limitations and extend the previous loss function-based estimation  
24 method for DCMs, we employ the generalized Bayesian (GB; Bissiri et al., 2016) method. The  
25 usual Bayesian parameter update determines the likelihood function and updates the model  
26 parameters in the likelihood with the observed dataset. By contrast, the GB method can express  
27 parameter updating with a dataset via loss functions. Therefore, Bayesian inference is applicable

1 to nonparametric-based estimation methods as well as to likelihood-based methods. Moreover,  
2 other benefits of the GB method are as follows. First, the GB method originally assumes the  $\mathcal{M}$ -  
3 open setting (Bernardo & Smith, 2009, Chap. 6), which implies that the GB method provides a  
4 valid inference even if the assumed model does not match the true data-generating mechanism.  
5 Various DCMs have been developed; however, selecting an appropriate item response function  
6 that expresses the true data-generating mechanism is not always possible. The GB method does  
7 not require an entire data-generating model but instead, sets a loss function related to the  
8 parameter sets of interest. This means that we do not need to find a correct data-generating  
9 model, which is always unknown and often misspecified. We expect the GB method to overcome  
10 the practical difficulties of DCMs' applications.

11         Second, the GB method allows the use of flexible loss functions and priors. The  
12 uncertainty of the parameters expressed in the loss functions is easily demonstrated in the  
13 generalized posteriors generated using the GB method. In other words, the GB method can  
14 handle the amount of uncertainty of attribute mastery estimates. Not only point estimates but also  
15 uncertainty variation is important for careful decisions of diagnostic evaluation. The GB method  
16 provides a useful tool for addressing the above problems, which both parametric and  
17 nonparametric methods have. Furthermore, the generalized posterior is easily obtained using a  
18 Markov chain Monte Carlo (MCMC) routine, such as the Metropolis-Hastings method. Third,  
19 we can control the relative importance between the dataset and the prior via the learning rate  
20 parameter. If the obtained data's quality is questionable, an inference that is completely  
21 dependent on the data may lead to inappropriate decisions. In such cases, the data's relative  
22 importance can be reduced. The learning rate parameter enables a more flexible inference.

23         Based on these discussions, we develop a GB method to overcome the limitations of the  
24 loss function-based estimation method for DCMs (C. Ma, de la Torre, et al., 2023). The  
25 remainder of this paper is organized as follows: The second section demonstrates the basic setup  
26 of the DCMs and the previous loss function-based estimation method. The third section provides  
27 the GB method's fundamentals and its application to DCMs based on their loss functions.

1 Therein, the MCMC algorithm for a generalized posterior is also discussed. The GB method's  
 2 mathematical properties under certain conditions are discussed in the fourth section. The fifth  
 3 and sixth sections comprise simulation and real data analysis examples of generalized Bayesian  
 4 inference for DCMs, wherein we compare previous nonparametric estimation methods in a  
 5 simulation study. Finally, the seventh section serves as the discussion, where the limitations of  
 6 the GB inference and future directions of DCMs' estimation methods are discussed.

## 7 **2. Model Setup and Previous Estimation Methods**

### 8 **2.1. Model Setup of DCMs**

9 First, we express an individual's attribute mastery pattern using a vector of length  
 10  $K$ ,  $\alpha_i \in \{0,1\}^K$ , where  $i \in \{1,2, \dots, I\}$ . The  $k$ -th element of the attribute mastery pattern vector  
 11  $\alpha_i$  is  $\alpha_{ik} \in \{0,1\}$ , where  $k \in \{1,2, \dots, K\}$ ; it takes one if individual  $i$  masters attribute  $k$ , and  
 12 otherwise, it takes 0. In this study, we assume unconditional attribute mastery patterns, where all  
 13 possible attribute mastery patterns and the number is  $L = 2^K$ . Therefore, the  $l \in \{1,2, \dots, L\}$ -th  
 14 attribute mastery pattern can be written as  $\alpha_l$ . The set of attribute mastery patterns for all  
 15 individuals is  $\mathcal{A} = \{\alpha_i\}_{i=1}^I$ . To define the parametric measurement model, we also need to  
 16 specify the diagnostic relationship between the attributes and item sets.

17 The diagnostic relationship between the attributes and test items is called the  $\mathbf{q}$ -vector,  
 18  $\mathbf{q}_j \in \{0,1\}^K \setminus \{\mathbf{0}_K\}$ , where  $j \in \{1,2, \dots, J\}$ ; if the  $k$ -th attribute is required for item  $j$ ,  $q_{jk} = 1$ ;  
 19 otherwise  $q_{jk} = 0$ . Additionally,  $\mathbf{0}_K$  is a vector of length  $K$  and all its elements are 0. Here  
 20 we assume there is no item with  $\mathbf{q}_j = \mathbf{0}_K$ . The Q-matrix (Tatsuoka, 1985) is a  $J \times K$  matrix  
 21 defined by  $(\mathbf{q}_1^\top, \mathbf{q}_2^\top, \dots, \mathbf{q}_J^\top)^\top$ .

22 Parametric DCMs define their measurement models using attribute mastery patterns and  
 23 a  $\mathbf{q}$ -vector. For example, one of the most general DCMs, known as the log linear cognitive  
 24 diagnostic model (LCDM; Henson et al., 2009) uses an item parameter vector  $\lambda_j =$   
 25  $(\lambda_{j0}, \lambda_{j1}, \dots, \lambda_{j12\dots K})^\top$ , and the measurement model of  $X_{ij} = 1$ , which is a conditional response  
 26 probability of individual  $i$  for item  $j$ , is

$$P(X_{ij} = 1 \mid \lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l) = \frac{\exp(f(\lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i))}{1 + \exp(f(\lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i))} \quad (1)$$

1 where  $f(\lambda_j, \boldsymbol{\alpha}_i)$  is:

$$\begin{aligned} f(\lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i) &= \log \frac{P(X_{ij} = 1 \mid \lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l)}{1 - P(X_{ij} = 1 \mid \lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l)} \\ &= \lambda_{j0} + \sum_{k=1}^K \lambda_{jk} q_{jk} \alpha_{ik} + \sum_{k=1}^K \sum_{k' < k} \lambda_{jk k'} q_{jk} q_{jk'} \alpha_{ik} \alpha_{ik'} + \cdots + \lambda_{j1 \dots K} \prod_{k=1}^K q_{jk} \alpha_{ik}. \end{aligned} \quad (2)$$

2 The LCDM has several parameters. The first parameter is the intercept  $\lambda_{j0}$ , which determines  
 3 the baseline correct item response probability. Attribute mastery patterns that do not master any  
 4 attributes requiring item  $j$  take response probability. The main effect parameters were  
 5  $\lambda_{j1}, \lambda_{j2}, \dots, \lambda_{jK}$ ; each parameter affected the correct item response probabilities with the  
 6 corresponding attributes. The first-order interaction parameter  $\lambda_{jk k'}$  is the effect of  
 7 simultaneously mastering the attributes  $k$  and  $k'$ . Similarly, we introduce the highest  
 8 interaction term as  $\lambda_{j12 \dots K}$ . General cognitive diagnosis models that are similar to LCDM have  
 9 also been proposed in the literature, including the generalized DINA (GDINA) model (de la  
 10 Torre, 2011) and general diagnostic model (GDM; von Davier, 2008).

11 As some attributes are not measured by item  $j$ , the number of estimated item  
 12 parameters under LCDM is  $2^{\sum_k q_{jk}} \leq 2^K$ . Moreover, notably, one-to-one mapping exists  
 13 between the LCDM item parameters and conditional item response probabilities (Rupp et al.,  
 14 2010). Therefore, it is convenient to use conditional item response probabilities to develop  
 15 parameter-estimation methods. The same strategy was adopted in previous studies (Yamaguchi  
 16 & Okada, 2020; Yamaguchi & Templin, 2022b), where DCMs are a restricted version of latent  
 17 class models (e.g., Rupp & Templin, 2008; Xu & Shang, 2018).

18 Therefore, let the correct item response probability be parameter  $\theta_{jl}$ :

$$\theta_{j, \boldsymbol{\alpha}_l} = P(X_{ij} = 1 \mid \lambda_j, \mathbf{q}_j, \boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l), \quad (3)$$



1 Additionally, the attribute mastery mixing parameters  $\pi_{\alpha_1}, \pi_{\alpha_2}, \dots, \pi_{\alpha_L} \in (0,1)$  are defined as  
 2  $\pi_{\alpha_l} = P(\alpha_i = \alpha_l)$ , satisfying  $\sum_l \pi_{\alpha_l} = 1$ . From this notation, the complete data likelihood  
 3 function of the LCDM is

$$\mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} | X) = \prod_{i=1}^I \prod_{j=1}^J \prod_{l=1}^L \left\{ \pi_l \theta_{j, \alpha_l}^{x_{ij}} (1 - \theta_{j, \alpha_l})^{1-x_{ij}} \right\}^{\mathcal{I}(\alpha_i = \alpha_l)}, \quad (4)$$

4 where  $X = \{x_{ij}\}_{i,j=1}^{N,J}$ ,  $\Theta = \{\theta_{jl}\}_{j,l=1}^{J,L}$ ,  $\boldsymbol{\pi} = (\pi_{\alpha_1}, \pi_{\alpha_2}, \dots, \pi_{\alpha_L})^\top$ , and  $\mathcal{I}(\cdot)$  is an indicator function.

5 We add some remarks on the correct item response probabilities for item  $j$ . First, as  
 6 mentioned previously, some attribute mastery patterns have the same item response probabilities  
 7 because of the setting of the  $\mathbf{q}$  vector. Moreover, some sub-models of the LCDM assume fewer  
 8 parameters than the general LCDM and have parsimonious model forms. The model settings for  
 9 each item can differ, but we assume that all test items have the same general LCDM form.

10 Second, the correct item response probabilities for item  $j$  exhibit an ordinal  
 11 relationship: These relationships are known as monotonicity constraints (Xu & Shang, 2018).  
 12 The formal expression of the monotonicity constraints proposed by Xu and Shang (2018) is

$$\max_{\boldsymbol{\alpha}: \boldsymbol{\alpha} \succcurlyeq \mathbf{q}_j} \theta_{j, \alpha} = \min_{\boldsymbol{\alpha}: \boldsymbol{\alpha} \preccurlyeq \mathbf{q}_j} \theta_{j, \alpha} \geq \theta_{j, \alpha'} \geq \theta_{j, \mathbf{0}_K}, \quad (5)$$

13 where we write  $\boldsymbol{\alpha} \succcurlyeq \mathbf{q}_j$  if  $\alpha_k \geq q_{jk}, \forall k$ ; otherwise,  $\boldsymbol{\alpha} \not\succeq \mathbf{q}_j$ . These constraints imply that the  
 14 patterns mastering all skills measured in item  $j$  ( $\boldsymbol{\alpha}: \boldsymbol{\alpha} \succcurlyeq \mathbf{q}_j$ ) should have the highest of all the  
 15 patterns. By contrast, all non-mastering patterns had the lowest correct item response probability.  
 16 The middle mastering patterns satisfying  $\boldsymbol{\alpha}' \not\succeq \mathbf{q}_j$  have response probabilities between these  
 17 two probabilities.

## 18 2.2. Loss Function-Based Parameter Estimation

19 This section introduces the loss function-based parameter estimation method proposed in C. Ma,  
 20 de la Torre, et al. (2023). First, we describe certain elements of the loss function-based method.

21 In this framework, we introduce the length  $J$  centroid parameter vector:  $\boldsymbol{\mu}_{\alpha_l} =$   
 22  $(\mu_{1, \alpha_l}, \mu_{2, \alpha_l}, \dots, \mu_{J, \alpha_l})^\top \in \mathbb{R}^J$ . Additionally, a penalty term for the mixing parameter  $\pi_{\alpha_l}$  is

1 introduced as  $h(\pi_{\alpha_l}) \in \mathbb{R}$ . Furthermore, an element-wise loss function taking item response  
 2 vector  $\mathbf{x}_i$  and a centroid parameter vector  $\boldsymbol{\mu}_{\alpha_l}$  is expressed as  $\ell(\mathbf{x}_i, \boldsymbol{\mu}_{\alpha_l})$ ; its codomain is real  
 3 positive number  $\mathbb{R}^+$ . The  $\ell(\mathbf{x}_i, \boldsymbol{\mu}_{\alpha_l})$  is the individual-level loss function. Therefore, the loss  
 4 function of the entire dataset is based on the individual-level loss function:

$$\mathcal{L}(\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}) = \sum_{l=1}^L \sum_{i:\alpha_i=\alpha_l} \{\ell(\mathbf{x}_i, \boldsymbol{\mu}_{\alpha_l}) + h(\pi_{\alpha_l})\}. \quad (6)$$

5 The second summation takes over the individuals with the attribute mastery pattern  $\alpha_l$ .

6 Parameter estimates are obtained to minimize the loss function defined above:

$$\{\hat{\mathcal{A}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\pi}}\} = \underset{\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}}{\operatorname{argmin}} \mathcal{L}(\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}). \quad (7)$$

7 Directly minimizing the above loss function is not easy; therefore, we use the iterative update  
 8 rule instead. In the estimation algorithm, we first set initial parameters  $\{\boldsymbol{\mu}^{(0)}, \boldsymbol{\pi}^{(0)}\}$ . When we  
 9 have parameter estimates at  $t$ -th iteration,  $\{\boldsymbol{\mu}^{(t)}, \boldsymbol{\pi}^{(t)}\}$ , the following update steps are repeated:

$$\begin{aligned} \text{Step 1: } \{\mathcal{A}^{(t+1)}\} &= \underset{\mathcal{A}}{\operatorname{argmin}} \mathcal{L}(\mathcal{A}, \boldsymbol{\mu}^{(t)}, \boldsymbol{\pi}^{(t)}), \\ \text{Step 2: } \{\boldsymbol{\mu}^{(t+1)}, \boldsymbol{\pi}^{(t+1)}\} &= \underset{\boldsymbol{\mu}, \boldsymbol{\pi}}{\operatorname{argmin}} \mathcal{L}(\mathcal{A}^{(t+1)}, \boldsymbol{\mu}, \boldsymbol{\pi}). \end{aligned} \quad (8)$$

10 If the predetermined convergence criterion is satisfied, for example  $\epsilon > 1 - \sum_i \left( \mathcal{J}(\boldsymbol{\alpha}_i^{(t+1)} = \right.$   
 11  $\left. \boldsymbol{\alpha}_i^{(t)}) \right) / I, 0 < \epsilon < 1$ , the update process is stopped, and the parameter estimates become output:

$$12 \quad \{\hat{\mathcal{A}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\pi}}\} = \{\mathcal{A}^{(t+1)}, \boldsymbol{\mu}^{(t+1)}, \boldsymbol{\pi}^{(t+1)}\}.$$

13 Many previous estimation methods can be viewed as special cases of the general loss  
 14 function formulation framework. The joint likelihood estimation of the parametric DCM is a first  
 15 example. In the following, we focus on the deterministic inputs noisy, “and” gate model (DINA  
 16 model; Junker & Sijtsma, 2001; MacReady & Dayton, 1977; Maris, 1999) as an example  
 17 because it is well-known and considered the most parsimonious DCM. The loss function further  
 18 used to obtain the MAP estimation, which is the negative of the sum of the log-likelihood and

1 log-prior density functions, is also presented. Subsequently, a nonparametric classification  
 2 method (NPC; Chiu & Douglas, 2013; Wang & Douglas, 2015) and generalized NPC (GNPC;  
 3 Chiu & Köhn, 2019; Chiu et al., 2018) are formulated under the above framework. Furthermore,  
 4 NPC and GNPC are extended to the GB framework in a later section.

5 The DINA model is the simplest and most fundamental DCM, which is a special case of  
 6 the LCDM. The DINA model assumes only the intercept and the highest interaction terms of the  
 7 LCDM item parameters. Let the subscript set of attributes measured by item  $j$  be  $\mathcal{K} =$   
 8  $\{k; q_{jk} = 1, k = 1, 2, \dots, K\}$  and let the LCDM kernel for the DINA model be reduced to

$$f(\boldsymbol{\lambda}_j, \mathbf{q}_j, \boldsymbol{\alpha}_i) = \lambda_{j0} + \lambda_{j\mathcal{K}} \prod_{k \in \mathcal{K}} \alpha_{ik}. \quad (9)$$

9 In the conventional DINA formulation, two-item response probabilities are represented by  
 10 estimating the  $g_j$  and slipping  $s_j$  parameters:

$$g_j = \frac{\exp(\lambda_{j0})}{1 + \exp(\lambda_{j0})}, \quad (10)$$

$$1 - s_j = \frac{\exp(\lambda_{j0} + \lambda_{j\mathcal{K}} \prod_{k \in \mathcal{K}} \alpha_{ik})}{1 + \exp(\lambda_{j0} + \lambda_{j\mathcal{K}} \prod_{k \in \mathcal{K}} \alpha_{ik})}. \quad (11)$$

11 The guessing parameter  $g_j$  indicates the chance level of a correct item response for attribute  
 12 mastery patterns that lack at least one attribute required by item  $j$ . The slipping parameter  $s_j$  is  
 13 the incorrect response probability of all mastering-attribute mastery patterns required by item  $j$ .  
 14 Both  $g_j$  and  $s_j$  can be represented as functions of the ideal responses

$$\eta_j^{DINA}(\boldsymbol{\alpha}_l) = \prod_{k=1}^K \alpha_{lk}^{q_{jk}}. \quad (12)$$

15 The ideal response represents the response of an individual who belongs to the  $l$ -th attribute  
 16 mastery pattern for item  $j$  without errors. Then, the  $g_j$  and  $s_j$  are represented as conditional  
 17 probabilities:

$$g_j = P(X_j = 1 \mid \eta_j^{DINA}(\boldsymbol{\alpha}_l) = 0), \quad (13)$$

$$s_j = P(X_j = 0 \mid \eta_j^{DINA}(\boldsymbol{\alpha}_l) = 1). \quad (14)$$

1 Using the item response probabilities, the centroid parameter under the DINA model is

$$\theta_{jl} = g_j^{1-\eta_j^{DINA}(\alpha_l)} (1 - s_j)^{\eta_j^{DINA}(\alpha_l)}. \quad (15)$$

2 Assuming a cross-entropy loss, which is  $-\log y$ , the likelihood-based loss function for the  
3 DINA model is

$$\mathcal{L}(\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}) = - \sum_{i=1}^I \sum_{l=1}^L \mathcal{J}(\alpha_i = \alpha_l) \left[ \sum_{j=1}^J \{x_{ij} \log \theta_{jl} + (1 - x_{ij}) \log(1 - \theta_{jl})\} + \log \pi_{\alpha_l} \right]. \quad (16)$$

4 We assume  $h(\pi_{\alpha_l}) = -\log \pi_{\alpha_l}$ . The loss function defined in Equation 16 is equivalent to a  
5 negative log complete likelihood function. Therefore, minimizing equation 16 corresponds to  
6 maximizing the likelihood function; the minimizers  $\{\hat{\mathcal{A}}, \hat{\boldsymbol{\mu}}, \hat{\boldsymbol{\pi}}\}$  can be considered as the  
7 maximum likelihood estimate.

8 Subsequently, we examine the NPC and GNPC methods. Following Chiu and Douglas  
9 (2013), the loss function in the NPC method is defined by the Hamming distance between the  
10 individual item response vector and the ideal response vector:

$$\ell(\mathbf{x}_i, \boldsymbol{\mu}_{\alpha_l}) = \sum_{j=1}^J \ell(x_{ij}, \mu_{j,\alpha_l}) = \sum_{j=1}^J |x_{ij} - \eta_j^{DINA}(\alpha_l)|. \quad (17)$$

11 In the NPC method, the centroid parameter is the ideal response  $\mu_{j,\alpha_l}$ . The NPC estimates are  
12 obtained to minimize Equation 17 for each individual:

$$\hat{\alpha}_i = \underset{\alpha_l}{\operatorname{argmin}} \sum_{j=1}^J |x_{ij} - \eta_j^{DINA}(\alpha_l)|, \forall i. \quad (18)$$

13 Clearly, the Hamming distance is a loss function, and the NPC method is a loss function-based  
14 estimation method.

15 As introduced in Chiu et al. (2018), the GNPC is a type of generalization that employs  
16 DINA- and deterministic inputs noisy, “or” gate (DINO; Templin & Henson, 2006)-type ideal  
17 responses to define a generalized ideal response. The DINO-type ideal response is

$$\eta_j^{DINO}(\alpha_l) = 1 - \prod_{k=1}^K (1 - \alpha_{lk})^{q_{jk}}, \quad (19)$$

1 and  $\eta_j^{DINO}(\alpha_l)$  becomes one if pattern  $l$  masters at least one attribute required for item  $j$ ;  
 2 otherwise, it becomes 0. The generalized ideal response is then defined as

$$\eta_j^{(w)}(\alpha_l) = w_{jl}\eta_j^{DINA}(\alpha_l) + (1 - w_{jl})\eta_j^{DINO}(\alpha_l), \quad (20)$$

3 where  $w_{jl} \in [0,1]$  is a weight parameter that determines an item's tendency. If the item is more  
 4 like DINA or conjunctive,  $w_{jl}$  is close to one. By contrast, a  $w_{jl}$  near zero means that the item  
 5 is DINO-like or disjunctive in nature. The GNPC assumes a Euclidean distance for its loss  
 6 function

$$d(x_j, \boldsymbol{\eta}^{(w)}(\alpha_l)) = \sum_{i=1}^I \mathcal{J}(\alpha_i = \alpha_l) (x_{ij} - \eta_j^{(w)}(\alpha_l))^2, \quad (21)$$

7 where  $\boldsymbol{\eta}^{(w)}(\alpha_l) = (\eta_1^{(w)}(\alpha_l), \eta_2^{(w)}(\alpha_l), \dots, \eta_j^{(w)}(\alpha_l))^T$ . The weight parameter is estimated via  
 8  $\hat{w}_{jl} = 1 - \sum_{i=1}^I \mathcal{J}(\alpha_i = \alpha_l) x_{ij} / \sum_{i=1}^I \mathcal{J}(\alpha_i = \alpha_l)$ . The loss function of the GNPC is

$$\mathcal{L}(\{\mathcal{A}, W\}) = \sum_{j=1}^J \sum_{l=1}^L d(x_j, \boldsymbol{\eta}^{(w)}(\alpha_l)) = \sum_{j=1}^J \sum_{l=1}^L \sum_{i=1}^I \mathcal{J}(\alpha_i = \alpha_l) (x_{ij} - \eta_j^{(w)}(\alpha_l))^2, \quad (22)$$

9 where  $W = \{w_{jl}\}_{j,l=1}^{J,L}$ . The GNPC requires iterative updates of weight  $w_{jl}$  and attribute  
 10 mastery patterns. The detailed update rule is described in Chiu et al. (2018). Note that if  $\eta_j^{DINA}$   
 11 and  $\eta_j^{DINO}$  are not distinguished for some items and attribute patterns, the weight value is fixed  
 12 to a value close to zero or one. See Chiu et al. (2018) for a detailed discussion.

13 As demonstrated above, parametric and nonparametric estimation methods can be  
 14 treated in a unified loss function-based framework (C. Ma, de la Torre, et al., 2023). However,  
 15 these loss function-based parameter estimates usually only provide point estimates, and  
 16 uncertainty quantification in the parameter estimates has been considered less serious.  
 17 Furthermore, different specifications of the measurement model precipitate significantly different  
 18 attribute mastery patterns (e.g., Li et al., 2016). However, assessing all possible measurement

1 models for all test items may be difficult. The GNPC is a promising estimation method that can  
 2 be used in varied situations, even when the measurement model is unknown. However, prior  
 3 knowledge of the weight parameters in the GNPC is often not considered. These problems can be  
 4 solved using the generalized Bayesian method introduced in the following section.

### 5 **3. Generalized Bayesian Method for DCMs**

#### 6 **3.1. Construction of the Generalized Posterior**

7 The GB method is a decision theory under a model misspecification situation (Bissiri  
 8 et al., 2016). In other words, the assumed model may not accurately represent the true data-  
 9 generating process, or the relationship between the model parameters and data may not be  
 10 described via the assumed model, which is known as the  $\mathcal{M}$ -open situation (Bissiri et al., 2016,  
 11 p. 1111). The GB method is a coherent belief update procedure that uses a loss function even in  
 12 the  $\mathcal{M}$ -open situation. Thus, the GB method extends the applicability of the typical Bayesian  
 13 methods, which require a likelihood function.

14 Let datasets and parameter sets be  $\mathbf{y}$  and  $\Theta$ , respectively. Additionally, the loss  
 15 function and prior distribution are  $\ell(\mathbf{y}; \Theta)$  and  $p(\Theta)$ . Then, the generalized posterior of the  
 16 parameter  $p(\Theta | \mathbf{y})$  is

$$p(\Theta | \mathbf{y}) \propto \exp(-\omega \ell(\mathbf{y}; \Theta)) p(\Theta), \quad (23)$$

17 where  $\omega$  is called the learning rate, a tuning parameter that controls a dataset's importance.  
 18 Methods for determining the learning rate are still being studied (Wu & Martin, 2023); notably,  
 19 no standard has been established thus far. The generalized posterior function is the result of  
 20 updating the prior distribution based on the loss function. If we select a negative log-likelihood  
 21 function for the loss function and  $\omega = 1$ , the generalized posterior becomes the usual Bayesian  
 22 posterior function.

23 Bissiri et al. (2016, pp. 1106-1107) discusses some of the validity requirements for loss  
 24 functions. First, the solution of the loss function must exist. Second, the loss function must  
 25 satisfy the following condition:

$$0 < \int \exp(-\ell(\mathbf{y}; \boldsymbol{\Theta})) p(\boldsymbol{\Theta}) d\boldsymbol{\Theta} < \infty. \quad (24)$$

1 Some major loss functions considered in this study, such as the Hamming distance, Euclid  
 2 distance, or cross-entropy loss, satisfy the above integral conditions. Additionally, Bissiri et al.  
 3 (2016, p. 1107) identify natural assumptions for deriving a generalized posterior from a loss  
 4 function. We should also point out that we only need to construct loss functions given a set of  
 5 data for only the parameter of interest to employ the GB method. In our manuscript, the GB  
 6 method employed the loss function for attribute mastery patterns. The loss function is based on  
 7 the GNPC: quadratic of Euclid distance. It satisfies the above conditions and is valid.

### 8 **3.3 General Form of the Generalized Bayesian Method for DCMs**

9 The general form of the GB method for DCMs can be expressed using Equations 6 and  
 10 23;

$$p(\{\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}\} | X) \propto \exp(-\omega \{\mathcal{L}(\{\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}\})\}) p(\boldsymbol{\mu}) p(\boldsymbol{\pi}),$$

$$\propto \exp\left(-\omega \sum_{l=1}^L \sum_{i:\alpha_i=\alpha_l} \{\ell(\mathbf{x}_i, \boldsymbol{\mu}_{\alpha_l}) + h(\pi_{\alpha_l})\}\right) p(\boldsymbol{\mu}) p(\boldsymbol{\pi}). \quad (25)$$

11 The penalty term was  $h(\pi_{\alpha_l}) = -\log \pi_{\alpha_l}$ .

12 Using the GNPC loss function defined in Equation 22 and adding a penalty term for the  
 13 mixing parameter  $\boldsymbol{\pi}$ , the generalized posterior is

$$p(\{\mathcal{A}, \boldsymbol{\mu}, \boldsymbol{\pi}\} | X) \propto \exp\left(-\omega \sum_{l=1}^L \sum_{i=1}^I \left\{ \mathcal{J}(\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l) \left[ \sum_{j=1}^J (x_{ij} - \eta_j^{(w)}(\boldsymbol{\alpha}_l))^2 \right] - \log \pi_{\alpha_l} \right\}\right) p(W) p(\boldsymbol{\pi}). \quad (26)$$

14 Notably, we treat weight  $W$  as a parameter and assume a prior instead of a centroid parameter  
 15  $\boldsymbol{\mu}$  because the centroid parameter  $\boldsymbol{\eta}^{(w)}(\boldsymbol{\alpha}_l)$  is determined by two ideal responses and weight  
 16 parameters; thus, it is natural. Priors for the mixing parameters and weight parameters are  
 17 assumed Dirichlet and Beta distributions:

$$p(\boldsymbol{\pi}) \propto \prod_{l=1}^L \pi_l^{\delta_l^0 - 1}, \quad (27)$$

$$p(W) \propto \prod_{j=1}^J \prod_{l=1}^L w_{jl}^{a_{jl}^0 - 1} (1 - w_{jl})^{b_{jl}^0 - 1}, \quad (28)$$

1 where  $\delta_1^0, \delta_2^0, \dots, \delta_L^0 \geq 0, \sum_l \delta_l^0 = 1$  and  $a_{jl}^0, b_{jl}^0 \geq 0$ .

2 The posterior was numerically obtained using MCMC techniques, such as Metropolis-  
3 Hastings within the Gibbs sampling method, or MCMC software, such as JAGS (Plummer,  
4 2003) or Stan (Carpenter et al., 2017). The conditional distribution of  $\boldsymbol{\alpha}_i$  is categorical:

$$p(\boldsymbol{\alpha}_i | \mathbf{x}_i, W, \boldsymbol{\pi}) \propto \prod_{l=1}^L r_{il}^{\mathcal{J}(\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l)},$$

$$r_{il} = \frac{\rho_{il}}{\sum_l \rho_{il}}, \quad (29)$$

$$\rho_{il} = \exp \left( -\omega \mathcal{J}(\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l) \left[ \sum_{j=1}^J (x_{ij} - \eta_j^{(w)}(\boldsymbol{\alpha}_l))^2 \right] - \log \pi_{\boldsymbol{\alpha}_l} \right).$$

5 The conditional distribution of the mixing parameters was a Dirichlet distribution:

$$p(\boldsymbol{\pi} | X) \propto \prod_{l=1}^L \pi_l^{\delta_l^* - 1},$$

$$\delta_l^* = \omega \sum_{i=1}^I \mathcal{J}(\boldsymbol{\alpha}_i = \boldsymbol{\alpha}_l) + \delta_l^0. \quad (30)$$

6 The conditional distribution of the weight parameter is not easily expressed; therefore, its  
7 MCMC update was performed using the Metropolis-Hastings method. The candidate was  
8 generated by a random walk using a uniform distribution:  $w_{jl}^{(\text{cand})} = w_{jl}^{(\text{now})} + u, u \sim$   
9  $\text{Unif}(-0.05, 0.05)$ . Using the above distributions and updating rules, the MCMC for a



1 generalized posterior is numerically approximated as follows: The mixing and weight parameters  
 2 were initialized as  $\{\boldsymbol{\pi}^{(0)}, W^{(0)}\}$  and the hyperparameters were set to  $\delta_l^0 = 1, \forall l$  and  $a_{jl}^0 =$   
 3  $2, b_{jl}^0 = 1, \forall j, l$  for example. Then, at the  $t$ -th MCMC iteration ( $t = 1, 2, \dots, T \in \mathbb{N}$ ),  $\boldsymbol{\alpha}_i^{(t+1)}$  is  
 4 generated from the categorical distribution expressed in Equation 29 with the  $t$ -th MCMC  
 5 sample of the parameter set at  $\{\boldsymbol{\alpha}^{(t)}, W^{(t)}\}$ . The  $(t + 1)$ -th MCMC sample of the mixing  
 6 parameter is generated from the Dirichlet distribution shown in Equation 30 using  $\mathcal{A}^{(t+1)}$ .  
 7  $W^{(t+1)}$  is obtained using the Metropolis-Hastings method.

8 Under the hyperparameter setting, the Dirichlet distribution for mixing parameters  
 9 becomes a uniform distribution. This represents a scenario in which we have no information  
 10 about the population attribute mastery ratio. In this means, the prior of the attribute mastery  
 11 pattern has almost no information. The mean and SD of the prior of the weight parameter are  
 12 0.667 and 0.236, respectively. Under this setting, interval  $[0.158, 0.987]$  covers 95% of the  
 13 support of the parameter. The data analyst expected the items in the test to have a slightly  
 14 conjunctive nature, which means the items behave more like in the DINA model than in the  
 15 DINO model. However, the expectation is not particularly strong because the interval covering  
 16 95% of the support of the parameter is wide. This interpretation indicates that the prior conveys  
 17 some information about the weight parameters.

#### 18 **4. Mathematical Properties of the Proposed Method: Consistency of the Maximum a** 19 **Posteriori Estimators**

20 First, we formally introduce the estimators under the GB framework and subsequently  
 21 discuss their statistical behaviors under certain conditions. The Appendix provides the full  
 22 proofs. In this work, we assume that the item responses were generated from the Bernoulli  
 23 distribution with parameter  $\Theta$  defined by Equation 3; the attribute mastery patterns were  
 24 generated from a categorical distribution with a mixing parameter  $\boldsymbol{\pi}$ . Although several  
 25 alternatives exist, MAP estimation provides a relatively natural and simple choice. Furthermore,  
 26 MAP estimators of the GB method  $(\hat{\mathcal{A}}, \hat{\Theta}, \hat{\boldsymbol{\pi}})$  are estimators of the true parameters

1  $(\mathcal{A}^0, \Theta^0, \boldsymbol{\pi}^0)$  in the data-generating process. These are obtained by minimizing the loss function  
 2 of  $(\mathcal{A}, \Theta, \boldsymbol{\pi})$  under the constraint imposed by the Q matrix, as follows:

$$\mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} | X) = \sum_{i=1}^I \left( \sum_{j=1}^J \ell(X_{ij}, \theta_{j,\alpha_i}) + h(\boldsymbol{\pi}_{\alpha_i}) \right) + \sum_{j,l} \log f_{j,\alpha_l}(\theta_{j,\alpha_l}) + \sum_{l=1}^L \log g_{\alpha_l}(\boldsymbol{\pi}_{\alpha_l}), \quad (31)$$

3 where  $h(\cdot)$  is a continuous nonincreasing regularization function of the proportion parameters  
 4  $\boldsymbol{\pi}_{\alpha}$ , often taken as  $h(\boldsymbol{\pi}) = -\log \boldsymbol{\pi}$ ;  $f_{j,\alpha}$  and  $g_{\alpha}$  are the prior density functions of  $\theta_{j,\alpha}$  and  $\boldsymbol{\pi}_{\alpha}$ ,  
 5 respectively. Note that we consider a model sequence indexed by  $(I, J)$ , where both  $I$  and  $J$   
 6 tend to infinity, while  $K$  is held constant.

7 Several regularity conditions are required to ensure the consistency of MAP estimators.  
 8 The first assumption is as follows.

9 **Assumption 1.** *There exists  $\delta_1, \delta_2 > 0$  such that*

$$10 \quad \min_{1 \leq j \leq J} \left\{ \min_{\alpha_l \circ q_j^0 \neq \alpha_{l'} \circ q_j^0} (\theta_{j,\alpha_l}^0 - \theta_{j,\alpha_{l'}}^0)^2 \right\} \geq \delta_1,$$

$$11 \quad \text{and } \delta_2 \leq \min_{j,\alpha} \theta_{j,\alpha}^0 < \max_{j,\alpha} \theta_{j,\alpha}^0 \leq 1 - \delta_2.$$

12 The first condition in Assumption 1 serves as an identification condition for local latent classes  
 13 at each item level. The gap denoted by  $\delta$  measures the separation between the latent classes,  
 14 thereby quantifying the signals' strength. The second condition in Assumption 1 keeps the true  
 15 parameters away from the boundaries of the parameter space to prevent unusual behaviors of the  
 16 element-wise loss.

17 Assumption 2 pertains to the discrete structures of  $\mathbf{Q}$  and is expressed as the following.

18 **Assumption 2.** *All proportion parameters  $\boldsymbol{\pi}_{\alpha}$  are strictly greater than zero, and there exist*  
 19  $\{\delta_j\} \subset (0, \infty)$  *such that*

$$\min_{1 \leq k \leq K} \frac{1}{J} \sum_{j=1}^J \mathcal{J}\{\mathbf{q}_j^0 = \mathbf{e}_k\} \geq \delta_J. \quad (32)$$

1 This assumption holds that  $\mathbf{Q}$  includes an increasing number of identity submatrices,  $\mathbf{I}_K$ , as  $J$   
 2 grows. Notably, by attaching the subscript  $J$  to the lower bound (32) in Assumption (2), we  
 3 allow it to decrease to zero as  $J$  approaches infinity. As the following theorems show, if the rate  
 4 at which variable  $\delta_J$  decreases meets certain mild requirements, the consistency of  $(\hat{\mathcal{A}}, \hat{\Theta})$  can  
 5 be ensured.

6 The subsequent assumption concerns the element-wise loss function  $\ell$ .

7 **Assumption 3.** *The loss function  $\ell(X, \theta)$  is twice continuously differentiable in  $\theta$  on  $(0,1)$   
 8 and  $\exists b_L > b_U > 0$  such that  $b_L \leq \partial_{\theta^2} \ell(R, \theta) \leq b_U$  for  $\theta$  in a compact subset of  $(0,1)$ . The  
 9 total loss (31) is minimized at class means given the subjects' membership, as in,  $\hat{\theta}_{j,\alpha} =$   
 10  $\sum_{i=1}^I \mathcal{J}\{\hat{\alpha}_i = \alpha\} X_{ij} / \sum_{i=1}^I \mathcal{J}\{\hat{\alpha}_i = \alpha\}$ .*

11 Assumption 3 imposes smoothness conditions on the element-wise loss function, rendering it  
 12 convex. The upper bound of the second derivative is necessary to control the remaining term in  
 13 the expansion of the first-order condition, and the lower bound allows us to quantify the  
 14 estimator drift caused by the given priors. For the sample average assumption, we can verify that  
 15 both  $\ell^2$  and cross-entropy loss functions satisfy Assumption 3.

16 Assumption 4 states that the true parameters minimize the element-wise loss functions  
 17 and quantify the deviations when  $\theta$  is not a true parameter. This assumption is expressed as  
 18 follows:

19 Assumption 4. *There exist constants  $\eta \geq 2, c > 0$  such that*

$$\mathbb{E}[\ell(X_{ij}, \theta)] - \mathbb{E}[\ell(X_{ij}, \theta_{j,\alpha_i}^0)] \geq c |\theta - \theta_{j,\alpha_i}^0|^\eta. \quad (33)$$

1 Assumption 4 holds for both the  $\ell^2$  loss and the cross-entropy loss.

2 Assumption 5 is a technical assumption that allows us to control the effects of prior  
3 distributions on the estimators.

4 **Assumption 5.**  $h(\cdot)$  in (31) is a continuous nonincreasing function of the proportion  
5 parameters, and  $C > c > 0$  exists such that for any  $j$  and  $\alpha, C > f_{j,\alpha}, g_\alpha > c$  on a compact  
6 parameter subspace of  $(0,1)$ .

7 We can verify that the Dirichlet and Beta distributions satisfy this assumption.

8 Under the aforementioned regularity conditions, we demonstrate the consistency  
9 properties of the GB method with constraints for different attribute mastery patterns  $\alpha_l$  and  
10  $\alpha_{l'}, l \neq l'$  (C. Ma, de la Torre, et al., 2023; Xu, 2017):

$$(\alpha_l \circ \mathbf{q}_j = \alpha_{l'} \circ \mathbf{q}_j) \implies (\theta_{j,\alpha_l} = \theta_{j,\alpha_{l'}}), \quad (34)$$

11 where  $\alpha \circ \mathbf{q}_j = (\alpha_1 \cdot q_{j1}, \dots, \alpha_K \cdot q_{jK})$  denotes the element-wise product of binary vectors  $\alpha$   
12 and  $\mathbf{q}_j$ . This implies that the item response parameter  $\theta_{j,\alpha}$  depends only on whether the  
13 attribute mastery pattern  $\alpha$  contains the required attributes  $\mathcal{K}_j := \{k \in [K]; q_{jk} = 1\}$  for item  
14  $j$ .

15 Based on the above five assumptions, we can derive consistent results for the GB  
16 method. The following main theorem first validates the clustering consistency of the GB method  
17 under the constraint (34), providing a bound for its convergence rate in recovering the attribute  
18 mastery patterns.

19 **Theorem 1** (Clustering Consistency). Consider  $(\hat{\mathcal{A}}, \hat{\Theta}, \hat{\boldsymbol{\pi}}) = \arg \min_{(\mathcal{A}, \Theta, \boldsymbol{\pi})} \mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} | X)$   
20 under the constraint (34). When  $I, J \rightarrow \infty$  jointly, suppose  $\sqrt{J} = O(I^{1-c})$  for some small  
21 constant  $c \in (0,1)$ . Under Assumption 1 to Assumption 5. the clustering error rate is

$$\frac{1}{I} \sum_{i=1}^I \mathcal{I}\{\hat{\alpha}_i \neq \alpha_i^0\} = o_p \left( \frac{(\log J)^{\xi/\eta}}{\delta_J(J)^{1/\eta}} \right), \quad (35)$$

1 where for a small positive constant  $\xi > 0$ .

2 Theorem 1 bounds the error of the estimator  $\hat{\mathcal{A}}$ , which establishes the clustering consistency of  
 3 the MAP estimators of the GB method, allowing the rate  $\delta_J$  to go to zero. Notably, the scaling  
 4 condition only assumes that  $J$  goes to infinity jointly with  $I$ , but at a slower rate.

5 The following result demonstrates that the MAP estimator of the item parameters can be  
 6 uniformly estimated consistently as  $I, J \rightarrow \infty$ :

7 **Theorem 2** (Item Parameters Consistency). *Under Assumptions 1 to 5 and the scaling conditions*  
 8 *given in Theorem 1, we have the following uniform consistency result for all  $j \in [J]$  and  $\alpha \in$*   
 9  *$\{0,1\}^K$ :*

$$\max_{j,\alpha} |\hat{\theta}_{j,\alpha} - \theta_{j,\alpha}^0| = o_p \left( \frac{1}{\sqrt{I^{1-\hat{c}}}} \right) + o_p \left( \frac{(\log J)^{\xi/\eta}}{\delta_J(J)^{1/\eta}} \right), \quad (36)$$

10 where  $\hat{c}$  and  $\xi$  are small positive constants.

11 On the first error term, the condition  $\pi_\alpha > 0$  for all  $\alpha \in \{0,1\}^K$  ensures that with probability  
 12 one, there are enough samples within each class to provide accurate estimates of item  
 13 parameters. Notably,  $\hat{c}$  the first error term arises because the number of parameters approaches  
 14 infinity jointly with the sample size  $I$ , which causes a slight deviation from the optimal error rate  
 15 of  $O_p(1/\sqrt{I})$ . The maximum deviation  $\max_{j,\alpha} |\hat{\theta}_{j,\alpha} - \theta_{j,\alpha}^0|$  is also affected by the classification  
 16 error. This is indicated in the second error term  $o_p((\log J)^\xi / \delta_J \sqrt{J})$ .

17 We can easily establish the consistency of the mixing parameter estimator  $\hat{\pi}$ . When  
 18  $h(\pi) = -\log \pi$ , the mixing parameters will be estimated as the sample average form  
 19  $\sum_i \mathcal{I}\{\alpha^0 = \alpha\} / I$ , which converge in probability to  $\pi_\alpha^0$  because of the clustering consistency.

20 **Corollary 1** (Proportion Parameters Consistency). *Under Assumptions 1 to 5 and the scaling*  
 21 *conditions given in Theorem 1, when  $h(\pi)$  is taken as  $-\log \pi$ , we have  $\hat{\pi}_\alpha \xrightarrow{P} \pi_\alpha^0$ .*

22

## 5. Simulation Study

This section compares the previous (G)NPC and the corresponding GB methods using the loss functions in NPC and GNPC, named as GBNPC and GBGNPC, respectively. This simulation study primarily aims to assess the behavior of the GB method's parameter estimates under finite small sample and item situations. As the GBNPC and GNPC are based on loss function in nonparametric methods, the most interesting parameters are attribute mastery patterns. In this simulation study, we mainly focus on the comparisons of the point estimates from these methods. To represent the uncertainty of the estimates, we also present attribute mastery probabilities using the GBGNPC and GBNPC methods, which indicate the benefit of the proposed method against the nonparametric methods.

The code for this simulation study is available on the Open Science Framework (OSF) webpage: <https://osf.io/sau6j/>.

### 5.1. Simulation Settings

Five factors are manipulated in the simulations. All factors had two conditions; hence,  $2^5 = 32$  simulation settings were used. The first factor was the data-generating model: DINA or general DCM (e.g., LCDM). The DINA model condition is a simpler data-generating situation, whereas the general DCM model is more complex. The second factor was the Q-matrix; four or five attributes are listed in Tables 1 and 2. Table 1 contains 19 items: eight simple items (i.e., measuring only one attribute), six items measuring two attributes, five items requiring three attributes, and the most complex item measuring all four attributes. Table 2 lists 30 items: eight simple items, ten items measuring two attributes, and ten items measuring three attributes.

Sample size was the third factor, with 30 or 300 participants assumed. The sample size setting of 30 participants mimicked classroom size. The sample size of 300 participants was 10 times larger than that of other classroom settings. The fourth condition was attribute correlation: independent ( $\rho = 0$ ) or highly correlated ( $\rho = 0.8$ ). The independent attribute condition was unrealistic but represented an ideal condition. The highly correlated condition was more realistic because the DCMs application indicated a high correlation among attributes (e.g., von Davier,

1 2008). The fifth condition was item quality. The high-item-quality condition indicates a high  
 2 description of all non-mastering attributes and all perfectly mastering attributes. In the high-item-  
 3 quality condition, the correct response probability of the all non-mastering pattern was 0.1 and  
 4 that of the all-mastering pattern was 0.9. On the other hand, in the low-item-quality condition,  
 5 the corresponding probabilities were 0.3 and 0.7. The correct item response probabilities of the  
 6 intermediate mastering patterns are generated based on Yamaguchi and Templin (2022b) or  
 7 Yamaguchi and Templin (2022a).

8           The data generation process used herein was similar to those in previous studies, such  
 9 as Chiu and Douglas (2013), Yamaguchi and Templin (2022b), and Yamaguchi and Templin  
 10 (2022a). First, for each individual, we generated a continuous latent variable vector  $\hat{\alpha}_i =$   
 11  $(\hat{\alpha}_{i1}, \dots, \hat{\alpha}_{iK})^T$  from  $K$ -dimensional normal distributions with zero means and compound  
 12 symmetry covariance with a correlation of 0 or 0.8, and variances of 1. Subsequently, the  
 13 continuous latent variable vector  $\hat{\alpha}_{i1}$  was converted into an attribute mastery pattern. More  
 14 precisely, if  $\hat{\alpha}_{ik}$  was greater than  $\Phi(k/(1+K))^{-1}$ ,  $\alpha_{ik} = 1$ ; otherwise,  $\alpha_{ik} = 0$ , where  $\Phi(\cdot)$   
 15  $^{-1}$  is the inverse cumulative normal distribution function. The simulated item responses were  
 16 randomly generated using these attribute mastery patterns, an assumed data-generating model  
 17 (DINA or general DCM), and item response probabilities. As mentioned in the previous section,  
 18 the parameters of the priors in the GB method were set to  $a_{jl}^0 = 2$ ,  $b_{jl}^0 = 1$ ,  $\forall j, l$  and  $\delta_l^0 = 1$ ,  $\forall l$ .  
 19 The step size of the Metropolis update was fixed at 0.05. A one-chain MCMC with 1,000  
 20 iterations was employed. The first 500 iterations are discarded as the burn-in period; therefore,  
 21 500 MCMC samples were used to approximate the posterior distributions.

22           The main target parameter is attribute masteries, and they are categorical latent  
 23 variables. However, common MCMC convergence criteria, such as Gelman-Rubin's  $\hat{R}$ , are for  
 24 continuous variables, which means the indicators may not be applicable to categorical variables.  
 25 Therefore, performing a convergence check of categorical variables in MCMC is not easy in this  
 26 context. Instead of directly checking for the convergence of attribute mastery, we calculated the  
 27 average correlations of the attribute mastery probabilities, which we estimated for the first and

1 second halves of the MCMC iterations after the burn-in period. If the estimated results of the  
 2 attribute mastery probabilities with the first half after the burn-in period are consistent with those  
 3 of the later MCMC iterations, we consider the attribute mastery results to be stable.

4 The attribute mastery pattern of the  $i$ -individual was calculated based on the posterior  
 5 attribute mastery probabilities. If the probability of the  $k$ -th attribute was greater than 0.5, the  
 6 attribute was considered mastered. Each estimation method was evaluated using two attribute  
 7 mastery recovery indices: attribute level agreement ratio (AAR) and pattern level agreement ratio  
 8 (PAR). AAR and PAR were calculated as follows:

$$\text{AAR}_k = \frac{1}{IM} \sum_{m=1}^M \sum_{i=1}^I \mathcal{J}(\hat{\alpha}_{ik}^{(m)} = \alpha_{ik}^{(\text{True})}), \forall k \quad (37)$$

$$\text{PAR} = \frac{1}{IM} \sum_{m=1}^M \sum_{i=1}^I \mathcal{J}(\hat{\alpha}_i^{(m)} = \alpha_i^{(\text{True})}), \quad (38)$$

9 where  $\hat{\alpha}_i^{(m)} = (\hat{\alpha}_{i1}^{(m)}, \hat{\alpha}_{i2}^{(m)}, \dots, \hat{\alpha}_{iK}^{(m)})^T$  is an estimate of the attribute mastery pattern for  
 10 individual  $i$  in the  $m$ -th simulation, and  $\alpha_i^{(\text{True})} = (\alpha_{i1}^{(\text{True})}, \alpha_{i2}^{(\text{True})}, \dots, \alpha_{iK}^{(\text{True})})^T$  is the true  
 11 attribute vector of individual  $i$ , where  $M$  is the total number of simulations, which is  $M =$   
 12 100.

### 13 5.2. Results

14 Table 3 shows the results, which indicate correlations greater than 0.98. Therefore, this  
 15 result can be interpreted as an indication that the MCMC iterations were stable and attribute  
 16 mastery can be estimated from the MCMC samples after the burn-in period.

17 Figures 1 and 2 present the simulation results of the DINA data generation with four-  
 18 and five-attribute Q-matrix conditions, respectively. In this simulation, the AARs and PARs of  
 19 the two Q-matrix conditions demonstrated similar tendencies; therefore, our discussion here  
 20 focuses on the four-attribute Q-matrix condition. The high item quality conditions presented in  
 21 the four left panels of Figure 1 indicated that all four estimation methods provide high AARs and  
 22 PARs. The low-item-quality conditions presented in the four right panels of Figure 1 indicated



1 lower AARs and PARs than the high-item-quality conditions, and the low-item-quality  
2 conditions exhibited some differences among the four estimation methods. Attribute correlations  
3 under low-item-quality conditions affected AARs and PARs more significantly. Furthermore,  
4 GBNPC and GBGNPC demonstrated higher AARs and PARs than the corresponding NPC and  
5 GNPC methods under the 30-sample size, 0.8 attribute correlation, and low-item-quality  
6 conditions. Interestingly, GBNPC exhibited the highest AARs and PARs under the 300-sample  
7 size, 0.8 attribute correlation, and low-item-quality conditions. Moreover, under these conditions,  
8 GBGNPC had similar AARs and PARs to the NPC, and the GNPC produced the least optimal  
9 result.

10           Figures 3 and 4 present the results of the general DCM data generation with four- and  
11 five-attribute Q-matrix conditions, respectively. Again, the AARs and PARs of the two Q-matrix  
12 conditions exhibited similar patterns; hereafter, we predominantly focus on results of the four-  
13 attribute Q-matrix conditions. Under high-item-quality conditions, GNPC and GBGNPC  
14 outperformed NPC and GBNPC. Furthermore, under high-attribute correlation conditions,  
15 GBGNPC was superior to GNPC; the same pattern was observed between GBNPC and NPC  
16 under the same conditions. In the low-item-quality conditions presented in the four right panels  
17 of Figure 3, GBGNPC and GBNPC tended to have higher AARs and PARs than GNPC and  
18 NPC. In particular, a sample size of 300, a high attribute correlation, and low-item-quality  
19 conditions indicated better AARs and PARs for GBGNPC and GBNPC than GNPC or NPC.

20           We also checked attribute mastery probabilities of the GBGNPC and GBNPC methods  
21 that represented uncertainty of parameter estimates. Figures 5 and 6 represent box plots of  
22 average attribute mastery probabilities of the four- and five-attribute conditions under the DINA  
23 model-based data-generating process. Interestingly, the GBNPC method tended to show higher  
24 average attribute mastery probabilities than the GBGNPC. The differences between the  
25 GBGNPC and GBNPC methods were relatively small in the first attribute but the discrepancy  
26 became larger as the attribute number increased. The later attributes were more difficult to  
27 master and the number of individuals mastering them was small. These tendencies also occurred

1 in the general data-generating process situations, which are shown in Figures 7 and 8. These  
2 posterior probabilities of attribute mastering represent estimation uncertainty, so we can carefully  
3 check the attribute mastery status. For example, the attribute mastery probabilities around cut-off  
4 values might represent indeterminacy of mastery or non-mastery. Such uncertainty quantification  
5 results cannot be obtained through the GNPC or NPC methods.

6 In summary, NPC and GBNPC tended to have higher AARs and PARs under DINA  
7 data generation, low-item-quality conditions, and high attribute correlations. However,  
8 GBGNPC was sometimes similar to NPC under the DINA data generation conditions, whereas  
9 GNPC was the least optimal. By contrast, under general DCM data generation conditions,  
10 GBGNPC and GNPC performed better than GBNPC and GNPC for high-quality items. For low-  
11 quality items, GBGNPC and GBNPC performed better. Based on these results, GBGNPC  
12 appears the optimal choice for attribute mastery estimation. If the DINA type item response  
13 mechanism is confirmed, GBNPC is the optimal choice among the four estimation methods from  
14 the perspective of attribute recovery.

15 The possible reason for the superiority of the GBGNPC over the GNPC is prior  
16 settings. In our simulation setting, sample sizes were relatively small in the situations in which  
17 the nonparametric methods were employed. Under such conditions, estimation of weight  
18 parameters might be difficult for the GNPC, especially in the low-item-quality conditions. The  
19 GBGNPC, on the other hand, assumed priors for the weight parameters, and the prior conveyed  
20 information of item characteristics and succeeded in estimating attribute mastery patterns.  
21 Another reason may be that the GBGNPC can deal with uncertainty in parameter estimation.  
22 This means that the GNPC uses parameter estimates to minimize the loss function, which simply  
23 selects the attribute mastery pattern that provides the minimum value of loss function without  
24 considering the second or third best attribute mastery patterns. By contrast, the GBGNPC can  
25 consider and use the second-best attribute mastery pattern for estimating attribute mastery  
26 probabilities. If these considerations are correct, even if we use non-information priors for the



1 GB method and nonparametric methods do not contain model parameters and only estimate  
 2 attribute mastery patterns that relate to individuals. Therefore, the attribute mastery patterns in a  
 3 training data set are not contained in a test data set. Exploring the appropriate quantitative  
 4 evaluation for the GB method in the DCMs is an important direction for future research.

### 5 **6.1. Data Analysis Settings**

6 The Examination of the Certificate of Proficiency in English (ECPE) data were  
 7 selected as an example. ECPE data have been analyzed in various previous studies, such as  
 8 Templin and Hoffman (2013) and Templin and Bradshaw (2014). The ECPE data contained  
 9 2,922 responses for 28 items. Table 4 presents a  $28 \times 3$  Q-matrix that assumes three attributes:  
 10 Morphosyntactic ( $\alpha_1$ ), cohesive ( $\alpha_2$ ), and lexical rules ( $\alpha_3$ ). The settings of the GB methods  
 11 were the same as those used in previous simulations. One difference was that we employed  
 12 GNPC and NPC estimates as initial values for GBGNPC and GBNPC. The data analysis code  
 13 can be obtained from the OSF webpage <https://osf.io/sau6j/>.

### 14 **5.2. Results**

15 The same correlations as in the simulation study were calculated. Again, the  
 16 correlations of the three attributes with the GBNPC and GBGNPC methods were all greater than  
 17 0.99. This indicated that the MCMC iterations for attribute mastery were stable.

18 Table 6 lists the frequencies and ratios of the attribute mastery patterns for the four  
 19 estimation methods. Several differences are observed in Table 6. First, GBGNPC and GBNPC  
 20 estimated the pattern (001) to be lower than GNPC and NPC estimates. Second, as pattern (011)  
 21 indicates, the frequency of pattern (011) for GBGNPC was the highest (1203), that for the GNPC  
 22 was the second (955), that for the NPC was the third (522), and that for the GBNPC was the last  
 23 (386). The GBGNPC and GBNPC produced lower frequencies than the GNPC and NPC for  
 24 patterns (100), (101), and (110). The final difference is indicated in pattern (111). GBGNPC and  
 25 GNPC had relatively smaller numbers than GBNPC and NPC.

26 Table 5 shows the means and *SDs* of the attribute mastery probabilities for the  
 27 GBGNPC and GBNPC methods. The attribute mastery probability for the first attribute

1 (Morphosyntactic rules) of GBGNPC was Mean = .551( $SD = .388$ ) and that of GBNPC was  
2 Mean = .807( $SD = .326$ ). The discrepancy was the largest among the three attributes. The  
3 attribute mastery probabilities for the second (cohesive rules) and third (lexical rules) attributes  
4 using the GBGNPC and GBNPC methods were higher than 0.90 so these attributes tended to be  
5 mastered.

6 Table 7 shows the estimated attribute mastery patterns of GBGNPC and GNPC. A  
7 large portion of the GBGNPC pattern (011) corresponds to patterns (001), (001), and (010) of the  
8 GNPC. Furthermore, patterns (011), (100), (101), and (110) of the GNPC correspond to pattern  
9 (111) of the GBGNPC. From these results, the GBGNPC tended to overestimate the number of  
10 attributes compared with the GNPC.

11 Table 8 presents the GBNPC's and NPC's estimated attribute mastery patterns. The  
12 results in Table 8 are similar to those of GBGNPC and GNPC. For example, patterns (000),  
13 (001), (010), and (010) with the NPC are sometimes estimated as pattern (011) in GBGNPC.  
14 Furthermore, patterns (000) to (110) in the NPC were classified as pattern (111) in the  
15 GBGNPC. Therefore, the GBNPC overestimates the number of attributes compared with the  
16 NPC.

17 We checked individual differences between the GBGNPC and GNPC methods. Table  
18 9 shows that some individuals indicated the largest pattern discrepancy of attribute mastery  
19 between GBGNPC and GNPC methods. The GBGNPC and the GNPC provided  $\alpha = (0,1,1)$   
20 and  $\alpha = (1,0,0)$ , respectively. The response patterns did not indicate systematic tendency but  
21 the sum scores of the individuals ranged from 11-15, which meant they could answer more than  
22 half of the test items. The maximum sub-scores for attributes one, two, and three were 13, 6, and  
23 18, respectively, so the individuals in Table 9 received half points out of the maximum total sub-  
24 scores. In addition, the sum scores of the individuals ranged from 11 – 15, which is about half  
25 the maximum sum-score of 28. Thus, the pattern (1,0,0) might saliently underestimate the  
26 latent attributes, making the pattern (0,1,1) possibly more likely. Furthermore, some attribute  
27 mastery probabilities were close to the cut-off value 0.5. For example, the mastery probability of



1 knowledge. The proposed estimation method can be used for any type of loss function and has  
2 great flexibility. This study's contribution is that the proposed method provides a novel approach  
3 for estimating the DCMs' parameters. The GB method is flexible because we can select any type  
4 of loss function and consider the uncertainty of the parameter estimation. Furthermore, the  
5 proposed method relaxes the assumption of the typical Bayesian method, which requires a  
6 likelihood function. The theoretical analysis revealed consistent results for the proposed GB  
7 method under mild regularity conditions. Additionally, the simulation study revealed that the GB  
8 method improved attribute mastery recoveries compared to previous nonparametric methods.  
9 The real data example indicated that the proposed GB method with the nonparametric loss  
10 function tended to overestimate attribute mastery compared to the nonparametric methods.

11           The theoretical results not only guarantee the consistency of the MAP estimation  
12 results, but also give convergence rate results, which is helpful in characterizing the finite sample  
13 estimate errors. All these results are new to the literature and provide theoretical justification for  
14 using the nonparametric methods and the proposed GB approach. Moreover, the theoretical  
15 results in the paper are established for the general loss function under the proposed assumptions.  
16 It covers popular loss functions, such as the GNPC and log-likelihood loss functions, which are  
17 used in C. Ma, de la Torre, et al. (2023).

18           One interesting future research problem is to establish consistent results for other  
19 Bayesian estimators, such as expected a posteriori (EAP). However, this is a more challenging  
20 question as it involves deriving the limiting distribution of the Bayesian posterior distribution.  
21 Intuitively, given our theoretical results of MAP, EAP would also be consistent, but technically  
22 this is not easy to determine and needs the development of new mathematical tools. Moreover,  
23 Assumption 2 may be further relaxed to allow for some latent attribute mastery patterns that do  
24 not exist in the population. In particular, if we know which attribute mastery patterns have zero  
25 probability, such as in hierarchical DCMs, then our theoretical results would still apply.  
26 However, if this information is unknown, while some latent attribute mastery patterns have zero

1 probability, the model itself may have some identifiability issues under the nonparametric DCMs  
2 setting. This is another interesting topic for future study.

3         Another future research direction is to explore how to determine the learning rate from  
4 data, especially under the  $\mathcal{M}$ -open setting. Intuitively, the learning rate controls the relative  
5 importance between prior information and the loss function. We can set a relatively small value  
6 for the learning rate if we have enough prior information about the attribute mastery distribution  
7 and use several new items whose nature we do not know. In this case, we put relatively great  
8 importance on the prior information rather than the obtained data. However, it may not be  
9 realistic to set the learning rate greater than one. Such a high learning rate would amplify the  
10 effect of the loss function but might indicate an overreliance on the data. It may not be suitable  
11 for the  $\mathcal{M}$ -open setting that the data-generating process is unknown. Therefore, we need to  
12 explore how to determine the learning rate from data.

13         As mentioned previously, no scholarly agreement exists regarding how to determine  
14 the learning rate, which is an important topic for future research especially in the DCM context.  
15 In particular, data-driven learning rate determination procedures were studied in Wu and Martin  
16 (2023), where several selection methods such as the SafeBayes algorithm based on the  
17 cumulative log-loss (Grünwald & van Ommen, 2017), information gain perspective (Holmes &  
18 Walker, 2017), modified weighted likelihood bootstrap approach (Lyddon et al., 2019), and the  
19 approximate achievement of nominal frequentist coverage probability (Syring & Martin, 2019)  
20 were compared. However, all of these methods have different foundations, and we need to  
21 explore which one is most appropriate for the DCM context.

22         Another topic that requires further investigation is model data fit evaluation. From our  
23 understanding, the GB method avoids explicit model representation in the framework. Therefore,  
24 the model evaluation scheme is not included in the procedure of the GB method. This is also true  
25 for the GBGNPC method proposed in this study. Therefore, future research needs to explore  
26 what kind of statistics can be used for model data fit. In particular, previously developed  
27 methods of model data fit assessment in psychometrics and Bayesian data analysis could be



1 employed in our setting. Following Sinharay et al. (2006), discrepancy measures such as  
2 observed score distribution, point biserial correlation, and statistical measures of association  
3 among the item pairs could be used for posterior predictive model checking (PPMC). For further  
4 details on PPMC methods for Bayesian networks and IRT models, see also Sinharay (2006) and  
5 Sinharay (2016). Moreover, PPMC for person fit (Sinharay, 2015) would also provide an  
6 important measure to assess the model fit for the attribute mastery patterns at the personal level,  
7 which is often of interest in cognitive diagnosis.

8           As a final note about the choice of estimation methods, it is necessary to consider  
9 estimation time. The GB method employs an MCMC procedure, so it has a longer estimation  
10 time than that of the nonparametric methods. In our simulation, the estimation times were less  
11 than ten seconds, so it is not irritatingly time consuming. However, if we need immediate  
12 feedback, the time difference between the two kinds of methods may be crucial. We also need to  
13 consider estimation time for the requirement of real data analysis.

14

15

### References

- 16 Bernardo, J. M., & Smith, A. F. (2009). *Bayesian theory*. John Wiley & Sons.
- 17 Bissiri, P. G., Holmes, C. C., & Walker, S. G. (2016). A general framework for updating belief  
18 distributions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*,  
19 78(5), 1103-1130. <https://doi.org/10.1111/rssb.12158>
- 20 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M., Brubaker,  
21 M. A., Guo, J., Li, P., & Riddell, A. (2017). Stan: A probabilistic programming language.  
22 *Journal of Statistical Software*, 76(1), 1-32. <https://doi.org/10.18637/jss.v076.i01>
- 23 Chen, Y., Liu, J., Xu, G., & Ying, Z. (2015). Statistical analysis of *Q*-matrix based diagnostic  
24 classification models. *Journal of the American Statistical Association*, 110, 850-866.  
25 <https://doi.org/10.1080/01621459.2014.934827>

- 1 Chiu, C.-Y., & Douglas, J. (2013). A nonparametric approach to cognitive diagnosis by  
2 proximity to ideal response patterns. *Journal of Classification*, *30*, 225-250.  
3 <https://doi.org/10.1007/s00357-013-9132-9>
- 4 Chiu, C.-Y., & Köhn, H.-F. (2019). Consistency theory for the general nonparametric  
5 classification method. *Psychometrika*, *84*, 830-845. [https://doi.org/10.1007/s11336-019-](https://doi.org/10.1007/s11336-019-09660-x)  
6 [09660-x](https://doi.org/10.1007/s11336-019-09660-x)
- 7 Chiu, C.-Y., Köhn, H.-F., Zheng, Y., & Henson, R. (2016). Joint maximum likelihood estimation  
8 for diagnostic classification models. *Psychometrika*, *81*, 1069-1092.  
9 <https://doi.org/10.1007/s11336-016-9534-9>
- 10 Chiu, C.-Y., Sun, Y., & Bian, Y. (2018). Cognitive diagnosis for small educational programs:  
11 The general nonparametric classification method. *Psychometrika*, *83*, 355-375.  
12 <https://doi.org/10.1007/s11336-017-9595-4>
- 13 Culpepper, S. A. (2015). Bayesian estimation of the DINA model with Gibbs sampling. *Journal*  
14 *of Educational and Behavioral Statistics*, *40*(5), 454-476.  
15 <https://doi.org/10.3102/1076998615595403>
- 16 de la Torre, J. (2009). DINA model and parameter estimation: A didactic. *Journal of Educational*  
17 *and Behavioral Statistics*, *34*(1), 115-130. <https://doi.org/10.3102/1076998607309474>
- 18 de la Torre, J. (2011). The generalized DINA model framework. *Psychometrika*, *76*, 179-199.  
19 <https://doi.org/10.1007/s11336-011-9207-7>
- 20 Grünwald, P., & van Ommen, T. (2017). Inconsistency of Bayesian inference for misspecified  
21 linear models, and a proposal for repairing it. *Bayesian Analysis*, *12*(4), 1069-1103.  
22 <https://doi.org/10.1214/17-BA1085>
- 23 Henson, R. A., Templin, J. L., & Willse, J. T. (2009). Defining a family of cognitive diagnosis  
24 models using log-linear models with latent variables. *Psychometrika*, *74*, 191-210.  
25 <https://doi.org/10.1007/s11336-008-9089-5>
- 26 Holmes, C. C., & Walker, S. G. (2017). Assigning a value to a power likelihood in a general  
27 Bayesian model. *Biometrika*, *104*(2), 497-503. <https://doi.org/10.1093/biomet/asx010>

- 1 Junker, B. W., & Sijtsma, K. (2001). Cognitive assessment models with few assumptions, and  
2 connections with nonparametric item response theory. *Applied Psychological Measurement*,  
3 25(3), 258-272. <https://doi.org/10.1177/01466210122032064>
- 4 Li, H., Hunter, C. V., & Lei, P.-W. (2016). The selection of cognitive diagnostic models for a  
5 reading comprehension test. *Language Testing*, 33(3), 391-409.  
6 <https://doi.org/10.1177/0265532215590848>
- 7 Lyddon, S. P., Holmes, C. C., & Walker, S. G. (2019). General Bayesian updating and the loss-  
8 likelihood bootstrap. *Biometrika*, 106(2), 465-478. <https://doi.org/10.1093/biomet/asz006>
- 9 Ma, C., de la Torre, J., & Xu, G. (2023). Bridging parametric and nonparametric methods in  
10 cognitive diagnosis. *Psychometrika*, 88(1), 51-75. <https://doi.org/10.1007/s11336-022->  
11 09878-2
- 12 Ma, C., Ouyang, J., & Xu, G. (2023). Learning latent and hierarchical structures in cognitive  
13 diagnosis models. *Psychometrika*, 88(1), 175-207. <https://doi.org/10.1007/s11336-022->  
14 09867-5
- 15 Ma, W., & Jiang, Z. (2021). Estimating cognitive diagnosis models in small samples: Bayes  
16 modal estimation and monotonic constraints. *Applied Psychological Measurement*, 45(2),  
17 95-111. <https://doi.org/10.1177/0146621620977681>
- 18 MacReady, G. B., & Dayton, C. M. (1977). The use of probabilistic models in the assessment of  
19 mastery. *Journal of Educational Statistics*, 2(2), 99-120. <https://doi.org/10.2307/1164802>
- 20 Maris, E. (1999). Estimating multiple classification latent class models. *Psychometrika*, 64, 187-  
21 212. <https://doi.org/10.1007/BF02294535>
- 22 Plummer, M. (2003). JAGS: A program for analysis of Bayesian graphical models using Gibbs  
23 sampling. In *Proceedings of the 3rd International Workshop on Distributed Statistical*  
24 *Computing* (Vol. 124, pp. 1-10). <https://www.r-project.org/conferences/DSC->  
25 2003/Drafts/Plummer.pdf

- 1 Rupp, A. A., & Templin, J. L. (2008). Unique characteristics of diagnostic classification models:  
2 A comprehensive review of the current state-of-the-art. *Measurement*, 6(4), 219-262.  
3 <https://doi.org/10.1080/15366360802490866>
- 4 Rupp, A. A., Templin, J. L., & Henson, R. A. (2010). *Diagnostic measurement: Theory,*  
5 *methods, and applications*. Guilford Press.
- 6 Sinharay, S. (2006). Model diagnostics for Bayesian networks. *Journal of Educational and*  
7 *Behavioral Statistics*, 31(1), 1-33. <https://doi.org/10.3102/10769986031001001>
- 8 Sinharay, S. (2015). Assessment of person fit for mixed-format tests. *Journal of Educational and*  
9 *Behavioral Statistics*, 40(4), 343-365. <https://doi.org/10.3102/1076998615589128>
- 10 Sinharay, S. (2016). Bayesian model fit and model comparison. In W. J. van der Linden (Ed.),  
11 *Handbook of item response theory, volume 2: Statistical tools* (pp. 379-394). Chapman and  
12 Hall/CRC. <https://doi.org/10.1201/b19166>
- 13 Sinharay, S., Johnson, M. S., & Stern, H. S. (2006). Posterior predictive assessment of item  
14 response theory models. *Applied Psychological Measurement*, 30(4), 298-321.  
15 <https://doi.org/10.1177/0146621605285517>
- 16 Syring, N., & Martin, R. (2019). Calibrating general posterior credible regions. *Biometrika*,  
17 106(2), 479-486. <https://doi.org/10.1093/biomet/asy054>
- 18 Tatsuoka, K. K. (1985). A probabilistic model for diagnosing misconceptions by the pattern  
19 classification approach. *Journal of Educational Statistics*, 10(1), 55-73.  
20 <https://doi.org/10.3102/10769986010001055>
- 21 Templin, J. L., & Bradshaw, L. (2014). Hierarchical diagnostic classification models: A family  
22 of models for estimating and testing attribute hierarchies. *Psychometrika*, 79, 317-339.  
23 <https://doi.org/10.1007/s11336-013-9362-0>
- 24 Templin, J. L., & Henson, R. A. (2006). Measurement of psychological disorders using cognitive  
25 diagnosis models. *Psychological Methods*, 11(3), 287-305. [https://doi.org/10.1037/1082-](https://doi.org/10.1037/1082-989X.11.3.287)  
26 989X.11.3.287

- 1 Templin, J. L., & Hoffman, L. (2013). Obtaining diagnostic classification model estimates using  
2 Mplus. *Educational Measurement: Issues and Practice*, 32(2), 37-50.  
3 <https://doi.org/10.1111/emip.12010>
- 4 von Davier, M. (2008). A general diagnostic model applied to language testing data. *British*  
5 *Journal of Mathematical and Statistical Psychology*, 61(2), 287-307.  
6 <https://doi.org/10.1348/000711007X193957>
- 7 von Davier, M., & Lee, Y.-S. (2019). *Handbook of diagnostic classification models: Models and*  
8 *model extensions, applications, software packages*. Springer.
- 9 Wang, S., & Douglas, J. (2015). Consistency of nonparametric classification in cognitive  
10 diagnosis. *Psychometrika*, 80, 85-100. <https://doi.org/10.1007/s11336-013-9372-y>
- 11 Wu, P.-S., & Martin, R. (2023). A comparison of learning rate selection methods in generalized  
12 Bayesian inference. *Bayesian Analysis*, 18(1), 105-132. <https://doi.org/10.1214/21-BA1302>
- 13 Xu, G. (2017). Identifiability of restricted latent class models with binary responses. *The Annals*  
14 *of Statistics*, 45(2), 675-707. <https://doi.org/10.1214/16-AOS1464>
- 15 Xu, G., & Shang, Z. (2018). Identifying latent structures in restricted latent class models. *Journal*  
16 *of the American Statistical Association*, 113(523), 1284-1295.  
17 <https://doi.org/10.1080/01621459.2017.1340889>
- 18 Yamaguchi, K., & Okada, K. (2020). Variational Bayes inference for the DINA model. *Journal*  
19 *of Educational and Behavioral Statistics*, 45(5), 569-597.  
20 <https://doi.org/10.3102/1076998620911934>
- 21 Yamaguchi, K., & Templin, J. L. (2022a). Direct estimation of diagnostic classification model  
22 attribute mastery profiles via a collapsed Gibbs sampling algorithm. *Psychometrika*, 87(4),  
23 1390-1421. <https://doi.org/10.1007/s11336-022-09857-7>
- 24 Yamaguchi, K., & Templin, J. L. (2022b). A Gibbs sampling algorithm with monotonicity  
25 constraints for diagnostic classification models. *Journal of Classification*, 39(1), 24-54.  
26 <https://doi.org/10.1007/s00357-021-09392-7>  
27

1  
2  
3  
4

Table 1.

*The four-attribute Q-matrix*

Item	Attribute			
	1	2	3	4
1	1	0	0	0
2	0	1	0	0
3	0	0	1	0
4	0	0	0	1
5	1	0	0	0
6	0	1	0	0
7	0	0	1	0
8	0	0	0	1
9	1	1	0	0
10	1	0	1	0
11	1	0	0	1
12	0	1	1	0
13	0	1	0	1
14	0	0	1	1
15	1	1	1	0
16	1	1	0	1
17	1	0	1	1
18	0	1	1	1
19	1	1	1	1

1

Table 2.

*The five-attribute Q-matrix*

Item	Attribute					Item	Attribute				
	1	2	3	4	5		1	2	3	4	5
1	1	0	0	0	0	16	0	1	0	1	0
2	0	1	0	0	0	17	0	1	0	0	1
3	0	0	1	0	0	18	0	0	1	1	0
4	0	0	0	1	0	19	0	0	1	0	1
5	0	0	0	0	1	20	0	0	0	1	1
6	1	0	0	0	0	21	1	1	1	0	0
7	0	1	0	0	0	22	1	1	0	1	0
8	0	0	1	0	0	23	1	1	0	0	1
9	0	0	0	1	0	24	1	0	1	1	0
10	0	0	0	0	1	25	1	0	1	0	1
11	1	1	0	0	0	26	1	0	0	1	1
12	1	0	1	0	0	27	0	1	1	1	0
13	1	0	0	1	0	28	0	1	1	0	1
14	1	0	0	0	1	29	0	1	0	1	1
15	0	1	1	0	0	30	0	0	1	1	1

2

3

1  
2

Table 3.

*The average correlations of attribute mastery probabilities estimated by first and second halves of MCMC iterations after the burn-in period*

Data generating model	Sample size	Attribute correlation	Item quality	GBGNPC		GBNPC	
				Three attributes	Four attributes	Three attributes	Four attributes
DINA	30	0	High	.994	.994	.998	.998
			Low	.984	.983	.993	.994
		0.8	High	.996	.995	.998	.998
			Low	.984	.983	.995	.995
	300	0	High	.997	.998	.998	.999
			Low	.992	.993	.993	.995
		0.8	High	.999	.999	.999	.999
			Low	.992	.993	.995	.996
General	30	0	High	.990	.989	.996	.997
			Low	.981	.980	.994	.995
		0.8	High	.994	.994	.998	.998
			Low	.982	.982	.995	.995
	300	0	High	.997	.997	.997	.997
			Low	.992	.993	.994	.995
		0.8	High	.999	.999	.999	.999
			Low	.991	.993	.996	.996

*Note:* GBGNPC: generalized Bayesian method with generalized nonparametric loss function, GBNPC: generalized Bayesian method with nonparametric loss function.

3



1  
2

Table 4.

*The Q-matrix of ECPE data*

Item	Attribute			Item	Attribute		
	Morphosyntactic rules: $\alpha_1$	Cohesive rules: $\alpha_2$	Lexical rules: $\alpha_3$		Morphosyntactic rules: $\alpha_1$	Cohesive rules: $\alpha_2$	Lexical rules: $\alpha_3$
1	1	1	0	15	0	0	1
2	0	1	0	16	1	0	1
3	1	0	1	17	0	1	1
4	0	0	1	18	0	0	1
5	0	0	1	19	0	0	1
6	0	0	1	20	1	0	1
7	1	0	1	21	1	0	1
8	0	1	0	22	0	0	1
9	0	0	1	23	0	1	0
10	1	0	0	24	0	1	0
11	1	0	1	25	1	0	0
12	1	0	1	26	0	0	1
13	1	0	0	27	1	0	0
14	1	0	0	28	0	0	1

3  
4

1

Table 5.  
*Means and SDs of posterior attribute mastery probabilities  
 for GBGNPC and GBNPC methods*

Estimation method	Attribute	Mean	SD
GBGNPC	Morphosyntactic rules: $\alpha_1$	.551	.388
	Cohesive rules: $\alpha_2$	.985	.077
	Lexical rules: $\alpha_3$	.939	.198
GBNPC	Morphosyntactic rules: $\alpha_1$	.807	.326
	Cohesive rules: $\alpha_2$	.978	.103
	Lexical rules: $\alpha_3$	.949	.187

2

3

1  
2

Table 6.

*Frequencies and ratios of the estimated attribute mastery patterns with the four estimation methods*

Pattern	GBGNPC		GBNPC		GNPC		NPC	
	Frequency	Ratio	Frequency	Ratio	Frequency	Ratio	Frequency	Ratio
000	24	.008	35	.012	29	.010	44	.015
001	2	.001	3	.001	155	.053	91	.031
010	88	.030	64	.022	88	.030	82	.028
011	1201	.411	384	.131	955	.327	522	.179
100	3	.001	3	.001	38	.013	27	.009
101	0	.000	0	.000	82	.028	87	.030
110	45	.015	36	.012	157	.054	96	.033
111	1559	.534	2397	.820	1418	.485	1973	.675

*Note:* GBGNPC: generalized Bayesian method with generalized nonparametric loss function, GB-NPC: generalized Bayesian method with nonparametric loss function, GNPC: generalized nonparametric method, NPC: nonparametric method.

3  
4

1  
2

Table 7.

*Contingency table of the estimated attribute mastery patterns by GBGNPC and GNPC*

GBGNPC	GNPC							
	000	001	010	011	100	101	110	111
000	18	1	0	0	5	0	0	0
001	0	1	0	0	1	0	0	0
010	7	1	54	0	9	0	17	0
011	4	152	34	924	11	6	41	29
100	0	0	0	0	3	0	0	0
101	0	0	0	0	0	0	0	0
110	0	0	0	0	5	0	40	0
111	0	0	0	31	4	76	59	1389

*Note:* GBGNPC: generalized Bayesian method with generalized nonparametric loss function,  
GNPC: generalized nonparametric method.

3  
4

1  
2

Table 8.

*Contingency table of the estimated attribute mastery patterns by GBNPC and NPC*

GBNPC	NPC							
	000	001	010	011	100	101	110	111
000	26	5	1	0	3	0	0	0
001	0	3	0	0	0	0	0	0
010	10	0	37	0	8	0	9	0
011	7	54	32	278	3	0	10	0
100	0	0	0	0	3	0	0	0
101	0	0	0	0	0	0	0	0
110	0	0	2	0	3	0	31	0
111	1	29	10	244	7	87	46	1973

*Note:* GBNPC: generalized Bayesian method with nonparametric loss function, NPC: nonparametric method.

3  
4  
5

1

Table 9.

*Individual differences in estimated patterns for GBGNPC and GNPC methods, response patterns, sum- and sub-scores, and attribute mastery probabilities*

ID	Attribute mastery pattern		Response pattern	Sum-score	Sub-score			Attribute mastery probability		
	GBGNPC	GNPC			$\alpha_1$	$\alpha_2$	$\alpha_3$	$\alpha_1$	$\alpha_2$	$\alpha_3$
813	011	100	1000100110100100000011101100	11	5	3	6	.130	.882	.516
1060	011	100	1000011101011000111010001101	14	7	3	9	.418	.956	.864
2378	011	100	1110110111110010010000001011	15	7	3	9	.420	.996	.982
2607	011	100	1000110000101000110010101100	11	5	3	7	.110	.874	.556

*Note:* GBGNPC: generalized Bayesian method with a generalized nonparametric loss function, GNPC: generalized Bayesian method with a nonparametric loss function.

2

1  
2

Table 10.

*Generalized posterior of attribute mastery pattern by GBNPC and NPC*

Estimation method	ID	Attribute mastery pattern							
		000	100	010	110	001	101	011	111
GBGNPC	813	.102	.016	<b>.264</b>	.102	0	0	<b>.504</b>	.012
	1060	.008	.022	.020	.086	.006	.008	<b>.548</b>	<b>.302</b>
	2378	.002	0	.004	.012	0	.002	<b>.574</b>	<b>.406</b>
	2607	.104	.012	<b>.246</b>	.082	.010	0	<b>.530</b>	.016
GBNPC	813	<b>.282</b>	.012	<b>.316</b>	<b>.250</b>	.006	0	.106	.028
	1060	.022	.008	.004	.016	0	.004	.066	<b>.880</b>
	2378	.006	.002	.004	.016	.002	.002	.072	<b>.896</b>
	2607	<b>.520</b>	.028	.082	.054	.022	0	<b>.244</b>	.050

*Note:* GBGNPC: generalized Bayesian method with a generalized nonparametric loss function, GBNPC: generalized Bayesian method with a nonparametric loss function.

3  
4  
5  
6

Figure 1.

Simulation results of the DINA data generation with four-attribute Q-matrix conditions

1  
2  
3  
5

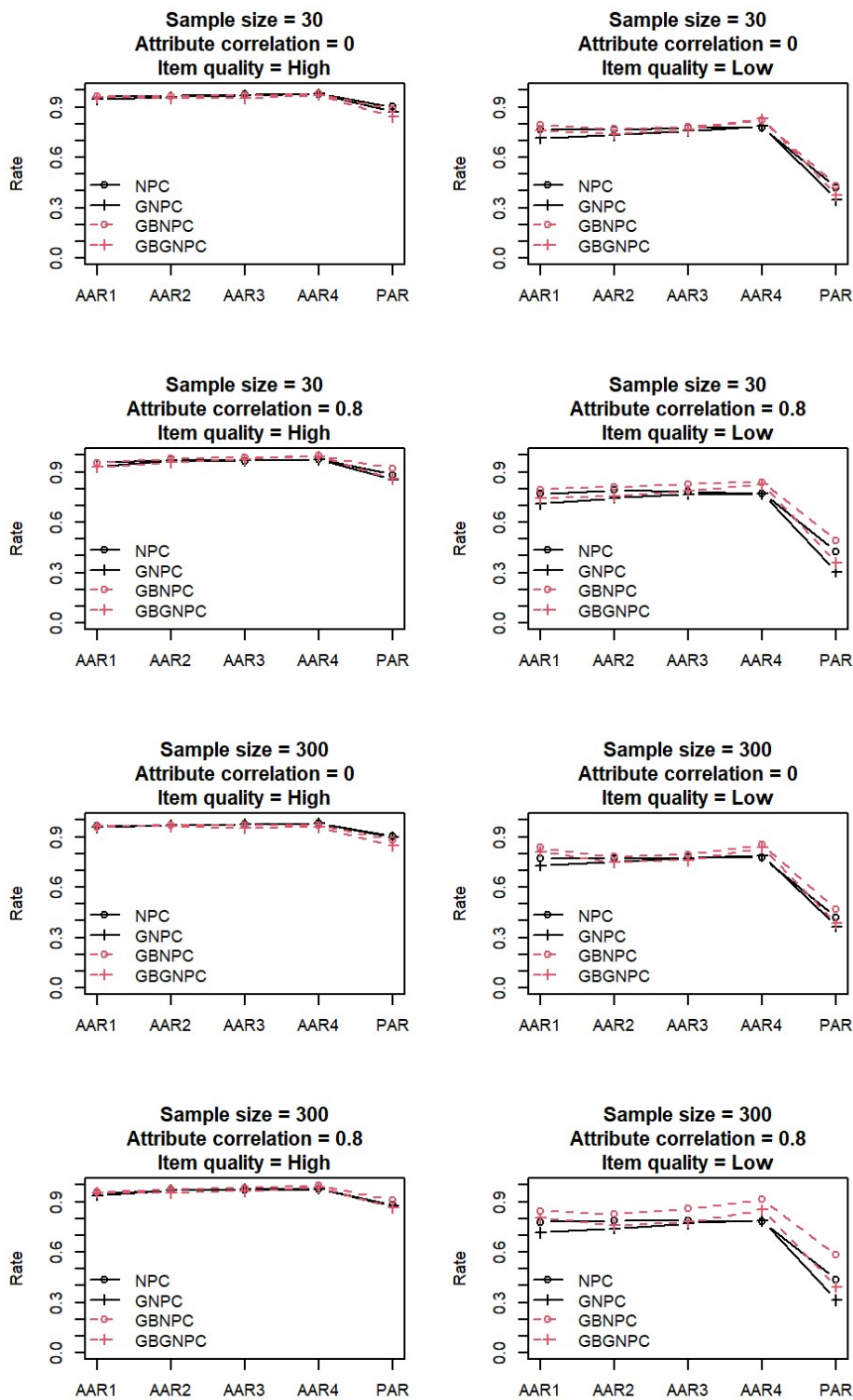
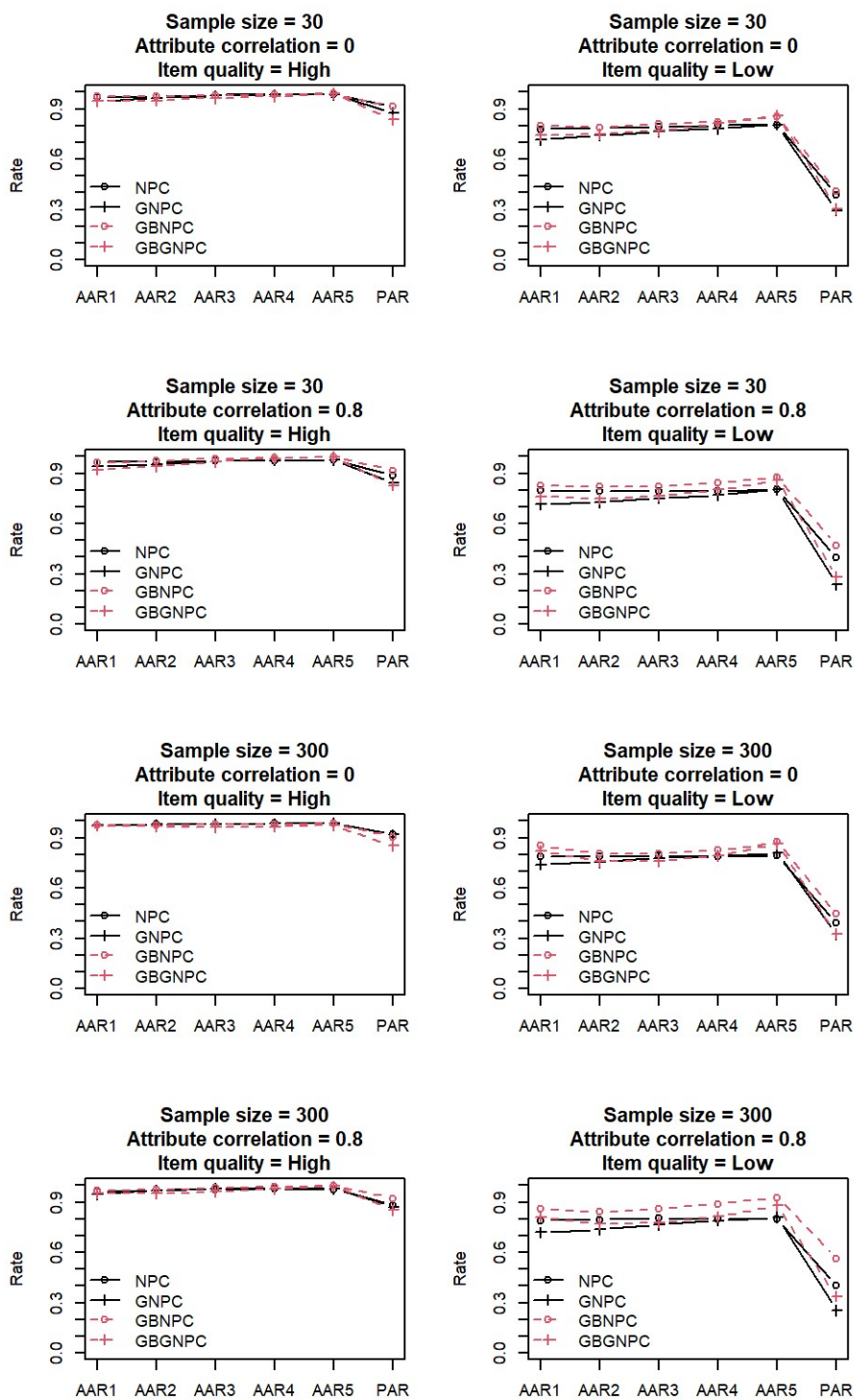




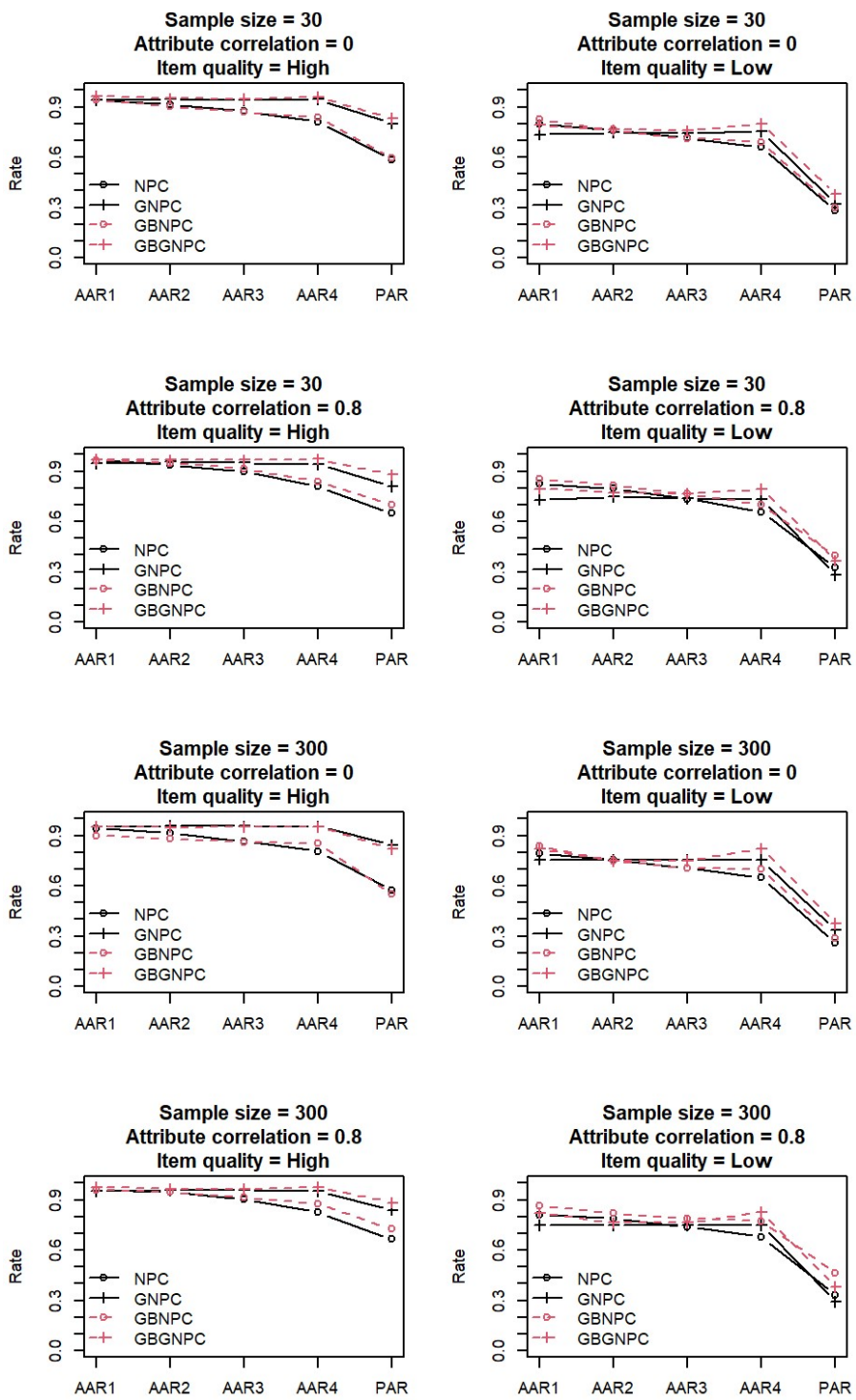
Figure 2.

Simulation results of the DINA data generation with five-attribute Q-matrix conditions



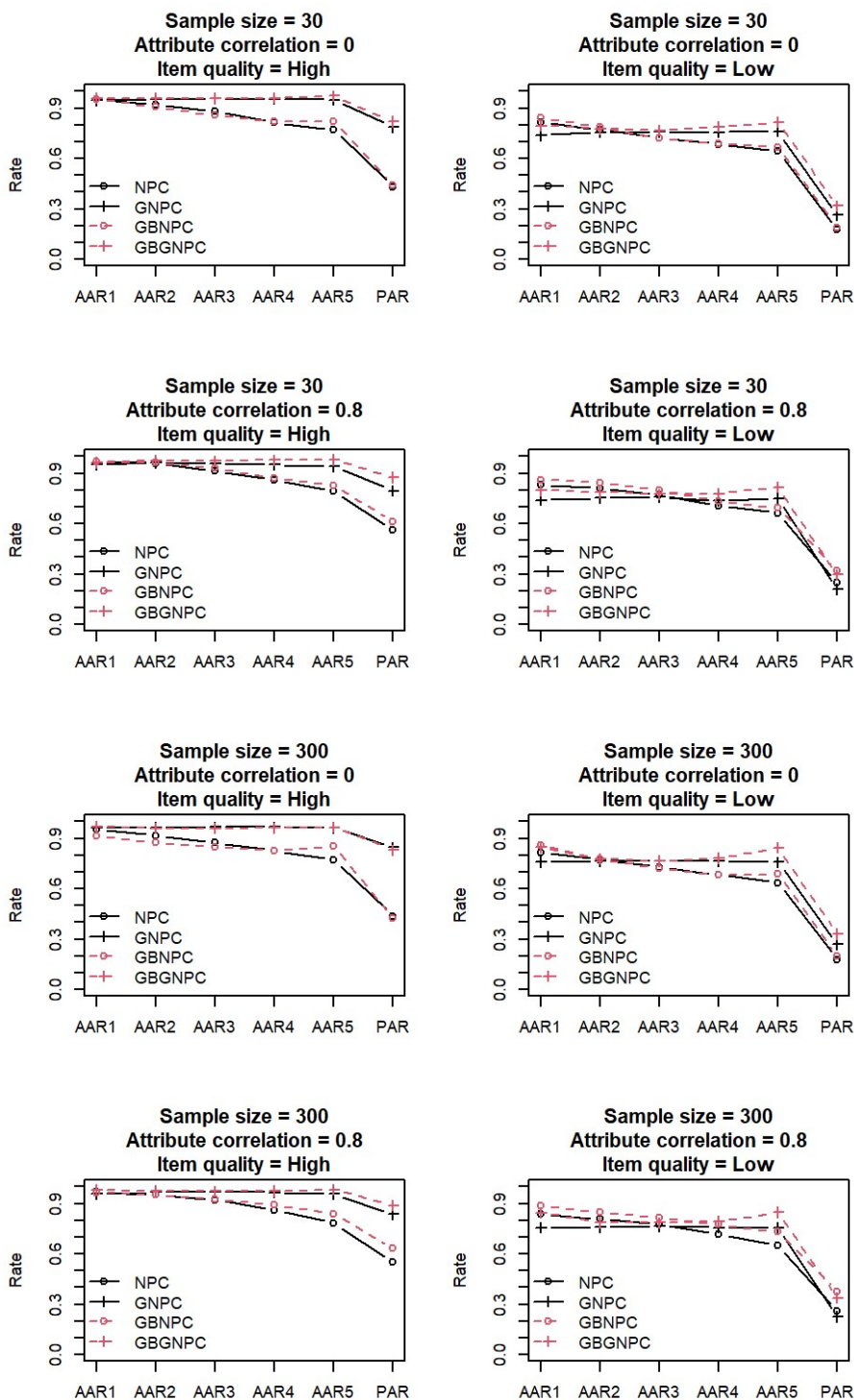
1  
2  
3  
4  
5  
6

Figure 3.  
*Simulation results of the general DCM data generation with four-attribute Q-matrix conditions*

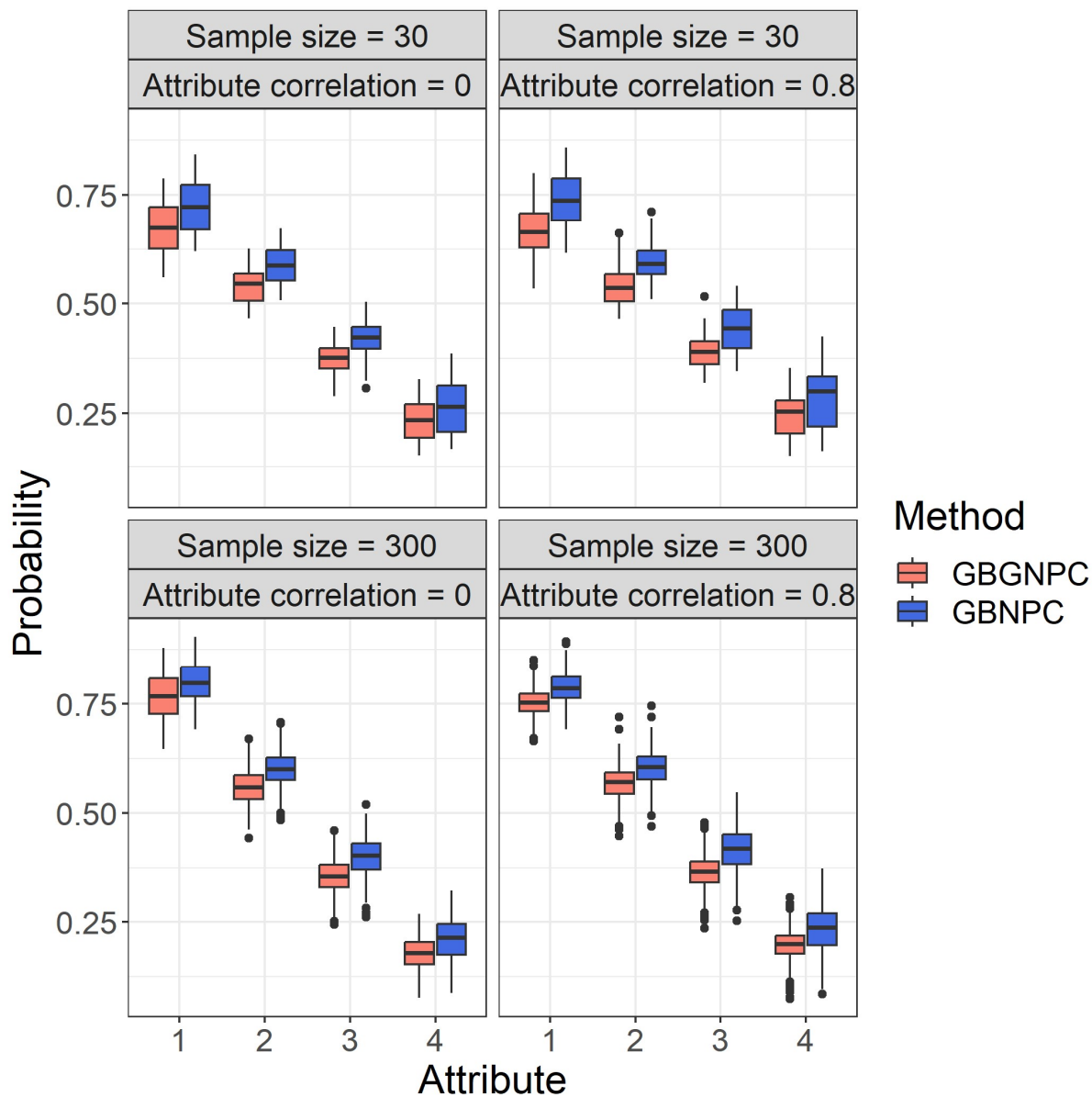


1  
2  
3  
4

Figure 4.  
*Simulation results of the general DCM data generation with five-attribute Q-matrix conditions*

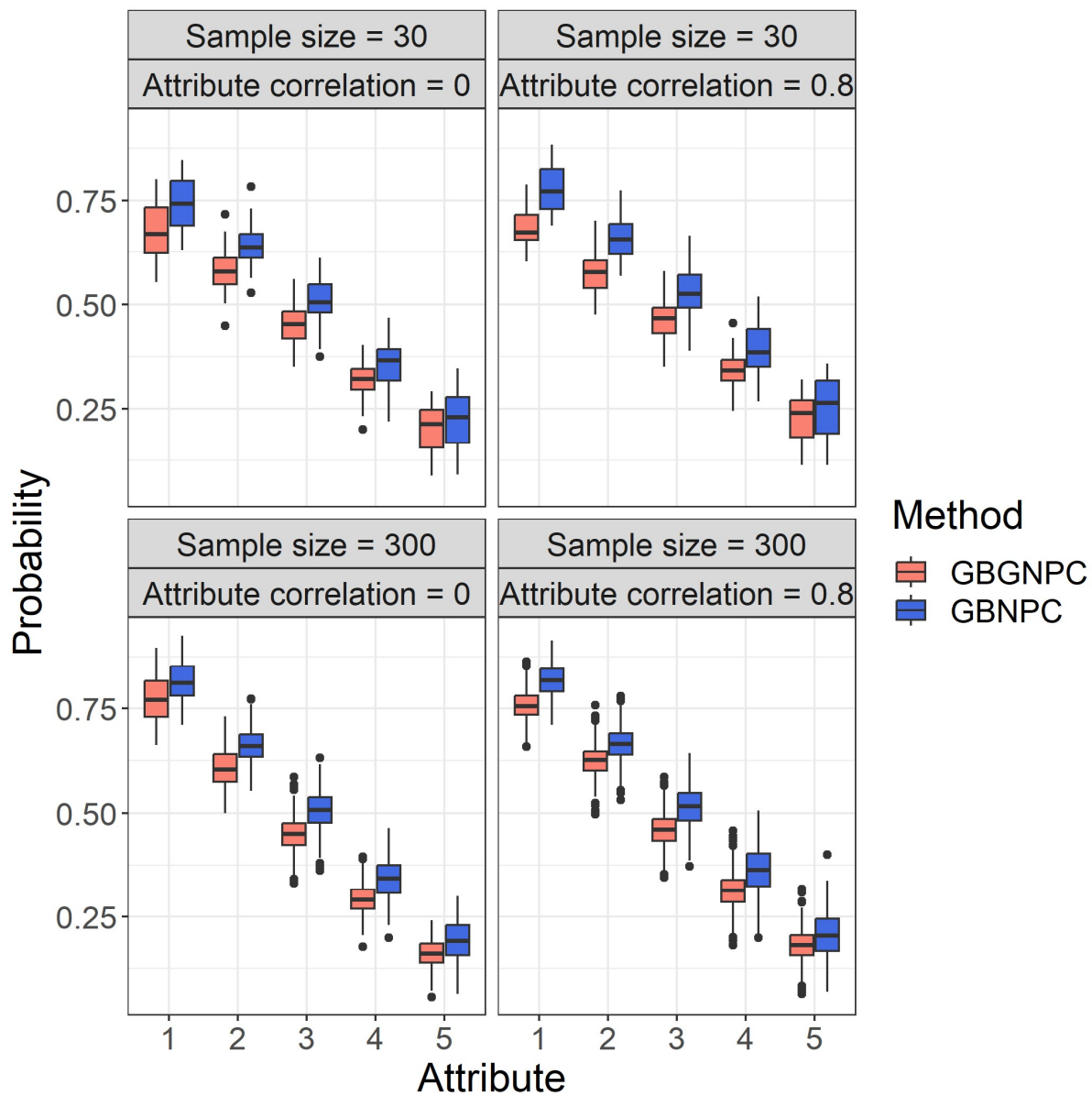


1 Figure 5.  
 2 *Box plots of attribute mastery probabilities of the DINA data generation with four-attribute Q-*  
 3 *matrix conditions*  
 4



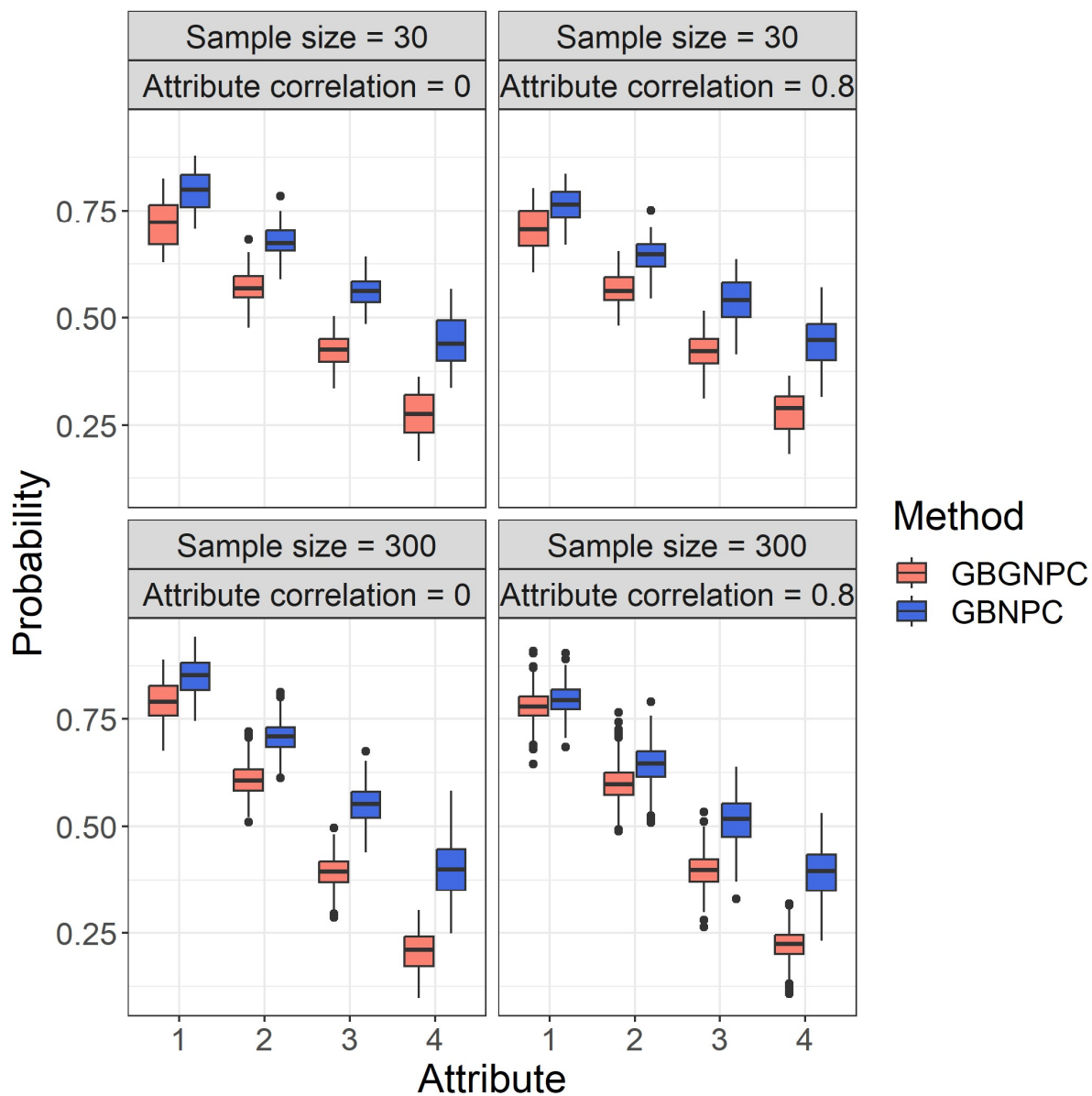
5  
 6

1 Figure 6.  
 2 Box plots of attribute mastery probabilities of the DINA data generation with five-attribute  $Q$ -  
 3 matrix conditions



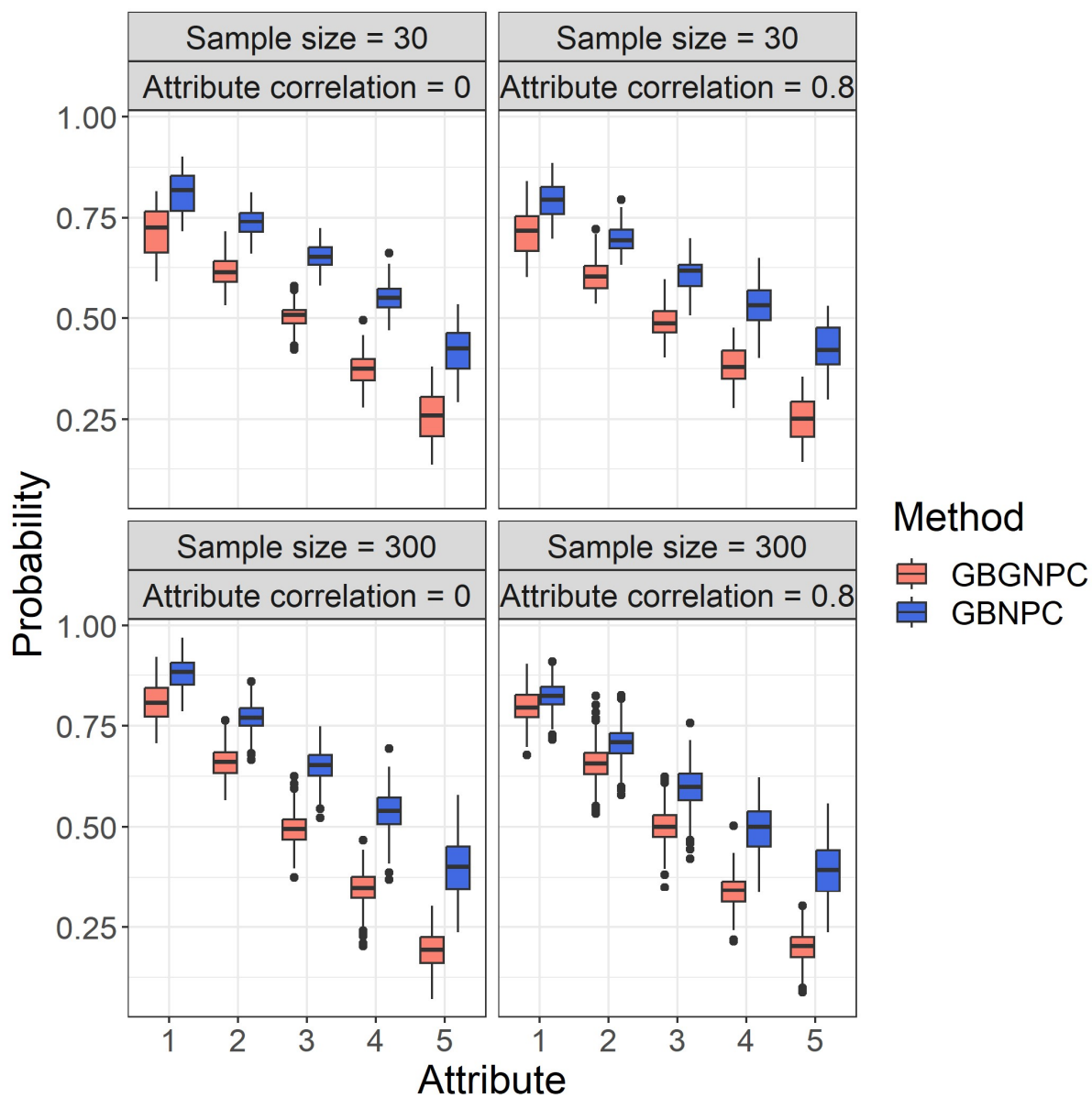
4  
 5  
 6

1 Figure 7.  
 2 *Box plots of attribute mastery probabilities of the general data generation with four-attribute Q-*  
 3 *matrix conditions*



4  
 5  
 6  
 7  
 8

1 Figure 8.  
2 *Box plots of attribute mastery probabilities of the general data generation with five-attribute Q-*  
3 *matrix conditions*



4  
5  
6

## Appendix

### Proofs of Theorems 1 and 2

#### *Preparation for the Proofs*

In this Appendix, we provide some basic tools and introduce helpful notations for the proofs of Theorems 1 and 2. The proofs are presented in the subsequent sections.

Motivated by the constraint (34), we introduce the concept of a "local" latent class at the item level. Considering item  $j$  with  $q$ -vector  $\mathbf{q}_j$ , the constraint (34) divides the collection of attribute mastery profiles  $\alpha$ , which is  $\{0,1\}^K$ , based on an equivalence relationship where  $\alpha_l \sim_j \alpha_{l'}$  is defined by  $\alpha_l \circ \mathbf{q}_j = \alpha_{l'} \circ \mathbf{q}_j$ ; here the subscript  $\sim_j$  emphasizes that the equivalence relationship is determined by the  $j$ -th item  $\mathbf{q}_j$ . On this basis, we introduce a function  $\xi: \{0,1\}^K \times \{0,1\}^K \rightarrow \mathbb{N}$  where  $\xi(\mathbf{q}_j, \alpha_l) = \xi(\mathbf{q}_j, \alpha_{l'})$  is equivalent to  $\alpha_l \circ \mathbf{q}_j = \alpha_{l'} \circ \mathbf{q}_j$ . This function assigns numbers to the equivalent classes induced by item  $j$  based on specific rules. In the following context, we refer to  $\xi(\mathbf{q}_j^0, \alpha)$  as the local latent class of  $\alpha$  induced by item  $j$ . It is straightforward to verify that the number of local latent classes induced by item  $j$ , denoted by  $|\xi(\mathbf{q}_j, \{0,1\}^K)|$ , is equal to  $L_j = 2^{K_j}$ . Here,  $K_j = \sum_{k=1}^K q_{jk}^0$  represents the number of latent attributes required for item  $j$ ; consequently, the range of function  $\xi$  satisfies  $\xi(\mathbf{q}_j, \{0,1\}^K) = [L_j] := \{1, \dots, L_j\}$ . As the local latent classes are identified up to permutations on  $[L_j]$ , owing to their categorical nature, the mapping rules between  $\xi(\mathbf{q}_j, \{0,1\}^K)$  and  $[L_j]$  need not be completely specified in our discussion.

For brevity, we use the general notation  $\mathbf{Z} = (z_{ij})$  to denote the collection of local latent classes for all items  $j \in [J]$  and subjects  $i \in [I]$ , where  $z_{ij}$  represents  $\xi(\mathbf{q}_j^0, \alpha_i)$ . Given that  $\xi(\mathbf{q}_j^0, \alpha_l) = \xi(\mathbf{q}_j^0, \alpha_{l'})$  implies  $\theta_{j,\alpha_l} = \theta_{j,\alpha_{l'}}$  by the definition of  $\xi$ , we express  $\theta_{j,\alpha_i}$  as  $\theta_{j,z_{ij}}$  to directly incorporate the constraint (34) into the loss function (31). For notational simplicity, we may write  $\theta_{j,z_{ij}}$  as  $\theta_{j,z_i}$ . Consequently, we define

$$P_{ij} = P(X_{ij} = 1) = \theta_{j,z_i}^0. \quad (\text{A.1})$$

Then, the loss function (31) can be rewritten as:



$$\mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} \mid X) = \sum_{i=1}^I \left( h(\boldsymbol{\pi}_{\alpha_i}) + \sum_{j=1}^J \ell(X_{ij}, \theta_{j,z_i}) \right) + \sum_{ja} \log f_{ja}(\theta_{ja}) + \sum_{\alpha} \log g_{\alpha}(\boldsymbol{\pi}_{\alpha}), \quad (\text{A.2})$$

1 where  $a \in [L_j]$ . Observe that  $X_{ij}^2 = X_{ij}$ , and  $\mathbb{E}[X_{ij}] = P_{ij}$ , we denote the expectation of the  
 2 above  $\mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} \mid X)$  by  $\bar{\mathcal{L}}(\mathcal{A}, \Theta, \boldsymbol{\pi}) := \mathbb{E}[\mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} \mid X)]$ .

3 Notably,  $\mathbf{Z} = (z_{ij})$  is determined only by  $\mathcal{A}$  because  $\mathbf{Q}^0$  is known. In the subsequent  
 4 context, the quantities determined by the latent attribute profiles  $\mathcal{A}$  are sometimes denoted by  
 5 the superscript  $\mathcal{A}$  to emphasize their relationships with  $\mathcal{A}$ . Considering an arbitrary  $\mathcal{A}$ , we  
 6 denote it as

$$\mathcal{L}(\mathcal{A}) = \inf_{\Theta, \boldsymbol{\pi}} \mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi}^{(\mathcal{A})} \mid X) = \mathcal{L}(\mathcal{A}, \hat{\Theta}^{(\mathcal{A})}, \hat{\boldsymbol{\pi}}^{(\mathcal{A})} \mid X), \quad (\text{A.3})$$

$$\bar{\mathcal{L}}(\mathcal{A}) = \bar{\mathcal{L}}(\mathcal{A}, \bar{\Theta}^{(\mathcal{A})}, \hat{\boldsymbol{\pi}}^{(\mathcal{A})}), \quad (\text{A.4})$$

7 where  $(\hat{\Theta}^{(\mathcal{A})}, \hat{\boldsymbol{\pi}}^{(\mathcal{A})}) := \operatorname{argmin}_{\Theta, \boldsymbol{\pi}} \mathcal{L}(\mathcal{A}, \Theta, \boldsymbol{\pi} \mid X)$  and the definition of  $\bar{\Theta}^{(\mathcal{A})}$  is provided later.  
 8 Notably,  $(\bar{\Theta}^{(\mathcal{A})}, \hat{\boldsymbol{\pi}}^{(\mathcal{A})})$  may not minimize  $\bar{\mathcal{L}}(\mathcal{A}, \Theta, \boldsymbol{\pi})$  for a given  $\mathcal{A}$ . Then, under any  
 9 realization of  $\mathcal{A}$ , if the prior distribution of  $\Theta$  is uniform, the following equations hold for any  
 10 local latent class  $a \in [L_j]$ :

$$\hat{\theta}_{ja}^{(\mathcal{A})} = \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^{(\mathcal{A})} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^{(\mathcal{A})} = a\}}, \quad \bar{\theta}_{ja}^{(\mathcal{A})} = \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^{(\mathcal{A})} = a\} P_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^{(\mathcal{A})} = a\}}. \quad (\text{A.5})$$

11 To derive (A.5), note the sum  $\sum_{j=1}^J \sum_{i=1}^I \ell(X_{ij}, \theta_{j,z_i})$  equals the sum  
 12  $\sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \ell(X_{ij}, \theta_{ja})$ . When estimating  $\hat{\theta}_{ja}$ , we focus on minimizing  $\sum_{z_i=a} \ell(X_{ij}, \theta_{ja})$ .  
 13 By substituting  $\mathbb{E}[X_{ij}] = P_{ij}$  into (A.5), we find that  $\mathbb{E}[\hat{\theta}_{ja}] = \bar{\theta}_{ja}$  holds for any  $(j, a)$ . In the  
 14 following section, we use the second formula in (A.5) to define  $\bar{\Theta}^{(\mathcal{A})}$  given in (A.4). When the  
 15 prior distribution is not a uniform distribution,  $\hat{\theta}_{ja}$  is obtained from minimizing  
 16  $\sum_{z_i=a} \ell(X_{ij}, \theta_{ja}) + \log f_{ja}(\theta_{ja})$ , where  $f_{ja}(\theta_{ja})$  is the prior density of  $\theta_{ja}$ . To avoid

1 ambiguity, we denote  $\hat{\theta}_{ja} := \operatorname{argmin}_{\theta} \sum_{z_i=a} \ell(X_{ij}, \theta_{ja}) + \log f_{ja}(\theta_{ja})$ , and  $\tilde{\theta}_{ja} :=$   
 2  $\operatorname{argmin}_{\theta} \sum_{z_i=a} \ell(X_{ij}, \theta_{ja})$ . It is clear that  $\mathbb{E}[\hat{\theta}_{ja}] = \tilde{\theta}_{ja}$ .

3 Before discussing the details of our proof, we provide technical remarks to simplify the  
 4 discussion. Notably, although we assume that the latent attributes have proportion parameters  
 5  $\boldsymbol{\pi}^0$ , they are still treated as unknown but fixed parameters that need to be estimated. As all the  
 6 proportion parameters  $\pi_{\alpha}^0$  are strictly greater than zero, with the probability converging to  
 7 1,  $\epsilon_1 > 0$  exists such that  $\min_{\alpha} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i = \boldsymbol{\alpha}\} \geq I\epsilon_1$ . Subsequently, we use this fact  
 8 interchangeably with the first condition in Assumption 2.

9 The second point concerns the compact parameter space specified in Assumptions 2 and  
 10 3. Some loss functions may exhibit unusual behavior near the boundary of the parameter space.  
 11 Although Assumption 2 confines the true item parameters to a compact subset within  $(0,1)$ , the  
 12 estimated item responses can still approach zero or one, making theoretical analysis more  
 13 difficult. For any pair  $(j, a)$ ,  $\theta_{ja}$  lies within  $[\delta_2, 1 - \delta_2]$ . We add a condition to  $\mathcal{A}$ , stating that  
 14 there exists an  $\epsilon_2 > 0$  such that for each  $\boldsymbol{\alpha}$ , the sum  $\sum_{i=1}^I \mathcal{J}\{\hat{\boldsymbol{\alpha}}_i = \boldsymbol{\alpha}\}$  is at least  $I\epsilon_2$ . With a  
 15 probability approaching one, this constraint is satisfied by the true latent attribute mastery  
 16 patterns  $\mathcal{A}^0$ . With this constraint, for any pair  $(j, a)$ , the probability that  $|\hat{\theta}_{ja} - \theta_{ja}|$  exceeds  $t$   
 17 can be bounded by  $2\exp(-I\epsilon_2 t^2)$  using Hoeffding's inequality. Thus, the probability that  
 18  $\max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}|$  exceeds  $t$  is less than  $J2^{K+1}\exp(-I\epsilon_2 t^2)$ . Based on the scaling condition in  
 19 Theorem 1,  $\max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}| = o_p(1)$ , implying that with probability converging to 1, all the  
 20  $\hat{\theta}_{ja}$  values fall within  $[\delta_2/2, 1 - \delta_2/2]$ . Based on this result, we assume that in the later content,  
 21 the estimators  $(\hat{\mathcal{A}}, \hat{\Theta})$  are obtained by minimizing the total loss (A.2), under the constraints that  
 22  $\min_{\alpha} \sum_i \mathcal{J}\{\hat{\boldsymbol{\alpha}}_i = \boldsymbol{\alpha}\} \geq I\epsilon_2$  and  $\hat{\Theta} \subset [\delta_3, 1 - \delta_3]$ , for two small positive constants  $\epsilon_2, \delta_3 > 0$ .

23 The third comment concerns how to quantify the effect of prior density  $f_{ja}$  on the  
 24 corresponding estimator  $\hat{\theta}_{ja}$ . Actually, under the smoothness and shape constraints given in  
 25 Assumption 5 and Assumption 3, the additional term  $\log f_{ja}(\theta_{ja})$  might cause the estimator  
 26  $\hat{\theta}_{ja}$  to have a  $O_p(1/\sqrt{I})$  level drift from the sample average form  $\tilde{\theta}_{ja}$  given in (A.6). By  
 27 considering the Taylor expansion formula, we have

$$\begin{aligned}
& \log f_{ja}(\theta_{ja}) + \sum_{z_{ij}=a} \ell(X_{ij}, \theta_{ja}) \\
&= \log f_{ja}(\theta_{ja}) + \sum_{z_{ij}=a} \ell(X_{ij}, \hat{\theta}_{ja}) + \left( \sum_{z_{ij}=a} \partial_{\theta} \ell(X_{ij}, \hat{\theta}_{ja}) \right) (\theta_{ja} - \hat{\theta}_{ja}) + \frac{1}{2} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\theta_{ja} - \hat{\theta}_{ja})^2 \\
&= \log f_{ja}(\theta_{ja}) + \sum_{z_{ij}=a} \ell(X_{ij}, \hat{\theta}_{ja}) + \frac{1}{2} \sum_{z_{ij}=a} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\theta_{ja} - \hat{\theta}_{ja})^2,
\end{aligned}$$

- 1 where  $\hat{\theta}_{ja}$  is between  $\hat{\theta}_{ja}$  and  $\theta_{ja}$  according to the mean value theorem; the second equality  
2 holds owing to Assumption 3. According to the above equation, we can find that  $\hat{\theta}_{ja} =$   
3  $\operatorname{argmin}_{\theta \in [\delta_3, 1-\delta_3]} \log f_{ja}(\theta) + \sum_{z_{ij}=a} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\theta - \hat{\theta}_{ja})^2 / 2$ . Note that if we take  $\theta =$   
4  $\hat{\theta}_{ja}$ ,  $\log f_{ja}(\theta) + \sum_{z_{ij}=a} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\theta - \hat{\theta}_{ja})^2 / 2 = \log f_{ja}(\hat{\theta}_{ja})$ . Based on Assumption 5,  
5 there exists a constant  $C > 0$  such that  $|\log f_{ja}(\hat{\theta}_{ja})| \leq \sup_{\theta \in [\delta_3, 1-\delta_3]} |\log f_{ja}(\theta)| < C$ ,  
6 implying the following:

$$\begin{aligned}
2C &\geq \frac{1}{2} \sum_{z_{ij}=a} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\hat{\theta}_{ja} - \hat{\theta}_{ja})^2 \\
&\geq \frac{b_L}{2} \sum_{z_{ij}=a} (\hat{\theta}_{ja} - \hat{\theta}_{ja})^2 = \frac{b_L i_{ja}}{2} (\hat{\theta}_{ja} - \hat{\theta}_{ja})^2,
\end{aligned}$$

- 7 where  $i_{ja} := \sum_{i=1}^I \mathcal{J}\{z_{ij}^{(\mathcal{A})} = a\}$ . Thus, a constant  $\hat{C} > 0$  exists such that for any pair  $(j, a)$ , we  
8 have

$$(\hat{\theta}_{ja} - \hat{\theta}_{ja})^2 \leq \hat{C} i_{ja}^{-1}.$$

- 9 This inequality will be used several times afterwards. In the theoretical analysis of the estimators  
10 above, uniform bounds related to the quantities of element-wise loss  $\ell(\cdot, \cdot)$  and prior densities  
11  $f_{ja}$  are frequently used. The existence of these uniform bounds requires restricting the parameter  
12 space of the item response parameters to a compact subspace. Therefore, discussing the compact  
13 parameter subspace of the item parameters is necessary.

14 **Outline of the first half of the proof**

1 **Step 1:** Express the upper bound of  $|\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})|$  in terms of  $(b/2) \cdot \left(\sum_{j=1}^J \sum_{a=1}^{L_j} i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2\right) + |\mathbb{E}[Y] - Y| + O_p(J)$ , where  $Y := \sum_i \sum_j \ell(X_{ij}, \bar{\theta}_{j,z_i}^{(\mathcal{A})})$  depending on  $Y$  and  $\bar{\Theta}^{(\mathcal{A})}$  under  $\mathcal{A}, b$  is the upper bound of the second order derivative of the  $\ell(\cdot, \cdot)$ .

2 **Step 2:** Bound  $\sum_j \sum_i i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2$  and  $|Y - \mathbb{E}[Y]|$  separately to obtain a uniform convergence rate  $\sup_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| = o_p(\delta_{IJ})$

3 **Step 3:** Based on the definition of  $\hat{\mathcal{A}}$ , it follows that  $0 \leq \bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0) \leq 2 \sup_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| = o_p(\delta_{IJ})$ , which controls the deviation  $\bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0)$ .

4 In some classical statistical inference contexts, consistent results for the parameters of interest are typically established through the uniform convergence of random functions associated with these parameters. For instance, if  $\sup_{\theta \in \Theta} |\hat{\ell}(\theta) - \ell(\theta)| \xrightarrow{P} 0$ , and if we further assume that  $\ell$  has a unique minimum  $\hat{\theta}$  on  $\Theta$ ,  $\operatorname{argmin}_{\theta} \hat{\ell}(\theta) =: \hat{\theta} \xrightarrow{P} \hat{\theta}$  under some regularity conditions. The regular conditions may vary across settings. Considering  $\mathcal{A}$  as the parameter to be estimated, the primary aim of the first three steps is to demonstrate that  $\mathcal{A}$  minimizes the expected loss and establishes a uniform convergence result for its random loss function of  $\mathcal{A}$ .

#### 15 *Outline of the second half of the proof*

16 **Step 4:** Define  $I_{a,b}^j = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} \mathcal{J}\{z_{ij} = b\}$ ,  $a, b \in [L_j]$  to represent the samples with the wrong local latent class assignments. Derive some upper bounds for the quantities based on  $I_{a,b}^j$  using  $\bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0)$  with the help of the identification assumptions.

17 **Step 5:** Bound the  $\sum_{i=1}^I \mathcal{J}\{\hat{\alpha}_i \neq \alpha^0\}$  using the quantities based on  $I_{a,b}^j$  with the help of the discrete structure of the Q-matrix, then obtain the desired classification error rate.

18 Assumptions 1-5 are the regularity conditions for achieving clustering consistency based on the uniform convergence results established in the first half of the proof. We have provided further details regarding the assumptions in later proofs.

24

#### 25 *First Half of the Proof of Theorem 1*

1 **Step 1.** The idea of decomposing  $\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})$  is to consider

$$\ell(X_{ij}, \hat{\theta}_{j,z_i}) - \mathbb{E}[\ell(X_{ij}, \bar{\theta}_{j,z_i})] = \left( \ell(X_{ij}, \hat{\theta}_{j,z_i}) - \ell(X_{ij}, \bar{\theta}_{j,z_i}) \right) + \left( \ell(X_{ij}, \bar{\theta}_{j,z_i}) - \mathbb{E}[\ell(X_{ij}, \bar{\theta}_{j,z_i})] \right).$$

2 The variability in the first term of the right-hand side mainly emerges from the fluctuation in  
 3  $|\hat{\theta}_{ja} - \bar{\theta}_{ja}|$ , while the randomness in the second term is attributable to the stochastic nature of  
 4  $X_{ij}$ .

5 **Lemma 1.** Let  $(X_{ij}; 1 \leq i \leq I, 1 \leq j \leq J)$  denote independent Bernoulli trials with parameters  
 6  $(P_{ij}; 1 \leq i \leq I, 1 \leq j \leq J)$ . In a general latent class model, given arbitrary latent attribute  
 7 mastery patterns  $\mathcal{A}$ ,

$$|\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| \leq \frac{b}{2} \cdot \left( \sum_{j=1}^J \sum_{a=1}^{L_j} i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \right) + |Y - \mathbb{E}(Y)| + O_p(J), \quad (\text{A.6})$$

8 where  $Y = \sum_{j=1}^J \sum_{i=1}^I \ell(X_{ij}, \bar{\theta}_{j,z_i})$  is a random variable depending on  $\mathcal{A}$  and  $L_j$  denotes the  
 9 number of the distinct local latent classes induced by  $\mathbf{q}_j$  for item  $j$ .

10 **Proof.** By noting the decomposition that we mentioned at the beginning of Step 1,  $|Y - \mathbb{E}[Y]|$  is  
 11 easy to check. It is sufficient for us to prove that

$$0 \leq \sum_i \sum_j \left( \ell(X_{ij}, \bar{\theta}_{j,z_i}) - \ell(X_{ij}, \hat{\theta}_{j,z_i}) \right) \leq \frac{b}{2} \cdot \left( \sum_{j=1}^J \sum_{a=1}^{L_j} i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \right) + O_p(J).$$

12 The first inequality is clear by the definition of  $\hat{\theta}_{j,z_i}$  (minimizing the loss). For the second part,  
 13 using the mean-value theorem for second-order derivatives, we obtain

$$\begin{aligned}
& \sum_i \sum_j \left( \ell(X_{ij}, \tilde{\theta}_{j,z_i}) - \ell(X_{ij}, \hat{\theta}_{j,z_i}) \right) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \ell(X_{ij}, \tilde{\theta}_{ja}) - \ell(X_{ij}, \hat{\theta}_{ja}) \right) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \partial_{\theta} \ell(X_{ij}, \hat{\theta}_{ja}) (\tilde{\theta}_{ja} - \hat{\theta}_{ja}) + \frac{1}{2} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\hat{\theta}_{ja} - \tilde{\theta}_{ja})^2 \right). \quad (\text{A.7}) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \frac{1}{2} \partial_{\theta^2} \ell(X_{ij}, \hat{\theta}_{ja}) (\tilde{\theta}_{ja} - \hat{\theta}_{ja})^2 \right) \\
&\leq \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \frac{b_U}{2} (\tilde{\theta}_{ja} - \hat{\theta}_{ja})^2 \right)
\end{aligned}$$

- 1 where  $\tilde{\theta}_{ja}$  is between  $\hat{\theta}_{ja}$  and  $\tilde{\theta}_{ja}$  according to the mean value theorem. The third equality
- 2 holds true since by Assumption 3, we have

$$\sum_{z_i=a} \partial_{\theta} \ell(X_{ij}, \hat{\theta}_{ja}) = 0.$$

- 3 Similarly, using  $(\hat{\theta}_{ja} - \tilde{\theta}_{ja})^2 \leq \hat{C} i_{ja}^{-1}$ , we have

$$\begin{aligned}
& \sum_i \sum_j \left( \ell(X_{ij}, \hat{\theta}_{j,z_i}) - \ell(X_{ij}, \tilde{\theta}_{j,z_i}) \right) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \ell(X_{ij}, \hat{\theta}_{ja}) - \ell(X_{ij}, \tilde{\theta}_{ja}) \right) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \partial_{\theta} \ell(X_{ij}, \tilde{\theta}_{ja}) (\hat{\theta}_{ja} - \tilde{\theta}_{ja}) + \frac{1}{2} \partial_{\theta^2} \ell(X_{ij}, \tilde{\theta}_{ja}) (\tilde{\theta}_{ja} - \hat{\theta}_{ja})^2 \right) \\
&= \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \frac{1}{2} \partial_{\theta^2} \ell(X_{ij}, \tilde{\theta}_{ja}) (\tilde{\theta}_{ja} - \hat{\theta}_{ja})^2 \right) \\
&\leq \sum_{j=1}^J \sum_{a=1}^{L_j} \sum_{z_i=a} \left( \frac{b_U}{2} (\hat{\theta}_{ja} - \tilde{\theta}_{ja})^2 \right) = \sum_{j=1}^J \sum_{a=1}^{L_j} \frac{b_U \hat{C}}{2} \leq (b_U \hat{C} 2^K) J,
\end{aligned}$$

1 which concludes the proof of this lemma. ■

2 **Lemma 2.** *The following event happens with a probability of at least  $1 - \delta$ ,*

$$\max_{\mathcal{A}} \left\{ \sum_{j=1}^J \sum_{a=1}^{L_j} i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \right\} < \frac{1}{2} \left( I \log 2^K + J 2^K \log \left( \frac{I}{2^K} + 1 \right) - \log \delta \right).$$

3 **Proof.** Under any realization of  $\mathcal{A}$ , each  $\hat{\theta}_{ja}$  is an average of  $i_{ja}$  independent Bernoulli  
4 random variables  $x_{1j}, \dots, x_{i_{ja}j}$  with mean  $\bar{\theta}_{ja}$ . By applying the Hoeffding inequality, we have

$$P(\hat{\theta}_{ja} \geq \bar{\theta}_{ja} + t) \leq \exp(-2i_{ja}t^2), \quad P(\hat{\theta}_{ja} \leq \bar{\theta}_{ja} - t) \leq \exp(-2i_{ja}t^2). \quad (\text{A.8})$$

5 Notably, considering a fixed  $\mathcal{A}$ , each  $\hat{\theta}_{ja}$  can take values only in the finite set  
6  $\{0, 1/i_{ja}, 2/i_{ja}, \dots, 1\}$  of cardinality  $i_{ja} + 1$ . We denote this range of  $\hat{\theta}_{ja}$  by  $\hat{\Theta}^{ja}$  and the range

7 of the matrix  $\hat{\Theta} = (\hat{\theta}_{ja})$  by  $\hat{\Theta}$ . Subsequently,  $P(\hat{\theta}_{ja} = v) \leq \exp(-2i_{ja}(v - \bar{\theta}_{ja})^2)$  for any

8  $v \in \hat{\Theta}^{ja}$ . As each of the  $J \times 2^K$  entries in  $\hat{\Theta}$ ,  $\hat{\theta}_{ja}$  can independently take on  $i_{ja} + 1$  different  
9 values, there is  $|\hat{\Theta}| = \prod_j \prod_{a=1}^{L_j} (i_{ja} + 1)$  with constraint  $\sum_{a=1}^{L_j} i_{ja} = I$ . As  $L_j = 2^{Kj} \leq 2^K$ , we

10 have  $\prod_{a=1}^{L_j} (i_{ja} + 1) \leq (1 + I/2^K)^{2^K}$ . Denote  $\hat{\Theta}_\epsilon = \{\hat{\Theta} \in \hat{\Theta} : \sum_j \sum_a i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \geq \epsilon\}$ ,  $\hat{\Theta}_\epsilon \subseteq$

11  $\hat{\Theta}$ , and

$$\begin{aligned} & P\left(\sum_j \sum_{a=1}^{L_j} i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \geq \epsilon\right) = \sum_{\hat{\Theta} \in \hat{\Theta}_\epsilon} P(\hat{\Theta} = \hat{\Theta}) \\ & \leq \sum_{\hat{\Theta} \in \hat{\Theta}_\epsilon} \prod_j \prod_a \exp(-2i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2) \\ & = \sum_{\hat{\Theta} \in \hat{\Theta}_\epsilon} \exp\left(-2i_{ja} \sum_j \sum_a (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2\right) \\ & \leq \sum_{\hat{\Theta} \in \hat{\Theta}_\epsilon} \exp(-2\epsilon) \leq |\hat{\Theta}_\epsilon| e^{-2\epsilon}. \end{aligned} \quad (\text{A.9})$$

12 The above result holds for any fixed  $\mathcal{A}$  when we apply a union bound over all the  $(2^K)^J$   
13 possible assignments of  $\mathcal{A}$  to obtain

$$P\left(\max_{\mathcal{A}} \left\{ \sum_j \sum_a i_{ja} (\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \right\} \geq \epsilon\right) \leq 2^{KI} \left(\frac{I}{2^K} + 1\right)^{J2^K} e^{-2\epsilon}. \quad (\text{A.10})$$

1 Take  $\delta = 2^{KI} \left(\frac{I}{2^K} + 1\right)^{J2^K} e^{-2\epsilon}$ ; then,  $2\epsilon = I \log 2^K + J2^K \log(1 + I/2^K) - \log \delta$ . This  
 2 concludes the proof of lemma 2. ■

3 **Lemma 3.** Define the random variable  $Y = \sum_i \sum_j \ell(X_{ij}, \bar{\theta}_{j,z_i}^{(\mathcal{A})})$ , and denote  $Y_{ij} = \ell(X_{ij}, \bar{\theta}_{j,z_i})$ .

4 Note that  $\bar{\theta}_{ja} \in [\delta_2, 1 - \delta_2]$  and  $\ell(\cdot, \cdot)$  are continuous on  $\theta$  in  $(0,1)$ . Since continuous  
 5 functions on the compact set are bounded, a constant  $U > 0$  exists such that  $|\ell(X_{ij}, \bar{\theta}_{j,z_i})| \leq$   
 6  $U, \forall(i, j)$ . By applying Hoeffding's inequality to bound  $|Y - \mathbb{E}[Y]|$  for any realization of  $\mathcal{A}$ , we  
 7 have:

$$P(|Y - \mathbb{E}[Y]| \geq \epsilon) \leq 2 \exp\left\{-\frac{\epsilon^2}{(4U^2)IJ}\right\}. \quad (\text{A.11})$$

8 With the help of Lemma 2 and Lemma 3, subsequently, we prove the following proposition:

9 **Proposition 1.** Under the following scaling for some small positive constant  $c > 0$ ,

$$\sqrt{J} = O(I^{1-c})$$

10 we have  $\max_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| = o_p(\delta_{IJ})$  where  $\delta_{IJ} = I\sqrt{J}(\log J)^\epsilon$  for a small positive  $\epsilon >$   
 11 0.

12 **Proof.** First, note that under the given scaling condition,  $J = o(I\sqrt{J})$ . Combining the results of  
 13 Lemma 2 and Lemma 3, since there are  $(2^K)^I$  possible assignments of  $\mathcal{A}$ , we apply the union  
 14 bound to obtain

15

16

17



$$\begin{aligned}
& P\left(\max_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| \geq 3\epsilon\delta_{IJ}\right) \\
& \leq (2^K)^I P\left[\left\{\sum_j \sum_a i_{ja}(\hat{\theta}_{ja} - \bar{\theta}_{ja})^2 \geq \epsilon\delta_{IJ}\right\} \cup \{|Y - \mathbb{E}[Y]| \geq \epsilon\delta_{IJ}\}\right] \\
& + P\left(\max_{\mathcal{A}} \sum_i \sum_j \left(\ell(X_{ij}, \hat{\theta}_{j,z_i}) - \ell(X_{ij}, \bar{\theta}_{j,z_i})\right) \geq J(\log)^\epsilon\right) \tag{A.12} \\
& \leq \exp\left(I \log(2^K) + J2^K \log\left(\frac{I}{2^K} + 1\right) - 2\epsilon\delta_{IJ}\right) + 2 \exp\left(I \log(2^K) - \frac{\epsilon^2 \delta_{IJ}^2}{4U^2 IJ}\right) \\
& + P\left(\max_{\mathcal{A}} \sum_i \sum_j \left(\ell(X_{ij}, \hat{\theta}_{j,z_i}) - \ell(X_{ij}, \bar{\theta}_{j,z_i})\right) \geq J(\log)^\epsilon\right).
\end{aligned}$$

1 For the third term, note that the following inequality holds for any given  $\mathcal{A}$ :

$$\sum_i \sum_j \left(\ell(X_{ij}, \hat{\theta}_{j,z_i}) - \ell(X_{ij}, \bar{\theta}_{j,z_i})\right) \leq (b_U \hat{C} 2^K)J$$

2 For the second term on the right-hand side of the aforementioned display to converge to zero, we  
3 set  $\delta_{IJ} = I\sqrt{J}(\log J)^\epsilon$  for a small positive constant  $\epsilon$ . Moreover, under this  $\delta_{IJ}$ , for the first  
4 term to converge to zero as  $I, J$  increase, the scaling  $\sqrt{J} = O(I^{1-c})$  given in the theorem results  
5 in  $P(\max_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| \geq \epsilon\delta_{IJ}) = o(1)$ , which implies the result in Proposition 1.

6 ■

7 **Step 3.** (A.5) implies that  $\bar{\theta}_{j,z_i}^{(\mathcal{A}^0)} = P_{ij}$ , which means that if we plug into the true latent attribute  
8 mastery pattern  $\mathcal{A}^0$ , the estimators will be the corresponding true parameters. According to this  
9 property, the following lemma indicates that  $\mathcal{A}^0$  minimizes expected loss.

10 **Lemma 4.** *By Assumption 4,  $\mathbb{E}[\ell(X_{ij}, \bar{\theta}_{j,z_i})] - \mathbb{E}[\ell(X_{ij}, \theta_{j,z_i}^0)] \geq c(\bar{\theta}_{j,z_i} - \theta_{j,z_i}^0)^\eta$  for some*  
11  *$\eta \geq 2, c > 0$ , then we have*

$$\bar{\mathcal{L}}(\mathcal{A}) - \bar{\mathcal{L}}(\mathcal{A}^0) \geq c \cdot \left( \sum_i \sum_j (P_{ij} - \bar{\theta}_{j,z_i})^\eta \right) \geq 0. \quad (\text{A.13})$$

1 Notably, while Lemma 4 holds for any  $\mathcal{A}$ , it also holds for the estimator  $\hat{\mathcal{A}}$ , then

$$0 \leq \bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0) = [\bar{\mathcal{L}}(\hat{\mathcal{A}}) - \mathcal{L}(\hat{\mathcal{A}})] + [\mathcal{L}(\hat{\mathcal{A}}) - \mathcal{L}(\mathcal{A}^0)] + [\mathcal{L}(\mathcal{A}^0) - \bar{\mathcal{L}}(\mathcal{A}^0)]. \quad (\text{A.14})$$

2 As  $\hat{\mathcal{A}} = \operatorname{argmin}_{\mathcal{A}} \mathcal{L}(\mathcal{A})$ , we have  $\mathcal{L}(\hat{\mathcal{A}}) - \mathcal{L}(\mathcal{A}^0) \leq 0$ . Substituting this into A.14, we can  
3 derive that

$$0 \leq \bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0) \leq 2 \sup_{\mathcal{A}} |\bar{\mathcal{L}}(\mathcal{A}) - \mathcal{L}(\mathcal{A})| = o_p(\delta_{IJ})$$

#### 4 **Second Half of the Proof of Theorem 1**

5 By applying Hölder's inequality, we have

$$(IJ)^{1-\frac{\eta}{2}} \left( \sum_i \sum_j (P_{ij} - \bar{\theta}_{j,z_i})^2 \right)^{\frac{\eta}{2}} \leq \sum_i \sum_j (P_{ij} - \bar{\theta}_{j,z_i})^\eta = o_p(\delta_{IJ}).$$

6 By letting  $(IJ)^{1-\eta/2} (S)^{\eta/2} = \delta_{IJ}$ , we can check that  $\sum_i \sum_j (P_{ij} - \bar{\theta}_{j,z_i})^2 = o_p(S)$  where  $S :=$   
7  $I(J)^{1-1/\eta} (\log J)^{2\hat{\epsilon}/\eta}$ . In the following, we derive a lower bound for  $\sum_i \sum_j (P_{ij} - \bar{\theta}_{j,z_i})^2$  because  
8 it is easier to work with than  $(P_{ij} - \bar{\theta}_{j,z_i})^\eta$ .

9 **Step 4.** Motivated by Assumption 2, we define  $\mathcal{J} := \{j \in [J]; \exists k \in [K] \text{ s.t. } \mathbf{q}_j^0 = \mathbf{e}_k\}$ , which  
10 represents the set of all items  $j$  that depend on only one latent attribute. Notably,  $\forall j \in$   
11  $\mathcal{J}, |\{\boldsymbol{\alpha} \circ \mathbf{q}_j^0; \boldsymbol{\alpha} \in \{0,1\}^K\}| = 2$ , as  $\mathbf{q}_j$  only contains one required latent attribute, then  $\xi(\mathbf{q}_j^0, \boldsymbol{\alpha}) \in$   
12  $\{1,2\}$  for all  $j \in \mathcal{J}$ . Without loss of generality, we assume that if  $\boldsymbol{\alpha} \circ \mathbf{q}_j^0 \neq \mathbf{0}$ , then let  
13  $\xi(\mathbf{q}_j^0, \boldsymbol{\alpha}) = 2$ , otherwise, let  $\xi(\mathbf{q}_j^0, \boldsymbol{\alpha}) = 1$ . Further, we assume that  $\theta_{j,2}^0 > \theta_{j,1}^0, \forall j \in \mathcal{J}$ , which  
14 aligns with the concept that subjects possessing the required latent attribute tend to perform  
15 better. For any  $j \in \mathcal{J}$ , define

$$I_{a,b}^j = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} \mathcal{J}\{\hat{z}_{ij} = b\}, \quad (a, b) \in \{1, 2\}^2. \quad (\text{A.15})$$

- 1 Note  $P_{ij} = \mathcal{J}\{z_{ij}^0 = 2\}\theta_{j,2}^0 + \mathcal{J}\{z_{ij}^0 = 1\}\theta_{j,1}^0$  and  $I_{2,2}^j + I_{1,2}^j = \sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 2\}$ ,  $I_{2,1}^j + I_{1,1}^j =$
- 2  $\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 1\}$ . By using (A.5), there are

$$\begin{aligned} \bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} &= \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 2\} P_{ij}}{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 2\}} \\ &= \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 2\} (\mathcal{J}\{z_{ij}^0 = 2\} \theta_{j,2}^0 + \mathcal{J}\{z_{ij}^0 = 1\} \theta_{j,1}^0)}{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = 2\}} \\ &= \frac{I_{2,2}^j \theta_{j,2}^0 + I_{1,2}^j \theta_{j,1}^0}{I_{2,2}^j + I_{1,2}^j}; \\ \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})} &= \frac{I_{2,1}^j \theta_{j,2}^0 + I_{1,1}^j \theta_{j,1}^0}{I_{2,1}^j + I_{1,1}^j}. \end{aligned} \quad (\text{A.16})$$

- 3 Under  $\hat{\mathcal{A}}$ , we impose a natural constraint  $\bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} > \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})}, \forall j \in \mathcal{J}$  on  $\hat{\mathcal{A}}$  for identifiability
- 4 purpose. This constraint does not change the previous results as  $\theta_{j,2}^0 > \theta_{j,1}^0$  allows  $\mathcal{L}(\hat{\mathcal{A}}) -$
- 5  $\mathcal{L}(\mathcal{A}^0) \leq 0$  in (A.14) still holds; thus,  $\bar{\mathcal{L}}(\hat{\mathcal{A}}) - \bar{\mathcal{L}}(\mathcal{A}^0) = o_p(\delta_{IJ})$  still holds under this
- 6 constraint. Combining  $\bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} > \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})}$  and  $\theta_{j,2}^0 > \theta_{j,1}^0$ , there is

$$\begin{aligned} \bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} > \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})} &\Leftrightarrow (I_{2,2}^j I_{1,1}^j - I_{1,0}^j I_{0,1}^j) \theta_{j,2}^0 > (I_{2,2}^j I_{1,1}^j - I_{1,0}^j I_{0,1}^j) \theta_{j,1}^0 \\ &\Leftrightarrow I_{2,2}^j I_{1,1}^j > I_{2,1}^j I_{1,2}^j. \end{aligned} \quad (\text{A.17})$$

- 7 From (A.15), we can obtain

$$\begin{aligned} \left| \theta_{j,1}^0 - \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})} \right| &= \frac{I_{2,1}^j (\theta_{j,2}^0 - \theta_{j,1}^0)}{I_{2,1}^j + I_{1,1}^j}, \quad \left| \theta_{j,2}^0 - \bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} \right| = \frac{I_{1,2}^j (\theta_{j,2}^0 - \theta_{j,1}^0)}{I_{2,2}^j + I_{1,2}^j}, \\ \left| \theta_{j,2}^0 - \bar{\theta}_{j,1}^{(\hat{\mathcal{A}})} \right| &= \frac{I_{1,1}^j (\theta_{j,2}^0 - \theta_{j,1}^0)}{I_{2,1}^j + I_{1,1}^j}, \quad \left| \theta_{j,1}^0 - \bar{\theta}_{j,2}^{(\hat{\mathcal{A}})} \right| = \frac{I_{2,2}^j (\theta_{j,2}^0 - \theta_{j,1}^0)}{I_{2,2}^j + I_{1,2}^j}. \end{aligned}$$

1 Therefore,

$$\begin{aligned}
& \sum_{j=1}^J \sum_{i=1}^I (P_{ij} - \bar{\theta}_{j,\hat{z}_i})^2 \geq \sum_{j \in \mathcal{J}} \sum_{i=1}^I (P_{ij} - \bar{\theta}_{j,\hat{z}_i})^2 \\
&= \sum_{j \in \mathcal{J}} \left( I_{1,1}^j (\theta_{j,1}^0 - \bar{\theta}_{j,1}^{(\mathcal{A})})^2 + I_{2,1}^j (\theta_{j,2}^0 - \bar{\theta}_{j,1}^{(\mathcal{A})})^2 + I_{1,2}^j (\theta_{j,1}^0 - \bar{\theta}_{j,2}^{(\mathcal{A})})^2 + I_{2,2}^j (\theta_{j,2}^0 - \bar{\theta}_{j,2}^{(\mathcal{A})})^2 \right) \\
&= \sum_{j \in \mathcal{J}} \left( \frac{I_{1,1}^j (I_{2,1})^2 + I_{2,1}^j (I_{1,1})^2}{(I_{2,1}^j + I_{1,1}^j)^2} + \frac{I_{1,2}^j (I_{2,2})^2 + I_{2,2}^j (I_{1,2})^2}{(I_{2,2}^j + I_{1,2}^j)^2} \right) (\theta_{j,2}^0 - \theta_{j,1}^0)^2 \\
&= \sum_{j \in \mathcal{J}} \left( \frac{I_{2,1}^j I_{1,1}^j}{I_{2,1}^j + I_{1,1}^j} + \frac{I_{2,2}^j I_{1,2}^j}{I_{2,2}^j + I_{1,2}^j} \right) (\theta_{j,2}^0 - \theta_{j,1}^0)^2 \tag{A.18} \\
&\geq \delta \sum_{j \in \mathcal{J}} \left( \frac{I_{2,1}^j I_{1,1}^j}{I_{2,1}^j + I_{1,1}^j} + \frac{I_{2,2}^j I_{1,2}^j}{I_{2,2}^j + I_{1,2}^j} \right) \\
&\geq \frac{1}{2} \delta \sum_{j \in \mathcal{J}} (\min\{I_{2,1}^j, I_{1,1}^j\} + \min\{I_{2,2}^j, I_{1,2}^j\}).
\end{aligned}$$

2 The second inequality holds since by Assumption 1,  $(\theta_{j,2}^0 - \theta_{j,1}^0)^2 \geq \delta$ . One ideal scenario is  
3 that for most  $j \in \mathcal{J}$ ,  $\min\{I_{2,1}^j, I_{1,1}^j\} + \min\{I_{2,2}^j, I_{1,2}^j\} = I_{2,1}^j + I_{1,2}^j = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 \neq \hat{z}_{ij}\}$ ; thus the  
4 misclassification error for the local latent classes could be bounded relatively tight. The  
5 following result confirms this intuition.

6 **Lemma 5.** Define the following random set depending on the estimated latent attribute mastery  
7 patterns  $\hat{\mathcal{A}}$  under constraint  $\bar{\theta}_{j,2}^{(\mathcal{A})} > \bar{\theta}_{j,1}^{(\mathcal{A})}, \forall j \in \mathcal{J}$ :

$$\mathcal{J}_0 = \{j \in \mathcal{J}; I_{2,1}^j < I_{1,1}^j, I_{1,2}^j < I_{2,2}^j\};$$

$$\mathcal{J}_1 = \{j \in \mathcal{J}; I_{2,1}^j < I_{1,1}^j, I_{1,2}^j > I_{2,2}^j\};$$

$$\mathcal{J}_2 = \{j \in \mathcal{J}; I_{2,1}^j > I_{1,1}^j, I_{1,2}^j < I_{2,2}^j\};$$

8 then under Assumption 1 and Assumption 2, there are  $|\mathcal{J}_1| = o_p(S/I), |\mathcal{J}_2| = o_p(S/I)$ .

9 **Proof.** If  $j \in \mathcal{J}_1$ ,  $\min\{I_{2,1}^j, I_{1,1}^j\} + \min\{I_{2,2}^j, I_{1,2}^j\} = I_{2,1}^j + I_{2,2}^j = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = 2\}$ . Under  
10 Assumption 2,

$$\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = 2\} \geq I\epsilon$$

1 then

$$\begin{aligned} & P\left(|\mathcal{J}_1| \geq \frac{S}{\delta I}\right) \\ & \leq P\left(\sum_{j \in \mathcal{J}_1} I_{2,1}^j + I_{2,2}^j \geq \frac{S}{\delta I} \cdot I\epsilon\right) \\ & \leq P\left(\sum_i \sum_j (P_{ij} - \bar{\theta}_{j,\hat{z}_i})^2 \geq \frac{\epsilon S}{2}\right). \end{aligned}$$

2 By noting  $\sum_i \sum_j (P_{ij} - \bar{\theta}_{j,\hat{z}_i})^2 = o_p(S)$ , then  $|\mathcal{J}_1| = o_p(S/I)$ . Similar arguments yield  $|\mathcal{J}_2| =$   
3  $o_p(S/I)$ , which concludes the proof of Lemma 5. ■

4 Note (A.17) implies that  $\min\{I_{2,1}^j, I_{1,1}^j\} + \min\{I_{2,2}^j, I_{1,2}^j\} \neq I_{1,1}^j + I_{2,2}^j, \forall j \in \mathcal{J}$ , thus  $\mathcal{J} =$   
5  $\mathcal{J}_0 \cup \mathcal{J}_1 \cup \mathcal{J}_2$ . Lemma 5 implies that when  $\delta_j$  goes to 0 with a mild rate, the number of elements  
6 in  $\mathcal{J}_0$  dominates the number of elements in  $\mathcal{J}_1 \cup \mathcal{J}_2$ ; thus for most  $j \in \mathcal{J}$ ,  $\min\{I_{2,1}^j, I_{1,1}^j\} +$   
7  $\min\{I_{2,2}^j, I_{1,2}^j\}$  should be  $I_{2,1}^j + I_{1,2}^j = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 \neq \hat{z}_{ij}\}$ , which represents the number of subjects  
8 with the incorrectly assigned local latent classes.

9 **Step 5.** (A.18) implies that  $\sum_i \sum_j (P_{ij} - \bar{\theta}_{j,\hat{z}_i})^2 \geq \delta \sum_{j \in \mathcal{J}_0} \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 \neq \hat{z}_{ij}\}/2$ . Next we focus on  
10 obtaining a lower bound of  $\sum_{j \in \mathcal{J}_0} \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 \neq \hat{z}_{ij}\}$  to control the classification error rate  
11  $I^{-1} \sum_{i=1}^I \mathcal{J}\{\alpha_i^0 \neq \hat{\alpha}_i\}$ .

12 Motivated by Assumption 2, for each latent attribute  $k$ , denote  $j_k^1$  as the smallest  
13 integer  $j$  such that item  $j$  has a  $\mathbf{q}$ -vector  $\mathbf{e}_k$ , and denote  $j_k^2$  as the second smallest integer  $j$   
14 such that  $\mathbf{q}_j = \mathbf{e}_k$ , etc. For positive integer  $m$ , denote

$$\mathcal{B}^m = \{j_1^m, \dots, j_K^m\}. \tag{A.19}$$

1 For each  $k \in \{1, \dots, K\}$ , denote

$$J_{\min} = \min_{1 \leq k \leq K} |\{j \in \mathcal{J}_0; \mathbf{q}_j^0 = \mathbf{e}_k\}|, \quad \hat{J}_{\min} = \min_{1 \leq k \leq K} |\{j \in \mathcal{J}; \mathbf{q}_j^0 = \mathbf{e}_k\}|. \quad (\text{A.20})$$

2 Then, we have that  $\mathcal{B}^m \cap \mathcal{B}^l = \emptyset$  for any  $m \neq l$ , thus

$$\begin{aligned} & \sum_{i=1}^I \sum_{j \in \mathcal{J}_0} \mathcal{J}\{\xi(\mathbf{q}_j^0, \boldsymbol{\alpha}_i^0) \neq \xi(\mathbf{q}_j^0, \hat{\boldsymbol{\alpha}}_i)\} \\ & \geq \sum_{i=1}^I \sum_{m=1}^{J_{\min}} \sum_{j \in \mathcal{B}^m} \mathcal{J}\{\xi(\mathbf{q}_j^0, \boldsymbol{\alpha}_i^0) \neq \xi(\mathbf{q}_j^0, \hat{\boldsymbol{\alpha}}_i)\} \\ & = J_{\min} \sum_{i=1}^I \sum_{k=1}^K \mathcal{J}\{\xi(\mathbf{e}_k, \boldsymbol{\alpha}_i^0) \neq \xi(\mathbf{e}_k, \hat{\boldsymbol{\alpha}}_i)\}. \end{aligned} \quad (\text{A.21})$$

3 The last inequality holds since  $\sum_{k=1}^K \mathcal{J}\{\xi(\mathbf{e}_k, \boldsymbol{\alpha}_i^0) \neq \xi(\mathbf{e}_k, \hat{\boldsymbol{\alpha}}_i)\} \geq \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\}$ . Note (A.21)

4 implies  $o_p(S/I) \geq J_{\min} I^{-1} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\}$ . For simplicity,

$$\gamma_J = \frac{S}{IJ} = J^{-\frac{1}{\eta}} (\log J)^{\frac{\xi}{\eta}}.$$

5 Note that (32) in Assumption 2 implies that  $|\mathcal{J}|/J \geq \hat{J}_{\min}/J \geq \delta_J$  and  $J_{\min} \geq \hat{J}_{\min} - |\mathcal{J}_1 \cup \mathcal{J}_2|$ ;

6 by plugging these results into  $o_p(S/I) \geq J_{\min} I^{-1} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\}$ , we can obtain

$$o_p\left(\frac{S}{I}\right) + |\mathcal{J}_1 \cup \mathcal{J}_2| \geq \frac{\hat{J}_{\min}}{I} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\} \geq \frac{J\delta_J}{I} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\}.$$

7 From Lemma 5, we have  $|\mathcal{J}_i| = o_p(S/I)$  for  $i = 1, 2$ , which implies that  $|\mathcal{J}_1 \cup \mathcal{J}_2| = o_p(S/I)$ .

8 By substituting this into the above inequality, we can conclude that

$$o_p\left(\frac{S}{I}\right) \geq \frac{J\delta_J}{I} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\},$$

9 which is equivalent to  $I^{-1} \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\} = o_p(\gamma_J/\delta_J)$ . The proof of this theorem is complete.

1           The inequality (32) in Assumption 2 bridges between the misclassification error for the  
 2 local latent classes and the misclassification error for the latent attribute mastery patterns  $\hat{\mathcal{A}}$  by  
 3 using the inequality  $\sum_{k=1}^K \mathcal{J}\{\xi(\mathbf{e}_k, \boldsymbol{\alpha}_i^0) \neq \xi(\mathbf{e}_k, \hat{\boldsymbol{\alpha}}_i)\} \geq \mathcal{J}\{\boldsymbol{\alpha}_i^0 \neq \hat{\boldsymbol{\alpha}}_i\}$ .

#### 4 **Proof of Theorem 2**

5 For notational simplicity, denote  $i_{ja}^0 = \sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}$ . Thus, Assumption 2 implies that  $\forall \boldsymbol{\alpha} \in$   
 6  $\{0,1\}^K, \sum_{i=1}^I \mathcal{J}\{\boldsymbol{\alpha}_i^0 = \boldsymbol{\alpha}\} \geq \epsilon I$  and

$$i_{ja}^0 \geq \frac{2^K}{2^{K_j}} I \epsilon \geq I \epsilon. \quad (\text{A.22})$$

7 Recall that

$$\hat{\theta}_{ja} = \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\}}.$$

8 Rewrite  $\theta_{ja}^0$  as similar form

$$\theta_{ja}^0 = \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} \theta_{ja}^0}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} = \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} P_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}}.$$

9 By triangle inequality, we have

$$\begin{aligned} & \max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}^0| \\ &= \max_{j,a} \left| \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\}} - \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} P_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} \right| \\ &\leq \max_{j,a} \left| \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\}} - \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} \right| \\ &\quad + \max_{j,a} \left| \frac{\sum_{i=1}^I \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} - \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} \right| \\ &\quad + \max_{j,a} \left| \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} X_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} - \frac{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\} P_{ij}}{\sum_{i=1}^I \mathcal{J}\{z_{ij}^0 = a\}} \right| \\ &\equiv \mathcal{J}_1 + \mathcal{J}_2 + \mathcal{J}_3. \end{aligned}$$

10 Thereafter, we analyze these three terms separately. For the first term,

$$\begin{aligned}
\mathcal{J}_1 &\leq \max_{j,a} \left( \sum_i \mathcal{J}\{\hat{z}_{ij} = a\} X_{ij} \right) \cdot \frac{\sum_i |\mathcal{J}\{\hat{z}_{ij} = a\} - \mathcal{J}\{z_{ij}^0 = a\}|}{i_{ja}^0 \sum_i \mathcal{J}\{\hat{z}_{ij} = a\}} \\
&\leq \max_{j,a} \frac{\sum_i |\mathcal{J}\{\hat{z}_{ij} = a\} - \mathcal{J}\{z_{ij}^0 = a\}|}{i_{ja}^0} \\
&\leq \frac{1}{\epsilon l} \sum_i \mathcal{J}\{\alpha_i^0 \neq \hat{\alpha}_i\} = o_p\left(\frac{\gamma_J}{\delta_J}\right).
\end{aligned}$$

1 The last inequality holds since  $\forall j \in [J], j \in [L_j], \sum_i |\mathcal{J}\{\hat{z}_{ij} = a\} - \mathcal{J}\{z_{ij}^0 = a\}| \leq \sum_i \mathcal{J}\{\alpha_i^0 \neq \hat{\alpha}_i\}$ .

2 For the second term, we have

$$\mathcal{J}_2 = \max_{j,a} \frac{\sum_i |X_{ij}(\mathcal{J}\{\hat{z}_{ij} = a\} - \mathcal{J}\{z_{ij}^0 = a\})|}{i_{ja}^0} \leq \max_{j,a} \frac{\sum_i |\mathcal{J}\{\hat{z}_{ij} = a\} - \mathcal{J}\{z_{ij}^0 = a\}|}{i_{ja}^0}.$$

3 For the same reason as  $\mathcal{J}_1 \xrightarrow{P} 0$ , we can also conclude that  $\mathcal{J}_2 = o_p(\gamma_J/\delta_J)$ ; thus,  $\mathcal{J}_1 + \mathcal{J}_2 =$

4  $o_p(\gamma_J/\delta_J)$ . For the third term, we apply Hoeffding's inequality for bounded random variables

5 and obtain

$$P\left(\frac{\sum_i \mathcal{J}\{z_{ij}^0 = a\}(X_{ij} - P_{ij})}{i_{ja}^0} \geq t\right) \leq 2\exp(-2i_{ja}^0 t^2) \leq 2\exp(-2\epsilon l t^2).$$

6 Note the number of  $(j, a)$  pairs less than or equal to  $J2^K$  under Assumption 2, we have for

7  $\forall t > 0$ ,

$$P(\mathcal{J}_3 \geq t) \leq J2^{K+1} \exp(-2\epsilon l t^2). \tag{A.23}$$

8 Notably,  $2^{K+1}$  remains a constant as  $K$  is fixed. By choosing  $t = 1/\sqrt{I^{1-\tilde{c}}}$  for a small  $\tilde{c} > 0$ ,

9 the tail probability in A.23 converges to zero when the scaling condition  $\sqrt{J} = O(I^{1-\tilde{c}})$  holds.

10 This implies that  $\mathcal{J}_3 = o_p(1/\sqrt{I^{1-\tilde{c}}})$ . Bringing together the preceding results, we have

$$\max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}^0| = o_p\left(\frac{\gamma_J}{\delta_J}\right) + o_p\left(\frac{1}{\sqrt{I^{1-\tilde{c}}}}\right).$$

11 Note for any  $(j, a)$ , we have  $(\hat{\theta}_{ja} - \hat{\theta}_{ja})^2 \leq \tilde{c} i_{ja}^{-1}$  and  $i_{ja} \geq I\epsilon_2$  with the probability

12 approaching 1; thus  $\max_{j,a} |\hat{\theta}_{ja} - \hat{\theta}_{ja}| = O_p(I^{-1/2})$ . Therefore



$$\max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}^0| \leq \max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}^0| + \max_{j,a} |\hat{\theta}_{ja} - \theta_{ja}^0| = o_p\left(\frac{\gamma_J}{\delta_J}\right) + o_p\left(\frac{1}{\sqrt{I^{1-\tilde{c}}}}\right).$$

1

■