

Supplementary Material for “Learning Latent and Hierarchical Structures in Cognitive Diagnosis Models”

Chenchen Ma, Jing Ouyang, and Gongjun Xu

Department of Statistics, University of Michigan

In the supplementary material, we provide the proof of the main theorem, the derivations for the penalized EM algorithm and a sensitivity analysis of our algorithm with varying upper bounds for the number of latent classes.

1 Proof for Theorem 3

In this section, we provide the proof of Theorem 3.

Proof. We first introduce some notations. For two sequences $\{a_N\}$ and $\{b_N\}$, we denote $a_N \lesssim b_N$ if $a_N = O(b_N)$, and $a_N \asymp b_N$ if $a_N \lesssim b_N$ and $b_N \lesssim a_N$. We use $(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)$ to denote the true model parameter and use $(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)$ to denote the oracle MLE obtained by assuming the number of latent attributes, the hierarchical structure, the Q -matrix and the item-level diagnostic models are known. Let $(\hat{\boldsymbol{\pi}}^*, \hat{\boldsymbol{\Theta}}^*)$ be the MLE obtained by directly optimizing log-likelihood (9) and $(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}})$ be the estimator obtained by optimizing the regularized log-likelihood (10). We define $\hat{\boldsymbol{\pi}}_{\rho_N} := \{\hat{\pi}_m : \hat{\pi}_m > \rho_N, m \in [M]\}$ and $\hat{\boldsymbol{\Theta}}_{\rho_N} := \{\hat{\theta}_{j,m} : \hat{\pi}_m > \rho_N, j \in [J], m \in [M]\}$, the model parameters corresponding to the selected latent classes. Let M be the upper bound for the number of latent classes, M_0 be the true number of latent classes, and $\hat{M} = |\{m : \hat{\pi}_m > \rho_N, m \in [M]\}|$ be the estimated number of latent classes. Without loss of generality, let $\hat{\boldsymbol{\pi}}_{\text{full}}^0 = (\hat{\boldsymbol{\pi}}^0, \mathbf{0}_{M-M_0})$. For the true item parameter matrix $\boldsymbol{\Theta}^0$, we defined the set of identical item parameter pairs $S^0 = \{(j, k_1, k_2) : \theta_{j,k_1}^0 = \theta_{j,k_2}^0, 1 \leq k_1 < k_2 \leq M_0\}$. Similarly, for $(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}})$ we define $\hat{S} = \{(j, k_1, k_2) : \hat{\theta}_{j,k_1} = \hat{\theta}_{j,k_2}, 1 \leq k_1 < k_2 \leq M, \hat{\pi}_{k_1} > \rho_N, \hat{\pi}_{k_2} > \rho_N\}$. We say $\hat{S} \sim S^0$ if there exists a column permutation σ of $\hat{\boldsymbol{\Theta}}$ such that $\hat{S}_\sigma = S^0$.

The probability $\mathbb{P}(\hat{M} \neq M_0)$ can be decomposed into two parts:

$$\mathbb{P}(\hat{M} \neq M_0) = \mathbb{P}(\hat{M} < M_0) + (\hat{M} > M_0). \quad (1)$$

Similarly, the probability $\mathbb{P}(\hat{S} \neq S^0)$ can be decomposed into three parts:

$$\mathbb{P}(\hat{S} \neq S^0) = \mathbb{P}(\hat{M} < M_0) + (\hat{M} > M_0) + \mathbb{P}(\hat{S} \neq S^0, \hat{M} = M_0). \quad (2)$$

In the following, we will bound each part in (1) and (2) respectively. Therefore, we will consider three cases below:

1. overfitted case: $\hat{M} > M_0$,
2. underfitted case: $\hat{M} < M_0$,
3. $\hat{M} = M_0$ but $\hat{S} \neq S^0$.

The objective function is

$$G_N(\boldsymbol{\pi}, \boldsymbol{\Theta}) = \frac{l_N(\boldsymbol{\pi}, \boldsymbol{\Theta}; \mathcal{R})}{N} - \frac{\lambda_N^{(1)}}{N} \sum_{k=1}^M \log_{[\rho_N]} \pi_k - \frac{\lambda_N^{(2)}}{N} \sum_{j=1}^J \mathcal{J}_{\tau, \rho_N}(\boldsymbol{\theta}_j), \quad (3)$$

where $\log_{[\rho_N]} \pi_k = \log \pi_k \cdot \mathbb{I}(\pi_k > \rho_N) + \log \rho_N \cdot \mathbb{I}(\pi_k \leq \rho_N)$. Let $\log_{[\rho_N]}(\boldsymbol{\pi}) = \sum_{k=1}^M \log_{[\rho_N]} \pi_k$.

First consider the overfitted case where $\hat{M} > M_0$. The event $\{G_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) > G_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)\}$ implies that

$$\begin{aligned} & \frac{1}{N} \sum_{i=1}^N \log \left[\frac{\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k})^{1-R_{ij}}}{\sum_{k=1}^M \hat{\pi}_k^0 \prod_{j=1}^J (\hat{\theta}_{j,k}^0)^{R_{ij}} (1 - \hat{\theta}_{j,k}^0)^{1-R_{ij}}} \right] \\ & > \frac{\lambda_N^{(1)}}{N} \{ \log_{[\rho_N]}(\hat{\boldsymbol{\pi}}) - \log_{[\rho_N]}(\hat{\boldsymbol{\pi}}_{full}^0) \} + \frac{\lambda_N^{(2)}}{N} \left\{ \sum_{j=1}^J \mathcal{J}_{\tau, \rho_N}(\hat{\boldsymbol{\theta}}_j) - \sum_{j=1}^J \mathcal{J}_{\tau, \rho_N}(\hat{\boldsymbol{\theta}}_j^0) \right\} \end{aligned} \quad (4)$$

$$:= J_1 + J_2.$$

For the RHS of (4), we have $J_1 \gtrsim N^{-1} \lambda_N^{(1)} |\log \rho_N|$ and $J_2 \gtrsim -N^{-1} \lambda_N^{(2)} \tau J M^2$. Since $\lambda_N^{(2)} \tau = o(\lambda_N^{(1)} |\log \rho_N|)$, we have $\text{RHS} \gtrsim N^{-1} \lambda_N^{(1)} |\log \rho_N|$.

For the LHS of (4), we have

$$\begin{aligned}
\text{LHS of (4)} &= \frac{1}{N} \log \left[\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k})^{1-R_{ij}} \right] - \frac{1}{N} \log \left[\sum_{k=1}^M \hat{\pi}_k^0 \prod_{j=1}^J (\hat{\theta}_{j,k}^0)^{R_{ij}} (1 - \hat{\theta}_{j,k}^0)^{1-R_{ij}} \right] \\
&\leq \frac{1}{N} \log \left[\sum_{k=1}^M \hat{\pi}_k^* \prod_{j=1}^J (\hat{\theta}_{j,k}^*)^{R_{ij}} (1 - (\hat{\theta}_{j,k}^*))^{1-R_{ij}} \right] - \frac{1}{N} \log \left[\sum_{k=1}^M \hat{\pi}_k^0 \prod_{j=1}^J (\hat{\theta}_{j,k}^0)^{R_{ij}} (1 - \hat{\theta}_{j,k}^0)^{1-R_{ij}} \right] \\
&\lesssim N^{-\delta},
\end{aligned}$$

where the last inequality follows from Assumption 1. When $N^{1-\delta}/|\log(\rho_N)| = o(\lambda_N^{(1)})$, we have $N^{-\delta} = o(N^{-1}\lambda_N^{(1)}|\log \rho_N|)$, which implies that the event described in (4) will happen with probability tending to zero. Therefore we have $\mathbb{P}(\hat{M} > M_0) \rightarrow 0$ as $N \rightarrow \infty$. That is to say, with the appropriate choice of tuning parameters, the extent that the log-penalty part favors a smaller model would dominate the extent that the likelihood part favors a larger model in the overfitted case.

Now consider the under-fitted case where $\hat{M} < M_0$. We need to bound

$$\mathbb{P}\left(\sup_{\hat{M} < M_0} [G_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - G_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)] > 0 \right). \quad (5)$$

We follow a similar argument to Shen et al. (2012). More specifically, since

$$\mathbb{P}\left(\sup_{\hat{M} < M_0} [G_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - G_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)] > 0 \right) \leq \sum_{m=1}^{M_0-1} \mathbb{P}\left(\sup_{\hat{M}=m} [G_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - G_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)] > 0 \right), \quad (6)$$

we will bound each term in the RHS of (6). By the large deviation inequality in Theorem 1 of Wong and Shen (1995), we have

$$\begin{aligned}
&\mathbb{P}\left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq \epsilon_N^2} \left[\frac{1}{N} l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - \frac{1}{N} l_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right] > -\epsilon_N^2 \right) \\
&\leq \mathbb{P}\left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq \epsilon_N^2} \left[\frac{1}{N} l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - \frac{1}{N} l_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right] > -\epsilon_N^2 \right) \leq \exp(-N\epsilon_N^2),
\end{aligned} \quad (7)$$

where $h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) = \sum_{\mathbf{R} \in \{0,1\}^J} [\mathbb{P}(\mathbf{R} | \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}})^{1/2} - \mathbb{P}(\mathbf{R} | \boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)^{1/2}]^2$ is the Hellinger distance. From the remark in Wong and Shen (1995), the inequality (7) holds for any $t > \epsilon_N$.

To use this large deviation inequality, we need to introduce the notion of bracketing Hellinger metric entropy $H(t, \mathcal{B}_m)$, which characterizes the size of the local parameter space. Consider the local parameter space $\mathcal{B}_m = \{(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) : \hat{M} = m \leq M_0, h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \leq 2\epsilon_N^2\}$, then $H(t, \mathcal{B}_m)$ is defined as the logarithm of the cardinality of the t-bracketing of \mathcal{B}_m of the smallest size. Specifically, following the definition in Shen et al. (2012), consider a bracket covering $S(t, m) = \{f_1^l, f_1^u, \dots, f_m^l, f_m^u\}$ such that $\max_{1 \leq j \leq m} \|f_j^u - f_j^l\|_2 \leq t$ and for any $f \in \mathcal{B}_m$, there is some j such that $f_j^l \leq f \leq f_j^u$ almost surely. Then $H(t, \mathcal{B}_m)$ is defined as $\log(\min\{m : S(t, m)\})$. Following Lemma 3 in Gu and Xu (2019), for any $2^{-4}\epsilon < t < \epsilon$, there is

$$H(t, \mathcal{B}_m) \lesssim M_0 \log M \log(2\epsilon/t). \quad (8)$$

Next we need to verify the conditions in Wong and Shen (1995). Let's take $\epsilon_N = \sqrt{M_0 \log M/N}$ and verify the entropy integral condition in Theorem 1 of Wong and Shen (1995) for such ϵ_N . The integral of bracketing Hellinger metric entropy on the interval $[2^{-8}\epsilon_N^2, \sqrt{2}\epsilon_N]$ satisfies the following inequality

$$\begin{aligned} \int_{2^{-8}\epsilon_N^2}^{\sqrt{2}\epsilon_N} H^{1/2}(t, \mathcal{B}_m) dt &\leq \int_{2^{-8}\epsilon_N^2}^{\sqrt{2}\epsilon_N} \sqrt{M_0 \log M \log(2\epsilon_N/t)} dt \\ &= \sqrt{M_0 \log M} \int_{\sqrt{\log \sqrt{2}}}^{\sqrt{\log \frac{2^9}{\epsilon_N}}} 4\epsilon_N u^2 e^{-u^2} du \\ &= \sqrt{M_0 \log M} \cdot 2\epsilon_N \int_{\log \sqrt{2}}^{\log \frac{2^9}{\epsilon_N}} \sqrt{u} e^{-u} du \\ &\lesssim \sqrt{N} \epsilon_N^2. \end{aligned}$$

Note that $\epsilon_N = o(1)$ as $N \rightarrow \infty$.

Following the proof in Gu and Xu (2019), there exists a constant c_0 , for some small constant $t > \epsilon_N$, we have

$$C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) := \inf_{(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}): \hat{M} \leq M_0} \left\{ \frac{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0))}{\max(M_0 - \hat{M}, 1)} \right\} \geq c_0 \gtrsim t^2 > \epsilon_N^2.$$

Moreover, for $\hat{M} = m < M_0$, there is $h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq (M_0 - m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)$. In order

to have the probability of the event (4) go to zero in the under-fitted case, the log-penalty term should not be too large such that the likelihood part is dominated by the log-penalty term that favors a smaller model. Here we take $\lambda_N^{(1)} = o(N \log \rho_N|^{-1})$. Then for (6) we have

$$\begin{aligned}
& \text{RHS of (6)} \\
& \leq \sum_{m=1}^{M_0-1} \mathbb{P} \left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq (M_0-m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0), \hat{M}=m} [G_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - G_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)] > 0 \right) \\
& \leq \sum_{m=1}^{M_0-1} \mathbb{P} \left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq (M_0-m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0), \hat{M}=m} [l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - l_N(\hat{\boldsymbol{\pi}}^0, \hat{\boldsymbol{\Theta}}^0)] > -\frac{\lambda_N^{(1)} M_0 |\log \rho_N|}{N} \right) \\
& \leq \sum_{m=1}^{M_0-1} \mathbb{P} \left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq (M_0-m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0), \hat{M}=m} [l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - l_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)] > -\frac{\lambda_N^{(1)} M_0 |\log \rho_N|}{N} \right) \\
& \leq \sum_{m=1}^{M_0-1} \mathbb{P} \left(\sup_{h^2((\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}), (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)) \geq (M_0-m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0), \hat{M}=m} [l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) - l_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)] > -(M_0 - m)C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right) \\
& \leq \sum_{m=1}^{M_0-1} \exp \left(-c_2 N (M_0 - m) C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right) \\
& \leq c_3 \exp \left(-c_2 N C_{\min}(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right).
\end{aligned}$$

Therefore we have $\mathbb{P}(\hat{M} < M_0) \rightarrow 0$ as $N \rightarrow \infty$. So far we have proved (12) in Theorem 3,

$$\mathbb{P}(\hat{M} \neq M_0) = \mathbb{P}(\hat{M} < M_0) + \mathbb{P}(\hat{M} > M_0) \rightarrow 0.$$

Finally we consider the third case where $\hat{M} = M_0$ but $\hat{S} \approx S^0$. The argument is similar to the proof of Proposition 2 in [Xu and Shang \(2018\)](#). We first show $(\hat{\boldsymbol{\pi}}_{\rho_N}, \hat{\boldsymbol{\Theta}}_{\rho_N})$ converge to $(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)$ with rate $N^{-1/2}$. For $(\boldsymbol{\pi}, \boldsymbol{\Theta})$ with $(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N})$ in a small neighborhood of $(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)$,

$$\begin{aligned}
G'_N(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}) & := \frac{l_N(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}; \mathcal{R})}{N} - \frac{\lambda_N^{(1)}}{N} \sum_{k:\pi_k > \rho_N} \log \pi_k - \frac{\lambda_N^{(2)}}{N} \sum_{j=1}^J \mathcal{J}_{\tau, \rho_N}(\boldsymbol{\theta}_j) \\
& = \frac{l_N(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}; \mathcal{R})}{N} - O(\lambda_N^{(1)} N^{-1} |\log \rho_N|) - O(\lambda_N^{(2)} \tau N^{-1}),
\end{aligned}$$

converges uniformly to the same limit of $l_N(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}; \mathcal{R})/N$ by the uniform law of large number, since $\lambda_N^{(1)} N^{-1} |\log \rho_N| \rightarrow 0$ and $\lambda_N^{(2)} \tau N^{-1} \rightarrow 0$. We use $G_0(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N})$ to denote the limit process, which is the expectation of the negative log-likelihood of a single observation. By Taylor's

expansion, we have $G_0(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}) - G_0(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) = O(\|(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}) - (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)\|^2)$.

For the log-likelihood function $l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}; \mathcal{R}) = \sum_{i=1}^N \log(\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}}))$, we have

$$\begin{aligned}
& \frac{1}{N} |l_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}; \mathcal{R}) - l_N(\hat{\boldsymbol{\pi}}_{\rho_N}, \hat{\boldsymbol{\Theta}}_{\rho_N}; \mathcal{R})| \\
& \leq \frac{1}{N} \sum_{i=1}^N \left| \log \left(\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}}) \right) - \log \left(\sum_{k:\hat{\pi}_k > \rho_N} \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}}) \right) \right| \\
& \leq \frac{1}{N} \sum_{i=1}^N \frac{|(\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}})) - (\sum_{k:\hat{\pi}_k > \rho_N} \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}}))|}{\sqrt{(\sum_{k=1}^M \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}})) \times (\sum_{k:\hat{\pi}_k > \rho_N} \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}}))}} \quad (9) \\
& \leq \frac{1}{N} \sum_{i=1}^N \frac{(M - \hat{M})\rho_N}{\sum_{k:\hat{\pi}_k > \rho_N} \hat{\pi}_k \prod_{j=1}^J \hat{\theta}_{j,k}^{R_{ij}} (1 - \hat{\theta}_{j,k}^{1-R_{ij}})} \\
& = O(\rho_N) = O(N^{-d}), \quad d \geq 1, \quad (10)
\end{aligned}$$

where inequality (9) follows from an upper bound for log function. Specifically, for $x \geq 1$, we know $\log x \leq (x - 1)/\sqrt{x}$, and thus for $0 < x \leq y$, we have $\log y - \log x \leq (y - x)/\sqrt{xy}$. From (10), $G'_N(\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}) = G'_N(\hat{\boldsymbol{\pi}}_{\rho_N}, \hat{\boldsymbol{\Theta}}_{\rho_N}) + O(N^{-d}) \geq G'_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)$ and thus $G'_N(\hat{\boldsymbol{\pi}}_{\rho_N}, \hat{\boldsymbol{\Theta}}_{\rho_N}) > G'_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) - O(N^{-d}) \geq G'_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) - O(N^{-1})$. Since $N^{-1/2}\lambda_N^{(1)} \rightarrow 0$ and $N^{-1/2}\lambda_N^{(2)}\tau \rightarrow 0$, then for sufficiently small ζ , by Taylor's expansion,

$$\mathbb{E} \left(\sup_{\|(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}) - (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0)\| \leq \zeta} G'_N(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}; \mathcal{R}) - G_0(\boldsymbol{\pi}_{\rho_N}, \boldsymbol{\Theta}_{\rho_N}) - G'_N(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0; \mathcal{R}) + G_0(\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) \right) = O(\zeta N^{-1/2}).$$

By Theorem 3.2.5 in [Van Der Vaart and Wellner \(1996\)](#), we have $(\hat{\boldsymbol{\pi}}_{\rho_N}, \hat{\boldsymbol{\Theta}}_{\rho_N}) - (\boldsymbol{\pi}^0, \boldsymbol{\Theta}^0) = O_p(N^{-1/2})$.

We next show selection consistency of S^0 . If true item parameters $\theta_{j,k_1}^0 \neq \theta_{j,k_2}^0$, then from the above convergence result, we know $\hat{\theta}_{j,k_1} \rightarrow \theta_{j,k_1}^0$ and $\hat{\theta}_{j,k_2} \rightarrow \theta_{j,k_2}^0$, and thus $\hat{\theta}_{j,k_1} \neq \hat{\theta}_{j,k_2}$ in probability. If true item parameters $\theta_{j,k_1}^0 = \theta_{j,k_2}^0$ but $\hat{\theta}_{j,k_1} \neq \hat{\theta}_{j,k_2}$, by the Karush-Kuhn-Tucker (KKT) conditions, we have $N^{-1/2} \partial l_N(\boldsymbol{\pi}, \boldsymbol{\Theta}; \mathcal{R}) / \partial \theta_{j,k_1} |_{(\boldsymbol{\pi}, \boldsymbol{\Theta}) = (\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}})} = N^{-1/2} \lambda_N^{(2)} \rightarrow \infty$ in probability. However $N^{-1/2} \partial l_N(\boldsymbol{\pi}, \boldsymbol{\Theta}; \mathcal{R}) / \partial \theta_{j,k_1} |_{(\boldsymbol{\pi}, \boldsymbol{\Theta}) = (\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}})} = O_p(1)$. Therefore, if $\theta_{j,k_1}^0 = \theta_{j,k_2}^0$, we have $\hat{\theta}_{j,k_1} = \hat{\theta}_{j,k_2}$ in probability, which proved the selection consistency that $\mathbb{P}(\hat{S} \approx S^0) \rightarrow 0$ as $N \rightarrow \infty$. ■

2 Derivations of PEM Algorithm

In this section, we give detailed derivations of the penalized EM algorithm in Section 4.1. First let's introduce a new variable $\mathbf{d} = (d_{jkl}, j = 1, \dots, J, 1 \leq k < l \leq M)$ to be the differences of the item parameters for each item. Then our problem becomes

$$\begin{aligned} \min_{\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}} \quad & G(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) \\ \text{s.t.} \quad & d_{jkl} = \theta_{jk} - \theta_{jl} \\ & j = 1, \dots, J, 1 \leq k < l \leq M. \end{aligned} \tag{11}$$

By using the difference convex property of the truncated Lasso penalty, we can decompose the objective function into two parts:

$$G(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) = G_1(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) - G_2(\mathbf{d}), \tag{12}$$

where

$$G_1(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) = -\frac{1}{N}Q(\boldsymbol{\pi}, \boldsymbol{\Theta} | \boldsymbol{\pi}^{(c)}, \boldsymbol{\Theta}^{(c)}) + \tilde{\lambda}_1 \sum_{k=1}^M \log \pi_k + \tilde{\lambda}_2 \sum_{j=1}^J \sum_{1 \leq k < l \leq M} |d_{jkl}|, \tag{13}$$

$$G_2(\mathbf{d}) = \tilde{\lambda}_2 \sum_{j=1}^J \sum_{1 \leq k < l \leq M} (|d_{jkl} - \tau|)_+. \tag{14}$$

Then we construct a sequence of upper approximation of $G(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d})$ iteratively by replacing $G_2(\mathbf{d})$ at iteration $c + 1$ with its piecewise affine minorization:

$$G_2^{(c)}(\mathbf{d}) = G_2(\hat{\mathbf{d}}^{(c)}) + \tilde{\lambda}_2 \sum_{j=1}^J \sum_{1 \leq k < l \leq M} (|d_{jkl}| - |\hat{d}_{jkl}^{(c)}|) \cdot \mathbb{I}(|\hat{d}_{jkl}^{(c)}| \geq \tau), \tag{15}$$

at the current estimate $\hat{\mathbf{d}}^{(c)}$, which lead to an upper convex approximation:

$$\begin{aligned} G^{(c+1)}(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) = & -\frac{1}{N}Q(\boldsymbol{\pi}, \boldsymbol{\Theta} | \boldsymbol{\pi}^{(c)}, \boldsymbol{\Theta}^{(c)}) + \tilde{\lambda}_1 \sum_{k=1}^M \log \pi_k \\ & + \tilde{\lambda}_2 \sum_{j=1}^J \sum_{1 \leq k < l \leq M} |d_{jkl}| \cdot \mathbb{I}(|\hat{d}_{jkl}^{(c)}| < \tau) \end{aligned}$$

$$+ \tilde{\lambda}_2 \tau \sum_{j=1}^J \sum_{1 \leq k < l \leq M} \mathbb{I}(|\hat{d}_{jkl}^{(c)}| \geq \tau).$$

Now we can apply ADMM. At iteration $c + 1$, the augmented Lagrangian is

$$L_\gamma(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}, \mathbf{y}) = G^{(c+1)}(\boldsymbol{\pi}, \boldsymbol{\Theta}, \mathbf{d}) + \sum_{j=1}^J \sum_{1 \leq k < l \leq M} y_{jkl} \cdot (d_{jkl} - (\theta_{jk} - \theta_{jl})) + \frac{\gamma}{2} \sum_{j=1}^J \sum_{1 \leq k < l \leq M} |d_{jkl} - (\theta_{jk} - \theta_{jl})|^2, \quad (16)$$

where y_{jkl} 's are the dual variables and γ is a nonnegative penalty parameter. Then ADMM (Boyd et al., 2011) consists of the following iterations:

$$\begin{aligned} \boldsymbol{\pi}^{(c+1)} &= \underset{\boldsymbol{\pi}}{\operatorname{argmin}} L_\gamma(\boldsymbol{\pi}, \boldsymbol{\Theta}^{(c)}, \mathbf{d}^{(c)}, \mathbf{y}^{(c)}), \\ \boldsymbol{\Theta}^{(c+1)} &= \underset{\boldsymbol{\Theta}}{\operatorname{argmin}} L_\gamma(\boldsymbol{\pi}^{(c+1)}, \boldsymbol{\Theta}, \mathbf{d}^{(c)}, \mathbf{y}^{(c)}), \\ \mathbf{d}^{(c+1)} &= \underset{\mathbf{d}}{\operatorname{argmin}} L_\gamma(\boldsymbol{\pi}^{(c+1)}, \boldsymbol{\Theta}^{(c+1)}, \mathbf{d}, \mathbf{y}^{(c)}), \\ y_{jkl}^{(c+1)} &= y_{jkl}^{(c)} + \gamma(d_{jkl}^{(c+1)} - (\theta_{jk}^{(c+1)} - \theta_{jl}^{(c+1)})), \quad j = 1, \dots, J, 1 \leq k < l \leq M. \end{aligned}$$

Using the scaled Lagrangian multiplier $\mu_{jkl} = y_{jkl}/\gamma$ and defining the residual $r_{jkl} = d_{jkl} - (\theta_{jk} - \theta_{jl})$, we have:

$$\begin{aligned} & y_{jkl} \cdot (d_{jkl} - (\theta_{jk} - \theta_{jl})) + \frac{\gamma}{2} |d_{jkl} - (\theta_{jk} - \theta_{jl})|^2 \\ &= y_{jkl} \cdot r_{jkl} + \frac{\gamma}{2} r_{jkl}^2 \\ &= \frac{\gamma}{2} (r_{jkl} + (1/\gamma)y_{jkl})^2 - \frac{1}{2\gamma} \mu_{jkl}^2 \\ &= \frac{\gamma}{2} (r_{jkl} + \mu_{jkl})^2 - \frac{1}{2\gamma} \mu_{jkl}^2. \end{aligned}$$

Then using the scaled dual variable, we can express ADMM as:

$$\begin{aligned} \boldsymbol{\pi}^{(c+1)} &= \underset{\boldsymbol{\pi}}{\operatorname{argmin}} G^{(c+1)}(\boldsymbol{\pi}, \boldsymbol{\Theta}^{(c)}, \mathbf{d}^{(c)}), \\ \boldsymbol{\Theta}^{(c+1)} &= \underset{\boldsymbol{\Theta}}{\operatorname{argmin}} G^{(c+1)}(\boldsymbol{\pi}^{(c+1)}, \boldsymbol{\Theta}, \mathbf{d}^{(c)}) + \frac{\gamma}{2} \sum_{j=1}^J \sum_{1 \leq k < l \leq M} (d_{jkl}^{(c)} - (\theta_{jk}^{(c)} - \theta_{jl}^{(c)}) + \mu_{jkl}^{(c)}), \\ \mathbf{d}^{(c+1)} &= \underset{\mathbf{d}}{\operatorname{argmin}} G^{(c+1)}(\boldsymbol{\pi}^{(c+1)}, \boldsymbol{\Theta}^{(c+1)}, \mathbf{d}) + \frac{\gamma}{2} \sum_{j=1}^J \sum_{1 \leq k < l \leq M} (d_{jkl} - (\theta_{jk}^{(c+1)} - \theta_{jl}^{(c+1)}) + \mu_{jkl}^{(c)}), \end{aligned}$$

$$\mu_{jkl}^{(c+1)} = \mu_{jkl}^{(c)} + d_{jkl}^{(c+1)} - (\theta_{jk}^{(c+1)} - \theta_{jl}^{(c+1)}), \quad j = 1, \dots, J, 1 \leq k < l \leq M.$$

Specifically, we get the following updates:

$$(1) \quad \pi_k^{(c+1)} = \frac{\sum_{i=1}^N s_{ik}^{(c+1)}/N - \tilde{\lambda}_1}{1 - M\tilde{\lambda}_1}, \quad \text{where } s_{ik}^{(c+1)} = \frac{\pi_k^{(c)} \varphi_k(\mathbf{R}_i; \boldsymbol{\Theta}_k^{(c)})}{\sum_{k'}^{(c)} \pi_{k'}^{(c)} \varphi_{k'}(\mathbf{R}_i; \boldsymbol{\Theta}_{k'}^{(c)})},$$

$$(2) \quad \hat{\theta}_{jk}^{(c+1)} = \underset{\theta_{jk}}{\operatorname{argmin}} \left\{ -\frac{\sum_{i=1}^N s_{ik}^{(c)} R_{ij}}{N} \log \theta_{jk} - \frac{\sum_{i=1}^N s_{ik}^{(c)} (1 - R_{ij})}{N} \log(1 - \theta_{jk}) \right. \\ \left. + \frac{\gamma}{2} \sum_{l>k} (d_{jkl}^{(c)} - (\theta_{jk} - \hat{\theta}_{jl}^{(c)}) + \hat{\mu}_{jkl}^{(c)})^2 \right. \\ \left. + \frac{\gamma}{2} \sum_{l<k} (d_{jlk}^{(c)} - (\hat{\theta}_{jl}^{(c+1)} - \theta_{jk}) + \hat{\mu}_{jlk}^{(c)})^2 \right\}$$

$$(3) \quad \hat{d}_{jkl}^{(c+1)} = \begin{cases} \hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)} - \hat{\mu}_{jkl}^{(c)}, & \text{if } |\hat{d}_{jkl}^{(c)}| \geq \tau \\ \operatorname{ST}(\hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)} - \hat{\mu}_{jkl}^{(c)}; \tilde{\lambda}_2/\gamma), & \text{if } |\hat{d}_{jkl}^{(c)}| < \tau, \text{ where } \operatorname{ST}(x; \gamma) = (|x| - \gamma)_+ x / |x| \end{cases},$$

$$(4) \quad \hat{\mu}_{jkl}^{(c+1)} = \hat{\mu}_{jkl}^{(c)} + \hat{d}_{jkl}^{(c+1)} - (\hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)}).$$

Note that the objective in step (2) is convex in θ_{jk} , therefore we use gradient descent to perform the minimization.

3 PEM Algorithm with Missing Values

In this section, we present the penalized EM algorithm with missing values. Here we use a mask matrix $M \in \{0, 1\}^{N \times J}$ to indicate the locations of the missing values, where $M_{i,j} = 0$ means the i th subject's response to the j th item is missing. The details of the algorithm is summarized in Algorithm 1.

Algorithm 1: Penalized EM with missing data

Data: Binary response matrix $\mathcal{R} = (R_{i,j})_{N \times J}$ and the mask matrix $\mathbf{M} = (M_{ij})_{N \times J}$ indicating missing values.

Set hyperparameters $\tilde{\lambda}_1, \tilde{\lambda}_2, \tau, \gamma$ and ρ .

Set an upper bound of the number of latent classes L .

Initialize parameters $\boldsymbol{\pi}, \boldsymbol{\Theta}$, and the conditional expectations \mathbf{s} .

while *not converged* **do**

 In the $(c+1)$ th iteration,

for $(i, k) \in [N] \times [L]$ **do**

$$\left[s_{ik}^{(c+1)} = \frac{\pi_k^{(c)} \varphi_k(\mathbf{R}_i; \boldsymbol{\theta}_k^{(c)})}{\sum_{k'}^{(c)} \pi_{k'}^{(c)} \varphi_{k'}(\mathbf{R}_i; \boldsymbol{\theta}_{k'}^{(c)})}, \quad \varphi(\mathbf{r}_i; \boldsymbol{\theta}_k) = \prod_{j=1}^J (\theta_{jk}^{R_{ij}} (1 - \theta_{kj})^{1-R_{ij}})^{m_{ij}} \right]$$

for $k \in [L]$ *and* $\pi_k^{(c)} > \rho$ **do**

$$\left[\pi_k^{(c+1)} = \frac{\sum_{i=1}^N s_{ik}^{(c+1)} / N - \tilde{\lambda}_1}{1 - L\tilde{\lambda}_1}. \right]$$

for $(j, k) \in [J] \times [L]$ *and* $\pi_k^{(c+1)} > \rho$ **do**

$$\left[\theta_{jk}^{(c+1)} = \underset{\theta_{jk}}{\operatorname{argmin}} \left\{ - \frac{\sum_{i=1}^N s_{ik}^{(c)} R_{ij} m_{ij}}{\sum_{i=1}^N m_{ij}} \log \theta_{jk} - \frac{\sum_{i=1}^N s_{ik}^{(c)} (1 - i_j) m_{ij}}{\sum_{i=1}^N m_{ij}} \log(1 - \theta_{jk}) \right. \right. \\ \left. \left. + \frac{\gamma}{2} \sum_{l>k} (\hat{d}_{jkl}^{(c)} - (\theta_{jk} - \hat{\theta}_{jl}^{(c)}) + \hat{\mu}_{jkl}^{(c)})^2 \right. \right. \\ \left. \left. + \frac{\gamma}{2} \sum_{l<k} (\hat{d}_{jlk}^{(c)} - (\hat{\theta}_{jl}^{(c+1)} - \theta_{jk}) + \hat{\mu}_{jlk}^{(c)})^2 \right\} \right]$$

for $j \in [J], k, l \in [L], k < l$ *and* $\pi_k^{(c+1)} > \rho, \pi_l^{(c+1)} > \rho$ **do**

$$\left[\hat{d}_{jkl}^{(c+1)} = \begin{cases} \hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)} - \hat{\mu}_{jkl}^{(c)}, & \text{if } |\hat{d}_{jkl}^{(c)}| \geq \tau \\ \operatorname{ST}(\hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)} - \hat{\mu}_{jkl}^{(c)}, \tilde{\lambda}_2/\gamma), & \text{if } |\hat{d}_{jkl}^{(c)}| < \tau \end{cases} \right. \\ \left. \hat{\mu}_{jkl}^{(c+1)} = \hat{\mu}_{jkl}^{(c)} + \hat{d}_{jkl}^{(c+1)} - (\hat{\theta}_{jk}^{(c+1)} - \hat{\theta}_{jl}^{(c+1)}). \right]$$

Output: $\{\hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\Theta}}, \hat{\mathbf{s}}\}$

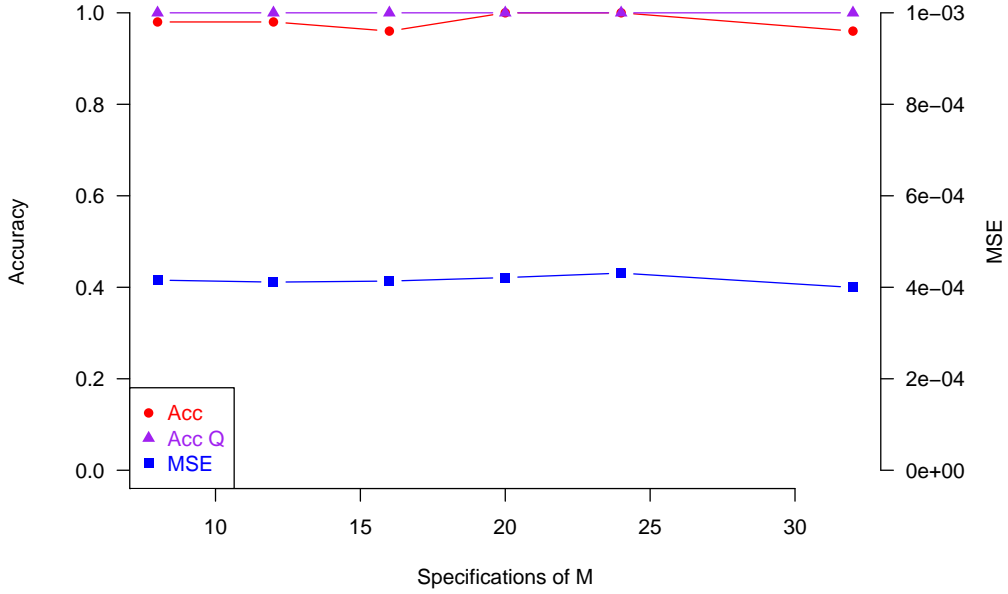
4 Sensitivity Analysis

In this section, we conduct the sensitivity analysis of our algorithm by investigating the effects of different inputs of M , the upper bound of the number of latent classes, on the simulation results.

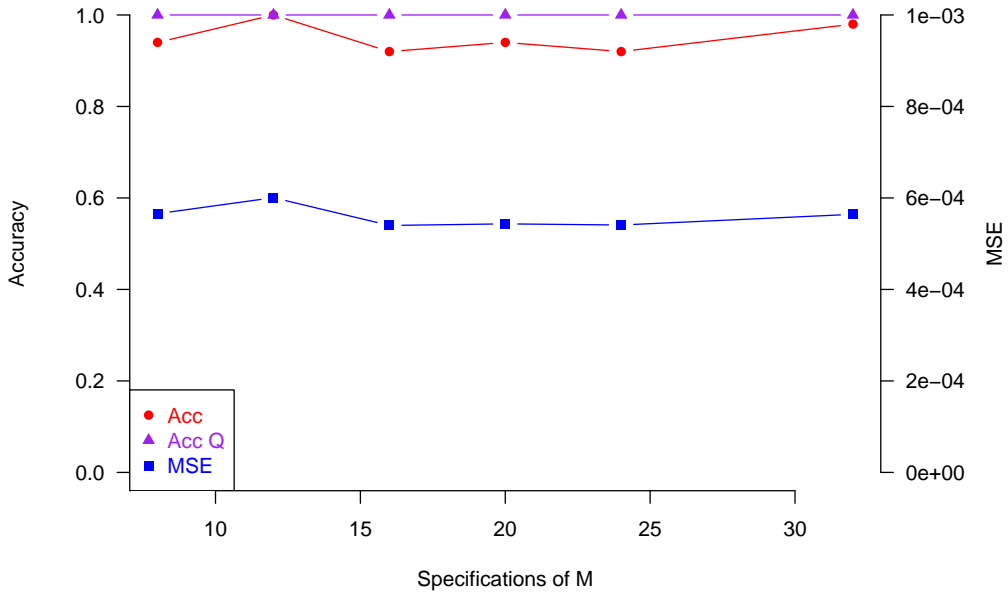
In particular, we focus on two simulation settings: (1) DINA model with linear hierarchical structure, $N = 500$ and $r = 0.1$; (2) GDINA model with linear hierarchical structure, $N = 1000$ and $r = 0.1$. Both two settings have $K = 4$ latent attributes and $J = 30$ test items, and run

for 50 repetitions. We keep the parameter generation process and the hyperparameter tuning strategy consistent with the simulation studies in the main article. In this sensitivity analysis, we fit our proposed method with various $M = \{8, 12, 16, 20, 24, 32\}$ in the two simulations settings. The evaluation results in DINA and GDINA settings are based on metrics $\text{Acc}(\hat{M})$, $\text{Acc}(\hat{\mathbf{P}})$, $\text{Acc}(\hat{\mathcal{E}})$, $\text{MSE}(\hat{\Theta})$ and $\text{Acc}(\hat{Q})$. Consistent with the simulation studies in the main article, the $\text{Acc}(\hat{M})$, $\text{Acc}(\hat{\mathbf{P}})$ and $\text{Acc}(\hat{\mathcal{E}})$ are calculated for all the cases; $\text{MSE}(\hat{\Theta})$ is calculated for the cases when the number of latent classes is correctly selected; $\text{Acc}(\hat{Q})$ is calculated for the cases when the hierarchical structure is successfully recovered. The results are plotted in Figure 1.

From the simulation results in Figure 1, we see our proposed method is robust to the different specifications of M , in terms of all metrics. Among cases with different M , our method achieves a high accuracy in estimating the number of latent classes, and in recovering the partial orders, the hierarchical structures, the item parameter matrix, and the Q -matrix. In terms of computation time, the average running time is 0.36 seconds and 1.12 seconds for DINA and GDINA, respectively, per repetition per set of tuning hyperparameters.



(a)



(b)

Figure 1: Sensitivity analysis results. (a) DINA results; (b) GDINA results. The red curve captures the $\text{Acc}(\hat{M})$, $\text{Acc}(\hat{P})$, $\text{Acc}(\hat{E})$, the blue curve captures $\text{MSE}(\hat{\Theta})$ and the purple curve captures the $\text{Acc}(\hat{Q})$ for various M .

References

- Boyd, S., Parikh, N., Chu, E., Peleato, B., and Eckstein, J. (2011). Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends[®] in Machine learning*, 3(1):1–122.
- Gu, Y. and Xu, G. (2019). Learning attribute patterns in high-dimensional structured latent attribute models. *Journal of Machine Learning Research*, 20(2019).
- Shen, X., Pan, W., and Zhu, Y. (2012). Likelihood-based selection and sharp parameter estimation. *Journal of the American Statistical Association*, 107(497):223–232.
- Van Der Vaart, A. W. and Wellner, J. A. (1996). Weak convergence. In *Weak convergence and empirical processes*, pages 16–28. Springer.
- Wong, W. H. and Shen, X. (1995). Probability inequalities for likelihood ratios and convergence rates of sieve MLEs. *The Annals of Statistics*, 23(2):339–362.
- Xu, G. and Shang, Z. (2018). Identifying latent structures in restricted latent class models. *Journal of the American Statistical Association*, 113(523):1284–1295.