# Supplement to "Structure-Preserving Embedding of Multi-layer Networks"

**Anonymous Author(s)**
Affiliation
Address
`email`

1 This supplement file contains the appendixes of the paper "Structure-Preserving Embedding of Multi-
2 layer Networks". In Appendix A, we summarize the projected gradient descent algorithm developed
3 in Section 2 of the paper. Appendix B contains the detailed cross-validation procedure in selecting
4 the tuning parameter $\lambda_n$. Additional simulation studies are provided in Appendix C. In Appendix D,
5 we provide an eigenvalue plot of the WAT dataset discussed in Section 4.2 of the paper. All technical
6 proofs and necessary lemmas are included in Appendix E.

## A Summary of the projected gradient descent algorithm

8 For easy of presentation, we denote the projection result from Step 1 to Step 3 discussed in Section
9 2.3 as $P_{\Omega_{\boldsymbol{\alpha}} \times \Omega_{\boldsymbol{\beta}}}(\tilde{\boldsymbol{\alpha}}, \tilde{\boldsymbol{\beta}})$. The developed projected gradient descent algorithm can be summarized in
10 Algorithm 1.

---

**Algorithm 1:** Projected gradient descent (PGD)

**Input** : Adjacency tensor $\boldsymbol{\mathcal{A}}$, sparsity factor $s_n$, number of communities $K$, embedding
dimension $R$, constraint parameter $\xi$, tuning parameter $\lambda_n$, learning rate $\eta$, number of
iterations $T$.

**Output** : Estimators of $\boldsymbol{\alpha}$ and $\boldsymbol{\beta}$, estimated vertex community memberships and community
centers.

1 Initialize $\boldsymbol{\alpha}^{(0)}$, $\boldsymbol{\beta}^{(0)}$ and hence obtain $\boldsymbol{Z}^{(0)}$, $\boldsymbol{C}^{(0)}$ by $(1+\delta)$-approximation $K$-means algorithm.
Set t=0.

2 **while** $t < T$ **do**

3     $\tilde{\boldsymbol{\alpha}}^{(t+1)} = \boldsymbol{\alpha}^{(t)} - \eta \nabla_{\boldsymbol{\alpha}} \mathcal{L}_\lambda(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\beta}^{(t)}; \boldsymbol{\mathcal{A}})$, $\tilde{\boldsymbol{\beta}}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \eta \nabla_{\boldsymbol{\beta}} \mathcal{L}_\lambda(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\beta}^{(t)}; \boldsymbol{\mathcal{A}})$;

4     $(\boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\beta}^{(t+1)}) = P_{\Omega_{\boldsymbol{\alpha}} \times \Omega_{\boldsymbol{\beta}}}(\tilde{\boldsymbol{\alpha}}^{(t+1)}, \tilde{\boldsymbol{\beta}}^{(t+1)})$;

5     Apply $(1+\delta)$-approximation $K$-means algorithm to $\boldsymbol{\alpha}^{(t+1)}$ to obtain $\boldsymbol{Z}^{(t+1)}$ and $\boldsymbol{C}^{(t+1)}$.

6     **if** $\frac{|\mathcal{L}_\lambda(\boldsymbol{\alpha}^{(t+1)}, \boldsymbol{\beta}^{(t+1)}; \boldsymbol{\mathcal{A}}) - \mathcal{L}_\lambda(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\beta}^{(t)}; \boldsymbol{\mathcal{A}})|}{\mathcal{L}_\lambda(\boldsymbol{\alpha}^{(t)}, \boldsymbol{\beta}^{(t)}; \boldsymbol{\mathcal{A}})} < 10^{-6}$ **then**

7       break.

8     **end**

9     t = t+1.

10 **end**

---

## B Selecting $\lambda_n$

12 In this appendix, we provide the detailed tuning procedure for selecting $\lambda_n$. Specifically, let $\Lambda =$
13 $\{\lambda_{n1}, ..., \lambda_{nQ}\}$ be the set of $Q$ candidates of $\lambda_n$, $p_0$ be the fraction of training data, and $\kappa$ be the
14 number of repetitions. For each repetition $\kappa_0$, we first sample the training data from the adjacency
15 tensor $\boldsymbol{\mathcal{A}}$ such that $a_{i,j,m}$ will be sampled independently with probability $p_0$, for any $i \leq j$. Denote

16  $\Delta$ be the index set of the training data. For each candidate $\lambda_{nq} \in \Lambda$, we apply Algorithm 1 to solve
17  for $(\boldsymbol{\alpha}^{\kappa_0,q}, \boldsymbol{\beta}^{\kappa_0,q}) \in \Omega_{\boldsymbol{\alpha}} \times \Omega_{\boldsymbol{\beta}}$ that minimizes

$$\frac{1}{|\Delta|} \sum_{(i,j,m) \in \Delta} L(\theta_{i,j,m}, a_{i,j,m}) + \lambda_{nq} J(\boldsymbol{\alpha}), \tag{1}$$

where $|\Delta|$ is the cardinality of $\Delta$. We then evaluate the negative log-likelihood over the held-out set

$$l^{\kappa_0,q} = \frac{1}{|\Delta^c|} \sum_{(i,j,m) \in \Delta^c} L(\theta_{i,j,m}^{\kappa_0,q}, a_{i,j,m}),$$

where $\Delta^c$ is the complement of $\Delta$ and $\theta_{i,j,m}^{\kappa_0,q} = \mathcal{I} \times_1 (\boldsymbol{\alpha}_{i,.}^{\kappa_0,q})^T \times_2 (\boldsymbol{\alpha}_{j,.}^{\kappa_0,q})^T \times_3 (\boldsymbol{\beta}_{m,.}^{\kappa_0,q})^T$. Finally, we select $\lambda_n$ from $\Lambda$ such that it minimizes the averaged held-out loss over $\kappa$ repetitions; that is, $\lambda_n = \lambda_{nq^*}$ with

$$q^* = \arg \min_{q \in [Q]} \frac{1}{\kappa} \sum_{\kappa_0 = 1}^{\kappa} l^{\kappa_0,q}.$$

18  We remark that when solving (1), one needs to replace $\mathcal{T}$ by $\mathcal{T} * \mathcal{B}$, to obtain the corresponding
19  gradients associated with the training data in the PGD algorithm. Herein, $\mathcal{B} \in \{0,1\}^{n \times n \times M}$ is the
20  binary indicator tensor associated with $\Delta$ such that $\mathcal{B}_{i,j,m} = 1$ if and only if $(i,j,m) \in \Delta$. Similarly,
21  when estimating $s_n$ inside the cross-validation process by equation (7) labeled in the paper, one need
22  to replace the coefficient $\frac{1}{nM}$ by $\frac{1}{nMp_0}$ and $\mathcal{A}$ by $\mathcal{A} * \widetilde{\mathcal{B}}$, where $\widetilde{\mathcal{B}}$ is a symmetrization version of $\mathcal{B}$
23  such that $\widetilde{\mathcal{B}}_{i,j,m} = \widetilde{\mathcal{B}}_{j,i,m} = \mathcal{B}_{i,j,m}$, for $i \leq j$, $m \in [M]$.

## 24  C  Additional simulation studies

25  As network gets sparser or community sizes gets more unbalanced, it becomes more difficult to
26  differentiate vertices community memberships based on the observed multi-layer network. In this
27  Appendix, we provide additional simulation studies of two scenarios. In Scenario I, we study the
28  performances of TLSM and its competitors on networks with various sparsity, while in Scenario
29  II, we study the performances of TLSM and its competitors on networks with various levels of
30  unbalanced structures.

31  **Scenario I:** The multi-layer network generating process is the same as that descried in the paper,
32  except that we vary $(n, s_n) \in \{200, 400\} \times \{0.025i : i \in [8]\}$ and fix $(M, K) = (5, 4)$. The
33  averaged Hamming errors with 95% confidence intervals over 50 replications of all methods are
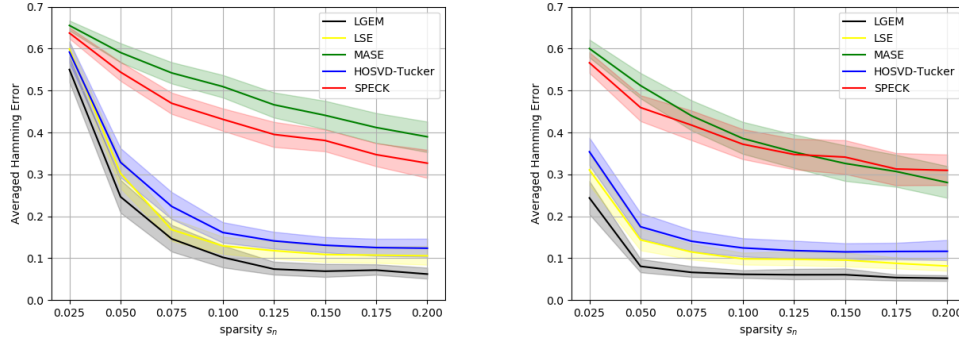34  plotted in Figure 1.



,

Figure 1: The averaged Hamming errors with 95% confidence intervals over 50 replications against various values of $s_n$ in Scenario I with $n = 200$ (Left) and 400 (Right).

35  **Scenario II:** The multi-layer network generating process is the same as that descried in the paper,
36  except that we generate $\psi \sim \text{Multi}(1, \boldsymbol{\pi})$ and vary $n \in \{200, 400\}$ while fixing $(M, K) = (5, 4)$,

where $\boldsymbol{\pi} = (\pi_1, \pi_2, \pi_3, \pi_4) = (0.25+p, 0.25+p, 0.25-p, 0.25-p)$ with $p \in \{1/24, 1/12, 1/8, 1/6\}$. The averaged Hamming errors with 95% confidence intervals over 50 replications of all methods are plotted in Figure 2.
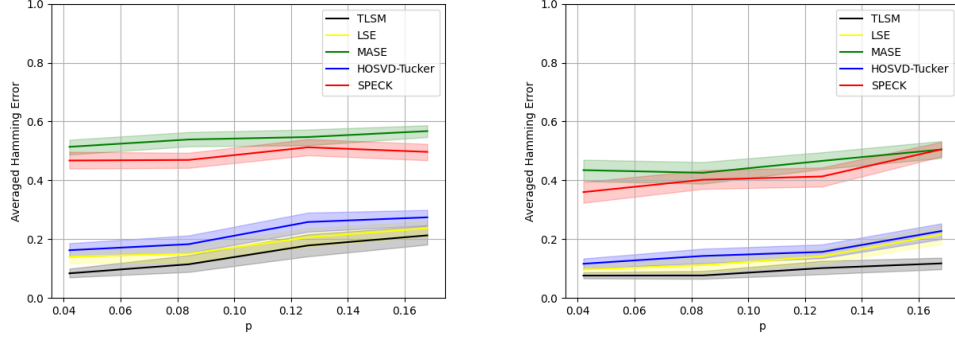


Figure 2: The averaged Hamming error with 95% confidence interval over 50 replications against various values of $p$ in Scenario II with $n = 200$ (Left) and $400$(Right).

It is evident that TLSM consistently outperforms the other competitors in both scenarios. In Scenario I, as $s_n$ becomes larger, the averaged hamming errors of all methods decrease as expected, and TLSM and LSE perform the best even for relatively small $s_n$. In Scenario II, the averaged hamming errors of all methods increase gradually when the networks get more and more unbalanced, whereas TLSM appears to be more robust against the unbalancedness.

## D  Eigenvalue plot of the WAT dataset

In this appendix, we provide a leading singular value plot of the mode-1 matricization of the WAT dataset as in Figure 3. Note that the mode-1 matricization of a tensor is to unfold it into a matrix by stacking its mode-1 fibers as the columns of its matricization. It is clear from Figure 3 that the 7th leading singular value of the mode-1 matricization of the WAT dataset is an elbow point, which suggests there are 6 potential communities among the vertices. We hence set $K = 6$ in our analysis at Section 4.2. Such an eigen-gap investigation approach has been popularly employed to determine the number of communities for a network data in literature [1, 4] when it is unknown.
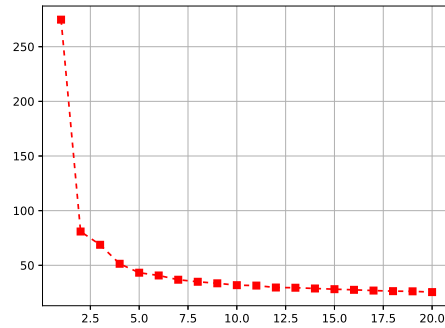


Figure 3: The first 20 leading singular values of the mode-1 matricization of the WAT dataset.

## E  Technical proofs

All Technical proofs and necessary lemmas are included in this appendix.

3

To begin with, we define the followings. For a Bernoulli random variable $Y$ with expectation $p = s_n(1 + \exp(-\theta))^{-1}$ and probability mass function $p(y; \theta)$, the discrete Hellinger distance between $p(y; \theta)$ and $p(y; \theta^*)$ is defined as

$$d(\theta, \theta^*) = \left[ \left( p^{1/2} - (p^*)^{1/2} \right)^2 + \left( (1-p)^{1/2} - (1-p^*)^{1/2} \right)^2 \right]^{1/2},$$

and the deviation of $\boldsymbol{\Theta}$ from $\boldsymbol{\Theta}^*$ can be assessed by the averaged squared Hellinger distance,

$$D^2(\boldsymbol{\Theta}, \boldsymbol{\Theta}^*) = \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} d^2(\theta_{i,j,m}, \theta^*_{i,j,m}).$$

In the proof of the main result, we will use the follow inequality several times.

**Lemma 1.** *Let $\boldsymbol{\mathcal{I}}$ be the order three $R$-dimensional identity matrix. For any matrix $\boldsymbol{A} \in \mathbb{R}^{n \times R}$, $\boldsymbol{B} \in \mathbb{R}^{n \times R}$ and $\boldsymbol{C} \in \mathbb{R}^{M \times R}$, we have*

$$||\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}||_F \leq \min\{\sqrt{M}||\boldsymbol{C}||_{vec(\infty)}, ||\boldsymbol{C}||_F\}||\boldsymbol{A}||_F||\boldsymbol{B}||_F,$$

*where $||\boldsymbol{C}||_{vec(\infty)}$ is the $l_\infty$-norm of the vectorization of $\boldsymbol{C}$.*

**Proof of Lemma 1**. The general Hölder inequality yields that the absolute value of the $(i_1, i_2, i_3)$-th entry of $\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}$ is upper bounded as

$$\left| (\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C})_{i_1,i_2,i_3} \right| = \left| \boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A}^T_{i_1,.} \times_2 \boldsymbol{B}^T_{i_2,.} \times_3 \boldsymbol{C}^T_{i_3,.} \right| \leq ||\boldsymbol{A}_{i_1,.}|| ||\boldsymbol{B}_{i_2,.}|| ||\boldsymbol{C}_{i_3,.}||_\infty.$$

Consequently,

$$||\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}||_F^2 = \sum_{i_1,i_2,i_3} \left| (\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C})_{i_1,i_2,i_3} \right|^2 \tag{2}$$

$$\leq M||\boldsymbol{C}||^2_{vec(\infty)} \sum_{i_1,i_2} ||\boldsymbol{A}_{i_1,.}||^2 ||\boldsymbol{B}_{i_2,.}||^2 = M||\boldsymbol{C}||^2_{vec(\infty)} ||\boldsymbol{A}||_F^2 ||\boldsymbol{B}||_F^2.$$

Besides, the Cauchy-Schwarz inequality implies that the absolute value of the $(i_1, i_2, i_3)$-th entry of $\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}$ is upper bounded as

$$\left| (\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C})_{i_1,i_2,i_3} \right| = \left| \sum_j \boldsymbol{A}_{i_1,j} \boldsymbol{B}_{i_2,j} \boldsymbol{C}_{i_3,j} \right| \leq ||\boldsymbol{A}_{i_1,.} * \boldsymbol{B}_{i_2,.}|| ||\boldsymbol{C}_{i_3,.}||,$$

where $\boldsymbol{A}_{i_1,.} * \boldsymbol{B}_{i_2,.}$ is the Hadamard product between $\boldsymbol{A}_{i_1,.}$ and $\boldsymbol{B}_{i_2,.}$. Note that

$$||\boldsymbol{A}_{i_1,.} * \boldsymbol{B}_{i_2,.}|| = \sqrt{\sum_j \boldsymbol{A}^2_{i_1,j} \boldsymbol{B}^2_{i_2,j}} \leq \sqrt{||\boldsymbol{A}_{i_1,.}||^2 ||\boldsymbol{B}_{i_2,.}||^2} = ||\boldsymbol{A}_{i_1,.}|| ||\boldsymbol{B}_{i_2,.}||,$$

which leads to

$$\left| (\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C})_{i_1,i_2,i_3} \right| \leq ||\boldsymbol{A}_{i_1,.}|| ||\boldsymbol{B}_{i_2,.}|| ||\boldsymbol{C}_{i_3,.}||.$$

It then follows that

$$||\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{A} \times_2 \boldsymbol{B} \times_3 \boldsymbol{C}||_F^2 \leq \sum_{i_1,i_2,i_3} ||\boldsymbol{A}_{i_1,.}||^2 ||\boldsymbol{B}_{i_2,.}||^2 ||\boldsymbol{C}_{i_3,.}||^2 = ||\boldsymbol{A}||_F^2 ||\boldsymbol{B}||_F^2 ||\boldsymbol{C}||_F^2. \tag{3}$$

Finally, the desired result immediately follows form (2) and (3). $\qquad\square$

**Proof of Proposition 1**. Denote $S = \{\boldsymbol{\Theta} \in \Omega | KL(\boldsymbol{\Theta}^* || \boldsymbol{\Theta}) \geq 4\epsilon_n\}$. Let

$$I := P\Big( \sup_S \big( \mathcal{L}_\lambda(\boldsymbol{\Theta}^*; \boldsymbol{\mathcal{A}}) - \mathcal{L}_\lambda(\boldsymbol{\Theta}; \boldsymbol{\mathcal{A}}) \big) \geq -\epsilon_n \Big)$$

$$= P\Big( \sup_S \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} \big( L(\theta^*_{i,j,m}; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m}) \big) + \lambda_n \big( J(\boldsymbol{\Theta}^*) - J(\boldsymbol{\Theta}) \big) \geq -\epsilon_n \Big).$$

We now decompose $S$ as follows. Let $S_u = \{\boldsymbol{\Theta} \in \Omega | 2^{u+1}\epsilon_n \leq KL(\boldsymbol{\Theta}^* || \boldsymbol{\Theta}) < 2^{u+2}\epsilon_n\}$, for $u = 1, 2, \dots$ . It immediately follows that $S = \bigcup_{u=1}^{+\infty} S_u$, and then

$$I \leq \sum_{u=1}^{+\infty} P\Big( \sup_{S_u} \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} \big( L(\theta^*_{i,j,m}; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m}) \big) + \lambda_n \big( J(\boldsymbol{\Theta}^*) - J(\boldsymbol{\Theta}) \big) \geq -\epsilon_n \Big)$$

$$:= \sum_{u=1}^{+\infty} I_u.$$

69 Define an empirical process $\nu_{n,M}(\Theta, \mathcal{A}) = \frac{1}{\varphi(n,M)} \sum_{m=1}^{M} \sum_{i \le j} \left( L(\theta_{i,j,m}^*; a_{i,j,m}) - \right.$

70 $L(\theta_{i,j,m}; a_{i,j,m}) - \mathbb{E}\left( L(\theta_{i,j,m}^*; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m}) \right) \right)$, for some independent but not iden-

71 tical data. It then follows that

$$I_u \le P\left( \sup_{S_u} \nu_{n,M}(\Theta, \mathcal{A}) \ge \inf_{S_u} \left( KL(\Theta^* || \Theta) + \lambda_n (J(\Theta) - J(\Theta^*)) \right) - \epsilon_n \right).$$

Since $\inf_{S_u} \left( KL(\Theta^* || \Theta) + \lambda_n (J(\Theta) - J(\Theta^*)) \right) - \epsilon_n \ge 2^{u+1}\epsilon_n - \epsilon_n - \epsilon_n \ge 2^u \epsilon_n$ and Lemma 2 shows that $\mathbb{E} \sup_{S_u} \nu_{n,M}(\Theta, \mathcal{A}) \le 2^{u-1}\epsilon_n$ when $n$ is large enough, we have

$$I_u \le P\left( \sup_{S_u} \nu_{n,M}(\Theta, \mathcal{A}) \ge 2^u \epsilon_n \right) \le P\left( \sup_{S_u} \nu_{n,M}(\Theta, \mathcal{A}) \ge \mathbb{E} \sup_{S_u} \nu_{n,M}(\Theta, \mathcal{A}) + 2^{u-1}\epsilon_n \right).$$

72 Let $Y$ be a Bernoulli random variable with expectation $p = s_n(1 + \exp(-\theta))^{-1}$, we have

$$\mathbb{E}\left( L(\theta; Y) - L(\theta^*; Y) \right) = -2p^* \log \left( (\frac{p}{p^*})^{1/2} \right) - 2(1 - p^*) \log \left( (\frac{1-p}{1-p^*})^{1/2} \right)$$

$$\ge -2p^* \left( \frac{p^{1/2}}{(p^*)^{1/2}} - 1 \right) - 2(1 - p^*) \left( \frac{(1-p)^{1/2}}{(1-p^*)^{1/2}} - 1 \right)$$

$$= \left( p^{1/2} - (p^*)^{1/2} \right)^2 + \left( (1-p)^{1/2} - (1-p^*)^{1/2} \right)^2,$$

73 where $p^* = s_n(1 + \exp(-\theta^*))^{-1}$. It immediately follows that $D^2(\Theta, \Theta^*) \le KL(\Theta^* || \Theta)$. More-

74 over, by Lagrange's mean value theorem, we further have

$$\mathbb{E}\left( L(\theta; Y) - L(\theta^*; Y) \right)^2$$

$$= 4p^* \left( \log(p^{1/2}) - \log((p^*)^{1/2}) \right)^2 + 4(1 - p^*) \left( \log((1-p)^{1/2}) - \log((1-p^*)^{1/2}) \right)^2$$

$$= 4p^* \eta_1^{-1} \left( (p^*)^{1/2} - p^{1/2} \right)^2 + 4(1 - p^*)(1 - \eta_2)^{-1} \left( (1-p^*)^{1/2} - (1-p)^{1/2} \right)^2,$$

75 where $\eta_1$ and $\eta_2$ are some real numbers between $p$ and $p^*$. Since $(1 - \xi)s_n \le p, p^* \le \xi s_n$, we

76 have $p^* \eta_1^{-1} \le \frac{\xi}{1-\xi}$ and $(1 - p^*)(1 - \eta_2)^{-1} \le \frac{1 - (1-\xi)s_n}{1 - \xi s_n} \le \frac{\xi}{1-\xi}$, which leads to $\mathbb{E}\left( L(\theta; Y) - \right.$

77 $\left. L(\theta^*; Y) \right)^2 \le \frac{4\xi}{1-\xi} d^2(\theta, \theta^*)$. On the set $S_u$, we have $KL(\Theta^* || \Theta) < 2^{u+2}\epsilon_n$. Therefore, the

78 variance of $\nu_{n,M}(\Theta, \mathcal{A})$ can be bounded as

$$Var\left( \nu_{n,M}(\Theta, \mathcal{A}) \right) \le \frac{4}{\varphi^2(n,M)} \sum_{m=1}^{M} \sum_{i \le j} \mathbb{E}\left( L(\theta_{i,j,m}; a_{i,j,m}) - L(\theta_{i,j,m}^*; a_{i,j,m}) \right)^2$$

$$\le \frac{4\xi}{(1-\xi)\varphi(n,M)} D^2(\Theta, \Theta^*) \le \frac{4\xi}{(1-\xi)\varphi(n,M)} KL(\Theta^* || \Theta) < \frac{\xi 2^{u+4}\epsilon_n}{(1-\xi)\varphi(n,M)}.$$

79 Also note that $|L(\theta; Y) - L(\theta^*; Y)|$ can be upper bounded as

$$\left| L(\theta; Y) - L(\theta^*; Y) \right| \le \max\left\{ \left| \log \frac{1 + \exp(-\theta)}{1 + \exp(-\theta^*)} \right|, \left| \log \frac{1 - s_n(1 + \exp(-\theta))^{-1}}{1 - s_n(1 + \exp(-\theta^*))^{-1}} \right| \right\} \le \log 2 + \frac{\xi}{1-\xi},$$

80 where the last inequality comes from the fact that $|\theta| \le \frac{\xi}{1-\xi}$. It follows that

$$\frac{1}{2(\log 2 + \frac{\xi}{1-\xi})} \left( L(\theta_{i,j,m}^*; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m}) - \mathbb{E}\left( L(\theta_{i,j,m}^*; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m}) \right) \right) \in [-1, 1].$$

5

81 Denote $\tilde{\nu}_{n,M}(\boldsymbol{\theta}; \mathcal{A}) = \varphi(n,M)\nu_{n,M}(\boldsymbol{\Theta}; \mathcal{A})/(2\log 2 + \frac{2\xi}{1-\xi})$. By the concentration inequality in
82 Theorem 1.1 of [2], we have

$$I_u \leq \exp\left(-\frac{\left(\varphi(n,M)2^{u-1}\epsilon_n/(2\log 2 + \frac{2\xi}{1-\xi})\right)^2}{2\left(2\mathbb{E}\sup_{S_u}\tilde{\nu}_{n,M}(\boldsymbol{\Theta},\mathcal{A}) + \sup_{S_u} Var\left(\tilde{\nu}_{n,M}(\boldsymbol{\Theta},\mathcal{A})\right)\right) + 3\frac{\varphi(n,M)}{2(\log 2 + R^{3/2}\xi^2)} * 2^{u-1}\epsilon_n}\right)$$

$$< \exp\left(-\frac{\left(\varphi(n,M)2^{u-1}\epsilon_n/(2\log 2 + \frac{2\xi}{1-\xi})\right)^2}{2\left(\frac{\varphi(n,M)}{(\log 2 + \frac{\xi}{1-\xi})} * 2^{u-1}\epsilon_n + \frac{\xi\varphi(n,M)}{4(1-\xi)(\log 2 + R^{3/2}\xi^2)^2} * 2^{u+4}\epsilon_n\right) + 3\frac{\varphi(n,M)}{2(\log 2 + \frac{\xi}{1-\xi})} * 2^{u-1}\epsilon_n}\right)$$

$$= \exp\left(-\frac{2^u\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2}\right).$$

Denote $\zeta = \exp\left(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2}\right)$. We have

$$I \leq \sum_{u=1}^{+\infty}\exp\left(-\frac{2^u\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2}\right) \leq \sum_{u=1}^{+\infty}\zeta^u = \frac{\zeta}{1-\zeta}.$$

83 As a result, $I \leq (1+I)\zeta \leq 2\zeta$. $\qquad\qquad\square$

84 **Lemma 2.** *Let the set $S_u$ and the empirical process $\nu_{n,M}(\boldsymbol{\Theta}; \mathcal{A})$ be defined in the proof of Proposition*
85 *1. If $\frac{(n+M)R}{\varphi(n,M)\epsilon_n}\log\sqrt{\frac{1}{\epsilon_n}} \leq c_1$, for some constant $c_1$ that depends on $\xi$ only, then for any $u = 1, 2, ...,$*
86 *we have $E\left(\sup_{S_u}\nu_{n,M}(\boldsymbol{\Theta}, \mathcal{A})\right) \leq 2^{u-1}\epsilon_n$.*

87 **Proof of Lemma 2.** Denote $f(\theta_{i,j,m}; a_{i,j,m}) = L(\theta^*_{i,j,m}; a_{i,j,m}) - L(\theta_{i,j,m}; a_{i,j,m})$, and hence
88 $\nu_{n,M}(\boldsymbol{\Theta}, \mathcal{A}) = \varphi^{-1}(n,M)\sum_{m=1}^{M}\sum_{i\leq j}\left(f(\theta_{i,j,m}; a_{i,j,m}) - \mathbb{E}f(\theta_{i,j,m}; a_{i,j,m})\right)$. Let $\mathcal{A}' =$
89 $(a'_{i,j,m})$ be an independent copy of $\mathcal{A}$ and $\boldsymbol{\tau} = (\tau_{i,j,m})$ be a collection of independent Rademacher
90 random variables. By the standard symmetrization argument, we have

$$\mathbb{E}_{\mathcal{A}}\sup_{S_u}\nu_{n,M}(\boldsymbol{\Theta}, \mathcal{A}) \leq \frac{1}{\varphi(n,M)}\mathbb{E}_{\mathcal{A},\mathcal{A}'}\sup_{S_u}\sum_{m=1}^{M}\sum_{i\leq j}\left(f(\theta_{i,j,m}; a_{i,j,m}) - f(\theta_{i,j,m}; a'_{i,j,m})\right)$$

$$\leq \frac{2}{\varphi(n,M)}\mathbb{E}_{\mathcal{A},\boldsymbol{\tau}}\sup_{S_u}\left|\sum_{m=1}^{M}\sum_{i\leq j}\tau_{i,j,m}f(\theta_{i,j,m}; a_{i,j,m})\right|.$$

91 Denote $X(\boldsymbol{\Theta}; \mathcal{A}) = \varphi^{-1/2}(n,M)\sum_{m=1}^{M}\sum_{i\leq j}\tau_{i,j,m}f(\theta_{i,j,m}; a_{i,j,m})$ as the conditional
92 Rademacher process. For any $\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)} \in S_u$, and $\omega \in \mathbb{R}$, we have $\mathbb{E}_{\boldsymbol{\tau}|\mathcal{A}}\exp\left(\omega\left(X(\boldsymbol{\Theta}^{(1)}; \mathcal{A}) - \right.\right.$
93 $\left.\left.X(\boldsymbol{\Theta}^{(2)}; \mathcal{A})\right)\right) \leq \exp\left(\frac{1}{2}\omega^2\rho^2(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)}; \mathcal{A})\right)$, where

$$\rho^2(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)}; \mathcal{A}) = \frac{1}{\varphi(n,M)}\sum_{m=1}^{M}\sum_{i\leq j}\left(f(\theta^{(1)}_{i,j,m}; a_{i,j,m}) - f(\theta^{(2)}_{i,j,m}; a_{i,j,m})\right)^2,$$

94 showing that $X(\boldsymbol{\Theta}; \mathcal{A})$ is a sub-Gaussian process with respect to $\rho$ when $\mathcal{A}$ is given. Thus, by
95 Theorem 3.11 of [3], there exists a positive constant $c_4$, such that

$$\varphi^{-1/2}(n,M)\mathbb{E}_{\mathcal{A},\boldsymbol{\tau}}\sup_{S_u}\left|\sum_{m=1}^{M}\sum_{i\leq j}\tau_{i,j,m}f(\theta_{i,j,m}; a_{i,j,m})\right| \leq \frac{c_4}{2}\mathbb{E}_{\mathcal{A}}\int_0^{\text{diam}(S_u)}H^{1/2}(\varepsilon; S_u, \rho)d\varepsilon,$$

where $\text{diam}(S_u)$ is the diameter of $S_u$ and $H(\varepsilon; S_u, \rho)$ is the metric entropy. Note that
$\left|\frac{dL(\theta_{i,j,m}; a_{i,j,m})}{d\theta_{i,j,m}}\right| = \left|\frac{\exp(-\theta_{i,j,m})}{1-s_n+\exp(-\theta_{i,j,m})}(p_{i,j,m} - a_{i,j,m})\right| < 1$. Thus, both $L(\theta_{i,j,m}; a_{i,j,m})$ and
$f(\theta_{i,j,m}; a_{i,j,m})$ are Lipschitz continuous with Lipschitz constant 1. Thus, for any $\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)} \in S_u$,
we have

$$\rho^2(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)}; \mathcal{A}) \leq \frac{1}{\varphi(n,M)}\sum_{m=1}^{M}\sum_{i\leq j}|\theta^{(1)}_{i,j,m} - \theta^{(2)}_{i,j,m}|^2 \leq \frac{1}{\varphi(n,M)}||\boldsymbol{\Theta}^{(1)} - \boldsymbol{\Theta}^{(2)}||_F^2.$$

6

96    By the triangle inequality and Lemma 1,

$$
\rho(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)}; \boldsymbol{\mathcal{A}}) \leq \frac{1}{\varphi^{1/2}(n, M)} \Big( ||\boldsymbol{\mathcal{I}} \times_1 (\boldsymbol{\alpha}^{(1)} - \boldsymbol{\alpha}^{(2)}) \times_2 \boldsymbol{\alpha}^{(1)} \times_3 \boldsymbol{\beta}^{(1)}||_F
$$

$$
+ ||\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{\alpha}^{(2)} \times_2 (\boldsymbol{\alpha}^{(1)} - \boldsymbol{\alpha}^{(2)}) \times_3 \boldsymbol{\beta}^{(1)}||_F + ||\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{\alpha}^{(2)} \times_2 \boldsymbol{\alpha}^{(2)} \times_3 (\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^{(2)}||_F \Big)
$$

$$
\leq \frac{\left( \sqrt{\min\{M, R\}} ||\boldsymbol{\alpha}^{(1)} - \boldsymbol{\alpha}^{(2)}||_F (||\boldsymbol{\alpha}^{(1)}||_F + ||\boldsymbol{\alpha}^{(2)}||_F) + \min\{2\sqrt{M}, ||\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^{(2)}||_F\} ||\boldsymbol{\alpha}^{(2)}||_F^2 \right)}{\varphi^{1/2}(n, M)}
$$

$$
\leq \frac{n \log \frac{\xi}{1-\xi}}{\varphi^{1/2}(n, M)} \Big( 2\sqrt{\min\{M, R\}} ||\frac{1}{\sqrt{\log \frac{\xi}{1-\xi}}} (\boldsymbol{\alpha}^{(1)} - \boldsymbol{\alpha}^{(2)})||_F + \sqrt{R} \min\{2\sqrt{\frac{M}{R}}, \frac{1}{\sqrt{R}} ||\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^{(2)}||_F\} \Big).
$$

97    This leads to

$$
H(\varepsilon; S_u, \rho) \leq H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{4n\sqrt{\min\{M, R\}} \log \frac{\xi}{1-\xi}}; B(nR), ||\cdot|| \Big) + H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{2n\sqrt{R} \log \frac{\xi}{1-\xi}}; B^h(MR), h \Big),
$$

98    where $B(nR)$ is the unit ball with respect to the $l_2$-norm in $\mathbb{R}^{nR}$, $B^h(MR)$ is the Euclidean
99    ball in $\mathbb{R}^{MR}$ with radius $\min\{2\sqrt{\frac{M}{R}}, 1\}$, $h$ is a truncated distance such that $h(\boldsymbol{\beta}^{(1)}, \boldsymbol{\beta}^{(2)}) = $
100    $\min\{2\sqrt{\frac{M}{R}}, \frac{1}{\sqrt{R}} ||\boldsymbol{\beta}^{(1)} - \boldsymbol{\beta}^{(2)}||_F\}$, and $H(\cdot; \cdot, \cdot)$ is the metric entropy.
101    In the case that $2\sqrt{M} \geq \sqrt{R}$, we have

$$
H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{2n\sqrt{R} \log \frac{\xi}{1-\xi}}; B^h(MR), h \Big) = H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{2n\sqrt{R} \log \frac{\xi}{1-\xi}}; B(MR), ||\cdot|| \Big)
$$

$$
\leq MR \log \frac{6n\sqrt{R} \log \frac{\xi}{1-\xi}}{\varphi^{1/2}(n, M)\epsilon} \leq MR \log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon},
$$

102    where $B(MR)$ is the unit ball with respect to the $l_2$ norm in $\mathbb{R}^{MR}$. In the case that $2\sqrt{M} < \sqrt{R}$, we
103    have

$$
H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{2n\sqrt{R} \log \frac{\xi}{1-\xi}}; B^h(MR), h \Big) = H\Big( \frac{\varphi^{1/2}(n, M)\varepsilon}{2n\sqrt{R} \log \frac{\xi}{1-\xi}} \cdot \frac{\sqrt{R}}{2\sqrt{M}}; B(MR), ||\cdot|| \Big) \leq MR \log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon}.
$$

104    Thus, $H(\varepsilon; S_u, \rho)$ can be bounded as

$$
H(\varepsilon; S_u, \rho) \leq nR \log \frac{12n\sqrt{\min\{M, R\}} \log \frac{\xi}{1-\xi}}{\varphi^{1/2}(n, M)\epsilon} + MR \log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon} \leq (n+M)R \log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon}.
$$

105    By concavity,

$$
\mathbb{E}_{\boldsymbol{\mathcal{A}}} \sup_{S_u} \nu_{n,M}(\boldsymbol{\Theta}; \boldsymbol{\mathcal{A}}) \leq \frac{c_4}{\varphi^{1/2}(n, M)} \mathbb{E}_{\boldsymbol{\mathcal{A}}} \int_0^{\mathrm{diam}(S_u)} \sqrt{(n+M)R \log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon}} \, d\varepsilon
$$

$$
\leq c_4 \sqrt{\frac{(n+M)R}{\varphi(n, M)}} \int_0^{\sqrt{\mathbb{E}_{\boldsymbol{\mathcal{A}}} \mathrm{diam}^2(S_u)}} \sqrt{\log \frac{12\sqrt{2} \log \frac{\xi}{1-\xi}}{\epsilon}} \, d\varepsilon.
$$

106    Furthermore, according to the same argument of bounding the variance of $\nu_{n,M}(\boldsymbol{\Theta}, \boldsymbol{\mathcal{A}})$, we
107    have $\mathbb{E}_{\boldsymbol{\mathcal{A}}} \rho^2(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^{(2)}; \boldsymbol{\mathcal{A}}) \leq 2\big( \mathbb{E}_{\boldsymbol{\mathcal{A}}} \rho^2(\boldsymbol{\Theta}^{(1)}, \boldsymbol{\Theta}^*; \boldsymbol{\mathcal{A}}) + \mathbb{E}_{\boldsymbol{\mathcal{A}}} \rho^2(\boldsymbol{\Theta}^{(2)}, \boldsymbol{\Theta}^*; \boldsymbol{\mathcal{A}}) \big) \leq 2\big( \frac{\xi}{1-\xi} 2^{u+2} \epsilon_n +$

108    $\frac{\xi}{1-\xi}2^{u+2}\epsilon_n) = \frac{\xi}{1-\xi}2^{u+4}\epsilon_n$, implying that $\mathbb{E}_{\boldsymbol{\mathcal{A}}}\text{diam}^2(S_u) \leq \frac{\xi}{1-\xi}2^{u+4}\epsilon_n$. Thus,

$$\mathbb{E}_{\boldsymbol{\mathcal{A}}}\sup_{S_u}\nu_{n,M}(\boldsymbol{\Theta},\boldsymbol{\mathcal{A}}) \leq c_4\sqrt{\frac{(n+M)R}{\varphi(n,M)}}\int_0^{\sqrt{\frac{\xi}{1-\xi}2^{(u+4)\epsilon_n}}}\sqrt{\log\frac{12\sqrt{2}\log\frac{\xi}{1-\xi}}{\epsilon}}d\varepsilon$$

$$\leq \frac{12\sqrt{2}c_4\sqrt{(n+M)R}\log\frac{\xi}{1-\xi}}{\sqrt{\varphi(n,M)\log\frac{12\sqrt{2}\log\frac{\xi}{1-\xi}}{\sqrt{\frac{\xi}{1-\xi}2^{u+4}\epsilon_n}}}}\int_{12\sqrt{2}\log\frac{\xi}{1-\xi}/\sqrt{\frac{\xi}{1-\xi}2^{u+4}\epsilon_n}}^{+\infty}\frac{\log\varepsilon}{\varepsilon^2}d\varepsilon$$

$$= c_4\sqrt{\frac{2^{u+4}(n+M)R\frac{\xi}{1-\xi}\epsilon_n}{\varphi(n,M)\log\frac{12\sqrt{2}\log\frac{\xi}{1-\xi}}{\sqrt{\frac{\xi}{1-\xi}2^{u+4}\epsilon_n}}}}\Big(1+\log\frac{12\sqrt{2}\log\frac{\xi}{1-\xi}}{\sqrt{\frac{\xi}{1-\xi}2^{u+4}\epsilon_n}}\Big)$$

$$\leq c_5\sqrt{\frac{2^{u+4}(n+M)R\epsilon_n}{\varphi(n,M)}\log\sqrt{\frac{1}{\epsilon_n}}},$$

109    for some positive constant $c_5$ that depends on $\xi$ only. Finally,

$$\mathbb{E}_{\boldsymbol{\mathcal{A}}}\sup_{S_u}\nu_{n,M}(\boldsymbol{\Theta},\boldsymbol{\mathcal{A}}) \leq 4\sqrt{2}c_5\sqrt{\frac{(n+M)R}{\varphi(n,M)\epsilon_n}\log\sqrt{\frac{1}{\epsilon_n}}}\cdot 2^{u-1}\epsilon_n \leq 2^{u-1}\epsilon_n, \tag{4}$$

110    where the second inequality follows from the condition that $\frac{(n+M)R}{\varphi(n,M)\epsilon_n}\log\sqrt{\frac{1}{\epsilon_n}} \leq c_1$ with $c_1$ taking

111    to be $\frac{1}{32c_5^2}$.      $\square$

112    **Proof of Theorem 1**. By definition of $\widehat{\boldsymbol{\Theta}}$, it follows from Proposition 1 that

$$P(D^2(\widehat{\boldsymbol{\Theta}},\boldsymbol{\Theta}^*) \geq 4\epsilon_n) \leq P\big(KL(\boldsymbol{\Theta}^*||\widehat{\boldsymbol{\Theta}}) \geq 4\epsilon_n\big)$$

$$\leq P\big(\sup_{\{\boldsymbol{\Theta}\in\Omega|KL(\boldsymbol{\Theta}^*||\boldsymbol{\Theta})\geq 4\epsilon_n\}}\mathcal{L}_\lambda(\boldsymbol{\Theta}^*)-\mathcal{L}_\lambda(\boldsymbol{\Theta}) \geq -\epsilon_n\big)$$

$$\leq 2\exp\Big(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\Big).$$

113    That is, with probability at least $1-2\exp\Big(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\Big)$, $D^2(\widehat{\boldsymbol{\Theta}},\boldsymbol{\Theta}^*) \leq 4\epsilon_n$.

114    Next, we bound the $F$-norm of the different between $\widehat{\boldsymbol{\Theta}}$ and $\boldsymbol{\Theta}^*$. Let $g(x) = \log\frac{x^2}{s_n-x^2}$. By
115    Lagrange's mean value theorem, for any $\boldsymbol{\Theta}\in\Omega$,

$$|\theta_{i,j,m}-\theta^*_{i,j,m}| = |g(p_{i,j,m}^{1/2})-g((p^*_{i,j,m})^{1/2})| \leq \max\{\frac{2}{\sqrt{(1-\xi)s_n\xi}},\frac{2}{\sqrt{\xi s_n(1-\xi)}}\}|p_{ijm}^{1/2}-(p^*_{ijm})^{1/2}|.$$

Moreover, $\xi > 1/2$ implies that $\max\{\frac{2}{\sqrt{(1-\xi)s_n\xi}},\frac{2}{\sqrt{\xi s_n(1-\xi)}}\} = \frac{2}{\sqrt{\xi s_n(1-\xi)}}$. It then follows that
$\frac{1}{\varphi(n,M)}\sum_{m=1}^M\sum_{i\leq j}(\theta_{ijm}-\theta^*_{ijm})^2 \leq \frac{4}{s_n\xi(1-\xi)^2}D^2(\boldsymbol{\Theta},\boldsymbol{\Theta}^*)$. Particularly, for the estimator $\widehat{\boldsymbol{\Theta}}$, we
have

$$\frac{1}{n^2M}\|\widehat{\boldsymbol{\Theta}}-\boldsymbol{\Theta}^*\|_F^2 \leq \frac{8\varphi(n,M)}{n^2Ms_n\xi(1-\xi)^2}D^2(\widehat{\boldsymbol{\Theta}};\boldsymbol{\Theta}^*) \leq \frac{8}{s_n\xi(1-\xi)^2}D^2(\widehat{\boldsymbol{\Theta}};\boldsymbol{\Theta}^*) \leq \frac{32\epsilon_n}{s_n\xi(1-\xi)^2},$$

116    whit probability at least $1-2\exp\Big(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\Big)$.      $\square$

117    **Lemma 3.** *Under the conditions of Theorem 1 and Assumption B, then there exists an absolute*
118    *constant $c_3$ that depends on $\xi$ only, such that*

$$\frac{1}{n\sqrt{M}}\|\boldsymbol{\mathcal{I}}\times_1\hat{\boldsymbol{Z}}\hat{\boldsymbol{C}}\times_2\hat{\boldsymbol{Z}}\hat{\boldsymbol{C}}\times_3\hat{\boldsymbol{\beta}}-\boldsymbol{\mathcal{I}}\times_1\boldsymbol{Z}^*\boldsymbol{C}^*\times_2\boldsymbol{Z}^*\boldsymbol{C}^*\times_3\boldsymbol{\beta}^*\|_F$$

$$\leq\Big(\frac{4\sqrt{2}}{(1-\xi)\sqrt{\xi}}+c_3\sqrt{\frac{(1+\delta)\min\{M,R\}}{M}}\Big)\sqrt{\epsilon_ns_n^{-1}},$$

119    *with probability at least $1-2\exp\Big(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\Big)-n^{-2}$.*

8

120 **Proof of Lemma 3**. We first provide a probabilistic upper bound for $J(\hat{\boldsymbol{\alpha}})$. Note that

$$\mathcal{L}(\boldsymbol{\Theta}^*, \boldsymbol{\mathcal{A}}) = \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} L(\theta_{i,j,m}^*; a_{i,j,m})$$

$$= \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} \left( a_{i,j,m} \log \frac{1 - p_{i,j,m}^*}{p_{i,j,m}^*} + \log \frac{1}{1 - p_{i,j,m}^*} \right).$$

Denote $X_{i,j,m} = a_{i,j,m} \log \frac{1-p_{i,j,m}^*}{p_{i,j,m}^*} + \log \frac{1}{1-p_{i,j,m}^*}$, for $i \leq j$, $m \in [M]$. It follows that $\mathcal{L}(\boldsymbol{\Theta}^*, \boldsymbol{\mathcal{A}})$ is the average of $\varphi(n, M)$ independent two-value random variables with $|X_{i,j,m}| \leq c_6 \log \frac{1}{s_n}$, $\mathbb{E}X_{i,j,m} \leq c_6 s_n \log \frac{1}{s_n}$ and $\mathbb{E}X_{i,j,m}^2 \leq c_6 s_n (\log \frac{1}{s_n})^2$, where $c_6$ is a constant that depends on $\xi$ only. By Bernstein inequality, for any $t > 0$,

$$P\left( \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} (X_{i,j,m} - \mathbb{E}X_{i,j,m}) > t \right) \leq \exp\left\{ -\frac{\frac{1}{2}\varphi^2(n, M)t^2}{c_5 \varphi(n, M) s_n (\log \frac{1}{s_n})^2 + c_5 \varphi(n, M) t \log \frac{1}{s_n}/3} \right\}.$$

Taking $t = \sqrt{6c_6} \varphi^{-1/2}(n, M) s_n^{1/2} (\log \frac{1}{s_n})(\log n)^{1/2}$, with probability at least $1 - n^{-2}$, we have

$$\lambda_n J(\hat{\boldsymbol{\alpha}}) < \mathcal{L}_\lambda(\widehat{\boldsymbol{\Theta}}; \boldsymbol{\mathcal{A}}) \leq \mathcal{L}_\lambda(\boldsymbol{\Theta}^*; \boldsymbol{\mathcal{A}}) + \epsilon_n \leq \frac{1}{\varphi(n, M)} \sum_{m=1}^{M} \sum_{i \leq j} \mathbb{E}X_{i,j,m} + t + \epsilon_n \leq c_6 s_n \log \frac{1}{s_n} + t + \epsilon_n.$$

121 Clearly $t = o(s_n \log \frac{1}{s_n})$ and $\epsilon_n = o(s_n \log \frac{1}{s_n})$. Thus, the assumption $\lambda_n \epsilon_n s_n^{-2} (\log s_n^{-1})^{-1} \geq c_2$

122 immediately implies that $J(\hat{\boldsymbol{\alpha}}) \leq \frac{(c_7-1)^2}{4} \epsilon_n s_n^{-1}$, for some constant $c_7 > 1$, with probability at least

123 $1 - n^{-2}$.

124 We now turn to bound the difference between $\boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_2 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_3 \hat{\boldsymbol{\beta}}$ and $\boldsymbol{\mathcal{I}} \times_1 \boldsymbol{Z}^*\boldsymbol{C}^* \times_2 \boldsymbol{Z}^*\boldsymbol{C}^* \times_3 \boldsymbol{\beta}^*$.

125 Applying the triangle inequality and Lemma 1 yields that

$$\frac{1}{n\sqrt{M}} \left\| \boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{\alpha}} \times_2 \hat{\boldsymbol{\alpha}} \times_3 \hat{\boldsymbol{\beta}} - \boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_2 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_3 \hat{\boldsymbol{\beta}} \right\|_F$$

$$\leq \frac{1}{n\sqrt{M}} \|\hat{\boldsymbol{\alpha}} - \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}}\|_F (\|\hat{\boldsymbol{\alpha}}\|_F + \|\hat{\boldsymbol{Z}}\hat{\boldsymbol{C}}\|_F) \min\{\sqrt{M}, \sqrt{R}\}$$

$$\leq 2\sqrt{\frac{\log \frac{\xi}{1-\xi}}{M}} \sqrt{(1+\delta)J(\hat{\boldsymbol{\alpha}})} \min\{\sqrt{M}, \sqrt{R}\} \tag{5}$$

$$\leq (c_7 - 1)\sqrt{\frac{\min\{M, R\}}{M}} \sqrt{(1+\delta)\epsilon_n s_n^{-1} \log \frac{\xi}{1-\xi}},$$

126 with probability at least $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) - n^{-2}$. Similarly,

$$\frac{1}{n\sqrt{M}} \left\| \boldsymbol{\mathcal{I}} \times_1 \boldsymbol{\alpha}^* \times_2 \boldsymbol{\alpha}^* \times_3 \boldsymbol{\beta}^* - \boldsymbol{\mathcal{I}} \times_1 \boldsymbol{Z}^*\boldsymbol{C}^* \times_2 \boldsymbol{Z}^*\boldsymbol{C}^* \times_3 \boldsymbol{\beta}^* \right\|_F \tag{6}$$

$$\leq 2\sqrt{\frac{\min\{M, R\}}{M}} \sqrt{J(\boldsymbol{\alpha}^*) \log \frac{\xi}{1-\xi}} = o\left( \sqrt{\frac{\min\{M, R\}}{M}} \sqrt{\epsilon_n s_n^{-1} \log \frac{\xi}{1-\xi}} \right),$$

127 where the equality follows from $\lambda_n J(\boldsymbol{\alpha}^*) \leq \epsilon_n$ and Assumption B. Finally, by (5), (6) and Theorem

128 1, we have

$$\frac{1}{n\sqrt{M}} \|\boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_2 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_3 \hat{\boldsymbol{\beta}} - \boldsymbol{\mathcal{I}} \times_1 \boldsymbol{Z}^*\boldsymbol{C}^* \times_2 \boldsymbol{Z}^*\boldsymbol{C}^* \times_3 \boldsymbol{\beta}^* \|_F$$

$$\leq \frac{1}{n\sqrt{M}} \|\boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_2 \hat{\boldsymbol{Z}}\hat{\boldsymbol{C}} \times_3 \hat{\boldsymbol{\beta}} - \widehat{\boldsymbol{\Theta}}\|_F + \frac{1}{n\sqrt{M}} \|\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}^*\|_F$$

$$+ \frac{1}{n\sqrt{M}} \|\boldsymbol{\Theta}^* - \boldsymbol{\mathcal{I}} \times_1 \boldsymbol{Z}^*\boldsymbol{C}^* \times_2 \boldsymbol{Z}^*\boldsymbol{C}^* \times_3 \boldsymbol{\beta}^* \|_F$$

$$\leq \left( \frac{4\sqrt{2}}{(1-\xi)\sqrt{\xi}} + c_7 \sqrt{\frac{\min\{M, R\}}{M}} \sqrt{(1+\delta)\log \frac{\xi}{1-\xi}} \right) \sqrt{\epsilon_n s_n^{-1}}.$$

9

129 The desired result follows by taking $c_3 = c_7 \sqrt{\log \frac{\xi}{1-\xi}}$.

130 $\square$

**Lemma 4.** *Let* $\widehat{\boldsymbol{\mathcal{B}}} = \boldsymbol{\mathcal{I}} \times_1 \hat{\boldsymbol{C}} \times_2 \hat{\boldsymbol{C}} \times_3 \hat{\boldsymbol{\beta}}$ *be the estimation counterpart of* $\boldsymbol{\mathcal{B}}^*$. *Denote* $\boldsymbol{\mathcal{M}}^* = \boldsymbol{\mathcal{B}}^* \times_2 \boldsymbol{Z}^*$ *and* $\widehat{\boldsymbol{\mathcal{M}}} = \widehat{\boldsymbol{\mathcal{B}}} \times_2 \hat{\boldsymbol{Z}}$. *Under the conditions of Lemma 3 and Assumption A and C, then with probability at least* $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) - n^{-2}$, *the following event F holds. F: for any* $k \in [K]$, *there exists an unique* $k' \in [K]$, *such that*

$$\frac{1}{\sqrt{nM}} \left\| \widehat{\boldsymbol{\mathcal{M}}}_{k',\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{k,\cdot,\cdot} \right\|_F = o(\gamma_n \sqrt{\frac{n_{\min} K}{n}}).$$

131 **Proof of Lemma 4.** Denote $F_0$ be the event that there exists $k \in [K]$ such that $\frac{1}{\sqrt{nM}} \| \widehat{\boldsymbol{\mathcal{M}}}_{k',\cdot,\cdot} -$

132 $\boldsymbol{\mathcal{M}}^*_{k,\cdot,\cdot} \|_F \geq \gamma_n \sqrt{\frac{n_{\min} K}{n}}$, for any $k' \in [K]$ and sufficiently large $n$ and $M$. It follows that

$$\frac{1}{n\sqrt{M}} \left\| \widehat{\boldsymbol{\mathcal{B}}} \times_1 \hat{\boldsymbol{Z}} \times_2 \hat{\boldsymbol{Z}} - \boldsymbol{\mathcal{B}}^* \times_1 \boldsymbol{Z}^* \times_2 \boldsymbol{Z}^* \right\|_F \geq \frac{1}{n\sqrt{M}} \left( \sum_{i \in N^*_k} \left\| (\widehat{\boldsymbol{\mathcal{M}}} \times_1 \hat{\boldsymbol{Z}})_{i,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{k,\cdot,\cdot} \right\|^2_F \right)^{1/2}$$

$$\geq \gamma_n \sqrt{\frac{n_{\min} K}{n}} \sqrt{\frac{n_k}{n}} \geq \frac{\gamma_n n_{\min} \sqrt{K}}{n}.$$

133 However, it follows from Lemma 3 and the Assumption C that

$$P(F_0) \leq P\left( \frac{1}{n\sqrt{M}} \left\| \widehat{\boldsymbol{\mathcal{B}}} \times_1 \hat{\boldsymbol{Z}} \times_2 \hat{\boldsymbol{Z}} - \boldsymbol{\mathcal{B}}^* \times_1 \boldsymbol{Z}^* \times \boldsymbol{Z}^* \right\|_F \geq \frac{\gamma_n n_{\min} \sqrt{K}}{n} \right)$$

$$\leq 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) + n^{-2}.$$

134 Therefore, with probability at least $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) - n^{-2}$, $F^c_0$, the complement of

135 $F_0$ holds; that is; the existence holds with high probability.

We now prove the uniqueness under $F^c_0$. Assume there exist $k_1 \neq k_2 \in [K]$ such that $\frac{1}{\sqrt{nM}} \| \widehat{\boldsymbol{\mathcal{M}}}_{k_i,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{k,\cdot,\cdot} \|_F = o(\gamma_n \sqrt{\frac{n_{\min} K}{n}})$, for $i \in [2]$. By existence, there exists $a \in [K]$ and $b_1 \neq b_2 \in [K]$ such that $\frac{1}{\sqrt{nM}} \| \widehat{\boldsymbol{\mathcal{M}}}_{a,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{b_j,\cdot,\cdot} \|_F = o(\gamma_n \sqrt{\frac{n_{\min} K}{n}})$, for $j \in [2]$. The triangle inequality implies that

$$\frac{1}{\sqrt{nM}} \left\| \boldsymbol{\mathcal{M}}^*_{b_1,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{b_2,\cdot,\cdot} \right\|_F \leq \frac{1}{\sqrt{nM}} \left( \left\| \widehat{\boldsymbol{\mathcal{M}}}_{a,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{b_1,\cdot,\cdot} \right\|_F + \left\| \widehat{\boldsymbol{\mathcal{M}}}_{a,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{b_2,\cdot,\cdot} \right\|_F \right) = o(\gamma_n \sqrt{\frac{n_{\min} K}{n}}).$$

On the other hand, Assumption A implies that

$$\frac{1}{\sqrt{nM}} \left\| \boldsymbol{\mathcal{M}}^*_{b_1,\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{b_2,\cdot,\cdot} \right\|_F \geq \sqrt{\frac{n_{\min}}{nM}} \left\| \boldsymbol{\mathcal{B}}^*_{b_1,\cdot,\cdot} - \boldsymbol{\mathcal{B}}^*_{b_2,\cdot,\cdot} \right\|_F \geq \sqrt{\frac{n_{\min} K}{n}} \gamma_n,$$

136 which is a contradiction. Hence, $F^c_0$ also implies uniqueness, showing that $F$ holds with probability

137 at least $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) - n^{-2}$. $\square$

138 **Proof of Theorem 2.** Based on Lemma 4, with probability at least $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) -$

139 $n^{-2}$, there exists a permutation $\pi^* \in S_K$ such that for each $k \in [K]$, $\frac{1}{\sqrt{nM}} \| \widehat{\boldsymbol{\mathcal{M}}}_{\pi^*(k),\cdot,\cdot} - \boldsymbol{\mathcal{M}}^*_{k,\cdot,\cdot} \|_F =$

140 $o(\gamma_n \sqrt{\frac{n_{\min} K}{n}})$. It then suffices to show that with probability at least $1 - 2\exp\left( -\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi} + 28\log 2} \right) -$

141 $n^{-2}$, it holds true that $\min_{\pi \in S_K} \frac{1}{n} \sum_{i=1}^{n} \mathbf{1}\{\psi^*_i \neq \pi(\hat{\psi}_i)\} \leq \frac{c^2_\xi n\epsilon_n}{n_{\min} K \gamma^2_n s_n}$. Let $\hat{N}_k = \{i : \hat{\psi}_i = k\}$, for

142 $k \in [K]$. Note that

$$\min_{\pi \in S_K} \sum_{i=1}^{n} \mathbf{1}\{\psi^*_i \neq \pi(\hat{\psi}_i)\} = \min_{\pi \in S_K} \sum_{k=1}^{K} |N^*_k \setminus \hat{N}_{\pi^{-1}(k)}| = \min_{\pi \in S_K} \sum_{k=1}^{K} |N^*_k \setminus \hat{N}_{\pi(k)}|,$$

143 where the last equality follows from the fact that $\pi^{-1}$ is also a permutation in $S_K$. It then suffices to

144 show that with probability at least $1 - 2\exp\left(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\right) - n^{-2}$, $\frac{1}{n}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq$

145 $\frac{c_\xi^2 n\epsilon_n}{n_{\min}K\gamma_n^2 s_n}$ for the particular permutation $\pi^*$. Let $F$ denote the same event in Lemma 4. In fact, by

146 Lemma 3, we have

$$P\Big(\frac{1}{n}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{c_\xi^2 n\epsilon_n}{n_{\min}K\gamma_n^2 s_n}|F\Big) = P\Big(\frac{n_{\min}K\gamma_n^2}{n^2}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| = \frac{c_\xi^2\epsilon_n}{s_n}|F\Big)$$

$$\geq P\Big(\frac{n_{\min}K\gamma_n^2}{n^2}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{1}{n^2 M}\big\|\widehat{\mathcal{B}}\times_1\hat{\mathbf{Z}}\times_2\hat{\mathbf{Z}} - \mathcal{B}^*\times_1\mathbf{Z}^*\times_2\mathbf{Z}^*\big\|_F^2|F\Big)+$$

$$P\Big(\frac{1}{n^2 M}\big\|\widehat{\mathcal{B}}\times_1\hat{\mathbf{Z}}\times_2\hat{\mathbf{Z}} - \mathcal{B}^*\times_1\mathbf{Z}^*\times_2\mathbf{Z}^*\big\|_F^2 \leq \frac{c_\xi^2\epsilon_n}{s_n}|F\Big) - 1$$

$$= P\Big(\frac{n_{\min}K\gamma_n^2}{n^2}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{1}{n^2 M}\big\|\widehat{\mathcal{B}}\times_1\hat{\mathbf{Z}}\times_2\hat{\mathbf{Z}} - \mathcal{B}^*\times_1\mathbf{Z}^*\times_2\mathbf{Z}^*\big\|_F^2|F\Big),$$

147 where the last equality comes from the fact that the event $F$ is based on the resultant inequality of

148 Lemma 3. Furthermore, note that

$$\frac{1}{n^2 M}\big\|\widehat{\mathcal{B}}\times_1\hat{\mathbf{Z}}\times_2\hat{\mathbf{Z}} - \mathcal{B}^*\times_1\mathbf{Z}^*\times_2\mathbf{Z}^*\big\|_F^2 \geq \frac{1}{n^2 M}\sum_{k=1}^{K}\sum_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\big\|(\widehat{\mathcal{M}}\times_1\hat{\mathbf{Z}})_{i,.,.} - \mathcal{M}_{k,.,.}^*\big\|_F^2$$

$$\geq \frac{1}{n^2 M}\sum_{k=1}^{K}\sum_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\left(\frac{\big\|\mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^* - \mathcal{M}_{k,.,.}^*\big\|_F^2}{2} - \big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2\right)$$

$$\geq \sum_{k=1}^{K}\frac{|N_k^*\setminus\hat{N}_{\pi^*(k)}|}{n^2 M}\min_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\left(\frac{1}{2}\big\|\mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^* - \mathcal{M}_{k,.,.}^*\big\|_F^2 - \big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2\right)$$

$$\geq \sum_{k=1}^{K}\frac{|N_k^*\setminus\hat{N}_{\pi^*(k)}|}{n^2 M}\left(\frac{1}{2}n_{\min}MK\gamma_n^2 - \max_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2\right).$$

149 Here we use the fact that $\min_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\frac{1}{nM}\big\|\mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^* - \mathcal{M}_{k,.,.}^*\big\|_F^2 \geq \frac{n_{\min}K}{n}\gamma_n^2$ according to

150 Assumption A. Consequently,

$$P\Big(\frac{1}{n}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{c_\xi^2 n\epsilon_n}{n_{\min}K\gamma_n^2 s_n}|F\Big)$$

$$\geq P\Big(\frac{n_{\min}K\gamma_n^2}{n^2}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq$$

$$\sum_{k=1}^{K}\frac{|N_k^*\setminus\hat{N}_{\pi^*(k)}|}{n^2 M}\Big(\frac{1}{2}n_{\min}MK\gamma_n^2 - \max_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2\Big)|F\Big)$$

$$\geq P\Big(\bigcap_{k=1}^{K}\big(\big\{\frac{n_{\min}K\gamma_n^2}{n^2} \leq \frac{n_{\min}K\gamma_n^2}{2n^2} - \max_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\frac{1}{n^2 M}\big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2\big\}\bigcap F\big)\Big)$$

$$\geq P\Big(\bigcap_{k=1}^{K}\big(\big\{\max_{i\in N_k^*\setminus\hat{N}_{\pi^*(k)}}\frac{1}{nM}\big\|\widehat{\mathcal{M}}_{\hat{\psi}_i,.,.} - \mathcal{M}_{(\pi^*)^{-1}(\hat{\psi}_i),.,.}^*\big\|_F^2 = o(\frac{n_{\min}K\gamma_n^2}{n})\big\}\bigcap F\big)\Big) = 1,$$

where the last equality is suggested by Lemma 4. Finally, by the definition of conditional probability,

$$P\Big(\frac{1}{n}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{c_\xi^2 n\epsilon_n}{n_{\min}K\gamma_n^2 s_n}\Big) = P\Big(\frac{1}{n}\sum_{k=1}^{K}|N_k^* \setminus \hat{N}_{\pi^*(k)}| \leq \frac{c_\xi^2 n\epsilon_n}{n_{\min}K\gamma_n^2 s_n}\Big|F\Big) \cdot P(F)$$

$$\geq 1 - 2\exp\Big(-\frac{\varphi(n,M)\epsilon_n}{156\frac{\xi}{1-\xi}+28\log 2}\Big) - n^{-2},$$

and thus the desired consistency result follows immediately. □

## References

[1] Pengsheng Ji and Jiashun Jin. Coauthorship and citation networks for statisticians. *The Annals of Applied Statistics*, 10(4):1779–1812, 2016.

[2] Thierry Klein and Emmanuel Rio. Concentration around the mean for maxima of empirical processes. *Ann. Prob.*, 33(3):1060–1077, 2005.

[3] Vladimir Koltchinskii. *Oracle Inequalities in Empirical Risk Minimization and Sparse Recovery Problems: Ecole d'Eté de Probabilités de Saint-Flour XXXVIII-2008*, volume 2033. Springer Science & Business Media, 2011.

[4] Karl Rohe, Tai Qin, and Bin Yu. Co-clustering directed graphs to discover asymmetries and directional communities. *Proceedings of the National Academy of Sciences*, 113(45):12679–12684, 2016.