

Non-convex Statistical Optimization for Sparse Tensor Graphical Model (Supplementary Material)

In this supplementary note, we provide the proofs of our main theorems in §A, prove the key lemmas in §B, list the auxiliary lemmas in §C, and illustrate additional simulation results in §D.

A Proof of main theorems

Proof of Theorem 3.1: To ease the presentation, we show that Theorem 3.1 holds when $K = 3$. The proof can be easily generalized to the case with $K > 3$.

We first simplify the population log-likelihood function. Note that when $\mathcal{T} \sim \text{TN}(\mathbf{0}; \Sigma_1^*, \Sigma_2^*, \Sigma_3^*)$, Lemma 1 of [9] implies that $\text{vec}(\mathcal{T}) \sim \text{N}(\text{vec}(\mathbf{0}); \Sigma_3^* \otimes \Sigma_2^* \otimes \Sigma_1^*)$. Therefore,

$$\begin{aligned} \mathbb{E}\{\text{tr}[\text{vec}(\mathcal{T})\text{vec}(\mathcal{T})^\top (\Omega_3 \otimes \Omega_2 \otimes \Omega_1)]\} &= \text{tr}[(\Sigma_3^* \otimes \Sigma_2^* \otimes \Sigma_1^*)(\Omega_3 \otimes \Omega_2 \otimes \Omega_1)] \\ &= \text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2) \text{tr}(\Sigma_1^* \Omega_1), \end{aligned}$$

where the second equality is due to the properties of kronecker product that $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$ and $\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A})\text{tr}(\mathbf{B})$. Therefore, the population log-likelihood function can be rewritten as

$$q(\Omega_1, \Omega_2, \Omega_3) = \frac{\text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2) \text{tr}(\Sigma_1^* \Omega_1)}{m_1 m_2 m_3} - \frac{1}{m_1} \log |\Omega_1| - \frac{1}{m_2} \log |\Omega_2| - \frac{1}{m_3} \log |\Omega_3|.$$

Taking derivative of $q(\Omega_1, \Omega_2, \Omega_3)$ with respect to Ω_1 while fixing Ω_2 and Ω_3 , we have

$$\nabla_1 q(\Omega_1, \Omega_2, \Omega_3) = \frac{\text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2)}{m_1 m_2 m_3} \Sigma_1^* - \frac{1}{m_1} \Omega_1^{-1}.$$

Setting it as zero leads to $\Omega_1 = m_2 m_3 [\text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2)]^{-1} \Omega_1^*$. This is indeed a minimizer of $q(\Omega_1, \Omega_2, \Omega_3)$ when fixing Ω_2 and Ω_3 , since the second derivative $\nabla_1^2 q(\Omega_1, \Omega_2, \Omega_3) = m_1^{-1} \Omega_1^{-1} \otimes \Omega_1^{-1}$ is positive definite. Therefore, we have

$$M_1(\Omega_2, \Omega_3) = \frac{m_2 m_3}{\text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2)} \Omega_1^*. \quad (\text{A.1})$$

Therefore, $M_1(\Omega_2, \Omega_3)$ equals to the true parameter Ω_1^* up to a constant. The computations of $M_2(\Omega_1, \Omega_3)$ and $M_3(\Omega_1, \Omega_2)$ follow from the same argument. This ends the proof of Theorem 3.1. ■

Proof of Theorem 3.4: To ease the presentation, we show that (3.5) holds when $K = 3$. The proof of the case when $K > 3$ is similar. We focus on the proof of the statistical error for the sample minimization function $\widehat{M}_1(\Omega_2, \Omega_3)$.

By definition, $\widehat{M}_1(\Omega_2, \Omega_3) = \arg\min_{\Omega_1} q_n(\Omega_1, \Omega_2, \Omega_3) = \arg\min_{\Omega_1} L(\Omega_1)$, where

$$L(\Omega_1) = \frac{1}{m_1} \text{tr}(\mathbf{S}_1 \Omega_1) - \frac{1}{m_1} \log |\Omega_1| + \lambda_1 \|\Omega_1\|_{1, \text{off}},$$

with the sample covariance matrix

$$\mathbf{S}_1 = \frac{1}{m_2 m_3 n} \sum_{i=1}^n \mathbf{V}_i \mathbf{V}_i^\top \text{ with } \mathbf{V}_i = [\mathcal{T}_i \times \{\mathbb{1}_{m_1}, \Omega_2^{1/2}, \Omega_3^{1/2}\}]_{(1)}.$$

For some constant $H > 0$, we define the set of convergence

$$\mathbb{A} := \left\{ \Delta \in \mathbb{R}^{m_1 \times m_1} : \Delta = \Delta^\top, \|\Delta\|_F = H \sqrt{\frac{(m_1 + s_1) \log m_1}{n m_2 m_3}} \right\}.$$

The key idea is to show that

$$\inf_{\Delta \in \mathbb{A}} \{L(M_1(\Omega_2, \Omega_3) + \Delta) - L(M_1(\Omega_2, \Omega_3))\} > 0, \quad (\text{A.2})$$

with high probability. To understand it, note that the function $L(M_1(\Omega_2, \Omega_3) + \Delta) - L(M_1(\Omega_2, \Omega_3))$ is convex in Δ . In addition, since $\widehat{M}_1(\Omega_2, \Omega_3)$ minimizes $L(\Omega_1)$, we have

$$L(\widehat{M}_1(\Omega_2, \Omega_3)) - L(M_1(\Omega_2, \Omega_3)) \leq L(M_1(\Omega_2, \Omega_3)) - L(M_1(\Omega_2, \Omega_3)) = 0.$$

If we can show (A.2), then the minimizer $\widehat{\Delta} = \widehat{M}_1(\Omega_2, \Omega_3) - M_1(\Omega_2, \Omega_3)$ must be within the interior of the ball defined by \mathbb{A} , and hence $\|\widehat{\Delta}\|_F \leq H\sqrt{(m_1 + s_1) \log m_1 / (nm_2m_3)}$. Similar technique is applied in vector-valued graphical model literature [25].

To show (A.2), we first decompose $L(M_1(\Omega_2, \Omega_3) + \Delta) - L(M_1(\Omega_2, \Omega_3)) = I_1 + I_2 + I_3$, where

$$\begin{aligned} I_1 &:= \frac{1}{m_1} \text{tr}(\Delta \mathbf{S}_1) - \frac{1}{m_1} \{ \log |M_1(\Omega_2, \Omega_3) + \Delta| - \log |M_1(\Omega_2, \Omega_3)| \}, \\ I_2 &:= \lambda_1 \{ \| [M_1(\Omega_2, \Omega_3) + \Delta]_{\mathbb{S}_1} \|_1 - \| [M_1(\Omega_2, \Omega_3)]_{\mathbb{S}_1} \|_1 \}, \\ I_3 &:= \lambda_1 \{ \| [M_1(\Omega_2, \Omega_3) + \Delta]_{\mathbb{S}_1^c} \|_1 - \| [M_1(\Omega_2, \Omega_3)]_{\mathbb{S}_1^c} \|_1 \}. \end{aligned}$$

It is sufficient to show $I_1 + I_2 + I_3 > 0$ with high probability. To simplify the term I_1 , we employ the Taylor expansion of $f(t) = \log |M_1(\Omega_2, \Omega_3) + t\Delta|$ at $t = 0$ to obtain

$$\begin{aligned} & \log |M_1(\Omega_2, \Omega_3) + \Delta| - \log |M_1(\Omega_2, \Omega_3)| \\ &= \text{tr} \{ [M_1(\Omega_2, \Omega_3)]^{-1} \Delta \} - [\text{vec}(\Delta)]^\top \left[\int_0^1 (1 - \nu) \mathbf{M}_\nu^{-1} \otimes \mathbf{M}_\nu^{-1} d\nu \right] \text{vec}(\Delta), \end{aligned}$$

where $\mathbf{M}_\nu := M_1(\Omega_2, \Omega_3) + \nu\Delta \in \mathbb{R}^{m_1 \times m_1}$. This leads to

$$I_1 = \underbrace{\frac{1}{m_1} \text{tr}(\{ \mathbf{S}_1 - [M_1(\Omega_2, \Omega_3)]^{-1} \} \Delta)}_{I_{11}} + \underbrace{\frac{1}{m_1} [\text{vec}(\Delta)]^\top \left[\int_0^1 (1 - \nu) \mathbf{M}_\nu^{-1} \otimes \mathbf{M}_\nu^{-1} d\nu \right] \text{vec}(\Delta)}_{I_{12}}.$$

For two symmetric matrices \mathbf{A}, \mathbf{B} , it is easy to see that $|\text{tr}(\mathbf{A}\mathbf{B})| = |\sum_{i,j} \mathbf{A}_{i,j} \mathbf{B}_{i,j}|$. Based on this observation, we decompose I_{11} into two parts: those in the set $\mathbb{S}_1 = \{(i, j) : [\Omega_1^*]_{i,j} \neq 0\}$ and those not in \mathbb{S}_1 . That is, $|I_{11}| \leq I_{111} + I_{112}$, where

$$\begin{aligned} I_{111} &:= \frac{1}{m_1} \left| \sum_{(i,j) \in \mathbb{S}_1} \{ \mathbf{S}_1 - [M_1(\Omega_2, \Omega_3)]^{-1} \}_{i,j} \Delta_{i,j} \right|, \\ I_{112} &:= \frac{1}{m_1} \left| \sum_{(i,j) \notin \mathbb{S}_1} \{ \mathbf{S}_1 - [M_1(\Omega_2, \Omega_3)]^{-1} \}_{i,j} \Delta_{i,j} \right|. \end{aligned}$$

Bound I_{111} : For two matrices \mathbf{A}, \mathbf{B} and a set \mathbb{S} , we have

$$\left| \sum_{(i,j) \in \mathbb{S}} \mathbf{A}_{i,j} \mathbf{B}_{i,j} \right| \leq \max_{i,j} |\mathbf{A}_{i,j}| \left| \sum_{(i,j) \in \mathbb{S}} \mathbf{B}_{i,j} \right| \leq \sqrt{|\mathbb{S}|} \max_{i,j} |\mathbf{A}_{i,j}| \|\mathbf{B}\|_F,$$

where the second inequality is due to the Cauchy-Schwarz inequality and the fact that $\sum_{(i,j) \in \mathbb{S}} \mathbf{B}_{i,j}^2 \leq \|\mathbf{B}\|_F^2$. Therefore, we have

$$\begin{aligned} I_{111} &\leq \frac{\sqrt{s_1 + m_1}}{m_1} \cdot \max_{i,j} \left| \{ \mathbf{S}_1 - [M_1(\Omega_2, \Omega_3)]^{-1} \}_{i,j} \right| \|\Delta\|_F \\ &\leq C \sqrt{\frac{(m_1 + s_1) \log m_1}{nm_1^2 m_2 m_3}} \|\Delta\|_F = \frac{CH \cdot (m_1 + s_1) \log m_1}{nm_1 m_2 m_3}, \end{aligned} \quad (\text{A.3})$$

where (A.3) is from Lemma B.2, the definition of $M_1(\Omega_2, \Omega_3)$ in (A.1), and the fact that $\Delta \in \mathbb{A}$.

Bound I_{12} : For any vector $\mathbf{v} \in \mathbb{R}^p$ and any matrix $\mathbf{A} \in \mathbb{R}^{p \times p}$, the variational form of Rayleigh quotients implies $\lambda_{\min}(\mathbf{A}) = \min_{\|\mathbf{x}\|=1} \mathbf{x}^\top \mathbf{A} \mathbf{x}$ and hence $\lambda_{\min}(\mathbf{A}) \|\mathbf{v}\|^2 \leq \mathbf{v}^\top \mathbf{A} \mathbf{v}$. Setting $\mathbf{v} = \text{vec}(\Delta)$ and $\mathbf{A} = \int_0^1 (1 - \nu) \mathbf{M}_\nu^{-1} \otimes \mathbf{M}_\nu^{-1} d\nu$ leads to

$$I_{12} \geq \frac{1}{m_1} \|\text{vec}(\Delta)\|_2^2 \int_0^1 (1 - \nu) \lambda_{\min}(\mathbf{M}_\nu^{-1} \otimes \mathbf{M}_\nu^{-1}) d\nu.$$

Moreover, by the property of kronecker product, we have

$$\lambda_{\min}(\mathbf{M}_\nu^{-1} \otimes \mathbf{M}_\nu^{-1}) = [\lambda_{\min}(\mathbf{M}_\nu^{-1})]^2 = [\lambda_{\max}(\mathbf{M}_\nu)]^{-2}.$$

In addition, by definition, $\mathbf{M}_\nu = M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3) + \nu \boldsymbol{\Delta}$, and hence we have

$$\lambda_{\max}[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3) + \nu \boldsymbol{\Delta}] \leq \lambda_{\max}[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)] + \lambda_{\max}(\nu \boldsymbol{\Delta}).$$

Therefore, we can bound I_{12} from below, that is,

$$\begin{aligned} I_{12} &\geq \frac{\|\text{vec}(\boldsymbol{\Delta})\|_2^2}{2m_1} \min_{0 \leq \nu \leq 1} [\lambda_{\max}[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)] + \lambda_{\max}(\nu \boldsymbol{\Delta})]^{-2} \\ &\geq \frac{\|\text{vec}(\boldsymbol{\Delta})\|_2^2}{2m_1} [\|M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)\|_2 + \|\boldsymbol{\Delta}\|_2]^{-2}. \end{aligned}$$

On the boundary of \mathbb{A} , it holds that $\|\boldsymbol{\Delta}\|_2 \leq \|\boldsymbol{\Delta}\|_F = o(1)$. Moreover, according to (A.1), we have

$$\|M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)\|_2 = \left| \frac{m_2 m_3}{\text{tr}(\boldsymbol{\Sigma}_3^* \boldsymbol{\Omega}_3) \text{tr}(\boldsymbol{\Sigma}_2^* \boldsymbol{\Omega}_2)} \right| \|\boldsymbol{\Omega}_1^*\|_2 \leq \frac{100}{81} \|\boldsymbol{\Sigma}_1^*\|_2 \leq \frac{1.5}{C_1}, \quad (\text{A.4})$$

where the first inequality is due to

$$\begin{aligned} \text{tr}(\boldsymbol{\Sigma}_3^* \boldsymbol{\Omega}_3) &= \text{tr}[\boldsymbol{\Sigma}_3^*(\boldsymbol{\Omega}_3 - \boldsymbol{\Omega}_3^*) + \mathbf{1}_{m_3}] \geq m_3 - |\text{tr}[\boldsymbol{\Sigma}_3^*(\boldsymbol{\Omega}_3 - \boldsymbol{\Omega}_3^*)]| \\ &\geq m_3 - \|\boldsymbol{\Sigma}_3^*\|_F \|\boldsymbol{\Omega}_3 - \boldsymbol{\Omega}_3^*\|_F \geq m_3(1 - \alpha \|\boldsymbol{\Sigma}_3^*\|_2 / \sqrt{m_3}) \geq 0.9m_3, \end{aligned}$$

for sufficiently large m_3 . Similarly, it holds that $\text{tr}(\boldsymbol{\Sigma}_2^* \boldsymbol{\Omega}_2) \geq 0.9m_2$. The second inequality in (A.4) is due to Condition 3.2. This together with the fact that $\|\text{vec}(\boldsymbol{\Delta})\|_2 = \|\boldsymbol{\Delta}\|_F = o(1) \leq 0.5/C_1$ for sufficiently large n imply that

$$I_{12} \geq \frac{\|\text{vec}(\boldsymbol{\Delta})\|_2^2}{2m_1} \left(\frac{C_1}{2} \right)^2 = \frac{C_1^2 H^2}{8} \cdot \frac{(m_1 + s_1) \log m_1}{nm_1 m_2 m_3}, \quad (\text{A.5})$$

which dominates the term I_{111} for sufficiently large H .

Bound I_2 : To bound I_2 , we apply the triangle inequality and then connect the ℓ_1 matrix norm with its Frobenius norm to obtain the final bound. Specifically, we have

$$|I_2| \leq \lambda_1 \|[\boldsymbol{\Delta}]_{\mathbb{S}_1}\|_1 = \lambda_1 \sum_{(i,j) \in \mathbb{S}_1} |\boldsymbol{\Delta}_{i,j}| \leq \lambda_1 \sqrt{(s_1 + m_1) \sum_{(i,j) \in \mathbb{S}_1} \boldsymbol{\Delta}_{i,j}^2} \leq \lambda_1 \sqrt{s_1 + m_1} \|\boldsymbol{\Delta}\|_F,$$

where the first inequality is from triangle inequality, the second inequality is due to the Cauchy-Schwarz inequality by noting that $s_1 = |\mathbb{S}_1| - m_1$, and the last inequality is due to the definition of Frobenius norm. By Condition 3.3, $\lambda_1 \leq C_2 \sqrt{\log m_1 / (nm_1^2 m_2 m_3)}$. Therefore,

$$|I_2| \leq C_2 H \cdot \frac{(m_1 + s_1) \log m_1}{nm_1 m_2 m_3},$$

which is dominated by I_{12} for sufficiently large H according to (A.5).

Bound $I_3 - |I_{112}|$: We show $I_3 - |I_{112}| > 0$. According to (A.1), we have that $M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)$ equals $\boldsymbol{\Omega}_1^*$ up to a non-zero coefficient. Therefore, for any entry $(i, j) \in \mathbb{S}_1^c$, we have $[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)]_{i,j} = 0$. This implies that

$$I_3 = \lambda_1 \sum_{(i,j) \in \mathbb{S}_1^c} \{ |[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)]_{i,j} + \boldsymbol{\Delta}_{i,j}| - |[M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)]_{i,j}| \} = \lambda_1 \sum_{(i,j) \in \mathbb{S}_1^c} |\boldsymbol{\Delta}_{i,j}|.$$

This together with the expression of I_{112} and the bound in Lemma B.2 leads to

$$\begin{aligned} I_3 - I_{112} &= \sum_{(i,j) \in \mathbb{S}_1^c} \left\{ \lambda_1 - m_1^{-1} \{ \mathbf{S}_1 - [M_1(\boldsymbol{\Omega}_2, \boldsymbol{\Omega}_3)]^{-1} \}_{i,j} \right\} |\boldsymbol{\Delta}_{i,j}| \\ &\geq \left(\lambda_1 - C \sqrt{\frac{\log m_1}{nm_1^2 m_2 m_3}} \right) \sum_{(i,j) \in \mathbb{S}_1^c} |\boldsymbol{\Delta}_{i,j}| > 0, \end{aligned}$$

as long as $1/C_2 > C$ for some constant C , which is valid for sufficient small C_2 in Condition 3.3.

Combining all these bounds together, we have, for any $\Delta \in \mathbb{A}$, with high probability,

$$L(M_1(\Omega_2, \Omega_3) + \Delta) - L(M_1(\Omega_2, \Omega_3)) \geq I_{12} - I_{111} - |I_2| + I_3 - I_{112} > 0,$$

which ends the proof Theorem 3.4. \blacksquare

Proof of Theorem 3.5: We show it by connecting the one-step convergence result in Theorem 3.1 and the statistical error result in Theorem 3.4. We show the case when $K = 3$. The proof of the $K > 3$ case is similar. We focus on the proof of the estimation error $\|\hat{\Omega}_1 - \Omega_1^*\|_F$.

To ease the presentation, in the following derivation we remove the superscript in the initializations $\Omega_2^{(0)}$ and $\Omega_3^{(0)}$ and use Ω_2 and Ω_3 instead. According to the procedure in Algorithm 1, we have

$$\begin{aligned} \|\hat{\Omega}_1 - \Omega_1^*\|_F &= \left\| \frac{\widehat{M}_1(\Omega_2, \Omega_3)}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} - \frac{\widehat{M}_1(\Omega_2, \Omega_3)}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} \right\|_F \\ &\leq \left\| \frac{\widehat{M}_1(\Omega_2, \Omega_3)}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} - \frac{M_1(\Omega_2, \Omega_3)}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} \right\|_F + \left\| \frac{M_1(\Omega_2, \Omega_3)}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} - \frac{M_1(\Omega_2, \Omega_3)}{\|M_1(\Omega_2, \Omega_3)\|_F} \right\|_F \\ &\leq \frac{2}{\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F} \|\widehat{M}_1(\Omega_2, \Omega_3) - M_1(\Omega_2, \Omega_3)\|_F, \end{aligned}$$

where the last inequality is due to the triangle inequality $\|a\| - \|b\| \leq \|a - b\|$ and the summation of two parts. We next bound $\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F$. By triangle inequality,

$$\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F \geq \|M_1(\Omega_2, \Omega_3)\|_F - \|M_1(\Omega_2, \Omega_3) - \widehat{M}_1(\Omega_2, \Omega_3)\|_F \geq 2^{-1} \|M_1(\Omega_2, \Omega_3)\|_F,$$

since $\|M_1(\Omega_2, \Omega_3) - \widehat{M}_1(\Omega_2, \Omega_3)\|_F = o_P(1)$ as shown in Theorem 3.4. Moreover, by the Cauchy-Schwarz inequality, we have

$$\text{tr}(\Sigma_2^* \Omega_2) \leq \|\Sigma_2^*\|_F \|\Omega_2\|_F \leq m_2 \|\Sigma_2^*\|_2 \|\Omega_2\|_2 \leq 2m_2/C_1,$$

due to Condition 3.2 and the fact that $\Omega_2 \in \mathbb{B}(\Omega_2^*)$. Similarly, we have $\text{tr}(\Sigma_3^* \Omega_3) \leq 2m_3/C_1$. This together with the expression of $M_1(\Omega_2, \Omega_3)$ in (A.1) imply that $\|\widehat{M}_1(\Omega_2, \Omega_3)\|_F \geq C_1^2/4$ and hence

$$\|\hat{\Omega}_1 - \Omega_1^*\|_F \leq \frac{8}{C_1^2} \|\widehat{M}_1(\Omega_2, \Omega_3) - M_1(\Omega_2, \Omega_3)\|_F = O_P\left(\sqrt{\frac{m_1(m_1 + s_1) \log m_1}{nm_1 m_2 m_3}}\right),$$

according to Theorem 3.4. This ends the proof Theorem 3.5. \blacksquare

Proof of Theorem 3.9: We prove it by transferring the optimization problem to an equivalent primal-dual problem and then applying the convergence results of [27] to obtain the desirable rate of convergence.

Given the sample covariance matrix $\widehat{\mathbf{S}}_k$ defined in Lemma B.3, according to (2.3), for each $k = 1, \dots, K$, the optimization problem has a unique solution $\widehat{\Omega}_k$ which satisfies the following Karush-Kuhn-Tucker (KKT) conditions

$$\widehat{\mathbf{S}}_k - \widehat{\Omega}_k + m_k \lambda_k \widehat{\mathbf{Z}}_k = 0, \quad (\text{A.6})$$

where $\widehat{\mathbf{Z}}_k \in \mathbb{R}^{m_k \times m_k}$ belongs to the sub-differential of $\|\Omega_k\|_{1, \text{off}}$ evaluated at $\widehat{\Omega}_k$, that is,

$$[\widehat{\mathbf{Z}}_k]_{i,j} := \begin{cases} 0, & \text{if } i = j \\ \text{sign}([\widehat{\Omega}_k]_{i,j}) & \text{if } i \neq j \text{ and } [\widehat{\Omega}_k]_{i,j} \neq 0 \\ \in [-1, +1] & \text{if } i \neq j \text{ and } [\widehat{\Omega}_k]_{i,j} = 0. \end{cases}$$

Following [27], we construct the primary-dual witness solution $(\widetilde{\Omega}_k, \widetilde{\mathbf{Z}}_k)$ such that

$$\widetilde{\Omega}_k := \underset{\Omega_k \succ 0, \Omega_k = \Omega_k^\top, [\Omega_k]_{\mathcal{S}_k^c} = 0}{\text{argmin}} \left\{ \text{tr}(\widehat{\mathbf{S}}_k \Omega_k) - \log |\Omega_k| + m_k \lambda_k \|\Omega_k\|_{1, \text{off}} \right\},$$

where the set \mathbb{S}_k refers to the set of true non-zero edges of Ω_k^* . Therefore, by construction, the support of the dual estimator $\hat{\Omega}_k$ is a subset of the true support, i.e., $\text{supp}(\hat{\Omega}_k) \subseteq \text{supp}(\Omega_k^*)$. We then construct $\tilde{\mathbf{Z}}_k$ as the sub-differential $\hat{\mathbf{Z}}_k$ and then for each $(i, j) \in \mathbb{S}_k^c$, we replace $[\tilde{\mathbf{Z}}_k]_{i,j}$ with $([\tilde{\Omega}_k^{-1}]_{i,j} - [\hat{\mathbf{S}}_k]_{i,j})/(m_k \lambda_k)$ to ensure that $(\tilde{\Omega}_k, \tilde{\mathbf{Z}}_k)$ satisfies the optimality condition (A.6).

Denote $\Delta := \tilde{\Omega}_k - \Omega_k^*$ and $R(\Delta) := \tilde{\Omega}_k^{-1} - \Omega_k^{*-1} + \Omega_k^{*-1} \Delta \tilde{\Omega}_k^{-1}$. According to Lemma 4 of [27], in order to show the strict dual feasibility $\tilde{\Omega}_k = \hat{\Omega}_k$, it is sufficient to prove

$$\max\{\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty, \|R(\Delta)\|_\infty\} \leq \frac{\alpha_k m_k \lambda_k}{8},$$

with α_k defined in Condition 3.7. As assumed in Condition 3.3, the tuning parameter satisfies $1/C_2 \sqrt{\log m_k / (nm m_k)} \leq \lambda_k \leq C_2 \sqrt{\log m_k / (nm m_k)}$ for some constant $C_2 > 0$ and hence $\alpha_k m_k \lambda_k / 8 \geq C_3 \sqrt{m_k \log m_k / (nm)}$ for some constant $C_3 > 0$.

In addition, according to Lemma B.3, we have

$$\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty = O_P \left(\max_{j=1, \dots, K} \sqrt{\frac{(m_j + s_j) \log m_j}{nm}} \right).$$

Under the assumption that $s_j = O(m_j)$ for $j = 1, \dots, K$ and $m_1 \asymp m_2 \asymp \dots \asymp m_K$, we have

$$\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty = O_P \left(\sqrt{\frac{m_k \log m_k}{nm}} \right).$$

Therefore, there exists a sufficiently small constant C_2 such that $\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty \leq \alpha_k m_k \lambda_k / 8$.

Moreover, according to Lemma 5 of [27], $\|R(\Delta)\|_\infty \leq 1.5 d_k \|\Delta\|_\infty^2 \kappa_{\Sigma_k^*}^3$ as long as $\|\Delta\|_\infty \leq (3 \kappa_{\Sigma_k^*} d_k)^{-1}$. According to Lemma 6 of [27], if we can show

$$r := 2 \kappa_{\Gamma_k^*} (\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty + m_k \lambda_k) \leq \min \left\{ \frac{1}{3 \kappa_{\Sigma_k^*} d_k}, \frac{1}{\kappa_{\Sigma_k^*}^3 \kappa_{\Gamma_k^*} d_k} \right\},$$

then we have $\|\Delta\|_\infty \leq r$. By Condition 3.8, $\kappa_{\Gamma_k^*}$ and $\kappa_{\Sigma_k^*}$ are bounded. Therefore, $\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty + m_k \lambda_k$ is in the same order of $\sqrt{m_k \log m_k / (nm)}$, which is in a smaller order of d_k^{-1} by the assumption of d_k in Condition 3.8. Therefore, we have shown that $\|R(\Delta)\|_\infty \leq m_k \lambda_k$ for a sufficiently small constant C_2 .

Combining above two bounds, we achieve the strict dual feasibility $\tilde{\Omega}_k = \hat{\Omega}_k$. Therefore, we have $\text{supp}(\hat{\Omega}_k) \subseteq \text{supp}(\Omega_k^*)$ and moreover,

$$\|\hat{\Omega}_k - \Omega_k^*\|_\infty = \|\Delta\|_\infty = O_P \left(\sqrt{\frac{m_k \log m_k}{nm}} \right).$$

This ends the proof of Theorem 3.9. ■

B Proof of key lemmas

The first key lemma establishes the rate of convergence of the difference between a sample-based quadratic form and its expectation. This new concentration result is also of independent interest.

Lemma B.1. Assume i.i.d. data $\mathbf{X}, \mathbf{X}_1, \dots, \mathbf{X}_n \in \mathbb{R}^{p \times q}$ follows the matrix-variate normal distribution such that $\text{vec}(\mathbf{X}_i) \sim N(\mathbf{0}; \Psi^* \otimes \Sigma^*)$ with $\Psi^* \in \mathbb{R}^{q \times q}$ and $\Sigma^* \in \mathbb{R}^{p \times p}$. Assume that $0 < C_1 \leq \lambda_{\min}(\Sigma^*) \leq \lambda_{\max}(\Sigma^*) \leq 1/C_1 < \infty$ and $0 < C_2 \leq \lambda_{\min}(\Psi^*) \leq \lambda_{\max}(\Psi^*) \leq 1/C_2 < \infty$ for some positive constants C_1, C_2 . For any symmetric and positive definite matrix $\Omega \in \mathbb{R}^{p \times p}$, we have

$$\max_{i,j} \left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \Omega \mathbf{X}_i - \frac{1}{p} \mathbb{E}(\mathbf{X}^\top \Omega \mathbf{X}) \right\}_{i,j} = O_P \left(\sqrt{\frac{\log q}{np}} \right).$$

Proof of Lemma B.1: Consider a random matrix \mathbf{X} following the matrix normal distribution such that $\text{vec}(\mathbf{X}) \sim \mathcal{N}(\mathbf{0}; \Psi^* \otimes \Sigma^*)$. Let $\Lambda^* = \Psi^{*-1}$ and $\Omega^* = \Sigma^{*-1}$. Let $\mathbf{Y} := (\Omega^*)^{1/2} \mathbf{X} (\Lambda^*)^{1/2}$. According to the properties of matrix normal distribution [30], \mathbf{Y} follows a matrix normal distribution such that $\text{vec}(\mathbf{Y}) \sim \mathcal{N}(\mathbf{0}; \mathbb{1}_q \otimes \mathbb{1}_p)$, that is, all the entries of \mathbf{Y} are i.i.d. standard Gaussian random variables. Next we rewrite the term $\mathbf{X}^\top \Omega \mathbf{X}$ by \mathbf{Y} and then simplify it. Simple algebra implies that

$$\mathbf{X}^\top \Omega \mathbf{X} = (\Lambda^*)^{-1/2} \mathbf{Y}^\top (\Omega^*)^{-1/2} \Omega (\Omega^*)^{-1/2} \mathbf{Y} (\Lambda^*)^{-1/2}.$$

When Ω is symmetric and positive definite, the matrix $\mathbf{M} := (\Omega^*)^{-1/2} \Omega (\Omega^*)^{-1/2} \in \mathbb{R}^{p \times p}$ is also symmetric and positive definite with Cholesky decomposition $\mathbf{U}^\top \mathbf{U}$, where $\mathbf{U} \in \mathbb{R}^{p \times p}$. Therefore,

$$\mathbf{X}^\top \Omega \mathbf{X} = (\Lambda^*)^{-1/2} \mathbf{Y}^\top \mathbf{U}^\top \mathbf{U} \mathbf{Y} (\Lambda^*)^{-1/2}.$$

Moreover, denote the column of the matrix $(\Lambda^*)^{-1/2}$ as $(\Lambda^*)_{(j)}^{-1/2}$ and denote its row as $(\Lambda^*)_i^{-1/2}$ for $i, j = 1, \dots, q$. Define the standard basis $\mathbf{e}_i \in \mathbb{R}^q$ as the vector with 1 in its i -th entry and 0 in all the rest entries. The (s, t) -th entry of matrix $\mathbf{X}^\top \Omega \mathbf{X}$ can be written as

$$\{\mathbf{X}^\top \Omega \mathbf{X}\}_{s,t} = \mathbf{e}_s^\top \mathbf{X}^\top \Omega \mathbf{X} \mathbf{e}_t = (\Lambda^*)_s^{-1/2} \mathbf{Y}^\top \mathbf{U}^\top \mathbf{U} \mathbf{Y} (\Lambda^*)_{(t)}^{-1/2}.$$

For the sample matrices $\mathbf{X}_1, \dots, \mathbf{X}_n$, we apply similar transformation that $\mathbf{Y}_i = (\Omega^*)^{1/2} \mathbf{X}_i (\Lambda^*)^{1/2}$. We apply the above derivation to the sample-based quadratic term $\mathbf{X}_i^\top \Omega \mathbf{X}_i$. Let $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_n) \in \mathbb{R}^{p \times n}$ with $\mathbf{a}_i = \mathbf{U} \mathbf{Y}_i (\Lambda^*)_s^{-1/2} \in \mathbb{R}^p$ and $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_n) \in \mathbb{R}^{p \times n}$ with $\mathbf{b}_i = \mathbf{U} \mathbf{Y}_i (\Lambda^*)_t^{-1/2} \in \mathbb{R}^p$. Then we have

$$\begin{aligned} \left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \Omega \mathbf{X}_i \right\}_{s,t} &= \frac{1}{np} \sum_{i=1}^n \mathbf{a}_i^\top \mathbf{b}_i = \frac{1}{np} \sum_{i=1}^n \sum_{j=1}^p \mathbf{A}_{i,j} \mathbf{B}_{i,j} \\ &= \frac{1}{4np} \sum_{i=1}^n \sum_{j=1}^p \{(\mathbf{A}_{i,j} + \mathbf{B}_{i,j})^2 - (\mathbf{A}_{i,j} - \mathbf{B}_{i,j})^2\} \\ &= \frac{1}{4np} \{ \|\text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B})\|_2^2 + \|\text{vec}(\mathbf{A}) - \text{vec}(\mathbf{B})\|_2^2 \}. \quad (\text{B.1}) \end{aligned}$$

Next we derive the explicit form of $\text{vec}(\mathbf{A})$ and $\text{vec}(\mathbf{B})$ in (B.1). Remind that $(\Lambda^*)_s^{-1/2}$ is a vector of length q . By the property of matrix products, we can rewrite $\mathbf{a}_i = [(\Lambda^*)_s^{-1/2} \otimes \mathbf{U}] \text{vec}(\mathbf{Y}_i)$, where \otimes is the Kronecker product. Therefore, we have

$$\begin{aligned} \text{vec}(\mathbf{A}) &= [\mathbb{1}_n \otimes (\Lambda^*)_s^{-1/2} \otimes \mathbf{U}] \mathbf{t} := \mathbf{Q}_1 \mathbf{t}, \\ \text{vec}(\mathbf{B}) &= [\mathbb{1}_n \otimes (\Lambda^*)_t^{-1/2} \otimes \mathbf{U}] \mathbf{t} := \mathbf{Q}_2 \mathbf{t}, \end{aligned}$$

where $\mathbf{t} = \{\text{vec}(\mathbf{Y}_1)^\top, \dots, \text{vec}(\mathbf{Y}_n)^\top\}^\top \in \mathbb{R}^{npq}$ is a vector with npq i.i.d. standard normal entries. Here $\mathbf{Q}_1 := \mathbb{1}_n \otimes (\Lambda^*)_s^{-1/2} \otimes \mathbf{U}$ and $\mathbf{Q}_2 := \mathbb{1}_n \otimes (\Lambda^*)_t^{-1/2} \otimes \mathbf{U}$ with $\mathbf{Q}_1, \mathbf{Q}_2 \in \mathbb{R}^{np \times npq}$. By the property of multivariate normal distribution, we have

$$\begin{aligned} \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B}) &\sim \mathcal{N}(\mathbf{0}; (\mathbf{Q}_1 + \mathbf{Q}_2)(\mathbf{Q}_1 + \mathbf{Q}_2)^\top) := \mathcal{N}(\mathbf{0}; \mathbf{H}_1), \\ \text{vec}(\mathbf{A}) - \text{vec}(\mathbf{B}) &\sim \mathcal{N}(\mathbf{0}; (\mathbf{Q}_1 - \mathbf{Q}_2)(\mathbf{Q}_1 - \mathbf{Q}_2)^\top) := \mathcal{N}(\mathbf{0}; \mathbf{H}_2). \end{aligned}$$

Next, we bound the spectral norm of two matrices \mathbf{H}_1 and \mathbf{H}_2 . By the property of matrix norm and the fact that one matrix and its transpose matrix have the same spectral norm, we have

$$\|\mathbf{H}_1\|_2 \leq \|\mathbf{Q}_1 \mathbf{Q}_1^\top\|_2 + 2\|\mathbf{Q}_1 \mathbf{Q}_2^\top\|_2 + \|\mathbf{Q}_2 \mathbf{Q}_2^\top\|_2,$$

then we bound each of these three terms individually. According to the definition of \mathbf{Q}_1 and the property of matrix Kronecker products, we have

$$\begin{aligned} \mathbf{Q}_1 \mathbf{Q}_1^\top &= [\mathbb{1}_n \otimes (\Lambda^*)_s^{-1/2} \otimes \mathbf{U}] [\mathbb{1}_n \otimes (\Lambda^*)_s^{-1/2} \otimes \mathbf{U}]^\top \\ &= \mathbb{1}_n \otimes (\Lambda^*)_s^{-1/2} [(\Lambda^*)_s^{-1/2}]^\top \otimes \mathbf{M}, \end{aligned}$$

where the last equality is due to the fact that $(\mathbf{C}_1 \otimes \mathbf{C}_2)^\top = \mathbf{C}_1^\top \otimes \mathbf{C}_2^\top$ and $(\mathbf{C}_1 \otimes \mathbf{C}_2)(\mathbf{C}_3 \otimes \mathbf{C}_4) = (\mathbf{C}_1 \mathbf{C}_3) \otimes (\mathbf{C}_2 \mathbf{C}_4)$ for any matrices $\mathbf{C}_1, \dots, \mathbf{C}_4$ such that the matrix multiplications $\mathbf{C}_1 \mathbf{C}_3$ and

$\mathbf{C}_2\mathbf{C}_4$ are valid. Moreover, we also use the Cholesky decomposition of \mathbf{M} , i.e., $\mathbf{M} = \mathbf{U}^\top \mathbf{U}$. Remind that $(\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_s^{-1/2}]^\top \in \mathbb{R}$, therefore, the spectral norm $\mathbf{Q}_1\mathbf{Q}_1^\top$ can be written as

$$\begin{aligned}\|\mathbf{Q}_1\mathbf{Q}_1^\top\|_2 &= |(\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_s^{-1/2}]^\top| \cdot \|\mathbf{1}_n\|_2 \|\mathbf{M}\|_2 \\ &\leq \|\Psi^*\|_2 \|\mathbf{M}\|_2 \leq (1 + \alpha/C_1)/C_2.\end{aligned}$$

Here the first inequality is because $\|\mathbf{1}_n\|_2 = 1$ and

$$\begin{aligned}|(\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_s^{-1/2}]^\top| &= \|[(\mathbf{\Lambda}^*)_s^{-1/2}]^\top (\mathbf{\Lambda}^*)_s^{-1/2}\|_2 \leq \max_j \|[(\Psi^*)_j^{1/2}]^\top (\Psi^*)_j^{1/2}\|_2 \\ &\leq \left\| \sum_{j=1}^q [(\Psi^*)_j^{1/2}]^\top (\Psi^*)_j^{1/2} \right\|_2 = \|\Psi^*\|_2,\end{aligned}$$

and the second inequality is because $\|\Psi^*\|_2 \leq 1/C_2$ and

$$\begin{aligned}\|\mathbf{M}\|_2 &= \left\| (\mathbf{\Omega}^*)^{-1/2} \mathbf{\Omega} (\mathbf{\Omega}^*)^{-1/2} \right\|_2 = \|(\mathbf{\Omega}^*)^{-1/2} (\mathbf{\Omega} - \mathbf{\Omega}^*) (\mathbf{\Omega}^*)^{-1/2} + \mathbf{1}_p\|_2 \\ &\leq \|(\mathbf{\Omega}^*)^{-1/2}\|_2^2 \|\mathbf{\Omega} - \mathbf{\Omega}^*\|_2 + 1 \leq \|\Sigma^*\|_2 \|\mathbf{\Omega} - \mathbf{\Omega}^*\|_F + 1 \leq 1 + \alpha/C_1.\end{aligned}$$

Similarly, we have $\|\mathbf{Q}_2\mathbf{Q}_2^\top\|_2 \leq (1 + \alpha/C_1)/C_2$. For $\|\mathbf{Q}_1\mathbf{Q}_2^\top\|_2$, similar arguments imply that

$$\mathbf{Q}_1\mathbf{Q}_2^\top = \mathbf{1}_n \otimes (\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_t^{-1/2}]^\top \otimes \mathbf{M},$$

and hence its spectral norm is bounded as

$$\begin{aligned}\|\mathbf{Q}_1\mathbf{Q}_2^\top\|_2 &= |(\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_t^{-1/2}]^\top| \cdot \|\mathbf{1}_n\|_2 \|\mathbf{M}\|_2 \\ &\leq \|\Psi^*\|_2 \|\mathbf{M}\|_2 \leq (1 + \alpha/C_1)/C_2,\end{aligned}$$

where the first inequality is because the above derivation and the Cauchy-Schwarz inequality. Specifically, let $\Psi^* = (\Psi_{i,j}^*)$, we have

$$\begin{aligned}|(\mathbf{\Lambda}^*)_s^{-1/2}[(\mathbf{\Lambda}^*)_t^{-1/2}]^\top| &= \sqrt{(\Psi^*)_s[(\Psi^*)_t]^\top} = \left[\sum_{j=1}^q \Psi_{s,j}^* \Psi_{t,j}^* \right]^{1/2} \\ &\leq \left\{ \left(\sum_{j=1}^q \Psi_{s,j}^{*2} \right) \left(\sum_{j=1}^q \Psi_{t,j}^{*2} \right) \right\}^{1/4} \leq \sqrt{\|\Psi^*\|_2 \|\Psi^*\|_2} \leq C_2^{-1}.\end{aligned}$$

Applying the same techniques to $\|\mathbf{H}_2\|_2$, we have

$$\|\mathbf{H}_1\|_2 \leq 4(1 + \alpha/C_1)/C_2, \quad (\text{B.2})$$

$$\|\mathbf{H}_2\|_2 \leq 4(1 + \alpha/C_1)/C_2. \quad (\text{B.3})$$

Next, we apply Lemma C.3 to bound the (s, t) -th entry of the differential matrix between the sample-based term and its expectation. Denote $\rho_{s,t} := [p^{-1}\mathbb{E}(\mathbf{X}^\top \mathbf{\Omega} \mathbf{X})]_{s,t}$. According to the derivation in (B.1), we have

$$\begin{aligned}&\left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \mathbf{\Omega} \mathbf{X}_i - \frac{1}{p} \mathbb{E}(\mathbf{X}^\top \mathbf{\Omega} \mathbf{X}) \right\}_{s,t} \\ &= \left[\frac{1}{4np} \sum_{i,j} (a_{ij} + b_{ij})^2 - \frac{\Delta_{s,t} + \rho_{s,t}}{2} \right] - \left[\frac{1}{4np} \sum_{i,j} (a_{ij} - b_{ij})^2 - \frac{\Delta_{s,t} - \rho_{s,t}}{2} \right], \quad (\text{B.4})\end{aligned}$$

where $\Delta_{s,t}$ is defined as

$$\Delta_{s,t} := \mathbb{E} \left\{ (4np)^{-1} \sum_{i,j} [(a_{ij} + b_{ij})^2 + (a_{ij} - b_{ij})^2] \right\}.$$

Furthermore, according to the definition of $\rho_{s,t}$ and (B.1), we have $\mathbb{E} \{ (4np)^{-1} \sum_{i=1}^n \sum_{j=1}^p [(a_{ij} + b_{ij})^2 - (a_{ij} - b_{ij})^2] \} = \rho_{s,t}$. Therefore, we have

$$\mathbb{E} \left\{ (4np)^{-1} \sum_{i,j} (a_{ij} + b_{ij})^2 \right\} = \frac{\Delta_{s,t} + \rho_{s,t}}{2}, \quad (\text{B.5})$$

$$\mathbb{E} \left\{ (4np)^{-1} \sum_{i,j} (a_{ij} - b_{ij})^2 \right\} = \frac{\Delta_{s,t} - \rho_{s,t}}{2}. \quad (\text{B.6})$$

Therefore, (B.4) implies that, for any $\delta > 0$,

$$\begin{aligned} & \mathbb{P} \left[\left| \left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \boldsymbol{\Omega} \mathbf{X}_i - \frac{1}{p} \mathbb{E}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}) \right\}_{s,t} \right| \geq \delta \right] \leq \\ & \underbrace{\mathbb{P} \left[\left| \frac{1}{np} \sum_{i,j} (a_{ij} + b_{ij})^2 - 2(\boldsymbol{\Delta}_{s,t} + \rho_{s,t}) \right| > 2\delta \right]}_{I_1} + \underbrace{\mathbb{P} \left[\left| \frac{1}{np} \sum_{i,j} (a_{ij} - b_{ij})^2 - 2(\boldsymbol{\Delta}_{s,t} - \rho_{s,t}) \right| > 2\delta \right]}_{I_2}. \end{aligned}$$

Remind that $\sum_{i=1}^n \sum_{j=1}^p (a_{ij} + b_{ij})^2 = \text{vec}(\mathbf{A}) + \text{vec}(\mathbf{B}) \sim \mathcal{N}(0; \mathbf{H}_1)$ and $\sum_{i=1}^n \sum_{j=1}^p (a_{ij} - b_{ij})^2 = \text{vec}(\mathbf{A}) - \text{vec}(\mathbf{B}) \sim \mathcal{N}(0; \mathbf{H}_2)$. According to (B.5) and (B.6), we apply Lemma C.3 to obtain

$$\begin{aligned} I_1 & \leq 2 \exp \left\{ -\frac{np}{2} \left(\frac{\delta}{2\|\mathbf{H}_1\|_2} - \frac{2}{\sqrt{np}} \right)^2 \right\} + 2 \exp(-np/2), \\ I_2 & \leq 2 \exp \left\{ -\frac{np}{2} \left(\frac{\delta}{2\|\mathbf{H}_2\|_2} - \frac{2}{\sqrt{np}} \right)^2 \right\} + 2 \exp(-np/2). \end{aligned}$$

Finally, in order to derive the convergence rate of the maximal difference over all index (s, t) , we employ the max sum inequality. That is, for random variables x_1, \dots, x_n , we have $\mathbb{P}(\max_i x_i \geq t) \leq \sum_{i=1}^n \mathbb{P}(x_i \geq t) \leq n \max_i \mathbb{P}(x_i \geq t)$. This together with (B.2) and (B.3) imply that

$$\begin{aligned} & \mathbb{P} \left[\max_{(s,t)} \left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \boldsymbol{\Omega} \mathbf{X}_i - \frac{1}{p} \mathbb{E}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}) \right\}_{s,t} \geq \delta \right] \\ & \leq 4q^2 \exp \left\{ -\frac{np}{2} \left[\frac{\delta C_1 C_2}{8(C_1 + \alpha)} - \frac{2}{\sqrt{np}} \right]^2 \right\} + 4q^2 \exp(-np/2). \end{aligned} \quad (\text{B.7})$$

Let $\delta = 8(C_1 + \alpha)(C_1 C_2)^{-1} [4\sqrt{\log q/(np)} + 3(np)^{-1/2}]$ in (B.7) which satisfies the condition in Lemma C.3 since $\delta > 2(np)^{-1/2}$ when q is sufficiently large. Therefore, we obtain the desirable conclusion that, with high probability,

$$\max_{(s,t)} \left\{ \frac{1}{np} \sum_{i=1}^n \mathbf{X}_i^\top \boldsymbol{\Omega} \mathbf{X}_i - \frac{1}{p} \mathbb{E}(\mathbf{X}^\top \boldsymbol{\Omega} \mathbf{X}) \right\}_{s,t} = O_P \left(\sqrt{\frac{\log q}{np}} \right).$$

This ends the proof of Lemma B.1. ■

Lemma B.2. Assume i.i.d. tensor data $\mathcal{T}, \mathcal{T}_1, \dots, \mathcal{T}_n \in \mathbb{R}^{m_1 \times m_2 \times \dots \times m_K}$ follows the tensor normal distribution $\text{TN}(\mathbf{0}; \boldsymbol{\Sigma}_1^*, \dots, \boldsymbol{\Sigma}_K^*)$. Assume Condition 3.2 holds. For any symmetric and positive definite matrices $\boldsymbol{\Omega}_j \in \mathbb{R}^{m_j \times m_j}, j \neq k$, we have

$$\mathbb{E}[\mathbf{S}_k] = \frac{m_k [\prod_{j \neq k} \text{tr}(\boldsymbol{\Sigma}_j^* \boldsymbol{\Omega}_j)]}{m} \boldsymbol{\Sigma}_k^*,$$

for $\mathbf{S}_k = \frac{m_k}{nm} \sum_{i=1}^n \mathbf{V}_i \mathbf{V}_i^\top$ with $\mathbf{V}_i = [\mathcal{T}_i \times \{\boldsymbol{\Omega}_1^{1/2}, \dots, \boldsymbol{\Omega}_{k-1}^{1/2}, \mathbf{1}_{m_k}, \boldsymbol{\Omega}_{k+1}^{1/2}, \dots, \boldsymbol{\Omega}_K^{1/2}\}]_{(k)}$ and $m = \prod_{k=1}^K m_k$. Moreover, we have

$$\max_{s,t} \left\{ \mathbf{S}_k - \frac{m_k [\prod_{j \neq k} \text{tr}(\boldsymbol{\Sigma}_j^* \boldsymbol{\Omega}_j)]}{m} \boldsymbol{\Sigma}_k^* \right\}_{s,t} = O_P \left(\sqrt{\frac{m_k \log m_k}{nm}} \right). \quad (\text{B.8})$$

Proof of Lemma B.2: The proof follows by carefully examining the distribution of \mathbf{V}_i and then applying Lemma B.1. We only show the case with $K = 3$ and $k = 1$. The extension to a general K follows similarly.

According to the property of mode- k tensor multiplication, we have $\mathbf{V}_i = [\mathcal{T}_i]_{(1)} (\boldsymbol{\Omega}_3^{1/2} \otimes \boldsymbol{\Omega}_2^{1/2})$, and hence

$$\begin{aligned} \mathbf{S}_1 &= \frac{1}{nm_2 m_3} \sum_{i=1}^n [\mathcal{T}_i]_{(1)} (\boldsymbol{\Omega}_3^{1/2} \otimes \boldsymbol{\Omega}_2^{1/2}) (\boldsymbol{\Omega}_3^{1/2} \otimes \boldsymbol{\Omega}_2^{1/2}) [\mathcal{T}_i]_{(1)}^\top \\ &= \frac{1}{nm_2 m_3} \sum_{i=1}^n [\mathcal{T}_i]_{(1)} (\boldsymbol{\Omega}_3 \otimes \boldsymbol{\Omega}_2) [\mathcal{T}_i]_{(1)}^\top. \end{aligned}$$

When tensor $\mathcal{T}_i \sim \text{TN}(\mathbf{0}; \Sigma_1^*, \Sigma_2^*, \Sigma_3^*)$, the property of mode- k tensor multiplication shown in Proposition 2.1 in [31] implies that

$$[\mathcal{T}_i]_{(1)} \in \mathbb{R}^{m_1 \times (m_2 m_3)} \sim \text{MN}(\mathbf{0}; \Sigma_1^*, \Sigma_3^* \otimes \Sigma_2^*),$$

where $\text{MN}(\mathbf{0}; \Sigma_1^*, \Sigma_3^* \otimes \Sigma_2^*)$ is the matrix-variate normal [32] such that the row covariance matrix of $[\mathcal{T}_i]_{(1)}$ is Σ_1^* and the column covariance matrix of $[\mathcal{T}_i]_{(1)}$ is $\Sigma_3^* \otimes \Sigma_2^*$. Therefore, in order to show (B.8), according to Lemma B.1, it is sufficient to show

$$\mathbb{E}[\mathbf{S}_1] = \frac{\text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2)}{m_2 m_3} \Sigma_1^*. \quad (\text{B.9})$$

According to the distribution of $[\mathcal{T}_i]_{(1)}$, we have

$$\mathbf{V}_i \sim \text{MN}\left(\mathbf{0}; \Sigma_1^*, (\Omega_3^{1/2} \otimes \Omega_2^{1/2})(\Sigma_3^* \otimes \Sigma_2^*)(\Omega_3^{1/2} \otimes \Omega_2^{1/2})\right),$$

and hence

$$\mathbf{V}_i^\top \sim \text{MN}\left(\mathbf{0}; (\Omega_3^{1/2} \otimes \Omega_2^{1/2})(\Sigma_3^* \otimes \Sigma_2^*)(\Omega_3^{1/2} \otimes \Omega_2^{1/2}), \Sigma_1^*\right).$$

Therefore, according to Lemma C.1, we have

$$\mathbb{E}[\mathbf{V}_i \mathbf{V}_i^\top] = \Sigma_1^* \text{tr}[(\Omega_3 \otimes \Omega_2)(\Sigma_3^* \otimes \Sigma_2^*)] = \Sigma_1^* \text{tr}(\Sigma_3^* \Omega_3) \text{tr}(\Sigma_2^* \Omega_2),$$

which implies (B.9) according to the definition of \mathbf{S}_1 . Finally, applying Lemma B.1 to \mathbf{S}_1 leads to the desirable result. This ends the proof of Lemma B.2. \blacksquare

The following lemma establishes the rate of convergence of the sample covariance matrix in max norm.

Lemma B.3. Assume i.i.d. tensor data $\mathcal{T}, \mathcal{T}_1, \dots, \mathcal{T}_n \in \mathbb{R}^{m_1 \times m_2 \times \dots \times m_K}$ follows the tensor normal distribution $\text{TN}(\mathbf{0}; \Sigma_1^*, \dots, \Sigma_K^*)$, and assume Condition 3.2 holds. Let $\hat{\Omega}_j \in \mathbb{R}^{m_j \times m_j}, j \neq k$, be the estimated precision matrix from Algorithm 1 with iteration number $T = 1$. Denote the k -th sample covariance matrix as

$$\hat{\mathbf{S}}_k = \frac{m_k}{nm} \sum_{i=1}^n \hat{\mathbf{V}}_i \hat{\mathbf{V}}_i^\top,$$

with $m = \prod_{k=1}^K m_k$ and $\hat{\mathbf{V}}_i := [\mathcal{T}_i \times \{\hat{\Omega}_1^{1/2}, \dots, \hat{\Omega}_{k-1}^{1/2}, \mathbf{1}_{m_k}, \hat{\Omega}_{k+1}^{1/2}, \dots, \hat{\Omega}_K^{1/2}\}]_{(k)}$. We have

$$\max_{s,t} [\hat{\mathbf{S}}_k - \Sigma_k^*]_{s,t} = O_P \left(\max_{j=1, \dots, K} \sqrt{\frac{(m_j + s_j) \log m_j}{nm}} \right). \quad (\text{B.10})$$

Proof of Lemma B.3: The proof follows by decomposing the $\hat{\mathbf{S}}_k - \Sigma_k^*$ into two parts and then applying Lemma B.2 and Theorem 3.5 for each part to bound the final error.

Note that the triangle inequality implies that

$$\|\hat{\mathbf{S}}_k - \Sigma_k^*\|_\infty \leq \underbrace{\left\| \hat{\mathbf{S}}_k - \frac{m_k [\prod_{j \neq k} \text{tr}(\Sigma_j^* \hat{\Omega}_j)]}{m} \Sigma_k^* \right\|_\infty}_{I_1} + \underbrace{\left\| \frac{m_k [\prod_{j \neq k} \text{tr}(\Sigma_j^* \hat{\Omega}_j)]}{m} \Sigma_k^* - \Sigma_k^* \right\|_\infty}_{I_2}.$$

Note that here the covariance matrix $\hat{\mathbf{S}}_k$ is constructed based on the estimators $\hat{\Omega}_j, j \neq k$. According to (B.8) in Lemma B.2, we have

$$I_1 = O_P \left(\sqrt{\frac{m_k \log m_k}{nm}} \right).$$

The remainder part is to bound the error I_2 . Note that $\text{tr}(\Sigma_j^* \Omega_j^*) = \text{tr}(\mathbf{1}_{m_j}) = m_j$. Therefore,

$$I_2 = \underbrace{\left\| \frac{m_k}{m} \left[\prod_{j \neq k} \text{tr}(\Sigma_j^* \hat{\Omega}_j) - \prod_{j \neq k} \text{tr}(\Sigma_j^* \Omega_j^*) \right] \right\|_\infty}_{I_3} \|\Sigma_k^*\|_\infty.$$

Given that $\|\Sigma_k^*\|_\infty = O_P(1)$, it is sufficient to bound the coefficient I_3 . We only demonstrate the proofs with $K = 3$ and $k = 1$. The extension to a general K follows similarly. In this case, we have

$$\begin{aligned} I_3 &= \frac{m_1}{m} \left| \text{tr}(\Sigma_2^* \widehat{\Omega}_2) \text{tr}(\Sigma_3^* \widehat{\Omega}_3) - \text{tr}(\Sigma_2^* \Omega_2^*) \text{tr}(\Sigma_3^* \Omega_3^*) \right| \\ &\leq \left| \frac{\text{tr}(\Sigma_2^* \widehat{\Omega}_2) \text{tr}[\Sigma_3^* (\widehat{\Omega}_3 - \Omega_3^*)]}{m_2 m_3} \right| + \left| \frac{\text{tr}[\Sigma_2^* (\widehat{\Omega}_2 - \Omega_2^*)] \text{tr}(\Sigma_3^* \Omega_3^*)}{m_2 m_3} \right|. \end{aligned}$$

According to the proof of Theorem 3.5, we have $C_1 \leq \text{tr}(\Sigma_j^* \Omega_j^*)/m_j \leq 1/C_1$ for any $j = 1, \dots, K$ and some constant $C_1 > 0$. Moreover, we have $\text{tr}(\Sigma_3^* \Omega_3^*) = m_3$. Therefore, we have

$$I_3 \leq \left| \frac{\text{tr}[\Sigma_3^* (\widehat{\Omega}_3 - \Omega_3^*)]}{m_3} \right| + \left| \frac{\text{tr}[\Sigma_2^* (\widehat{\Omega}_2 - \Omega_2^*)]}{m_2} \right|.$$

Here $\text{tr}[\Sigma_j^* (\widehat{\Omega}_j - \Omega_j^*)] \leq \|\Sigma_j^*\|_F \|\widehat{\Omega}_j - \Omega_j^*\|_F \leq \sqrt{m_j} \|\Sigma_j^*\|_2 \|\widehat{\Omega}_j - \Omega_j^*\|_F$. According to Condition 3.2, $\|\Sigma_j^*\|_2 = O_P(1)$. This together with Theorem 3.5 implies that

$$I_3 = O_P \left(\sqrt{\frac{(m_3 + s_3) \log m_3}{nm}} + \sqrt{\frac{(m_2 + s_2) \log m_2}{nm}} \right).$$

By generalizing it to a general K and k , we have that

$$I_3 = O_P \left(\max_{j \neq k} \sqrt{\frac{(m_j + s_j) \log m_j}{nm}} \right),$$

and hence

$$\|\widehat{\Sigma}_k - \Sigma_k^*\|_\infty = O_P \left(\sqrt{\frac{m_k \log m_k}{nm}} + \max_{j \neq k} \sqrt{\frac{(m_j + s_j) \log m_j}{nm}} \right),$$

which leads to the desirable result. This ends the proof of Lemma B.3. \blacksquare

C Auxiliary lemmas

Lemma C.1. Assume a random matrix $\mathbf{X} \in \mathbb{R}^{p \times q}$ follows the matrix-variate normal distribution such that $\text{vec}(\mathbf{X}) \sim N(\mathbf{0}; \Psi^* \otimes \Sigma^*)$ with $\Psi^* \in \mathbb{R}^{q \times q}$ and $\Sigma^* \in \mathbb{R}^{p \times p}$. Then for any symmetric and positive definite matrix $\Omega \in \mathbb{R}^{p \times p}$, we have $\mathbb{E}(\mathbf{X}^\top \Omega \mathbf{X}) = \Psi^* \text{tr}(\Omega \Sigma^*)$.

Proof of Lemma C.1: Since the matrix Ω is symmetric and positive definite, it has the Cholesky decomposition $\Omega = \mathbf{V}^\top \mathbf{V}$, where \mathbf{V} is upper triangular with positive diagonal entries. Let $\mathbf{Y} := \mathbf{V} \mathbf{X}$ and denote the j -th row of matrix \mathbf{Y} as $\mathbf{y}_j = (y_{j,1}, \dots, y_{j,q})$. We have $\mathbb{E}(\mathbf{X}^\top \Omega \mathbf{X}) = \mathbb{E}(\mathbf{Y}^\top \mathbf{Y}) = \sum_{j=1}^p \mathbb{E}(\mathbf{y}_j^\top \mathbf{y}_j)$. Here $\mathbf{y}_j = \mathbf{v}_j \mathbf{X}$ with \mathbf{v}_j the j -th row of \mathbf{V} . Denote the i -th column of matrix \mathbf{X} as $\mathbf{x}_{(i)}$, we have $y_{j,i} = \mathbf{v}_j \mathbf{x}_{(i)}$. Therefore, the (s, t) -th entry of $\mathbb{E}(\mathbf{y}_j^\top \mathbf{y}_j)$ is

$$[\mathbb{E}(\mathbf{y}_j^\top \mathbf{y}_j)]_{(s,t)} = \mathbb{E}[\mathbf{v}_j \mathbf{x}_{(s)} \mathbf{v}_j \mathbf{x}_{(t)}] = \mathbf{v}_j \mathbb{E}[\mathbf{x}_{(s)} \mathbf{x}_{(t)}^\top] \mathbf{v}_j^\top = \mathbf{v}_j \Psi_{s,t}^* \Sigma^* \mathbf{v}_j^\top,$$

where $\Psi_{s,t}^*$ is the (s, t) -th entry of Ψ^* . The last equality is due to $\text{vec}(\mathbf{X}) = (\mathbf{x}_{(1)}^\top, \dots, \mathbf{x}_{(q)}^\top)^\top \sim N(\mathbf{0}; \Psi^* \otimes \Sigma^*)$. Therefore, we have

$$\mathbb{E}(\mathbf{X}^\top \Omega \mathbf{X}) = \sum_{j=1}^p \mathbb{E}(\mathbf{y}_j^\top \mathbf{y}_j) = \Psi^* \sum_{j=1}^p \mathbf{v}_j \Sigma^* \mathbf{v}_j^\top = \Psi^* \text{tr} \left(\sum_{j=1}^p \mathbf{v}_j^\top \mathbf{v}_j \Sigma^* \right) = \Psi^* \text{tr}(\Omega \Sigma^*).$$

This ends the proof of Lemma C.1. \blacksquare

The following lemma is stated by [24].

Lemma C.2. Let random variables $x_1, \dots, x_n \in \mathbb{R}$ be i.i.d. drawn from standard normal $N(0; 1)$ and denote $\mathbf{x} = (x_1, \dots, x_n)^\top \in \mathbb{R}^n$ be a random vector. For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ with Lipschitz constant L , that is, for any vectors $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{R}^n$, there exists $L \geq 0$ such that $|f(\mathbf{v}_1) - f(\mathbf{v}_2)| \leq L\|\mathbf{v}_1 - \mathbf{v}_2\|_2$. Then, for any $t > 0$, we have

$$\mathbb{P}\{|f(\mathbf{x}) - \mathbb{E}[f(\mathbf{x})]| > t\} \leq 2 \exp\left(-\frac{t^2}{2L^2}\right).$$

The following lemma is useful for the proof of Lemma B.1. A similar statement was given in Lemma I.2 of [33].

Lemma C.3. Suppose that a d -dimensional Gaussian random vector $\mathbf{y} \sim N(0; \mathbf{Q})$, Then, for any $t > 2/\sqrt{d}$, we have

$$\mathbb{P}\left[\frac{1}{d}|\|\mathbf{y}\|_2^2 - \mathbb{E}(\|\mathbf{y}\|_2^2)| > 4t\|\mathbf{Q}\|_2\right] \leq 2 \exp\left\{-\frac{d(t - 2/\sqrt{d})^2}{2}\right\} + 2 \exp(-d/2).$$

Proof of Lemma C.3: Note that $\mathbb{E}(\|\mathbf{y}\|_2^2) \leq [\mathbb{E}(\|\mathbf{y}\|_2)]^2$ and hence

$$\|\mathbf{y}\|_2^2 - \mathbb{E}(\|\mathbf{y}\|_2^2) \leq [\|\mathbf{y}\|_2 - \mathbb{E}(\|\mathbf{y}\|_2)][\|\mathbf{y}\|_2 + \mathbb{E}(\|\mathbf{y}\|_2)].$$

The term $(\|\mathbf{y}\|_2 - \mathbb{E}(\|\mathbf{y}\|_2))$ can be bounded via the concentration inequality in Lemma C.2 by noting that $\|\mathbf{y}\|_2$ is a Lipschitz function of Gaussian random vector \mathbf{y} . The term $\|\mathbf{y}\|_2 + \mathbb{E}(\|\mathbf{y}\|_2)$ can also be bounded by the large deviation bound since \mathbf{y} is a Gaussian random vector. This ends the proof of Lemma C.3. \blacksquare

D Additional simulation results

In this section, we explain the details in generating the true precision matrices and then show additional numerical results.

Triangle: For each $k = 1, \dots, K$, we construct the covariance matrix $\Sigma_k \in \mathbb{R}^{m_k \times m_k}$ such that its (i, j) -th entry is $[\Sigma_k]_{i,j} = \exp(-|h_i - h_j|/2)$ with $h_1 < h_2 < \dots < h_{m_k}$. The difference $h_i - h_{i-1}$ with $i = 2, \dots, m_k$ is generated independently and identically from $\text{Unif}(0.5, 1)$. This generated covariance matrix mimics the autoregressive process of order one, i.e., AR(1). We set $\Omega_k^* = \Sigma_k^{-1}$.

Nearest neighbor: For each $k = 1, \dots, K$, we construct the precision matrix $\Omega_k \in \mathbb{R}^{m_k \times m_k}$ directly from a four nearest-neighbor network. We first randomly pick m_k points from a unit square and compute all pairwise distances among the points. We then search for the four nearest-neighbors of each point and a pair of symmetric entries in the precision matrix Ω_k that has a random chosen value from $[-1, -0.5] \cup [0.5, 1]$. To ensure its positive definite property, we let the final precision matrix as $\Omega_k^* = \Omega_k + (|\lambda_{\min}(\Omega_k)| + 0.2) \cdot \mathbf{1}_{m_k}$, where $\lambda_{\min}(\cdot)$ refers to the smallest eigenvalue.

The additional error criterions for comparison are the averaged estimation errors in Frobenius norm and max norm, i.e.,

$$\frac{1}{K} \sum_{k=1}^K \|\hat{\Omega}_k - \Omega_k^*\|_F, \quad \frac{1}{K} \sum_{k=1}^K \|\hat{\Omega}_k - \Omega_k^*\|_\infty.$$

Note that these two criterions are only available to the P-MLE method and our Tlasso. The direct Glasso method estimate the whole Kronecker product and hence could not produce the estimator for each precision matrix.

Remind that, as we show in Theorem 3.5 and Theorem 3.9, the estimation error for the k -th precision matrix is $O_p(\sqrt{m_k(m_k + s_k) \log m_k / (nm)})$ in Frobenius norm or $O_p(\sqrt{m_k \log m_k / (nm)})$ in max norm, where $m = m_1 m_2 m_3$ in this example. These theoretical findings are supported by the numerical results in Figure 2. In particular, as sample size n increases from Scenario s1 to s2, the estimation errors in both Frobenius norm and max norm expectedly decrease. From Scenario s1 to s3, one dimension m_1 increases from 10 to 100, and other dimensions m_2, m_3 decrease from 10 to 5, in which case the averaged estimation error in max norm is decreasing, while the error in Frobenius norm increases due to its additional $\sqrt{m_k + s_k}$ effect. Moreover, compared to the P-MLE

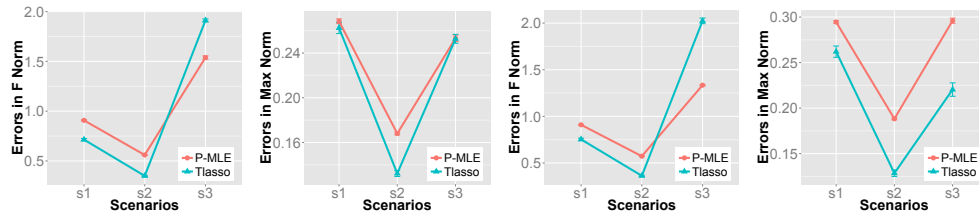


Figure 2: Averaged estimation errors of the precision matrices in Frobenius norm and max norm of each method in Simulations 1&2, respectively. The left two plots are for Simulation 1, and the right two are for Simulation 2.

method, our Tlasso is better in Scenarios s1 and s2 and is worse in Scenario s3 in Frobenius norm. However, in terms of the max norm, our Tlasso delivers significant better performance in 4 scenarios and comparable results in the rest 2 scenarios.

References

- [1] S. Rendle and L. Schmidt-Thieme. Pairwise interaction tensor factorization for personalized tag recommendation. In *International Conference on Web Search and Data Mining*, 2010.
- [2] G.I. Allen. Sparse higher-order principal components analysis. In *International Conference on Artificial Intelligence and Statistics*, 2012.
- [3] J. Zahn, S. Poosala, A. Owen, D. Ingram, et al. AGEMAP: A gene expression database for aging in mice. *PLOS Genetics*, 3:2326–2337, 2007.
- [4] T. Cai, W. Liu, and H.H. Zhou. Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation. *Annals of Statistics*, 2015.
- [5] C. Leng and C.Y. Tang. Sparse matrix graphical models. *Journal of the American Statistical Association*, 107:1187–1200, 2012.
- [6] J. Yin and H. Li. Model selection and estimation in the matrix normal graphical model. *Journal of Multivariate Analysis*, 107:119–140, 2012.
- [7] T. Tsiligkaridis, A. O. Hero, and S. Zhou. On convergence of Kronecker graphical Lasso algorithms. *IEEE Transactions on Signal Processing*, 61:1743–1755, 2013.
- [8] S. Zhou. Gemini: Graph estimation with matrix variate normal instances. *Annals of Statistics*, 42:532–562, 2014.
- [9] S. He, J. Yin, H. Li, and X. Wang. Graphical model selection and estimation for high dimensional tensor data. *Journal of Multivariate Analysis*, 128:165–185, 2014.
- [10] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Symposium on Theory of Computing*, pages 665–674, 2013.
- [11] P. Netrapalli, P. Jain, and S. Sanghavi. Phase retrieval using alternating minimization. In *Advances in Neural Information Processing Systems*, pages 2796–2804, 2013.
- [12] J. Sun, Q. Qu, and J. Wright. Complete dictionary recovery over the sphere. *arXiv:1504.06785*, 2015.
- [13] S. Arora, R. Ge, T. Ma, and A. Moitra. Simple, efficient, and neural algorithms for sparse coding. *arXiv:1503.00778*, 2015.
- [14] A. Anandkumar, R. Ge, D. Hsu, S. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- [15] W. Sun, J. Lu, H. Liu, and G. Cheng. Provable sparse tensor decomposition. *arXiv:1502.01425*, 2015.
- [16] S. Zhe, Z. Xu, X. Chu, Y. Qi, and Y. Park. Scalable nonparametric multiway data analysis. In *International Conference on Artificial Intelligence and Statistics*, 2015.
- [17] S. Zhe, Z. Xu, Y. Qi, and P. Yu. Sparse bayesian multiview learning for simultaneous association discovery and diagnosis of alzheimer’s disease. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- [18] T. Kolda and B. Bader. Tensor decompositions and applications. *SIAM Review*, 51:455–500, 2009.
- [19] R. Tibshirani. Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, 58:267–288, 1996.
- [20] M. Yuan and Y. Lin. Model selection and estimation in the gaussian graphical model. *Biometrika*, 94:19–35, 2007.
- [21] J. Friedman, H. Hastie, and R. Tibshirani. Sparse inverse covariance estimation with the graphical Lasso. *Biostatistics*, 9:432–441, 2008.
- [22] A. J. Rothman, P. J. Bickel, E. Levina, and J. Zhu. Sparse permutation invariant covariance estimation. *Electronic Journal of Statistics*, 2:494–515, 2008.
- [23] W. Sun, J. Wang, and Y. Fang. Consistent selection of tuning parameters via variable selection stability. *Journal of Machine Learning Research*, 14:3419–3440, 2013.

- [24] M. Ledoux and M. Talagrand. *Probability in Banach Spaces: Isoperimetry and Processes*. Springer, 2011.
- [25] J. Fan, Y. Feng, and Y. Wu. Network exploration via the adaptive Lasso and scad penalties. *Annals of Statistics*, 3:521–541, 2009.
- [26] P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine Learning Research*, 7:2541–2567, 2006.
- [27] P. Ravikumar, M.J. Wainwright, G. Raskutti, and B. Yu. High-dimensional covariance estimation by minimizing ℓ_1 -penalized log-determinant divergence. *Electronic Journal of Statistics*, 5:935–980, 2011.
- [28] Z. Wang, H. Liu, and T. Zhang. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *Annals of Statistics*, 42:2164–2201, 2014.
- [29] T. Zhao, H. Liu, K. Roeder, J. Lafferty, and L. Wasserman. The huge package for high-dimensional undirected graph estimation in R. *Journal of Machine Learning Research*, 13:1059–1062, 2012.
- [30] A. Gupta and D. Nagar. *Matrix variate distributions*. Chapman and Hall/CRC Press, 2000.
- [31] P. Hoff. Separable covariance arrays via the Tucker product, with applications to multivariate relational data. *Bayesian Analysis*, 6:179–196, 2011.
- [32] A.P. Dawid. Some matrix-variate distribution theory: Notational considerations and a bayesian application. *Biometrika*, 68:265–274, 1981.
- [33] S. Negahban and M.J. Wainwright. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *Annals of Statistics*, 39:1069–1097, 2011.