

# A Consistent Method for Graph Based Anomaly Localization

## – Supplemental Material –

---

**Satoshi Hara**      **Tetsuro Morimura**      **Toshihiro Takahashi**      **Hiroki Yanagisawa**  
 IBM Research – Tokyo, Japan, {satohara, tetsuro, e30137, yanagis}@jp.ibm.com

**Taiji Suzuki**  
 Tokyo Institute of Technology & PRESTO, JST, Japan, s-taiji@is.titech.ac.jp

## 6 Settings for the Synthetic Data

### 6.1 Generation of Synthetic Matrices

Let  $\Lambda_A$  and  $\Lambda_B$  be a pair of matrices. For our simulation in Section 4.1, we generated matrices that satisfy Assumption 1. Specifically, we generated matrices that simultaneously satisfy three conditions:

- positive-definite:  $\Lambda_u \succeq \nu I_d$  for  $u \in \{A, B\}$ ,
- sparse:  $\Lambda_{u,ij} = 0$  for  $(i, j) \in \Omega_u$  with  $u \in \{A, B\}$ ,
- homogeneous:  $\Lambda_{A,ij} = \Lambda_{B,ij}$  for  $(i, j) \in \Pi$ ,

where  $\nu$  is some positive constant, and  $\Omega_A, \Omega_B$ , and  $\Pi$  are properly chosen subsets of  $\{1, 2, \dots, d\}^2$ . In particular, all pairs of healthy variable indices have to be involved in  $\Pi$  so that the homogeneity is satisfied. In addition, based on the sparseness and homogeneity conditions,  $\Omega_A \cap \Omega_B \subseteq \Pi$  must be satisfied. We note that, although the first condition is not necessary for general adjacency matrices, we used it since both the covariance and precision matrices we used in the simulation are positive-definite matrices.

The matrix generation procedure starts by specifying  $\nu, \Omega_A, \Omega_B, \Pi$ , and two reference matrices  $C_A, C_B \in \mathbb{R}^{d \times d}$ . Here we chose  $C_A$  and  $C_B$  to be symmetric but not necessarily to be positive-definite. We then derive  $\Lambda_A$  and  $\Lambda_B$  as the solution to the problem:

$$\begin{aligned}
 & \min_{\{\Lambda_u\}_{u \in \{A, B\}}} \frac{1}{2} \sum_{u \in \{A, B\}} \|\Lambda_u - C_u\|_F^2, \\
 & \text{s.t. } \Lambda_u \succeq \nu I_d \quad (u \in \{A, B\}), \\
 & \quad \Lambda_{u,ij} = 0 \quad \text{for } (i, j) \in \Omega_u \quad (u \in \{A, B\}), \\
 & \quad \Lambda_{A,ij} = \Lambda_{B,ij} \quad \text{for } (i, j) \in \Pi,
 \end{aligned}$$

where  $\|\cdot\|_F$  denotes the Frobenius-norm of the matrix. This problem corresponds to searching for the matrices

$\Lambda_A$  and  $\Lambda_B$  that are closest to the reference matrices  $C_A$  and  $C_B$  under the specified conditions. The problem is convex and we can solve it by using ADMM [24]. We first rewrite the problem into an equivalent form:

$$\begin{aligned}
 & \min_{\{X_u\}_{u \in \{A, B\}}, \{Y_u\}_{u \in \{A, B\}}} \frac{1}{2} \sum_{u \in \{A, B\}} (\|X_u - C_u\|_F^2 \\
 & \quad + \delta_{\Omega_u}(X_u) + \tilde{\delta}(Y_u)), \\
 & \text{s.t. } X_u - Y_u + \nu I_d = 0 \quad (u \in \{A, B\}), \\
 & \quad X_{A,ij} - X_{B,ij} = 0 \quad \text{for } (i, j) \in \Pi,
 \end{aligned}$$

where  $\delta_{\Omega_u}(X_u)$  and  $\tilde{\delta}(Y_u)$  are the indicator functions defined as

$$\begin{aligned}
 \delta_{\Omega_u}(X_u) & := \begin{cases} 0 & \text{if } X_{u,ij} = 0 \text{ for all } (i, j) \in \Omega_u, \\ \infty & \text{otherwise,} \end{cases} \\
 \tilde{\delta}(Y_u) & := \begin{cases} 0 & \text{if } Y_u \succeq 0, \\ \infty & \text{otherwise,} \end{cases}
 \end{aligned}$$

and we have  $\Lambda_u = Y_u$  as the solution for  $u \in \{A, B\}$ . Let  $Z_u \in \mathbb{R}^{d \times d}$  ( $u \in \{A, B, O\}$ ) be the matrix of Lagrange multipliers. We then define the Augmented Lagrangian (AL) function as

$$\begin{aligned}
 \mathcal{L}_\beta(X, Y, Z) & = \frac{1}{2} \sum_{u \in \{A, B\}} (\|X_u - C_u\|_F^2 + \delta_{\Omega_u}(X_u) + \tilde{\delta}(Y_u)) \\
 & \quad + \frac{\beta}{2} \left( \sum_{u \in \{A, B\}} \|X_u - Y_u + \nu I_d + \frac{1}{\beta} Z_u\|_F^2 \right. \\
 & \quad \left. + \|E_\Pi \odot (X_A - X_B) + \frac{1}{\beta} Z_O\|_F^2 \right),
 \end{aligned}$$

where we set  $X := \{X_u\}_{u \in \{A, B\}}$ ,  $Y := \{Y_u\}_{u \in \{A, B\}}$ , and  $Z := \{Z_u\}_{u \in \{A, B, O\}}$  to simplify the notation,  $\odot$  denotes the Hadamard product of matrices, and  $E_\Pi$  is an indicator matrix of the set  $\Pi$  defined as

$$E_{\Pi,ij} := \begin{cases} 1 & \text{if } (i, j) \in \Pi, \\ 0 & \text{otherwise.} \end{cases}$$

The optimization procedure of ADMM is defined using this AL function. Specifically, we repeat these steps until one of the termination criteria is fulfilled:

$$\begin{cases} X^{(k+1)} \in \operatorname{argmin}_X \mathcal{L}_\beta(X, Y^{(k)}, Z^{(k)}), \\ Y^{(k+1)} \in \operatorname{argmin}_Y \mathcal{L}_\beta(X^{(k+1)}, Y, Z^{(k)}), \\ Z_u^{(k+1)} = Z_u^{(k)} + \beta(X_u^{(k+1)} - Y_u^{(k+1)} + \nu I_d) \quad (u \in \{A, B\}) \\ Z_O^{(k+1)} = Z_O^{(k)} + \beta E_\Pi \odot (X_A^{(k+1)} - X_B^{(k+1)}). \end{cases}$$

The first step, the update of  $X$ , can be decomposed into individual problems on each  $(i, j)$ th entry given as

$$\begin{aligned} \min_{X_{A,ij}, X_{B,ij}} \frac{1}{2} \sum_{u \in \{A, B\}} \{ & (X_{u,ij} - C_{u,ij})^2 + \beta (X_{u,ij} - P_{u,ij}^{(k)})^2 \} \\ & + \frac{\beta}{2} \{ E_{\Pi,ij} (X_{A,ij} - X_{B,ij}) + \frac{1}{\beta} Z_{O,ij} \}^2, \end{aligned}$$

s.t.  $X_{u,ij} = 0$  for  $(i, j) \in \Omega_u$  ( $u \in \{A, B\}$ ),

where  $P_u^{(k)}$  is the matrix defined as

$$P_u^{(k)} = Y_u^{(k)} - \frac{1}{\beta} Z_u^{(k)} - \nu I_d \quad (u \in \{A, B\}).$$

When  $E_{\Pi,ij} = 0$  or  $(i, j) \notin \Pi$ , the problem can be further reduced into the individual problems on  $X_{A,ij}$  and  $X_{B,ij}$ . Hence, we have the solution

$$X_{u,ij}^{(k+1)} = \begin{cases} \frac{1}{1+\beta} (C_{u,ij} + \beta P_{u,ij}^{(k)}) & \text{if } (i, j) \notin \Omega_u, \\ 0 & \text{otherwise.} \end{cases}$$

For the case of  $E_{\Pi,ij} = 1$  and  $(i, j) \notin \Omega_u$  ( $u \in \{A, B\}$ ), we have the solution

$$\begin{aligned} \begin{bmatrix} X_{A,ij} \\ X_{B,ij} \end{bmatrix} &= \frac{1}{(1+\beta)(1+3\beta)} \\ &\times \begin{bmatrix} 1+2\beta & \beta \\ \beta & 1+2\beta \end{bmatrix} \begin{bmatrix} C_{A,ij} + \beta P_{A,ij}^{(k)} - Z_{O,ij}^{(k)} \\ C_{B,ij} + \beta P_{B,ij}^{(k)} + Z_{O,ij}^{(k)} \end{bmatrix}. \end{aligned}$$

The update problem for  $Y$  can be decomposed into individual problems on  $Y_A$  and  $Y_B$ , which are given as

$$\begin{aligned} \min_{Y_u} \frac{1}{2} \|Y_u - Q_u^{(k)}\|_F^2, \quad \text{s.t. } Y_u \succeq 0 \quad (u \in \{A, B\}), \\ Q_u^{(k)} = X_u^{(k+1)} + \frac{1}{\beta} Z_u^{(k)} + \nu I_d. \end{aligned}$$

This problem is equivalent to the Euclidean projection of the matrix  $Q_u^{(k)}$  onto the positive semidefinite cone. Hence, this can be computed analytically using the eigenvalue decomposition. Here we assume all of the matrices are symmetric, which can be assured by initializing all of the matrices to be symmetric. Let  $Q_u^{(k)} = UDU^\top$  be the eigenvalue decomposition with  $D = \operatorname{diag}(\sigma_1, \sigma_2, \dots, \sigma_d)$ . We then have  $Y_u^{(k+1)} = U\tilde{D}U^\top$  with  $D = \operatorname{diag}(\tilde{\sigma}_1, \tilde{\sigma}_2, \dots, \tilde{\sigma}_d)$  where  $\tilde{\sigma}_i = \max(\sigma_i, 0)$ .

## 6.2 Parameter Settings

This subsection explains how the sparsity patterns  $\Omega_A$  and  $\Omega_B$  and the shared pattern  $\Pi$  are chosen. For the sparsity pattern, we randomly pick index pairs so that the average size of  $\Omega_u$  is 3,000 when  $d = 100$ . This corresponds to choosing 70% of the matrix entries to be zero. For  $d = 200$ , they are set to be 8,000 (80%). The set  $\Pi$ , the shared pattern, is chosen as  $\Pi = (\mathcal{I} \times \mathcal{I}) \cup (\Omega_A \cap \Omega_B) \cup \Pi_+$ . The first set is all pairs of healthy variables and the second set corresponds to the shared zero entries between the two matrices. The set  $\Pi_+$  specifies the common non-zero entries, so that  $\Pi_+ \subseteq \Omega_A^c \cap \Omega_B^c$  holds where  $c$  denotes the compliment of the set. For the construction of  $\Pi_+$ , we consider the subset of  $\Omega_A^c \cap \Omega_B^c$  defined as

$$\Pi_0 := (\mathcal{I} \times \mathcal{I})^c \cap (\Omega_A^c \cap \Omega_B^c),$$

which means that  $\Pi_0$  specifies the set of index pairs whose corresponding edges are non-zeros, and hence connected to anomalous variables. We then randomly pick index pairs  $(i, j) \in \Pi_0$  and add to  $\Pi_+$ . We set the size  $|\Pi|$  to be 70% of  $|\Pi_0|$  for both  $d = 100$  and 200.

Across all of the settings, the value of  $\nu$  is set as  $10^{-3}$ . Here is how the reference matrices  $C_A$  and  $C_B$  are generated. We first generate random matrices  $L_A, L_B \in \mathbb{R}^{d \times \lceil \sqrt{d} \rceil}$  where each entry of each of the matrices is generated from a standard Gaussian distribution  $\mathcal{N}(0, 1)$ . We then set  $C_A = L_A L_A^\top$  and  $C_B = L_B L_B^\top$ , and rescale them so that their diagonals are one. Note that the generated matrices have at most rank  $\lceil \sqrt{d} \rceil$  which implies that they are rank deficient. Therefore, the resulting covariance and precision matrices encourage variables to have higher dependencies with other variables.

## 6.3 Generation of Matrices with a Concentrated Anomaly Pattern

In Section 4.1, Figure 6, we used a concentrated anomaly pattern for the simulation. The synthetic matrices were generated by modifying the generated matrices in Section 6.1. Here, we denote by  $M_{\mathcal{I} \times \mathcal{J}}$  the sub-matrix consisting of the components of a matrix  $M$  with indices in  $\mathcal{I} \times \mathcal{J}$ . We choose a subset  $\mathcal{H} \subseteq \mathcal{I}$  and update the matrix  $\Lambda_A$  by  $\Lambda_{A, \mathcal{H} \times \mathcal{I}} \leftarrow a \Lambda_{A, \mathcal{H} \times \mathcal{I}}$  and  $\Lambda_{A, \mathcal{I} \times \mathcal{H}} \leftarrow a \Lambda_{A, \mathcal{I} \times \mathcal{H}}$  where  $a > 1$  is some large positive value. The matrix  $\Lambda_B$  is updated in the same manner. Since the modified matrices  $\Lambda_A$  and  $\Lambda_B$  are no longer positive definite, we add an identity matrix  $\Lambda_A \leftarrow \Lambda_A + cI_d$  and  $\Lambda_B \leftarrow \Lambda_B + cI_d$  so that the matrices become positive definite. In our simulation in Section 4.1, we set the size of  $\mathcal{H}$  to be 5 and  $a = 10$ .

## 7 Proofs of Theorems

Here, we provide the proofs of the theorems in the manuscript. We first enumerate key lemmas relevant to a GGM estimator that are needed in the proofs, and then present the proofs of theorems.

### 7.1 Key Lemmas on GGM estimators

Suppose that we have two groups of i.i.d. samples  $\{x_i^{(A)}\}_{i=1}^{n_A} \sim p_A$  and  $\{x_i^{(B)}\}_{i=1}^{n_B} \sim p_B$  ( $x_i^{(A)} \in \mathbb{R}^d$  and  $x_i^{(B)} \in \mathbb{R}^d$ ). To estimate the precision matrices  $\Lambda_A$  and  $\Lambda_B$  corresponding to  $p_A$  and  $p_B$ , respectively, we use the following graphical-Lasso type estimators:

$$\hat{\Lambda}_A := \operatorname{argmin}_{\Lambda \succ 0} -\log \det(\Lambda) + \operatorname{tr}[\hat{\Sigma}_A \Lambda] + \lambda_n^{(A)} \|\Lambda\|_{1,\text{off}},$$

$$\hat{\Lambda}_B := \operatorname{argmin}_{\Lambda \succ 0} -\log \det(\Lambda) + \operatorname{tr}[\hat{\Sigma}_B \Lambda] + \lambda_n^{(B)} \|\Lambda\|_{1,\text{off}},$$

where  $\hat{\Sigma}_A := \sum_{i=1}^{n_A} x_i^{(A)} x_i^{(A)\top} / n_A$ ,  $\hat{\Sigma}_B := \sum_{i=1}^{n_B} x_i^{(B)} x_i^{(B)\top} / n_B$ , and  $\|\Lambda\|_{1,\text{off}} := \sum_{i \neq j} |\Lambda_{ij}|$ .

Here we write  $\Sigma_A = \Lambda_A^{-1}$ , and  $\Sigma_B = \Lambda_B^{-1}$ . We define a linear operator  $\Phi_A : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$  as  $\Phi_A \Lambda = \Sigma_A \operatorname{tr}[\Sigma_A \Lambda]$  (in other words,  $\Phi_A = \Sigma_A \otimes \Sigma_A$ ), and  $\Phi_B : \mathbb{R}^{d \times d} \rightarrow \mathbb{R}^{d \times d}$  in the same manner.

Now let  $S_A := \{(i, j) \mid \Lambda_{A,ij} \neq 0\} \cup \{(1, 1), \dots, (d, d)\}$  and similarly  $S_B := \{(i, j) \mid \Lambda_{B,ij} \neq 0\} \cup \{(1, 1), \dots, (d, d)\}$ . Let  $s_A := \max_{i=1, \dots, p} |\{j \mid \Lambda_{A,ij} \neq 0\}|$  and  $s_B := \max_{i=1, \dots, p} |\{j \mid \Lambda_{B,ij} \neq 0\}|$ . For matrix  $M \in \mathbb{R}^{p \times p}$  ( $p \in \mathbb{N}$ ), we denote by  $M_{\mathcal{I} \times \mathcal{J}}$  the sub-matrix consisting of the components of  $M$  with indices in  $\mathcal{I} \times \mathcal{J}$ , and let  $\|M\|_{\infty, \infty}$  be an operator norm of  $M$  as an operator from  $\ell_\infty$  to  $\ell_\infty$ , that is,  $\|M\|_{\infty, \infty} := \max_{i=1, \dots, p} \sum_{j=1}^p |M_{ij}|$ .

Using these notations, we define

$$K_{\Sigma_A} := \|\Sigma_A\|_{\infty, \infty},$$

$$K_{\Psi_A} := \|(\Psi_{A, S \times S})^{-1}\|_{\infty, \infty} = \|((\Sigma_A \otimes \Sigma_A)_{S \times S})^{-1}\|_{\infty, \infty},$$

and define  $K_{\Sigma_B}$  and  $K_{\Psi_B}$  in the same manner.

Finally, we assume the following conditions on the distributions  $p_A$  and  $p_B$ .

#### Assumption 2

- (Sub-Gaussianity) Both  $p_A$  and  $p_B$  have zero means and sub-Gaussian tails:  $\exists \sigma > 0, \forall t \in \mathbb{R}$ ,

$$\mathbb{E}_{p_A}[\exp(tX)] \leq \exp(\sigma^2 t^2),$$

$$\mathbb{E}_{p_B}[\exp(tX)] \leq \exp(\sigma^2 t^2).$$

- (Incoherence) There exists a real number  $\alpha \in$

$(0, 1]$  such that

$$\|\Psi_{A, S^c \times S}(\Psi_{A, S \times S})^{-1}\|_{\infty, \infty} \leq 1 - \alpha,$$

$$\|\Psi_{B, S^c \times S}(\Psi_{B, S \times S})^{-1}\|_{\infty, \infty} \leq 1 - \alpha.$$

Now let  $n := \min\{n_A, n_B\}$  and we have the following lemma.

**Lemma 1** Under Assumption 2, if the sample size  $n$  satisfies

$$n > C_3 s_A^2 (1 + 3/\alpha)^2 (2 \log(2d/\eta) + \log 4),$$

$$C_3 := \{48\sqrt{2}(1 + 4\sigma^2) \max_i(\Sigma_{A,ii}) \max\{K_{\Sigma_A} K_{\Psi_A}, K_{\Sigma_A}^3 K_{\Psi_A}^2\}\}^2,$$

and the regularization parameter  $\lambda_n^{(A)}$  satisfies

$$\lambda_n^{(A)} = (8/\alpha) \sqrt{128(1 + 4\sigma^2)^2 \max_i(\Sigma_{A,ii})^2} \sqrt{\frac{2 \log(2d/\eta)}{n}},$$

then with probability greater than  $1 - \eta/2$  ( $\forall \eta > 0$ ), the estimated  $\hat{\Lambda}_A$  satisfies the bound,

$$\|\hat{\Lambda}_A - \Lambda_A\|_{\infty} \leq 16\sqrt{2}(1 + 4\sigma^2) \max_i(\Sigma_{A,ii}) (1 + 8\alpha^{-1})$$

$$\times K_{\Psi_A} \sqrt{\frac{2 \log(2d/\eta) + \log 4}{n}},$$

and the non-zero elements of  $\hat{\Lambda}_A$  are included in  $\Lambda_A$ ,

$$\{(i, j) \mid \hat{\Lambda}_{A,ij} \neq 0\} \subseteq S_A.$$

(Proof of Lemma 1) The statement is a modification of Corollary 1 of Ravikumar et al. [10]. Assumption 2 yields the conditions assumed in Theorem 1 and Corollary 1 (sub-Gaussian condition) [10], thus we can apply these theorems. The assertion is immediately proven by replacing  $\bar{\delta}_f(n, p^\tau)$  and  $\bar{n}_f(n, p^\tau)$  with  $\bar{\delta}_f(n, \eta/(2d^2))$  and  $\bar{n}_f(n, \eta/(2d^2))$  in the proof of Corollary 1 in the paper. If we replace  $d^\tau$  with  $\eta/(2d^2)$ , then the tail probability of the noise can be evaluated as

$$P[\|\hat{\Sigma} - \Sigma^*\|_{\infty} \geq \bar{\delta}_f(n, \eta/(2d^2))] \leq \frac{\eta}{2},$$

by Lemma 8 and its proof [10]. Finally, we notice  $\log(2d^2/\eta) \leq 2 \log(2d/\eta)$ , and obtain the assertion.  $\square$

We have the same statement for the estimation of  $\Lambda_B$ .

**Lemma 2** Under the conditions assumed in Lemma 1,  $\hat{\Gamma} - \Gamma$  satisfies the following inequality with probability greater than  $1 - \eta$  ( $\eta \in (0, 1)$ ):

$$\|\hat{\Gamma} - \Gamma\|_{\infty} \leq C_1 \sqrt{\frac{2 \log(2d/\eta) + \log 4}{n}},$$

$$C_1 := 16\sqrt{2}(1 + 4\sigma^2)(1 + 8\alpha^{-1})$$

$$\times (\max_i(\Sigma_{A,ii}^*) K_{\Gamma_A^*} + \max_i(\Sigma_{B,ii}^*) K_{\Gamma_B^*}).$$

In particular, we have

$$\begin{aligned} \|\hat{\Gamma} - \Gamma\|_{S_\infty} &\leq \|\hat{\Gamma} - \Gamma\|_F \\ &\leq C_1 \sqrt{\frac{d \min\{s_A + s_B, d\} \{2 \log(2d/\eta) + \log 4\}}{n}}, \end{aligned}$$

with probability greater than  $1 - \eta$ , where  $\|\cdot\|_{S_\infty}$  is the spectrum norm and  $\|\cdot\|_F$  is the Frobenius norm.

(Proof of Lemma 2) Since

$$\begin{aligned} & - |\hat{\Lambda}_{A,ij} - \Lambda_{A,ij} - \hat{\Lambda}_{B,ij} + \Lambda_{B,ij}| \\ & \leq |\hat{\Lambda}_{A,ij} - \hat{\Lambda}_{B,ij}| - |\Lambda_{A,ij} - \Lambda_{B,ij}| \\ & \leq |\hat{\Lambda}_{A,ij} - \Lambda_{A,ij} - \hat{\Lambda}_{B,ij} + \Lambda_{B,ij}|, \end{aligned}$$

we obtain

$$\begin{aligned} \|\|\hat{\Gamma} - \Gamma\|_\infty &\leq \|\|\hat{\Lambda}_A - \Lambda_A - \hat{\Lambda}_B + \Lambda_B\|_\infty \\ &\leq \|\|\hat{\Lambda}_A - \Lambda_A\|_\infty + \|\|\hat{\Lambda}_B - \Lambda_B\|_\infty. \end{aligned}$$

Thus, applying the bound derived in Lemma 1 to both  $\|\|\hat{\Lambda}_A - \Lambda_A\|_\infty$  and  $\|\|\hat{\Lambda}_B - \Lambda_B\|_\infty$ , we obtain the first assertion.

As for the second assertion, since  $\hat{\Gamma} - \Gamma$  has at most  $\min\{d(s_A + s_B), d^2\}$  non-zero components by the second assertion of Lemma 1, the Frobenius norm between  $\hat{\Gamma} - \Gamma$  has the bound,

$$\begin{aligned} \|\hat{\Gamma} - \Gamma\|_F &\leq \sqrt{\min\{d(s_A + s_B), d^2\} \|\|\hat{\Gamma} - \Gamma\|_\infty^2} \\ &= \sqrt{\min\{d(s_A + s_B), d^2\} \|\|\hat{\Gamma} - \Gamma\|_\infty}, \end{aligned}$$

which gives the assertion.  $\square$

This lemma gives a bound on the discrepancy between  $\lambda_{\min}(\hat{\Gamma})$  and  $\lambda_{\min}(\Gamma)$  as follows.

**Lemma 3** *Under the conditions assumed in Lemma 1, we have the following bound with probability greater than  $1 - \eta$ ,*

$$\begin{aligned} & |\lambda_{\min}(\hat{\Gamma}) - \lambda_{\min}(\Gamma)| \\ & \leq C_1 \sqrt{\frac{d \min\{s_A + s_B, d\} \{2 \log(2d/\eta) + \log 4\}}{n}}. \end{aligned}$$

(Proof of Lemma 3) For a real symmetric matrix  $Q \in \mathbb{R}^{d \times d}$ , let  $\lambda_i(Q)$  be the  $i$ -th smallest eigenvalue ( $\lambda_1(Q) \leq \lambda_2(Q) \leq \dots \leq \lambda_d(Q)$ ). For all real symmetric matrices  $Q, R \in \mathbb{R}^{d \times d}$ , the well-known Hoffman and Wielandt inequality (see Theorem 6.3.5 and its corollaries in Horn et al. [25]) yields

$$\sqrt{\sum_{i=1}^d (\lambda_i(Q) - \lambda_i(R))^2} \leq \|Q - R\|_F.$$

Then applying Lemma 2 to the right-hand side of this inequality and noticing the relation  $|\lambda_{\min}(Q) - \lambda_{\min}(R)| \leq \sqrt{\sum_{i=1}^d (\lambda_i(Q) - \lambda_i(R))^2}$ , we obtain the assertion.  $\square$

## 7.2 Proof of Theorem 1

The theorem is true for  $k = d$  and  $k = 0$  since  $\mathcal{I} = \hat{\mathcal{I}} = \{1, 2, \dots, d\}$  and  $\mathcal{I} = \hat{\mathcal{I}} = \emptyset$  hold, respectively. Therefore, we only need to consider the case when  $1 \leq k \leq d - 1$ .

Let  $\epsilon := \|\|\hat{\Gamma} - \Gamma\|_\infty$  and  $f(\mathcal{K}, \mathcal{K}'; M) := \sum_{i \in \mathcal{K}, j \in \mathcal{K}'} M_{ij}$  for a matrix  $M$ . We also set the index sets  $\mathcal{P}$ ,  $\mathcal{Q}$ , and  $\mathcal{R}$  as  $\mathcal{P} := \mathcal{I} \setminus \hat{\mathcal{I}}$ ,  $\mathcal{Q} := \hat{\mathcal{I}} \setminus \mathcal{I}$ , and  $\mathcal{R} := \mathcal{I} \cap \hat{\mathcal{I}}$ . We now have

$$\begin{aligned} f(\mathcal{I}, \mathcal{I}; \hat{\Gamma}) - f(\hat{\mathcal{I}}, \hat{\mathcal{I}}; \hat{\Gamma}) &= f(\mathcal{P}, \mathcal{P}; \hat{\Gamma}) + 2f(\mathcal{P}, \mathcal{R}; \hat{\Gamma}) - f(\mathcal{Q}, \mathcal{Q}; \hat{\Gamma}) - 2f(\mathcal{Q}, \mathcal{R}; \hat{\Gamma}) \\ &\leq f(\mathcal{P}, \mathcal{P}; \Gamma) + 2f(\mathcal{P}, \mathcal{R}; \Gamma) - f(\mathcal{Q}, \mathcal{Q}; \Gamma) - 2f(\mathcal{Q}, \mathcal{R}; \Gamma) \\ &\quad + \epsilon(|\mathcal{P}|^2 + |\mathcal{Q}|^2 + 2|\mathcal{P}||\mathcal{R}| + 2|\mathcal{Q}||\mathcal{R}|) \\ &\leq f(\mathcal{I}, \mathcal{I}; \Gamma) - f(\hat{\mathcal{I}}, \hat{\mathcal{I}}; \Gamma) + \epsilon(k^2 + d^2), \end{aligned}$$

where, in the first inequality, we used the fact that  $|f(\mathcal{K}, \mathcal{K}'; \hat{\Gamma}) - f(\mathcal{K}, \mathcal{K}'; \Gamma)| \leq \epsilon |\mathcal{K}| |\mathcal{K}'|$ , and in the second inequality, we used  $|\mathcal{P}|^2 + |\mathcal{Q}|^2 + 2|\mathcal{P}||\mathcal{R}| + 2|\mathcal{Q}||\mathcal{R}| \leq (|\mathcal{P}| + |\mathcal{R}|)^2 + (|\mathcal{Q}| + |\mathcal{R}|)^2 = |\mathcal{I}|^2 + |\hat{\mathcal{I}}|^2 \leq k^2 + d^2$ . Since this inequality is valid for all  $\hat{\mathcal{I}} \neq \mathcal{I}$ , we have

$$f(\mathcal{I}, \mathcal{I}; \hat{\Gamma}) - f(\hat{\mathcal{I}}, \hat{\mathcal{I}}; \hat{\Gamma}) \leq \epsilon(k^2 + d^2) - h.$$

From the assumption that  $\mathcal{I}$  is unique, we have  $h > 0$ . If  $\epsilon < h/(k^2 + d^2)$ , then the right hand side becomes negative implying  $\mathcal{I}$  is the minimizer of (1), which proves the claim.  $\square$

## 7.3 Proof of Theorem 2

The proof of the theorem immediately follows from the next two lemmas:

**Lemma 4** *Let  $\hat{A}_{\mu, \mathcal{I} \times \mathcal{I}} \in \mathbb{R}^{k \times k}$  and  $\hat{A}_{\mu, \mathcal{J} \times \mathcal{I}} \in \mathbb{R}^{(d-k) \times k}$  be sub-matrices of  $\hat{A}_\mu$  indexed by  $\mathcal{I}$  and  $\mathcal{J}$ . Suppose the following conditions hold for some  $\tau, \tau' > 0$ :*

$$\hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}^{-1} \mathbf{1}_k \geq \tau \mathbf{1}_k, \quad (5)$$

$$\hat{A}_{\mu, \mathcal{J} \times \mathcal{I}} \mathbf{1}_k \geq \frac{1 + \tau'}{\tau} \mathbf{1}_{d-k}. \quad (6)$$

Then  $\tilde{\mathcal{I}}_0 = \mathcal{I}$ .

**Lemma 5** *Suppose the conditions in Theorem 2 hold true. Then there exists  $\tau, \tau' > 0$  that satisfy the conditions (5) and (6).*

(Proof of Lemma 4) Since the problem (3) is a convex quadratic programming with a positive-definite matrix  $\hat{A}_\mu$ , the KKT conditions given here are both necessary

and sufficient for the optimality of the solution:

$$\hat{A}_\mu \tilde{\mathbf{s}} - \tilde{\boldsymbol{\zeta}} - \tilde{\nu} \mathbf{1}_d = \mathbf{0}_d, \quad (7)$$

$$\tilde{\mathbf{s}}, \tilde{\boldsymbol{\zeta}} \geq \mathbf{0}_d, \quad (8)$$

$$\tilde{\mathbf{s}} \odot \tilde{\boldsymbol{\zeta}} = \mathbf{0}_d, \quad (9)$$

$$\mathbf{1}_d^\top \tilde{\mathbf{s}} - 1 = 0, \quad (10)$$

where  $\tilde{\boldsymbol{\zeta}}$  and  $\tilde{\nu}$  are the dual parameters, and  $\odot$  denotes the Hadamard product.

Here, let  $\mathbf{t}_{\mathcal{I}}$ ,  $\mathbf{t}_{\mathcal{J}}$ ,  $\boldsymbol{\xi}_{\mathcal{I}}$ , and  $\boldsymbol{\xi}_{\mathcal{J}}$  be sub-vectors of  $\mathbf{t}, \boldsymbol{\xi} \in \mathbb{R}^d$  indexed by  $\mathcal{I}$  and  $\mathcal{J}$ . We define vectors  $\mathbf{t}$  and  $\boldsymbol{\xi}$  by  $\mathbf{t}_{\mathcal{I}} = \hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}^{-1} \mathbf{1}_k$ ,  $\mathbf{t}_{\mathcal{J}} = \mathbf{0}_{d-k}$ ,  $\boldsymbol{\xi}_{\mathcal{I}} = \mathbf{0}_k$ ,  $\boldsymbol{\xi}_{\mathcal{J}} = \hat{A}_{\mu, \mathcal{J} \times \mathcal{I}} \hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}^{-1} \mathbf{1}_k - \mathbf{1}_{d-k}$ . We also define a scalar  $\gamma$  by  $\gamma = 1/\mathbf{1}_k^\top \hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}^{-1} \mathbf{1}_k$ . We show that setting  $\tilde{\mathbf{s}} \leftarrow \gamma \mathbf{t}$ ,  $\tilde{\boldsymbol{\zeta}} \leftarrow \gamma \boldsymbol{\xi}$ , and  $\tilde{\nu} \leftarrow \gamma$  satisfies the conditions (7)–(10). From the definitions of  $\mathbf{t}, \boldsymbol{\xi}$ , and  $\gamma$ , the conditions (7), (9), and (10) are obvious, and there remains only the condition (8) to be verified. The condition  $\tilde{\mathbf{s}} \geq \mathbf{0}_d$  is guaranteed from the condition (5). By combining (5) and (6), we also have

$$\boldsymbol{\xi}_{\mathcal{J}} \geq \tau \hat{A}_{\mu, \mathcal{J} \times \mathcal{I}} \mathbf{1}_k - \mathbf{1}_{d-k} \geq \tau' \mathbf{1}_{d-k} > \mathbf{0}_{d-k},$$

which guarantees  $\tilde{\boldsymbol{\zeta}} \geq \mathbf{0}_d$ . This result indicates the vector  $\tilde{\mathbf{s}}$  defined here is the optimal solution to the problem (3) with  $\tilde{\mathcal{I}}_0 = \mathcal{I}$ , which completes the proof.  $\square$

(Proof of Lemma 5) Let  $\epsilon := \|\hat{A}_\mu - A_\mu\|_\infty$ . Since  $B_\delta \leq 1/2$  for any  $\delta > 0$ , we have  $\epsilon < \mu/2d$ .

We first consider the condition (5). We note that, from Assumption 1,  $\Gamma_{\mathcal{I} \times \mathcal{I}} = 0_{|\mathcal{I}| \times |\mathcal{I}|}$  holds and thus  $A_{\mu, \mathcal{I} \times \mathcal{I}} = \mu I_k$ . Hence,  $\hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}$  can be expressed as  $\hat{A}_{\mu, \mathcal{I} \times \mathcal{I}} = \mu I_k + E$  with  $\|E\|_\infty \leq \epsilon$ . From Ravikumar et al. [10, pp. 972], we have

$$(\mu I_k + E)^{-1} \mathbf{1}_k = \frac{1}{\mu} \mathbf{1}_k - \frac{1}{\mu^2} E \mathbf{1}_k + \frac{1}{\mu^3} E^2 J \mathbf{1}_k,$$

where  $J = \sum_{m=0}^{\infty} (-1)^m (\mu^{-1} E)^m$ . From  $\|E\|_\infty \leq \epsilon$ , we have  $\|E^2\|_{\infty, \infty} \leq d^2 \epsilon^2$  and  $\|J\|_{\infty, \infty} \leq \sum_{m=0}^{\infty} \|\mu^{-1} E\|_{\infty, \infty}^m \leq \sum_{m=0}^{\infty} (\mu^{-1} d \epsilon)^m \leq 2$  where we used  $\epsilon < \mu/2d$  for the last inequality. Hence, we have

$$\|E^2 J \mathbf{1}_k\|_\infty \leq \|E^2 J\|_{\infty, \infty} \leq \|E^2\|_{\infty, \infty} \|J\|_{\infty, \infty} \leq 2d^2 \epsilon^2.$$

Thus, we can conclude

$$(\hat{A}_{\mu, \mathcal{I} \times \mathcal{I}}^{-1} \mathbf{1}_k)_i \geq \frac{1}{\mu} - \frac{d}{\mu^2} \epsilon - \frac{2d^2}{\mu^3} \epsilon^2,$$

where  $(\cdot)_i$  denotes the  $i$ th entry of the vector in the parenthesis. In addition, the right hand side of this inequality becomes positive when  $\epsilon < \mu/2d$ , which is assured by the assumption. This indicates that we can

choose  $\tau = 1/\mu - d\epsilon/\mu^2 - 2d^2\epsilon^2/\mu^3$  and the condition (5) holds.

We now turn to the condition (6). We have

$$(\hat{A}_{\mu, \mathcal{J} \times \mathcal{I}} \mathbf{1}_k)_i \geq (A_{\mu, \mathcal{J} \times \mathcal{I}} \mathbf{1}_k)_i - k\epsilon \geq (1 + \delta)\mu - k\epsilon,$$

from the assumption in Theorem 2, and we can choose  $\tau'$  as  $\tau' = \tau\{(1 + \delta)\mu - k\epsilon\} - 1$ . The condition  $\tau' > 0$  is assured when

$$\begin{aligned} \tau' &= \delta - \frac{k}{\mu} \epsilon - \left( \frac{d}{\mu} \epsilon + 2 \frac{d^2}{\mu^2} \epsilon^2 \right) \left( 1 + \delta - \frac{k}{\mu} \epsilon \right) \\ &\geq \delta - (2 + \delta) \frac{d}{\mu} \epsilon - 2(1 + \delta) \frac{d^2}{\mu^2} \epsilon^2 > 0, \end{aligned}$$

which is guaranteed by the assumption  $\epsilon < B_\delta \mu/d$ , and therefore the condition (6) holds.  $\square$

#### 7.4 Proof of Theorem 3

From Lemma 2, we have  $\|\hat{\Gamma} - \Gamma\|_\infty < h/(k^2 + d^2)$  when  $\eta > 4d \exp\{-h^2 n/2C_1^2(k^2 + d^2)^2\}$ , which proves the claim.  $\square$

#### 7.5 Proof of Theorem 4

From Lemmas 2 and 3, with probability greater than  $1 - \eta$ , we have

$$\begin{aligned} \|\hat{A}_\mu - A_\mu\|_\infty &\leq \|\hat{\Gamma} - \Gamma\|_\infty + |\lambda_{\min}(\hat{\Gamma}) - \lambda_{\min}(\Gamma)| \\ &\leq C_1 \left( 1 + \sqrt{d \min\{s_A + s_B, d\}} \right) \\ &\quad \times \sqrt{\frac{2 \log(2d/\eta) + \log 4}{n}} \\ &= C_2 \sqrt{\frac{2 \log(2d/\eta) + \log 4}{n}}, \end{aligned}$$

where  $C_2 := C_1 \left( 1 + \sqrt{d \min\{s_A + s_B, d\}} \right)$ .

From Lemma 2, we have  $\|\hat{A}_\mu - A_\mu\|_\infty < B_\delta \mu/d$  when  $\eta > 4d \exp\{-B_\delta^2 \mu^2 n/2C_2^2 d^2\}$ , which proves the claim.  $\square$

## References

- [24] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, Distributed optimization and statistical learning via the alternating direction method of multipliers. *Foundations and Trends in Machine Learning* 3(1):1122, 2011.
- [25] R. A. Horn, and C. R. Johnson, *Matrix Analysis*. Cambridge University Press, New York, 1985.