

## A Proof of Main Theorems

In this section, we provide the proof for our main theories.

We start by defining some notations. Note that the estimator in (3.7) can be rewritten as

$$\hat{\Delta} = \underset{\Delta \in \mathbb{R}^{d \times d}}{\operatorname{argmin}} \ell(\Delta) + \mathcal{G}_\lambda(\Delta), \quad (\text{A.1})$$

where  $\ell(\Delta) = 1/2 \operatorname{tr}(\Delta \hat{\Sigma}_Y \Delta \hat{\Sigma}_X) - \operatorname{tr}(\Delta(\hat{\Sigma}_X - \hat{\Sigma}_Y))$ ,  $\mathcal{G}_\lambda(\Delta)$  is the nonconvex penalty defined in Section 3 and  $\lambda$  is a non-negative regularization parameter. By the definition and decomposition of nonconvex penalty in Section 3, we can written the estimator as

$$\hat{\Delta} = \underset{\Delta \in \mathbb{R}^{d \times d}}{\operatorname{argmin}} \tilde{\ell}_\lambda(\Delta) + \lambda \|\Delta\|_{1,1}, \quad (\text{A.2})$$

where  $\tilde{\ell}_\lambda(\Delta) = \ell(\Delta) + \mathcal{H}_\lambda(\Delta)$ , and  $\mathcal{H}_\lambda(\Delta) = \sum_{j,k=1}^d h_\lambda(\Delta_{jk})$  is the concave part of  $\mathcal{G}(\Delta)$ .

To simplify the proof, we further make some transformations on the notations. By some linear algebra identities [13], we have  $\operatorname{tr}(\mathbf{A}^\top \mathbf{B}) = \operatorname{vec}(\mathbf{A})^\top \operatorname{vec}(\mathbf{B})$  and  $\operatorname{tr}(\mathbf{A}^\top \mathbf{B} \mathbf{C} \mathbf{D}^\top) = \operatorname{vec}(\mathbf{A})^\top (\mathbf{D} \otimes \mathbf{B}) \operatorname{vec}(\mathbf{C})$  for any matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}$  and  $\mathbf{D}$  with commensurate dimensions. Using these identities, we can rewrite the quasi log likelihood in (3.6) as

$$\mathcal{L}(\beta) = \frac{1}{2} \beta^\top \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}^\top \beta, \quad (\text{A.3})$$

where  $\beta = \operatorname{vec}(\Delta) \in \mathbb{R}^{d^2}$ ,  $\hat{\mathbf{Q}} = \hat{\Sigma}_X \otimes \hat{\Sigma}_Y \in \mathbb{R}^{d^2 \times d^2}$  and  $\hat{\mathbf{b}} = \operatorname{vec}(\hat{\Sigma}_X - \hat{\Sigma}_Y) \in \mathbb{R}^{d^2}$ . Then the estimator in (A.1) can be rewritten as

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{d^2}}{\operatorname{argmin}} \mathcal{L}(\beta) + \mathcal{G}_\lambda(\beta), \quad (\text{A.4})$$

where  $\mathcal{L}(\beta) = 1/2 \beta^\top \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}^\top \beta$  is the counterpart of loss function  $\ell(\Delta) = 1/2 \operatorname{tr}(\Delta \hat{\Sigma}_Y \Delta \hat{\Sigma}_X) - \operatorname{tr}(\Delta(\hat{\Sigma}_X - \hat{\Sigma}_Y))$ ,  $\mathcal{G}_\lambda(\beta) = \sum_{i=1}^{d^2} g_\lambda(\beta_i)$  is the nonconvex penalty defined in Section 3 and  $\lambda$  is a non-negative regularization parameter. Therefore, the optimization problem in (A.2) turns to be

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^{d^2}}{\operatorname{argmin}} \tilde{\mathcal{L}}_\lambda(\beta) + \lambda \|\beta\|_1, \quad (\text{A.5})$$

where  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{H}_\lambda(\beta)$ , and  $\mathcal{H}_\lambda(\beta) = \sum_{i=1}^{d^2} h_\lambda(\beta_i)$  is the concave part of  $\mathcal{G}(\beta)$ .

Denote  $\operatorname{vec}(S) := \operatorname{supp}(\beta^*)$ , where  $\beta^* = \operatorname{vec}(\Delta^*)$  and  $S = \operatorname{supp}(\Delta^*)$  is the support of the true differential graph. Finally, the vectorized oracle estimator of  $\beta^*$  in (4.1) turns to be

$$\hat{\beta}_O = \underset{\operatorname{supp}(\beta) \subseteq \operatorname{vec}(S)}{\operatorname{argmin}} \mathcal{L}(\beta), \quad (\text{A.6})$$

where  $\mathcal{L}(\beta) = \frac{1}{2} \beta^\top \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}^\top \beta$ .

Now, we are ready to prove our main results. In order to make the proof concise, we first prove Theorem 4.6, followed which we prove Theorem 4.4. Note that the proof of Theorem 4.4 relies on the proof of Theorem 4.6.

*Proof of Theorem 4.6.* Suppose  $\hat{\mathbf{z}} \in \partial \|\hat{\beta}\|_1$ . In particular, the estimator  $\hat{\beta}$  in (A.5) satisfies optimality condition for unconstrained problem

$$\langle \hat{\beta} - \beta', \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}) + \lambda \hat{\mathbf{z}} \rangle \leq 0, \quad (\text{A.7})$$

for any  $\beta'$ .

First, we want to show that there exists some  $\hat{\mathbf{z}}_O \in \partial \|\hat{\beta}_O\|_1$ , such that  $\hat{\mathbf{z}}_O$  satisfies the optimality condition as follows

$$\langle \hat{\beta}_O - \beta', \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O \rangle \leq 0, \quad (\text{A.8})$$

for any  $\beta'$ . Since  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{H}_\lambda(\beta)$ , we have

$$\begin{aligned} \langle \hat{\beta}_O - \beta', \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O \rangle &= \underbrace{\sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i \cdot (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i}_{(i)} \\ &+ \underbrace{\sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i}_{(ii)}. \end{aligned} \quad (\text{A.9})$$

For term (i) in (A.9), by Lemma B.3, we have with probability at least  $1 - 3/s$  that

$$\|\hat{\beta}_O - \beta^*\|_\infty \leq C \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\frac{\log s}{n}},$$

where  $C$  is an absolute constant. Recall the assumption on entry magnitude of  $\beta^*$ , i.e.,  $\min_{i \in \text{vec}(S)} |\beta_i^*| \geq \nu + C \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\log s/n}$ , we have with probability at least  $1 - 3/s$  that

$$\begin{aligned} \min_{i \in \text{vec}(S)} |(\hat{\beta}_O)_i| &= \min_{i \in \text{vec}(S)} |(\hat{\beta}_O - \beta^* + \beta^*)_i| \geq \min_{i \in \text{vec}(S)} \{ |(\beta^*)_i| - |(\hat{\beta}_O - \beta^*)_i| \} \\ &\geq - \max_{i \in \text{vec}(S)} |(\hat{\beta}_O - \beta^*)_i| + \min_{i \in \text{vec}(S)} |(\beta^*)_i|. \end{aligned} \quad (\text{A.10})$$

The right hand side of (A.10) can be further lower bounded by

$$\min_{i \in S} |(\hat{\beta}_O)_i| \geq -C \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\frac{\log s}{n}} + \nu + C \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\frac{\log s}{n}}.$$

Following condition (a) in Assumption 4.3 for  $\mathcal{G}(\beta)$ , we have

$$(\nabla \mathcal{H}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i = (\nabla \mathcal{G}(\hat{\beta}_O))_i = g'_\lambda((\hat{\beta}_O)_i) = 0,$$

for  $i \in \text{vec}(S)$ . Hence we have

$$\begin{aligned} \sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i &= \sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i \cdot (\nabla \mathcal{L}(\hat{\beta}_O) + \nabla \mathcal{H}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i, \\ &= \sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i \cdot (\nabla \mathcal{L}(\hat{\beta}_O))_i. \end{aligned}$$

Recall that  $\hat{\beta}_O$  is the global solution to the problem in (A.6). Hence we have  $\hat{\beta}_O$  satisfies the optimality condition as follows

$$\sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i (\nabla \mathcal{L}(\hat{\beta}_O))_i \leq 0,$$

which leads to

$$\sum_{i \in \text{vec}(S)} (\hat{\beta}_O - \beta')_i (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i \leq 0. \quad (\text{A.11})$$

For term (ii) in (A.9), notice that  $(\hat{\beta}_O)_i = 0$  for  $i \in \text{vec}(S)^c$ . By the regularity condition (c), we have

$$(\nabla \mathcal{H}_\lambda(\hat{\beta}_O))_i = h'_\lambda((\hat{\beta}_O)_i) = 0,$$

for  $i \in \text{vec}(S)^c$ . This leads to

$$\begin{aligned} \sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i &= \sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \mathcal{L}(\hat{\beta}_O) + \nabla \mathcal{H}_\lambda(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i, \\ &= \sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \mathcal{L}(\hat{\beta}_O) + \lambda \hat{\mathbf{z}}_O)_i. \end{aligned}$$

Since  $\nabla \mathcal{L}(\beta) = \hat{\mathbf{Q}}\beta - \hat{\mathbf{b}}$  and note that  $\mathbf{Q}^*\beta^* = \mathbf{b}^*$ , we have

$$\begin{aligned} \|\nabla \mathcal{L}(\hat{\beta}_O)\|_\infty &= \|\hat{\mathbf{Q}}\hat{\beta}_O - \hat{\mathbf{Q}}\beta^* + \hat{\mathbf{Q}}\beta^* - \mathbf{Q}^*\beta^* + \mathbf{b}^* - \hat{\mathbf{b}}\|_\infty \\ &\leq \|\hat{\mathbf{Q}}\|_1 \cdot \|\hat{\beta}_O - \beta^*\|_\infty + \|\hat{\mathbf{Q}} - \mathbf{Q}^*\|_{\infty, \infty} \cdot \|\beta^*\|_1 + \|\mathbf{b}^* - \hat{\mathbf{b}}\|_\infty \\ &\leq \|\hat{\mathbf{Q}}\|_1 \cdot C\theta_X^2\theta_Y^2\sigma_X\sigma_Y M \sqrt{\frac{\log s}{n}} + \|\beta^*\|_1 \cdot \sqrt{5}\pi \sqrt{\frac{\log d}{n}} + 6\pi \sqrt{\frac{\log d}{n}}, \end{aligned}$$

where in the last inequality the first term is due to Lemma B.3, the second term is from (B.3), and the last term is due to Lemma C.1 and (B.6). In addition, we have  $\|\hat{\mathbf{Q}}\|_1 \leq \|\hat{\Sigma}_X\|_1 \cdot \|\hat{\Sigma}_Y\|_1 \leq 4\|\Sigma_X^*\|_1 \cdot \|\Sigma_Y^*\|_1$  when  $n$  is sufficient large, and thus  $\|\hat{\mathbf{Q}}\|_1 \leq 4\sigma_X\sigma_Y$  by Assumption 4.1. By Assumption 4.2, we have  $\|\beta^*\|_1 \leq M$ . Therefore, for any  $i \in \text{vec}(S)^c$ , we obtain

$$|(\nabla \mathcal{L}(\hat{\beta}_O))_i| \leq \|\nabla \mathcal{L}(\hat{\beta}_O)\|_\infty \leq C_0\theta_X^2\theta_Y^2\sigma_X\sigma_Y M \sqrt{\frac{\log d}{n}},$$

where  $C_0$  is an absolute constant. By Assumption , it follows that  $|(\nabla \mathcal{L}(\hat{\beta}_O))_i| \leq \lambda/2$  for any  $i \in \text{vec}(S)^c$ . Since we have  $\hat{\mathbf{z}}_O \in \partial\|\hat{\beta}_O\|_1$ , hence  $|(\hat{\mathbf{z}}_O)_i| \leq 1$  for  $i \in \text{vec}(S)^c$ . By setting  $(\hat{\mathbf{z}}_O)_i = -(\nabla \mathcal{L}(\hat{\beta}_O))_i/\lambda$  for  $i \in \text{vec}(S)^c$ , we can enforce the following equality to hold

$$(\nabla \mathcal{L}(\hat{\beta}_O) + \lambda\hat{\mathbf{z}}_O)_i = 0,$$

for  $i \in \text{vec}(S)^c$ . Hence, we have

$$\sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda\hat{\mathbf{z}}_O)_i = \sum_{i \in \text{vec}(S)^c} (\hat{\beta}_O - \beta')_i \cdot (\nabla \mathcal{L}(\hat{\beta}_O) + \lambda\hat{\mathbf{z}}_O)_i = 0. \quad (\text{A.12})$$

By using (A.11) and (A.12), we obtain (A.8).

Now we are ready to provide proof on  $\hat{\beta} = \hat{\beta}_O$ . Recall that  $\text{supp}(\hat{\beta}_O) = \text{vec}(S)$ , and Lemma B.2 shows that under suitable condition, we have

$$\tilde{\mathcal{L}}_\lambda(\hat{\beta}) \geq \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \langle \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O), \hat{\beta} - \hat{\beta}_O \rangle + \frac{\rho - \zeta_-}{2} \|\hat{\beta} - \hat{\beta}_O\|_2^2, \quad (\text{A.13})$$

$$\tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) \geq \tilde{\mathcal{L}}_\lambda(\hat{\beta}) + \langle \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}), \hat{\beta}_O - \hat{\beta} \rangle + \frac{\rho - \zeta_-}{2} \|\hat{\beta}_O - \hat{\beta}\|_2^2, \quad (\text{A.14})$$

hold with high probability.

By convexity of  $\ell_1$  norm  $\|\cdot\|_1$ , we have following two inequality hold

$$\lambda\|\hat{\beta}\|_1 \geq \lambda\|\hat{\beta}_O\|_1 + \lambda\langle \hat{\beta} - \hat{\beta}_O, \hat{\mathbf{z}}_O \rangle, \quad (\text{A.15})$$

$$\lambda\|\hat{\beta}_O\|_1 \geq \lambda\|\hat{\beta}\|_1 + \lambda\langle \hat{\beta}_O - \hat{\beta}, \hat{\mathbf{z}} \rangle. \quad (\text{A.16})$$

By adding Equations (A.13)-(A.16), we have

$$0 \geq \underbrace{\langle \hat{\beta}_O - \hat{\beta}, \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}) + \lambda\hat{\mathbf{z}} \rangle}_{(a)} + \underbrace{\langle \hat{\beta} - \hat{\beta}_O, \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda\hat{\mathbf{z}}_O \rangle}_{(b)} + (\rho - \zeta_-) \|\hat{\beta} - \hat{\beta}_O\|_2^2.$$

Recall that  $\hat{\beta}$  satisfies the optimality condition

$$\langle \hat{\beta} - \hat{\beta}_O, \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}) + \lambda\hat{\mathbf{z}} \rangle \leq 0,$$

hence we have term (a)  $\geq 0$ .

Similarly, by (A.8), we have

$$\langle \hat{\beta}_O - \hat{\beta}, \nabla \tilde{\mathcal{L}}_\lambda(\hat{\beta}_O) + \lambda\hat{\mathbf{z}}_O \rangle \leq 0,$$

which leads to term (b)  $\geq 0$ . Therefore, we have  $(\rho - \zeta_-) \|\hat{\beta} - \hat{\beta}_O\|_2^2 \leq 0$ , which implies  $\hat{\beta} = \hat{\beta}_O$ . Thus we can conclude that, under suitable condition, the proposed estimator  $\hat{\beta}$  is the oracle estimator  $\hat{\beta}_O$ , which exactly recover the true support of  $\beta^*$  with probability at least  $1 - 3/s$ .  $\square$

Next, we are able to prove Theorem 4.4.

*Proof of Theorem 4.4.* Recall that in Theorem 4.6, we have proved that under certain conditions  $\hat{\beta} = \hat{\beta}_O$  holds. Then by Lemma B.3 we have

$$\|\hat{\beta} - \beta^*\|_\infty = \|\hat{\beta}_O - \beta^*\|_\infty \leq 6\pi\theta_X\theta_Y\sqrt{\frac{\log s}{n}} + 2\sqrt{10}\pi\theta_X^2\theta_Y^2\sigma_X\sigma_Y M\sqrt{\frac{\log s}{n}}$$

holds with probability at least  $1 - 3/s$ , where the second term dominates the first one.

Next, we just need to bound  $\|\hat{\beta}_O - \beta^*\|_2$ . By definition in (A.6), we have

$$\hat{\beta}_O = \underset{\text{supp}(\beta) \subseteq \text{vec}(S)}{\text{argmin}} \frac{1}{2} \beta^\top \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}^\top \beta.$$

By definition we have  $\hat{\mathbf{Q}} = \hat{\Sigma}_X \otimes \hat{\Sigma}_Y \in \mathbb{R}^{d^2 \times d^2}$ . For any  $j, k, p, q = 1, \dots, d$ , we use  $\hat{\mathbf{Q}}_{(j,k,p,q)}$  to denote the entry in  $\hat{\mathbf{Q}}$  that is obtained from the product of the  $(j, k)$ -th entry in  $\hat{\Sigma}_X$  and the  $(p, q)$ -th entry in  $\hat{\Sigma}_Y$ . Specifically, we have

$$\begin{aligned} \hat{\mathbf{Q}}_{(j,k,p,q)} &= \hat{\Sigma}_X^{jk} \hat{\Sigma}_Y^{pq} = \sin\left(\frac{\pi}{2} \hat{\tau}_{jk}^X\right) \sin\left(\frac{\pi}{2} \hat{\tau}_{pq}^Y\right) \\ &= \frac{1}{2} \cos\left(\frac{\pi}{2} (\hat{\tau}_{jk}^X - \hat{\tau}_{pq}^Y)\right) - \frac{1}{2} \cos\left(\frac{\pi}{2} (\hat{\tau}_{jk}^X + \hat{\tau}_{pq}^Y)\right). \end{aligned}$$

Furthermore, we define  $\hat{\mu}_{jk;pq} = \hat{\tau}_{jk}^X - \hat{\tau}_{pq}^Y$  and  $\hat{\mu}'_{jk;pq} = \hat{\tau}_{jk}^X + \hat{\tau}_{pq}^Y$ . All the notations above can be easily extended to  $\mathbf{Q}^*$ . Then we have

$$\begin{aligned} \hat{\mathbf{Q}}_{(j,k,p,q)} - \mathbf{Q}_{(j,k,p,q)}^* &= \frac{1}{2} \left( \cos\left(\frac{\pi}{2} \hat{\mu}_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) \right) \\ &\quad + \frac{1}{2} \left( \cos\left(\frac{\pi}{2} \hat{\mu}'_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}'^*\right) \right). \end{aligned}$$

We only need to bound the first term, and the second term is very similar and the bound should be exactly the same. Note that

$$\begin{aligned} \cos\left(\frac{\pi}{2} \hat{\mu}_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) &= -\frac{\pi}{2} \sin\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) (\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*) \\ &\quad - \frac{\pi^2}{8} \cos\left(\frac{\pi}{2} \tilde{\mu}_{jk;pq}\right) (\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*)^2, \end{aligned}$$

where  $\tilde{\mu}_{jk;pq}$  lies between  $\hat{\mu}_{jk;pq}$  and  $\mu_{jk;pq}^*$ . Let  $\hat{\mathbf{L}} \in \mathbb{R}^{d^2 \times d^2}$  be the matrix with the same structure as  $\hat{\mathbf{Q}}$  whose  $(j, k, p, q)$ -th entry is  $\cos(\pi/2 \hat{\mu}_{jk;pq})$ . Similar notations are defined for  $\mathbf{L}^*$  and  $\tilde{\mathbf{L}}$ . Then for any  $\mathbf{x} \in \mathbb{R}^{d^2}$  we have

$$|\mathbf{x}^\top (\hat{\mathbf{Q}} - \mathbf{Q}^*) \mathbf{x}| \leq \frac{\pi}{2} \left| \mathbf{x}^\top \left[ \sin\left(\frac{\pi}{2} \hat{\mathbf{L}}\right) \circ (\hat{\mathbf{L}} - \mathbf{L}^*) \right] \mathbf{x} \right| + \frac{\pi^2}{8} \left| \mathbf{x}^\top \left[ \cos\left(\frac{\pi}{2} \tilde{\mathbf{L}}\right) \circ (\hat{\mathbf{L}} - \mathbf{L}^*) \circ (\hat{\mathbf{L}} - \mathbf{L}^*) \right] \mathbf{x} \right|.$$

Recall the results in Lemma C.2 and Lemma C.3, following a similar proof we can show that with probability at least  $1 - 1/s$

$$\sup_{\|\mathbf{x}\| \leq 1} |\mathbf{x}^\top (\hat{\mathbf{Q}}_{SS} - \mathbf{Q}_{SS}^*) \mathbf{x}| \leq 2\pi^2 \frac{s \log s}{n} + 8\pi\sqrt{C} \sqrt{\frac{\log s + s \log 9}{n}}, \quad (\text{A.17})$$

where  $C$  is an absolute constant. We have the closed form solution  $\beta^* = \mathbf{Q}^{*-1} \mathbf{b}^*$ . It follows that

$$\begin{aligned} \|\hat{\beta}_O - \beta^*\|_2 &= \|\hat{\mathbf{Q}}_{SS}^{-1} \hat{\mathbf{b}} - \mathbf{Q}^{*-1} \mathbf{b}^*\|_2 = \|\hat{\mathbf{Q}}_{SS}^{-1} \hat{\mathbf{b}} - \hat{\mathbf{Q}}_{SS}^{-1} \mathbf{b}^* + \hat{\mathbf{Q}}_{SS}^{-1} \mathbf{b}^* - \mathbf{Q}^{*-1} \mathbf{b}^*\|_2 \\ &\leq \underbrace{\|\hat{\mathbf{Q}}_{SS}^{-1}\|_2 \cdot \|\hat{\mathbf{b}} - \mathbf{b}^*\|_2}_{(i)} + \underbrace{\|\hat{\mathbf{Q}}_{SS}^{-1} - \mathbf{Q}_{SS}^{*-1}\|_2 \cdot \|\mathbf{b}^*\|_2}_{(ii)}, \end{aligned} \quad (\text{A.18})$$

where we use the fact that  $\text{vec}(S) = \text{supp}(\beta^*) = \text{supp}(\mathbf{Q}^{*-1} \mathbf{b}^*)$ . Note that  $\|\mathbf{b}^*\|_2 = \|\Sigma_X^* - \Sigma_Y^*\|_F = \|\Sigma_X^* \Delta^* \Sigma_Y^*\|_F \leq \|\Sigma_X^*\|_2 \cdot \|\Sigma_Y^*\|_2 \cdot \|\Delta^*\|_F \leq M/(\kappa_1 \kappa_2)$ , here we used the fact that

$\|\Delta^*\|_F \leq \|\Delta^*\|_{1,1} \leq M$  and  $\lambda_{\max}(\Sigma_X^*) \leq 1/\kappa_1$  by Assumption 4.1 and 4.2. Then term (ii) in (A.18) can be bounded as

$$\begin{aligned} \|\widehat{\mathbf{Q}}_{SS}^{-1} - \mathbf{Q}_{SS}^{*-1}\|_2 \cdot \|\mathbf{b}_S^*\|_2 &\leq \frac{M}{\kappa_1 \kappa_2} \cdot \|\widehat{\mathbf{Q}}_{SS}^{-1}\|_2 \cdot \|\widehat{\mathbf{Q}}_{SS} - \mathbf{Q}_{SS}^*\|_2 \cdot \|\mathbf{Q}_{SS}^{*-1}\|_2 \\ &\leq \frac{4\pi^2 M}{\kappa_1^3 \kappa_2^3} \frac{s \log s}{n} + \frac{16\pi M C}{\kappa_1^3 \kappa_2^3} \sqrt{\frac{\log s + s \log 9}{n}}, \end{aligned} \quad (\text{A.19})$$

where the second inequality uses the bound in (A.17) and the fact that  $\|\mathbf{Q}_{SS}^{*-1}\|_2 \leq 1/(\kappa_1 \kappa_2)$  by Assumption 4.1 and  $\|\widehat{\mathbf{Q}}_{SS}^{-1}\|_2 \leq 2\|\mathbf{Q}_{SS}^{*-1}\|_2$  when  $n$  is sufficient large.

For term (i) in (A.18), with probability at least  $1 - 1/d - 2/d^2$  we have  $\|\widehat{\Sigma}_X - \Sigma_X^*\|_2 \leq 2\pi^2 d \log d/n + 8\pi\sqrt{C}\sqrt{(\log d + d \log 9)/n}$  by Lemma C.3. The dominating term is  $\sqrt{d/n}$ . Similar bound for  $\widehat{\Sigma}_Y$  holds. It immediately implies

$$\|\widehat{\mathbf{Q}}_{SS}^{-1}\|_2 \cdot \|\widehat{\mathbf{b}} - \mathbf{b}^*\|_S \leq \frac{2}{\kappa_1 \kappa_2} (\|\widehat{\Sigma}_X - \Sigma_X^*\|_S + \|\widehat{\Sigma}_Y - \Sigma_Y^*\|_S) \leq \frac{C_1}{\kappa_1 \kappa_2} \sqrt{\frac{s}{n}}, \quad (\text{A.20})$$

where  $C_1$  is an absolute constant. Submitting (A.19) and (A.20) into (A.18), we obtain

$$\|\widehat{\beta}_O - \beta^*\|_2 \leq \frac{C_1}{\kappa_1 \kappa_2} \sqrt{\frac{s}{n}} + \frac{4\pi^2 M}{\kappa_1^3 \kappa_2^3} \frac{s \log s}{n} + \frac{64(1 + \sqrt{5})\pi M}{\kappa_1^3 \kappa_2^3} \sqrt{\frac{\log s + s \log 9}{n}},$$

which holds with probability at least  $1 - 1/d - 1/d^{2.5}$ . In Theorem 4.6, we have proved that under certain conditions  $\widehat{\beta} = \widehat{\beta}_O$  holds, which further implies that

$$\begin{aligned} \|\widehat{\beta} - \beta^*\|_2 &\leq \frac{C_1(\kappa_1^2 + \kappa_2^2)}{\kappa_1^3 \kappa_2^3} \sqrt{\frac{s}{n}} + \frac{4\pi^2 \sigma_X \sigma_Y M}{\kappa_1^3 \kappa_2^3} \frac{s \log s}{n} + \frac{64(1 + \sqrt{5})\pi M}{\kappa_1^3 \kappa_2^3} \sqrt{\frac{\log s + s \log 9}{n}} \\ &\leq \frac{C_2 M}{\kappa_1 \kappa_2} \sqrt{\frac{s}{n}} \end{aligned}$$

holds with probability at least  $1 - 1/s - 1/s^{2.5} \geq 1 - 2/s$ , where  $C_2$  is an absolute constant.  $\square$

Finally, we are ready to prove Theorem 4.7.

*Proof of Theorem 4.7.* Let  $\mathbf{z} \in \partial\|\beta\|_1$  and  $\widehat{\mathbf{z}} \in \partial\|\widehat{\beta}\|_1$  denote the subgradient. Recall that,  $\widehat{\beta}$  is the global solution to (A.5). Hence we have

$$\langle \widehat{\beta} - \beta', \nabla \widetilde{\mathcal{L}}_\lambda(\widehat{\beta}) + \lambda \widehat{\mathbf{z}} \rangle \leq 0, \quad (\text{A.21})$$

for any  $\beta'$ . By Lemma B.2, under suitable condition, with high probability, we have

$$\widetilde{\mathcal{L}}_\lambda(\widehat{\beta}) \geq \widetilde{\mathcal{L}}_\lambda(\beta^*) + \langle \nabla \widetilde{\mathcal{L}}_\lambda(\beta^*), \widehat{\beta} - \beta^* \rangle + \frac{\rho - \zeta_-}{2} \|\widehat{\beta} - \beta^*\|_2^2, \quad (\text{A.22})$$

$$\widetilde{\mathcal{L}}_\lambda(\beta^*) \geq \widetilde{\mathcal{L}}_\lambda(\widehat{\beta}) + \langle \nabla \widetilde{\mathcal{L}}_\lambda(\widehat{\beta}), \beta^* - \widehat{\beta} \rangle + \frac{\rho - \zeta_-}{2} \|\beta^* - \widehat{\beta}\|_2^2. \quad (\text{A.23})$$

By convexity of  $\ell_1$  norm  $\|\cdot\|_1$ , we have

$$\lambda \|\widehat{\beta}\|_1 \geq \lambda \|\beta^*\|_1 + \lambda \langle \widehat{\beta} - \beta^*, \mathbf{z}^* \rangle, \quad (\text{A.24})$$

$$\lambda \|\beta^*\|_1 \geq \lambda \|\widehat{\beta}\|_1 + \lambda \langle \beta^* - \widehat{\beta}, \widehat{\mathbf{z}} \rangle. \quad (\text{A.25})$$

Adding up (A.22) to (A.25), we have

$$0 \geq \langle \beta^* - \widehat{\beta}, \nabla \widetilde{\mathcal{L}}_\lambda(\widehat{\beta}) + \lambda \widehat{\mathbf{z}} \rangle + \langle \widehat{\beta} - \beta^*, \nabla \widetilde{\mathcal{L}}_\lambda(\beta^*) + \lambda \mathbf{z}^* \rangle + (\rho - \zeta_-) \|\widehat{\beta} - \beta^*\|_2^2.$$

Meanwhile, (A.21) leads to

$$\langle \nabla \widetilde{\mathcal{L}}_\lambda(\widehat{\beta}) + \lambda \widehat{\mathbf{z}}, \beta^* - \widehat{\beta} \rangle \geq 0.$$

Hence, we have

$$(\rho - \zeta_-) \|\hat{\beta} - \beta^*\|_2^2 \leq \langle \beta^* - \hat{\beta}, \nabla \tilde{\mathcal{L}}_\lambda(\beta^*) + \lambda \mathbf{z}^* \rangle. \quad (\text{A.26})$$

Recall that  $\tilde{\mathcal{L}}_\lambda(\beta)$  is restricted strongly convex provided  $\rho = \kappa_1 \kappa_2 / 2$ . With  $\zeta_- \leq \kappa_1 \kappa_2 / 4$  and (A.26), we have

$$\begin{aligned} \frac{\kappa_1 \kappa_2}{4} \|\hat{\beta} - \beta^*\|_2^2 &\leq (\rho - \zeta_-) \|\hat{\beta} - \beta^*\|_2^2 \\ &\leq \langle \nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*, \beta^* - \hat{\beta} \rangle \\ &\leq \sum_{i=1}^{d^2} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i|. \end{aligned} \quad (\text{A.27})$$

Now, we decompose (A.27) into three parts:  $i \in \text{vec}(S)^c$ ,  $i \in S_1$  and  $i \in S_2$ , where we define  $S_1 = \{i \mid |(\beta^*)_i| \geq \nu\}$  and  $S_2 = \{i \mid |(\beta^*)_i| < \nu\}$ .

**Case 1:** For  $i \in \text{vec}(S)^c$ , based on regularity condition (c) in Assumption 4.3, we have

$$(\nabla \mathcal{H}_\lambda(\beta^*))_i = h'_\lambda(\beta_i^*) = h'_\lambda(0) = 0.$$

Recall that we have  $|(\nabla \mathcal{L}(\beta^*))_i| \leq CM \sqrt{\log d/n} = \lambda/2$  according to Lemma B.4. Hence,

$$|(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*))_i| \leq \frac{\lambda}{2}.$$

Since  $\mathbf{z}^* \in \partial \|\beta^*\|_1$ , we have  $|z_i^*| \leq 1$  and thus  $\lambda z_i^* \in [-\lambda, \lambda]$ . Therefore, for any  $i \in \text{vec}(S)^c$ , by definition of subgradient of  $\mathbf{z}^*$  we can always find a  $z_i^*$  such that

$$|(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| = 0.$$

This leads to

$$\sum_{i \in \text{vec}(S)^c} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| = 0. \quad (\text{A.28})$$

**Case 2:** For  $i \in S_1$ , we have  $|\beta_i^*| \geq \nu$ . By condition (a) in Assumption 4.3 on  $\mathcal{G}(\beta) = \mathcal{H}_\lambda(\beta) + \lambda \|\beta\|_1$ , we have

$$(\nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i = g'_\lambda(\beta_i^*) = 0,$$

which implies

$$\sum_{i \in S_1} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| = \sum_{i \in S_1} |[\nabla \mathcal{L}(\beta^*)]_i| \cdot |(\beta^* - \hat{\beta})_i|.$$

Hence by Cauchy's inequality we have

$$\sum_{i \in S_1} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| \leq \|(\nabla \mathcal{L}(\beta^*))_{S_1}\|_2 \cdot \|(\beta^* - \hat{\beta})_{S_1}\|_2.$$

Since  $\nabla \mathcal{L}(\beta) = \hat{\mathbf{Q}}\beta - \hat{\mathbf{b}}$  and note that  $\mathbf{Q}^* \beta^* = \mathbf{b}^*$ , we have

$$\begin{aligned} \|[\nabla \mathcal{L}(\beta^*)]_{S_1}\|_2 &= \|[\hat{\mathbf{Q}}\beta^* - \mathbf{Q}^* \beta^* + \mathbf{b}^* - \hat{\mathbf{b}}]_{S_1}\|_2 \leq \|[\hat{\mathbf{Q}} - \mathbf{Q}^*]_{S_1 S_1}\|_2 \cdot \|[\beta^*]_{S_1}\|_2 + \|[\mathbf{b}^* - \hat{\mathbf{b}}]_{S_1}\|_2 \\ &\leq \sqrt{5\pi} M \sqrt{\frac{\log s_1}{n}} + 4\sqrt{3\pi} \sqrt{\frac{s_1}{n}} \end{aligned}$$

holds with probability at least  $1 - 2/s_1 - 1/s_1 = 1 - 3/s_1$ , where the first term in the last inequality is due to (B.3) and  $\|\beta^*\|_2 \leq \|\beta^*\|_1 \leq M$  by Assumption 4.2, and the second term is from Lemma C.4. Thus, we obtain

$$\sum_{i \in S_1} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| \leq 4\sqrt{3\pi} M \sqrt{\frac{s_1}{n}} \cdot \|\beta^* - \hat{\beta}\|_2. \quad (\text{A.29})$$

**Case 3:** For  $i \in S_2$ , we have  $|\beta_i^*| \leq \nu$ . By condition (d) in Assumption 4.3, we have

$$\max_{i \in S_2} |(\nabla \mathcal{H}_\lambda(\beta^*))_i| \leq \max_{i \in S_2} |h'_\lambda(\beta_i^*)| \leq \max_{1 \leq i \leq d^2} |h'_\lambda(\beta_i^*)| \leq \lambda.$$

Since  $\mathbf{z}^* \in \partial \|\beta^*\|_1$ , we have  $|z_i^*| \leq 1$ . Therefore, for  $i \in S_2$ , the following results hold

$$\begin{aligned} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| &\leq |(\nabla \mathcal{L}(\beta^*))_i| + |(\nabla \mathcal{H}_\lambda(\beta^*))_i| + \lambda |(\mathbf{z}^*)_i| \\ &\leq |(\nabla \mathcal{L}(\beta^*))_i| + 2\lambda. \end{aligned}$$

Again, by Lemma B.4  $\|\nabla \mathcal{L}(\beta^*)\|_\infty \leq \lambda/2$  holds with probability at least  $1 - 3/d$ , we obtain

$$\begin{aligned} \sum_{i \in S_2} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| &\leq (\|\nabla \mathcal{L}(\beta^*)\|_\infty + 2\lambda) \sum_{i \in S_2} |(\beta^* - \hat{\beta})_i| \\ &\leq \frac{5}{2} \lambda \sum_{i \in S_2} |(\beta^* - \hat{\beta})_i|. \end{aligned}$$

Hence we have

$$\sum_{i \in S_2} |(\nabla \mathcal{L}(\beta^*) + \nabla \mathcal{H}_\lambda(\beta^*) + \lambda \mathbf{z}^*)_i| \cdot |(\beta^* - \hat{\beta})_i| \leq \frac{5}{2} \lambda \sqrt{s_2} \|(\beta^* - \hat{\beta})_{S_2}\|_2 \leq \frac{5}{2} \lambda \sqrt{s_2} \|\beta^* - \hat{\beta}\|_2. \quad (\text{A.30})$$

Adding up (A.28) (A.29) and (A.30), and substituting the right term in (A.27), we obtain

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{16\sqrt{3}\pi M}{\kappa_1 \kappa_2} \sqrt{\frac{s_1}{n}} + \frac{10\sqrt{s_2}\lambda}{\kappa_1 \kappa_2} \leq \frac{16\sqrt{3}\pi M}{\kappa_1 \kappa_2} \sqrt{\frac{s_1}{n}} + \frac{10\pi M C}{\kappa_1 \kappa_2} \sqrt{\frac{s_2 \log d}{n}} \quad (\text{A.31})$$

holds with probability at least  $1 - 3/d - 3/s_1 \geq 1 - 6/s_1$ .  $\square$

## B Lemmas in the Proof of Main Theorems

**Lemma B.1.** Under Assumptions 4.1, the loss function  $\mathcal{L}(\beta) = 1/2 \beta^\top \hat{\mathbf{Q}} \beta - \beta^\top \hat{\mathbf{b}}$  is strongly convex with constant  $\kappa_1 \kappa_2 / 2$ .

*Proof of Lemma B.1.* Note that  $\nabla \mathcal{L}(\beta) = \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}$ , we have

$$\langle \nabla \mathcal{L}(\beta) - \nabla \mathcal{L}(\beta'), \beta - \beta' \rangle = (\beta - \beta')^\top \hat{\mathbf{Q}} (\beta - \beta').$$

Then we get

$$\begin{aligned} \min_{\beta, \beta' \in \text{vec}(S)} (\beta - \beta')^\top \hat{\mathbf{Q}} (\beta - \beta') &= \min_{\beta, \beta' \in \text{vec}(S)} (\beta - \beta')^\top (\hat{\mathbf{Q}} - \mathbf{Q}^* + \mathbf{Q}^*) (\beta - \beta') \\ &\geq \lambda_{\min}(\mathbf{Q}^*) \|\beta - \beta'\|_2^2 - \max_{\beta, \beta' \in \text{vec}(S)} (\beta - \beta')^\top (\hat{\mathbf{Q}} - \mathbf{Q}^*) (\beta - \beta') \\ &\geq \lambda_{\min}(\mathbf{Q}^*) \|\beta - \beta'\|_2^2 - \|\hat{\mathbf{Q}} - \mathbf{Q}^*\|_2 \cdot \|\beta - \beta'\|_2^2. \end{aligned}$$

By Assumption 4.1, we have  $\lambda_{\min}(\mathbf{Q}^*) = \lambda_{\min}(\Sigma_X^*) \lambda_{\min}(\Sigma_Y^*) = \kappa_1 \kappa_2$ . For the second term, we have

$$\begin{aligned} \|\hat{\mathbf{Q}} - \mathbf{Q}^*\|_2 &= \|\hat{\Sigma}_X \otimes \hat{\Sigma}_Y - \Sigma_X^* \otimes \Sigma_Y^* + \Sigma_X^* \otimes \hat{\Sigma}_Y + \Sigma_X^* \otimes \Sigma_Y^*\|_2 \\ &\leq \|\hat{\Sigma}_X - \Sigma_X^*\|_2 \cdot \|\hat{\Sigma}_Y - \Sigma_Y^*\|_2 + \|\Sigma_X^*\|_2 \cdot \|\hat{\Sigma}_Y - \Sigma_Y^*\|_2 + \|\Sigma_Y^*\|_2 \cdot \|\hat{\Sigma}_X - \Sigma_X^*\|_2 \\ &\leq \|\hat{\Sigma}_X - \Sigma_X^*\|_2 \cdot \|\hat{\Sigma}_Y - \Sigma_Y^*\|_2 + \frac{1}{\kappa_1} \cdot \|\hat{\Sigma}_Y - \Sigma_Y^*\|_2 + \frac{1}{\kappa_2} \cdot \|\hat{\Sigma}_X - \Sigma_X^*\|_2, \end{aligned}$$

where the second inequality is due to Assumption 4.1. By Lemma C.1, we have

$$\|\hat{\Sigma}_X - \Sigma_X^*\|_{\infty, \infty} \leq 3\pi \sqrt{\frac{\log d}{n}}, \quad \|\hat{\Sigma}_Y - \Sigma_Y^*\|_{\infty, \infty} \leq 3\pi \sqrt{\frac{\log d}{n}}$$

with probability at least  $1 - d^{-5/2}$ . Therefore, we have

$$\|(\hat{\mathbf{Q}} - \mathbf{Q}^*)_S\|_2 \leq 3\pi \sqrt{\frac{s \log d}{n}} \left( \frac{1}{\kappa_1} + \frac{1}{\kappa_2} \right) + \frac{9\pi^2 s \log d}{n}$$

with high probability. When  $n$  is sufficient large, we have

$$\min_{\beta, \beta' \in C} (\beta - \beta')^\top \hat{\mathbf{Q}} (\beta - \beta') \geq \frac{\kappa_1 \kappa_2}{2} \|\beta - \beta'\|_2^2.$$

This immediately implies  $\mathcal{L}(\beta)$  is restrictively strongly convex with constant  $\kappa_1 \kappa_2 / 2$ .  $\square$

Lemma (B.1) shows that, with high probability,  $\mathcal{L}(\beta)$  is a strongly convex function with modulus  $\rho = \kappa_1 \kappa_2 / 2 > 0$ . In (A.5) we defined  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{H}_\lambda(\beta)$ , where  $\mathcal{L}(\beta) = 1/2 \beta^\top \hat{\mathbf{Q}} \beta - \beta^\top \mathbf{b}$ ,  $\mathcal{H}_\lambda(\beta) = \sum_{i=1}^{d-1} h_\lambda(\beta_i) = \mathcal{G}_\lambda(\beta) - \lambda \|\beta\|_1$ . We now show that, with high probability,  $\tilde{\mathcal{L}}_\lambda(\beta)$  is strongly convex.

**Lemma B.2** (Restricted Strongly Convex). Let  $S = \text{supp}(\beta^*)$ . Given  $n \geq C_1 s \log d$  and appropriate parameter in nonconvex penalty  $\mathcal{G}_\lambda(\beta)$ ,  $\tilde{\mathcal{L}}_\lambda(\beta')$  is strongly convex.

$$\tilde{\mathcal{L}}_\lambda(\beta') \geq \tilde{\mathcal{L}}_\lambda(\beta^*) + \langle \nabla \tilde{\mathcal{L}}_\lambda(\beta^*), \beta' - \beta^* \rangle + \frac{\rho - \zeta_-}{2} \|\beta' - \beta^*\|_2^2,$$

holds with probability at least  $1 - C' \exp(-Cn)$ .

*Proof.* Recall that  $\mathcal{H}_\lambda(\beta)$  is the concave part of  $\tilde{\mathcal{L}}_\lambda(\beta)$ , which implies  $-\mathcal{H}_\lambda(\beta)$  is convex. Meanwhile, recall that  $\mathcal{H}_\lambda(\beta) = \sum_{i=1}^d h_\lambda(\beta_i)$ , where  $h_\lambda(\beta_i) + \zeta_- / 2 \beta_i^2$  is convex by Assumption 4.3. Hence we have

$$h_\lambda(\beta'_i) + \frac{\zeta_-}{2} \beta'^2_i \geq h_\lambda(\beta^*_i) + \frac{\zeta_-}{2} \beta^{*2}_i + (h'_\lambda(\beta^*_i) + \zeta_- \beta^*_i)(\beta'_i - \beta^*_i),$$

and

$$\mathcal{H}_\lambda(\beta') + \frac{\zeta_-}{2} \|\beta'\|_2^2 \geq \mathcal{H}_\lambda(\beta^*) + \frac{\zeta_-}{2} \|\beta^*\|_2^2 + \langle \nabla \mathcal{H}_\lambda(\beta^*) + \zeta_- \beta^*, \beta' - \beta^* \rangle.$$

This immediately implies

$$\mathcal{H}_\lambda(\beta') \geq \mathcal{H}_\lambda(\beta^*) + \langle \nabla \mathcal{H}_\lambda(\beta^*), \beta' - \beta^* \rangle - \frac{\zeta_-}{2} \|\beta' - \beta^*\|_2^2. \quad (\text{B.1})$$

Recall that by Lemma B.1, provided suitable condition,  $\mathcal{L}(\beta')$  is (w.h.p.) strongly convex. with modulus  $\rho = \kappa_1 \kappa_2 / 2$ , we have

$$\mathcal{L}(\beta') \geq \mathcal{L}(\beta^*) + \langle \nabla \mathcal{L}(\beta^*), \beta' - \beta^* \rangle + \frac{\rho}{2} \|\beta' - \beta^*\|_2^2. \quad (\text{B.2})$$

By the definition of  $\tilde{\mathcal{L}}_\lambda(\beta) = \mathcal{L}(\beta) + \mathcal{H}_\lambda(\beta)$ , adding (B.1) and (B.2) together, we obtain

$$\tilde{\mathcal{L}}_\lambda(\beta') \geq \tilde{\mathcal{L}}_\lambda(\beta^*) + \langle \nabla \tilde{\mathcal{L}}_\lambda(\beta^*), \beta' - \beta^* \rangle + \frac{\rho - \zeta_-}{2} \|\beta' - \beta^*\|_2^2,$$

holds with probability at least  $1 - C' \exp(-Cn)$ . Here,  $\rho = \kappa_1 \kappa_2 / 2$  and  $\zeta_-$  is depended on the nonconvex penalty. For example, in MCP penalty  $\zeta_- = 1/b$ . When  $\rho = \kappa_1 \kappa_2 / 2 > 1/b$ , the above equation leads to strongly convexity of  $\tilde{\mathcal{L}}_\lambda(\beta)$ , w.h.p. in the cone, provided suitable condition on  $n$ .  $\square$

**Lemma B.3.** Under Assumption 4.1, the oracle estimator  $\hat{\beta}_O$  in (A.6) satisfies

$$\|\hat{\beta}_O - \beta^*\|_\infty \leq 6\pi \theta_X \theta_Y \sqrt{\frac{\log s}{n}} + 2\sqrt{10} \pi \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\frac{\log s}{n}},$$

with probability at least  $1 - 3/s$ .

*Proof.* By definition in (A.6), we have

$$\hat{\beta}_O = \underset{\text{supp}(\beta) \subseteq \text{vec}(S)}{\text{argmin}} \frac{1}{2} \beta^\top \hat{\mathbf{Q}} \beta - \hat{\mathbf{b}}^\top \beta.$$

By definition we have  $\hat{\mathbf{Q}} = \hat{\Sigma}_X \otimes \hat{\Sigma}_Y \in \mathbb{R}^{d^2 \times d^2}$ . For any  $j, k, p, q = 1, \dots, d$ , we use  $\hat{\mathbf{Q}}_{(j,k,p,q)}$  to denote the entry in  $\hat{\mathbf{Q}}$  that is obtained from the product of the  $(j, k)$ -th entry in  $\hat{\Sigma}_X$  and the  $(p, q)$ -th entry in  $\hat{\Sigma}_Y$ . Specifically, we have

$$\begin{aligned} \hat{\mathbf{Q}}_{(j,k,p,q)} &= \hat{\Sigma}_X^{jk} \hat{\Sigma}_Y^{pq} = \sin\left(\frac{\pi}{2} \hat{\tau}_{jk}^X\right) \sin\left(\frac{\pi}{2} \hat{\tau}_{pq}^Y\right) \\ &= \frac{1}{2} \cos\left(\frac{\pi}{2} (\hat{\tau}_{jk}^X - \hat{\tau}_{pq}^Y)\right) - \frac{1}{2} \cos\left(\frac{\pi}{2} (\hat{\tau}_{jk}^X + \hat{\tau}_{pq}^Y)\right). \end{aligned}$$

Furthermore, we define  $\hat{\mu}_{jk;pq} = \hat{\tau}_{jk}^X - \hat{\tau}_{pq}^Y$  and  $\hat{\mu}'_{jk;pq} = \hat{\tau}_{jk}^X + \hat{\tau}_{pq}^Y$ . All the notations above can be easily extended to  $\mathbf{Q}^*$ . Then we have

$$\begin{aligned} \hat{\mathbf{Q}}_{(j,k,p,q)} - \mathbf{Q}_{(j,k,p,q)}^* &= \frac{1}{2} \left( \cos\left(\frac{\pi}{2} \hat{\mu}_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) \right) \\ &\quad + \frac{1}{2} \left( \cos\left(\frac{\pi}{2} \hat{\mu}'_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}'^*\right) \right). \end{aligned}$$

We only need to bound the first term, and the second term is very similar. Note that

$$\cos\left(\frac{\pi}{2} \hat{\mu}_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) = -\frac{\pi}{2} \sin\left(\frac{\pi}{2} \tilde{\mu}_{jk;pq}\right) (\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*),$$

where  $\tilde{\mu}_{jk;pq}$  lies between  $\hat{\mu}_{jk;pq}$  and  $\mu_{jk;pq}^*$ . To bound  $\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*$ , note that  $\hat{\tau}_{jk}^X, \hat{\tau}_{pq}^Y$  are sub-Gaussian random variables and  $|\hat{\tau}_{jk}^X|, |\hat{\tau}_{pq}^Y| \leq 1$ . Thus  $|\hat{\mu}_{jk;pq}| \leq 2$  and  $\hat{\mu}_{jk;pq}$  is also sub-Gaussian. In addition, we have  $\mathbb{E}(\hat{\mu}_{jk;pq}) = \mathbb{E}(\hat{\tau}_{jk}^X) - \mathbb{E}(\hat{\tau}_{pq}^Y) = \mu_{jk;pq}^*$ . Then by Hoeffding's inequality for U-statistics and applying union bound, we get

$$\mathbb{P}\left(\sup_{j,k,p,q} |\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*| > t\right) \leq 2d^4 e^{-\frac{nt^2}{4}}.$$

Take  $t = \sqrt{20 \log d/n}$ , we have that

$$\sup_{j,k,p,q} \left| \cos\left(\frac{\pi}{2} \hat{\mu}_{jk;pq}\right) - \cos\left(\frac{\pi}{2} \mu_{jk;pq}^*\right) \right| \leq \frac{\pi}{2} \sup_{j,k,p,q} |\hat{\mu}_{jk;pq} - \mu_{jk;pq}^*| \leq \frac{\pi}{2} \sqrt{\frac{20 \log d}{n}}$$

holds with probability at least  $1 - 2/d$ . It follows that

$$\|\hat{\mathbf{Q}} - \mathbf{Q}\|_{\infty, \infty} = \sup_{j,k,p,q} |\hat{\mathbf{Q}}_{(j,k,p,q)} - \mathbf{Q}_{(j,k,p,q)}^*| \leq \sqrt{5\pi} \sqrt{\frac{\log d}{n}} \quad (\text{B.3})$$

holds with probability at least  $1 - 2/d$ . We have the closed form solution of  $\hat{\beta}_O$  as  $\hat{\beta}_O = \hat{\mathbf{Q}}_{SS}^{-1} \hat{\mathbf{b}}$ . Then we have

$$\begin{aligned} \|\hat{\beta}_O - \beta^*\|_{\infty} &= \|\hat{\mathbf{Q}}_{SS}^{-1} \hat{\mathbf{b}} - \mathbf{Q}^{*-1} \mathbf{b}^*\|_{\infty} = \|\hat{\mathbf{Q}}_{SS}^{-1} \hat{\mathbf{b}} - \hat{\mathbf{Q}}_{SS}^{-1} \mathbf{b}^* + \hat{\mathbf{Q}}_{SS}^{-1} \mathbf{b}^* - \mathbf{Q}^{*-1} \mathbf{b}^*\|_{\infty} \\ &\leq \underbrace{\|\hat{\mathbf{Q}}_{SS}^{-1}\|_1 \cdot \|\hat{\mathbf{b}} - \mathbf{b}^*\|_{\infty}}_{(i)} + \underbrace{\|\hat{\mathbf{Q}}_{SS}^{-1} - \mathbf{Q}_{SS}^{*-1}\|_{\infty, \infty} \cdot \|\mathbf{b}^*\|_1}_{(ii)}. \end{aligned} \quad (\text{B.4})$$

For term (ii), we have

$$\begin{aligned} \|\hat{\mathbf{Q}}_{SS}^{-1} - \mathbf{Q}_{SS}^{*-1}\|_{\infty, \infty} &= \|\mathbf{Q}_{SS}^{*-1} (\hat{\mathbf{Q}}_{SS} - \mathbf{Q}_{SS}^*) \hat{\mathbf{Q}}_{SS}^{-1}\|_{\infty, \infty} \\ &\leq \|\mathbf{Q}_{SS}^{*-1}\|_1 \cdot \|\hat{\mathbf{Q}}_{SS} - \mathbf{Q}_{SS}^*\|_{\infty, \infty} \cdot \|\hat{\mathbf{Q}}_{SS}^{-1}\|_1. \end{aligned}$$

By Assumption we have  $\|\mathbf{Q}^{*-1}\|_1 = \|\Sigma_X^{*-1}\|_1 \cdot \|\Sigma_Y^{*-1}\|_1 \leq \theta_X \theta_Y$ . When  $n$  is sufficient large, we have  $\|\hat{\mathbf{Q}}^{-1}\|_1 \leq 2\|\mathbf{Q}^{*-1}\|_1$  by concentration. By (B.3) we have with probability at least  $1 - 2/s$  that

$$\|\hat{\mathbf{Q}}_{SS}^{-1} - \mathbf{Q}_{SS}^{*-1}\|_{\infty, \infty} \cdot \|\mathbf{b}^*\|_1 \leq 2\sqrt{10\pi} \theta_X^2 \theta_Y^2 \sigma_X \sigma_Y M \sqrt{\frac{\log s}{n}}, \quad (\text{B.5})$$

where we used the fact that  $\|\Sigma_X^* - \Sigma_Y^*\|_{1,1} \leq \sigma_X \sigma_Y M$  by Assumption 4.2. For term (i), we have with probability at least  $1 - s^{-2.5}$

$$\|[\hat{\mathbf{b}} - \mathbf{b}^*]_S\|_\infty \leq \|[\hat{\Sigma}_X - \Sigma_X^*]_S\|_{\infty, \infty} + \|[\hat{\Sigma}_Y - \Sigma_Y^*]_S\|_{\infty, \infty} \leq 6\pi \sqrt{\frac{\log s}{n}}, \quad (\text{B.6})$$

where the second inequality is due to Lemma C.1. Therefore, submitting (B.5) and (B.6) into (B.4), we obtain

$$\|\hat{\beta}_O - \beta^*\|_\infty \leq 6\pi\theta_X\theta_Y\sqrt{\frac{\log s}{n}} + 2\sqrt{10}\pi\theta_X^2\theta_Y^2\sigma_X\sigma_Y M\sqrt{\frac{\log s}{n}},$$

which holds with probability at least  $1 - 2/s - 1/s^{2.5} \geq 1 - 3/s$ .  $\square$

**Lemma B.4.** We have with probability at least  $1 - 3/d$  that

$$\|\mathcal{L}(\beta^*)\|_\infty \leq CM\sqrt{\frac{\log d}{n}},$$

where  $C$  is an absolute constants.

*Proof of Lemma B.4.* Since  $\nabla \mathcal{L}(\beta) = \hat{\mathbf{Q}}\beta - \hat{\mathbf{b}}$  and note that  $\mathbf{Q}^*\beta^* = \mathbf{b}^*$ , we have

$$\begin{aligned} \|\nabla \mathcal{L}(\beta^*)\|_\infty &= \|\hat{\mathbf{Q}}\beta^* - \mathbf{Q}^*\beta^* + \mathbf{b}^* - \hat{\mathbf{b}}\|_\infty \\ &\leq \|\hat{\mathbf{Q}} - \mathbf{Q}^*\|_{\infty, \infty} \cdot \|\beta^*\|_1 + \|\mathbf{b}^* - \hat{\mathbf{b}}\|_\infty \\ &\leq \sqrt{5}\pi M\sqrt{\frac{\log d}{n}} + 6\pi\sqrt{\frac{\log d}{n}} \end{aligned}$$

holds with probability at least  $1 - 2/d - 1/d^{2.5} \geq 1 - 3/d$ , where the first term in the last inequality is due to (B.3) and  $\|\beta^*\|_1 \leq M$  by Assumption 4.2, and the second term is due to Lemma C.1.  $\square$

## C Auxiliary Lemmas

**Lemma C.1.** [20] Given  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$  are i.i.d. random vectors following  $TE_d(\Sigma^*, \xi; f_1, f_2, \dots, f_d)$  and letting  $\hat{\Sigma}$  be the Kendall tau correlation matrix, we have that

$$\|\hat{\Sigma} - \Sigma^*\|_{\infty, \infty} \leq 3\pi\sqrt{\frac{\log d}{n}}$$

holds with probability at least  $1 - d^{-5/2}$ .

To prove the spectral norm error of  $\hat{\Sigma}$ , we first introduce the following bound for  $\hat{\mathbf{T}}$ , where  $\hat{T}_{jk} = \hat{\tau}_{jk}$ .

**Lemma C.2.** Suppose  $\delta \in (0, 1)$  satisfy  $\log(1/\delta) + d \log(9) \leq n$ . Then with probability  $1 - \delta$  it holds that

$$\sup_{\|\mathbf{x}\|_2 \leq 1} |\mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{x}| \leq 4\sqrt{C} \sqrt{\frac{\log(1/\delta) + d \log(9)}{n}},$$

where  $C$  is an absolute constant.

*Proof of Lemma C.2.* Let  $\theta = 4\sqrt{C} \sqrt{[\log(1/\delta) + s \log(12)]/n}$ . For any fixed  $\mathbf{x}$  with  $\|\mathbf{x}\|_2 \leq 1$  and any  $0 < t$ , by Markov's inequality

$$\mathbb{P}(\mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{x} > \theta) \leq \mathbb{E} \left[ \exp \left( t \cdot \mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{x} - t\theta \right) \right]. \quad (\text{C.1})$$

By Lemma D.1, we have

$$\mathbb{E} \left[ \exp \left( t \cdot \mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{x} \right) \right] \leq \exp \left( \frac{8Ct^2}{n} \right).$$

Submitting the above inequality into (C.1) and setting  $t = n\theta/(16C)$ , we obtain

$$\mathbb{P}\left(\mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{x} > \theta\right) \leq \exp\left(-\frac{n\theta^2}{16C}\right). \quad (\text{C.2})$$

Thus the error bound for  $\hat{\mathbf{T}}$  in spectral norm satisfies

$$\begin{aligned} \mathbb{P}(\|\hat{\mathbf{T}} - \mathbf{T}^*\|_2 > \theta) &= \mathbb{P}\left(\sup_{\|\mathbf{x}\|_2=1} \mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{x} > \theta\right) \\ &\leq \mathbb{P}\left((1-2\epsilon)^{-1} \sup_{\mathbf{x} \in \mathcal{N}_\epsilon} \mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{x} > \theta\right) \\ &\leq (1+2/\epsilon)^d \mathbb{P}(\mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{x} > (1-2\epsilon)\theta), \end{aligned}$$

where the first inequality is due to Lemma D.3 and the second one due to Lemma D.2. Take  $\epsilon = 1/4$  we obtain

$$\mathbb{P}(\|\hat{\mathbf{T}} - \mathbf{T}^*\|_2 > \theta) \leq 9^d \mathbb{P}(\mathbf{x}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{x} > \frac{\theta}{2}) \leq 9^d \exp\left(-\log(1/\delta) - d \log 9\right) = \delta.$$

where in the second inequality we used (C.2) and the definition of  $\theta$ . This completes the proof.  $\square$

Next, we relate  $\hat{\Sigma}$  to  $\hat{\mathbf{T}}$ . We have the following bound on the error of covariance estimator  $\hat{\Sigma}$ :

**Lemma C.3.** Assume that  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are i.i.d. random vectors following  $TE_d(\Sigma^*, \xi; f_1, f_2, \dots, f_d)$  and letting  $\hat{\Sigma}$  be the Kendall tau correlation matrix, we have

$$\|\hat{\Sigma} - \Sigma^*\|_2 \leq 2\pi^2 \frac{d \log d}{n} + 8\pi\sqrt{C} \sqrt{\frac{\log d + d \log 9}{n}}$$

holds with probability at least  $1 - 1/d - 2/d^2$ .

*Proof of Lemma C.3.* By definition in Section 3.2, we have  $\hat{\Sigma} = \cos(\pi/2\hat{\mathbf{T}})$ , where the  $\cos(\cdot)$  function is elementwise. By Taylor's theorem,

$$\hat{\Sigma} = \Sigma^* + \frac{\pi}{2} \cos\left(\frac{\pi}{2}\mathbf{T}^*\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) - \frac{\pi^2}{8} \sin\left(\frac{\pi}{2}\tilde{\mathbf{T}}\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \circ (\hat{\mathbf{T}} - \mathbf{T}^*),$$

where  $\tilde{\mathbf{T}}$  has entries  $\tilde{\tau}_{jk}$  which lies between  $\hat{\tau}_{jk}$  and  $\tau_{jk}^*$  for all  $j, k = 1, \dots, d$ . Here  $\circ$  is the Hadamard (elementwise) product for matrices. For any  $\mathbf{x} \in \mathbb{R}^d$  and  $\|\mathbf{x}\|_2 \leq 1$ ,

$$|\mathbf{x}^\top (\hat{\Sigma} - \Sigma^*)\mathbf{x}| \leq \underbrace{\frac{\pi}{2} \left| \mathbf{x}^\top \left[ \cos\left(\frac{\pi}{2}\mathbf{T}^*\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \right] \mathbf{x} \right|}_{(i)} + \underbrace{\frac{\pi^2}{8} \left| \mathbf{x}^\top \left[ \sin\left(\frac{\pi}{2}\tilde{\mathbf{T}}\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \right] \mathbf{x} \right|}_{(ii)}. \quad (\text{C.3})$$

We first bound the term (ii). Note that

$$\begin{aligned} \left| \mathbf{x}^\top \left[ \sin\left(\frac{\pi}{2}\tilde{\mathbf{T}}\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \right] \mathbf{x} \right| &\leq \|\mathbf{x}\|_1^2 \cdot \left\| \sin\left(\frac{\pi}{2}\tilde{\mathbf{T}}\right) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \circ (\hat{\mathbf{T}} - \mathbf{T}^*) \right\|_{\infty, \infty} \\ &\leq d \cdot \left\| \hat{\mathbf{T}} - \mathbf{T}^* \right\|_{\infty, \infty}^2, \end{aligned}$$

where the second inequality holds because  $|\sin(\pi/2\tilde{\tau}_{jk})| \leq 1$ , for all  $j, k = 1, \dots, d$ , and  $\|\mathbf{x}\|_1 \leq \sqrt{d}\|\mathbf{x}\|_2 \leq \sqrt{d}$ . Next, we bound term (i) in (C.3). By Lemma D.4, we can express  $\cos(\pi/2\mathbf{T}^*)$  as a convex combination,

$$\cos\left(\frac{\pi}{2}\mathbf{T}^*\right) = \sum_{i=1}^{\infty} a_i \mathbf{u}_i \mathbf{v}_i^\top,$$

where  $\mathbf{u}_i, \mathbf{v}_i \in \mathbb{R}^d$  satisfy  $\|\mathbf{u}_i\|_\infty, \|\mathbf{v}_i\|_\infty \leq 1$  for all  $i \geq 1$ , and the non-negative sequence  $a_1, a_2, \dots$  satisfy  $\sum_{i=1}^\infty a_i = 4$ . Then

$$\begin{aligned} \left| \mathbf{x}^\top \left[ \cos\left(\frac{\pi}{2} \mathbf{T}^*\right) \circ (\widehat{\mathbf{T}} - \mathbf{T}^*) \right] \mathbf{x} \right| &\leq \sum_{i=1}^\infty a_i \left| \mathbf{x}^\top \left[ \mathbf{u}_i \mathbf{v}_i^\top \circ (\widehat{\mathbf{T}} - \mathbf{T}^*) \right] \mathbf{x} \right| \\ &= \sum_{i=1}^\infty a_i \left| (\mathbf{x} \circ \mathbf{v})^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) (\mathbf{x} \circ \mathbf{u}) \right| \\ &\leq 4 \sup_{\mathbf{u}_0, \mathbf{v}_0} |\mathbf{v}_0^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \mathbf{u}_0|, \end{aligned}$$

where  $\mathbf{u}_0 = \mathbf{x} \circ \mathbf{u}$ ,  $\mathbf{v}_0 = \mathbf{x} \circ \mathbf{v}$ . Note that  $\|\mathbf{u}_0\|_2 \leq \|\mathbf{x}\|_2 \cdot \|\mathbf{u}\|_\infty \leq 1$ , similarly  $\|\mathbf{v}\|_2 \leq 1$ . Let  $\tilde{\mathbf{u}} = (\mathbf{u}_0 + \mathbf{v}_0)/2$  and  $\tilde{\mathbf{v}} = (\mathbf{u}_0 - \mathbf{v}_0)/2$ . Observe that we have  $\|\tilde{\mathbf{u}}\|_2 \leq (\|\mathbf{u}_0\|_2 + \|\mathbf{v}_0\|_2)/2 \leq 1$ , similarly  $\|\tilde{\mathbf{v}}\|_2 \leq 1$ . Therefore,

$$\begin{aligned} \sup_{\|\mathbf{u}_0\|_2, \|\mathbf{v}_0\|_2 \leq 1} |\mathbf{v}_0^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \mathbf{u}_0| &= \frac{1}{2} \sup_{\|\tilde{\mathbf{u}}\|_2, \|\tilde{\mathbf{v}}\|_2 \leq 1} |\tilde{\mathbf{u}}^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \tilde{\mathbf{u}} - \tilde{\mathbf{v}}^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \tilde{\mathbf{v}}| \\ &\leq \sup_{\|\tilde{\mathbf{u}}\|_2 \leq 1} |\tilde{\mathbf{u}}^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \tilde{\mathbf{u}}|. \end{aligned}$$

Recall the bound in (C.3), we now obtain

$$\sup_{\|\mathbf{x}\|_2 \leq 1} \left| \mathbf{x}^\top (\widehat{\Sigma} - \Sigma^*) \mathbf{x} \right| \leq \frac{\pi^2}{8} \cdot d \|\widehat{\mathbf{T}} - \mathbf{T}\|_{\infty, \infty}^2 + 2\pi \sup_{\|\mathbf{x}\|_2 \leq 1} \left| \mathbf{x}^\top (\widehat{\mathbf{T}} - \mathbf{T}^*) \mathbf{x} \right|. \quad (\text{C.4})$$

Since  $\widehat{\tau}_{jk}$  is a U-statistic, and its kernel is a bounded function between  $-1$  and  $1$  and  $\mathbb{E} \widehat{\tau}_{jk} = \tau_{jk}$ . Then by Hoeffding's inequality for U-statistics, we obtain

$$\mathbb{P} \left( \sup_{j,k} |\widehat{\tau}_{jk} - \tau_{jk}| > t \right) \leq 2d^2 e^{-\frac{nt^2}{4}}. \quad (\text{C.5})$$

Choose  $t = 4\sqrt{\log d/n}$ , we have  $\|\widehat{\mathbf{T}} - \mathbf{T}\|_{\infty, \infty} \leq 4\sqrt{\log d/n}$  with probability at least  $1 - 2/d^2$ . Plugging the bound in Lemma C.2 and (C.5) into (C.4), we obtain that

$$\|\widehat{\Sigma} - \Sigma^*\|_2 \leq 2\pi^2 \frac{d \log d}{n} + 8\pi \sqrt{C} \sqrt{\frac{\log d + d \log(9)}{n}}$$

holds with probability at least  $1 - 1/d - 2/d^2$ , where we set  $\delta = 1/d$  in Lemma C.2 and  $C$  is an absolute constant.  $\square$

**Lemma C.4.** For vectors  $\widehat{\mathbf{b}}, \mathbf{b}^* \in \mathbb{R}^{d^2}$  with entries  $\widehat{b}_j = \sin(\widehat{\tau}_j)$ , and a index set  $S$  with  $|S| = s$ , we have

$$\|[\widehat{\mathbf{b}} - \mathbf{b}^*]_S\|_2 \leq 4\sqrt{3}\pi \sqrt{\frac{s}{n}}$$

holds with probability at least  $1 - 1/s$ .

*Proof of Lemma C.4.* By definition  $\widehat{\mathbf{b}} = \text{vec}(\widehat{\Sigma}_X - \widehat{\Sigma}_Y)$ ,

$$\|\widehat{\mathbf{b}} - \mathbf{b}^*\|_2 \leq \|\text{vec}(\widehat{\Sigma}_X) - \text{vec}(\Sigma_X^*)\|_2 + \|\text{vec}(\widehat{\Sigma}_Y) - \text{vec}(\Sigma_Y^*)\|_2.$$

Denote  $\widehat{\mathbf{x}} = \text{vec}(\widehat{\Sigma}_X)$ ,  $\mathbf{x}^* = \text{vec}(\Sigma_X^*)$  and  $\widehat{\mathbf{y}} = \text{vec}(\widehat{\Sigma}_Y)$ ,  $\mathbf{y}^* = \text{vec}(\Sigma_Y^*)$ . We only need to bound  $\|\widehat{\mathbf{x}} - \mathbf{x}^*\|_2$ . By definition in Section 3.2, we have  $\widehat{\Sigma} = \cos(\pi/2 \widehat{\mathbf{T}})$  and  $\widehat{\mathbf{x}} = \cos(\pi/2 \widehat{\boldsymbol{\tau}})$ , where the  $\cos(\cdot)$  function is elementwise. Here we use  $\boldsymbol{\tau} = \text{vec}(\widehat{\mathbf{T}})$  to denote the vectorized Tau statistic. By Taylor's theorem,

$$\widehat{\mathbf{x}} = \mathbf{x}^* + \frac{\pi}{2} \cos\left(\frac{\pi}{2} \boldsymbol{\tau}^*\right) \circ (\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) - \frac{\pi^2}{8} \sin\left(\frac{\pi}{2} \widetilde{\boldsymbol{\tau}}\right) \circ (\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \circ (\widehat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*),$$

where  $\tilde{\tau}$  has entries  $\tilde{\tau}_j$  which lies between  $\hat{\tau}_j$  and  $\tau_j^*$  for all  $j = 1, \dots, d^2$  and  $\circ$  is the Hadamard (elementwise) product. For any  $\mathbf{u} \in \mathbb{R}^s$  and  $\|\mathbf{u}\|_2 \leq 1$ ,

$$|\mathbf{u}^\top (\hat{\mathbf{x}} - \mathbf{x})_S| \leq \underbrace{\frac{\pi}{2} \left| \mathbf{u}^\top \left[ \cos \left( \frac{\pi}{2} \boldsymbol{\tau}^* \right) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \right]_S \right|}_\text{(i)} + \underbrace{\frac{\pi^2}{8} \left| \mathbf{u}^\top \left[ \sin \left( \frac{\pi}{2} \tilde{\boldsymbol{\tau}} \right) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \right]_S \right|}_\text{(ii)}. \quad (\text{C.6})$$

We first bound the term (i). Note that

$$\begin{aligned} \left| \mathbf{u}^\top \left[ \sin \left( \frac{\pi}{2} \tilde{\boldsymbol{\tau}} \right) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \right]_S \right| &\leq \|\mathbf{u}\|_1 \cdot \left\| \left[ \sin \left( \frac{\pi}{2} \tilde{\boldsymbol{\tau}} \right) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \right]_S \right\|_\infty \\ &\leq \sqrt{s} \cdot \|\hat{\mathbf{T}} - \mathbf{T}^*\|_{\infty, \infty}^2, \end{aligned}$$

where the second inequality holds because  $|\sin(\pi/2 \tilde{\tau}_j)| \leq 1$ , for all  $j \in S$ , and  $\|\mathbf{u}\|_1 \leq \sqrt{s} \|\mathbf{u}\|_2 \leq \sqrt{s}$ . By (C.5), we have  $\|\hat{\mathbf{T}} - \mathbf{T}^*\|_\infty \leq 4\sqrt{\log s/n}$ .

Next, we bound term (ii) in (C.6). Note that

$$\left| \mathbf{u}^\top \left[ \cos \left( \frac{\pi}{2} \boldsymbol{\tau}^* \right) \circ (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*) \right]_S \right| = \left| \left[ \mathbf{u} \circ \cos \left( \frac{\pi}{2} \boldsymbol{\tau}_S^* \right) \right]^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S \right| = |\mathbf{u}_1^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S|,$$

where  $\mathbf{u}_1 \in \mathbb{R}^s$  is a constant vector with  $\|\mathbf{u}_1\|_2 \leq 1$ . Since  $\hat{\tau}_j$  is a U-statistic, and its kernel is a bounded function between  $-1$  and  $1$  and  $\mathbb{E}\hat{\tau}_j = \tau_j^*$ . Thus  $\tau_j - \tau_j^*$  are centered sub-Gaussian random variables and  $\|\tau_j - \tau_j^*\|_{\psi_2} \leq 2$ . Then by Hoeffding's inequality, we obtain

$$\mathbb{P}(|\mathbf{u}_1^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S| > t) \leq e^{-\frac{nt^2}{4}}. \quad (\text{C.7})$$

By Lemma D.2 and Lemma D.3

$$\begin{aligned} \mathbb{P}\left(\sup_{\mathbf{u}_1 \in \mathbb{R}^s, \|\mathbf{u}_1\|_2 \leq 1} |\mathbf{u}_1^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S| > t\right) &\leq \mathbb{P}\left(\sup_{\mathbf{u}_1 \in S_\epsilon} |\mathbf{u}_1^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S| > (1 - \epsilon)^{-1}t\right) \\ &\leq (1 + 2/\epsilon)^s \mathbb{P}\left(|\mathbf{u}_1^\top (\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*)_S| > (1 - \epsilon)^{-1}t\right) \\ &\leq 5^s \exp\left(-\frac{nt^2}{16}\right), \end{aligned}$$

where we set  $\epsilon = 1/2$  for the  $\epsilon$ -net of  $s$ -sphere. Choose  $t = 4\sqrt{3s/n}$  and then we have  $\|[\hat{\boldsymbol{\tau}} - \boldsymbol{\tau}^*]_S\|_2 \leq 4\sqrt{3s/n}$  with probability at least  $1 - 1/s$ . Therefore, we obtain

$$\|[\hat{\mathbf{x}} - \mathbf{x}^*]_S\|_2 = \max_{\|\mathbf{u}\|_2=1} |\mathbf{u}^\top (\hat{\mathbf{x}} - \mathbf{x})_S| \leq 2\pi^2 \frac{\sqrt{s \log s}}{n} + 2\sqrt{3}\pi \sqrt{\frac{s}{n}},$$

and it follows that

$$\|[\hat{\mathbf{b}} - \mathbf{b}^*]_S\|_2 \leq 4\sqrt{3}\pi \sqrt{\frac{s}{n}}$$

holds with probability at least  $1 - 1/s$ .  $\square$

## D Additional Lemmas

**Lemma D.1.** Assume that  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are i.i.d. random vectors following  $TE_d(\boldsymbol{\Sigma}^*, \xi; f_1, f_2, \dots, f_d)$ .  $\mathbf{T}^*$  is the Kendall's tau matrix defined in Section 3.2, and  $\hat{\mathbf{T}}$  is the Kendall's tau estimator in (3.1). For fixed  $\mathbf{x}$  with  $\|\mathbf{u}\|_2 \leq 1$ , for any  $t \leq nC/2$ , we have

$$\mathbb{E}[\exp(t \cdot \mathbf{u}^\top (\hat{\mathbf{T}} - \mathbf{T}^*)\mathbf{u})] \leq \exp\left(\frac{8Ct^2}{n}\right).$$

*Proof.* Denote  $S_n$  the group of permutations of  $[n]$ . For a fixed  $\mathbf{u} \in \mathbb{R}^d$  and a permutation  $\sigma \in S_n$ , define

$$Z_{\sigma, i} = \mathbf{u}^\top \left( \text{sign}((\mathbf{X}_{\sigma(i)} - \mathbf{X}_{\sigma_{i+n/2}})(\mathbf{X}_{\sigma(i)} - \mathbf{X}_{\sigma_{i+n/2}})^\top) - \mathbf{T}^* \right) \mathbf{u}, \quad i = 1, \dots, n.$$

Observe that

$$\mathbf{u}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{u} = \frac{1}{n!} \sum_{\sigma \in S_n} \frac{2}{n} \sum_{i \in [n/2]} Z_{\sigma,i}, \quad (\text{D.1})$$

and that for any fixed  $\sigma \in S_n$ , the  $Z_{\sigma,i}$ 's are i.i.d. distributed for  $i = 1, \dots, n/2$ . We denote the identical distribution as

$$\begin{aligned} \tilde{Z} &= \mathbf{u}^\top \left( \text{sign}((\mathbf{X}_i - \mathbf{X}_{i+n/2})(\mathbf{X}_i - \mathbf{X}_{i+n/2})^\top) - \mathbf{T}^* \right) \mathbf{u} \\ &= (\mathbf{u}^\top \text{sign}(\mathbf{X}_i - \mathbf{X}_{i+n/2}))^2 - \mathbb{E}(\mathbf{u}^\top \text{sign}(\mathbf{X}_i - \mathbf{X}_{i+n/2}))^2. \end{aligned}$$

Note that  $|\mathbf{u}^\top \text{sign}(\mathbf{X}_{ik} - \mathbf{X}_{i+n/2,k})| \leq \|\mathbf{u}\|_2 \leq 1$  for any  $k = 1, \dots, d$ . So  $\mathbf{u}^\top \text{sign}(\mathbf{X}_{ik} - \mathbf{X}_{i+n/2,k})$  is sub-Gaussian and  $\tilde{Z}$  is a centered sub-exponential random variable with  $\|\tilde{Z}\|_{\psi_1} \leq 2$ . By Bernstein-type inequality [30], we have

$$\mathbb{E} \exp(t\tilde{Z}) \leq e^{4Ct^2}, \quad (\text{D.2})$$

where  $C$  is an absolute constant. It immediately implies that

$$\begin{aligned} \mathbb{E}[\exp(t \cdot \mathbf{u}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{u})] &= \mathbb{E}\left[\exp\left(\frac{1}{n!} \sum_{\sigma \in S_n} \frac{2t}{n} \sum_{i \in [n/2]} Z_{\sigma,i}\right)\right] \\ &\leq \frac{1}{n!} \sum_{\sigma \in S_n} \mathbb{E}\left[\exp\left(\frac{2t}{n} \sum_{i \in [n/2]} Z_{\sigma,i}\right)\right], \end{aligned}$$

where the inequality is due to Jensen's inequality. Since for any fixed  $\sigma \in S_n$ ,  $Z_{\sigma,i}$ 's are i.i.d. distributed and equal to  $\tilde{Z}$  in distribution. We have

$$\begin{aligned} \mathbb{E}[\exp(t \cdot \mathbf{u}^\top (\hat{\mathbf{T}} - \mathbf{T}^*) \mathbf{u})] &\leq \left( \mathbb{E}\left[\exp\left(\frac{2t}{n} \tilde{Z}\right)\right] \right)^{n/2} \\ &\leq \exp\left(\frac{8Ct^2}{n}\right), \end{aligned}$$

where the second inequality is due to (D.2).  $\square$

The following lemma is about covering numbers of the sphere.

**Lemma D.2.** [30] The unit Euclidean sphere  $S^{n-1}$  equipped with the Euclidean metric satisfies for every  $\epsilon > 0$  that

$$|\mathcal{N}_\epsilon| \leq \left(1 + \frac{2}{\epsilon}\right)^n,$$

where  $\mathcal{N}_\epsilon$  is the  $\epsilon$ -net of  $S^{n-1}$ .

The following lemma is about how to compute the spectral norm on a  $\epsilon$ -net.

**Lemma D.3.** [30] Let  $\mathbf{A}$  be a symmetric  $n \times n$  matrix. For some  $\epsilon \in [0, 1/2)$ , let  $\mathcal{N}_\epsilon$  be an  $\epsilon$ -net of the unit sphere  $S^{n-1}$  in  $\mathbb{R}^n$ . Then

$$\|\mathbf{A}\|_2 = \sup_{\mathbf{x} \in S^{n-1}} |\mathbf{x}^\top \mathbf{A} \mathbf{x}| \leq (1 - 2\epsilon)^{-1} \sup_{\mathbf{x} \in \mathcal{N}_\epsilon} |\mathbf{x}^\top \mathbf{A} \mathbf{x}|.$$

**Lemma D.4.** [33, 2] There exist vectors  $\mathbf{x}_1, \mathbf{x}_2, \dots$  and  $\mathbf{y}_1, \mathbf{y}_2, \dots$  with  $\|\mathbf{x}_i\|_\infty, \|\mathbf{y}_i\|_\infty \leq 1$  for all  $i \geq 1$  and a non-negative sequence  $a_1, a_2, \dots$  with  $\sum_{i=1}^\infty a_i = 4$ , such that  $\cos(\pi/2\mathbf{T}) = \sum_{i=1}^\infty a_i \mathbf{x}_i \mathbf{y}_i^\top$ .

## E Additional Experimental Results

In this section, we present the simulation results of Gaussian differential graph model, which is a special case of the semiparametric differential graph models. Note that in the Gaussian case, the

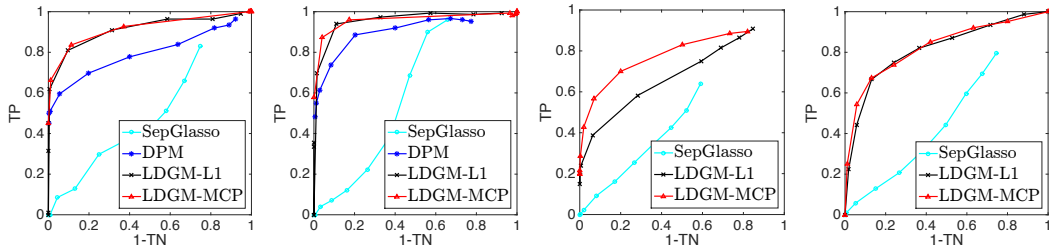
Table 3: Comparisons of estimation errors in Frobenius norm  $\|\hat{\Delta} - \Delta^*\|_F$  for Gaussian differential graph models. N/A means the algorithm did not output the solution in one day.

Methods	$n = 100, d = 100$		$n = 200, d = 400$	
	Setting 1	Setting 2	Setting 1	Setting 2
SepGlasso	33.0574 $\pm$ 0.4551	56.8891 $\pm$ 0.1778	70.1670 $\pm$ 0.4316	84.9336 $\pm$ 0.0025
DPM	23.5676 $\pm$ 0.7222	39.4366 $\pm$ 0.3814	N/A	N/A
LDGM-L1	14.0990 $\pm$ 0.6233	32.1872 $\pm$ 0.4237	29.1737 $\pm$ 0.4597	44.4980 $\pm$ 0.5482
LDGM-MCP	12.4052 $\pm$ 0.5758	28.7305 $\pm$ 0.3477	27.8458 $\pm$ 0.5843	38.7960 $\pm$ 0.3976

Table 4: Comparisons of estimation errors in infinity norm  $\|\hat{\Delta} - \Delta^*\|_{\infty, \infty}$  for Gaussian differential graph models. N/A means the algorithm did not output the solution in one day.

Methods	$n = 100, d = 100$		$n = 200, d = 400$	
	Setting 1	Setting 2	Setting 1	Setting 2
SepGlasso	3.8932 $\pm$ 0.1362	5.1321 $\pm$ 0.0102	4.1205 $\pm$ 0.1081	3.8786 $\pm$ 0.0369
DPM	3.1945 $\pm$ 0.0291	4.4132 $\pm$ 0.1060	N/A	N/A
LDGM-L1	2.7127 $\pm$ 0.0364	4.1265 $\pm$ 0.3595	2.2423 $\pm$ 0.1490	3.0224 $\pm$ 0.1088
LDGM-MCP	2.6549 $\pm$ 0.1648	3.5277 $\pm$ 0.0609	2.0638 $\pm$ 0.0388	2.3904 $\pm$ 0.1831

inputs for all the methods are the sample covariance matrices  $\hat{\Sigma}_X$  and  $\hat{\Sigma}_Y$  instead of the Kendall’s tau based correlation matrices. The ROC curves by averaging the results over 10 repetitions for Gaussian differential graph models are shown in Figure 4, from which we can see our estimator (LDGM-MCP) outperforms the others in all settings. In addition, LDGM-L1 as a special case of our estimator also performs better than DPM and SepGlasso. SepGlasso’s performance is poor since it highly depends on the sparsity of both individual graphs. When  $n > 100$ , the DPM method failed to output the solution in one day. The average results over 10 replicates for estimation in terms of Frobenius norm and infinity norm are summarized in Tables 3 and 4 respectively. Our estimator again achieves smaller error than the other baselines in all settings. In addition, LDGM-L1 also performs better than DPM and SepGlasso.



(a) Setting 1:  $n=100, d=100$  (b) Setting 2:  $n=100, d=100$  (c) Setting 1:  $n=200, d=400$  (d) Setting 2:  $n=200, d=400$

Figure 4: ROC curves for Gaussian differential graph models of all the 4 methods. There are two settings of graph structure. Note that DPM is not scalable to  $d = 400$ .