

## Supplementary material for: Mutual Transfer Learning for Massive Data

Ching-Wei Cheng Xingye Qiao Guang Cheng

The supplementary material is organized as follows:

- In Section S.1, we give the formulas for the penalty functions we have used.
- In Section S.2, we prove Theorem 2.1.
- In Section S.3, we prove Theorem 4.1.
- In Section S.4, we derive the bounds related to random matrices, which are needed in Section S.3.
- In Section S.5, we present the detailed derivation of the ADMM Algorithm 1.
- In Section S.6, we prove Proposition 3.1.
- In Section S.7, we prove Theorem 4.2.
- In Section S.8, we prove (ii) in Corollary 4.1.
- In Section S.9, we present technical lemmas.
- In Section S.10, we define the signal-to-noise ratio (SNR) for the proposed method.
- In Section S.11, we present additional simulation and real data results.

### S.1. Penalty Functions

More specifically, for  $t > 0$ , the  $L_1$ , MCP, SCAD and TLP are respectively given by

$$p^{L_1}(t; \lambda) = \lambda t, \quad (\text{S.1})$$

$$p_\gamma^{\text{MCP}}(t; \lambda) = \lambda \int_0^t \left(1 - \frac{x}{\lambda\gamma}\right)_+ dx, \quad \gamma > 1, \quad (\text{S.2})$$

$$p_\gamma^{\text{SCAD}}(t; \lambda) = \lambda \int_0^t \min \left\{ 1, \frac{(\gamma - \frac{x}{\lambda})_+}{\gamma - 1} \right\} dx, \quad \gamma > 2, \quad (\text{S.3})$$

and

$$p_\gamma^{\text{TLP}}(t; \lambda) = \lambda \min(t, \lambda\gamma), \quad \gamma > 0. \quad (\text{S.4})$$

### S.2. Proof of Theorem 2.1

To prove  $Q_N^{\text{CD}}(\beta, \Theta) - Q_N(\beta, \Theta)$  is a constant, first note that

$$Q_N^{\text{CD}}(\beta, \Theta) - Q_N(\beta, \alpha) = L_N^{\text{CD}}(\beta, \alpha) - L_N(\beta, \Theta),$$

where

$$L_N^{\text{CD}}(\beta, \Theta) = \frac{1}{2N} \sum_{i=1}^M \begin{pmatrix} \check{\beta}_i - \beta \\ \check{\theta}_i - \theta_i \end{pmatrix}^\top (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \check{\beta}_i - \beta \\ \check{\theta}_i - \theta_i \end{pmatrix}, \quad (\text{S.5})$$

$$L_N(\beta, \Theta) = \frac{1}{2N} \sum_{i=1}^M (\mathbf{y}_i - \mathbf{x}_i\beta - \mathbf{z}_i\theta_i)^\top \mathbf{W}_i(\mathbf{y}_i - \mathbf{x}_i\beta - \mathbf{z}_i\theta_i). \quad (\text{S.6})$$

Define  $\mathbf{r}_i(\boldsymbol{\beta}, \boldsymbol{\theta}_i) = (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \boldsymbol{\beta} - \boldsymbol{\beta}_0 \\ \boldsymbol{\theta}_i - \boldsymbol{\theta}_{i0} \end{pmatrix}$ , where  $\boldsymbol{\beta}_0$  and  $\boldsymbol{\theta}_{i0}$  are true values of  $\boldsymbol{\beta}$  and  $\boldsymbol{\theta}_i$ , respectively. Then we can write

$$\begin{aligned} (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \check{\boldsymbol{\beta}}_i - \boldsymbol{\beta} \\ \check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i \end{pmatrix} &= (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \check{\boldsymbol{\beta}}_i - \boldsymbol{\beta}_0 \\ \check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_{i0} \end{pmatrix} - (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \boldsymbol{\beta} - \boldsymbol{\beta}_0 \\ \boldsymbol{\theta}_i - \boldsymbol{\theta}_{i0} \end{pmatrix} \\ &= (\mathbf{x}_i, \mathbf{z}_i) [(\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i)]^{-1} (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i) - \mathbf{r}_i(\boldsymbol{\beta}, \boldsymbol{\theta}_i) \\ &= \mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)} \mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i) - \mathbf{r}_i(\boldsymbol{\beta}, \boldsymbol{\theta}_i) \end{aligned}$$

where  $\mathbf{V}_i$  denote the squared root matrix of  $\mathbf{W}_i$  such that  $\mathbf{W}_i = \mathbf{V}_i^2$ . In addition, we can also write

$$\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i = \mathbf{y}_i - (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \boldsymbol{\beta}_0 \\ \boldsymbol{\theta}_{i0} \end{pmatrix} - (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \boldsymbol{\beta} - \boldsymbol{\beta}_0 \\ \boldsymbol{\theta}_i - \boldsymbol{\theta}_{i0} \end{pmatrix} = \mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i - \mathbf{r}_i(\boldsymbol{\beta}, \boldsymbol{\theta}_i).$$

Consequently, we have

$$\begin{aligned} (\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \check{\boldsymbol{\beta}}_i - \boldsymbol{\beta} \\ \check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i \end{pmatrix} - (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i) &= (\mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)} \mathbf{V}_i - \mathbf{I}_{n_i}) (\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i) \\ &= -\mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i). \end{aligned}$$

Note that  $\mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp$  is the projection matrix onto the orthogonal complement of the column space of  $\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)$  such that  $\mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \mathbf{x}_i = \mathbf{0}$  and  $\mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \mathbf{z}_i = \mathbf{0}$ . Hence we have

$$(\mathbf{x}_i, \mathbf{z}_i) \begin{pmatrix} \check{\boldsymbol{\beta}}_i - \boldsymbol{\beta} \\ \check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i \end{pmatrix} - (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i) = -\mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i.$$

and therefore we can write

$$\begin{aligned} L_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\Theta}) &= \frac{1}{2N} \sum_{i=1}^M (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i - \mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i)^\top \mathbf{W}_i (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i - \mathbf{V}_i^{-1} \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i) \\ &= L_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) - \frac{1}{N} \sum_{i=1}^M (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \boldsymbol{\theta}_i)^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i + \frac{1}{2N} \sum_{i=1}^M \boldsymbol{\varepsilon}_i^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i \\ &= L_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) - \frac{1}{N} \sum_{i=1}^M \mathbf{y}_i^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i + \frac{1}{2N} \sum_{i=1}^M \boldsymbol{\varepsilon}_i^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp \mathbf{V}_i \boldsymbol{\varepsilon}_i, \end{aligned}$$

where the last identity follows from  $\mathbf{x}_i^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp = \mathbf{0}$  and  $\mathbf{z}_i^\top \mathbf{V}_i \mathbf{P}_{\mathbf{V}_i(\mathbf{x}_i, \mathbf{z}_i)}^\perp = \mathbf{0}$ . Hence we conclude  $Q_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\Theta}) - Q_N(\boldsymbol{\beta}, \boldsymbol{\alpha}) = L_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\alpha}) - L_N(\boldsymbol{\beta}, \boldsymbol{\Theta})$  is a constant.

### S.3. Proof of Theorem 4.1

To prove (i), we write

$$\begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} = (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{W} (\mathbf{Z} \mathbf{U} + \boldsymbol{\varepsilon}),$$

where  $\mathbf{F} = (\mathbf{X}, \mathbf{Z} \mathbf{A})$ , and hence

$$\begin{aligned} \left\| \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} \right\| &\leq \left\| (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \right\| \left\| \mathbf{F}^\top \mathbf{W} (\mathbf{Z} \mathbf{U} + \boldsymbol{\varepsilon}) \right\| \\ &\leq \sqrt{p + Sq} \left\| (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \right\| \left\| \mathbf{F}^\top \mathbf{W} (\mathbf{Z} \mathbf{U} + \boldsymbol{\varepsilon}) \right\|_\infty. \end{aligned}$$

We derive the upper bounds for  $\|(\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1}\|$  and  $\|\mathbf{F}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty$  under the events

$$E_{CB} = \left\{ \max_{1 \leq i \leq p} \|[\mathbf{X}]_{\cdot j}\| \leq \sqrt{5c_f N} \text{ and } \max_{1 \leq s \leq S} \max_{1 \leq k \leq q} \|[\tilde{\mathbf{Z}}_s]_{\cdot k}\| \leq \sqrt{5c_f g_{\max}} \right\}$$

and

$$E_{EB} = \left\{ \underline{C}_f \lesssim \lambda_{\min}(\mathbf{f}_i^\top \mathbf{f}_i / n_i) \leq \lambda_{\max}(\mathbf{f}_i^\top \mathbf{f}_i / n_i) \lesssim \bar{C}_f, i = 1, \dots, M \right\},$$

where  $\mathbf{f}_i = (\mathbf{x}_i, \mathbf{z}_i)$  and  $n_{\min} = \min_{1 \leq i \leq M} n_i$ . According to Sections S.4.1 and S.4.2, we have  $P(E_{CB}) \geq 1 - (p + Sq)e^{-g_{\max}}$  and  $P(E_{EB}) \geq 1 - 2Me^{-n_{\min}}$ .

We first look at  $\|(\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1}\|$ . By definition, we can write

$$\mathbf{W} = [\sigma_\varepsilon^2 \mathbf{I}_N + \mathbf{Z}(\mathbf{I}_M \otimes \boldsymbol{\Psi})\mathbf{Z}^\top]^{-1} = \sigma_\varepsilon^{-2} \left[ \mathbf{I}_N - \sigma_\varepsilon^{-2} \mathbf{Z}(\mathbf{I}_M \otimes \boldsymbol{\Psi}^{-1} + \sigma_\varepsilon^{-2} \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \right],$$

where the second identity holds by applying the Woodbury identity (S.41). It then follows that

$$\begin{aligned} \|(\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1}\| &= \lambda_{\min}^{-1}(\mathbf{F}^\top \mathbf{W} \mathbf{F}) \\ &= \sigma_\varepsilon^2 \lambda_{\min}^{-1} \left[ \mathbf{F}^\top \mathbf{F} - \sigma_\varepsilon^{-2} \mathbf{F}^\top \mathbf{Z}(\mathbf{I}_M \otimes \boldsymbol{\Psi}^{-1} + \sigma_\varepsilon^{-2} \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{F} \right] \\ &\leq \sigma_\varepsilon^2 \left\{ \lambda_{\min}(\mathbf{F}^\top \mathbf{F}) + \sigma_\varepsilon^2 \lambda_{\min} \left[ \mathbf{F}^\top \mathbf{Z}(\mathbf{I}_M \otimes \boldsymbol{\Psi}^{-1} + \sigma_\varepsilon^{-2} \mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{F} \right] \right\}^{-1} \\ &\leq \sigma_\varepsilon^2 \lambda_{\min}^{-1}(\mathbf{F}^\top \mathbf{F}) \\ &\lesssim \underline{C}_f^{-1} \sigma_\varepsilon^2 g_{\min}^{-1}, \end{aligned} \tag{S.7}$$

where the last inequality holds from (S.14).

It remains to show the upper bound for  $\|\mathbf{F}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|$ . Recall that  $\mathbf{F} = (\mathbf{X}, \mathbf{Z} \mathbf{A})$ . Let  $\mathbf{V}_i$  denote the square root matrix of  $\mathbf{W}_i$  such that  $\mathbf{W}_i = \mathbf{V}_i^2$ , and then define  $\mathbf{V} = \text{diag}[(\mathbf{V}_i)_{i=1, \dots, M}]$  and  $\mathbb{V}_s = \text{diag}[(\mathbf{V}_i)_{i:L_i=s}]$ . It is easy to see from (S.15) that  $\|\mathbf{V}\| = \|\mathbf{W}\|^{1/2} = \max_{1 \leq i \leq M} \|\mathbf{W}_i\|^{1/2} \leq \sigma_\varepsilon^{-1}$  and  $\|\mathbb{V}_s\| = \max_{i:L_i=s} \|\mathbf{W}_i\|^{1/2} \leq \sigma_\varepsilon^{-1}$ . Let  $\tilde{\mathbf{Z}}_s = (\mathbf{z}_i^\top)_{i:L_i=s}^\top$  denote the stacked design matrix for the  $s$ -th subgroup. Under the events  $E_{CB}$  and  $E_{EB}$ , we have

$$\|\mathbf{V}[\mathbf{X}]_{\cdot j}\| \leq \|\mathbf{V}\| \|[\mathbf{X}]_{\cdot j}\| \leq \sigma_\varepsilon^{-1} \sqrt{5c_f N} \quad \text{and} \quad \|\mathbb{V}_s[\tilde{\mathbf{Z}}_s]_{\cdot k}\| \leq \|\mathbb{V}_s\| \|[\tilde{\mathbf{Z}}_s]_{\cdot k}\| \leq \sigma_\varepsilon^{-1} \sqrt{5c_f g_{\max}}$$

for all  $j = 1, \dots, p, s = 1, \dots, S, k = 1, \dots, q$ . By union bound and the tail bounds in (S.18) and (S.20), it holds

$$\begin{aligned} P\left(\|\mathbf{X}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty > t\right) &\leq \sum_{j=1}^p P\left(\|[\mathbf{X}]_{\cdot j}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\| > t\right) \\ &\leq 4p \exp\left(-\frac{c_e (\tau \wedge \sigma_\varepsilon^2) \sigma_\varepsilon^2 t^2}{5c_f N}\right), \end{aligned} \tag{S.8}$$

$$\begin{aligned} P\left(\|(\mathbf{Z} \mathbf{A})^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty > t\right) &\leq \sum_{s=1}^S \sum_{k=1}^q P\left(\|[\tilde{\mathbf{Z}}_s]_{\cdot k}^\top \mathbb{W}_s(\mathbb{Z}_s \mathbf{U}_s + \boldsymbol{\mathcal{E}}_s)\| > t\right) \\ &\leq 4Sq \exp\left(-\frac{c_e (\tau \wedge \sigma_\varepsilon^2) \sigma_\varepsilon^2 t^2}{5c_f g_{\max}}\right), \end{aligned} \tag{S.9}$$

where  $\mathbb{Z}_s = \text{diag}[(\mathbf{z}_i)_{i:L_i=s}]$ ,  $\mathbf{U}_s = (\mathbf{u}_i^\top)_{i:L_i=s}^\top$ ,  $\mathbb{W}_s = \text{diag}[(\mathbf{W}_i)_{i:L_i=s}]$ ,  $\boldsymbol{\mathcal{E}}_s = (\boldsymbol{\varepsilon}_i^\top)_{i:L_i=s}^\top$ , and  $\tau = \lambda_{\min}(\boldsymbol{\Psi})$ . Since  $\mathbf{F} = (\mathbf{X}, \mathbf{Z} \mathbf{A})$ , it follows that

$$\begin{aligned} P\left(\|\mathbf{F}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty > t\right) &\leq P\left(\|\mathbf{X}^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty > t\right) + P\left(\|(\mathbf{Z} \mathbf{A})^\top \mathbf{W}(\mathbf{Z} \mathbf{U} + \boldsymbol{\mathcal{E}})\|_\infty > t\right) \\ &\leq 2p \exp\left(-\frac{c_e (\tau \wedge \sigma_\varepsilon^2) \sigma_\varepsilon^2 t^2}{5c_f N}\right) + 2Sq \exp\left(-\frac{c_e (\tau \wedge \sigma_\varepsilon^2) \sigma_\varepsilon^2 t^2}{5c_f g_{\max}}\right) \\ &\leq 4(p + Sq) \exp\left(-\frac{c_e (\tau \wedge \sigma_\varepsilon^2) \sigma_\varepsilon^2 t^2}{5c_f N}\right). \end{aligned}$$

Taking  $t = \sqrt{5}c_e^{-1/2}c_f^{1/2}(\tau \wedge \sigma_\varepsilon^2)^{-1/2}\sigma_\varepsilon^{-1}\sqrt{N \log N}$  hence yields

$$\|\mathbf{F}^\top \mathbf{W}(\mathbf{Z}\mathbf{U} + \boldsymbol{\varepsilon})\|_\infty \leq \sqrt{5}c_e^{-1/2}c_f^{1/2}(\tau \wedge \sigma_\varepsilon^2)^{-1}\sigma_\varepsilon^{-1/2}\sqrt{N \log N} \quad (\text{S.10})$$

with probability at least  $1 - (p + Sq)(4N^{-1} + e^{-g_{\max}}) - 2Me^{-n_{\min}}$ .

In summary, combining the results in (S.7) and (S.10) leads to

$$\begin{aligned} \left\| \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} \right\| &\leq \sqrt{p + Sq} \left\| (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \right\| \|\mathbf{F}^\top \mathbf{W}(\mathbf{Z}\mathbf{U} + \boldsymbol{\varepsilon})\|_\infty \\ &\leq \sqrt{5}c_e^{-1/2}c_f^{1/2}\underline{C}_f^{-1}(\tau \wedge \sigma_\varepsilon^2)^{-1/2}\sigma_\varepsilon g_{\min}^{-1} \sqrt{(p + Sq)N \log N}, \end{aligned}$$

with probability at least  $1 - (p + Sq)(4N^{-1} + e^{-g_{\max}}) - 2Me^{-n_{\min}}$ .

To prove (ii), we show the oracle estimator satisfies the Lindeberg-Feller condition as follows. Note that for any  $\mathbf{a}_N \in \mathbb{R}^{p+Sq}$  with  $\|\mathbf{a}_N\| = 1$ ,

$$\mathbf{a}_N^\top \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} = \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \mathbf{F}^\top \mathbf{W}(\mathbf{Z}\mathbf{U} + \boldsymbol{\varepsilon}).$$

It is then clear that

$$\mathbb{E} \left[ \mathbf{a}_N^\top \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} \right] = \mathbf{0}$$

and

$$\begin{aligned} \sigma_N^2(\mathbf{a}_N) &= \text{Var} \left[ \mathbf{a}_N^\top \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} - \boldsymbol{\beta}_0 \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} - \boldsymbol{\alpha}_0 \end{pmatrix} \right] = \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \mathbf{a}_N \\ &\geq \lambda_{\min} [(\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1}] = \|\mathbf{F}^\top \mathbf{W} \mathbf{F}\|^{-1} \geq \|\mathbf{F}\|^{-2} \|\mathbf{W}\|^{-1} \gtrsim \overline{C}_f \sigma_\varepsilon^2 N^{-1}, \end{aligned}$$

where the last inequality is due to (S.13) and (S.16). Let  $\xi_{il} = [\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)]_l$  denote the  $l$ -th coordinate of  $\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)$ , and  $\mathbf{F}_i$  the  $n_i \times (p + Sq)$  submatrix of  $\mathbf{F}$  that collects the corresponding rows of the  $i$ -th unit. Then for any  $\epsilon > 0$ , we can write

$$\begin{aligned} \sigma_N^{-2}(\mathbf{a}_N) &\sum_{i=1}^M \sum_{l=1}^{n_i} \mathbb{E} \left\{ \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^2 \mathbb{1}_{\left| \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right| > \epsilon \sigma_N(\mathbf{a}_N)} \right\} \\ &\leq \overline{C}_f^{-1} \sigma_\varepsilon^{-2} N \sum_{i=1}^M \sum_{l=1}^{n_i} \left\{ \mathbb{E} \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^4 \right\}^{1/2} \\ &\quad \times \left\{ P \left( \left| \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right| > \epsilon \sigma_N(\mathbf{a}_N) \right) \right\}^{1/2}. \end{aligned}$$

We are going to show the above quantity is  $o(1)$ .

From (S.19) we can see that  $\xi_{il}$  has sub-Gaussian tails, and thus there exists a constant  $c_3 > 0$  such that  $\mathbb{E}[\|\xi_{il}\|^4] \leq c_3$  and  $\mathbb{E}[\|\xi_{il}\|^2] \leq c_3$ . By the definition of  $\mathbf{F}_i$  and Lemma S.9.2, we have  $\|[\mathbf{F}_i]_l\|^2 \leq 5c_4(p + Sq)$  with probability at least  $1 - e^{-(p+Sq)}$  for some constant  $c_4 > 0$ . It then follows that

$$\begin{aligned} &\left\{ \mathbb{E} \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^4 \right\}^{1/2} \\ &\leq \|\mathbf{a}_N\|^2 \left\| (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \right\|^2 \|\mathbf{V}_i\|^2 \|[\mathbf{F}_i]_l\|^2 \left\{ \mathbb{E}[\|\xi_{il}\|^4] \right\}^{1/2} \\ &\leq 5c_3^{1/2} c_4 \underline{C}_f^{-2} \overline{C}_f \sigma_\varepsilon^2 (p + Sq) g_{\min}^{-2} \\ &= O \left[ (p + Sq) g_{\min}^{-2} \right], \end{aligned}$$

where the second inequality holds based on (S.7) and (S.15) under the event  $E_{EB}$ . In addition, similar argument gives

$$\begin{aligned} \mathbb{E} \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^2 &\leq \|\mathbf{a}_N\|^2 \left\| (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} \right\|^2 \|\mathbf{V}_i\|^2 \|\mathbf{F}_i\|_l^2 \mathbb{E} [\|\xi_{il}\|^2] \\ &= O \left[ (p + Sq) g_{\min}^{-2} \right]. \end{aligned}$$

Then by Markov inequality, it holds

$$\begin{aligned} P \left( \left| \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right| > \epsilon \sigma_N(\mathbf{a}_N) \right) &\leq \frac{\mathbb{E} \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^2}{\epsilon^2 \sigma_N^2(\mathbf{a}_N)} \\ &= O \left[ (p + Sq) N g_{\min}^{-2} \right]. \end{aligned}$$

It then can be concluded from the above results that

$$\begin{aligned} \sigma_N^{-2}(\mathbf{a}_N) \sum_{i=1}^M \sum_{l=1}^{n_i} \mathbb{E} \left\{ \left[ \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right]^2 \mathbb{1}_{\left| \mathbf{a}_N^\top (\mathbf{F}^\top \mathbf{W} \mathbf{F})^{-1} [\mathbf{V}_i \mathbf{F}_i]_l^\top \xi_{il} \right| > \epsilon \sigma_N(\mathbf{a}_N)} \right\} \\ = O \left[ (p + Sq)^2 N^3 g_{\min}^{-4} \right], \end{aligned}$$

which is  $o(1)$  provided that  $g_{\min} \gg (p + Sq)^{1/2} N^{3/4}$ . Accordingly, the asymptotic normality assertion (12) follows from the Lindeberg-Feller Central Limit Theorem.

## S.4. Properties of Random Matrices

In this section, we derive the bounds for the  $L_2$  norms and eigenvalues of the random design matrices under Assumption 4.1. In addition, we derive the non-asymptotic upper bound for  $\mathbf{V}(\mathbf{Z}\mathbf{U} + \mathcal{E})$ , where  $\mathbf{V}$  denotes the square root matrix of  $\mathbf{W}$  such that  $\mathbf{W} = \mathbf{V}^2$ , under Assumption 4.2.

### S.4.1. $L_2$ -norm Bounds for Columns of Design Matrices

By union bound and (S.42), there exists a constant  $c_1 > 0$  such that

$$P \left( \max_{1 \leq j \leq p} \|\mathbf{X}_{\cdot j}\| > \sqrt{5c_1 N} \right) \leq \sum_{j=1}^p P \left( \|\mathbf{X}_{\cdot j}\| > \sqrt{5c_1 N} \right) \leq p e^{-N}$$

and similarly

$$P \left( \max_{1 \leq s \leq S} \max_{1 \leq k \leq q} \|\tilde{\mathbf{Z}}_s\|_{\cdot k} > \sqrt{5c_1 g_{\max}} \right) \leq S q e^{-g_{\max}},$$

where  $\mathbb{Z}_s = (\mathbf{z}_i^\top)_{i:L_i=s}^\top$ .

We define an event

$$E_{CB} = \left\{ \max_{1 \leq j \leq p} \|\mathbf{X}_{\cdot j}\| \leq \sqrt{5c_1 N} \quad \text{and} \quad \max_{1 \leq s \leq S} \max_{1 \leq k \leq q} \|\tilde{\mathbf{Z}}_s\|_{\cdot k} \leq \sqrt{5c_1 g_{\max}} \right\} \quad (\text{S.11})$$

with  $P(E_{CB}) \geq 1 - (p + Sq)e^{-g_{\max}}$  for later use. (The subscript ‘‘CB’’ stands for ‘‘column bounds.’’)

### S.4.2. Eigenvalue Bounds for Design Matrices

In this section, we first derive the eigenvalue bounds for  $\mathbf{x}_i^\top \mathbf{x}_i$  and  $\mathbf{z}_i^\top \mathbf{z}_i$  for  $i = 1, \dots, M$ , and then those for  $(\mathbf{X}, \mathbf{AZ})^\top (\mathbf{X}, \mathbf{AZ})$ .

Let  $\mathbf{f}_i = (\mathbf{x}_i, \mathbf{z}_i)$ . Since the rows of  $\mathbf{f}_i$  are independently drawn from  $\mathbf{F}$ , then the rows of  $\mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i$  are isotropic, i.e.,  $\text{Cov}(\mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} [\mathbf{f}_i]_k) = \mathbf{I}$ . Following Theorem 5.39 of Vershynin (2012), it then can be shown that, on an event with probability at least  $1 - 2 \exp(-c_2 t^2)$ , we have

$$\begin{aligned} \sqrt{n_i} - C\sqrt{p+q} - t &\leq \lambda_{\min}^{1/2} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i^\top \mathbf{f}_i \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \right) \\ &\leq \lambda_{\max}^{1/2} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i^\top \mathbf{f}_i \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \right) \leq \sqrt{n_i} + C\sqrt{p+q} + t, \end{aligned}$$

where  $C, c_2$  are some positive constants. By taking  $t = c_2^{-1/2} \sqrt{n_i}$ , with probability at least  $1 - 2e^{-n_i}$  we have

$$\begin{aligned} \lambda_{\min}^{1/2}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) &= n_i^{-1/2} \lambda_{\min}^{1/2} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{1/2} \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i^\top \mathbf{f}_i \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{1/2} \right) \\ &\geq n_i \lambda_{\min} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{1/2} \right) \lambda_{\min}^{1/2} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i^\top \mathbf{f}_i \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \right) \\ &\geq C_f^{1/2} \left( 1 - C \sqrt{\frac{p+q}{n_i}} - c_2^{-1/2} \right) \gtrsim (1 - c_2^{-1/2}) C_f^{1/2} \end{aligned}$$

and similarly

$$\begin{aligned} \lambda_{\max}^{1/2}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) &\leq n_i^{-1/2} \lambda_{\max} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{1/2} \right) \lambda_{\max}^{1/2} \left( \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \mathbf{f}_i^\top \mathbf{f}_i \mathbb{E}[\mathbf{F}\mathbf{F}^\top]^{-1/2} \right) \\ &\leq C_f^{-1/2} \left( 1 + C \sqrt{\frac{p+q}{n_i}} + c_2^{-1/2} \right) \lesssim (1 + c_2^{-1/2}) C_f^{-1/2}. \end{aligned}$$

In summary, it holds

$$(1 - c_2^{-1/2}) C_f \lesssim \lambda_{\min}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) \leq \lambda_{\max}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) \lesssim (1 + c_2^{-1/2}) C_f^{-1},$$

with probability at least  $1 - 2e^{-n_i}$  for  $i = 1, \dots, M$ .

By the fact that the largest (smallest) eigenvalue of a submatrix with fewer columns is bounded above (below) from the largest (smallest) eigenvalue of the full matrix, we have

$$\begin{aligned} (1 - c_2^{-1/2}) C_f &\lesssim \lambda_{\min}(\mathbf{x}_i^\top \mathbf{x}_i/n_i) \leq \lambda_{\max}(\mathbf{x}_i^\top \mathbf{x}_i/n_i) \lesssim (1 + c_2^{-1/2}) C_f^{-1}, \\ (1 - c_2^{-1/2}) C_f &\lesssim \lambda_{\min}(\mathbf{z}_i^\top \mathbf{z}_i/n_i) \leq \lambda_{\max}(\mathbf{z}_i^\top \mathbf{z}_i/n_i) \lesssim (1 + c_2^{-1/2}) C_f^{-1}, \end{aligned}$$

with probability at least  $1 - 2e^{-n_i}$  for each  $i = 1, \dots, M$ .

Define the event

$$E_{EB} = \left\{ \underline{C}_f \lesssim \lambda_{\min}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) \leq \lambda_{\max}(\mathbf{f}_i^\top \mathbf{f}_i/n_i) \lesssim \bar{C}_f, i = 1, \dots, M \right\} \quad (\text{S.12})$$

with  $P(E_{EB}) \geq 1 - 2Me^{-n_{\min}}$ , where  $\underline{C}_f = (1 - c_2^{-1/2}) C_f$ ,  $\bar{C}_f = (1 + c_2^{-1/2}) C_f^{-1}$  and  $n_{\min} = \min_{1 \leq i \leq M} n_i$ . Then, on the event  $E_{EB}$ , it follows that

$$\begin{aligned} \lambda_{\max}(\mathbf{F}^\top \mathbf{F}) &\leq \max \left\{ \lambda_{\max}(\mathbf{X}^\top \mathbf{X}), \lambda_{\max}[(\mathbf{Z}\mathbf{A})^\top (\mathbf{Z}\mathbf{A})] \right\} \\ &\leq \max \left\{ \sum_{i=1}^M \lambda_{\max}(\mathbf{x}_i^\top \mathbf{x}_i), \max_{1 \leq s \leq S} \sum_{i:L_i=s} \lambda_{\max}(\mathbf{z}_i^\top \mathbf{z}_i) \right\} \\ &\lesssim \max \{ \bar{C}_f N, \bar{C}_f g_{\max} \} = \bar{C}_f N \end{aligned} \quad (\text{S.13})$$

and

$$\begin{aligned} \lambda_{\min}(\mathbf{F}^\top \mathbf{F}) &\geq \min \left\{ \lambda_{\min}(\mathbf{X}^\top \mathbf{X}), \lambda_{\min}[(\mathbf{Z}\mathbf{A})^\top (\mathbf{Z}\mathbf{A})] \right\} \\ &\geq \min \left\{ \sum_{i=1}^M \lambda_{\min}(\mathbf{x}_i^\top \mathbf{x}_i), \min_{1 \leq s \leq S} \sum_{i:L_i=s} \lambda_{\min}(\mathbf{z}_i^\top \mathbf{z}_i) \right\} \\ &\gtrsim \min \{ \underline{C}_f N, \underline{C}_f g_{\min} \} = \underline{C}_f g_{\min}. \end{aligned} \quad (\text{S.14})$$

### S.4.3. Weight Matrix

It is easy to see that

$$\lambda_{\max}(\mathbf{W}_i) = \lambda_{\min}^{-1}(\sigma_\varepsilon^2 \mathbf{I}_{n_i} + \mathbf{z}_i \boldsymbol{\Psi} \mathbf{z}_i^\top) \leq \frac{1}{\sigma_\varepsilon^2 + \lambda_{\min}(\mathbf{z}_i \boldsymbol{\Psi} \mathbf{z}_i^\top)} = \sigma_\varepsilon^{-2}, \quad (\text{S.15})$$

where  $\lambda_{\min}(\mathbf{z}_i \boldsymbol{\Psi} \mathbf{z}_i^\top) = 0$  since  $q < n_i$ . Note that this is an exact result without any probabilistic nor asymptotic statement. Moreover, since  $\mathbf{W} = \text{diag}[(\mathbf{W}_i)_{i=1, \dots, M}]$ , we can also conclude

$$\lambda_{\max}(\mathbf{W}) = \max_{1 \leq i \leq M} \lambda_{\max}(\mathbf{W}_i) \leq \sigma_\varepsilon^{-2}. \quad (\text{S.16})$$

#### S.4.4. Tail bounds of Weighted Random Effects and Noises

On the event  $E_{EB}$ , we show the non-asymptotic tail bounds for  $\mathbf{V}(\mathbf{Z}\mathbf{U} + \boldsymbol{\mathcal{E}})$ . For  $\mathbf{V}\mathbf{Z}\mathbf{U}$ , we first check  $\|\mathbf{Z}^\top \mathbf{V}\|$  as follows. By definition, we can write

$$\mathbf{W}_i = (\sigma_\varepsilon^2 \mathbf{I}_{n_i} + \mathbf{z}_i \boldsymbol{\Psi} \mathbf{z}_i^\top)^{-1} = \sigma_\varepsilon^{-2} [\mathbf{I}_{n_i} - \sigma_\varepsilon^{-2} \mathbf{z}_i (\boldsymbol{\Psi}^{-1} + \sigma_\varepsilon^{-2} \mathbf{z}_i^\top \mathbf{z}_i)^{-1} \mathbf{z}_i^\top],$$

where the second identity holds by applying the Woodbury identity (S.41). It then follows that

$$\begin{aligned} (\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i)^{-1} &= \sigma_\varepsilon^2 [\mathbf{z}_i^\top \mathbf{z}_i - \sigma_\varepsilon^{-2} \mathbf{z}_i^\top \mathbf{z}_i (\boldsymbol{\Psi}^{-1} + \sigma_\varepsilon^{-2} \mathbf{z}_i^\top \mathbf{z}_i)^{-1} \mathbf{z}_i^\top \mathbf{z}_i]^{-1} \\ &= \sigma_\varepsilon^2 (\mathbf{z}_i^\top \mathbf{z}_i)^{-1} + \boldsymbol{\Psi}^{-1}, \end{aligned}$$

where the second and identity holds by applying Woodbury identity (S.41) on the inverse term. It then holds

$$\begin{aligned} \lambda_{\max}(\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i) &= \lambda_{\min}^{-1}[(\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i)^{-1}] \\ &= \sigma_\varepsilon^{-2} \lambda_{\min}^{-1}[(\mathbf{z}_i^\top \mathbf{z}_i)^{-1} + \sigma_\varepsilon^{-2} \boldsymbol{\Psi}^{-1}] \\ &\leq \sigma_\varepsilon^{-2} \{ \lambda_{\min}[(\mathbf{z}_i^\top \mathbf{z}_i)^{-1}] + \sigma_\varepsilon^{-2} \lambda_{\min}(\boldsymbol{\Psi}^{-1}) \}^{-1} \\ &\leq \tau^{-1}. \end{aligned}$$

Accordingly, we have

$$\|\mathbf{Z}^\top \mathbf{V}\| \leq \|\mathbf{V}\mathbf{Z}\| = \lambda_{\max}^{1/2}(\mathbf{Z}^\top \mathbf{W}\mathbf{Z}) = \max_{1 \leq i \leq M} \lambda_{\max}^{1/2}(\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i) \leq \tau^{-1/2}. \quad (\text{S.17})$$

Under events  $E_{CB}$  and  $E_{EB}$ , it follows for any  $\mathbf{a} \in \mathbb{R}^N$  that

$$\begin{aligned} P(|\mathbf{a}^\top \mathbf{V}\mathbf{Z}\mathbf{U}| > \|\mathbf{a}\|t | \mathbf{Z}) &\leq P\left(|\mathbf{a}^\top \mathbf{V}\mathbf{Z}\mathbf{U}| > \|\mathbf{Z}^\top \mathbf{V}\mathbf{a}\| \frac{t}{\|\mathbf{V}\mathbf{Z}\|} \middle| \mathbf{Z}\right) \\ &\leq 2 \exp\left(-\frac{c_e t^2}{\|\mathbf{V}\mathbf{Z}\|^2}\right) \leq 2 \exp(-c_e \tau t^2). \end{aligned}$$

By law of iterated expectations, we have

$$\begin{aligned} P(|\mathbf{a}^\top \mathbf{V}\mathbf{Z}\mathbf{U}| > \|\mathbf{a}\|t) &= \mathbb{E}_{\mathbf{Z}} [P(|\mathbf{a}^\top \mathbf{V}\mathbf{Z}\mathbf{U}| > \|\mathbf{a}\|t | \mathbf{Z})] \\ &\leq \mathbb{E}_{\mathbf{Z}} \{2 \exp[-c_e \tau t^2]\} = 2 \exp[-c_e \tau t^2]. \end{aligned}$$

For  $\mathbf{V}\boldsymbol{\mathcal{E}}$ , since  $\|\mathbf{V}\| \leq \sigma_\varepsilon^{-1}$ , similar argument yields

$$P(|\mathbf{a}^\top \mathbf{V}\boldsymbol{\mathcal{E}}| > \|\mathbf{a}\|t) \leq 2 \exp[-c_e \sigma_\varepsilon^2 t^2].$$

In summary, it holds

$$\begin{aligned} P(|\mathbf{a}^\top \mathbf{V}(\mathbf{Z}\mathbf{U} + \boldsymbol{\mathcal{E}})| > \|\mathbf{a}\|t) &\leq P(|\mathbf{a}^\top \mathbf{V}\mathbf{Z}\mathbf{U}| > \|\mathbf{a}\|t) + P(|\mathbf{a}^\top \mathbf{V}\boldsymbol{\mathcal{E}}| > \|\mathbf{a}\|t) \\ &\leq 4 \exp(-c_e (\tau \wedge \sigma_\varepsilon^2) t^2), \end{aligned} \quad (\text{S.18})$$

where  $\tau = \lambda_{\min}(\boldsymbol{\Psi})$ .

In addition, similar argument results in

$$P(|\mathbf{a}_i^\top \mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)| > \|\mathbf{a}_i\|t) \leq 4 \exp\{-c_e (\tau \wedge \sigma_\varepsilon^2) t^2\} \quad (\text{S.19})$$

for any  $\mathbf{a}_i \in \mathbb{R}^{n_i}$ , where  $\mathbf{V}_i$  is the square root matrix of  $\mathbf{W}_i$  such that  $\mathbf{W}_i = \mathbf{V}_i^2$ . Moreover, let  $\mathbb{V}_s = \text{diag}[(\mathbf{V}_i)_{i:L_i=s}]$ ,  $\mathbb{U}_s = (\mathbf{u}_i)_{i:L_i=s}^\top$  and  $\boldsymbol{\mathcal{E}}_s = (\boldsymbol{\varepsilon}_i)_{i:L_i=s}^\top$ . Then we have

$$P(|\mathbf{b}_s^\top \mathbb{V}_s(\mathbb{Z}_s \mathbb{U}_s + \boldsymbol{\mathcal{E}}_s)| > \|\mathbf{b}_s\|t) \leq 4 \exp\{-c_e (\tau \wedge \sigma_\varepsilon^2) t^2\} \quad (\text{S.20})$$

for any  $\mathbf{b}_s \in \mathbb{R}^{g_s}$ .

## S.5. Derivation of Algorithm 1

The augmented Lagrangian of (9) can then be formulated by

$$\begin{aligned} L(\boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\delta}, \boldsymbol{\nu}) = & L_0(\boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\delta}) \\ & + \sum_{1 \leq i < j \leq M} \boldsymbol{\nu}_{ij}^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij}) \\ & + \frac{\kappa}{2} \sum_{1 \leq i < j \leq M} \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij}\|^2, \end{aligned} \quad (\text{S.21})$$

where the dual variables  $\boldsymbol{\nu} = (\boldsymbol{\nu}_{ij}^\top)_{i < j}^\top$  are Lagrange multipliers and  $\kappa > 0$  is the penalty parameter.

The minimizer of the augmented Lagrangian (S.21) can be solved by developing an ADMM algorithm. Given the values of  $\boldsymbol{\Theta}^{[k]}$ ,  $\boldsymbol{\delta}^{[k]}$  and  $\boldsymbol{\nu}^{[k]}$  at the  $k$ -th iteration, we update  $\boldsymbol{\beta}$  using

$$\begin{aligned} \boldsymbol{\beta}^{[k+1]} = & \left( \sum_{i=1}^M \mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i \right)^{-1} \times \\ & \left( \sum_{i=1}^M \mathbf{x}_i^\top \mathbf{W}_i \left[ \mathbf{x}_i \check{\boldsymbol{\beta}}_i + z_i (\check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i^{[k]}) \right] \right). \end{aligned} \quad (\text{S.22})$$

Next, given  $\boldsymbol{\beta}^{[k+1]}$ ,  $\boldsymbol{\delta}^{[k]}$  and  $\boldsymbol{\nu}^{[k]}$ , the updating formula for  $\boldsymbol{\Theta}$  is

$$\begin{aligned} \boldsymbol{\Theta}^{[k+1]} = & (\mathbf{Z}^\top \mathbf{W} \mathbf{Z} + \kappa N \mathbf{B} \mathbf{B}^\top)^{-1} \times \\ & [\mathbf{Z}^\top \mathbf{W} \mathbf{Z} \check{\boldsymbol{\Theta}} + \mathbf{D}_{zx} (\check{\boldsymbol{\beta}} - \mathbf{1}_M \otimes \boldsymbol{\beta}^{[k+1]}) \\ & + N \mathbf{B} (\kappa \boldsymbol{\delta}^{[k]} - \boldsymbol{\nu}^{[k]})], \end{aligned} \quad (\text{S.23})$$

where  $\check{\boldsymbol{\Theta}} = (\check{\boldsymbol{\theta}}_i^\top)_{i=1, \dots, M}^\top$ ,  $\check{\boldsymbol{\beta}} = (\check{\boldsymbol{\beta}}_i^\top)_{i=1, \dots, M}^\top$ ,  $\mathbf{D}_{zx} = \text{diag} [(z_i^\top \mathbf{W}_i \mathbf{x}_i)_{i=1, \dots, M}]$  and  $\mathbf{B}$  is a matrix such that  $\mathbf{B}^\top \boldsymbol{\Theta} = (\boldsymbol{\theta}_i^\top - \boldsymbol{\theta}_j^\top)_{i < j}^\top$ . Note that the form of  $\mathbf{B}$  depends on how we arrange  $\boldsymbol{\delta}_{ij}$ 's in  $\boldsymbol{\delta}$ ; see Remark S.5.1.

To update  $\boldsymbol{\delta}$ , we use (S.29), (S.30), (S.31) and (S.32) derived below for  $L_1$ , MCP, SCAD and TLP, respectively. Finally, updating the Lagrange multipliers  $\boldsymbol{\nu}$  takes a standard form in general ADMM algorithms,

$$\boldsymbol{\nu}^{[k+1]} = \boldsymbol{\nu}^{[k]} + \kappa (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+1]} - \boldsymbol{\delta}^{[k+1]}). \quad (\text{S.24})$$

The above discussions are summarized in Algorithm 1.

As the initial values of  $(\boldsymbol{\Theta}^{[0]}, \boldsymbol{\delta}^{[0]}, \boldsymbol{\nu}^{[0]})$ , we use  $\boldsymbol{\Theta}^{[0]} = \hat{\boldsymbol{\Theta}}$ ,  $\boldsymbol{\delta}^{[0]} = \mathbf{B}^\top \hat{\boldsymbol{\Theta}}$  and  $\boldsymbol{\nu}^{[0]} = \mathbf{0}$  in our implementation. We also fix  $\kappa = 1$  because the ADMM method can be shown to converge for all values of the penalty parameter  $\kappa > 0$  (e.g., Eckstein (2012); Ghadimi et al. (2015)).

### Detailed derivation for (S.29), (S.30), (S.31) and (S.32).

Given the values of  $\boldsymbol{\delta}^{[k]}$  and  $\boldsymbol{\nu}^{[k]}$  in the  $k$ -th iteration, the ADMM that solves the augmented Lagrangian (S.21) goes as follows:

$$(\boldsymbol{\beta}^{[k+1]}, \boldsymbol{\Theta}^{[k+1]}) = \arg \min_{\boldsymbol{\Theta}} L(\boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\delta}^{[k]}, \boldsymbol{\nu}^{[k]}), \quad (\text{S.25})$$

$$\boldsymbol{\delta}^{[k+1]} = \arg \min_{\boldsymbol{\delta}} L(\boldsymbol{\beta}^{[k+1]}, \boldsymbol{\Theta}^{[k+1]}, \boldsymbol{\delta}, \boldsymbol{\nu}^{[k]}), \quad (\text{S.26})$$

and  $\boldsymbol{\nu}^{[k+1]}$  is updated using (S.24).



Firstly, we rewrite the augmented Lagrangian (S.21) by expanding the first term as

$$\begin{aligned}
 L(\boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\delta}, \boldsymbol{\nu}) &= \frac{1}{2N} \sum_{i=1}^M (\check{\boldsymbol{\beta}}_i - \boldsymbol{\beta})^\top \mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i (\check{\boldsymbol{\beta}}_i - \boldsymbol{\beta}) + \frac{1}{2N} \sum_{i=1}^M (\check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i)^\top \mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i (\check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i) \\
 &\quad + \frac{1}{N} \sum_{i=1}^M (\check{\boldsymbol{\beta}}_i - \boldsymbol{\beta})^\top \mathbf{x}_i^\top \mathbf{W}_i \mathbf{z}_i (\check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i) + \sum_{1 \leq i < j \leq M} p_\gamma (\|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|; \lambda) \\
 &\quad + \sum_{1 \leq i < j \leq M} \boldsymbol{\nu}_{ij}^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij}) + \frac{\kappa}{2} \sum_{1 \leq i < j \leq M} \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij}\|^2 \\
 &= \frac{1}{2N} (\check{\boldsymbol{\beta}} - \mathbf{1}_M \otimes \boldsymbol{\beta})^\top \mathbf{D}_x (\check{\boldsymbol{\beta}} - \mathbf{1}_M \otimes \boldsymbol{\beta}) + \frac{1}{2N} (\check{\boldsymbol{\theta}} - \boldsymbol{\Theta})^\top \mathbf{Z}^\top \mathbf{W} \mathbf{Z} (\check{\boldsymbol{\theta}} - \boldsymbol{\Theta}) \\
 &\quad + \frac{1}{N} (\check{\boldsymbol{\theta}} - \boldsymbol{\Theta})^\top \mathbf{D}_{zx} (\check{\boldsymbol{\beta}} - \mathbf{1}_M \otimes \boldsymbol{\beta}) + \sum_{i < j} p_\gamma (\|\boldsymbol{\delta}_{ij}\|; \lambda) \\
 &\quad + (\mathbf{B}^\top \boldsymbol{\Theta} - \boldsymbol{\delta})^\top \boldsymbol{\nu} + \frac{\kappa}{2} (\mathbf{B}^\top \boldsymbol{\Theta} - \boldsymbol{\delta})^\top (\mathbf{B}^\top \boldsymbol{\Theta} - \boldsymbol{\delta}),
 \end{aligned} \tag{S.27}$$

where  $\mathbf{D}_x = \text{diag} [(\mathbf{x}_i^\top \mathbf{W}_i \mathbf{x}_i)_{i=1, \dots, M}]$ ,  $\mathbf{D}_{zx} = \text{diag} [(\mathbf{z}_i^\top \mathbf{W}_i \mathbf{x}_i)_{i=1, \dots, M}]$  and  $\mathbf{B}$  is a matrix such that  $\mathbf{B}^\top \boldsymbol{\Theta} = (\boldsymbol{\theta}_i^\top - \boldsymbol{\theta}_j^\top)_{i < j}$ . Given  $\boldsymbol{\Theta}^{[k]}$ ,  $\boldsymbol{\delta}^{[k]}$  and  $\boldsymbol{\nu}^{[k]}$ , solving the equation using (S.27)

$$\frac{\partial L(\boldsymbol{\beta}, \boldsymbol{\Theta}^{[k]}, \boldsymbol{\delta}^{[k]}, \boldsymbol{\nu}^{[k]})}{\partial \boldsymbol{\beta}} = -\frac{1}{N} \sum_{i=1}^M \mathbf{x}_i^\top \mathbf{W}_i \left[ \mathbf{x}_i (\check{\boldsymbol{\beta}}_i - \boldsymbol{\beta}) + \mathbf{z}_i (\check{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i^{[k]}) \right] = 0$$

yields the updating formula (S.27) for  $\boldsymbol{\beta}$ . Similarly, given  $\boldsymbol{\beta}^{[k+1]}$ ,  $\boldsymbol{\delta}^{[k]}$  and  $\boldsymbol{\nu}^{[k]}$ , setting the partial derivative w.r.t.  $\boldsymbol{\Theta}$  using (S.28),

$$\frac{\partial L(\boldsymbol{\beta}^{[k+1]}, \boldsymbol{\Theta}, \boldsymbol{\delta}^{[k]}, \boldsymbol{\nu}^{[k]})}{\partial \boldsymbol{\Theta}} = -\frac{1}{N} \left[ \mathbf{Z}^\top \mathbf{W} \mathbf{Z} (\check{\boldsymbol{\theta}} + \boldsymbol{\Theta}) - \mathbf{D}_{zx} (\check{\boldsymbol{\beta}} - \mathbf{1}_M \otimes \boldsymbol{\beta}^{[k+1]}) \right] + \mathbf{B} \boldsymbol{\nu}^{[k]} + \kappa \mathbf{B} (\mathbf{B}^\top \boldsymbol{\Theta} - \boldsymbol{\delta}^{[k]})$$

equal to zero leads to the updating formula (S.23) for  $\boldsymbol{\Theta}$ .

Now, given  $\boldsymbol{\beta}^{[k+1]}$ ,  $\boldsymbol{\Theta}^{[k+1]}$  and  $\boldsymbol{\nu}^{[k]}$ , we derive the updating formula for  $\boldsymbol{\delta}$ . By omitting irrelevant terms w.r.t.  $\boldsymbol{\delta}$ , (S.27) results in

$$\sum_{i < j} \left[ p_\gamma (\|\boldsymbol{\delta}_{ij}\|; \lambda) + \frac{\kappa}{2} \|\boldsymbol{\zeta}_{ij}^{[k+1]} - \boldsymbol{\delta}_{ij}\|^2 \right],$$

where  $\boldsymbol{\zeta}_{ij}^{[k+1]} = \boldsymbol{\theta}_i^{[k+1]} - \boldsymbol{\theta}_j^{[k+1]} + \kappa^{-1} \boldsymbol{\nu}_{ij}^{[k]}$ . Its minimizer is taken as the update of  $\boldsymbol{\delta}$ . In particular, the solution for the  $L_1$  penalty is

$$\boldsymbol{\delta}_{ij}^{[k+1]} = \text{ST} (\boldsymbol{\zeta}_{ij}^{[k+1]}; \lambda/\kappa), \tag{S.29}$$

where  $\text{ST}(\mathbf{v}; t) = (1 - t/\|\mathbf{v}\|)_+ \mathbf{v}$  is the soft thresholding operator. For the MCP given in (S.2) with  $\gamma > \kappa^{-1}$ , the solution is

$$\boldsymbol{\delta}_{ij}^{[k+1]} = \begin{cases} \frac{\text{ST}(\boldsymbol{\zeta}_{ij}^{[k+1]}; \lambda/\kappa)}{1 - 1/(\gamma\kappa)} & \text{if } \|\boldsymbol{\zeta}_{ij}^{[k+1]}\| \leq \gamma\lambda, \\ \boldsymbol{\zeta}_{ij}^{[k+1]} & \text{if } \|\boldsymbol{\zeta}_{ij}^{[k+1]}\| > \gamma\lambda. \end{cases} \tag{S.30}$$

For the SCAD penalty given in (S.3) with  $\gamma > \kappa^{-1} + 1$ , the solution is

$$\boldsymbol{\delta}_{ij}^{[k+1]} = \begin{cases} \text{ST}(\boldsymbol{\zeta}_{ij}^{[k+1]}; \lambda/\kappa) & \text{if } \|\boldsymbol{\zeta}_{ij}^{[k+1]}\| \leq \lambda + \lambda/\kappa, \\ \frac{\text{ST}(\boldsymbol{\zeta}_{ij}^{[k+1]}; \gamma\lambda/[(\gamma-1)\kappa])}{1 - 1/[(\gamma-1)\kappa]} & \text{if } \lambda + \lambda/\kappa < \|\boldsymbol{\zeta}_{ij}^{[k+1]}\| \leq \gamma\lambda, \\ \boldsymbol{\zeta}_{ij}^{[k+1]} & \text{if } \|\boldsymbol{\zeta}_{ij}^{[k+1]}\| > \gamma\lambda. \end{cases} \tag{S.31}$$

For the TLP given in (S.4) with  $\gamma > \kappa^{-1}$ , the solution is

$$\delta_{ij}^{[k+1]} = \begin{cases} \text{ST}(\zeta_{ij}^{[k+1]}; \lambda/\kappa) & \text{if } \|\zeta_{ij}^{[k+1]}\| \leq \gamma\lambda, \\ \zeta_{ij}^{[k+1]} & \text{if } \|\zeta_{ij}^{[k+1]}\| > \gamma\lambda. \end{cases} \quad (\text{S.32})$$

---

**Algorithm 1** ADMM for CD Fusion Approach
 

---

- 1: **Input:** Initial values of  $\Theta^{[0]}$ ,  $\delta^{[0]}$  and  $\nu^{[0]}$
  - 2: **for**  $k = 1, 2, \dots$  **do**
  - 3:   Compute  $\beta^{[k+1]}$  using (S.22)
  - 4:   Compute  $\Theta^{[k+1]}$  using (S.23)
  - 5:   Compute  $\delta^{[k+1]}$  using (S.29), (S.30), (S.31) or (S.32)
  - 6:   Compute  $\nu^{[k+1]}$  using (S.24)
  - 7:   **if** Convergence criterion is met **then**
  - 8:     Break
  - 9:   **else**
  - 10:     $k \leftarrow k + 1$
  - 11:   **end if**
  - 12: **end for**
  - 13: **Input:**  $\hat{\beta}(\lambda) = \beta^{[k+1]}$  and  $\hat{\Theta}(\lambda) = \Theta^{[k+1]}$
- 

**Remark S.5.1.** Here we stress that, though  $\delta = (\delta_{ij}^\top)_{i < j}^\top$ ,  $\nu = (\nu_{ij}^\top)_{i < j}^\top$  and the matrix  $B$  such that  $B^\top \Theta = (\theta_i^\top - \theta_j^\top)_{i < j}^\top$  are just auxiliary parameters in the ADMM algorithm, one has to be cautious with matching their forms. More specifically,  $\nu_{ij}$ 's are the corresponding Lagrangian multipliers for  $\delta_{ij}$ , which are made to match  $\theta_i - \theta_j$ . For example, suppose  $M = 4$ , and one may have

$$\delta = \begin{pmatrix} \delta_{12} \\ \delta_{13} \\ \delta_{14} \\ \delta_{23} \\ \delta_{24} \\ \delta_{34} \end{pmatrix}, \quad \nu = \begin{pmatrix} \nu_{12} \\ \nu_{13} \\ \nu_{14} \\ \nu_{23} \\ \nu_{24} \\ \nu_{34} \end{pmatrix} \quad \text{and} \quad B^\top = \begin{bmatrix} I_q & -I_q & & & & \\ I_q & & -I_q & & & \\ I_q & & & -I_q & & \\ & I_q & -I_q & & -I_q & \\ & I_q & & -I_q & & \\ & & I_q & -I_q & & \end{bmatrix} \quad \text{s.t.} \quad B^\top \Theta = \begin{pmatrix} \theta_1 - \theta_2 \\ \theta_1 - \theta_3 \\ \theta_1 - \theta_4 \\ \theta_2 - \theta_3 \\ \theta_2 - \theta_4 \\ \theta_3 - \theta_4 \end{pmatrix}.$$

However, it is also possible to set

$$\delta = \begin{pmatrix} \delta_{12} \\ \delta_{23} \\ \delta_{34} \\ \delta_{13} \\ \delta_{24} \\ \delta_{14} \end{pmatrix}, \quad \nu = \begin{pmatrix} \nu_{12} \\ \nu_{23} \\ \nu_{34} \\ \nu_{13} \\ \nu_{24} \\ \nu_{14} \end{pmatrix} \quad \text{and} \quad B^\top = \begin{bmatrix} I_q & -I_q & & & & \\ & I_q & -I_q & & & \\ & & I_q & -I_q & & \\ I_q & & & -I_q & & \\ & I_q & & -I_q & & \\ I_q & & & -I_q & & \end{bmatrix} \quad \text{s.t.} \quad B^\top \Theta = \begin{pmatrix} \theta_1 - \theta_2 \\ \theta_2 - \theta_3 \\ \theta_3 - \theta_4 \\ \theta_1 - \theta_3 \\ \theta_2 - \theta_4 \\ \theta_1 - \theta_4 \end{pmatrix}.$$

As can be seen that the order of  $\nu_{ij}$ 's and the form of  $B$  both depend on the order of  $\delta_{ij}$ 's. Although it is not necessary to have guidelines to make their forms unique, while matching them correspondingly is crucial in implementation.

## S.6. Proof of Proposition 3.1

Here we show  $\lim_{k \rightarrow \infty} \|r^{[k]}\|^2 = 0$ . Define

$$\begin{aligned} f^{[k+1]} &= \inf_{B^\top \Theta^{[k+1]} - \delta = 0} \left\{ L_N^{\text{CD}}(\beta^{[k+1]}, \Theta^{[k+1]}) + \sum_{i < j} p_\gamma(\|\delta\|; \lambda) \right\} \\ &= \inf_{B^\top \Theta^{[k+1]} - \delta = 0} L(\beta^{[k+1]}, \Theta^{[k+1]}, \delta, \nu^{[k]}), \end{aligned}$$

where  $L_N^{\text{CD}}(\boldsymbol{\beta}, \boldsymbol{\Theta})$  is defined in (S.5). By the definition of  $\boldsymbol{\delta}^{[k+1]}$  as in (S.26), we have

$$L(\boldsymbol{\beta}^{[k+1]}, \boldsymbol{\Theta}^{[k+1]}, \boldsymbol{\delta}^{[k+1]}, \boldsymbol{\nu}^{[k]}) \leq f^{[k+1]}.$$

Let  $t$  be an integer. Since  $\boldsymbol{\nu}^{[k+t-1]} = \boldsymbol{\nu}^{[k]} + \kappa \sum_{i=1}^{t-1} (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+i]} - \boldsymbol{\delta}^{[k+i]})$ , it follows that

$$\begin{aligned} L(\boldsymbol{\beta}^{[k+t]}, \boldsymbol{\Theta}^{[k+t]}, \boldsymbol{\delta}^{[k+t]}, \boldsymbol{\nu}^{[k+t-1]}) &= L_N^{\text{CD}}(\boldsymbol{\beta}^{[k+t]}, \boldsymbol{\Theta}^{[k+t]}) + \sum_{i < j} p_\gamma(\|\boldsymbol{\delta}^{[k+t]}\|; \lambda) \\ &\quad + (\boldsymbol{\nu}^{[k+t-1]})^\top (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+t]} - \boldsymbol{\delta}^{[k+t]}) + \frac{\kappa}{2} \left\| \mathbf{B}^\top \boldsymbol{\Theta}^{[k+t]} - \boldsymbol{\delta}^{[k+t]} \right\|^2 \\ &= L_N^{\text{CD}}(\boldsymbol{\beta}^{[k+t]}, \boldsymbol{\Theta}^{[k+t]}) + \sum_{i < j} p_\gamma(\|\boldsymbol{\delta}^{[k+t]}\|; \lambda) \\ &\quad + (\boldsymbol{\nu}^{[k]})^\top (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+t]} - \boldsymbol{\delta}^{[k+t]}) \\ &\quad + \kappa \sum_{i=1}^{t-1} (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+i]} - \boldsymbol{\delta}^{[k+i]})^\top (\mathbf{B}^\top \boldsymbol{\Theta}^{[k+t]} - \boldsymbol{\delta}^{[k+t]}) \\ &\quad + \frac{\kappa}{2} \left\| \mathbf{B}^\top \boldsymbol{\Theta}^{[k+t]} - \boldsymbol{\delta}^{[k+t]} \right\|^2 \\ &\leq f^{[k+t]}. \end{aligned}$$

We are going to show that  $\lim_{k \rightarrow \infty} \|\mathbf{r}^{[k]}\|^2 = 0$  is a necessary condition for the inequality

$$\lim_{k \rightarrow \infty} L(\boldsymbol{\beta}^{[k+t]}, \boldsymbol{\Theta}^{[k+t]}, \boldsymbol{\delta}^{[k+t]}, \boldsymbol{\nu}^{[k+t-1]}) \leq \lim_{k \rightarrow \infty} f^{[k+t]}$$

for all  $t \geq 0$ .

By applying the results in Theorem 4.1 of (Tseng, 2001), the sequence  $(\boldsymbol{\beta}^{[k]}, \boldsymbol{\Theta}^{[k]}, \boldsymbol{\delta}^{[k]})$  has a limit point, denoted by  $(\boldsymbol{\beta}^{[*]}, \boldsymbol{\Theta}^{[*]}, \boldsymbol{\delta}^{[*]})$ , since the objective function  $L(\boldsymbol{\beta}, \boldsymbol{\Theta}, \boldsymbol{\delta}, \boldsymbol{\nu})$  is differentiable w.r.t.  $(\boldsymbol{\beta}, \boldsymbol{\Theta})$  and is convex w.r.t.  $\boldsymbol{\delta}$ . Define

$$f^{[*]} = \lim_{k \rightarrow \infty} f^{[k+1]} = \lim_{k \rightarrow \infty} f^{[k+t]} = \inf_{\mathbf{B}^\top \boldsymbol{\Theta}^{[*]} - \boldsymbol{\delta}^{[*]} = \mathbf{0}} \left\{ L_N^{\text{CD}}(\boldsymbol{\beta}^{[*]}, \boldsymbol{\Theta}^{[*]}) + \sum_{i < j} p_\gamma(\|\boldsymbol{\delta}^{[*]}\|; \lambda) \right\}.$$

For all  $t \geq 0$ , we have

$$\begin{aligned} \lim_{k \rightarrow \infty} L(\boldsymbol{\beta}^{[k+t]}, \boldsymbol{\Theta}^{[k+t]}, \boldsymbol{\delta}^{[k+t]}, \boldsymbol{\nu}^{[k+t-1]}) &= L_N^{\text{CD}}(\boldsymbol{\beta}^{[*]}, \boldsymbol{\Theta}^{[*]}) + \sum_{i < j} p_\gamma(\|\boldsymbol{\delta}_{ij}^{[*]}\|; \lambda) \\ &\quad + \lim_{k \rightarrow \infty} (\boldsymbol{\nu}^{[k]})^\top (\mathbf{B}^\top \boldsymbol{\Theta}^{[*]} - \boldsymbol{\delta}^{[*]}) + (t - \frac{1}{2})\kappa \left\| \mathbf{B}^\top \boldsymbol{\Theta}^{[*]} - \boldsymbol{\delta}^{[*]} \right\|^2 \\ &\leq f^{[*]}. \end{aligned}$$

Note that  $f^{[*]}$  is a constant by definition, and thus we must have  $\lim_{k \rightarrow \infty} \|\mathbf{r}^{[k]}\|^2 = \|\mathbf{B}^\top \boldsymbol{\Theta}^{[*]} - \boldsymbol{\delta}^{[*]}\|^2 = 0$  to make the last inequality hold for all  $t \geq 0$ .

## S.7. Proof of Theorem 4.2

In this section, we prove the results in Theorem 4.2. We start with introducing the following notations. The objective functions for the proposed and oracle estimators can be written respectively by

$$Q_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) = L_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) + P_N(\boldsymbol{\Theta}) \quad \text{and} \quad Q_N^G(\boldsymbol{\beta}, \boldsymbol{\alpha}) = L_N^G(\boldsymbol{\beta}, \boldsymbol{\alpha}) + P_N^G(\boldsymbol{\alpha}),$$

where

$$\begin{aligned} L_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) &= \frac{1}{2N} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Theta})^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\boldsymbol{\Theta}), \quad P_N(\boldsymbol{\Theta}) = \lambda \sum_{1 \leq i < j \leq N} \rho(\|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|), \\ L_N^G(\boldsymbol{\beta}, \boldsymbol{\alpha}) &= \frac{1}{2N} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{A}\boldsymbol{\alpha})^\top \mathbf{W} (\mathbf{Y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{Z}\mathbf{A}\boldsymbol{\alpha}), \quad P_N^G(\boldsymbol{\Theta}) = \lambda \sum_{1 \leq s < s' \leq S} M_s M_{s'} \rho(\|\boldsymbol{\alpha}_s - \boldsymbol{\alpha}_{s'}\|), \end{aligned}$$

$\rho(t) = \lambda^{-1}p_\gamma(t; \lambda)$  and  $M_s$  is the number of units in the  $s$ -th subgroup. Recall that the set

$$\mathcal{M}_G = \{\Theta \in \mathbb{R}^{Mq} : \theta_i = \theta_j \text{ for any } L_i = L_j = s, s = 1, \dots, S\}$$

is the collection of stacked heterogeneous effect vector  $\Theta = (\theta_i^\top)_{i=1, \dots, M}^\top$  with only  $S$  distinct values of  $\theta_i$ 's. For each  $\theta \in \mathcal{M}_G$ , there exists a  $Mq \times Sq$  label matrix  $\mathbb{A}$  and  $\alpha \in \mathbb{R}^{Sq}$  such that  $\Theta = \mathbb{A}\alpha$  and  $\alpha = (\mathbb{A}^\top \mathbb{A})^{-1} \mathbb{A}^\top \Theta$ . Hence, let  $T : \mathcal{M}_G \rightarrow \mathbb{R}^{Sq}$  be the mapping such that  $T(\Theta)$  is the  $(Sq)$ -vector stacking the  $S$  different  $\theta_i$ 's in  $\Theta$ , and its inverse mapping  $T^{-1} : \mathbb{R}^{Sq} \rightarrow \mathcal{M}_G$  is well-defined. Let  $T^* : \mathbb{R}^{Mq} \rightarrow \mathbb{R}^{Sq}$  be the mapping such that  $T^*(\Theta) = \{M_s^{-1} \sum_{i:L_i=s} \theta_i\}_{s=1}^S$ . For any  $\Theta \in \mathbb{R}^{Mq}$ , define  $\Theta^* = T^{-1}(T^*(\Theta))$ .

By the above definitions, it clearly follows that  $T(\Theta) = T^*(\Theta)$  and  $P_N(\Theta) = P_N^G(T(\Theta))$  for every  $\Theta \in \mathcal{M}_G$ , and that  $P_N(T^{-1}(\alpha)) = P_N^G(\alpha)$  for every  $\alpha \in \mathbb{R}^{Sq}$ . Moreover, we can write

$$\begin{aligned} Q_N(\beta, \Theta) &= Q_N^G(\beta, T(\Theta)) \quad \text{for every } \Theta \in \mathcal{M}_G, \\ Q_N^G(\beta, \alpha) &= Q_N(\beta, T^{-1}(\alpha)) \quad \text{for every } \alpha \in \mathbb{R}^{Sq}. \end{aligned} \quad (\text{S.33})$$

Consider a neighborhood of  $(\beta_0, \Theta_0)$ :

$$\Upsilon = \left\{ \beta \in \mathbb{R}^p, \Theta \in \mathbb{R}^{Mq} : \|\beta - \beta_0\| \leq \phi_N, \max_{1 \leq i \leq M} \|\theta_i - \theta_{i0}\| \leq \phi_N \right\}.$$

By Theorem 4.1, there is an event  $E_1$  with  $P(E_1) \geq 1 - (p + Sq)(4N^{-1} + e^{-g_{\max}}) - 2Me^{-n_{\min}}$  such that  $(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}}) \in \Upsilon$  on event  $E_1$ . We show that  $(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$  is a strict local minimizer of the objective function (5) with probability tending to 1 through the following two steps.

- (i) On the event  $E_1$ ,  $Q_N(\beta, \Theta^*) > Q_N(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$  for any  $(\beta, \Theta) \in \Upsilon$  and  $(\beta, \Theta^*) \neq (\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$ .
- (ii) There is an event  $E_2$  with  $P(E_2) \geq 1 - Me^{-g_{\min}}$ . On  $E_1 \cap E_2$ , there is a neighborhood of  $(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$ , denoted by  $\Upsilon_N$ , such that  $Q_N(\beta, \Theta) \geq Q_N(\beta, \Theta^*)$  for any  $(\beta, \Theta) \in \Upsilon_N \cap \Upsilon$  for sufficiently large  $N$ .

By results (i) and (ii), we have  $Q_N(\beta, \Theta) > Q_N(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$  for any  $(\beta, \Theta) \in \Upsilon_N \cap \Upsilon$  and  $(\beta, \Theta^*) \neq (\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$  on the event  $E_1 \cap E_2$ , so that  $(\hat{\beta}_{\text{OR}}, \hat{\Theta}_{\text{OR}})$  is a strict local minimizer of  $Q_N(\beta, \Theta)$  given in (5) over the event  $E_1 \cap E_2$  with  $P(E_1 \cap E_2) \geq 1 - (p + Sq)(4N^{-1} + e^{-g_{\max}}) - M(2e^{-n_{\min}} + e^{-g_{\min}})$  for sufficiently large  $N$ .

We first show  $P_N^G(T^*(\Theta)) = K_N$  for any  $\Theta \in \Upsilon$ , where  $K_N$  is a constant independent of  $\Theta$ . Let  $T^*(\Theta) = \alpha = (\alpha_s^\top)_{s=1, \dots, S}^\top$ . Since  $\rho(t) = \lambda^{-1}p_\gamma(t; \lambda)$  is constant for  $t \geq a\lambda$  according to Assumption 4.3, it suffices to show that  $\|\alpha_s - \alpha_{s'}\| > a\lambda$  for all  $s, s' \in \{1, \dots, S\}$ . Since

$$\begin{aligned} \max_{1 \leq s \leq S} \|\alpha_s - \alpha_{s0}\| &= \max_{1 \leq s \leq S} \left\| \frac{1}{M_s} \sum_{i:L_i=s} (\theta_i - \theta_{i0}) \right\| \\ &\leq \max_{1 \leq s \leq S} \frac{1}{M_s} \sum_{i:L_i=s} \|\theta_i - \theta_{i0}\| \\ &\leq \max_{1 \leq s \leq S} \|\theta_i - \theta_{i0}\| \\ &\leq \phi_N, \end{aligned} \quad (\text{S.34})$$

then for all  $s, s' \in \{1, \dots, S\}$  we have

$$\|\alpha_s - \alpha_{s'}\| \geq \|\alpha_{s0} - \alpha_{s'0}\| - 2 \max_{1 \leq s \leq S} \|\alpha_s - \alpha_{s0}\| \geq \Delta_N - 2\phi_N > a\lambda, \quad (\text{S.35})$$

where the last inequality follows from the assumption that  $\Delta_N \gg a\lambda \gg \phi_N$ . Accordingly, we have  $P_N^G(T^*(\Theta)) = K_N$  for some constant  $K_N > 0$  and hence  $Q_N^G(\beta, T^*(\Theta)) = L_N^G(\beta, T^*(\Theta)) + K_N$  for all  $(\beta, \Theta) \in \Upsilon$ . Since  $(\hat{\beta}_{\text{OR}}, \hat{\alpha}_{\text{OR}})$  is the unique global minimizer of  $L_N^G(\beta, \alpha)$ , we have  $L_N^G(\beta, T^*(\Theta)) > L_N^G(\hat{\beta}_{\text{OR}}, \hat{\alpha}_{\text{OR}})$  for all  $(\beta, T^*(\Theta)) \neq (\hat{\beta}_{\text{OR}}, \hat{\alpha}_{\text{OR}})$ .

It then follows that  $Q_N^{\mathcal{G}}(\beta, T^*(\Theta)) > Q_N^{\mathcal{G}}(\widehat{\beta}_{\text{OR}}, \widehat{\alpha}_{\text{OR}})$  for all  $T^*(\Theta) \neq \widehat{\alpha}_{\text{OR}}$ . Let  $\widehat{\Theta}_{\text{OR}} = \mathbf{A}\widehat{\alpha}_{\text{OR}}$ , and by (S.33) we have  $Q_N^{\mathcal{G}}(\widehat{\beta}_{\text{OR}}, \widehat{\alpha}_{\text{OR}}) = Q_N(\widehat{\beta}_{\text{OR}}, \widehat{\Theta}_{\text{OR}})$  and  $Q_N^{\mathcal{G}}(\beta, T^*(\Theta)) = Q_N(\beta, T^{-1}(T^*(\Theta))) = Q_N(\beta, \Theta^*)$ . Consequently, we have  $Q_N(\beta, \Theta^*) > Q_N(\widehat{\beta}_{\text{OR}}, \widehat{\Theta}_{\text{OR}})$  for all  $(\beta, \Theta) \in \Upsilon$  and  $(\beta, \Theta^*) \neq (\widehat{\beta}_{\text{OR}}, \widehat{\Theta}_{\text{OR}})$ , and the result (i) follows.

Next we prove the result in (ii). For positive sequence  $\{t_N\}$ , let

$$\Upsilon_N = \{\Theta : \max_{1 \leq i \leq M} \|\theta_i - \widehat{\Theta}_{\text{OR},i}\| \leq t_N\}.$$

For  $(\beta, \Theta) \in \Upsilon_N \cap \Upsilon$ , Taylor's expansion leads to

$$Q_N(\beta, \Theta) - Q_N(\beta, \Theta^*) = \underbrace{-N^{-1} \sum_{i=1}^M (\mathbf{y}_i - \mathbf{x}_i \beta - \mathbf{z}_i \tilde{\theta}_i)^\top \mathbf{W}_i \mathbf{z}_i (\theta_i - \theta_i^*)}_{\Gamma_1} + \underbrace{\sum_{i=1}^M \frac{\partial P_N(\tilde{\Theta})}{\partial \theta_i} (\theta_i - \theta_i^*)}_{\Gamma_2},$$

where  $\tilde{\Theta} = \varsigma \Theta + (1 - \varsigma) \Theta^*$  for some  $\varsigma \in (0, 1)$ .

We firstly deal with  $\Gamma_2$ . Observe that

$$\begin{aligned} \Gamma_2 &= \lambda \sum_{i \neq j} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top (\theta_i - \theta_i^*) \\ &= \lambda \sum_{i < j} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top (\theta_i - \theta_i^*) + \lambda \sum_{i < j} \frac{\rho'(\|\tilde{\theta}_j - \tilde{\theta}_i\|)}{\|\tilde{\theta}_j - \tilde{\theta}_i\|} (\tilde{\theta}_j - \tilde{\theta}_i)^\top (\theta_j - \theta_j^*) \\ &= \lambda \sum_{i < j} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top (\theta_i - \theta_i^*) - \lambda \sum_{i < j} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top (\theta_j - \theta_j^*) \\ &= \lambda \sum_{i < j} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top [(\theta_i - \theta_i^*) - (\theta_j - \theta_j^*)]. \end{aligned} \quad (\text{S.36})$$

When  $L_i = L_j = s$ , we have  $\theta_i^* = \theta_j^*$  and  $\tilde{\theta}_i - \tilde{\theta}_j = \varsigma(\theta_i - \theta_j)$ , and therefore

$$\begin{aligned} \Gamma_2 &= \lambda \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|) \|\theta_i - \theta_j\| \\ &\quad + \lambda \sum_{s < s'} \sum_{\substack{L_i=s \\ L_j=s'}} \frac{\rho'(\|\tilde{\theta}_i - \tilde{\theta}_j\|)}{\|\tilde{\theta}_i - \tilde{\theta}_j\|} (\tilde{\theta}_i - \tilde{\theta}_j)^\top [(\theta_i - \theta_i^*) - (\theta_j - \theta_j^*)]. \end{aligned}$$

We look at the second term in  $\Gamma_2$ . Since  $\Theta \in \Upsilon$  and  $\Theta^* = T^{-1}(\alpha) \in \mathcal{M}_{\mathcal{G}}$ , it holds

$$\max_{1 \leq i \leq M} \|\theta_i - \theta_{i0}\| \leq \phi_N \quad \text{and} \quad \max_{1 \leq i \leq M} \|\theta_i^* - \theta_{i0}\| = \max_{1 \leq s \leq S} \|\alpha_s - \alpha_{s0}\| \leq \phi_N,$$

where the inequality follows from (S.34). Since  $\tilde{\Theta} = \varsigma \Theta + (1 - \varsigma) \Theta^*$ , we have

$$\max_{1 \leq i \leq M} \|\tilde{\theta}_i - \theta_{i0}\| \leq \varsigma \max_{1 \leq i \leq M} \|\theta_i - \theta_{i0}\| + (1 - \varsigma) \max_{1 \leq i \leq M} \|\theta_i^* - \theta_{i0}\| \leq \varsigma \phi_N + (1 - \varsigma) \phi_N = \phi_N.$$

Then for  $L_i = s$  and  $L_j = s'$  with  $s \neq s'$ , it follows that

$$\|\tilde{\theta}_i - \tilde{\theta}_j\| \geq \min_{L_i=s, L_j=s'} \|\theta_{i0} - \theta_{j0}\| - 2 \max_{1 \leq i \leq M} \|\tilde{\theta}_i - \theta_{i0}\| \geq \Delta_N - 2\phi_N \gg a\lambda,$$

and accordingly  $\rho'(\|\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}_j\|) = 0$ . As a result, the second term in  $\Gamma_2$  is zero and we have

$$\Gamma_2 = \lambda \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \rho'(\|\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}_j\|) \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|.$$

For the first term in  $\Gamma_2$ , since  $\boldsymbol{\Theta} \in \Upsilon_N$ , the same reasoning as (S.34) yields

$$\max_{1 \leq i \leq M} \|\boldsymbol{\theta}_i^* - \hat{\boldsymbol{\theta}}_{\text{OR},i}\| = \max_{1 \leq i \leq M} \|\boldsymbol{\alpha}_s - \hat{\boldsymbol{\alpha}}_{\text{OR},s}\| \leq \max_{1 \leq i \leq M} \|\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_{\text{OR},i}\| \leq t_N.$$

Since  $\tilde{\boldsymbol{\theta}}_i = \varsigma \boldsymbol{\theta}_i + (1 - \varsigma) \boldsymbol{\theta}_i^*$  for some  $\varsigma \in (0, 1)$ , it holds  $\tilde{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i^* = \varsigma(\boldsymbol{\theta}_i - \boldsymbol{\theta}_i^*)$ . For  $L_i = L_j$ , due to  $\boldsymbol{\theta}_i^* = \boldsymbol{\theta}_j^*$ , we have

$$\begin{aligned} \max_{1 \leq i \leq M} \|\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}_j\| &\leq \max_{1 \leq i \leq M} \left( \|\tilde{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i^*\| + \|\tilde{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j^*\| \right) \\ &\leq 2 \max_{1 \leq i \leq M} \|\tilde{\boldsymbol{\theta}}_i - \boldsymbol{\theta}_i^*\| \\ &\leq 2 \max_{1 \leq i \leq M} \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_i^*\| \\ &\leq 2 \left( \max_{1 \leq i \leq M} \|\boldsymbol{\theta}_i - \hat{\boldsymbol{\theta}}_{\text{OR},i}\| + \max_{1 \leq i \leq M} \|\boldsymbol{\theta}_i^* - \hat{\boldsymbol{\theta}}_{\text{OR},i}\| \right) \\ &\leq 4t_N. \end{aligned}$$

The concavity of  $\rho(\cdot)$  then implies that  $\rho'(\|\tilde{\boldsymbol{\theta}}_i - \tilde{\boldsymbol{\theta}}_j\|) \geq \rho'(4t_N)$ . In summary, we have

$$\Gamma_2 \geq \lambda \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \rho'(4t_N) \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|. \quad (\text{S.37})$$

Now we analyze  $\Gamma_1$ . For  $i = 1, \dots, M$ , let

$$\mathbf{h}_i = (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \tilde{\boldsymbol{\theta}}_i)^\top \mathbf{W}_i \mathbf{z}_i = \mathbf{z}_i^\top \mathbf{W}_i \left[ \mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i + \mathbf{x}_i (\boldsymbol{\beta}_0 - \boldsymbol{\beta}) + \mathbf{z}_i (\boldsymbol{\theta}_{i0} - \tilde{\boldsymbol{\theta}}_i) \right].$$

Recall that  $M_s$  is the number of units in the  $s$ -th subgroup. By the definition of  $\boldsymbol{\theta}_i^*$ , it holds

$$\boldsymbol{\theta}_i - \boldsymbol{\theta}_i^* = \boldsymbol{\theta}_i - \sum_{j:L_j=L_i} \frac{\boldsymbol{\theta}_j}{M_s} = \sum_{j:L_j=L_i} \frac{(\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{M_s}.$$

Then we have

$$\begin{aligned}
 \Gamma_1 &= -N^{-1} \sum_{i=1}^M (\mathbf{y}_i - \mathbf{x}_i \boldsymbol{\beta} - \mathbf{z}_i \tilde{\boldsymbol{\theta}}_i)^\top \mathbf{W}_i \mathbf{z}_i (\boldsymbol{\theta}_i - \boldsymbol{\theta}_i^*) \\
 &= -N^{-1} \sum_{i=1}^M \mathbf{h}_i^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_i^*) \\
 &= -N^{-1} \sum_{s=1}^S \sum_{L_i=L_j=s} \frac{\mathbf{h}_i^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{M_s} \\
 &= -N^{-1} \sum_{s=1}^S \sum_{L_i=L_j=s} \frac{\mathbf{h}_i^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{2M_s} - \sum_{s=1}^S \sum_{L_i=L_j=s} \frac{\mathbf{h}_j^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{2M_s} \\
 &= -N^{-1} \sum_{s=1}^S \sum_{L_i=L_j=s} \frac{(\mathbf{h}_i - \mathbf{h}_j)^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{2M_s} \\
 &= -N^{-1} \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \frac{(\mathbf{h}_i - \mathbf{h}_j)^\top (\boldsymbol{\theta}_i - \boldsymbol{\theta}_j)}{M_s} \\
 &\geq -N^{-1} \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\| \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|. \tag{S.38}
 \end{aligned}$$

Combining (S.37) and (S.38) yields

$$\begin{aligned}
 Q_N(\boldsymbol{\beta}, \boldsymbol{\Theta}) - Q_N(\boldsymbol{\beta}, \boldsymbol{\Theta}^*) &= \Gamma_1 + \Gamma_2 \geq \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \left[ \lambda \rho'(4t_N) - N^{-1} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\| \right] \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\| \\
 &\gtrsim \sum_{s=1}^S \sum_{\substack{L_i=L_j=s \\ i < j}} \left[ \lambda - N^{-1} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\| \right] \|\boldsymbol{\theta}_i - \boldsymbol{\theta}_j\|,
 \end{aligned}$$

where the last inequality holds by taking  $t_N = o(1)$  so that  $\rho'(4t_N) \rightarrow 1$ . It is then left to show  $\lambda > N^{-1} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\|$  to conclude the results in (ii).

Recall that  $\mathbf{W}_i = \mathbf{V}_i^2$ . On the event  $E_1$ , we can write

$$\begin{aligned}
 \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\| &\leq 2 \max_{1 \leq i \leq M} \|\mathbf{h}_i\| \\
 &\leq 2 \max_{1 \leq i \leq M} \left\{ \|\mathbf{V}_i \mathbf{z}_i\| \|\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)\| + (\|\mathbf{z}_i^\top \mathbf{W}_i \mathbf{x}_i\| + \|\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i\|) \phi_N \right\} \\
 &\leq 2 \max_{1 \leq i \leq M} \left\{ \|\mathbf{V}_i \mathbf{z}_i\| \|\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)\| + (\|\mathbf{V}_i \mathbf{z}_i\| \|\mathbf{V}_i\| \|\mathbf{x}_i\| + \|\mathbf{z}_i^\top \mathbf{W}_i \mathbf{z}_i\|) \phi_N \right\} \\
 &\lesssim 2 \max_{1 \leq i \leq M} \left\{ \tau^{-1} \|\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)\| + (\bar{C}_f^{1/2} \sigma_\varepsilon^{-1} \tau^{-1} n_i^{1/2} + \tau^{-2}) \phi_N \right\},
 \end{aligned}$$

where the last inequality follows from (S.17) and (S.15) under events  $E_{CB}$  and  $E_{EB}$ . Moreover, since  $\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)$  has sub-Gaussian tails according to (S.19), applying union bound and (S.42) yields

$$P\left( \max_{1 \leq i \leq M} \|\mathbf{V}_i(\mathbf{z}_i \mathbf{u}_i + \boldsymbol{\varepsilon}_i)\| > \sqrt{5c_5 g_{\min}} \right) \leq M e^{-g_{\min}},$$

where  $c_5 > 0$  is a constant. Then we have

$$\begin{aligned} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\| &\lesssim 2 \left[ \sqrt{5}c_5^{1/2} \tau^{-1} g_{\min}^{1/2} + \left( \overline{C}_f^{1/2} \sigma_\varepsilon^{-1} \tau^{-1} \max_{1 \leq i \leq M} n_i^{1/2} + \tau^{-2} \right) \phi_N \right] \\ &\leq 2g_{\min}^{1/2} \left[ \sqrt{5}c_5^{1/2} \tau^{-1} + \left( \overline{C}_f^{1/2} \sigma_\varepsilon^{-1} \tau^{-1} + \tau^{-2} \right) \phi_N \right] \end{aligned}$$

with probability at least  $1 - Me^{-g_{\min}}$ . Hence, there is an event  $E_2$  with  $P(E_2) \geq 1 - Me^{-g_{\min}}$  and, on the event  $E_2$ , we have

$$\frac{N^{-1} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\|}{\phi_N} \lesssim \frac{2N^{-1} g_{\min}^{1/2} \left[ \sqrt{5}c_5^{1/2} \tau^{-1} + \left( \overline{C}_f^{1/2} \sigma_\varepsilon^{-1} \tau^{-1} + \tau^{-2} \right) \phi_N \right]}{\phi_N} = O(N^{-1} g_{\min}^{1/2}) = o(1),$$

which indicates that

$$\lambda \gg \phi_N \gg N^{-1} \max_{i,j} \|\mathbf{h}_i - \mathbf{h}_j\|. \quad (\text{S.39})$$

In summary, we conclude the results in (ii) and the proof is complete.

### S.8. Proof of (ii) in Corollary 4.1

Since the asymptotic equivalence between the proposed estimator and the oracle estimator has been shown in Theorem 4.2, it suffices to show, for any  $p$ -vector  $\mathbf{v}_p$  and  $q$ -vector  $\mathbf{v}_q$ ,

$$\mathbf{v}_p^\top \text{Cov}(\widehat{\boldsymbol{\beta}}_{\text{OR}}) \mathbf{v}_p \leq \mathbf{v}_p^\top \text{Cov}(\check{\boldsymbol{\beta}}_i) \mathbf{v}_p \quad \text{and} \quad \mathbf{v}_q^\top \text{Cov}(\widehat{\boldsymbol{\theta}}_{\text{OR},i}) \mathbf{v}_q \leq \mathbf{v}_q^\top \text{Cov}(\check{\boldsymbol{\theta}}_i) \mathbf{v}_q, \quad (\text{S.40})$$

for all  $i = 1, \dots, M$ . Recall that

$$\text{Cov} \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} \end{pmatrix} = \left[ (\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A}) \right]^{-1} \quad \text{and} \quad \text{Cov} \begin{pmatrix} \check{\boldsymbol{\beta}}_i \\ \check{\boldsymbol{\theta}}_i \end{pmatrix} = \left[ (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) \right]^{-1}.$$

With straightforward matrix algebra, we have

$$(\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A}) = \sum_{i=1}^M \mathbf{G}_{L_i}^\top (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) \mathbf{G}_{L_i},$$

where  $\mathbf{G}_s$  are the  $(p+q) \times (p+Sq)$  matrices such that  $(\boldsymbol{\beta}_0^\top, \boldsymbol{\theta}_{i0}^\top)^\top = \mathbf{G}_{L_i} (\boldsymbol{\beta}_0^\top, \boldsymbol{\alpha}_{s0}^\top)^\top$ , for any parameter dimensions.

Without loss of generality, assume  $M = 2$ . When  $S = 1$ , we have  $\mathbf{G}_{L_1} = \mathbf{G}_{L_2} = \mathbf{I}_{p+q}$  and  $\widehat{\boldsymbol{\theta}}_{1,0} = \widehat{\boldsymbol{\theta}}_{2,0} = \widehat{\boldsymbol{\alpha}}_0$ . Hence, we have for any  $(p+q)$ -vector  $\mathbf{a}$ ,

$$\begin{aligned} \mathbf{a}^\top \text{Cov} \begin{pmatrix} \widehat{\boldsymbol{\beta}}_{\text{OR}} \\ \widehat{\boldsymbol{\alpha}}_{\text{OR}} \end{pmatrix} \mathbf{a} &= \mathbf{a}^\top \left[ (\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A}) \right]^{-1} \mathbf{a} \\ &= \mathbf{a}^\top \left[ (\mathbf{x}_1, \mathbf{z}_1)^\top \mathbf{W}_1(\mathbf{x}_1, \mathbf{z}_1) + (\mathbf{x}_2, \mathbf{z}_2)^\top \mathbf{W}_2(\mathbf{x}_2, \mathbf{z}_2) \right]^{-1} \mathbf{a} \\ &\leq \mathbf{a}^\top \left[ (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) \right]^{-1} \mathbf{a} \\ &= \mathbf{a}^\top \text{Cov} \begin{pmatrix} \check{\boldsymbol{\beta}}_i \\ \check{\boldsymbol{\theta}}_i \end{pmatrix} \mathbf{a}, \quad \text{for } i = 1, 2. \end{aligned}$$

where the inequality follows from Lemma S.9.3.

When  $S = 2$ , we calculate the expressions of  $\text{Cov}(\check{\boldsymbol{\beta}}_i)$ ,  $\text{Cov}(\check{\boldsymbol{\theta}}_i)$ ,  $\text{Cov}(\widehat{\boldsymbol{\beta}}_{\text{OR}})$  and  $\text{Cov}(\widehat{\boldsymbol{\theta}}_{\text{OR},i})$  in the following so as to establish the inequalities in (S.40). Firstly, for the unit estimators, let  $\mathbf{D}_i = \text{Cov}(\check{\boldsymbol{\beta}}_i)$  and  $\mathbf{H}_i = \text{Cov}(\check{\boldsymbol{\theta}}_i)$ , and then we can write

$$\text{Cov} \begin{pmatrix} \check{\boldsymbol{\beta}}_i \\ \check{\boldsymbol{\theta}}_i \end{pmatrix} = \left[ (\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) \right]^{-1} = \begin{bmatrix} \mathbf{D}_i & \mathbf{B}_i \\ \mathbf{B}_i^\top & \mathbf{H}_i \end{bmatrix}, \quad i = 1, 2,$$



where  $B_i$  are some  $p \times q$  matrices.

For the oracle estimators, suppose  $L_1 = 1$  and  $L_2 = 2$ . Then we have  $(\beta_0^\top, \alpha_0^\top)^\top = (\beta_0^\top, \theta_{1,0}, \theta_{2,0})^\top$ ,

$$G_1 = \begin{bmatrix} I_p & O_{p \times q} & O_{p \times q} \\ O_{p \times q}^\top & I_q & O_{q \times q} \end{bmatrix} \quad \text{and} \quad G_2 = \begin{bmatrix} I_p & O_{p \times q} & O_{p \times q} \\ O_{p \times q}^\top & O_{q \times q} & I_q \end{bmatrix},$$

where  $O_{d_1 \times d_2}$  stands for the  $d_1 \times d_2$  zero matrix. By blockwise matrix inversion, we can write

$$(\mathbf{x}_i, \mathbf{z}_i)^\top \mathbf{W}_i(\mathbf{x}_i, \mathbf{z}_i) = \begin{bmatrix} D_i^{-1} + D_i^{-1} B_i K_i B_i^\top D_i^{-1} & -D_i^{-1} B_i K_i \\ -K_i B_i^\top D_i^{-1} & K_i \end{bmatrix}, \quad i = 1, 2,$$

where  $K_i = (H_i - B_i^\top D_i^{-1} B_i)^{-1}$ . Accordingly, it follows that

$$(\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A}) = \begin{bmatrix} \sum_{i=1}^2 \{D_i^{-1} + D_i^{-1} B_i K_i B_i^\top D_i^{-1}\} & -D_1^{-1} B_1 K_1 & -D_2^{-1} B_2 K_2 \\ -K_1 B_1^\top D_1^{-1} & K_1 & O_{q \times q} \\ -K_2 B_2^\top D_2^{-1} & O_{q \times q} & K_2 \end{bmatrix}.$$

By the blockwise matrix inversion, we have

$$\begin{aligned} \text{Cov} \begin{pmatrix} \widehat{\beta}_{\text{OR}} \\ \widehat{\alpha}_{\text{OR}} \end{pmatrix} &= [(\mathbf{X}, \mathbf{Z}\mathbf{A})^\top \mathbf{W}(\mathbf{X}, \mathbf{Z}\mathbf{A})]^{-1} \\ &= \begin{bmatrix} Q & QD_1^{-1}B_1 & -QD_2^{-1}B_2 \\ B_1^\top D_1^{-1}Q & H_1 - B_1 D_1^{-1} B_1 + B_1^\top D_1^{-1} Q D_1^{-1} B_1 & B_1^\top D_1^{-1} Q D_2^{-1} B_2 \\ -B_2^\top D_2^{-1}Q & B_2^\top D_2^{-1} Q D_1^{-1} B_1 & H_2 - B_2 D_2^{-1} B_2 + B_2^\top D_2^{-1} Q D_2^{-1} B_2 \end{bmatrix}, \end{aligned}$$

where  $Q = (D_1^{-1} + D_2^{-1})^{-1}$ . As a result, we have  $\text{Cov}(\widehat{\beta}_{\text{OR}}) = Q = (D_1^{-1} + D_2^{-1})^{-1}$  and  $\text{Cov}(\widehat{\theta}_{\text{OR},i}) = H_i - B_i D_i^{-1} B_i + B_i^\top D_i^{-1} Q D_i^{-1} B_i, i = 1, 2$ .

With the above results, it is then straightforward that, for any  $p$ -vector  $\mathbf{v}_p$ ,

$$\mathbf{v}_p^\top \text{Cov}(\widehat{\beta}_{\text{OR}}) \mathbf{v}_p = \mathbf{v}_p^\top (D_1^{-1} + D_2^{-1})^{-1} \mathbf{v}_p \leq \mathbf{v}_p^\top D_i \mathbf{v}_p = \mathbf{v}_p^\top \text{Cov}(\check{\beta}_i) \mathbf{v}_p, \quad \text{for } i = 1, 2,$$

where the inequality follows from Lemma S.9.3. In addition, for any  $q$ -vector  $\mathbf{v}_q$ ,

$$\mathbf{v}_q^\top \text{Cov}(\widehat{\theta}_{\text{OR},i}) \mathbf{v}_q = \mathbf{v}_q^\top [H_i - B_i^\top D_i^{-1} B_i + B_i^\top D_i^{-1} (D_1^{-1} + D_2^{-1})^{-1} D_i^{-1} B_i] \mathbf{v}_q.$$

For the third term of the R.H.S., it holds

$$\mathbf{v}_q^\top B_i^\top D_i^{-1} (D_1^{-1} + D_2^{-1})^{-1} D_i^{-1} B_i \mathbf{v}_q \leq \mathbf{v}_q^\top B_i^\top D_i^{-1} D_i D_i^{-1} B_i \mathbf{v}_q = \mathbf{v}_q^\top B_i^\top D_i^{-1} B_i \mathbf{v}_q,$$

where the inequality follows from Lemma S.9.3. Consequently, we have

$$\mathbf{v}_q^\top \text{Cov}(\widehat{\theta}_{\text{OR},i}) \mathbf{v}_q \leq \mathbf{v}_q^\top H_i \mathbf{v}_q = \text{Cov}(\check{\theta}_i), \quad i = 1, 2.$$

The proof is thus complete.

## S.9. Auxiliary Lemmas

In this section we provide technical lemmas, with a bit abuse of notations.

**Lemma S.9.1** (Woodbury matrix identity, Theorem 18.2.8 of Harville (2000)). *Let  $A, U, B$  and  $V$  denote  $n \times n, n \times m, m \times m$  and  $m \times n$  matrices, respectively. Suppose  $A$  and  $B$  are nonsingular. Then we have*

$$(A \pm UBV)^{-1} = A^{-1} \mp A^{-1}U (B^{-1} \pm VA^{-1}U)^{-1} VA^{-1}. \quad (\text{S.41})$$

**Lemma S.9.2** (Lemma 8 of Hsu et al. (2014)). *Suppose  $\xi$  is a sub-Gaussian random  $n$ -vector. For all symmetric positive semidefinite matrices  $M \succeq \mathbf{O}$  and all  $t > 0$ , we have*

$$P \left[ \xi^\top M \xi > c \left( \text{tr}(M) + 2\sqrt{\text{tr}(M^2)t} + 2\|M\|t \right) \right] \leq e^{-t},$$

where  $c > 0$  is an absolute constant. With  $M = \mathbf{I}$  and  $t = n$ , it follows that

$$P(\|\xi\|^2 > 5cn) \leq e^{-n}. \quad (\text{S.42})$$

**Lemma S.9.3** (Lemma A.3 of Liu et al. (2015)). *Suppose  $M_1$  and  $M_2$  are  $d \times d$  positive definite matrices. Then, for any  $d$ -vector  $\mathbf{v}$ ,*

$$(\mathbf{v}^\top M_1^{-1} \mathbf{v})^{-1} + (\mathbf{v}^{-1} M_1^\top \mathbf{v})^{-1} \leq \left[ \mathbf{v}^\top (M_1 + M_2)^{-1} \mathbf{v} \right]^{-1}. \quad (\text{S.43})$$

This implies that  $\mathbf{v}^\top (M_1 + M_2)^{-1} \mathbf{v} \leq \mathbf{v}^\top M_1^{-1} \mathbf{v}$ .

## S.10. Signal-to-Noise Ratio for the proposed MTL method

Here we define the signal-to-noise ratio (SNR) for the proposed method. In order to supply information of how strong the subgroup effects compared to the noise, we first compute the signal for different subgroups  $i$  and  $j$  as

$$\begin{aligned} \text{Signal}(i, j) &= \text{Var}[(\mathbf{X}^\top \beta_0 + \mathbf{Z}^\top \alpha_i) - (\mathbf{X}^\top \beta_0 + \mathbf{Z}^\top \alpha_j)] \\ &= \text{Var}[\mathbf{Z}^\top (\alpha_i - \alpha_j)] \\ &= (\alpha_i - \alpha_j)^\top \Sigma_{\mathbf{Z}} (\alpha_i - \alpha_j), \end{aligned}$$

and then define the desired signal as the minimal signal variance between different subgroups, i.e.,

$$\text{Signal} = \min_{1 \leq i < j \leq S} \text{Signal}(i, j).$$

For noise, we calculate

$$\begin{aligned} \text{Noise} &= \text{Var}[(\mathbf{Z}^\top \mathbf{u}_i + \varepsilon_i) - (\mathbf{Z}^\top \mathbf{u}_i + \varepsilon_i)] \\ &= \text{Var}[\mathbf{Z}^\top (\mathbf{u}_j - \mathbf{u}_j)] + 2\text{Var}(\varepsilon) = 2[\text{tr}(\mathbf{E}[\mathbf{Z}\mathbf{Z}^\top] \Psi) + \sigma_\varepsilon^2]. \end{aligned}$$

Accordingly, we have

$$\text{SNR} = \frac{\text{Signal}}{\text{Noise}} = \frac{\min_{1 \leq i < j \leq S} \text{Signal}(i, j)}{2[\text{tr}(\mathbf{E}[\mathbf{Z}\mathbf{Z}^\top] \Psi) + \sigma_\varepsilon^2]}.$$

## S.11. Additional Simulation and Real Data Results.

### S.11.1. Simulation Results for all the cases.

We present the all-case version of Table 2.

Table S.1. Evaluation of subgroup recovery. Complete version of Table 2.

	Method	$\hat{S}$		NMI	Perfect
		Mean (SD)	Median (Min,Max)		Recovery
Case 1 ( $S = 2$ )	MCP	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	SCAD	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	TLP	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	$L_1$	2.39 (0.6651)	2 (1, 4)	0.9476	0.62
	$K$ -Means	2.05 (0.2179)	2 (2, 3)	0.9896	0.95
Case 2 ( $S = 3$ )	MCP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	SCAD	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	TLP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	$L_1$	2.40 (1.0731)	3 (1, 4)	0.6467	0.55
	$K$ -Means	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
Case 3 ( $S = 2$ )	MCP	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	SCAD	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	TLP	2.00 (0.0000)	2 (2, 2)	0.9988	0.99
	$L_1$	2.33 (0.5870)	2 (1, 4)	0.9718	0.70
	$K$ -Means	2.01 (0.0995)	2 (2, 3)	0.9979	0.99
Case 4 ( $S = 3$ )	MCP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	SCAD	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	TLP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	$L_1$	2.43 (1.2248)	3 (1, 6)	0.6141	0.47
	$K$ -Means	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
Case 5 ( $S = 2$ )	MCP	2.00 (0.0000)	2 (2, 2)	1.0000	1.00
	SCAD	2.01 (0.1000)	2 (2, 3)	0.9990	0.99
	TLP	2.01 (0.1000)	2 (2, 3)	0.9990	0.99
	$L_1$	2.27 (0.5835)	2 (1, 4)	0.9807	0.77
	$K$ -Means	3.36 (1.0538)	4 (2, 5)	0.7877	0.36
Case 6 ( $S = 3$ )	MCP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	SCAD	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	TLP	3.00 (0.0000)	3 (3, 3)	1.0000	1.00
	$L_1$	4.09 (9.7152)	3 (1, 100)	0.9303	0.72
	$K$ -Means	3.76 (0.8261)	4 (3, 6)	0.9303	0.48
Case 7 ( $S = 5$ )	MCP	5.01 (0.1000)	5 (5, 6)	0.9996	0.98
	SCAD	5.00 (0.0000)	5 (5, 5)	0.9997	0.99
	TLP	5.01 (0.1000)	5 (5, 6)	0.9996	0.99
	$L_1$	63.24 (47.1983)	100 (1, 100)	0.5878	0.00
	$K$ -Means	5.00 (0.0000)	5 (5, 5)	0.9768	0.37
Case 8 ( $S = 5$ )	MCP	5.01 (0.1000)	5 (5, 6)	0.9998	0.99
	SCAD	5.02 (0.1407)	5 (5, 6)	0.9998	0.98
	TLP	5.02 (0.1407)	5 (5, 6)	0.9998	0.98
	$L_1$	97.26 (70.6762)	150 (1, 150)	0.6137	0.00
	$K$ -Means	5.43 (0.5875)	5 (5, 7)	0.9550	0.16
Case 9 ( $S = 7$ )	MCP	7.00 (0.0000)	7 (7, 7)	0.9999	0.99
	SCAD	7.26 (0.5245)	7 (7, 9)	0.9981	0.77
	TLP	7.22 (0.5041)	7 (7, 9)	0.9982	0.81
	$L_1$	150.00 (0.0000)	150 (150, 150)	0.6205	0.00
	$K$ -Means	7.00 (0.0000)	7 (7, 7)	0.9786	0.20

### S.11.2. NOAA Data Features

In our analysis, each climate division is taken as a data unit. The monthly average temperature is the response of interest, and we include precipitation (PCPN), Palmer Drought Severity Index (PDSI, based on the principles of a balance between moisture supply and demand without considering man-made changes), Palmer Hydrological Drought Index (PHDI, based on principles similar to PDSI, but with consideration of some man-made changes) and Palmer Z Index (ZNDX, a moisture anomaly index) into the covariate pool. Negative and positive values of PDSI, PHDI and ZNDX indicate dry and wet spells, respectively. See the database documentation<sup>4</sup> for more details about these drought indices. To account for seasonal effects, we include dummy variables for the four seasons (March, April and May as spring; June, July and August as summer;

<sup>4</sup>Available on <ftp://ftp.ncdc.noaa.gov/pub/data/cirs/climdiv/divisional-readme.txt>.

Table S.2. Average ASD( $\hat{\alpha}_{1,1}$ ) and ESD( $\hat{\alpha}_{1,1}$ ), based on replicates with perfect subgroup structure recovery.

	Average ASD( $\hat{\alpha}_{1,1}$ )				ESD( $\hat{\alpha}_{1,1}$ )			
	MCP	SCAD	TLP	Oracle	MCP	SCAD	TLP	Oracle
<b>Case 1</b>	0.1123	0.1123	0.1123	0.1205	0.1166	0.1166	0.1166	0.1166
<b>Case 2</b>	0.1508	0.1508	0.1508	0.1531	0.1569	0.1569	0.1569	0.1569
<b>Case 3</b>	0.1244	0.1244	0.1244	0.1200	0.1117	0.1117	0.1117	0.1177
<b>Case 4</b>	0.1591	0.1591	0.1591	0.1525	0.1728	0.1728	0.1728	0.1728
<b>Case 5</b>	0.0781	0.0781	0.0781	0.0823	0.0871	0.0871	0.0871	0.0871
<b>Case 6</b>	0.0988	0.0988	0.0988	0.1025	0.1131	0.1131	0.1131	0.1131
<b>Case 7</b>	0.1319	0.1318	0.1319	0.1380	0.1285	0.1147	0.1282	0.1285
<b>Case 8</b>	0.1069	0.1066	0.1167	0.1104	0.1023	0.1028	0.1034	0.1030
<b>Case 9</b>	0.1284	0.1280	0.1284	0.1339	0.1414	0.1348	0.1448	0.1409

September, October and November as Fall; December, January and February as winter), and the spring effect is taken as the baseline intercept term. In summary, there are 8 covariates in our analysis.

We next determine heterogeneous effects by observing the kernel densities of the ordinary least-squares (OLS) estimates of the candidate effects, as shown in Figure S.1. Intuitively, the distributions of heterogeneous effects are likely to form a multimodal or wide-spread shapes. As a result, we choose the intercept, PCPN and ZNDX as the heterogeneous effects (i.e.,  $q = 3$  and hence  $p = 5$ ).

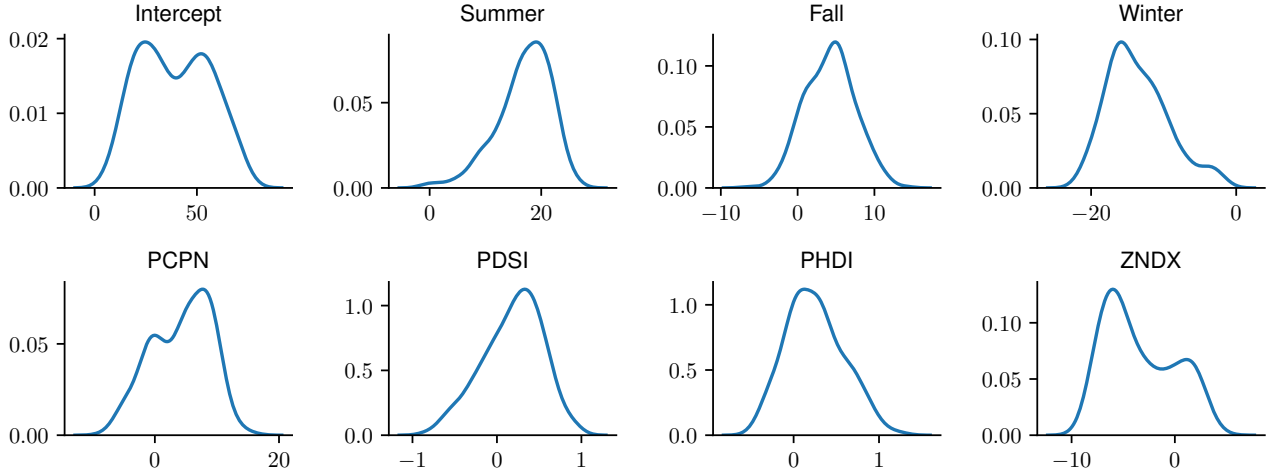


Figure S.1. Kernel densities of the 344 OLS estimates obtained from the nClimDiv database.

### S.11.3. Asymptotic covariance approximation

We assess how close the estimated asymptotic covariance is to its oracle and empirical counterparts. To this end, we observe  $\|\hat{\Sigma} - \hat{\Sigma}_{\text{OR}}\|_{\max}$ , where  $\hat{\Sigma}$  and  $\hat{\Sigma}_{\text{OR}}$  stand for the asymptotic covariance matrices of  $(\hat{\beta}^{\top}, \hat{\alpha}^{\top})^{\top}$  and  $(\hat{\beta}_{\text{OR}}^{\top}, \hat{\alpha}_{\text{OR}}^{\top})^{\top}$ , respectively. Furthermore, we compare the average asymptotic standard deviation (ASD) with the empirical standard deviation (ESD) to evaluate the finite-sample second moment approximation. Since the ASDs and ESDs of all coordinates in  $\hat{\alpha}$  behave similarly, we only show the results for  $\hat{\alpha}_{1,1}$ .

Results of  $\|\hat{\Sigma} - \hat{\Sigma}_{\text{OR}}\|_{\max}$ ,  $\text{ASD}(\hat{\alpha}_{1,1})$  and  $\text{ESD}(\hat{\alpha}_{1,1})$  are shown in Figure S.2 and Table S.2, based on the replicates with perfect subgroup structure recovery. It can be seen that the  $\hat{\Sigma}$ 's are fairly close to  $\hat{\Sigma}_{\text{OR}}$ 's, and the average  $\text{ASD}(\hat{\alpha}_{1,1})$ 's are decently comparable to their empirical counterpart  $\text{ESD}(\hat{\alpha}_{1,1})$ 's. Accordingly, we confirm that the empirical covariance matrix of the proposed estimator can be properly approximated by the asymptotic covariance matrix given in Corollary 4.1.

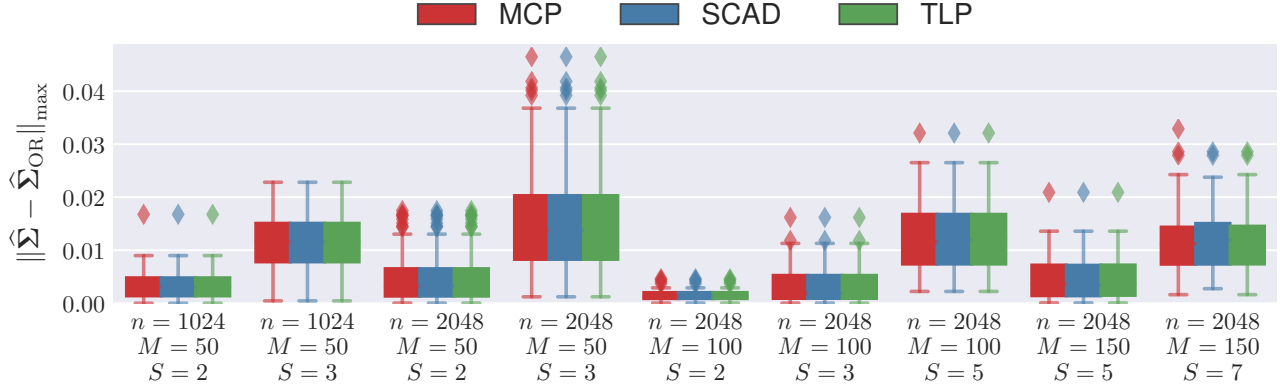


Figure S.2. Boxplots of  $\|\hat{\Sigma} - \hat{\Sigma}_{\text{OR}}\|_{\max}$ , based on replicates with perfect subgroup structure recovery.

#### S.11.4. Inferential Accuracy

In this section, we study the statistical inferential accuracy based on the asymptotic normality result in Corollary 4.1. Let  $\hat{\Sigma}_{\alpha} = (\mathbf{Z}\hat{\mathbf{A}})^{\top} \mathbf{W}(\mathbf{Z}\hat{\mathbf{A}}) - [(\mathbf{Z}\hat{\mathbf{A}})^{\top} \mathbf{W}\mathbf{X}](\mathbf{X}^{\top} \mathbf{W}\mathbf{X})^{-1}[\mathbf{X}^{\top} \mathbf{W}(\mathbf{Z}\hat{\mathbf{A}})]$  denote the estimated asymptotic covariance matrix for  $\hat{\alpha}$ , where  $\hat{\mathbf{A}}$  is the estimated label matrix  $\mathbf{A}$  by substituting the estimated subgroups for the true ones. For the first subgroup effect  $\alpha_1$ , a 95% confidence region can be constructed by

$$\text{CR}_{\alpha_1} = \left\{ \mathbf{v} \in \mathbb{R}^q : (\hat{\alpha}_1 - \mathbf{v})^{\top} \left( \mathbf{L}_1 \hat{\Sigma}_{\alpha} \mathbf{L}_1^{\top} \right)^{-1} (\hat{\alpha}_1 - \mathbf{v}) \leq \chi_q^2(0.95) \right\},$$

where  $\mathbf{L}_1 = [\mathbf{I}_q, \mathbf{O}_q, \dots, \mathbf{O}_q]_{q \times Sq}$  such that  $\mathbf{L}_1 \hat{\alpha} = \hat{\alpha}_1$  and  $\chi_q^2(0.95)$  is the 95% percentile of the  $\chi_q^2$  distribution. Table S.3 displays the empirical coverage probabilities obtained from replications with perfect subgroup recovery. It can be seen that the empirical coverages, including that by the oracle estimates, are close to the nominal level 95% except for Case 7, which again indicates that  $S = 5$  could be too large for  $M = 100$ .

Table S.3. Empirical coverage probabilities for  $\text{CR}_{\alpha_1}$ .

	MCP	SCAD	TLP	Oracle
<b>Case 1</b>	0.9600	0.9600	0.9596	0.9500
<b>Case 2</b>	0.9200	0.9200	0.9200	0.9200
<b>Case 3</b>	0.9394	0.9394	0.9394	0.9100
<b>Case 4</b>	0.9300	0.9300	0.9300	0.9400
<b>Case 5</b>	0.9400	0.9556	0.9596	0.9100
<b>Case 6</b>	0.9200	0.9200	0.9200	0.9300
<b>Case 7</b>	0.8800	0.8800	0.8800	0.8800
<b>Case 8</b>	0.9596	0.9468	0.9600	0.9500
<b>Case 9</b>	0.9286	0.9286	0.9310	0.9200

We further formulate a heterogeneous test  $H_0 : \alpha_1 = \alpha_2$  between the first two subgroups. Since  $\alpha_1 - \alpha_2 = \mathbf{L}_{12}\alpha$ , where  $\mathbf{L}_{12} = [\mathbf{I}_q, -\mathbf{I}_q, \mathbf{O}_q, \dots, \mathbf{O}_q]_{q \times Sq}$ , we define the  $F$  test statistic as

$$T = (\hat{\alpha}_1 - \hat{\alpha}_2)^{\top} \left( \mathbf{L}_{12} \hat{\Sigma}_{\alpha} \mathbf{L}_{12}^{\top} \right)^{-1} (\hat{\alpha}_1 - \hat{\alpha}_2) / q.$$

Under  $H_0$ , the test statistic  $T$  asymptotically follows the  $F_{q, N-p-\hat{S}q-2}$  distribution. For replications with perfect subgroup recovery, the  $p$ -values of the heterogeneity test are all less than 0.0001. This indicates that the estimated subgroups obtained from the proposed method cannot be further combined.