# Supplemental Materials for "A concave pairwise fusion approach to subgroup analysis"

Shujie Ma\*

Department of Statistics, University of California at Riverside and Jian Huang<sup>†</sup> Department of Statistics and Actuarial Science, University of Iowa

In this supplement, we give the technical proofs for Proposition 1 and Theorems 1-3. We also provide a detailed estimation procedure for model (2) based on the ADMM algorithm

in a way similar to that for model (1).

## A.1 Proof of Proposition 1

In this section we show the results in Proposition 1. By the definition of  $\eta^{(m+1)}$ , we have

$$L(\mu^{(m+1)}, \beta^{(m+1)}, \eta^{(m+1)}, v^{(m)}) \leq L(\mu^{(m+1)}, \beta^{(m+1)}, \eta, v^{(m)})$$

for any  $\eta$ . Define

$$f^{(m+1)} = \inf_{\substack{\Delta \mu^{(m+1)} - \eta = \mathbf{0} \\ \Delta \mu^{(m+1)} - \eta = \mathbf{0}}} \left\{ \frac{1}{2} \left\| \mathbf{y} - \boldsymbol{\mu}^{(m+1)} - \mathbf{X} \boldsymbol{\beta}^{(m+1)} \right\|^2 + \sum_{i < j} p_{\gamma}(|\eta_{ij}|, \lambda) \right\}$$
  
= 
$$\inf_{\substack{\Delta \mu^{(m+1)} - \eta = \mathbf{0}}} L(\boldsymbol{\mu}^{(m+1)}, \boldsymbol{\beta}^{(m+1)}, \boldsymbol{\eta}, \boldsymbol{v}^{(m)}).$$

Then

$$L(\boldsymbol{\mu}^{(m+1)}, \boldsymbol{\beta}^{(m+1)}, \boldsymbol{\eta}^{(m+1)}, \boldsymbol{v}^{(m)}) \leq f^{(m+1)}.$$

<sup>\*</sup>The research of Ma is supported in part by the U.S. NSF grant DMS-13-06972.

<sup>&</sup>lt;sup>†</sup>Corresponding author. The research of Huang is supported in part by the U.S. NSF grant DMS-12-08225.

Let t be an integer. Since  $\boldsymbol{v}^{(m+t-1)} = \boldsymbol{v}^{(m)} + \vartheta \sum_{i=1}^{t-1} (\boldsymbol{\Delta \mu}^{(m+i)} - \boldsymbol{\eta}^{(m+i)})$ , we have

$$\begin{split} & L(\boldsymbol{\mu}^{(m+t)}, \boldsymbol{\beta}^{(m+t)}, \boldsymbol{\eta}^{(m+t)}, \boldsymbol{v}^{(m+t-1)}) \\ &= \frac{1}{2} \left\| \mathbf{y} - \boldsymbol{\mu}^{(m+t)} - \mathbf{X} \boldsymbol{\beta}^{(m+t)} \right\|^2 + \boldsymbol{v}^{(m+t-1)\mathrm{T}} (\boldsymbol{\Delta} \boldsymbol{\mu}^{(m+t)} - \boldsymbol{\eta}^{(m+t)}) \\ &\quad + \frac{\vartheta}{2} || \boldsymbol{\Delta} \boldsymbol{\mu}^{(m+t)} - \boldsymbol{\eta}^{(m+t)} ||^2 + \sum_{i < j} p_{\gamma}(|\boldsymbol{\eta}_{ij}^{(m+t)}|, \boldsymbol{\lambda}) \\ &= \frac{1}{2} \left\| \mathbf{y} - \boldsymbol{\mu}^{(m+t)} - \mathbf{X} \boldsymbol{\beta}^{(m+t)} \right\|^2 + \boldsymbol{v}^{(m)\mathrm{T}} (\boldsymbol{\Delta} \boldsymbol{\mu}^{(m+t)} - \boldsymbol{\eta}^{(m+t)}) \\ &\quad + \vartheta \sum_{i=1}^{t-1} (\boldsymbol{\Delta} \boldsymbol{\mu}^{(m+i)} - \boldsymbol{\eta}^{(m+i)})^{\mathrm{T}} (\boldsymbol{\Delta} \boldsymbol{\mu}^{(m+t)} - \boldsymbol{\eta}^{(m+t)}) \\ &\quad + \frac{\vartheta}{2} || \boldsymbol{\Delta} \boldsymbol{\mu}^{(m+t)} - \boldsymbol{\eta}^{(m+t)} ||^2 + \sum_{i < j} p_{\gamma}(|\boldsymbol{\eta}_{ij}^{(m+t)}|, \boldsymbol{\lambda}) \\ &\leq f^{(m+t)}. \end{split}$$

Since the objective function  $L(\boldsymbol{\mu}, \boldsymbol{\beta}, \boldsymbol{\eta}, \boldsymbol{v})$  is differentiable with respect to  $(\boldsymbol{\mu}, \boldsymbol{\beta})$  and is convex with respect to  $\boldsymbol{\eta}$ , by applying the results in Theorem 4.1 of Tseng (1991), the sequence  $(\boldsymbol{\mu}^{(m)}, \boldsymbol{\beta}^{(m)}, \boldsymbol{\eta}^{(m)})$  has a limit point, denoted by  $(\boldsymbol{\mu}^*, \boldsymbol{\beta}^*, \boldsymbol{\eta}^*)$ . Then we have

$$f^* = \lim_{m \to \infty} f^{(m+1)} = \lim_{m \to \infty} f^{(m+t)} = \inf_{\Delta \mu^* - \eta = 0} \{ \frac{1}{2} \| \mathbf{y} - \mu^* - \mathbf{X} \boldsymbol{\beta}^* \|^2 + \sum_{i < j} p_{\gamma}(|\eta_{ij}|, \lambda) \},\$$

and for all  $t\geq 0$ 

$$\lim_{m \to \infty} L(\boldsymbol{\mu}^{(m+t)}, \boldsymbol{\beta}^{(m+t)}, \boldsymbol{\eta}^{(m+t)}, \boldsymbol{v}^{(m+t-1)})$$

$$= \frac{1}{2} \|\mathbf{y} - \boldsymbol{\mu}^* - \mathbf{X}\boldsymbol{\beta}^*\|^2 + \sum_{i < j} p_{\gamma}(|\eta_{ij}^*|, \lambda) + \lim_{m \to \infty} \boldsymbol{v}^{(m)\mathrm{T}}(\boldsymbol{\Delta}\boldsymbol{\mu}^* - \boldsymbol{\eta}^*) + (t - \frac{1}{2})\vartheta \|\boldsymbol{\Delta}\boldsymbol{\mu}^* - \boldsymbol{\eta}^*\|^2$$

$$\leq f^*.$$

Hence  $\lim_{m\to\infty} ||\mathbf{r}^{(m)}||^2 = r^* = ||\Delta \boldsymbol{\mu}^* - \boldsymbol{\eta}^*||^2 = 0.$ Since  $(\boldsymbol{\mu}^{(m+1)}, \boldsymbol{\beta}^{(m+1)})$  minimize  $L(\boldsymbol{\mu}, \boldsymbol{\beta}, \boldsymbol{\eta}^{(m)}, \boldsymbol{v}^{(m)})$  by definition, we have that

$$\partial L(\boldsymbol{\mu}^{(m+1)},\boldsymbol{\beta}^{(m+1)},\boldsymbol{\eta}^{(m)},\boldsymbol{v}^{(m)})/\partial \boldsymbol{\mu} = \mathbf{0},$$

and moreover,

$$\begin{array}{l} \partial L(\boldsymbol{\mu}^{(m+1)},\boldsymbol{\beta}^{(m+1)},\boldsymbol{\eta}^{(m)},\boldsymbol{v}^{(m)})/\partial\boldsymbol{\mu} \\ = & \boldsymbol{\mu}^{(m+1)} + \mathbf{X}\boldsymbol{\beta}^{(m+1)} - \mathbf{y} + \boldsymbol{\Delta}^{\mathrm{T}}\boldsymbol{v}^{(m)} + \boldsymbol{\Delta}^{\mathrm{T}}\vartheta(\boldsymbol{\Delta}\boldsymbol{\mu}^{(m+1)} - \boldsymbol{\eta}^{(m)}) \\ = & \boldsymbol{\mu}^{(m+1)} + \mathbf{X}\boldsymbol{\beta}^{(m+1)} - \mathbf{y} + \boldsymbol{\Delta}^{\mathrm{T}}(\boldsymbol{v}^{(m)} + \vartheta(\boldsymbol{\Delta}\boldsymbol{\mu}^{(m+1)} - \boldsymbol{\eta}^{(m)})) \\ = & \boldsymbol{\mu}^{(m+1)} + \mathbf{X}\boldsymbol{\beta}^{(m+1)} - \mathbf{y} + \boldsymbol{\Delta}^{\mathrm{T}}\boldsymbol{v}^{(m+1)} + \vartheta\boldsymbol{\Delta}^{\mathrm{T}}(\boldsymbol{\eta}^{(m+1)} - \boldsymbol{\eta}^{(m)}). \end{array}$$

The last step follows from  $\boldsymbol{v}^{(m+1)} = \boldsymbol{v}^{(m)} + \vartheta(\Delta \boldsymbol{\mu}^{(m+1)} - \boldsymbol{\eta}^{(m+1)})$ . Therefore,

$$\mathbf{s}^{(m+1)} = \vartheta \mathbf{\Delta}^{\mathrm{T}}(\boldsymbol{\eta}^{(m+1)} - \boldsymbol{\eta}^{(m)}) = -(\boldsymbol{\mu}^{(m+1)} + \mathbf{X}\boldsymbol{\beta}^{(m+1)} - \mathbf{y} + \mathbf{\Delta}^{\mathrm{T}}\boldsymbol{v}^{(m+1)}).$$

Since  $||\Delta \mu^* - \eta^*||^2 = 0$ ,

$$\lim_{m \to \infty} \partial L(\boldsymbol{\mu}^{(m+1)}, \boldsymbol{\beta}^{(m+1)}, \boldsymbol{\eta}^{(m)}, \boldsymbol{v}^{(m)}) / \partial \boldsymbol{\mu}$$
  
= 
$$\lim_{m \to \infty} \boldsymbol{\mu}^{(m+1)} + \mathbf{X} \boldsymbol{\beta}^{(m+1)} - \mathbf{y} + \boldsymbol{\Delta}^{\mathrm{T}} \boldsymbol{v}^{(m+1)}$$
  
= 
$$\boldsymbol{\mu}^{*} + \mathbf{X} \boldsymbol{\beta}^{*} - \mathbf{y} + \boldsymbol{\Delta}^{\mathrm{T}} \boldsymbol{v}^{*} = \mathbf{0}.$$

Therefore,  $\lim_{m\to\infty} \mathbf{s}^{(m+1)} = \mathbf{0}$ .

# A.2 Proof of Theorem 1

In this section we show the results in Theorem 1. Since for every  $\boldsymbol{\mu} \in \mathcal{M}_{\mathcal{G}}$ , it can be written as  $\boldsymbol{\mu} = \mathbf{Z}\boldsymbol{\alpha}$ , and hence  $\boldsymbol{\alpha} = \mathbf{D}^{-1}\mathbf{Z}^{\mathrm{T}}\boldsymbol{\mu}$ . Then  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}} = ((\mathbf{Z}\widehat{\boldsymbol{\alpha}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ , where

$$\begin{pmatrix} \widehat{\boldsymbol{\alpha}}^{or} \\ \widehat{\boldsymbol{\beta}}^{or} \end{pmatrix} = \arg\min_{\boldsymbol{\alpha} \in R^{K}, \boldsymbol{\beta} \in R^{p}} \frac{1}{2} ||\mathbf{y} - \mathbf{Z}\boldsymbol{\alpha} - \mathbf{X}\boldsymbol{\beta}||^{2} = [(\mathbf{Z}, \mathbf{X})^{\mathrm{T}} (\mathbf{Z}, \mathbf{X})]^{-1} (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} \mathbf{y}.$$

Then

$$\left( egin{array}{c} \widehat{oldsymbol{lpha}}^{or} - oldsymbol{lpha}^0 \\ \widehat{oldsymbol{eta}}^{or} - oldsymbol{eta}^0 \end{array} 
ight) = [(\mathbf{Z}, \mathbf{X})^{\mathrm{T}} (\mathbf{Z}, \mathbf{X})]^{-1} (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} oldsymbol{\epsilon},$$

where  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_n)^{\mathrm{T}}$  and  $\boldsymbol{\alpha}^0 = (\alpha_1^0, \dots, \alpha_K^0)^{\mathrm{T}}$ . Hence

$$\left\| \left( \begin{array}{c} \widehat{\boldsymbol{\alpha}}^{or} - \boldsymbol{\alpha}^{0} \\ \widehat{\boldsymbol{\beta}}^{or} - \boldsymbol{\beta}^{0} \end{array} \right) \right\|_{\infty} \leq \left\| \left[ (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} (\mathbf{Z}, \mathbf{X}) \right]^{-1} \right\|_{\infty} \left\| (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} \boldsymbol{\epsilon} \right\|_{\infty}.$$
(A.1)

By Condition (C1), we have  $\left\| \left[ \left( \mathbf{Z}, \mathbf{X} \right)^{\mathrm{T}} \left( \mathbf{Z}, \mathbf{X} \right) \right]^{-1} \right\| \leq C_{1}^{-1} \left| \mathcal{G}_{\min} \right|^{-1}$  and thus

$$\left\| \left[ \left( \mathbf{Z}, \mathbf{X} \right)^{\mathrm{T}} \left( \mathbf{Z}, \mathbf{X} \right) \right]^{-1} \right\|_{\infty} \leq \sqrt{K + p} C_{1}^{-1} \left| \mathcal{G}_{\min} \right|^{-1}.$$
(A.2)

Moreover

$$P(\left\| (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} \boldsymbol{\epsilon} \right\|_{\infty} > C\sqrt{n \log n}) \le P(\left\| \mathbf{Z}^{\mathrm{T}} \boldsymbol{\epsilon} \right\|_{\infty} > C\sqrt{n \log n}) + P(\left\| \mathbf{X}^{\mathrm{T}} \boldsymbol{\epsilon} \right\|_{\infty} > C\sqrt{n \log n}),$$

for some constant  $0 < C < \infty$ . By Condition (C3) and union bound,

$$P\left(\left\|\mathbf{Z}^{\mathrm{T}}\boldsymbol{\epsilon}\right\|_{\infty} > C\sqrt{n\log n}\right)$$

$$\leq \sum_{k=1}^{K} P(\left|\sum_{i\in\mathcal{G}_{k}}\epsilon_{i}\right| > C\sqrt{n\log n}) \leq \sum_{k=1}^{K} P(\left|\sum_{i\in\mathcal{G}_{k}}\epsilon_{i}\right| > \sqrt{|\mathcal{G}_{k}|}C\sqrt{\log n})$$

$$\leq 2K\exp(-c_{1}C^{2}\log n) = 2Kn^{-c_{1}C^{2}},$$

and by Conditions (C1) and (C3) and union bound,

$$P\left(\left\|\mathbf{X}^{\mathrm{T}}\boldsymbol{\epsilon}\right\|_{\infty} > C\sqrt{n\log n}\right)$$

$$\leq \sum_{j=1}^{p} P\left(\left|\mathbf{X}_{j}^{\mathrm{T}}\boldsymbol{\epsilon}\right| > \sqrt{n}C\sqrt{\log n}\right)$$

$$\leq 2p\exp(-c_{1}C^{2}\log n) = 2pn^{-c_{1}C^{2}}.$$

By the above results, we have

$$P(\left\| (\mathbf{Z}, \mathbf{X})^{\mathrm{T}} \boldsymbol{\epsilon} \right\|_{\infty} > C \sqrt{n \log n}) \le 2(K+p) n^{-c_1 C^2}.$$
(A.3)

Therefore, by (A.1), (A.2) and (A.3), we have with probability at least  $1 - 2(K + p)n^{-c_1C^2}$ ,

$$\left\| \left( \begin{array}{c} \widehat{\boldsymbol{\alpha}}^{or} - \boldsymbol{\alpha}^{0} \\ \widehat{\boldsymbol{\beta}}^{or} - \boldsymbol{\beta}^{0} \end{array} \right) \right\|_{\infty} \leq C C_{1}^{-1} \sqrt{K + p} \left| \mathcal{G}_{\min} \right|^{-1} \sqrt{n \log n},$$

and hence  $\|\widehat{\boldsymbol{\mu}}^{or} - \boldsymbol{\mu}^0\|_{\infty} = \|\widehat{\boldsymbol{\alpha}}^{or} - \boldsymbol{\alpha}^0\|_{\infty} \leq CC_1^{-1}\sqrt{K+p} |\mathcal{G}_{\min}|^{-1}\sqrt{n\log n}$ . The result (8) in Theorem 1 is proved by letting  $C = c_1^{-1/2}$ , and result (10) follows from Central Limit Theorem.

#### **A.3** Proof of Theorem 2

In this section we show the results in Theorem 2. Define

$$L_n(\boldsymbol{\mu}, \boldsymbol{\beta}) = \frac{1}{2} ||\mathbf{y} - \boldsymbol{\mu} - \mathbf{X}\boldsymbol{\beta}||^2, P_n(\boldsymbol{\mu}) = \lambda \sum_{i < j} \rho(|\mu_i - \mu_j|),$$
  

$$L_n^{\mathcal{G}}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{1}{2} ||\mathbf{y} - \mathbf{Z}\boldsymbol{\alpha} - \mathbf{X}\boldsymbol{\beta}||^2, P_n^{\mathcal{G}}(\boldsymbol{\alpha}) = \lambda \sum_{k < k'} |\mathcal{G}_k| |\mathcal{G}_{k'}| \rho(|\alpha_k - \alpha_{k'}|),$$

and let

$$Q_n(\boldsymbol{\mu},\boldsymbol{\beta}) = L_n(\boldsymbol{\mu},\boldsymbol{\beta}) + P_n(\boldsymbol{\mu}), Q_n^{\mathcal{G}}(\boldsymbol{\alpha},\boldsymbol{\beta}) = L_n^{\mathcal{G}}(\boldsymbol{\alpha},\boldsymbol{\beta}) + P_n^{\mathcal{G}}(\boldsymbol{\alpha}).$$

Let  $T: \mathcal{M}_{\mathcal{G}} \to \mathbb{R}^{K}$  be the mapping such that  $T(\boldsymbol{\mu})$  is the  $K \times 1$  vector whose  $k^{\text{th}}$  coordinate equals to the common value of  $\mu_{i}$  for  $i \in \mathcal{G}_{k}$ . Let  $T^{*}: \mathbb{R}^{n} \to \mathbb{R}^{K}$  be the mapping such that  $T^{*}(\boldsymbol{\mu}) = \{ |\mathcal{G}_{k}|^{-1} \sum_{i \in \mathcal{G}_{k}} \mu_{i} \}_{k=1}^{K}. \text{ Clearly, when } \boldsymbol{\mu} \in \mathcal{M}_{\mathcal{G}}, T(\boldsymbol{\mu}) = T^{*}(\boldsymbol{\mu}).$ By calculation, for every  $\boldsymbol{\mu} \in \mathcal{M}_{\mathcal{G}}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\alpha} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\mu} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\mu} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\mu} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) = P_{n}^{\mathcal{G}}(T(\boldsymbol{\mu})) \text{ and for every } \boldsymbol{\mu} \in \mathbb{R}^{K}, \text{ we have } P_{n}(\boldsymbol{\mu}) \in \mathbb{R}^{K}, \text{ we have } P$ 

have  $P_n(T^{-1}(\boldsymbol{\alpha})) = P_n^{\mathcal{G}}(\boldsymbol{\alpha})$ . Hence

$$Q_n(\boldsymbol{\mu},\boldsymbol{\beta}) = Q_n^{\mathcal{G}}(T(\boldsymbol{\mu}),\boldsymbol{\beta}), Q_n^{\mathcal{G}}(\boldsymbol{\alpha},\boldsymbol{\beta}) = Q_n(T^{-1}(\boldsymbol{\alpha}),\boldsymbol{\beta}).$$
(A.4)

Consider the neighborhood of  $(\boldsymbol{\mu}^0, \boldsymbol{\beta}^0)$ :

$$\Theta = \{\boldsymbol{\mu} \in \mathbb{R}^{n}, \boldsymbol{\beta} \in \mathbb{R}^{p} : \left\| ((\boldsymbol{\mu} - \boldsymbol{\mu}^{0})^{\mathrm{T}}, (\boldsymbol{\beta} - \boldsymbol{\beta}^{0})^{\mathrm{T}})^{\mathrm{T}} \right\|_{\infty} \le \phi_{n} \}$$

By the result in Theorem 1, there is an event  $E_1$  such that on the event  $E_1$ ,

$$\left\| ((\widehat{\boldsymbol{\mu}}^{or} - \boldsymbol{\mu}^0)^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or} - \boldsymbol{\beta}^0)^{\mathrm{T}})^{\mathrm{T}} \right\|_{\infty} \leq \phi_n,$$

and  $P(E_1^C) \leq 2(K+p)n^{-1}$ . Hence  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}} \in \Theta$  on the event  $E_1$ . For any  $\boldsymbol{\mu} \in \mathbb{R}^n$ , let  $\boldsymbol{\mu}^* = T^{-1}(T^*(\boldsymbol{\mu}))$ . We show that  $(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  is a strictly local minimizer of the objective function (3) with probability approaching 1 through the following two steps.

(i). On the event  $E_1$ ,  $Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta}) > Q_n(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  for any  $((\boldsymbol{\mu})^{\mathrm{T}}, (\boldsymbol{\beta})^{\mathrm{T}})^{\mathrm{T}} \in \Theta$  and  $((\boldsymbol{\mu}^*)^{\mathrm{T}}, (\boldsymbol{\beta})^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}.$ 

(ii). There is an event  $E_2$  such that  $P(E_2^C) \leq 2n^{-1}$ . On  $E_1 \cap E_2$ , there is a neighborhood of  $((\widehat{\boldsymbol{\mu}}^{or})^T, (\widehat{\boldsymbol{\beta}}^{or})^T)^T$ , denoted by  $\Theta_n$ , such that  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) \geq Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta})$  for any  $((\boldsymbol{\mu})^T, (\boldsymbol{\beta})^T)^T \in \Theta_n \cap \Theta$  for sufficiently large n.

Therefore, by the results in (i) and (ii), we have  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) > Q_n(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  for any  $((\boldsymbol{\mu})^{\mathrm{T}}, (\boldsymbol{\beta})^{\mathrm{T}})^{\mathrm{T}} \in \Theta_n \cap \Theta$  and  $((\boldsymbol{\mu})^{\mathrm{T}}, (\boldsymbol{\beta})^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ , so that  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$  is a strict local minimizer of  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta})$  given in (3) on the event  $E_1 \cap E_2$  with  $P(E_1 \cap E_2) \geq 1 - 2(K + p + 1)n^{-1}$  for sufficiently large n.

In the following we prove the result in (i). We first show  $P_n^{\mathcal{G}}(T^*(\boldsymbol{\mu})) = C_n$  for any  $\boldsymbol{\mu} \in \Theta$ , where  $C_n$  is a constant which does not depend on  $\boldsymbol{\mu}$ . Let  $T^*(\boldsymbol{\mu}) = \boldsymbol{\alpha} = (\alpha_1, \ldots, \alpha_K)^{\mathrm{T}}$ . It suffices to show that  $|\alpha_k - \alpha_{k'}| > a\lambda$  for all k and k'. Then by Condition (C2),  $\rho(|\alpha_k - \alpha_{k'}|)$ is a constant, and as a result  $P_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}))$  is a constant. Since

$$|\alpha_k - \alpha_{k'}| \ge |\alpha_k^0 - \alpha_{k'}^0| - 2||\boldsymbol{\alpha} - \boldsymbol{\alpha}^0||_{\infty},$$

and

$$||\boldsymbol{\alpha} - \boldsymbol{\alpha}^{0}||_{\infty} = \sup_{k} |\sum_{i \in \mathcal{G}_{k}} \mu_{i}/|\mathcal{G}_{k}| - \alpha_{k}^{0}| = \sup_{k} |\sum_{i \in \mathcal{G}_{k}} (\mu_{i} - \mu_{i}^{0})/|\mathcal{G}_{k}||$$

$$\leq \sup_{k} \sup_{i \in \mathcal{G}_{k}} |\mu_{i} - \mu_{i}^{0}| = ||\boldsymbol{\mu} - \boldsymbol{\mu}^{0}||_{\infty}, \qquad (A.5)$$

then for all k and k'

$$|\alpha_k - \alpha_{k'}| \ge |\alpha_k^0 - \alpha_{k'}^0| - 2||\boldsymbol{\mu} - \boldsymbol{\mu}^0||_{\infty} \ge b_n - 2\phi_n > a\lambda,$$

where the last inequality follows from the assumption that  $b_n > a\lambda \gg \phi_n$ . Therefore, we have  $P_n^{\mathcal{G}}(T^*(\boldsymbol{\mu})) = C_n$ , and hence  $Q_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}),\boldsymbol{\beta}) = L_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}),\boldsymbol{\beta}) + C_n$  for all  $(\boldsymbol{\mu}^{\mathrm{T}},\boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta$ . Since  $((\widehat{\boldsymbol{\alpha}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$  is the unique global minimizer of  $L_n^{\mathcal{G}}(\boldsymbol{\alpha},\boldsymbol{\beta})$ , then  $L_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}),\boldsymbol{\beta}) > L_n^{\mathcal{G}}(\widehat{\boldsymbol{\alpha}}^{or},\widehat{\boldsymbol{\beta}}^{or})$  for all  $(T^*(\boldsymbol{\mu})^{\mathrm{T}},\boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\alpha}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$  and thus  $Q_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}),\boldsymbol{\beta}) > Q_n^{\mathcal{G}}(\widehat{\boldsymbol{\alpha}}^{or},\widehat{\boldsymbol{\beta}}^{or})$  for all  $T^*(\boldsymbol{\mu}) \neq \widehat{\boldsymbol{\alpha}}^{or}$ . By (A.4), we have  $Q_n^{\mathcal{G}}(\widehat{\boldsymbol{\alpha}}^{or},\widehat{\boldsymbol{\beta}}^{or}) = Q_n(\widehat{\boldsymbol{\mu}}^{or},\widehat{\boldsymbol{\beta}}^{or})$  and  $Q_n^{\mathcal{G}}(T^*(\boldsymbol{\mu}),\boldsymbol{\beta}) = Q_n(T^{-1}(T^*(\boldsymbol{\mu})),\boldsymbol{\beta}) = Q_n(\boldsymbol{\mu}^*,\boldsymbol{\beta})$ . Therefore,  $Q_n(\boldsymbol{\mu}^*,\boldsymbol{\beta}) > Q_n(\widehat{\boldsymbol{\mu}}^{or},\widehat{\boldsymbol{\beta}}^{or})$  for all  $\boldsymbol{\mu}^* \neq \widehat{\boldsymbol{\mu}}^{or}$ , and the result in (i) is proved.

Next we prove the result in (ii). For a positive sequence  $t_n$ , let  $\Theta_n = \{ \boldsymbol{\mu} : || \boldsymbol{\mu} - \hat{\boldsymbol{\mu}}^{or} || \leq t_n \}$ . For  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta_n \cap \Theta$ , by Taylor's expansion, we have

$$Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) - Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta}) = \Gamma_1 + \Gamma_2,$$

where

$$\Gamma_{1} = -(\mathbf{y} - (\mathbf{I}_{n}, \mathbf{X})((\boldsymbol{\mu}^{m})^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}})^{\mathrm{T}}(\boldsymbol{\mu} - \boldsymbol{\mu}^{*}),$$
  
$$\Gamma_{2} = \sum_{i=1}^{n} \frac{\partial P_{n}(\boldsymbol{\mu}^{m})}{\partial \mu_{i}}(\mu_{i} - \mu_{i}^{*}),$$

in which  $\boldsymbol{\mu}^m = \varsigma \boldsymbol{\mu} + (1 - \varsigma) \boldsymbol{\mu}^*$  for some  $\varsigma \in (0, 1)$ . Moreover,

$$\Gamma_{2} = \lambda \sum_{\{j>i\}} \overline{\rho}(\mu_{i}^{m} - \mu_{j}^{m})(\mu_{i} - \mu_{i}^{*}) + \lambda \sum_{\{ji\}} \overline{\rho}(\mu_{i}^{m} - \mu_{j}^{m})(\mu_{i} - \mu_{i}^{*}) + \lambda \sum_{\{ii\}} \overline{\rho}(\mu_{i}^{m} - \mu_{j}^{m})(\mu_{i} - \mu_{i}^{*}) - \lambda \sum_{\{ii\}} \overline{\rho}(\mu_{i}^{m} - \mu_{j}^{m})\{(\mu_{i} - \mu_{i}^{*}) - (\mu_{j} - \mu_{j}^{*})\}.$$
(A.6)

When  $i, j \in \mathcal{G}_k$ ,  $\mu_i^* = \mu_j^*$ , and  $\mu_i^m - \mu_j^m$  has the same sign as  $\mu_i - \mu_j$ . Hence

$$\Gamma_{2} = \lambda \sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}, i < j\}} \rho'(|\mu_{i}^{m} - \mu_{j}^{m}|)|\mu_{i} - \mu_{j}| + \lambda \sum_{k < k'} \sum_{\{i \in \mathcal{G}_{k}, j' \in \mathcal{G}_{k'}\}} \overline{\rho}(\mu_{i}^{m} - \mu_{j}^{m})\{(\mu_{i} - \mu_{i}^{*}) - (\mu_{j} - \mu_{j}^{*})\}$$

As shown in (A.5),

$$||oldsymbol{\mu}^*-oldsymbol{\mu}^0||_\infty=||oldsymbol{lpha}-oldsymbol{lpha}^0||_\infty\leq ||oldsymbol{\mu}-oldsymbol{\mu}^0||_\infty$$

Since  $\boldsymbol{\mu}^m = \varsigma \boldsymbol{\mu} + (1 - \varsigma) \boldsymbol{\mu}^*$ ,

$$||\boldsymbol{\mu}^{m} - \boldsymbol{\mu}^{0}||_{\infty} \leq ||\boldsymbol{\mu} - \boldsymbol{\mu}^{0}||_{\infty} \leq \phi_{n}, \qquad (A.7)$$

and then for  $k \neq k', i \in \mathcal{G}_k, j \in \mathcal{G}_{k'}$ ,

$$\begin{aligned} |\mu_i^m - \mu_j^m| &\geq \min_{i \in \mathcal{G}_k, j \in \mathcal{G}_{k'}} |\mu_i^0 - \mu_j^0| - 2||\boldsymbol{\mu}^m - \boldsymbol{\mu}^0||_{\infty} \\ &\geq b_n - 2||\boldsymbol{\mu} - \boldsymbol{\mu}^0||_{\infty} \geq b_n - 2\phi_n > a\lambda, \end{aligned}$$

and thus  $\overline{\rho}(\mu_i^m - \mu_j^m) = 0$ . Therefore,

$$\Gamma_2 = \lambda \sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_k, i < j\}} \rho'(|\mu_i^m - \mu_j^m|)|\mu_i - \mu_j|.$$

Furthermore, by the same reasoning as (A.5), we have

$$|| \boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}^{or} ||_{\infty} \leq || \boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or} ||_{\infty}.$$

Then

$$\begin{aligned} |\mu_i^m - \mu_j^m| &\leq 2||\boldsymbol{\mu}^m - \boldsymbol{\mu}^*||_{\infty} \leq 2||\boldsymbol{\mu} - \boldsymbol{\mu}^*||_{\infty} \\ &\leq 2(||\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or}||_{\infty} + ||\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}^{or}||_{\infty}) \\ &\leq 4||\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or}||_{\infty} \leq 4t_n. \end{aligned}$$

Hence  $\rho'(|\mu_i^m - \mu_j^m|) \ge \rho'(4t_n)$  by concavity of  $\rho(\cdot)$ . As a result,

$$\Gamma_2 \ge \lambda \sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_k, i < j\}} \rho'(4t_n) |\mu_i - \mu_j|.$$
(A.8)

Let

$$\mathbf{w} = (w_1, \dots, w_n)^{\mathrm{T}} = \mathbf{y} - (\mathbf{I}_n, \mathbf{X})((\boldsymbol{\mu}^m)^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}}.$$

Then

$$\Gamma_{1} = -\mathbf{w}^{\mathrm{T}}(\boldsymbol{\mu} - \boldsymbol{\mu}^{*}) = -\sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}\}} \frac{w_{i}(\mu_{i} - \mu_{j})}{|\mathcal{G}_{k}|} 
= -\sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}\}} \frac{w_{i}(\mu_{i} - \mu_{j})}{2|\mathcal{G}_{k}|} - \sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}\}} \frac{w_{i}(\mu_{i} - \mu_{j})}{2|\mathcal{G}_{k}|} 
= -\sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}, i < j\}} \frac{(w_{j} - w_{i})(\mu_{j} - \mu_{i})}{2|\mathcal{G}_{k}|} 
= -\sum_{k=1}^{K} \sum_{\{i,j \in \mathcal{G}_{k}, i < j\}} \frac{(w_{j} - w_{i})(\mu_{j} - \mu_{i})}{|\mathcal{G}_{k}|}.$$
(A.9)

Since

$$\mathbf{w} = \boldsymbol{\epsilon} + \mathbf{X}(\boldsymbol{\beta}^0 - \boldsymbol{\beta}) + \boldsymbol{\mu}^0 - \boldsymbol{\mu}^m,$$

then

$$\max_{i,j} |w_j - w_i| \le 2||\mathbf{w}||_{\infty} \le 2||\boldsymbol{\epsilon}||_{\infty} + 2||\mathbf{X}||_{\infty}||\boldsymbol{\beta}^0 - \boldsymbol{\beta}||_{\infty} + 2||\boldsymbol{\mu}^0 - \boldsymbol{\mu}^m||_{\infty}.$$

Hence by (A.7) and Condition (C1),

$$\max_{i,j} |w_j - w_i| \le 2||\boldsymbol{\epsilon}||_{\infty} + 2C_2 p\phi_n + 2\phi_n$$

By Condition (C3),

$$P(||\boldsymbol{\epsilon}||_{\infty} > \sqrt{2c_1^{-1}}\sqrt{\log n}) \le \sum_{i=1}^n P(|\varepsilon_i| > \sqrt{2c_1^{-1}}\sqrt{\log n}) \le 2n^{-1}.$$

Thus there is an event  $E_2$  such that  $P(E_2^C) \leq 2n^{-1}$ , and on the event  $E_2$ ,

$$\max_{i,j} |w_j - w_i| \le 2\sqrt{2c_1^{-1}}\sqrt{\log n} + 2(C_2p + 1)\phi_n.$$
(A.10)

Hence

$$|\mathcal{G}_{\min}|^{-1} \max_{i,j} |w_j - w_i| \le |\mathcal{G}_{\min}|^{-1} \{ 2\sqrt{2c_1^{-1}} \sqrt{\log n} + 2(C_2p + 1)\phi_n \}$$

Since  $|\mathcal{G}_{\min}| \gg \sqrt{(K+p)n\log n}$  and p = o(n), then  $|\mathcal{G}_{\min}|^{-1}p = o(1)$ . Thus  $\lambda \gg \phi_n \gg |\mathcal{G}_{\min}|^{-1}2(C_2p+1)\phi_n$ . Moreover,  $\lambda \gg \phi_n \gg |\mathcal{G}_{\min}|^{-1}\sqrt{\log n}$ . Hence

$$\lambda \gg |\mathcal{G}_{\min}|^{-1} \max_{i,j} |w_j - w_i|.$$
(A.11)

Let  $t_n = o(1)$ , then  $\rho'(4t_n) \to 1$ . Therefore, by (A.8), (A.9), and (A.11),

$$Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) - Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta}) = \Gamma_1 + \Gamma_2$$
  

$$\geq \sum_{k=1}^K \sum_{\{i,j \in \mathcal{G}_k, i < j\}} [\lambda \rho'(4t_n) - |\mathcal{G}_{\min}|^{-1} \max_{i,j} |w_j - w_i|] |\mu_i - \mu_j| \ge 0,$$

for sufficiently large n, so that the result in (ii) is proved.

### A.4 Proof of Theorem 3

In this section we show the results in Theorem 3. The proofs of (12) and (13) follow the same arguments as the proof of Theorem 1 by letting  $\mathbf{Z} = \mathbf{1}_n$  and  $|\mathcal{G}_{\min}| = n$ , and thus they are omitted. Next, we will show (14). It follows similar procedures as the proof of Theorem 2 with the details given below. Let  $\mathcal{M}$  be the subspace of  $\mathbb{R}^n$ , defined as

$$\mathcal{M} = \{ \boldsymbol{\mu} \in \mathbb{R}^n : \mu_1 = \cdots = \mu_n \}.$$

For each  $\boldsymbol{\mu} \in \mathcal{M}$ , it can be written as  $\boldsymbol{\mu} = \mathbf{1}_n \alpha$ , where  $\alpha$  is the common value of  $\boldsymbol{\mu}$ . Let  $T : \mathcal{M} \to R$  be the mapping such that  $T(\boldsymbol{\mu})$  is the scalar that equals to the common value of  $\mu_i$ 's. Let  $T^* : R^n \to R$  be the mapping such that  $T^*(\boldsymbol{\mu}) = n^{-1} \sum_{i=1}^n \mu_i$ . Clearly, when  $\boldsymbol{\mu} \in \mathcal{M}, T(\boldsymbol{\mu}) = T^*(\boldsymbol{\mu})$ . Consider the neighborhood of  $(\boldsymbol{\mu}^0, \boldsymbol{\beta}^0)$ :

$$\Theta = \{\boldsymbol{\mu} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^p: \left\| ((\boldsymbol{\mu} - \boldsymbol{\mu}^0)^{\mathrm{T}}, (\boldsymbol{\beta} - \boldsymbol{\beta}^0)^{\mathrm{T}})^{\mathrm{T}} \right\|_{\infty} \le \phi_n \},\$$

where  $\phi_n = c_1^{-1/2} C_1^{-1} \sqrt{1+p} \sqrt{n^{-1} \log n}$ . By the result in (12), there is an event  $E_1$  such that on the even  $E_1$ ,

$$\left\| ((\widehat{\boldsymbol{\mu}}^{or} - \boldsymbol{\mu}^0)^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or} - \boldsymbol{\beta}^0)^{\mathrm{T}})^{\mathrm{T}} \right\|_{\infty} \leq \phi_n,$$

and  $P(E_1^C) \leq 2(1+p)n^{-1}$ . Hence  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}} \in \Theta$  on the event  $E_1$ . For any  $\boldsymbol{\mu} \in \mathbb{R}^n$ , let  $\boldsymbol{\mu}^* = T^{-1}(T^*(\boldsymbol{\mu}))$ . We show that  $(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  is a strictly local minimizer of the objective function (3) with probability approaching 1 through the following two steps.

(i). On the event  $E_1, Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta}) > Q_n(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  for any  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta$  and  $((\boldsymbol{\mu}^*)^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}.$ 

(ii). There is an event  $E_2$  such that  $P(E_2^C) \leq 2n^{-1}$ . On  $E_1 \cap E_2$ , there is a neighborhood of  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ , denoted by  $\Theta_n$ , such that  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) \geq Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta})$  for any  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta_n \cap \Theta$  for sufficiently large n.

Therefore, by the results in (i) and (ii), we have  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) > Q_n(\widehat{\boldsymbol{\mu}}^{or}, \widehat{\boldsymbol{\beta}}^{or})$  for any  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta_n \cap \Theta$  and  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ , so that  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$  is a strict local minimizer of  $Q_n(\boldsymbol{\mu}, \boldsymbol{\beta})$  on the event  $E_1 \cap E_2$  with  $P(E_1 \cap E_2) \geq 1 - 2(p+2)n^{-1}$  for sufficiently large n.

By the definition of  $((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ , we have  $\frac{1}{2} \sum_{i=1}^{n} (y_i - \boldsymbol{\mu}^* - \mathbf{x}_i^{\mathrm{T}} \boldsymbol{\beta})^2 > \frac{1}{2} \sum_{i=1}^{n} (y_i - \widehat{\boldsymbol{\mu}}^{or} - \mathbf{x}_i^{\mathrm{T}} \widehat{\boldsymbol{\beta}}^{or})^2$  for any  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta$  and  $((\boldsymbol{\mu}^*)^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \neq ((\widehat{\boldsymbol{\mu}}^{or})^{\mathrm{T}}, (\widehat{\boldsymbol{\beta}}^{or})^{\mathrm{T}})^{\mathrm{T}}$ . Moreover, since

 $p_{\gamma}(|\widehat{\mu}_{i}^{or}-\widehat{\mu}_{j}^{or}|,\lambda) = p_{\gamma}(|\mu_{i}^{*}-\mu_{j}^{*}|,\lambda) = 0 \text{ for } 1 \leq i,j \leq n, \text{ we have } Q_{n}(\boldsymbol{\mu}^{*},\boldsymbol{\beta}) = \frac{1}{2} \sum_{i=1}^{n} (y_{i}-\boldsymbol{\mu}^{*}-\mathbf{x}_{i}^{\mathrm{T}}\boldsymbol{\beta})^{2} \text{ and } Q_{n}(\widehat{\boldsymbol{\mu}}^{or},\widehat{\boldsymbol{\beta}}^{or}) = \frac{1}{2} \sum_{i=1}^{n} (y_{i}-\widehat{\boldsymbol{\mu}}^{or}-\mathbf{x}_{i}^{\mathrm{T}}\widehat{\boldsymbol{\beta}}^{or})^{2}. \text{ Therefore, } Q_{n}(\boldsymbol{\mu}^{*},\boldsymbol{\beta}) > Q_{n}(\widehat{\boldsymbol{\mu}}^{or},\widehat{\boldsymbol{\beta}}^{or}).$ 

Next we prove the result in (ii). For a positive sequence  $t_n$ , let  $\Theta_n = \{\boldsymbol{\mu} : ||\boldsymbol{\mu} - \hat{\boldsymbol{\mu}}^{or}|| \leq t_n\}$ . For  $(\boldsymbol{\mu}^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}} \in \Theta_n \cap \Theta$ , by Taylor's expansion, we have

$$Q_n(\boldsymbol{\mu},\boldsymbol{\beta}) - Q_n(\boldsymbol{\mu}^*,\boldsymbol{\beta}) = \Gamma_1 + \Gamma_2,$$

where

$$\Gamma_{1} = -(\mathbf{y} - (\mathbf{I}_{n}, \mathbf{X})((\boldsymbol{\mu}^{m})^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}})^{\mathrm{T}}(\boldsymbol{\mu} - \boldsymbol{\mu}^{*}),$$
  
$$\Gamma_{2} = \sum_{i=1}^{n} \frac{\partial P_{n}(\boldsymbol{\mu}^{m})}{\partial \mu_{i}}(\mu_{i} - \mu_{i}^{*}).$$

in which  $\mu^m = \varsigma \mu + (1 - \varsigma) \mu^*$  for some  $\varsigma \in (0, 1)$ . Moreover, by (A.6), we have

$$\Gamma_2 = \lambda \sum_{i < j} \overline{\rho}(\mu_i^m - \mu_j^m) \{ (\mu_i - \mu_i^*) - (\mu_j - \mu_j^*) \}$$
  
=  $\lambda \sum_{i < j} \rho'(|\mu_i^m - \mu_j^m|) |\mu_i - \mu_j|,$ 

where the second equality holds due to the fact that  $\mu_i^* = \mu_j^*$  and  $\mu_i^m - \mu_j^m$  has the same sign as  $\mu_i - \mu_j$ . Let  $T^*(\boldsymbol{\mu}) = \alpha$ . Following the same reasoning as the proof for (A.5), we have

$$||\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}^{or}||_{\infty} = |\alpha - \widehat{\alpha}^{or}| \leq ||\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or}||_{\infty}.$$

Then

$$\begin{aligned} |\mu_i^m - \mu_j^m| &\leq 2||\boldsymbol{\mu}^m - \boldsymbol{\mu}^*||_{\infty} \leq 2||\boldsymbol{\mu} - \boldsymbol{\mu}^*||_{\infty} \\ &\leq 2(||\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or}||_{\infty} + ||\boldsymbol{\mu}^* - \widehat{\boldsymbol{\mu}}^{or}||_{\infty}) \\ &\leq 4||\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}^{or}||_{\infty} \leq 4t_n. \end{aligned}$$

Hence  $\rho'(|\mu_i^m - \mu_j^m|) \ge \rho'(4t_n)$  by concavity of  $\rho(\cdot)$ . As a result,

$$\Gamma_2 \ge \lambda \sum_{i < j} \rho'(4t_n) |\mu_i - \mu_j|. \tag{A.12}$$

Then, by the same reasoning as the proof for (A.9), we have

$$\Gamma_1 = -\mathbf{w}^{\mathrm{T}}(\boldsymbol{\mu} - \boldsymbol{\mu}^*) = -n^{-1} \sum_{i < j} (w_j - w_i)(\mu_j - \mu_i), \qquad (A.13)$$

where  $\mathbf{w} = (w_1, \ldots, w_n)^{\mathrm{T}} = \mathbf{y} - (\mathbf{I}_n, \mathbf{X})((\boldsymbol{\mu}^m)^{\mathrm{T}}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}}$ . By the same reasoning as the proof for (A.10), we have that there is an event  $E_2$  such that  $P(E_2^C) \leq 2n^{-1}$ , and on the event  $E_2$ ,

$$\max_{i,j} |w_j - w_i| \le 2\sqrt{2c_1^{-1}}\sqrt{\log n} + 2(C_2p + 1)\phi_n.$$

Hence

$$n^{-1} \max_{i,j} |w_j - w_i| \le n^{-1} \{ 2\sqrt{2c_1^{-1}} \sqrt{\log n} + 2(C_2p + 1)\phi_n \}.$$

Since  $n^{-1}p = o(1)$ , then  $\lambda \gg \phi_n \gg n^{-1}2(C_2p+1)\phi_n$ . Moreover,  $\lambda \gg \phi_n \gg n^{-1}\sqrt{\log n}$ . Hence

$$\lambda \gg n^{-1} \max_{i,j} |w_j - w_i|. \tag{A.14}$$

Let  $t_n = o(1)$ , then  $\rho'(4t_n) \to 1$ . Therefore, by (A.12), (A.13), and (A.14),

$$Q_n(\boldsymbol{\mu}, \boldsymbol{\beta}) - Q_n(\boldsymbol{\mu}^*, \boldsymbol{\beta}) = \Gamma_1 + \Gamma_2$$
  
 
$$\geq \sum_{i < j} [\lambda \rho'(4t_n) - n^{-1} \max_{i,j} |w_j - w_i|] |\mu_i - \mu_j| \ge 0,$$

for sufficiently large n, so that the result in (ii) is proved.

### A.5 Estimation procedure for model (2)

We let  $\widetilde{\mathbf{x}}_i = (1, \mathbf{x}_i^{\mathrm{T}})^{\mathrm{T}}$  and  $\boldsymbol{\beta}^* = (\boldsymbol{\mu}, \boldsymbol{\beta}^{\mathrm{T}})^{\mathrm{T}}$ . The model (2) can be written as  $y_i = \mathbf{z}_i^{\mathrm{T}} \boldsymbol{\theta}_i + \widetilde{\mathbf{x}}_i^{\mathrm{T}} \boldsymbol{\beta}^* + \epsilon_i, i = 1, \ldots, n$ . Similar to the assumption for model (1), we assume that observations can be divided into K different subgroups with K < n. Let  $\boldsymbol{\mathcal{G}} = (\boldsymbol{\mathcal{G}}_1, \ldots, \boldsymbol{\mathcal{G}}_K)$  be a partition of  $\{1, \ldots, n\}$ , and we assume  $\boldsymbol{\theta}_i = \boldsymbol{\alpha}_k$  for all  $i \in \boldsymbol{\mathcal{G}}_k$ , where  $\boldsymbol{\alpha}_k$  is the common value for the  $\boldsymbol{\theta}_i$ 's from group  $\boldsymbol{\mathcal{G}}_k$ . Then the estimates of  $\boldsymbol{\theta} = (\boldsymbol{\theta}_1^{\mathrm{T}}, \ldots, \boldsymbol{\theta}_n^{\mathrm{T}})^{\mathrm{T}}$  and  $\boldsymbol{\beta}^*$  can be obtained by minimizing

$$Q_n(\boldsymbol{\theta}, \boldsymbol{\beta}^*; \lambda) = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{z}_i^{\mathrm{T}} \boldsymbol{\theta}_i - \widetilde{\mathbf{x}}_i^{\mathrm{T}} \boldsymbol{\beta}^*)^2 + \sum_{1 \le i < j \le n} p(||\boldsymbol{\theta}_i - \boldsymbol{\theta}_j||, \lambda), \quad (A.15)$$

where  $p(\cdot, \lambda)$  is a concave penalty function with a tuning parameter  $\lambda$ , such as MCP or SCAD as described in Section 2. Then for a given  $\lambda > 0$ , define

$$(\widehat{\boldsymbol{\theta}}(\lambda), \widehat{\boldsymbol{\beta}}^*(\lambda)) = \operatorname{argmin} Q_n(\boldsymbol{\theta}, \boldsymbol{\beta}^*; \lambda)$$

The penalty shrinks some of  $||\boldsymbol{\theta}_i - \boldsymbol{\theta}_j||$  to zero. Based on this, we can partition the treatment effects into subgroups. Specifically, let  $\hat{\lambda}$  be the value of the tuning parameter selected based on a data-driven procedure such as the BIC. For simplicity, write  $(\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\beta}}^*) \equiv (\hat{\boldsymbol{\theta}}(\lambda), \hat{\boldsymbol{\beta}}^*(\lambda))$ . Let  $\{\hat{\boldsymbol{\alpha}}_1, \dots, \hat{\boldsymbol{\alpha}}_{\hat{K}}\}$  be the distinct values of  $\hat{\boldsymbol{\theta}}$ . Let  $\hat{\mathcal{G}}_k = \{i : \hat{\boldsymbol{\theta}}_i = \hat{\boldsymbol{\alpha}}_k, 1 \leq i \leq n\}, 1 \leq k \leq \hat{K}$ . Then  $\{\hat{\mathcal{G}}_1, \dots, \hat{\mathcal{G}}_{\hat{K}}\}$  constitutes a partition of  $\{1, \dots, n\}$ . Then we apply our proposed ADMM algorithm to obtain the estimates of  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}^*$  described as follows.

We reparametrize by introducing a new set of parameters  $\delta_{ij} = \theta_i - \theta_j$ , and hence minimization of (A.15) is equivalent to the constraint optimization problem:

$$S(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}) = \frac{1}{2} \sum_{i=1}^{n} (y_i - \mathbf{z}_i^{\mathrm{T}} \boldsymbol{\theta}_i - \widetilde{\mathbf{x}}_i^{\mathrm{T}} \boldsymbol{\beta}^*)^2 + \sum_{i < j} p_{\gamma}(||\boldsymbol{\delta}_{ij}||, \lambda),$$
  
subject to  $\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij} = \mathbf{0},$ 

where  $\boldsymbol{\delta} = \{\boldsymbol{\delta}_{ij}^{\mathrm{T}}, i < j\}^{\mathrm{T}}$ . By the augmented Lagrangian method (ALM), the estimates of the parameters can be obtained by minimizing

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v}) = S(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}) + \sum_{i < j} \langle \boldsymbol{v}_{ij}, \boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij} \rangle + \frac{\vartheta}{2} \sum_{i < j} ||\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij}||^2,$$

where the dual variables  $\boldsymbol{v} = \{\boldsymbol{v}_{ij}^{\mathrm{T}}, i < j\}^{\mathrm{T}}$  are Lagrange multipliers and  $\vartheta$  is the penalty parameter. We then can obtain the estimators of  $(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v})$  through iterations by the ADMM.

For given  $(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{v})$ , the minimizer of  $L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v})$  with respect to  $\boldsymbol{\delta}_{ij}$  is unique and has a closed-form expression for the L<sub>1</sub>, MCP and SCAD penalties, respectively. Specifically, for given  $(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{v})$ , the minimization problem is the same as minimizing

$$\frac{\vartheta}{2}\sum_{i < j} ||\boldsymbol{\zeta}_{ij} - \boldsymbol{\delta}_{ij}||^2 + \sum_{i < j} p_{\gamma}(||\boldsymbol{\delta}_{ij}||, \lambda)$$

with respect to  $\delta_{ij}$ , where  $\zeta_{ij} = \theta_i - \theta_j + \vartheta^{-1} v_{ij}$ . Hence, the closed-form solution for the L<sub>1</sub> penalty is

$$\widehat{\boldsymbol{\delta}}_{ij} = S(\boldsymbol{\zeta}_{ij}, \lambda/\vartheta), \tag{A.16}$$

where  $S(\mathbf{z}, t) = (1 - t/||\mathbf{z}||)_{+}\mathbf{z}$  is the groupwise soft thresholding rule, and  $(x)_{+} = x$  if x > 0and 0, otherwise. For the MCP penalty with  $\gamma > 1/\vartheta$ , it is

$$\widehat{\boldsymbol{\delta}}_{ij} = \begin{cases} \frac{S(\boldsymbol{\zeta}_{ij}, \lambda/\vartheta)}{1 - 1/(\gamma\vartheta)} & \text{if } ||\boldsymbol{\zeta}_{ij}|| \le \gamma\lambda \\ \boldsymbol{\zeta}_{ij} & \text{if } ||\boldsymbol{\zeta}_{ij}|| > \gamma\lambda. \end{cases}$$
(A.17)

For the SCAD penalty with  $\gamma > 1/\vartheta + 1$ , it is

$$\widehat{\boldsymbol{\delta}}_{ij} = \begin{cases} ST(\boldsymbol{\zeta}_{ij}, \lambda/\vartheta) & \text{if } ||\boldsymbol{\zeta}_{ij}|| \leq \lambda + \lambda/\vartheta \\ \frac{ST(\boldsymbol{\zeta}_{ij}, \gamma\lambda/((\gamma-1)\vartheta))}{1 - 1/((\gamma-1)\vartheta)} & \text{if } \lambda + \lambda/\vartheta < ||\boldsymbol{\zeta}_{ij}|| \leq \gamma\lambda \\ \boldsymbol{\zeta}_{ij} & \text{if } ||\boldsymbol{\zeta}_{ij}|| > \gamma\lambda. \end{cases}$$
(A.18)

ADMM algorithm for (A.15). We now describe the computational algorithm based on the ADMM for minimizing (A.15). It consists of iteratively updating  $\theta, \beta^*, \delta$  and v. The main ingredients of the algorithm are as follows.

First, for a given  $(\boldsymbol{\delta}, \boldsymbol{v})$ , to obtain an update of  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}^*$ , we set the derivatives  $\partial L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v}) / \partial \boldsymbol{\theta}$  and  $\partial L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v}) / \partial \boldsymbol{\beta}^*$  to zero, where

$$L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}, \boldsymbol{v}) = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{z}_i^{\mathrm{T}} \boldsymbol{\theta}_i - \widetilde{\mathbf{x}}_i^{\mathrm{T}} \boldsymbol{\beta}^*)^2 + \frac{\vartheta}{2} \sum_{i < j} ||\boldsymbol{\theta}_i - \boldsymbol{\theta}_j - \boldsymbol{\delta}_{ij} + \vartheta^{-1} \boldsymbol{v}_{ij}||^2 + C$$
  
$$= \frac{1}{2} ||\mathbf{Z}\boldsymbol{\theta} + \widetilde{\mathbf{X}} \boldsymbol{\beta}^* - \mathbf{y}||^2 + \frac{\vartheta}{2} ||\mathbf{A}\boldsymbol{\beta} - \boldsymbol{\delta} + \vartheta^{-1} \boldsymbol{v}||^2 + C.$$

Here *C* is a constant independent of  $\boldsymbol{\theta}$  and  $\boldsymbol{\beta}^*$ ,  $\mathbf{y} = (y_1, \ldots, y_n)^{\mathrm{T}}$ ,  $\mathbf{Z} = \operatorname{diag}(\mathbf{z}_1^{\mathrm{T}}, \ldots, \mathbf{z}_n^{\mathrm{T}})$  and  $\widetilde{\mathbf{X}} = (\widetilde{\mathbf{x}}_1, \ldots, \widetilde{\mathbf{x}}_n)^{\mathrm{T}}$ . Moreover,  $e_i$  is the  $n \times 1$  vector whose  $i^{\mathrm{th}}$  element is 1 and the remaining ones are 0,  $\Delta = \{(e_i - e_j), i < j\}^{\mathrm{T}}$  and  $\mathbf{A} = \Delta \otimes \mathbf{I}_p$ , where  $\mathbf{I}_d$  denotes the  $d \times d$  identity matrix and  $\otimes$  denotes the Kronecker product.

Thus for given  $\boldsymbol{\delta}^{(m)}$  and  $\boldsymbol{v}^{(m)}$  at the  $m^{\text{th}}$  step, the updates  $\boldsymbol{\theta}^{(m+1)}$  and  $\boldsymbol{\beta}^{*(m+1)}$ , which are the minimizers of  $L(\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}^{(m)}, \boldsymbol{v}^{(m)})$ , are

$$\boldsymbol{\theta}^{(m+1)} = (\mathbf{Z}^{\mathrm{T}}(\mathbf{I}_{n} - \mathbf{Q}_{\widetilde{\mathbf{X}}})\mathbf{Z} + \vartheta \mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}[\mathbf{Z}^{\mathrm{T}}(\mathbf{I}_{n} - \mathbf{Q}_{\widetilde{\mathbf{X}}})\mathbf{y} + \vartheta \mathbf{A}^{\mathrm{T}}(\boldsymbol{\delta}^{(m)} - \vartheta^{-1}\boldsymbol{v}^{(m)})],$$
  
$$\boldsymbol{\beta}^{*(m+1)} = (\widetilde{\mathbf{X}}^{\mathrm{T}}\widetilde{\mathbf{X}})^{-1}\widetilde{\mathbf{X}}^{\mathrm{T}}(\mathbf{y} - \mathbf{Z}\boldsymbol{\theta}^{(m+1)}),$$

where  $\mathbf{Q}_{\widetilde{\mathbf{X}}} = \widetilde{\mathbf{X}} (\widetilde{\mathbf{X}}^{\mathrm{T}} \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^{\mathrm{T}}.$ 

Second, the update of  $\boldsymbol{\delta}_{ij}$  at the  $(m+1)^{\text{th}}$  iteration is obtained by the formula given in (A.16), (A.17) and (A.18), respectively, by the Lasso, MCP and SCAD penalties with  $\boldsymbol{\zeta}_{ij}$  replaced by  $\boldsymbol{\zeta}_{ij}^{(m+1)} = \boldsymbol{\beta}_i^{(m+1)} - \boldsymbol{\beta}_j^{(m+1)} + \vartheta^{-1} \boldsymbol{v}_{ij}^{(m+1)}$ .

Finally, the estimate of  $\boldsymbol{v}_{ij}$  is updated as

$$\boldsymbol{v}_{ij}^{(m+1)} = \boldsymbol{v}_{ij}^{(m)} + \vartheta(\boldsymbol{\beta}_i^{(m+1)} - \boldsymbol{\beta}_j^{(m+1)} - \boldsymbol{\delta}_{ij}^{(m+1)}).$$

We iteratively update the estimates of  $\boldsymbol{\theta}, \boldsymbol{\beta}^*, \boldsymbol{\delta}$  and  $\boldsymbol{v}$  until the stopping rule is met. We track the progress of the ADMM based on the primal residual  $\mathbf{r}^{(m+1)} = \mathbf{A}\boldsymbol{\theta}^{(m+1)} - \boldsymbol{\delta}^{(m+1)}$ . We stop the algorithm when  $\mathbf{r}^{(m+1)}$  is close to zero such that  $\|\mathbf{r}^{(m+1)}\| < \epsilon$  for some small value  $\epsilon$ .

# References

Tseng, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimization, Journal of Optimization Theory and Applications, 109, 475-494.