Supplemental Material for "Distributed Algorithms for High-Dimensional Statistical Inference and Structure Learning with Heterogeneous Data"

Hongru Zhao and Xiaotong Shen

School of Statistics, University of Minnesota, Twin Cities

S1Heterogeneous Linear Regression with Site-Specific Heteroskedastic Errors

We extend the heterogeneous linear model in Remark 1 to allow each site j its own error variance $N_{n_j}(0, \sigma_j^2 I_{n_j})$. The negative log-likelihood (up to constant) is $\frac{n_j}{2} \log \sigma_j^2 + \frac{1}{2\sigma_j^2} \|\mathbf{Y}_j - \mathbf{Y}_j\|_{2\sigma_j^2}$

 $\mathbf{X}_{j}\boldsymbol{\beta}_{0}-\mathbf{W}_{j}\boldsymbol{\beta}_{j}\big\|_{2}^{2}$. Profiling out σ_{j}^{2} or equivalently reparameterizing via weighted least squares, one obtains the site-wise loss

$$L_j(\boldsymbol{\beta}_0, \boldsymbol{\beta}_j) = \frac{1}{2} \| \mathbf{Y}_j' - \mathbf{X}_j' \boldsymbol{\beta}_0 - \mathbf{W}_j' \boldsymbol{\beta}_j \|_2^2$$
 (S1.1)

with $\mathbf{Y}_j' = \frac{\mathbf{Y}_j}{\sigma_j}$, $\mathbf{X}_j' = \frac{\mathbf{X}_j}{\sigma_j}$, $\mathbf{W}_j' = \frac{\mathbf{W}_j}{\sigma_j}$. This weighted least squares loss simply replaces the loss L_j in (5.25). To compute the CMLE for the parameters $\{\boldsymbol{\beta}_j\}_{j=0}^K$ and variances $\{\sigma_j^2\}_{j=1}^K$, we use a block coordinate descent (BCD) algorithm, see Algorithm S1.

By treating the number of sites K as a fixed parameter, the constrained likelihood-ratio testing procedure and Theorem 2 remain valid without modification.

S2Proof of Theorem 1

Proof. We will show that if $\kappa \geq |A_{H_0}^0|$, then $\{i \in [p] \backslash B : \beta_i^0 \neq 0\} \subset \{i \in [p] \backslash B : \widehat{\Gamma}_i \neq 0\}$ almost surely, where the estimators $\widehat{\Gamma}_{H_0}$ and $\widehat{\Gamma}_{H_1}$ are obtained from Algorithm 3. For H_0 , let $A^0 = \{i \in [p] \backslash B : \beta_i^0 \neq 0\}$ and $A^{[t]} = \{i \in [p] \backslash B : |\widetilde{\Gamma}_i^{[t]}| \ge \tau\}. \text{ For } H_1, \text{ let } A^0 = \{i \in [p] \backslash B : \beta_i^0 \ne 0\} \cup B$ and $A^{[t]} = \{i \in [p] \setminus B : |\widetilde{\Gamma}_i^{[t]}| \geq \tau\} \cup B$. Set $\widehat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} - \mathbf{X}\widehat{\boldsymbol{\beta}}^{ol}$, where $\widehat{\boldsymbol{\beta}}^{ol}$ is the

Algorithm S1 BCD for CMLE in Heterogeneous Linear Regression with Heteroskedastic Errors

- 1: Initialize: Set t = 0 and $\sigma_j^{2(0)} = 1$ for $j = 1, \dots, K$.
- 2: while not converged do
- **Parameter update:** Treat (S1.1) (with the current $\{\hat{\sigma}_i^{2(t)}\}\)$ as the site-wise loss and run Algorithm 2 to obtain $\{\widehat{\boldsymbol{\beta}}_{j}^{(t+1)}\}_{j=0}^{K}$. Variance update: For each $j=1,\ldots,K,$

$$\widehat{\sigma}_{j}^{2\,(t+1)} = \frac{1}{n_{j}} \left\| \mathbf{Y}_{j} - \mathbf{X}_{j} \, \widehat{\boldsymbol{\beta}}_{0}^{(t+1)} - \mathbf{W}_{j} \, \widehat{\boldsymbol{\beta}}_{j}^{(t+1)} \right\|_{2}^{2}.$$

- $t \leftarrow t + 1$
- 6: end while
- 7: Output: $\{\widehat{\boldsymbol{\beta}}_i\}_{i=0}^K$ and $\{\widehat{\sigma}_i^2\}_{i=1}^K$.

oracle estimate that minimizes $\|\boldsymbol{Y} - \mathbf{X}\boldsymbol{\beta}\|_2^2$ under the constraint $\boldsymbol{\beta}_{(A^0)^c} = 0$. Set $E = \{ \| \mathbf{X}^T \widehat{\boldsymbol{\varepsilon}} / n \|_{\infty} \le 0.5 \lambda \tau \} \cap \{ \| \boldsymbol{\beta}^0 - \widehat{\boldsymbol{\beta}}^{ol} \|_{\infty} \le 0.5 \tau \}.$

We will show that $A^0 \triangle A^{[t]}$ is eventually empty set on event E, which has a probability tending to 1, where \triangle denotes the symmetric difference. By the optimality criterion Lee and Lee (2005) for (4.12) and (4.13), we have

$$\langle \widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}, -\mathbf{X}^T (\boldsymbol{Y} - \mathbf{X} \widetilde{\boldsymbol{\Gamma}}^{[t]}) / n + \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^c}^{[t]} \right\|_{\mathbf{I}} \rangle \ge 0.$$
 (S2.2)

Rearranging the terms, we have

$$\begin{aligned}
& \left\| \mathbf{X} (\widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}) \right\|_{2}^{2} / n \\
& \leq \langle \widetilde{\boldsymbol{\Gamma}}^{[t]} - \widehat{\boldsymbol{\beta}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle \\
& = \langle \widetilde{\boldsymbol{\Gamma}}_{A^{0} \backslash A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \backslash A^{[t-1]}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle \\
& + \langle \widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \backslash A^{0}} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \backslash A^{0}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle \\
& + \langle \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle \\
& + \langle \widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \cap A^{0}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \cap A^{0}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle.
\end{aligned} \tag{S2.3}$$

Let \mathbf{X}_{A^0} denote the matrix, whose A^0 columns are the same as the A^0 columns of X and all other columns are zeros. Similarly, we can define $\mathbf{X}_{(A^0)^c}. \text{ Notice that } \mathbf{X} = \mathbf{X}_{A^0} + \mathbf{X}_{(A^0)^c}, \ \widehat{\boldsymbol{\beta}}_{(A^0)^c}^{ol} = \mathbf{0}, \ \widehat{\boldsymbol{\beta}}^{ol} = (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^\dagger \mathbf{X}_{A^0}^T \boldsymbol{Y},$

$$\widehat{\boldsymbol{Y}} = \mathbf{X}_{A^0} (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger} \mathbf{X}_{A^0}^T \boldsymbol{Y}, \ P_{A^0} = \mathbf{X}_{A^0} (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger} \mathbf{X}_{A^0}^T, \text{ as well as } \widehat{\boldsymbol{\varepsilon}} = \boldsymbol{Y} - \widehat{\boldsymbol{Y}} = (I - P_{A^0}) \boldsymbol{Y} = (I - P_{A^0}) \boldsymbol{\varepsilon}. \text{ Note that } \mathbf{X}^T \widehat{\boldsymbol{\varepsilon}} = \mathbf{X}_{(A^0)^c}^T \widehat{\boldsymbol{\varepsilon}}.$$

Thus, we know that the fourth term in the last equation of (S2.3) vanishes, since

$$\langle \widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \cap A^{0}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \cap A^{0}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]})^{c}}^{[t]} \right\|_{1} \rangle$$

$$= \langle \widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \cap A^{0}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \cap A^{0}}^{ol}, \mathbf{X}_{(A^{0})^{c}}^{T} \widehat{\boldsymbol{\varepsilon}} / n \rangle$$

$$= \langle \mathbf{X}_{(A^{0})^{c}} (\widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \cap A^{0}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \cap A^{0}}^{ol}), \widehat{\boldsymbol{\varepsilon}} / n \rangle = 0.$$
(S2.4)

Note that the third term in the last equation of (S2.3) satisfies

$$\langle \widetilde{\Gamma}_{(A^{[t-1]} \cup A^0)^c}^{[t]} - \widehat{\beta}_{(A^{[t-1]} \cup A^0)^c}^{ol}, \nabla \left\| \widetilde{\Gamma}_{(A^{[t-1]})^c}^{[t]} \right\|_1 \rangle = \left\| \widetilde{\Gamma}_{(A^{[t-1]} \cup A^0)^c}^{[t]} - \widehat{\beta}_{(A^{[t-1]} \cup A^0)^c}^{ol} \right\|_1,$$
(S2.5)

because $\widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^0)^c}^{ol} = 0$. Recalling (S2.3), we have

$$0 \leq \left\| \mathbf{X} (\widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}) \right\|_{2}^{2} / n$$

$$\leq \langle \widetilde{\boldsymbol{\Gamma}}_{A^{0} \backslash A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \backslash A^{[t-1]}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{A^{0} \backslash A^{[t-1]}}^{[t]} \right\|_{1} \rangle$$

$$+ \langle \widetilde{\boldsymbol{\Gamma}}_{A^{[t-1]} \backslash A^{0}} - \widehat{\boldsymbol{\beta}}_{A^{[t-1]} \backslash A^{0}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n \rangle$$

$$+ \langle \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol}, \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} / n - \lambda \tau \nabla \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} \right\|_{1} \rangle$$

$$\leq \left\| \widetilde{\boldsymbol{\Gamma}}_{A^{0} \triangle A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \triangle A^{[t-1]}}^{ol} \right\|_{1} \cdot (\left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n + \lambda \tau)$$

$$+ \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol} \right\|_{1} \cdot (\left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n - \lambda \tau).$$
(S2.6)

Rearranging inequality (S2.6) implies that

$$\begin{split} & \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol} \right\|_{1} \cdot (\lambda \tau - \left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n) \\ \leq & \left\| \widetilde{\boldsymbol{\Gamma}}_{A^{0} \triangle A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \triangle A^{[t-1]}}^{ol} \right\|_{1} \cdot (\left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n + \lambda \tau). \end{split}$$

Over event E, we have

$$\begin{split} \left\| \widetilde{\Gamma}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\beta}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol} \right\|_{1} &\leq 3 \left\| \widetilde{\Gamma}_{A^{0} \triangle A^{[t-1]}}^{[t]} - \widehat{\beta}_{A^{0} \triangle A^{[t-1]}}^{ol} \right\|_{1} \\ &\leq 3 \left\| \widetilde{\Gamma}_{A^{0} \cup A^{[t-1]}}^{[t]} - \widehat{\beta}_{A^{0} \cup A^{[t-1]}}^{ol} \right\|_{1}. \end{split}$$

Recall that $\kappa_1 = \left| \{ i \in [p] \backslash B; |\widetilde{\Gamma}_i^{[0]}| \geq \tau \} \right|$ and $\kappa_{max} = \max\{\kappa, \kappa_1\}$. Without loss of generality, we can assume that $\widetilde{\Gamma}_B^{[0]} = \mathbf{0}$. For the base case, we know that

$$|A^0 \triangle A^{[0]}| \le |A^0 \backslash B| + |A^{[0]} \backslash A^0| \le 2\kappa_{max}.$$

For the induction step, assume $|A^0 \triangle A^{[t-1]}| \leq 2\kappa_{max}$ on event E, for $t \geq 1$. We aim to show that $|A^0 \triangle A^{[t]}| \leq 2\kappa_{max}$ over event E. Applying Assumption 1 and (S2.6), over event E, we have

$$c_{1} \left\| \widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]} \right\|_{2}^{2} \leq \left\| \mathbf{X} (\widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}) \right\|_{2}^{2} / n$$

$$\leq \left\| \widetilde{\boldsymbol{\Gamma}}_{A^{0} \triangle A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \triangle A^{[t-1]}}^{ol} \right\|_{1} \cdot (\left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n + \lambda \tau)$$

$$+ \left\| \widetilde{\boldsymbol{\Gamma}}_{(A^{[t-1]} \cup A^{0})^{c}}^{[t]} - \widehat{\boldsymbol{\beta}}_{(A^{[t-1]} \cup A^{0})^{c}}^{ol} \right\|_{1} \cdot (\left\| \mathbf{X}^{T} \widehat{\boldsymbol{\varepsilon}} \right\|_{\infty} / n - \lambda \tau)$$

$$\leq \frac{3}{2} \lambda \tau \left\| \widetilde{\boldsymbol{\Gamma}}_{A^{0} \triangle A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^{0} \triangle A^{[t-1]}}^{ol} \right\|_{1}.$$
(S2.7)

By Cauchy-Schwarz inequality,

$$\left\| \widetilde{\Gamma}_{A^0 \triangle A^{[t-1]}}^{[t]} - \widehat{\boldsymbol{\beta}}_{A^0 \triangle A^{[t-1]}}^{ol} \right\|_{1}^{2} \le \left| A^0 \triangle A^{[t-1]} \right| \left\| \widetilde{\Gamma}^{[t]} - \widehat{\boldsymbol{\beta}}^{ol} \right\|_{2}^{2}. \tag{S2.8}$$

Combining (S2.7), (S2.8) and the third condition in Theorem 1,

$$\left\|\widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}\right\|_{2} \le \frac{3\lambda\tau}{2c_{1}}\sqrt{|A^{0}\triangle A^{[t-1]}|} \le \frac{\tau}{4}\sqrt{|A^{0}\triangle A^{[t-1]}|}.$$
 (S2.9)

Applying (S2.9) and induction assumption $|A^0 \triangle A^{[t-1]}| \leq 2\kappa_{max}$, we have

$$\left\|\widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]}\right\|_{2} / \tau \le \frac{1}{4} \sqrt{2\kappa_{max}} \le \frac{1}{2} \sqrt{\kappa_{max}} \le \sqrt{\kappa_{max}}.$$
 (S2.10)

Because for any $i \in A^{[t]} \setminus A^0$, $|\widetilde{\Gamma}_i^{[t]} - \widehat{\beta}_i^{ol}| = |\widetilde{\Gamma}_i^{[t]}| \ge \tau$, we have

$$\sqrt{|A^{[t]} \backslash A^0|} \le \left\| \widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]} \right\|_2 / \tau \le \sqrt{\kappa_{\max}}.$$

Thus,

$$|A^0 \triangle A^{[t]}| \le |A^0 \backslash B| + |A^{[t]} \backslash A^0| \le 2\kappa_{max}.$$

By induction, we already showed that over event E, $|A^0 \triangle A^{[t]}| \leq 2\kappa_{max}$ for any $0 \leq t \leq t_{max}$, where t_{max} denotes the total number of the iteration of Algorithm 3.

Note that Algorithm 3 is terminated at t if

$$\operatorname{supp}\{\widetilde{\Gamma}^{[t]}\}\backslash B = \operatorname{supp}\{\widetilde{\Gamma}^{[t-1]}\}\backslash B.$$

We will show that over event E, for any $t \geq 1$,

$$\sqrt{|A^0 \triangle A^{[t]}|} \leq \frac{1}{2} \sqrt{|A^0 \triangle A^{[t-1]}|}.$$

If $i \in A^0 \backslash A^{[t]}$, then we know that $\beta_i^0 \neq 0$ and $\widetilde{\Gamma}_i^{[t]} = 0$. Over event E, $|\beta_i^0 - \widehat{\beta}_i^{ol}| \leq 0.5\tau$. According to assumption $\kappa \geq |A_{H_0}^0|$ in Theorem 1, we know that

$$|\widetilde{\Gamma}_i^{[t]} - \widehat{\beta}_i^{ol}| \ge |\widetilde{\Gamma}_i^{[t]} - \beta_i^{0}| - |\beta_i^{0} - \widehat{\beta}_i^{ol}| \ge \tau - 0.5\tau = 0.5\tau. \tag{S2.11}$$

If $i \in A^{[t]} \setminus A^0$, then we know that $\beta_i^0 = \widehat{\beta}_i^{ol} = 0$ and $|\widetilde{\Gamma}_i^{[t]}| \ge \tau$, which implies $|\widetilde{\Gamma}_i^{[t]} - \widehat{\beta}_i^{ol}| = |\widetilde{\Gamma}_i^{[t]}| \ge \tau$.

Over event E, we obtain that for any $i \in A^0 \triangle A^{[t]}$, $|\widetilde{\Gamma}_i^{[t]} - \widehat{\beta}_i^{ol}| \ge \tau - 0.5\tau = 0.5\tau$. Combining (S2.9), we obtain

$$\sqrt{|A^0 \triangle A^{[t]}|} \le \frac{1}{0.5\tau} \left\| \widehat{\boldsymbol{\beta}}^{ol} - \widetilde{\boldsymbol{\Gamma}}^{[t]} \right\|_2 \le \frac{1}{2} \sqrt{|A^0 \triangle A^{[t-1]}|}. \tag{S2.12}$$

In conclusion, over event E, we obtain

$$\sqrt{|A^0 \triangle A^{[t]}|} \le \frac{1}{2^t} \sqrt{2\kappa_{max}}.$$
 (S2.13)

If $t \ge \lceil \frac{\log(2\kappa_{max})}{\log 4} \rceil$, $\sqrt{|A^0 \triangle A^{[t]}|} < 1$, i.e., $A^0 \triangle A^{[t]} = \emptyset$.

Set $t_{max} = \lceil \frac{\log(2\kappa_{max})}{\log 4} \rceil$. We already showed that over event E,

$$\{i \in [p] \backslash B : |\widetilde{\Gamma}_i^{[t_{max}]}| \ge \tau\} = \{i \in [p] \backslash B : \beta_i^0 \ne 0\}.$$

This means that for any $j \notin \{i \in [p] \setminus B : |\widetilde{\Gamma}_i^{[t_{max}]}| \geq \tau\}, |\widetilde{\Gamma}_j^{[t_{max}]}| < \tau$, and for any $k \in \{i \in [p] \setminus B : \widetilde{\Gamma}_i^{[t_{max}]} \geq \tau\}, |\widetilde{\Gamma}_k^{[t_{max}]}| \geq \tau$. Thus, over the event E

$$supp{\beta^{0}} B = \{i \in [p] B : \beta_{i}^{0} \neq 0\}$$

$$= \{i \in [p] B : |\widetilde{\Gamma}_{i}^{[t_{max}]}| \geq \tau\} \subset supp{\widehat{\Gamma}} B.$$
(S2.14)

Next, we aim to bound the probability of event E. Let $\mathbf{a} := \mathbf{X}^T \widehat{\boldsymbol{\varepsilon}} = (a_1, \dots, a_p)^T$. It is obvious that $\mathbf{a} = \mathbf{X}^T (I - P_{A^0}) \boldsymbol{\varepsilon} = \mathbf{X}_{(A^0)^c}^T (I - P_{A^0}) \boldsymbol{\varepsilon}$.

Recall the error $\varepsilon \sim N(0, \Sigma)$ with $\Sigma = \sigma^2 I$. By Assumption 2, for any $1 \le l \le p$,

$$\operatorname{Var}(a_l) = \boldsymbol{e}_l^T \mathbf{X}^T (I - P_{A^0}) \Sigma (I - P_{A^0}) \mathbf{X} \boldsymbol{e}_l \le \sigma^2 (\mathbf{X}^T (I - P_{A^0}) \mathbf{X})_{ll} \le n \sigma^2 c_2^2.$$
(S2.15)

By the upper-tail inequality for sub-gaussian distribution and the third condition in Theorem 1, we have

$$\mathbb{P}(\|\mathbf{X}^T \widehat{\boldsymbol{\varepsilon}}\|_{\infty} / n > 0.5\lambda\tau) = \mathbb{P}(\|\boldsymbol{a}\|_{\infty} / n > 0.5\lambda\tau)$$

$$\leq \sum_{l=1}^{p} \mathbb{P}(|\boldsymbol{a}|_{l} > 0.5n\lambda\tau) \leq 2p \exp\left(-\frac{1}{8} \frac{n\lambda^2\tau^2}{\sigma^2 c_2^2}\right) \leq \frac{2}{p^3 n^4}.$$
(S2.16)

Set $\boldsymbol{b} := \widehat{\boldsymbol{\beta}}^{ol} - \boldsymbol{\beta}^{0}$. A direct calculation implies that

$$\boldsymbol{b} = (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger} \mathbf{X}_{A^0}^T \boldsymbol{\varepsilon} = (b_1, \cdots, b_p)^T.$$

By Assumption 2, for any $1 \le l \le p$,

$$\operatorname{Var}(b_l) = \boldsymbol{e}_l^T (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger} \mathbf{X}_{A^0}^T \boldsymbol{\Sigma} \mathbf{X}_{A^0} (\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger} \boldsymbol{e}_l$$
$$\leq \sigma^2 ((\mathbf{X}_{A^0}^T \mathbf{X}_{A^0})^{\dagger})_{ll} \leq \frac{1}{n} \sigma^2 c_3^2 I(l \in A^0).$$

By Assumptions 2 and (4.17), we have

$$\mathbb{P}\left(\left\|\widehat{\boldsymbol{\beta}}^{ol} - \boldsymbol{\beta}^{0}\right\|_{\infty} > 0.5\tau\right) \leq \sum_{i \in A^{0}} \mathbb{P}(|\boldsymbol{b}|_{i} > 0.5\tau)$$

$$\leq 2|A^{0}| \exp\left(-\frac{1}{8} \frac{n\tau^{2}}{\sigma^{2} c_{3}^{2}}\right) \leq \frac{2(\kappa_{H_{0}}^{0} + |B|)}{p^{4} n^{4}} = o\left(\frac{1}{p^{4} n^{3}}\right). \tag{S2.17}$$

Case 1: $\kappa = |A_{H_0}^0|$. In this case, (S2.14) implies that over the event E

$$\operatorname{supp}\{\beta^{0}\}\backslash B = \operatorname{supp}\{\widehat{\Gamma}\}\backslash B. \tag{S2.18}$$

Thus, under the requirement for (τ, λ) in Theorem 1,

$$\mathbb{P}(\operatorname{supp}\{\widehat{\Gamma}\}\backslash B \neq A^0\backslash B)$$

$$= \mathbb{P}(\{\operatorname{supp}\{\widehat{\Gamma}\}\backslash B \neq A^0\backslash B\} \cap E) + \mathbb{P}(\{\operatorname{supp}\{\widehat{\Gamma}\}\backslash B \neq A^0\backslash B\} \cap E^c). \tag{S2.19}$$

Notice that $E \subset \{ \sup\{\widehat{\Gamma}\} \setminus B = A^{[t_{max}]} \setminus B \}$, and $E \subset \{A^{[t_{max}]} \setminus B = A^0 \setminus B \}$. Thus,

$$\mathbb{P}(\{\operatorname{supp}\{\widehat{\Gamma}\}\setminus B \neq A^0\setminus B\} \cap E) = 0$$
, and

$$\mathbb{P}(\operatorname{supp}\{\widehat{\Gamma}\}\backslash B \neq A^0\backslash B) \leq \mathbb{P}(E^c) \leq \frac{C'}{n^3}, \tag{S2.20}$$

where C' is a absolute constant.

Note that $\{\widehat{\boldsymbol{\Gamma}} = \widehat{\boldsymbol{\beta}}^{ol}\} = \{ \sup\{\widehat{\boldsymbol{\Gamma}}\} \setminus B = A^0 \setminus B \}$. By Borel-Cantelli lemma, we have $\{\widehat{\boldsymbol{\Gamma}} = \widehat{\boldsymbol{\beta}}^{ol}\}$ almost surely as $n \to \infty$.

It remains to show that $\widehat{\boldsymbol{\beta}}^{ol}$ is a global minimizer of (4.8) or (4.9) with high probability.

Applying Theorem 2 of Shen et al. (2013) and its proof, with the degree of separation condition 4, we obtain

$$\mathbb{P}\left(\widehat{\boldsymbol{\beta}}^{\ell_0} \neq \widehat{\boldsymbol{\beta}}^{ol}\right) \leq \frac{e+1}{e-1} \exp\left(-\frac{n}{18\sigma^2} \left(C_{\min} - 36\frac{\log p + \log n}{n}\sigma^2\right)\right)$$
$$\leq \frac{e+1}{e-1} \frac{1}{p^2 n^2},$$

where $\widehat{\boldsymbol{\beta}}^{\ell_0}$ denotes the global minimizer of (4.8) or (4.9). By Borel-Cantelli lemma, we have $\{\widehat{\boldsymbol{\Gamma}} = \widehat{\boldsymbol{\beta}}^{ol} = \widehat{\boldsymbol{\beta}}^{\ell_0}\}$ almost surely as $n \to \infty$.

Case 2: $\kappa \geq |A_{H_0}^0|$. Thus, under the requirement for (τ, λ) in Theorem 1,

$$\mathbb{P}(\operatorname{supp}\{\boldsymbol{\beta}^{0}\}\backslash B \not\subset \operatorname{supp}\{\widehat{\boldsymbol{\Gamma}}\}\backslash B)$$

$$=\mathbb{P}(\{\operatorname{supp}\{\boldsymbol{\beta}^{0}\}\backslash B \not\subset \operatorname{supp}\{\widehat{\boldsymbol{\Gamma}}\}\backslash B\} \cap E)$$

$$+\mathbb{P}(\{\operatorname{supp}\{\boldsymbol{\beta}^{0}\}\backslash B \not\subset \operatorname{supp}\{\widehat{\boldsymbol{\Gamma}}\}\backslash B\} \cap E^{c}).$$

By (S2.14), we know that $E \subset \{\sup\{\beta^0\} \setminus B \subset \sup\{\widehat{\Gamma}\} \setminus B\}$. Thus,

$$\mathbb{P}(\{\operatorname{supp}\{\boldsymbol{\beta}^0\}\backslash B\not\subset\operatorname{supp}\{\widehat{\boldsymbol{\Gamma}}\}\backslash B\}\cap E)=0$$
 and

$$\mathbb{P}\left(\operatorname{supp}\{\boldsymbol{\beta}^{0}\}\backslash B \not\subset \operatorname{supp}\{\widehat{\boldsymbol{\Gamma}}\}\backslash B\right) \leq \mathbb{P}(E^{c}) \leq \frac{C'}{n^{3}}.$$
 (S2.21)

By the Borel-Cantelli lemma, we have $\{\sup\{\beta^0\}\setminus B\subset \sup\{\widehat{\Gamma}\}\setminus B\}$ almost surely as $n\to\infty$.

Under H_0 , we know that $\sup\{\widehat{\boldsymbol{\beta}}_{H_0}^{ol}\}\subset \sup\{\boldsymbol{\beta}^0\}\setminus B$, which implies that $\{\sup\{\widehat{\boldsymbol{\beta}}_{H_0}^{ol}\}\subset \sup\{\widehat{\boldsymbol{\Gamma}}\}\setminus B\}$ almost surely as $n\to\infty$.

Under H_1 , we know that $\sup\{\widehat{\boldsymbol{\beta}}_{H_1}^{ol}\} \subset \sup\{\boldsymbol{\beta}^0\} \cup B$, which implies that $\{\sup\{\widehat{\boldsymbol{\beta}}_{H_1}^{ol}\} \subset \sup\{\widehat{\boldsymbol{\Gamma}}\} \cup B\}$ almost surely as $n \to \infty$.

S3Proof of Theorem 2

Proof. By Theorem 1, we have

$$\lim_{n \to \infty} \mathbb{P}(\{\widehat{\boldsymbol{\Gamma}}^{(0)} = \widehat{\boldsymbol{\beta}}_{H_0}^{ls}\} \cap \{\widehat{\boldsymbol{\Gamma}}^{(1)} = \widehat{\boldsymbol{\beta}}_{H_1}^{ls}\}) = 1, \tag{S3.22}$$

where $\widehat{\Gamma}^{(0)}$ and $\widehat{\Gamma}^{(1)}$ are obtained from Algorithm 3, and $\widehat{\boldsymbol{\beta}}_{H_0}^{ls}$ and $\widehat{\boldsymbol{\beta}}_{H_1}^{ls}$ are obtained from (4.10) and (4.11), respectively.

The remainder of the proof follows the argument in Theorem 2 of Zhu et al. (2020), which establishes the sampling distribution of $\Lambda_n(B)$.

S4Proof of Theorem 3

Proof. By Theorem 1, we have

$$\lim_{n \to \infty} \mathbb{P}(\widehat{\Gamma}^{(1)} = \widehat{\boldsymbol{\beta}}^{ls}) = 1, \tag{S4.23}$$

where $\widehat{\pmb{\Gamma}}^{(1)}$ is obtained from Algorithm 3 and $\widehat{\pmb{\beta}}^{ls}=\widehat{\pmb{\beta}}^{ls}_{H_1}$ is obtained from (4.11), supported on $A^0 \cup B$.

Thus, to show (5.26), it suffices to show

$$\sqrt{n}(\widehat{\boldsymbol{\beta}}_{B}^{ls} - \boldsymbol{\beta}_{B}^{0}) \xrightarrow{d} N(0, \Sigma),$$
(S4.24)

Let $t \in \mathbb{R}^B$ and $u \in \mathbb{R}^{A^0 \cup B}$ such that $u_B = t$ and $u_{B^c} = 0$. Note that

$$\widehat{\boldsymbol{\beta}}_{A^0 \cup B}^{ls} = (\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})^{\dagger} (\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B}) \boldsymbol{\beta}_{A^0 \cup B}^0 + (\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})^{\dagger} \mathbf{X}_{A^0 \cup B}^T \boldsymbol{\varepsilon}, \tag{S4.25}$$

where $\mathbf{X}_{A^0 \cup B}$ denotes the sub-matrix with columns of $A^0 \cup B$. Due to $u \in \{\boldsymbol{\xi} \in \mathbb{R}^{A^0 \cup B}; \boldsymbol{\xi}_i = 0, i \notin B\} \subset \mathcal{R}(\mathbf{X}_{A^0 \cup B}^T)$, we know that there exists vector v such that $u = \mathbf{X}_{A^0 \cup B}^T v$ and

$$(\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})(\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})^{\dagger} u$$

$$= \mathbf{X}_{A^0 \cup B}^T \left(\mathbf{X}_{A^0 \cup B}(\mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})^{\dagger} \mathbf{X}_{A^0 \cup B}^T \right) v = \mathbf{X}_{A^0 \cup B}^T v = u.$$

By direct calculation of the characteristic function of $\sqrt{n}(\hat{\boldsymbol{\beta}}_B^{ls} - \boldsymbol{\beta}_B^0)$, we

know that

$$\mathbb{E}e^{it^{T}\sqrt{n}(\widehat{\boldsymbol{\beta}}_{B}^{ls}-\boldsymbol{\beta}_{B}^{0})} = \mathbb{E}e^{iu^{T}\sqrt{n}(\widehat{\boldsymbol{\beta}}_{A^{0}\cup B}^{ls}-\boldsymbol{\beta}_{A^{0}\cup B}^{0})}$$

$$=\mathbb{E}_{\mathbf{X}}\mathbb{E}_{\boldsymbol{\varepsilon}}e^{iu^{T}\sqrt{n}(\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B})^{\dagger}\mathbf{X}_{A^{0}\cup B}^{T}\boldsymbol{\varepsilon}}$$

$$=\mathbb{E}_{\mathbf{X}}\exp\{-1/2\cdot\sigma^{2}u^{T}n(\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B})^{\dagger}\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B}(\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B})^{\dagger}u\}$$

$$=\mathbb{E}_{\mathbf{X}}\exp\{-1/2\cdot\sigma^{2}u^{T}(1/n\cdot\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B})^{\dagger}u\}$$

$$=\mathbb{E}_{\mathbf{X}}\exp\{-1/2\cdot\sigma^{2}t^{T}(1/n\cdot\mathbf{X}_{A^{0}\cup B}^{T}\mathbf{X}_{A^{0}\cup B})^{\dagger}_{B,B}t\}.$$

Under the assumption of Moore–Penrose inverse $\sigma^2 \left(\frac{1}{n} \mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B}\right)_{B,B}^{\dagger}$ converges in distribution to some positive semi-definite matrix Σ , and by applying the continuous mapping theorem, we obtain that

$$\mathbb{E}_{\mathbf{X}} \exp\{-1/2 \cdot \sigma^2 t^T (1/n \cdot \mathbf{X}_{A^0 \cup B}^T \mathbf{X}_{A^0 \cup B})_{B,B}^{\dagger} t\} \to e^{-1/2t^T \Sigma t}, \quad (S4.26)$$

as $n \to \infty$, for any $t \in \mathbb{R}^B$.

By Lévy's continuity theorem, we complete the proof of weak convergence (S4.24). Therefore, the proof of Theorem 3 is completed. \Box

S5 Top- κ index set Selection

Algorithm S2 Threshold-Based Top- κ index set Selection

- 1: Inputs:
 - Current parameter estimates $\{\widetilde{\Gamma}_{k,j}^{[t]}\}_{(k,j)\in\mathcal{S}}$. Target sparsity level κ .
 - Initial thresholds a, b such that

$$\left|\left\{(k,j): |\widetilde{\Gamma}_{k,j}^{[t]}| \leq a\right\}\right| < \kappa \quad \text{and} \quad \left|\left\{(k,j): |\widetilde{\Gamma}_{k,j}^{[t]}| \leq b\right\}\right| > \kappa.$$

- 2: Initialize: Set $c \leftarrow \frac{a+b}{2}$.
- 3: while b a is not sufficiently small do
- 4: **Site-Level Operations:** Each site counts how many local parameters satisfy $|\widetilde{\Gamma}_{k,j}^{[t]}| \leq c$ and sends *only* this count to a central location.
- 5: Aggregation and Update:
 - Sum the local counts with the center count to get the total count

$$\left|\left\{(k,j): |\widetilde{\Gamma}_{k,j}^{[t]}| \le c\right\}\right|.$$

- If the total is exactly κ , break the loop.
- If the total is $<\kappa$, set $a \leftarrow c$. Otherwise, set $b \leftarrow c$. Update $c \leftarrow \frac{a+b}{2}$.
- 6: end while
- 7: Finalize Threshold c^* : Let $c^* \leftarrow c$ if the loop ended early due to an exact match, or set $c^* \leftarrow a$ if we finished the binary search without an exact match.
- 8: Construct the ℓ_0 Projected Set:
 - Each site identifies which local parameters exceed $|\widetilde{\Gamma}_{k,j}^{[t]}| > c^*$.
 - Let C be the union of those parameter indices across all sites,

$$C = \left\{ (k, j) : |\widetilde{\Gamma}_{k, j}^{[t]}| \ge \text{threshold } c^* \right\}.$$

- If $|C| < \kappa$, select additional $\kappa |C|$ parameters with $|\widetilde{\Gamma}_{k,j}^{[t]}| \in (a,b)$ until $|C| = \kappa$.
- 9: Output: The projected index set C.

References

- Lee, G. M. and K. B. Lee (2005). Vector variational inequalities for nondifferentiable convex vector optimization problems. *Journal of Global Optimization* 32, 597–612.
- Shen, X., W. Pan, Y. Zhu, and H. Zhou (2013). On constrained and regularized high-dimensional regression. *Annals of the Institute of Statistical Mathematics* 65(5), 807–832.
- Zhu, Y., X. Shen, and W. Pan (2020). On high-dimensional constrained maximum likelihood inference. *Journal of the American Statistical Association* 115(529), 217–230.