

Supplementary Material for “Byzantine-tolerant Distributed Multiclass Sparse Discriminant Analysis”

Abstract

This document provides supplementary material to the article “Byzantine-tolerant Distributed Multiclass Sparse Discriminant Analysis” written by the same authors.

1 Proofs of Main Results

1.1 Proof of Theorem 3.1

Theorem 1.1. *Let $\max_{2 \leq k \leq K} \|\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*\|_1 = O_{\mathbb{P}}(a_n)$ and choose some sufficiently large positive constant η_1 such that*

$$\lambda_1 = \begin{cases} \eta_1 \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} \right), & \text{System I} \\ \eta_1 \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} + \frac{\alpha}{\sqrt{n}} + \frac{1}{n} \right), & \text{System II} \end{cases}$$

under conditions (C1), (C3) and (C5), we have

$$\max_{2 \leq k \leq K} \left\| \widehat{\boldsymbol{\theta}}_k^{(1)} - \boldsymbol{\theta}_k^* \right\|_2 = O_{\mathbb{P}}(\sqrt{s}\lambda_1), \quad (1.1)$$

and

$$\max_{2 \leq k \leq K} \left\| \widehat{\boldsymbol{\theta}}_k^{(1)} - \boldsymbol{\theta}_k^* \right\|_1 = O_{\mathbb{P}}(s\lambda_1). \quad (1.2)$$

Note that the initial error for the t -th iteration would be $\max_{2 \leq k \leq K} \|\widehat{\boldsymbol{\theta}}_k^{(t-1)} - \boldsymbol{\theta}_k^*\|_1$ by plug-in the initial estimator $\widehat{\boldsymbol{\theta}}_k^{(t-1)}$. Then we can obtain the ℓ_1 and ℓ_2 error bound of the $\widehat{\boldsymbol{\theta}}_k^{(t)}$ easily by induction according to (1.1) and (1.2) in Theorem 1.1. In the next, to show the proof of Theorem 1.1, we present several useful lemmas in the following.

Lemma 1.1. *For $\mathbf{x}_k \in \mathbb{R}^p, k = 1, \dots, K-1$ such that*

$$\sum_{k=1}^{K-1} \|\mathbf{x}_k\|_1 \leq 4\sqrt{s(K-1)} \left(\sum_{k=1}^{K-1} \|\mathbf{x}_k\|_2^2 \right)^{1/2},$$

we have

$$\sum_{k=1}^{K-1} \mathbf{x}_k^T \widehat{\Sigma}_{(0)} \mathbf{x}_k^T \geq L \sum_{k=1}^{K-1} \|\mathbf{x}_k\|_2^2,$$

holds with probability tending to 1.

Lemma 1.2. Let $\|\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*\|_1 = O_{\mathbb{P}}(a_n)$, then for $k = 2, \dots, K$ we have

$$\left\| \widehat{\Sigma}_{(0)}(\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \widehat{\mathbf{b}}_{k,0} \right\|_{\infty} = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} \right),$$

under System I.

Lemma 1.3. Let $\|\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*\|_1 = O_{\mathbb{P}}(a_n)$, then for $k = 2, \dots, K$ we have

$$\left\| \widehat{\Sigma}_{(0)}(\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \widetilde{\mathbf{b}}_{k,0} \right\|_{\infty} = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} + \frac{\alpha}{\sqrt{n}} + \frac{1}{n} \right),$$

under System II.

The proofs of above lemmas are relegated to Section 2 in the following.

Proof of Theorem 3.1. For simplicity of notations, we use $\widehat{\boldsymbol{\theta}}_k$ to denote $\widehat{\boldsymbol{\theta}}_k^{(1)}$. By the optimality of $(\widehat{\boldsymbol{\theta}}_2, \dots, \widehat{\boldsymbol{\theta}}_K)$, we have

$$\begin{aligned} & \frac{1}{2} \widehat{\boldsymbol{\theta}}_k^T \widehat{\Sigma}_{(0)} \widehat{\boldsymbol{\theta}}_k - \left(\widehat{\Sigma}_{(0)} \widehat{\boldsymbol{\theta}}_k^{(0)} - \mathbf{b}_{k,0} \right)^T \widehat{\boldsymbol{\theta}}_k + \lambda_1 \sum_{j=1}^p \left(\sum_{l=2}^K \widehat{\boldsymbol{\theta}}_{l,j}^2 \right)^{1/2} \\ & \leq \frac{1}{2} \boldsymbol{\theta}_k^{*T} \widehat{\Sigma}_{(0)} \boldsymbol{\theta}_k^* - \left(\widehat{\Sigma}_{(0)} \widehat{\boldsymbol{\theta}}_k^{(0)} - \mathbf{b}_{k,0} \right)^T \boldsymbol{\theta}_k^* + \lambda_1 \sum_{j=1}^p \left(\sum_{l \neq k} \widehat{\boldsymbol{\theta}}_{l,j}^2 + (\boldsymbol{\theta}_{k,j}^*)^2 \right)^{1/2}. \end{aligned}$$

By rearranging the terms above, we have

$$\begin{aligned} & \frac{1}{2} \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right)^T \widehat{\Sigma}_{(0)} \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right) \\ & \leq \left(\widehat{\Sigma}_{(0)}(\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \mathbf{b}_{k,0} \right)^T \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right) + \lambda_1 \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_1. \end{aligned} \tag{1.3}$$

According to Lemmas 1.2 and 1.3,

$$\left\| \widehat{\Sigma}_{(0)}(\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \mathbf{b}_{k,0} \right\|_{\infty} = \begin{cases} O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} \right), & \text{System I} \\ O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{N}} + a_n \sqrt{\frac{\log p}{n}} + \frac{\alpha}{\sqrt{n}} + \frac{1}{n} \right). & \text{System II} \end{cases}$$

Thus $\|\widehat{\Sigma}_{(0)} \boldsymbol{\theta}_k^* - \mathbf{b}_{k,0}\|_{\infty} \leq \lambda_1/2$ holds with high probability for some sufficiently large positive constant η_1 under both System I and II. Then (1.3) indicates that

$$\frac{1}{2} \sum_{k=2}^K \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right)^T \widehat{\Sigma}_{(0)} \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right) \leq \frac{3\lambda_1}{2} \sum_{k=2}^K \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_1. \tag{1.4}$$

By the optimality of $(\hat{\boldsymbol{\theta}}_2, \dots, \hat{\boldsymbol{\theta}}_K)$, we can also obtain that

$$\begin{aligned} & \frac{1}{2} \sum_{k=2}^K (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*)^\top \hat{\boldsymbol{\Sigma}}_{(0)} (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*) + \lambda_1 \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)}\|_2 \\ & \leq \sum_{k=2}^K (\hat{\boldsymbol{\Sigma}}_{(0)}(\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \mathbf{b}_{k,0})^\top (\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^*) + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2. \end{aligned} \quad (1.5)$$

Let $\mathbf{c}_k = \hat{\boldsymbol{\Sigma}}_{(0)}(\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \mathbf{b}_{k,0}$ and $\mathbf{c}_{(j)} = (\mathbf{c}_{2,j}, \dots, \mathbf{c}_{K,j})^\top$, then it follows from (1.5) and $\|\mathbf{c}_k\|_\infty \leq \lambda_1/2$ that

$$\begin{aligned} \lambda_1 \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)}\|_2 & \leq \sum_{j=1}^p \mathbf{c}_{(j)}^\top (\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*) + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2 \\ & \leq \sum_{j=1}^p \|\mathbf{c}_{(j)}\|_2 \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2 \\ & \leq \max_j \|\mathbf{c}_{(j)}\|_2 \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2 \\ & \leq \sqrt{K-1} \max_{2 \leq k \leq K} \|\mathbf{c}_k\|_\infty \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2 \\ & \leq \frac{\lambda_1}{2} \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 + \lambda_1 \sum_{j=1}^p \|\boldsymbol{\theta}_{(j)}^*\|_2, \end{aligned}$$

which implies that $\sum_{j \in S^c} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 \leq 3 \sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2$. In conjunction with the fact that

$$\left(\sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 \right)^2 \leq s \sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2^2,$$

we have

$$\begin{aligned} \sum_{j \in S^c} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_1 & \leq \sqrt{K-1} \sum_{j \in S^c} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 \leq 3\sqrt{K-1} \sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2 \\ & \leq 3\sqrt{s(K-1)} \left(\sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_2^2 \right)^{1/2} \\ & \leq 3\sqrt{s(K-1)} \left(\sum_{k=2}^K \|\hat{\boldsymbol{\theta}}_{k,S}^{(1)} - \boldsymbol{\theta}_{k,S}^*\|_2^2 \right)^{1/2}. \end{aligned}$$

Similarly, we have $\sum_{j \in S} \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_1 \leq \sqrt{s(K-1)} (\sum_{k=2}^K \|\hat{\boldsymbol{\theta}}_{k,S}^{(1)} - \boldsymbol{\theta}_{k,S}^*\|_2^2)^{1/2}$. It implies that

$$\sum_{k=2}^K \|\boldsymbol{\theta}_k^* - \hat{\boldsymbol{\theta}}_k\|_1 = \sum_{j=1}^p \|\hat{\boldsymbol{\theta}}_{(j)} - \boldsymbol{\theta}_{(j)}^*\|_1 \leq 4\sqrt{s(K-1)} \left(\sum_{k=2}^K \|\hat{\boldsymbol{\theta}}_k^{(1)} - \boldsymbol{\theta}_k^*\|_2^2 \right)^{1/2}. \quad (1.6)$$

By applying Lemma 1.1, we have

$$\frac{1}{2} \sum_{k=2}^K \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right)^T \widehat{\boldsymbol{\Sigma}}_{(0)} \left(\widehat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}_k^* \right) \geq L \sum_{k=2}^K \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_2^2. \quad (1.7)$$

Combining the inequalities (1.4), (1.6) and (1.7), we have

$$\left(\sum_{k=2}^K \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_2^2 \right)^{1/2} \leq \frac{6\sqrt{s(K-1)}}{L} \lambda_1,$$

and

$$\sum_{k=2}^K \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_1 \leq \frac{24s(K-1)}{L} \lambda_1.$$

It also indicates that

$$\max_{2 \leq k \leq K} \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_2 = O_{\mathbb{P}}(\sqrt{s}\lambda_1),$$

and

$$\max_{2 \leq k \leq K} \left\| \boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k \right\|_1 = O_{\mathbb{P}}(s\lambda_1).$$

□

1.2 Proof of Theorem 3.2

It suffices to prove Theorem 1.2 in the following.

Theorem 1.2. *Under conditions (C1)-(C5), with the same choice of λ_1 as in Theorem 1.1, we have $\widehat{S}^{(1)} \subseteq S$ holds with probability tending to 1 and $\widehat{\boldsymbol{\theta}}_k^{(1)}$ satisfies that*

$$\left\| \widehat{\boldsymbol{\theta}}_k^{(1)} - \boldsymbol{\theta}_k^* \right\|_{\infty} = O_{\mathbb{P}} \left(\left\| \boldsymbol{\Sigma}_{SS}^{-1} \right\|_{\infty} \lambda_1 \right). \quad (1.8)$$

Moreover, suppose that there exists a sufficiently large constant $C > 0$ such that

$$\theta_{\min}^* \geq C \left\| \boldsymbol{\Sigma}_{SS}^{-1} \right\|_{\infty} \lambda_1, \quad (1.9)$$

we have $\widehat{S}^{(1)} = S$ with probability tending to 1.

Lemma 1.4. *By partitioning $\boldsymbol{\Sigma}$ as*

$$\boldsymbol{\Sigma} = \begin{pmatrix} \boldsymbol{\Sigma}_{SS} & \boldsymbol{\Sigma}_{SS^c} \\ \boldsymbol{\Sigma}_{S^cS} & \boldsymbol{\Sigma}_{S^cS^c} \end{pmatrix},$$

and $\boldsymbol{\mu}_k$ according to sets S and S^c for $k = 2, \dots, K$ respectively, we have

$$\boldsymbol{\theta}_{k,S}^* = \boldsymbol{\Sigma}_{SS}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)_S, \quad (1.10)$$

and

$$(\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)_{S^c} = \boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)_S. \quad (1.11)$$

Proof of Theorem 1.2. Here we only prove the results in System II and the proof for System I is similar. First we define the oracle sub-problem as

$$\widehat{\boldsymbol{\theta}}_S^o = \arg \min_{\boldsymbol{\theta}_{k,S^c}=\mathbf{0}} \sum_{k=2}^K \left\{ \frac{1}{2} \boldsymbol{\theta}_k^T \widehat{\boldsymbol{\Sigma}}_{(0)} \boldsymbol{\theta}_k - \widetilde{\mathbf{b}}_{k,g-1}^T \boldsymbol{\theta}_k \right\} + \lambda_1 \sum_{j \in S} \|\boldsymbol{\theta}_{(j)}\|_2. \quad (1.12)$$

Once we show $\widehat{\boldsymbol{\theta}}_k^{(1)} = (\widehat{\boldsymbol{\theta}}_{k,S}^o, \mathbf{0})$ is the solution to (7), it is clear that $\widehat{S}^{(1)} \subseteq S$. According to the KKT condition, for any $j \in S$,

$$\begin{pmatrix} \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS} \widehat{\boldsymbol{\theta}}_{2,S}^o - (\widetilde{\mathbf{b}}_{2,0})_S \right)_j \\ \vdots \\ \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS} \widehat{\boldsymbol{\theta}}_{K,S}^o - (\widetilde{\mathbf{b}}_{K,0})_S \right)_j \end{pmatrix} + \lambda_1 \mathbf{Z}_j = 0, \quad (1.13)$$

and for any $j \notin S$,

$$\begin{pmatrix} \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS^c} \widehat{\boldsymbol{\theta}}_{2,S}^o - (\widetilde{\mathbf{b}}_{2,0})_{S^c} \right)_j \\ \vdots \\ \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS^c} \widehat{\boldsymbol{\theta}}_{K,S}^o - (\widetilde{\mathbf{b}}_{K,0})_{S^c} \right)_j \end{pmatrix} + \lambda_1 \mathbf{Z}_j = 0, \quad (1.14)$$

where $\mathbf{Z}_j \in \mathbb{R}^{K-1}$ is subgradient of $\|\boldsymbol{\theta}\|_2$ evaluated at $\widehat{\boldsymbol{\theta}}_{(j)}$. It suffices to show that

$$\lambda_1^{-1} \max_{j \in S^c} \left(\sum_{k=2}^K \left\{ \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\theta}}_{k,S}^o \right)_j - (\widetilde{\mathbf{b}}_{k,0})_j \right\}^2 \right)^{1/2} < 1, \quad (1.15)$$

holds with probability tending to 1. From equation (1.13), we have $\widehat{\boldsymbol{\theta}}_{k,S}^o = \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} ((\widetilde{\mathbf{b}}_{k,0})_S - \lambda_1 \widetilde{\mathbf{Z}}_{k,S})$ where $\widetilde{\mathbf{Z}}_{k,S} = (Z_{j,k} : j \in S) \in \mathbb{R}^s$ and $\sum_{k=2}^K (Z_{j,k})^2 = 1$. Note that

$$\begin{aligned} & \widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\theta}}_{k,S}^o - (\widetilde{\mathbf{b}}_{k,0})_{S^c} \\ &= \widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} ((\widetilde{\mathbf{b}}_{k,0})_S - \lambda_1 \widetilde{\mathbf{Z}}_{k,S}) - (\widetilde{\mathbf{b}}_{k,0})_{S^c} \\ &= \widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \left\{ \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS} - \boldsymbol{\Sigma}_{SS} \right) \left(\widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - \boldsymbol{\theta}_{k,S}^* \right) + \boldsymbol{\Sigma}_{SS} \widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - (\widetilde{\mathbf{d}}_{k,0})_S + (\widetilde{\boldsymbol{\mu}}_k - \widetilde{\boldsymbol{\mu}}_1)_S - \boldsymbol{\Sigma}_{SS} \boldsymbol{\theta}_{k,S}^* \right\} \\ &\quad - \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^c} - \boldsymbol{\Sigma}_{S^c S} \right) \left(\widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - \boldsymbol{\theta}_{k,S}^* \right) - \boldsymbol{\Sigma}_{SS} \widehat{\boldsymbol{\theta}}_{k,S}^{(0)} + (\widetilde{\mathbf{d}}_{k,0})_{S^c} + \boldsymbol{\Sigma}_{S^c S} \boldsymbol{\theta}_{k,S}^* - (\widetilde{\boldsymbol{\mu}}_k - \widetilde{\boldsymbol{\mu}}_1)_{S^c} \\ &\quad - \lambda_1 \widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \widetilde{\mathbf{Z}}_{k,S}. \end{aligned}$$

We denote

$$I_1 = \widehat{\boldsymbol{\Sigma}}_{(0),S^c} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \left\{ \left(\widehat{\boldsymbol{\Sigma}}_{(0),SS} - \boldsymbol{\Sigma}_{SS} \right) \left(\widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - \boldsymbol{\theta}_{k,S}^* \right) + \boldsymbol{\Sigma}_{SS} \widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - (\widetilde{\mathbf{d}}_{k,0})_S + (\widetilde{\boldsymbol{\mu}}_k - \widetilde{\boldsymbol{\mu}}_1)_S - \boldsymbol{\Sigma}_{SS} \boldsymbol{\theta}_{k,S}^* \right\},$$

and

$$I_2 = \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^c} - \boldsymbol{\Sigma}_{S^c S} \right) \left(\widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - \boldsymbol{\theta}_{k,S}^* \right) + \boldsymbol{\Sigma}_{SS} \widehat{\boldsymbol{\theta}}_{k,S}^{(0)} - (\widetilde{\mathbf{d}}_{k,0})_{S^c} - \boldsymbol{\Sigma}_{S^c S} \boldsymbol{\theta}_{k,S}^* + (\widetilde{\boldsymbol{\mu}}_k - \widetilde{\boldsymbol{\mu}}_1)_{S^c}.$$

Observe that

$$\begin{aligned}\widehat{\Sigma}_{(0),S^cS}\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{S^cS}\Sigma_{SS}^{-1} &= \left(\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right)\left(\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right)+\Sigma_{S^cS}\left(\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right) \\ &\quad +\Sigma_{SS}^{-1}\left(\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right),\end{aligned}$$

which implies

$$\begin{aligned}\left\|\widehat{\Sigma}_{(0),S^cS}\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{S^cS}\Sigma_{SS}^{-1}\right\|_\infty &\leq\left\|\left(\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right)\left(\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right)\right\|_\infty \\ &\quad +\left\|\Sigma_{S^cS}\left(\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right)\right\|_\infty+\left\|\Sigma_{SS}^{-1}\left(\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right)\right\|_\infty \\ &\leq s^{3/2}\left|\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right|_\infty\left\|\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right\|_2 \\ &\quad +s^{3/2}\left|\Sigma_{S^cS}\right|_\infty\left\|\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right\|_2+s^{3/2}\left|\Sigma_{SS}^{-1}\right|_\infty\left\|\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right\|_2.\end{aligned}$$

Using the inequalities (58a) and (58b) in [Wainwright \[2009\]](#), we have

$$\left\|\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{SS}^{-1}\right\|_2=O_{\mathbb{P}}\left(\sqrt{\frac{s}{n}}\right),$$

and

$$\left\|\widehat{\Sigma}_{(0),SS}-\Sigma_{SS}\right\|_2=O_{\mathbb{P}}\left(\sqrt{\frac{s}{n}}\right).$$

Combining with the fact $\left|\widehat{\Sigma}_{(0),S^cS}-\Sigma_{S^cS}\right|_\infty=O_{\mathbb{P}}(\sqrt{\log p/n})$, it yields

$$\left\|\widehat{\Sigma}_{(0),S^cS}\widehat{\Sigma}_{(0),SS}^{-1}-\Sigma_{S^cS}\Sigma_{SS}^{-1}\right\|_\infty=O_{\mathbb{P}}\left(s^{3/2}\sqrt{\frac{\log p+s}{n}}\right).$$

Owing to the fact that $\boldsymbol{\theta}_{k,S}=\Sigma_{SS}^{-1}(\boldsymbol{\mu}_k-\boldsymbol{\mu}_1)_S$ in Lemma 1.4, we have

$$(\widetilde{\boldsymbol{\mu}}_k-\widetilde{\boldsymbol{\mu}}_1)_S-\Sigma_{SS}\boldsymbol{\theta}_{k,S}^*=(\widetilde{\boldsymbol{\mu}}_k-\widetilde{\boldsymbol{\mu}}_1)_S-(\boldsymbol{\mu}_k-\boldsymbol{\mu}_1)_S.$$

It yields that

$$\left\|(\widetilde{\boldsymbol{\mu}}_k-\widetilde{\boldsymbol{\mu}}_1)_S-\Sigma_{SS}\boldsymbol{\theta}_{k,S}^*\right\|_\infty=O_{\mathbb{P}}\left(\sqrt{\frac{\log p}{N}}\right).$$

Moreover, note that

$$\left\|\left(\widehat{\Sigma}_{(0),SS}-\Sigma_{SS}\right)\left(\widehat{\boldsymbol{\theta}}_{k,S}^{(0)}-\boldsymbol{\theta}_{k,S}^*\right)\right\|_\infty\leq\left|\widehat{\Sigma}_{(0),SS}-\Sigma_{SS}\right|_\infty\left\|\widehat{\boldsymbol{\theta}}_{k,S}^{(0)}-\boldsymbol{\theta}_{k,S}^*\right\|_1=O_{\mathbb{P}}\left(\sqrt{\frac{\log p}{n}}a_n\right).$$

Then together with the assumption $\|\Sigma_{S^cS}\Sigma_{SS}^{-1}\|_\infty\leq\xi$ and Lemma 1.2 we have

$$\|I_1\|_\infty=O_{\mathbb{P}}\left(\sqrt{\frac{\log p}{N}}+\sqrt{\frac{\log p}{n}}a_n+\frac{\alpha}{\sqrt{n}}+\frac{1}{n}\right).$$

Similarly, we can show that

$$\|I_2\|_\infty = O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{N}} + \sqrt{\frac{\log p}{n}} a_n + \frac{\alpha}{\sqrt{n}} + \frac{1}{n} \right).$$

Owing to the choice of λ_1 in Theorem 1.1 and following the analysis above, there exists some positive constant C_1 such that for $k = 2, \dots, K$,

$$\begin{aligned} & \left\| \widehat{\boldsymbol{\theta}}_{k,S}^o - \boldsymbol{\theta}_{k,S}^* \right\|_\infty \\ & \leq \lambda_1 \left\| \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \widetilde{\mathbf{Z}}_{k,S} \right\|_\infty + \left\| \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} (\widehat{\boldsymbol{\Sigma}}_{(0),SS} - \boldsymbol{\Sigma}_{SS}) (\widehat{\boldsymbol{\theta}}_{k,S}^o - \boldsymbol{\theta}_{k,S}^*) \right\|_\infty \\ & \quad + \left\| \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} (\boldsymbol{\Sigma}_{SS} \widehat{\boldsymbol{\theta}}_{k,S}^o - (\tilde{\mathbf{d}}_{k,0})_S) \right\|_\infty + \left\| \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} ((\tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\mu}}_1)_S - \boldsymbol{\Sigma}_{SS} \boldsymbol{\theta}_{k,S}^*) \right\|_\infty \\ & \leq C_1 \left\| \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty \lambda_1, \end{aligned} \tag{1.16}$$

holds with probability tending to 1. Moreover, for $j \in S^c$, we have

$$\begin{aligned} \left| \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^cS} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \widetilde{\mathbf{Z}}_{k,S} \right)_j \right| & \leq \left\| \widehat{\boldsymbol{\Sigma}}_{(0),S^cS} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} - \boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty \left\| \widetilde{\mathbf{Z}}_{k,S} \right\|_\infty \\ & \quad + \left\| \boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty \left\| \widetilde{\mathbf{Z}}_{k,S} - \mathbf{Z}_{k,S}^* \right\|_\infty + \left| (\boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \mathbf{Z}_{k,S}^*)_j \right|, \end{aligned}$$

and

$$\begin{aligned} \left\| \widetilde{\mathbf{Z}}_{k,S} - \mathbf{Z}_{k,S}^* \right\|_\infty & = \max_{j \in S} \left| \frac{\widehat{\theta}_{kj}^o}{\|\boldsymbol{\theta}_{(j)}^o\|_2} - \frac{\theta_{kj}^*}{\|\boldsymbol{\theta}_{(j)}^*\|_2} \right| \\ & \leq \max_{j \in S} \frac{\left| \widehat{\theta}_{kj}^o - \theta_{kj}^* \right|}{\|\boldsymbol{\theta}_{(j)}^*\|_2} + \max_{j \in S} \left| \widehat{\theta}_{kj}^o \right| \frac{\left| \|\boldsymbol{\theta}_{(j)}^o\|_2 - \|\boldsymbol{\theta}_{(j)}^*\|_2 \right|}{\|\boldsymbol{\theta}_{(j)}^o\|_2 \|\boldsymbol{\theta}_{(j)}^*\|_2} \\ & \leq \max_{j \in S} \frac{\left| \widehat{\theta}_{kj}^o - \theta_{kj}^* \right|}{\|\boldsymbol{\theta}_{(j)}^*\|_2} + \max_{j \in S} \frac{\|\boldsymbol{\theta}_{(j)}^o - \boldsymbol{\theta}_{(j)}^*\|_2}{\|\boldsymbol{\theta}_{(j)}^*\|_2} \\ & \lesssim 2 \max_{2 \leq k \leq K} \left\| \widehat{\boldsymbol{\theta}}_{k,S}^o - \boldsymbol{\theta}_{k,S}^* \right\|_\infty / \theta_{\min}^*. \end{aligned}$$

Combining with the Conditions (C2) and inequality (1.16), with probability tending to 1 we have

$$\begin{aligned} & \lambda_1^{-2} \max_{j \in S^c} \sum_{k=2}^K \left\{ \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^cS} \widehat{\boldsymbol{\theta}}_{k,S}^o \right)_j - (\mathbf{b}_{k,g-1})_j \right\}^2 \leq \sum_{k=2}^K \left| \left(\widehat{\boldsymbol{\Sigma}}_{(0),S^cS} \widehat{\boldsymbol{\Sigma}}_{(0),SS}^{-1} \widetilde{\mathbf{Z}}_{k,S} \right)_j \right|^2 + o(1) \\ & \leq \max_{j \in S^c} \sum_{k=2}^K \left| (\boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \mathbf{Z}_{k,S}^*)_j \right|^2 + C_1^2 \left\| \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty^2 \left\| \boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty^2 (K-1) \lambda_1^2 / \theta_{\min}^{*2} + o(1) \\ & \leq 1 - \kappa + C_1^2 \left\| \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty^2 \left\| \boldsymbol{\Sigma}_{S^cS} \boldsymbol{\Sigma}_{SS}^{-1} \right\|_\infty^2 (K-1) \lambda_1^2 / \theta_{\min}^{*2} + o(1) \\ & \leq 1 - \kappa / 2, \end{aligned} \tag{1.17}$$

then we have shown the inequality (1.15) holds. Recall that with inequality (1.16), we have

$$\|\hat{\boldsymbol{\theta}}_k^o - \boldsymbol{\theta}_k^*\|_\infty \leq C_1 \|\Sigma_{SS}^{-1}\|_\infty \lambda_1,$$

holds with probability tending to 1. And note that $\hat{\boldsymbol{\theta}}_k^o$ is a solution to (7) with probability tending to 1, that is $\mathbb{P}(\hat{\boldsymbol{\theta}}_k^{(1)} = \hat{\boldsymbol{\theta}}_k^o) \rightarrow 1$. It yields

$$\|\hat{\boldsymbol{\theta}}_k^{(1)} - \boldsymbol{\theta}_k^*\|_\infty \leq C_1 \|\Sigma_{SS}^{-1}\|_\infty \lambda_1,$$

holds with probability tending to 1. If $\theta_{\min}^* \geq C \|\Sigma_{SS}^{-1}\|_\infty \lambda_1$ for some sufficiently large positive constant C , then $\hat{S}^{(1)} = S$ holds with probability tending to 1. In fact, the inequality (1.17) still holds if we choose sufficiently large C . Therefore, we have finished the proof of Theorem 1.2. \square

2 Proof of Auxiliary Lemmas

2.1 Proof or Lemma 1.1

Proof or Lemma 1.1. With probability tending to 1, there exists some sufficiently large positive constant L such that

$$\begin{aligned} \sum_{k=1}^{K-1} \mathbf{x}_k^T \hat{\boldsymbol{\Sigma}}_{(0)} \mathbf{x}_k &\geq \sum_{k=1}^{K-1} \mathbf{x}_k^T \boldsymbol{\Sigma} \mathbf{x}_k - \left| \hat{\boldsymbol{\Sigma}}_{(0)} - \boldsymbol{\Sigma} \right|_\infty \sum_{k=1}^{K-1} \|\mathbf{x}_k\|_1^2 \\ &\geq \sum_{k=1}^{K-1} \mathbf{x}_k^T \boldsymbol{\Sigma} \mathbf{x}_k - \left| \hat{\boldsymbol{\Sigma}}_{(0)} - \boldsymbol{\Sigma} \right|_\infty \left(\sum_{k=1}^{K-1} \|\mathbf{x}_k\|_1 \right)^2 \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}) \sum_{k=1}^{K-1} \|\mathbf{x}_k\|_2^2 - 16s(K-1) \left| \hat{\boldsymbol{\Sigma}}_{(0)} - \boldsymbol{\Sigma} \right|_\infty \sum_{k=1}^{K-1} \|\mathbf{x}_k\|_2^2 \\ &\geq L \sum_{k=1}^{K-1} \|\mathbf{x}_k\|_2^2, \end{aligned}$$

and the last inequality follows from the fact that $|\hat{\boldsymbol{\Sigma}}_{(0)} - \boldsymbol{\Sigma}|_\infty = O_{\mathbb{P}}(\sqrt{\log p/n})$ and $s\sqrt{\log p/n} = o(1)$. \square

2.2 Proof of Lemma 1.2

First note that,

$$\hat{\boldsymbol{\Sigma}}_{(0)}(\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \hat{\mathbf{b}}_{k,0} = (\hat{\boldsymbol{\Sigma}}_{(0)} - \hat{\boldsymbol{\Sigma}})(\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \hat{\boldsymbol{\Sigma}}\boldsymbol{\theta}_k^* + (\hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_1).$$

Due to the definition of $\widehat{\Sigma}$, we have

$$\begin{aligned}\widehat{\Sigma}\boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^* &= \frac{1}{N} \sum_{d=1}^K \sum_{\{i:Y_i=d\}} (\mathbf{X}_i - \widehat{\mu}_d)(\mathbf{X}_i - \widehat{\mu}_d)^T \boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^* \\ &= \frac{1}{N} \sum_{d=1}^K \sum_{\{i:Y_i=d\}} (\mathbf{X}_i - \boldsymbol{\mu}_d)(\mathbf{X}_i - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^* + \frac{1}{N} \sum_{d=1}^K N_d(\widehat{\mu}_d - \boldsymbol{\mu}_d)(\widehat{\mu}_d - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^*.\end{aligned}$$

We note that $\mathbf{X}_i^T \boldsymbol{\theta}_k^* \sim \mathcal{N}(\boldsymbol{\mu}_k^T \boldsymbol{\theta}_k^*, (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1))$ for i such that $Y_i = k$ and $k \neq 1$, which yields that $|(\widehat{\mu}_d - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^*| = O_{\mathbb{P}}(\Delta_{\max}/\sqrt{N})$. Let $\mathbf{D}_{di} = (\mathbf{X}_i - \boldsymbol{\mu}_d)(\mathbf{X}_i - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^*$, then $\mathbf{D}_{di,j}$ is sub-exponential variable with parameter $\sigma_{j,j}\Delta_k$. According to Bernstein's inequality for sub-exponential variable [Vershynin, 2018], we have

$$\left\| \frac{1}{N} \sum_{d=1}^K \sum_{\{i:Y_i=d\}} \mathbf{D}_{di} \right\|_{\infty} = O_{\mathbb{P}} \left(\Delta_{\max} \sqrt{\frac{\log p}{N}} \right).$$

It follows that

$$\begin{aligned}\|\widehat{\Sigma}\boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^*\|_{\infty} &\leq \left\| \frac{1}{N} \sum_{d=1}^K \sum_{\{i:Y_i=d\}} \mathbf{D}_{di} \right\|_{\infty} + \max_d \|(\widehat{\mu}_d - \boldsymbol{\mu}_d)(\widehat{\mu}_d - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^*\|_{\infty} \\ &\leq \left\| \frac{1}{N} \sum_{d=1}^K \sum_{\{i:Y_i=d\}} \mathbf{D}_{di} \right\|_{\infty} + \max_d \|(\widehat{\mu}_d - \boldsymbol{\mu}_d)\|_{\infty} |(\widehat{\mu}_d - \boldsymbol{\mu}_d)^T \boldsymbol{\theta}_k^*| \\ &\lesssim \Delta_{\max} \sqrt{\frac{\log p}{N}} + \sqrt{\frac{\log p}{N}} \frac{\Delta_{\max}}{\sqrt{N}}\end{aligned}$$

with high probability. It yields that

$$\begin{aligned}\|\widehat{\Sigma}_{(0)}(\widehat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \widehat{\boldsymbol{b}}_{k,0}\|_{\infty} &\leq \|\widehat{\Sigma}_{(0)} - \widehat{\Sigma}\|_{\infty} \|\boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k^{(0)}\|_1 + \|\widehat{\Sigma}\boldsymbol{\theta}_k^* - (\widehat{\boldsymbol{\mu}}_k - \widehat{\boldsymbol{\mu}}_1)\|_{\infty} \\ &\leq \|\widehat{\Sigma}_{(0)} - \Sigma\|_{\infty} \|\boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k^{(0)}\|_1 + \|\widehat{\Sigma} - \Sigma\|_{\infty} \|\boldsymbol{\theta}_k^* - \widehat{\boldsymbol{\theta}}_k^{(0)}\|_1 \\ &\quad + \|\widehat{\Sigma}\boldsymbol{\theta}_k^* - \Sigma\boldsymbol{\theta}_k^*\|_{\infty} + \|(\widehat{\boldsymbol{\mu}}_k - \widehat{\boldsymbol{\mu}}_1) - (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)\|_{\infty} \\ &\lesssim \sqrt{\frac{\log p}{n}} a_n + \sqrt{\frac{\log p}{N}} a_n + \Delta_{\max} \sqrt{\frac{\log p}{N}} + \sqrt{\frac{\log p}{N}} \\ &= O_{\mathbb{P}} \left(\sqrt{\frac{\log p}{n}} a_n + \sqrt{\frac{\log p}{N}} \right),\end{aligned}$$

where the third inequality follows from the basic bound $\|\widehat{\boldsymbol{\mu}}_k - \widehat{\boldsymbol{\mu}}\|_{\infty} = O_{\mathbb{P}}(\sqrt{\log p/N})$.

2.3 Proof of Lemma 1.3

Lemma 2.1 (Berry-Esseen inequality [Petrov, 1975]). *Let X_1, \dots, X_n are i.i.d random variables and suppose*

$$\mathbb{E}X_1 = 0, \quad \mathbb{E}X_1^2 = \sigma^2 > 0, \quad \mathbb{E}|X_1|^3 < \infty, \quad \varrho = \frac{\mathbb{E}|X_1|^3}{\sigma^3}.$$

Then for some absolute positive constants A

$$\sup_x \left| \mathbb{P} \left(\frac{1}{\sigma\sqrt{n}} \sum_{j=1}^n X_j < x \right) - \Phi(x) \right| \leq A \frac{\varrho}{\sqrt{n}}.$$

Proof of Lemma 1.3. Denote the Byzantine local machines by \mathcal{B} and $|\mathcal{B}| = \alpha M$. Let $Y_l = \sqrt{n}(\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1)_l / \sqrt{\sigma_{ll}}$ then

$$Y_l = \text{med} \{Y_{l,0}, Y_{l,1}, \dots, Y_{l,M}\}.$$

where $Y_{l,m} = \sqrt{n}(\tilde{\boldsymbol{\mu}}_1^{(m)} - \boldsymbol{\mu}_1)_l / \sqrt{\sigma_{ll}} \sim N(0, 1)$ if $m \notin \mathcal{B}$. Using the uniform bound and the fact that for any $t \in \mathbb{R}$

$$\left| \frac{1}{M+1} \sum_{m=0}^M \mathbb{I}(Y_{l,m} \geq t) - \frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} \mathbb{I}(Y_{l,m} \geq t) \right| \leq \alpha,$$

we have

$$\begin{aligned} & \mathbb{P} \left(\max_l \left| \frac{\sqrt{n}}{\sigma_{ll}} (\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1)_l \right| \geq u_n \right) \\ & \leq p \max_l \mathbb{P} \left(\left| \frac{\sqrt{n}}{\sigma_{ll}} (\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1)_l \right| \geq u_n \right) \\ & = p \max_l \mathbb{P} \left(\frac{1}{M+1} \sum_{m=0}^M \mathbb{I}(Y_{l,m} \geq u_n) \geq \frac{1}{2} \right) + p \max_l \mathbb{P} \left(\frac{1}{M+1} \sum_{m=0}^M \mathbb{I}(Y_{l,m} \leq -u_n) \leq \frac{1}{2} \right) \\ & = p \max_l \mathbb{P} \left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} \mathbb{I}(Y_{l,m} \geq u_n) - \mathbb{P}(Y_{l,m} \geq u_n) \geq \frac{1}{2} - \alpha - (1 - \Phi(u_n)) \right) \\ & \quad + p \max_l \mathbb{P} \left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} \mathbb{I}(Y_{l,m} \leq -u_n) - \mathbb{P}(Y_{l,m} \leq -u_n) \leq \frac{1}{2} + \alpha - \Phi(-u_n) \right). \end{aligned}$$

By Taylor expansion we have

$$\Phi(u_n) = \Phi(0) + \phi(0)u_n + o(u_n).$$

Thus

$$\begin{aligned} & \mathbb{P} \left(\max_l \left| \frac{\sqrt{n}}{\sigma_{ll}} (\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1)_l \right| \geq u_n \right) \\ & \leq p \max_l \mathbb{P} \left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} \mathbb{I}(Y_{l,m} \geq u_n) - \mathbb{P}(Y_{l,m} \geq u_n) \geq \phi(0)u_n + o(u_n) - \alpha \right) \\ & \quad + p \max_l \mathbb{P} \left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} \mathbb{I}(Y_{l,m} \leq -u_n) - \mathbb{P}(Y_{l,m} \leq -u_n) \leq \phi(0)u_n + o(u_n) + \alpha \right). \end{aligned}$$

Let $u_n = \rho'(\sqrt{\log p / ((1 - \alpha)M + 1)} + \alpha)$ for some sufficiently large positive constant ρ' and using Bernstein's inequality we have

$$\mathbb{P} \left(\max_l \frac{\sqrt{n}}{\sigma_{ll}} (\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1)_l \geq u_n \right) \leq 2p^{-1}.$$

And this means that

$$\|\tilde{\boldsymbol{\mu}}_1 - \boldsymbol{\mu}_1\|_\infty \lesssim \sqrt{\frac{\log p}{N}} + \frac{\alpha}{\sqrt{n}},$$

holds with at least probability $1 - 2p^{-1}$. Similarly, we can prove

$$\|\tilde{\boldsymbol{\mu}}_k - \boldsymbol{\mu}_k\|_\infty \lesssim \sqrt{\frac{\log p}{N}} + \frac{\alpha}{\sqrt{n}},$$

holds with at least probability $1 - 2p^{-1}$. For the second inequality, note that

$$(\tilde{\mathbf{d}}_{k,0})_l = \text{med} \left\{ \frac{1}{n} \sum_{d=1}^K \sum_{\{i \in \mathcal{H}_m : Y_i=d\}} (X_{il} - \hat{\mu}_{dl}^{(m)}) (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_d^{(m)})^\top \hat{\boldsymbol{\theta}}_k^{(0)} : m = 1, 2, \dots, M \right\},$$

where X_{il} is the l -th entry of \mathbf{X}_i , $\hat{\mu}_{dl}^{(m)}$ is the l -th entry of $\hat{\boldsymbol{\mu}}_d^{(m)}$. By straightforward calculation we can write

$$\begin{aligned} & \frac{1}{n} \sum_{d=1}^K \sum_{\{i \in \mathcal{H}_m : Y_i=d\}} (X_{il} - \hat{\mu}_{dl}^{(m)}) (\mathbf{X}_i - \hat{\boldsymbol{\mu}}_d^{(m)})^\top \hat{\boldsymbol{\theta}}_k^{(0)} \\ &= \frac{1}{n} \sum_{d=1}^K \sum_{\{i \in \mathcal{H}_m : Y_i=d\}} (X_{il} - \mu_{dl}) (\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} + \frac{1}{n} \sum_{d=1}^K n_k (\hat{\mu}_{dl}^{(m)} - \mu_{dl}) (\hat{\boldsymbol{\mu}}_d^{(m)} - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)}, \end{aligned}$$

and for $m \notin \mathcal{B}$

$$(\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} \sim \mathcal{N} \left(0, (\hat{\boldsymbol{\theta}}_k^{(0)})^\top \Sigma(\hat{\boldsymbol{\theta}}_k^{(0)}) \right), \quad i \in \mathcal{H}_m$$

where $(\hat{\boldsymbol{\theta}}_k^{(0)})^\top \Sigma(\hat{\boldsymbol{\theta}}_k^{(0)}) \lesssim \boldsymbol{\theta}_k^{*T} \Sigma \boldsymbol{\theta}_k^* \leq \Delta_{\max}^2$. Conditioning on $\hat{\boldsymbol{\theta}}_k^{(0)}$, we have

$$\mathbb{E} \left[(X_{il} - \mu_{dl}) (\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} \middle| \hat{\boldsymbol{\theta}}_k^{(0)} \right] = (\Sigma \hat{\boldsymbol{\theta}}_k^{(0)})_l,$$

and

$$\begin{aligned} \tilde{\sigma}_l^2 &:= \text{Var} \left[(X_{il} - \mu_{dl}) (\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} \middle| \hat{\boldsymbol{\theta}}_k^{(0)} \right] \\ &\leq (\mathbb{E}(X_{il} - \mu_{dl})^4)^{1/2} \left(\mathbb{E} \left((\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} \right)^4 \right)^{1/2} \lesssim 3\sigma_{ll}^2 \Delta_k^2, \end{aligned}$$

for $i \in \mathcal{H}_m$ and $m \notin \mathcal{B}$ and $\tilde{\sigma}_l^2 < \infty$ according to assumption. Let

$$W_{l,m} = \frac{1}{\sqrt{n}} \sum_{d=1}^K \sum_{\{i \in \mathcal{H}_m : Y_i=d\}} (X_{il} - \mu_{dl}) (\mathbf{X}_i - \boldsymbol{\mu}_d)^\top \hat{\boldsymbol{\theta}}_k^{(0)} - (\Sigma \hat{\boldsymbol{\theta}}_k^{(0)})_l,$$

and

$$V_{l,m} = \frac{1}{\sqrt{n}} \sum_{d=1}^K n_k (\hat{\mu}_{dl}^{(m)} - \mu_{dl}) (\hat{\boldsymbol{\mu}}_d^{(m)} - \boldsymbol{\mu}_d)^T \hat{\boldsymbol{\theta}}_k^{(0)},$$

then for $m \notin \mathcal{B}$

$$W_{l,m} \xrightarrow{d} N(0, \tilde{\sigma}_l^2) \text{ and } V_{l,m} = O_{\mathbb{P}}\left(\frac{\tilde{\sigma}_l}{\sqrt{n}}\right).$$

Denote

$$Z_{l,m} = \mathbb{I}(W_{l,m} + V_{l,m} \geq u_n) - \mathbb{P}(W_{l,m} + V_{l,m} \geq u_n),$$

and

$$Z_{l,m'} = \mathbb{I}(W_{l,m} + V_{l,m} \leq -u_n) - \mathbb{P}(W_{l,m} + V_{l,m} \leq -u_n).$$

Owing to the definition of sample median, we have

$$\begin{aligned} & \mathbb{P}\left(\max_l \sqrt{n} \left| \left(\tilde{\boldsymbol{d}}_{k,0} - \Sigma \hat{\boldsymbol{\theta}}_k^{(0)}\right)_l \right| \geq u_n\right) \leq p \max_l \mathbb{P}\left(\sqrt{n} \left| \left(\tilde{\boldsymbol{d}}_{k,0} - \Sigma \hat{\boldsymbol{\theta}}_k^{(0)}\right)_l \right| \geq u_n\right) \\ &= p \max_l \mathbb{P}\left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} Z_{l,m} \geq \frac{1}{2} - \alpha - \mathbb{P}(W_{l,m} + V_{l,m} \geq u_n)\right) \\ &+ p \max_l \mathbb{P}\left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} Z_{l,m'} \leq \frac{1}{2} + \alpha - \mathbb{P}(W_{l,m} + V_{l,m} \leq -u_n)\right). \end{aligned}$$

Using the fact $\mathbb{P}(W_{l,m} + V_{l,m} \leq u_n) = \mathbb{P}(W_{l,m}/\tilde{\sigma}_l + V_{l,m}/\tilde{\sigma}_l \leq u_n/\theta_l)$ we have

$$\begin{aligned} & \left| \mathbb{P}\left(\frac{W_{l,m}}{\tilde{\sigma}_l} + \frac{V_{l,m}}{\tilde{\sigma}_l} \leq \frac{u_n}{\tilde{\sigma}_l}\right) - \Phi\left(\frac{u_n}{\tilde{\sigma}_l}\right) \right| \\ & \leq \left| \mathbb{P}\left(\frac{W_{l,m}}{\tilde{\sigma}_l} \leq \frac{u_n}{\tilde{\sigma}_l} - \frac{V_{l,m}}{\tilde{\sigma}_l}\right) - \Phi\left(\frac{u_n}{\tilde{\sigma}_l} - \frac{V_{l,m}}{\tilde{\sigma}_l}\right) \right| + \left| \Phi\left(\frac{u_n}{\tilde{\sigma}_l} - \frac{V_{l,m}}{\tilde{\sigma}_l}\right) - \Phi\left(\frac{u_n}{\tilde{\sigma}_l}\right) \right| \\ & \leq \sup_{x \in \mathbb{R}} \left| \mathbb{P}\left(\frac{W_{l,m}}{\tilde{\sigma}_l} \leq x\right) - \Phi(x) \right| + \left| \Phi\left(\frac{u_n}{\tilde{\sigma}_l} - \frac{V_{l,m}}{\tilde{\sigma}_l}\right) - \Phi\left(\frac{u_n}{\tilde{\sigma}_l}\right) \right| \\ & = \sup_{x \in \mathbb{R}} \left| \mathbb{P}\left(\frac{W_{l,m}}{\tilde{\sigma}_l} \leq x\right) - \Phi(x) \right| + \phi\left(\frac{u_n}{\tilde{\sigma}_l}\right) \frac{V_{l,m}}{\tilde{\sigma}_l} + o(V_{l,m}) \\ & \lesssim \frac{1}{\sqrt{n}} + \frac{1}{\sqrt{n}}, \end{aligned}$$

where the last inequality follows from Berry-Esseen inequality and the normal density function $\phi(x)$ is bounded. It yields that

$$\begin{aligned} & \mathbb{P}\left(\max_l \sqrt{n} \left| \left(\tilde{\boldsymbol{d}}_{k,0} - \Sigma \hat{\boldsymbol{\theta}}_k^{(0)}\right)_l \right| \geq u_n\right) \\ & \leq p \max_l \mathbb{P}\left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} Z_{l,m} \geq \phi(0)u_n/\tilde{\sigma}_l - \alpha + O\left(\frac{1}{\sqrt{n}}\right) + o(u_n)\right) \\ & + p \max_l \mathbb{P}\left(\frac{1}{(1-\alpha)M+1} \sum_{m \notin \mathcal{B}} Z_{l,m'} \leq \phi(0)u_n/\tilde{\sigma}_l + \alpha + O\left(\frac{1}{\sqrt{n}}\right) + o(u_n)\right). \end{aligned}$$

Let $u_n = \rho''(\Delta_{\max}\sqrt{\log p/(M+1)} + 1/\sqrt{n} + \alpha)$ for some sufficiently large positive constant ρ'' , then by Bernstein's inequality we can prove that

$$\|\tilde{\mathbf{d}}_{k,0} - \Sigma\hat{\boldsymbol{\theta}}_k^{(0)}\|_\infty = O_{\mathbb{P}}\left(\Delta_{\max}\sqrt{\frac{\log p}{N}} + \frac{\alpha}{\sqrt{n}} + \frac{1}{n}\right).$$

By the definition of $\tilde{\mathbf{b}}_{k,0}$, we have

$$\begin{aligned} & \|\hat{\Sigma}_{(0)}(\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \tilde{\mathbf{b}}_{k,0}\|_\infty \\ &= \left\| (\hat{\Sigma}_{(0)} - \Sigma) (\hat{\boldsymbol{\theta}}_k^{(0)} - \boldsymbol{\theta}_k^*) - \tilde{\mathbf{d}}_{k,0} + \Sigma\hat{\boldsymbol{\theta}}_k^{(0)} - \Sigma\boldsymbol{\theta}_k^* + (\tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\mu}}_1) \right\|_\infty \\ &\leq \left\| (\hat{\Sigma}_{(0)} - \Sigma) (\boldsymbol{\theta}_k^* - \hat{\boldsymbol{\theta}}_k^{(0)}) \right\|_\infty + \left\| \tilde{\mathbf{d}}_{k,0} - \Sigma\hat{\boldsymbol{\theta}}_k^{(0)} \right\|_\infty + \|\boldsymbol{\mu}_k - \boldsymbol{\mu}_1 - (\tilde{\boldsymbol{\mu}}_k - \tilde{\boldsymbol{\mu}}_1)\|_\infty, \end{aligned}$$

then the results follow. \square

2.4 Proof of Lemma 1.4

Proof of Lemma 1.4. By the definition of the support set S and $\Sigma\boldsymbol{\theta}_k^* = \boldsymbol{\mu}_k - \boldsymbol{\mu}_1$ we have

$$\begin{pmatrix} (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)_S \\ (\boldsymbol{\mu}_k - \boldsymbol{\mu}_1)_{S^c} \end{pmatrix} = \begin{pmatrix} \Sigma_{SS} & \Sigma_{SS^c} \\ \Sigma_{S^c S} & \Sigma_{S^c S^c} \end{pmatrix} \begin{pmatrix} \boldsymbol{\theta}_{k,S}^* \\ \mathbf{0} \end{pmatrix},$$

then the results follow immediately. \square

References

- V. V. Petrov. *Sums of Independent Random Variables*, volume 82. Springer, New York, 1975.
- R. Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- M. J. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using ℓ_1 -constrained quadratic programming (lasso). *IEEE Transactions on Information Theory*, 55(5):2183–2202, 2009.