

## A. Proofs of the Main Results

In this section, we lay out the proofs of the two theorems in §4, which establish the statistical rates of convergence of our estimators.

### A.1. Proof of Theorem 4.2

*Proof.* Since  $\widehat{\beta}$  is the solution of the optimization problem in (3.4), the first-order optimality condition states that

$$\nabla L(\widehat{\beta}) + \lambda \xi = 0, \quad \text{where } \xi \in \partial \|\widehat{\beta}\|_1. \quad (\text{A.1})$$

Then the entries of  $\xi \in \mathbb{R}^d$  are given by

$$\xi_j = \text{sign}(\widehat{\beta}_j), \quad \forall j \in \text{supp}(\widehat{\beta}); \quad \xi_j \in [-1, 1], \quad \forall j \notin \text{supp}(\widehat{\beta}).$$

For any index set  $\mathcal{A} \subseteq [d]$  and  $z \in \mathbb{R}^d$ , we define the restriction of  $z$  to  $\mathcal{A}$ ,  $z_{\mathcal{A}} \in \mathbb{R}^d$ , by letting

$$[z_{\mathcal{A}}]_j = z_j \quad \text{if } j \in \mathcal{A}, \quad [z_{\mathcal{A}}]_j = 0 \quad \text{otherwise.}$$

Here  $[z_{\mathcal{A}}]_j$  is the  $j$ -th entry of  $z_{\mathcal{A}}$ . Let  $\mathcal{S} = \text{supp}(\beta^*)$ , then we can write  $\xi = \xi_{\mathcal{S}} + \xi_{\mathcal{S}^c}$ . For notational simplicity, in the sequel, we define  $\theta = \widehat{\beta} - \mu \cdot \beta^*$ . Thus by (A.1) it holds that

$$\begin{aligned} \langle \nabla L(\widehat{\beta}) - \nabla L(\mu\beta^*), \theta \rangle &= \langle -\lambda \cdot \xi - \nabla L(\mu\beta^*), \theta \rangle \\ &\leq \langle -\lambda \cdot \xi_{\mathcal{S}} - \lambda \cdot \xi_{\mathcal{S}^c}, \theta \rangle + \|\nabla L(\mu\beta^*)\|_{\infty} \cdot \|\theta\|_1. \end{aligned} \quad (\text{A.2})$$

By the definition of  $\xi$ , we have

$$\langle -\lambda \cdot \xi_{\mathcal{S}^c}, \widehat{\beta} - \mu\beta^* \rangle = -\lambda \cdot \|\widehat{\beta}\|_1. \quad (\text{A.3})$$

Moreover, since  $\|\xi\|_{\infty} \leq 1$ , Hölder's inequality implies that

$$\langle -\lambda \cdot \xi_{\mathcal{S}}, \theta \rangle \leq \|\theta_{\mathcal{S}}\|_1. \quad (\text{A.4})$$

Note that  $\nabla^2 L(\beta) = 2I_d$ . Combining (A.9), (A.3), and (A.4), we obtain

$$2\|\theta\|_2^2 = \langle \nabla L(\widehat{\beta}) - \nabla L(\mu\beta^*), \theta \rangle \leq -\lambda\|\theta_{\mathcal{S}^c}\|_1 + \lambda\|\theta_{\mathcal{S}}\|_1 + \|\nabla L(\mu\beta^*)\|_{\infty} \cdot \|\theta\|_1. \quad (\text{A.5})$$

For an upper bound of the right-hand side of (A.5), we apply the following lemma to obtain an upper bound on  $\|\nabla L(\mu\beta^*)\|_{\infty}$ .

**Lemma A.1** (Bound on  $\|\nabla L(\mu\beta^*)\|_{\infty}$ ). We set the truncation level in (4.1) as  $\tau = 2(M \cdot n / \log d)^{1/4}$ . Then we have

$$\mathbb{P}\left[\|\nabla L(\mu\beta^*)\|_{\infty} > 7\sqrt{M \cdot \log d/n}\right] \leq d^{-2}.$$

*Proof.* See §B.1 for a detailed proof. □

Thus by Lemma A.1 and the choice of  $\lambda$ , we have  $\lambda > 2\|\nabla L(\mu\beta^*)\|_{\infty}$  with probability at least  $1 - d^{-2}$ . This implies that

$$2\|\theta\|_2^2 \leq -\lambda/2 \cdot \|\theta_{\mathcal{S}^c}\|_1 + 3\lambda/2 \cdot \|\theta_{\mathcal{S}}\|_1 \leq 2\lambda \cdot \|\theta_{\mathcal{S}}\|_1. \quad (\text{A.6})$$

Since the leftmost term in (A.6) is nonnegative, we obtain  $\|\theta_{\mathcal{S}^c}\|_1 \leq 3 \cdot \|\theta_{\mathcal{S}}\|_1$ . In addition, since  $|\mathcal{S}| = s^*$ ,  $\|\theta_{\mathcal{S}}\|_1 \leq \sqrt{s^*} \cdot \|\theta_{\mathcal{S}}\|_2$ . Thus by (A.6) we have  $\|\theta\|_2 \leq \sqrt{s^*} \cdot \lambda$ . Moreover, we also have  $\|\theta_{\mathcal{S}}\|_1 \leq s^* \lambda$ , which further implies that

$$\|\theta\|_1 = \|\theta_{\mathcal{S}}\|_1 + \|\theta_{\mathcal{S}^c}\|_1 \leq 4 \cdot \|\theta_{\mathcal{S}}\|_1 \leq 4s^* \lambda.$$

Therefore, we conclude the proof. □

### A.2. Proof of Theorem 4.3

*Proof.* The proof of Theorem 4.3 is parallel to that of Theorem 4.2. Here the difference is to handle the nuclear norm regularization, instead of the  $\ell_1$ -penalty. Since  $\hat{\beta}$  is the solution of the optimization problem in (3.4), the first order optimality condition states that

$$L(\hat{\beta}) + \lambda \|\hat{\beta}\|_* \leq L(\mu\beta^*) + \lambda \|\mu\beta^*\|_*. \quad (\text{A.7})$$

To simplify the notation, we define  $\Theta = \hat{\beta} - \mu \cdot \beta^*$ . Since  $L$  is quadratic,

$$L(\hat{\beta}) - L(\mu\beta^*) = \langle \nabla L(\mu\beta^*), \Theta \rangle + 2\|\Theta\|_{\text{fro}}^2, \quad (\text{A.8})$$

where  $\nabla L$  takes values in  $\mathbb{R}^{d_1 \times d_2}$ . Then combining (A.7), (A.8), and Hölder's inequality, we have

$$\|\Theta\|_{\text{fro}}^2 \leq -\langle \nabla L(\mu\beta^*), \Theta \rangle + \lambda \|\mu\beta^*\|_* - \lambda \|\hat{\beta}\|_* \leq \|\nabla L(\mu\beta^*)\|_{\text{op}} \cdot \|\Theta\|_* + \lambda \|\mu\beta^*\|_* - \lambda \|\hat{\beta}\|_*. \quad (\text{A.9})$$

In the following, we focus on the term  $\|\mu\beta^*\|_* - \|\hat{\beta}\|_*$  in (A.9). Let  $U\Lambda^*V^\top$  be the singular value decomposition of  $\mu\beta^*$ , where  $U \in \mathbb{R}^{d_1 \times d_1}$  and  $V \in \mathbb{R}^{d_2 \times d_2}$  are orthogonal matrices, and  $\Lambda^* \in \mathbb{R}^{d_1 \times d_2}$  be formed by the singular values of  $\mu\beta^*$ . Moreover, since  $\text{rank}(\beta^*) = r^*$ ,  $\Lambda^*$  can be written in block form as

$$\Lambda^* = \begin{bmatrix} \Lambda_{11}^* & 0 \\ 0 & 0 \end{bmatrix}, \quad (\text{A.10})$$

where  $\Lambda_{11}^* \in \mathbb{R}^{r^* \times r^*}$  is a diagonal matrix whose diagonal elements are the nonzero singular values of  $\mu\beta^*$ . We define  $\Gamma = U^\top \Theta V$ , which can be written in block form as

$$\Gamma = \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & \Gamma_{22} \end{bmatrix},$$

where  $\Gamma_{11} \in \mathbb{R}^{r^* \times r^*}$ . In addition, we define matrices

$$\Gamma^{(1)} = \begin{bmatrix} 0 & 0 \\ 0 & \Gamma_{22} \end{bmatrix} \quad \text{and} \quad \Gamma^{(2)} = \begin{bmatrix} \Gamma_{11} & \Gamma_{12} \\ \Gamma_{21} & 0 \end{bmatrix}.$$

Then by (A.10) and triangle inequality of the nuclear norm, we have

$$\begin{aligned} \|\hat{\beta}\|_* &= \|\mu\beta^* + \Theta\|_* = \|U(\Lambda^* + \Gamma)V^\top\|_* = \|\Lambda^* + \Gamma\|_* \\ &\geq \|\Lambda^* + \Gamma^{(1)}\|_* - \|\Gamma^{(2)}\|_* = \|\Lambda^*\|_* + \|\Gamma^{(1)}\|_* - \|\Gamma^{(2)}\|_*, \end{aligned} \quad (\text{A.11})$$

where the last equality follows from the fact that  $\Lambda^* + \Gamma^{(1)}$  is block diagonal. Since  $\|\mu\beta^*\|_* = \|\Lambda^*\|_*$ , by (A.11) we obtain

$$\|\mu\beta^*\|_* - \|\hat{\beta}\|_* \leq \|\Gamma^{(2)}\|_* - \|\Gamma^{(1)}\|_*. \quad (\text{A.12})$$

In addition, triangle inequality implies that

$$\|\Theta\|_* = \|U\Gamma V^\top\|_* \leq \|\Gamma^{(1)}\|_* + \|\Gamma^{(2)}\|_*. \quad (\text{A.13})$$

Thus combining (A.11), (A.12), (A.13), we have

$$\|\Theta\|_{\text{fro}}^2 \leq (\|\nabla L(\mu\beta^*)\|_{\text{op}} + \lambda) \cdot \|\Gamma^{(2)}\|_* + (\|\nabla L(\mu\beta^*)\|_{\text{op}} - \lambda) \cdot \|\Gamma^{(1)}\|_*, \quad (\text{A.14})$$

We utilize the following lemma to obtain an upper bound of  $\|\nabla L(\mu\beta^*)\|_{\text{op}}$ .

**Lemma A.2** (Upper bound of  $\|\nabla L(\mu\beta^*)\|_{\text{op}}$ ). Let loss function  $L: \mathbb{R}^{d_1 \times d_2} \rightarrow \mathbb{R}$  be defined in (3.4) for the matrix setting. Setting  $\kappa = 2\sqrt{n \cdot \log(d_1 + d_2)} / \sqrt{(d_1 + d_2)M}$ , then it holds that

$$\mathbb{P}\left[\|\nabla L(\mu\beta^*)\|_{\text{op}} > 6\sqrt{(d_1 + d_2)/n}\right] \leq (d_1 + d_2)^{-2}.$$

*Proof.* See B.2 for a detailed proof.  $\square$

By Lemma A.2 and the choice of  $\lambda$ , we conclude that  $\lambda > 2 \cdot \|\nabla L(\mu\beta^*)\|_{\text{op}}$  with probability at least  $1 - (d_1 + d_2)^{-3}$ . Thus by (A.14) we have

$$\|\Theta\|_{\text{fro}}^2 \leq 3\lambda/2 \cdot \|\Gamma^{(2)}\|_* - \lambda/2 \cdot \|\Gamma^{(1)}\|_* \quad (\text{A.15})$$

which implies that  $\|\Gamma^{(1)}\|_* \leq 3 \cdot \|\Gamma^{(2)}\|_*$ . Moreover, by the subadditivity of rank, we obtain

$$\text{rank}(\Gamma^{(2)}) \leq \text{rank}\left(\begin{bmatrix} \Gamma_{11}/2 & \Gamma_{12} \\ 0 & 0 \end{bmatrix}\right) + \text{rank}\left(\begin{bmatrix} \Gamma_{11}/2 & 0 \\ \Gamma_{21} & 0 \end{bmatrix}\right) = 2r^*,$$

which implies that  $\|\Gamma^{(2)}\|_* \leq \sqrt{2r^*} \cdot \|\Gamma^{(2)}\|_{\text{fro}}$ . Then by (A.15) we obtain that  $\|\Theta\|_{\text{fro}} \leq 3/\sqrt{2} \cdot \sqrt{r^*} \cdot \lambda$ . Finally, by triangle inequality for the nuclear norm,

$$\|\Theta\|_* = \|\Gamma\|_* \leq \|\Gamma^{(1)}\|_* + \|\Gamma^{(2)}\|_* \leq 4 \cdot \|\Gamma^{(2)}\|_* \leq 4\sqrt{2r^*} \|\Gamma^{(2)}\|_{\text{fro}} = 12r^* \lambda.$$

Thus we conclude the proof of Theorem 4.3.  $\square$

## B. Proof of Auxiliary Results

### B.1. Proof of Lemma A.1

*Proof.* By definition of the loss function  $L$  in (3.4), we have

$$\nabla L(\mu\beta^*) = 2\mu\beta^* - \frac{2}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}(X_i) = \mathbb{E}[2Y_i \cdot S(X_i)] - \frac{2}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}(X_i).$$

By triangle inequality,

$$\|\nabla L(\mu\beta^*)\|_{\infty} \leq \left\| \mathbb{E}[2Y \cdot S(X)] - \mathbb{E}[2\tilde{Y} \cdot \tilde{S}(X)] \right\|_{\infty} + \left\| \mathbb{E}[2\tilde{Y} \cdot \tilde{S}(X)] - \frac{2}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}(X_i) \right\|_{\infty}. \quad (\text{B.1})$$

For any  $j \in [d]$ , by the definition of the truncated response  $\tilde{Y}$  and truncated score  $\tilde{S}$ , we obtain

$$\begin{aligned} |\mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] - \mathbb{E}[Y \cdot S_j(X)]| &\leq \left| \mathbb{E}\left\{ \tilde{Y} \cdot [\tilde{S}_j(X) - S_j(X)] \right\} \right| + |\mathbb{E}[(\tilde{Y} - Y) \cdot S_j(X)]| \\ &= \underbrace{|\mathbb{E}[\tilde{Y} \cdot S_j(X) \cdot \mathbb{1}_{\{|S_j(X)| > \tau\}}]|}_{a_1} + \underbrace{|\mathbb{E}[Y \cdot S_j(X) \cdot \mathbb{1}_{\{|Y| > \tau\}}]|}_{a_2}. \end{aligned} \quad (\text{B.2})$$

By Cauchy-Schwarz inequality, we have

$$\begin{aligned} a_1^2 &\leq \mathbb{E}[\tilde{Y}^2 S_j^2(X)] \cdot \mathbb{P}[|S_j(X)| \geq \tau] \\ &\leq \sqrt{\mathbb{E}(\tilde{Y}^4) \cdot \mathbb{E}[S_j^4(X)]} \cdot \mathbb{E}[S_j^4(X)] \cdot \tau^{-4} \\ &= M^2 \cdot \tau^{-4}, \end{aligned} \quad (\text{B.3})$$

where the second inequality follows from Chebyshev's inequality. Similarly, for  $a_2$  we have

$$\begin{aligned} a_2^2 &\leq \mathbb{E}[Y^2 S_j^2(X)] \cdot \mathbb{P}(|Y| \geq \tau) \\ &\leq \sqrt{\mathbb{E}(\tilde{Y}^4) \cdot \mathbb{E}[S_j^4(X)]} \cdot \mathbb{E}(Y^4) \cdot \tau^{-4} \\ &\leq M^2 \cdot \tau^{-4}. \end{aligned} \quad (\text{B.4})$$

Thus combining (B.2), (B.3), and (B.4), we conclude that

$$\left| \mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] - \mathbb{E}[Y \cdot S_j(X)] \right| \leq a_1 + a_2 \leq 2M \cdot \tau^{-2}$$

for all  $j \in [d]$ . Thus choosing  $\tau = 2(M \cdot n / \log d)^{1/4}$ , we have

$$\left\| \mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] - \mathbb{E}[Y \cdot S_j(X)] \right\|_\infty \leq 1/2 \cdot \sqrt{M \cdot \log d / n}. \quad (\text{B.5})$$

Furthermore, under Assumption 4.1, the variance of  $\tilde{Y} \cdot \tilde{S}_j(X)$  is bounded by

$$\text{Var}[\tilde{Y} \cdot \tilde{S}_j(X)] \leq \mathbb{E}[\tilde{Y}^2 \cdot \tilde{S}_j^2(X)] \leq \mathbb{E}[Y^2 \cdot S_j^2(X)] \leq \sqrt{\mathbb{E}(Y^4) \cdot \mathbb{E}[S_j^4(X)]} \leq M.$$

Thus for the second term in (B.1), since  $|\tilde{Y} \cdot \tilde{S}_j(X)| \leq \tau^2$ , by the Bernstein inequality in (Boucheron et al., 2013) (Theorem 2.10), for any  $j \in [d]$  and any  $t > 0$ , we have

$$\mathbb{P}\left\{ \left| \frac{1}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}_j(X_i) - \mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] \right| \geq \sqrt{\frac{2M \cdot t}{n}} + \frac{\tau^2 \cdot t}{3n} \right\} \leq \exp(-t). \quad (\text{B.6})$$

Taking union bound over  $j \in [t]$  in (B.6) yields

$$\mathbb{P}\left\{ \left\| \frac{1}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}_j(X_i) - \mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] \right\|_\infty \geq \sqrt{\frac{2M \cdot t}{n}} + \frac{\tau^2 \cdot t}{3n} \right\} \leq \exp(-t + \log d). \quad (\text{B.7})$$

Finally, we plug in  $\tau = 2(M \cdot n / \log d)^{1/4}$  and set  $t = 3 \log d$  in (B.7) to obtain that

$$\left\| \frac{1}{n} \sum_{i=1}^n \tilde{Y}_i \cdot \tilde{S}_j(X_i) - \mathbb{E}[\tilde{Y} \cdot \tilde{S}_j(X)] \right\|_\infty \leq (4 + \sqrt{6}) \sqrt{\frac{M \cdot \log d}{n}} \quad (\text{B.8})$$

with probability at least  $1 - d^{-2}$ . Finally, combining (B.1), (B.5), and (B.8), we conclude the proof.  $\square$

## B.2. Proof of Lemma A.2

*Proof.* For loss function  $L$  defined in (3.4) in the matrix setting, we have

$$\nabla L(\mu\beta^*) = 2\mu\beta^* - \frac{2}{\kappa \cdot n} \sum_{i=1}^n \psi[\kappa \cdot Y_i \cdot S(X_i)] = 2\mathbb{E}[Y \cdot S(X)] - \frac{2}{\kappa \cdot n} \sum_{i=1}^n \psi[\kappa \cdot Y_i \cdot S(X_i)]. \quad (\text{B.9})$$

Here the last equality follows from the generalized Stein's identity. In the sequel, we apply results in (Minsker, 2016) to bound  $\|\nabla L(\mu\beta^*)\|_{\text{op}}$ . To begin with, we first consider the operator norm of  $\mathbb{E}[Y^2 \cdot S(X)S(X)^\top] \in \mathbb{R}^{d_1 \times d_2}$  and  $\mathbb{E}[Y^2 \cdot S(X)^\top S(X)] \in \mathbb{R}^{d_2 \times d_1}$ . For notational simplicity, we denote by  $S_{j,\cdot}(\cdot) \in \mathbb{R}^{d_2}$ ,  $S_{\cdot,k}(\cdot) \in \mathbb{R}^{d_1}$  the  $j$ -th row and  $k$ -th column of the score function  $S(\cdot)$ , respectively. For any  $u \in \mathbb{R}^{d_1-1}$ , by Cauchy-Schwarz inequality we have

$$\mathbb{E}[Y^2 \cdot u^\top S(X)S(X)^\top u] = \sum_{k=1}^{d_2} \mathbb{E}\{[Y^2 \cdot S_{\cdot,k}(X)^\top u]^2\} \leq d_2 \cdot \sqrt{\mathbb{E}(Y^4) \cdot \mathbb{E}\{[S_{\cdot,1}(X)^\top u]^4\}}, \quad (\text{B.10})$$

where we use the fact that the entries of  $S(X)$  are i.i.d. Since  $\mathbb{E}[S_{ij}(X)] = 0$  and  $\mathbb{E}[S_{ij}^4(X)] \leq M$ , by Cauchy-Schwarz inequality we obtain that

$$\begin{aligned} \mathbb{E}\{[S_{\cdot,1}(X)^\top u]^4\} &= \sum_{j_1=1}^d \sum_{j_2=1}^d \mathbb{E}[S_{j_1,1}(X)^2 \cdot S_{j_2,1}(X)^2] \cdot u_{j_1}^2 u_{j_2}^2 \\ &\leq \sum_{j_1=1}^d \sum_{j_2=1}^d \sqrt{\mathbb{E}[S_{j_1,1}^4(X)] \cdot \mathbb{E}[S_{j_2,1}^4(X)]} \cdot u_{j_1}^2 u_{j_2}^2 \leq M \sum_{j_1=1}^d \sum_{j_2=1}^d u_{j_1}^2 u_{j_2}^2 = M. \end{aligned} \quad (\text{B.11})$$

Thus combining (B.10) and (B.11) we obtain that

$$\mathbb{E}[Y^2 \cdot u^\top S(X) S(X)^\top u] \leq d_2 \cdot M,$$

which implies that  $\|\mathbb{E}[Y^2 \cdot S(X) S(X)^\top]\|_{\text{op}} \leq d_2 \cdot M$ . Similarly, we obtain  $\|\mathbb{E}[Y^2 \cdot S(X)^\top S(X)]\|_{\text{op}} \leq d_1 \cdot M$ . Thus by Corollary 3.1 in (Minsker, 2016), we have

$$\mathbb{P}\left\{\left\|\frac{1}{\kappa \cdot n} \sum_{i=1}^n \psi[\kappa \cdot Y_i \cdot S(X_i)] - \mathbb{E}[Y \cdot S(X)]\right\|_{\text{op}} \geq \frac{t}{\sqrt{n}}\right\} \leq 2(d_1 + d_2) \exp[-\kappa t \sqrt{n} + \kappa^2 (d_1 + d_2) M / 2] \quad (\text{B.12})$$

for any  $t > 0$  and  $\kappa > 0$ . We set  $\kappa = 2\sqrt{n \cdot \log(d_1 + d_2)} / \sqrt{(d_1 + d_2)M}$  and  $t = \sqrt{(d_1 + d_2)M} \cdot s$  in (B.12), which implies that

$$\begin{aligned} \mathbb{P}\left\{\left\|\frac{1}{\kappa \cdot n} \sum_{i=1}^n \psi[\kappa \cdot Y_i \cdot S(X_i)] - \mathbb{E}[Y \cdot S(X)]\right\|_{\text{op}} \geq \sqrt{\frac{(d_1 + d_2)M}{n}} \cdot s\right\} \\ \leq 2(d_1 + d_2) \exp[-2\sqrt{\log(d_1 + d_2)} \cdot s + 2 \cdot \log(d_1 + d_2)]. \end{aligned} \quad (\text{B.13})$$

Now we set  $s = 3 \cdot \sqrt{\log(d_1 + d_2)}$ , which implies that the right-hand side of (B.13) is less than

$$2(d_1 + d_2) \exp[-6 \log(d_1 + d_2) + 2 \cdot \log(d_1 + d_2)] \leq (d_1 + d_2)^2 \cdot \exp[-4 \cdot \log(d_1 + d_2)] = (d_1 + d_2)^{-2}.$$

Therefore, combining (B.9) and (B.13) we conclude that

$$\|\nabla L(\mu\beta^*)\|_{\text{op}} \leq 6\sqrt{(d_1 + d_2) \cdot M/n} \cdot \sqrt{\log(d_1 + d_2)}$$

with probability at least  $1 - (d_1 + d_2)^{-2}$ , which concludes the proof. □