

## A Additional Experiments

In this section, we present additional experimental results to verify the linear convergence rate, sample complexity, and statistical rate of our proposed algorithm.

### A.1 Robust Matrix Sensing

Our data are generated from the same procedure as described before. In addition, we study the same experimental setting as before except we choose  $\alpha = r, \nu = 1, \beta = 0.1$ . Figure 2 summarized the experimental results for robust matrix sensing. Figure 2(a) and 2(c) illustrate the relative error  $\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2 / \|\mathbf{X}^*\|_F^2$  in log scale versus number of iterations. Note that, we only lay out results under setting  $d_1 = d_2 = 100, r = 3$  with number of observations  $n = 0.2 * d_1 d_2$  to avoid redundancy. These plots demonstrate the linear rate of convergence of our algorithm. Figure 2(b) demonstrates the sample complexity requirement to achieve exact recovery for low-rank structure in the noiseless setting. Note that we say  $\widehat{\mathbf{X}}$  achieves exact recovery if  $\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F / \|\mathbf{X}^*\|_F \leq 10^{-3}$ . It confirms our theoretical results regarding the sample complexity. The statistical error for the low-rank matrix is demonstrated in Figure 2(d), which is consistent with our result  $O(rd/n)$ .

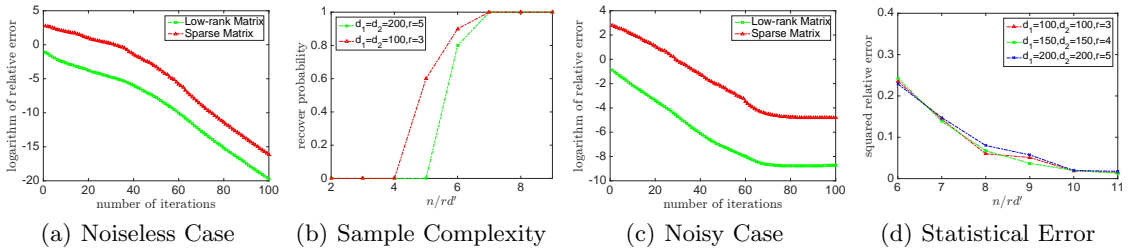


Figure 2: Experimental results for robust matrix sensing. (a),(c) Relative error in log scale vs. number of iterations in the noiseless and noisy settings respectively. (b) Recovering probability of low-rank matrix vs. scaled sample size in the noiseless setting. (d) Relative error vs. scaled sample size in the noisy setting.

### A.2 Robust PCA

We generate the data according to the same procedure as before. Furthermore, we consider the same experimental settings as robust matrix sensing except. The experimental results for robust PCA are summarized in Figure 3. In detail, Figures 3(a) and 3(c) report the squared estimation error  $\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F^2 / (d_1 d_2)$  in log scale versus number of iterations. Note that we only lay out the results under fully observed model with setting  $d_1 = d_2 = 200, r = 5$ , because other settings will give us similar plots, and we leave them out for simplicity. The results verify the linear convergence rate of our algorithm. In the noiseless setting, the sample complexity for achieving exactly recovery of the low-rank matrix is illustrated in Figure 3(b). The result of recovery probability indicates the sample complexity requirement  $n = O(rd \log d)$  for robust PCA. Finally, Figure 3(d) demonstrates the statistical error for the low-rank matrix, which is at the order  $O(rd \log d/n)$ . Although our theoretical results suggest  $O(r^2 d \log d)$  sample complexity and  $O(r^2 d \log d)$  statistical error, the simulation results indicate that both the sample complexity and the statistical error scale linearly with  $rd$ .

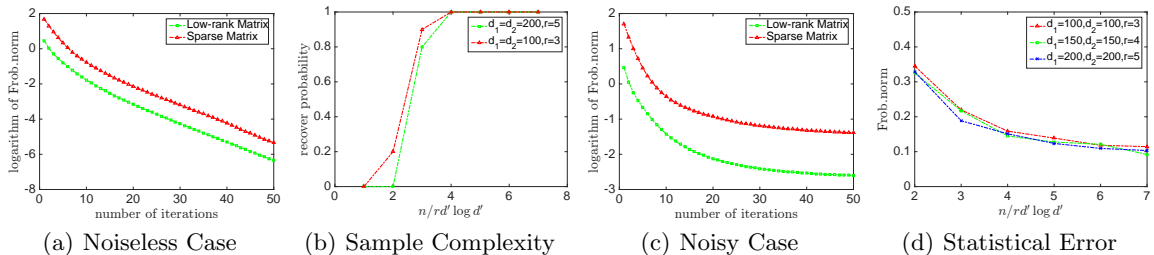


Figure 3: Experimental results for robust PCA. (a),(c) Squared estimation error in log scale vs. number of iterations in the noiseless and noisy settings respectively. (b) Recovering probability of low-rank matrix vs. scaled sample size in the noiseless setting. (d) Squared relative error vs. scaled sample size in the noisy setting.

## B Proof of the Main Theory

In this section, we establish the proof of our main theory. Before proceeding any further, we introduce the following notations. For any index set  $\Omega \subseteq [d_1] \times [d_2]$ , let  $\Omega_{i,*}$  and  $\Omega_{*,j}$  be the  $i$ -th row and  $j$ -th column of  $\Omega$  respectively. Denote the column and row space of  $\mathbf{A}$  by  $\text{col}(\mathbf{A})$  and  $\text{row}(\mathbf{A})$  respectively. Let the top  $d_1 \times r$  and bottom  $d_2 \times r$  matrices of any matrix  $\mathbf{A} \in \mathbb{R}^{(d_1+d_2) \times r}$  be  $\mathbf{A}_U$  and  $\mathbf{A}_V$  respectively. Let the nuclear norm of any matrix  $\mathbf{A}$  be  $\|\mathbf{A}\|_*$ . Denote  $\mathbf{Z} = [\mathbf{U}; \mathbf{V}] \in \mathbb{R}^{(d_1+d_2) \times r}$ , then according to (3.3), we reformulate the regularized objective function as follows

$$\tilde{F}_n(\mathbf{Z}, \mathbf{S}) = F_n(\mathbf{U}, \mathbf{V}, \mathbf{S}) = \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S}) + \frac{1}{8} \|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2. \quad (\text{B.1})$$

Therefore, the corresponding gradient regarding to  $\mathbf{Z}$  is as follows

$$\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S}) = \begin{bmatrix} \nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S}) + \frac{1}{2} \mathbf{U}(\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}) \\ \nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S}) + \frac{1}{2} \mathbf{V}(\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}) \end{bmatrix}. \quad (\text{B.2})$$

### B.1 Proof of Theorem 4.6

In order to prove Theorem 4.6, we need to make use of the following lemmas. Since both low-rank and sparse structures exist in our model, it is necessary to derive the convergence results for both structures. Lemma B.1, proved in Section C.1 characterizes the convergence of the low rank structure, while Lemma B.2, proved in Section C.2 corresponds to the convergence of the sparse structure.

**Lemma B.1** (Convergence for Low-Rank Structure). Suppose the sample loss function  $\mathcal{L}_n$  satisfies Conditions 4.2 and 4.4. Recall that  $\mathbf{X}^* = \mathbf{U}^* \mathbf{V}^{*\top}$  is the unknown rank- $r$  matrix that satisfies (3.1),  $\mathbf{S}^*$  is the unknown  $s$ -sparse matrix with at most  $\beta$ -fraction nonzero entries per row and column. There exist constants  $c_1, c_2$  and  $c_3$  such that if  $\mathbf{Z}^t \in \mathbb{B}(c_2 \sqrt{\sigma_r})$  with  $c_2 \leq \min\{1/4, \sqrt{\mu'_1/[10(L_1 + 1 + 8/\mu_2)]}\}$ , and we set the step size  $\eta = c_1/\sigma_1$  with  $c_1 \leq \min\{1/32, \mu_1/(192L_1^2)\}$ , then the output of Algorithm 1  $\mathbf{Z}^t = [\mathbf{U}^t; \mathbf{V}^t]$  satisfies

$$d^2(\mathbf{Z}^{t+1}, \mathbf{Z}^*) \leq \rho_1 d^2(\mathbf{Z}^t, \mathbf{Z}^*) - \frac{\eta \mu_1}{4} \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + \Gamma_1 \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + \Gamma_2 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2,$$

provided that  $\beta \leq 1/(c_3 \alpha r \kappa)$  with  $c_3 \geq 720(\gamma + 1)\mu_2/\mu'_1$ , where contraction parameter  $\rho_1 = 1 - \eta \mu'_1 \sigma_r / 40$ ,  $\mu'_1 = \min\{\mu_1, 2\}$ ,  $\Gamma_1 = 48\eta^2(1 + K)^2 \sigma_1 + \eta(\mu_2 + 4K^2/\mu_1)$ , and  $\Gamma_2 = 48\eta^2 r \sigma_1 + 2\eta(8r/\mu_1 + r/L_1)$ .

**Lemma B.2** (Convergence for Sparse Structure). Suppose the sample loss function  $\mathcal{L}_n$  satisfies Conditions 4.3 and 4.4. Recall that  $\mathbf{X}^*$  is the unknown rank- $r$  matrix,  $\mathbf{S}^*$  is the unknown  $s$ -sparse matrix. If we set the step size  $\tau \leq 1/(3L_2)$  and choose appropriate parameters  $\gamma, \gamma'$ , then the output of Algorithm 1 satisfies

$$\|\mathbf{S}^{t+1} - \mathbf{S}^*\|_F^2 \leq \rho_2 \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + \Gamma_3 \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + \Gamma_4 \|\mathbf{H}^t\|_F^2 + \Gamma_5 \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2.$$

Here,  $\rho_2$  is the contraction parameter satisfying  $\rho_2 = C(\gamma, \gamma') \cdot (1 - \mu_2 \tau / 4) < 1$ , where  $C(\gamma, \gamma')$  is defined in Theorem 4.6, and  $\Gamma_3, \Gamma_4$  and  $\Gamma_5$  are constants satisfying

$$\begin{aligned} \Gamma_3 &= C(\gamma, \gamma') \cdot \left( \frac{4\tau K^2}{\mu_2} + 3\tau^2(1 + K)^2 \right), & \Gamma_4 &= C(\gamma, \gamma') \cdot \frac{\tau(\gamma + 1)\beta \alpha r \sigma_1}{\mu_2}, \\ \Gamma_5 &= C(\gamma, \gamma') \cdot \left( \frac{4\tau(\gamma' + 1)s}{\mu_2} + 3\tau^2(2\gamma' + 1)s \right). \end{aligned}$$

Now we are ready to prove Theorem 4.6.

*Proof of Theorem 4.6.* Given a fixed step size  $\tau$ , we set  $\gamma, \gamma'$  such that  $\gamma' \geq 1 + 256/(\mu_2^2 \tau^2)$  and  $\gamma \geq \max\{5, 1 + 64^2/(\mu_2 \tau)^2\}$ , then we obtain

$$\rho_2 = \left( 1 + \sqrt{\frac{2}{\gamma - 1}} \right)^2 \cdot \left( 1 + \frac{2}{\sqrt{\gamma' - 1}} \right) \cdot \left( 1 - \frac{\mu_2 \tau}{4} \right) \leq 1 - \frac{\mu_2 \tau}{16}.$$

Consider iteration stage  $t$ . According to Lemmas B.1 and B.2, we have

$$\begin{aligned} d^2(\mathbf{Z}^{t+1}, \mathbf{Z}^*) + \frac{1}{\sigma_1} \|\mathbf{S}^{t+1} - \mathbf{S}^*\|_F^2 &\leq \left(\rho_1 + \frac{\Gamma_4}{\sigma_1}\right) \cdot d^2(\mathbf{Z}^t, \mathbf{Z}^*) + \frac{1}{\sigma_1} (\rho_2 + \Gamma_1 \sigma_1) \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 \\ &+ \left(-\frac{\eta\mu_1}{4} + \frac{\Gamma_3}{\sigma_1}\right) \cdot \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + \Gamma_2 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 + \frac{\Gamma_5}{\sigma_1} \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2. \end{aligned}$$

Recall the formula of  $\Gamma_1$  and  $\Gamma_3, \Gamma_4$  from Lemmas B.1 and B.2 respectively. Note that under condition  $\eta = c_1/\sigma_1$  and  $\beta = c_3/(\alpha r \kappa)$ , we can set  $c_3$  to be sufficiently small such that

$$\Gamma_4 = C(\gamma, \gamma') \cdot \frac{\tau(\gamma + 1)\beta\alpha r \sigma_1}{\mu_2} \leq \frac{c_1 \mu'_1 \sigma_r}{80},$$

where  $\mu'_1 = \min\{\mu_1, 2\}$ , which implies that  $\rho_1 + \Gamma_4/\sigma_1 \leq 1 - \eta\mu'_1\sigma_r/80$ . Besides, under condition that  $K$  is sufficiently small, we can set  $c_1 \leq \min\{\mu_2/50, \tau/96\}$  such that the following inequality holds

$$\Gamma_1 \sigma_1 = 48c_1^2(1+K)^2 + c_1 \left(\frac{4K^2}{\mu_1} + \mu_2\right) \leq 50c_1^2 + 2c_1\mu_2 \leq 3\mu_2 c_1 \leq \frac{\mu_2 \tau}{32}. \quad (\text{B.3})$$

Finally, consider the formula of  $\Gamma_3$ . Note that similarly we can set  $K$  to be small enough such that

$$\Gamma_3 = C(\gamma, \gamma') \cdot \left(\frac{4\tau K^2}{\mu_2} + 3\tau^2(1+K)^2\right) \leq 4\tau^2,$$

thus as long as  $\tau$  is sufficiently small, there exist  $c_1$  such that  $16\tau^2/\mu_1 \leq c_1 \leq \min\{\mu_2/50, \tau/96\}$ , which implies  $\Gamma_3 \leq c_1\mu_1/4$  while ensuring (B.3) holds as well. Therefore, we obtain

$$\begin{aligned} d^2(\mathbf{Z}^{t+1}, \mathbf{Z}^*) + \frac{1}{\sigma_1} \|\mathbf{S}^{t+1} - \mathbf{S}^*\|_F^2 &\leq \left(1 - \frac{\eta\mu'_1\sigma_r}{80}\right) \cdot d^2(\mathbf{Z}^t, \mathbf{Z}^*) + \frac{1}{\sigma_1} \left(1 - \frac{\mu_2\tau}{32}\right) \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 \\ &+ \Gamma_2 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 + \frac{\Gamma_5}{\sigma_1} \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2. \end{aligned}$$

For simplicity, we denote  $D(\mathbf{Z}^t, \mathbf{S}^t) = d^2(\mathbf{Z}^t, \mathbf{Z}^*) + \|\mathbf{S}^t - \mathbf{S}^*\|_F^2/\sigma_1$ , and  $\rho = \max\{1 - \eta\mu'_1\sigma_r/80, 1 - \mu_2\tau/32\} \in (0, 1)$ , then we have

$$D(\mathbf{Z}^{t+1}, \mathbf{S}^{t+1}) \leq \rho D(\mathbf{Z}^t, \mathbf{S}^t) + \Gamma_2 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 + \frac{\Gamma_5}{\sigma_1} \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2.$$

Recall the formula of  $\Gamma_2$  and  $\Gamma_5$  in Lemmas B.1 and B.2 respectively. Under Condition 4.5, we can always set the sample size  $n$  to be large enough such that

$$\Gamma_2 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 + \frac{\Gamma_5}{\sigma_1} \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2 \leq \Gamma_2 \epsilon_1^2(n, \delta) + \frac{\Gamma_5}{\sigma_1} \epsilon_2^2(n, \delta) \leq (1 - \rho) c_2^2 \sigma_r$$

holds with probability at least  $1 - \delta$ . Thus as long as  $D(\mathbf{Z}^0, \mathbf{S}^0) \leq c_2^2 \sigma_r$ , we have by induction  $D(\mathbf{Z}^t, \mathbf{S}^t) \leq c_2^2 \sigma_r$  for any  $t \geq 0$ , which implies  $\mathbf{Z}^t \in \mathbb{B}(c_2\sqrt{\sigma_r})$ , for any  $t \geq 0$ . Hence, we obtain

$$D(\mathbf{Z}^t, \mathbf{S}^t) \leq \rho^t D(\mathbf{Z}^0, \mathbf{S}^0) + \frac{\Gamma_2}{1 - \rho} \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 + \frac{\Gamma_5}{(1 - \rho)\sigma_1} \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2,$$

which completes the proof.  $\square$

## B.2 Proof of Theorem 4.8

In order to prove Theorem 4.8, we need to make use of the following lemma. Lemma B.3 characterizes a variation of regularity condition for the sample loss function  $\mathcal{L}_n$  with respect to the sparse structure, which is proved in Section C.3.

**Lemma B.3.** Suppose the sample loss function  $\mathcal{L}_n$  satisfies Condition 4.3. Given a fixed rank- $r$  matrix  $\mathbf{X}$ , for any sparse matrices  $\mathbf{S}_1, \mathbf{S}_2 \in \mathbb{R}^{d_1 \times d_2}$  with cardinality at most  $\gamma's$ , we have

$$\begin{aligned} \langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_1 - \mathbf{S}_2 \rangle &\geq \frac{\mu_2}{2} \|\mathbf{S}_1 - \mathbf{S}_2\|_F^2 \\ &+ \frac{1}{2L_2} \|\mathcal{P}_{\Omega}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2))\|_F^2, \end{aligned}$$

where  $\Omega \subseteq [d_1] \times [d_2]$  is an index set with cardinality at most  $\tilde{s}$  such that  $\text{supp}(\mathbf{S}_1) \subseteq \Omega$  and  $\mathcal{P}_{\Omega}$  is the projection operator onto  $\Omega$ .

*Proof of Theorem 4.8.* Consider a fixed iteration  $\ell$  in Algorithm 2. As for the sparse structure, we have

$$\mathbf{S}_{\ell+1} = \mathcal{H}_{\lambda s}(\mathbf{S}_{\ell} - \tau' \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}_{\ell})).$$

Denote  $\Omega' = \text{supp}(\mathbf{S}^*) \cup \text{supp}(\mathbf{S}_{\ell}) \cup \text{supp}(\mathbf{S}_{\ell+1})$ , then we have  $\lambda s \leq |\Omega'| \leq (2\lambda + 1)s$ . We further denote  $\tilde{\mathbf{S}}_{\ell+1} = \mathcal{P}_{\Omega'}(\mathbf{S}_{\ell} - \tau' \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}_{\ell}))$ , then we obtain  $\mathbf{S}_{\ell+1} = \mathcal{H}_{\lambda s}(\tilde{\mathbf{S}}_{\ell+1})$ . Thus, according to Lemma 3.3 in [35], we have

$$\|\mathbf{S}_{\ell+1} - \mathbf{S}^*\|_F^2 \leq \left(1 + \frac{2}{\sqrt{\lambda' - 1}}\right) \cdot \|\tilde{\mathbf{S}}_{\ell+1} - \mathbf{S}^*\|_F^2. \quad (\text{B.4})$$

Therefore, it is sufficient to upper bound  $\|\tilde{\mathbf{S}}_{\ell+1} - \mathbf{S}^*\|_F$  for the sparse structure. We have

$$\begin{aligned} \|\tilde{\mathbf{S}}_{\ell+1} - \mathbf{S}^*\|_F &= \|\mathbf{S}_{\ell} - \mathbf{S}^* - \tau' \mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}_{\ell}))\|_F \\ &\leq \underbrace{\|\mathbf{S}_{\ell} - \mathbf{S}^* - \tau' \mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}_{\ell}) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}^*))\|_F}_{I_1} \\ &\quad + \underbrace{\tau' \|\mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F}_{I_2} + \underbrace{\tau' \|\mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F}_{I_3}, \end{aligned} \quad (\text{B.5})$$

where the second inequality follows from the triangle inequality. As for the first term  $I_1$  in (B.5), according to Lemma B.3, we have

$$\begin{aligned} I_1^2 &\leq (1 - \mu_2 \tau') \cdot \|\mathbf{S}_{\ell} - \mathbf{S}^*\|_F^2 - \left(\frac{\tau'}{L_2} - \tau'^2\right) \cdot \|\mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}_{\ell}) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}^*))\|_F^2 \\ &\leq (1 - \mu_2 \tau') \cdot \|\mathbf{S}_{\ell} - \mathbf{S}^*\|_F^2, \end{aligned} \quad (\text{B.6})$$

provided that  $\tau' \leq 1/L_2$ . Consider the second term  $I_2$  in (B.5). Note that  $|\Omega'| \leq (2\lambda + 1)s$ , thus according to the definition of Frobenius norm, we have

$$\begin{aligned} I_2 &= \sup_{\|\mathbf{W}\|_F \leq 1} \langle \mathcal{P}_{\Omega'}(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}_{\ell} + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)), \mathbf{W} \rangle \\ &\leq \sup_{\|\mathbf{W}\|_F \leq 1} \{|\langle \mathbf{X}_{\ell} - \mathbf{X}^*, \mathcal{P}_{\Omega'}(\mathbf{W}) \rangle| + K \|\mathbf{X}_{\ell} - \mathbf{X}^*\|_F \cdot \|\mathcal{P}_{\Omega'}(\mathbf{W})\|_F\} \\ &\leq \|\mathbf{X}_{\ell} - \mathbf{X}^*\|_{\infty, \infty} \cdot \|\mathcal{P}_{\Omega'}(\mathbf{W})\|_{1,1} + K \|\mathbf{X}_{\ell} - \mathbf{X}^*\|_F \leq 4\zeta^* \sqrt{\lambda s} + K \|\mathbf{X}_{\ell} - \mathbf{X}^*\|_F, \end{aligned} \quad (\text{B.7})$$

where the first inequality follows from Condition 4.4, the second inequality holds because  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_{1,1} \cdot \|\mathbf{B}\|_{\infty, \infty}$  and  $\|\mathcal{P}_{\Omega'}(\mathbf{W})\|_F \leq \|\mathbf{W}\|_F \leq 1$ , and the last inequality is due to the fact that  $\|\mathbf{X}_{\ell}\|_{\infty, \infty} \leq \zeta^*$ ,  $\|\mathbf{X}^*\|_{\infty, \infty} \leq \zeta^*$  and the triangle inequality. And for the third term  $I_3$ , we have

$$I_3 \leq \sqrt{(2\lambda + 1)s} \cdot \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}. \quad (\text{B.8})$$

Therefore, plugging (B.6), (B.7) and (B.8) into (B.5), we obtain

$$\begin{aligned} \|\tilde{\mathbf{S}}_{\ell+1} - \mathbf{S}^*\|_F &\leq \sqrt{1 - \mu_2 \tau'} \cdot \|\mathbf{S}_{\ell} - \mathbf{S}^*\|_F + \tau' K \|\mathbf{X}_{\ell} - \mathbf{X}^*\|_F + 2\tau' \zeta^* \sqrt{s} \\ &\quad + \tau' \sqrt{(2\lambda + 1)s} \cdot \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}. \end{aligned} \quad (\text{B.9})$$

Hence, combining (B.4) and (B.9), we obtain the following result for sparse structure

$$\begin{aligned} \|\mathbf{S}_{\ell+1} - \mathbf{S}^*\|_F &\leq \left(1 + \frac{2}{\sqrt{\lambda-1}}\right) \cdot (\sqrt{1 - \mu_2 \tau'} \cdot \|\mathbf{S}_\ell - \mathbf{S}^*\|_F + \tau' K \|\mathbf{X}_\ell - \mathbf{X}^*\|_F) \\ &\quad + \tau' \left(1 + \frac{2}{\sqrt{\lambda-1}}\right) \cdot (4\zeta^* \sqrt{\lambda s} + \sqrt{3\lambda s} \cdot \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}). \end{aligned} \quad (\text{B.10})$$

Next, let us consider the low-rank structure. According to Algorithm 2, we have

$$\mathbf{X}_{\ell+1} = \mathcal{P}_{\lambda', \zeta^*}(\mathbf{X}_\ell - \eta' \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}_\ell + \mathbf{S}_\ell)),$$

where the projection operator  $\mathcal{P}_{\lambda', \zeta^*}$  is defined as

$$\mathcal{P}_{\lambda', \zeta^*}(\mathbf{X}) = \underset{\text{rank}(\mathbf{Y}) \leq \lambda' r, \|\mathbf{Y}\|_{\infty, \infty} \leq \zeta^*}{\text{argmin}} \|\mathbf{Y} - \mathbf{X}\|_F, \text{ for any } \mathbf{X} \in \mathbb{R}^{d_1 \times d_2}.$$

Let the singular value decomposition of  $\mathbf{X}_\ell, \mathbf{X}_{\ell+1}$  be  $\mathbf{X}_\ell = \bar{\mathbf{U}}^\ell \bar{\Sigma}^\ell \bar{\mathbf{V}}^{\ell \top}$  and  $\mathbf{X}_{\ell+1} = \bar{\mathbf{U}}^{\ell+1} \bar{\Sigma}^{\ell+1} \bar{\mathbf{V}}^{\ell+1 \top}$  respectively. Define the following subspace spanned by the column vectors of  $\bar{\mathbf{U}}^*, \bar{\mathbf{U}}^\ell$  and  $\bar{\mathbf{U}}^{\ell+1}$  as

$$\text{span}(\tilde{\mathbf{U}}) = \text{span}\{\bar{\mathbf{U}}^*, \bar{\mathbf{U}}^\ell, \bar{\mathbf{U}}^{\ell+1}\} = \text{col}(\bar{\mathbf{U}}^*) + \text{col}(\bar{\mathbf{U}}^\ell) + \text{col}(\bar{\mathbf{U}}^{\ell+1}),$$

where each column vector of  $\tilde{\mathbf{U}}$  is a basis vector of the above subspace. Similarly, we define the subspace spanned by the column vectors of  $\bar{\mathbf{V}}^*, \bar{\mathbf{V}}^\ell$  and  $\bar{\mathbf{V}}^{\ell+1}$  as

$$\text{span}(\tilde{\mathbf{V}}) = \text{span}\{\bar{\mathbf{V}}^*, \bar{\mathbf{V}}^\ell, \bar{\mathbf{V}}^{\ell+1}\} = \text{col}(\bar{\mathbf{V}}^*) + \text{col}(\bar{\mathbf{V}}^\ell) + \text{col}(\bar{\mathbf{V}}^{\ell+1}),$$

Note that  $\mathbf{X}^*$  has rank  $r$ ,  $\mathbf{X}_\ell$  and  $\mathbf{X}_{\ell+1}$  has rank at most  $\lambda' r$ , thus both  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  have at most  $(2\lambda' + 1)r$  columns. Moreover, we further define the following subspace

$$\mathcal{A} = \{\Delta \in \mathbb{R}^{d_1 \times d_2} \mid \text{row}(\Delta) \subseteq \text{span}(\tilde{\mathbf{V}}) \text{ and } \text{col}(\Delta) \subseteq \text{span}(\tilde{\mathbf{U}})\}.$$

Let  $\Pi_{\mathcal{A}}$  be the projection operator onto  $\mathcal{A}$ , then for any  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ , we have  $\Pi_{\mathcal{A}}(\mathbf{X}) = \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{X} \tilde{\mathbf{V}} \tilde{\mathbf{V}}^\top$ . Note that for any  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ , we have  $\text{rank}(\Pi_{\mathcal{A}}(\mathbf{X})) \leq (2\lambda' + 1)r$ , since  $\text{rank}(\mathbf{A}\mathbf{B}) \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})\}$ . Besides, we denote

$$\tilde{\mathbf{X}}_{\ell+1} = \mathbf{X}_\ell - \eta' \Pi_{\mathcal{A}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}_\ell + \mathbf{S}_\ell)).$$

Similar to the proof of Theorem 5.9 in [49], we have  $\mathbf{X}_{\ell+1}$  is actually the best rank- $\lambda' r$  approximation of  $\tilde{\mathbf{X}}_{\ell+1}$  satisfying the infinity norm constraint, or in other words,  $\mathbf{X}_{\ell+1} = \mathcal{P}_{\lambda', \zeta^*}(\tilde{\mathbf{X}}_{\ell+1})$ . Note that  $\mathcal{P}_{\lambda', \zeta^*}(\mathbf{X}^*) = \mathbf{X}^*$ , thus according to Lemma 3.18 in [35], we obtain

$$\|\mathbf{X}_{\ell+1} - \mathbf{X}^*\|_F^2 = \|\mathcal{P}_{\lambda', \zeta^*}(\tilde{\mathbf{X}}_{\ell+1}) - \mathbf{X}^*\|_F^2 \leq \left(1 + \frac{2}{\sqrt{\lambda' - 1}}\right) \cdot \|\tilde{\mathbf{X}}_{\ell+1} - \mathbf{X}^*\|_F^2. \quad (\text{B.11})$$

Thus, it suffices to bound the term  $\|\tilde{\mathbf{X}}_{\ell+1} - \mathbf{X}^*\|_F$ . Note that  $\mathbf{X}^* \in \mathcal{A}$ , thus according to the triangle inequality, we have

$$\begin{aligned} \|\tilde{\mathbf{X}}_{\ell+1} - \mathbf{X}^*\|_F &\leq \underbrace{\|\mathbf{X}_\ell - \mathbf{X}^* - \eta' \Pi_{\mathcal{A}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}_\ell + \mathbf{S}_\ell) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}_\ell))\|_F}_{I'_1} \\ &\quad + \eta' \underbrace{\|\Pi_{\mathcal{A}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}_\ell) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F}_{I'_2} + \eta' \underbrace{\|\Pi_{\mathcal{A}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F}_{I'_3}. \end{aligned} \quad (\text{B.12})$$

Consider  $I'_1$  in (B.12) first. According to Lemma B.2 in [50], we have

$$I'_1{}^2 \leq (1 - \eta' \mu_1) \cdot \|\mathbf{X}_\ell - \mathbf{X}^*\|_F^2, \quad (\text{B.13})$$

provided that  $\eta' \leq 1/L_1$ . As for the second term  $I'_2$  in (B.12), by the definition of Frobenius norm, we have

$$\begin{aligned} I'_2 &= \sup_{\|\mathbf{W}\|_F \leq 1} \langle \Pi_{\mathcal{A}_{3r}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}_\ell) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)), \mathbf{W} \rangle \\ &\leq (1+K) \cdot \|\mathbf{S}_\ell - \mathbf{S}^*\|_F \cdot \|\Pi_{\mathcal{A}}(\mathbf{W})\|_F \leq (1+K) \cdot \|\mathbf{S}_\ell - \mathbf{S}^*\|_F, \end{aligned} \quad (\text{B.14})$$

where the first inequality holds because of Condition 4.4. As for  $I'_3$ , we have

$$I'_3 \leq \sqrt{(2\lambda' + 1)r} \cdot \|\Pi_{\mathcal{A}}(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_2 \leq \sqrt{(2\lambda' + 1)r} \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2. \quad (\text{B.15})$$

Therefore, plugging (B.13), (B.14) and (B.15) into (B.12), we obtain

$$\|\tilde{\mathbf{X}}_{\ell+1} - \mathbf{X}^*\|_F \leq \sqrt{1 - \eta'\mu_1} \cdot \|\mathbf{X}_\ell - \mathbf{X}^*\|_F + \eta'(1+K) \cdot \|\mathbf{S}_\ell - \mathbf{S}^*\|_F + \eta'\sqrt{3\lambda'r'} \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2. \quad (\text{B.16})$$

Finally, combining (B.11) and (B.16), we obtain the following result for low rank structure

$$\begin{aligned} \|\mathbf{X}_{\ell+1} - \mathbf{X}^*\|_F &\leq \left(1 + \frac{2}{\sqrt{\lambda' - 1}}\right) \cdot \left(\sqrt{1 - \eta'\mu_1} \cdot \|\mathbf{X}_\ell - \mathbf{X}^*\|_F + \eta'(1+K) \cdot \|\mathbf{S}_\ell - \mathbf{S}^*\|_F\right) \\ &\quad + \eta' \left(1 + \frac{2}{\sqrt{\lambda' - 1}}\right) \cdot \sqrt{3\lambda'r'} \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2. \end{aligned} \quad (\text{B.17})$$

Hence, combining (B.10) and (B.17), we obtain

$$\begin{aligned} \|\mathbf{X}_{\ell+1} - \mathbf{X}^*\|_F + \|\mathbf{S}_{\ell+1} - \mathbf{S}^*\|_F &\leq \rho'_1 \|\mathbf{X}_\ell - \mathbf{X}^*\|_F + \rho'_2 \|\mathbf{S}_\ell - \mathbf{S}^*\|_F + 4\sqrt{\lambda'r'} \left(1 + \frac{2}{\sqrt{\lambda' - 1}}\right) \cdot \zeta^* \sqrt{s} \\ &\quad + \Gamma_1 \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2 + \Gamma_2 \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}, \end{aligned} \quad (\text{B.18})$$

where  $\Gamma_1 = \eta'(1 + 2/\sqrt{\lambda' - 1})\sqrt{3\lambda'r}$ ,  $\Gamma_2 = \tau'(1 + 2/\sqrt{\lambda - 1})\sqrt{3\lambda s}$ , and contraction parameter  $\rho'_1, \rho'_2$  are defined in Theorem 4.8. Note that we set  $\eta' = 1/(6\mu_1) \leq 1/L_1$ ,  $\tau' = 3/(4\mu_2) \leq 1/L_2$ , and we assume  $\mu_1 \geq 1/3$ . Then with sufficient large  $\lambda$  and  $\lambda'$  and structural Lipschitz gradient parameter  $K$  small enough, we could guarantee  $\rho'_1, \rho'_2 \in (0, 19/20)$ . Plugging in the definition of  $\zeta^* = c_0 \alpha r \kappa / \sqrt{d_1 d_2}$ , we complete the proof by induction.  $\square$

### B.3 Proof of Theorem 4.9

*Proof.* To prove Theorem 4.9, it is sufficient to verify the assumption  $D(\mathbf{Z}^0, \mathbf{S}^0) \leq c_4^2 \sigma_r$  in Theorem 4.6. Thus, according to Theorem 4.8, it is sufficient to make sure the right hand side of (4.2) is small enough.

As for the optimization error, i.e., the first term on the R.H.S. of (4.2), we can perform  $L \geq \log\{c\sigma_r/(2\|\mathbf{X}^*\|_F + 2\|\mathbf{S}^*\|_F)\}/\log(\rho')$  iterations in Algorithm 2 to make sure the optimization error is sufficiently small such that  $\rho'^L \cdot (\|\mathbf{X}^*\|_F + \|\mathbf{S}^*\|_F) \leq c\sigma_r/2$ , where  $c = \min\{1/2, c_4/4\}$ .

On the other hand, for the statistical error, i.e., the last three terms on the R.H.S. of (4.2), we assume  $s \leq cd_1 d_2 / (\alpha^2 r^2 \kappa^2)$ , where  $c$  is a small enough constant, and sample size  $n$  is sufficiently large such that  $\Gamma_1 \sqrt{r} \epsilon_1(n, \delta) + \Gamma_2 \sqrt{s} \epsilon_2(n, \delta) \leq c\sigma_r/4$ . Putting pieces together, we arrive at  $\|\mathbf{X}^0 - \mathbf{X}^*\|_F + \|\mathbf{S}^0 - \mathbf{S}^*\|_F \leq c \cdot \sigma_r$ . Finally, based on Lemma 5.14 in [46], the initial assumption that  $D(\mathbf{Z}^0, \mathbf{S}^0) \leq c_4^2 \sigma_r$  in Theorem 4.6 is satisfied, which completes the proof.  $\square$

## C Proofs of Technical Lemmas

### C.1 Proof of Lemma B.1

In order to prove Theorem B.1, we need to make use of the following lemmas. Lemma C.1 characterizes a local curvature property of the low-rank structure, which gives us the lower bound of the inner product term. We provide its proof in Section D.1. Lemma C.2, proved in Section D.2, characterizes a local smoothness property of the low-rank structure and gives us an upper bound of the Frobenius term.

**Lemma C.1** (Local Curvature Property for Low-Rank Structure). Suppose the sample loss function  $\mathcal{L}_n$  satisfies Conditions 4.2 and 4.4. Recall that  $\mathbf{X}^* = \mathbf{U}^* \mathbf{V}^{*\top}$  is the unknown rank- $r$  matrix that satisfies (3.1), and  $\mathbf{S}^*$  is the unknown  $s$ -sparse matrix. Let  $\mathbf{Z} \in \mathbb{R}^{(d_1+d_2) \times r}$  be any matrix with  $\mathbf{Z} = [\mathbf{U}; \mathbf{V}]$ , where  $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$ ,  $\mathbf{V} \in \mathbb{R}^{d_2 \times r}$

satisfy  $\|\mathbf{U}\|_{2,\infty} \leq 2\sqrt{\alpha r \sigma_1 / d_1}$  and  $\|\mathbf{V}\|_{2,\infty} \leq 2\sqrt{\alpha r \sigma_1 / d_2}$ . Let  $\mathbf{S} \in \mathbb{R}^{d_1 \times d_2}$  be any matrix with at most  $\beta'$ -fraction nonzero entries per row and column and satisfying  $\|\mathbf{S}\|_0 \leq s' \leq \tilde{s}$ . Denote the optimal rotation with respect to  $\mathbf{Z}$  by  $\mathbf{R} = \operatorname{argmin}_{\tilde{\mathbf{R}} \in \mathbb{Q}_r} \|\mathbf{Z} - \mathbf{Z}^* \tilde{\mathbf{R}}\|_F$ , and  $\mathbf{H} = \mathbf{Z} - \mathbf{Z}^* \mathbf{R}$ , then we have

$$\begin{aligned} \langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S}), \mathbf{H} \rangle &\geq \frac{\mu_1}{4} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{1}{16} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 + \left( \frac{\mu'_1}{20} \sigma_r - C \right) \cdot \|\mathbf{H}\|_F^2 - \left( \frac{L_1 + 1}{8} + \frac{1}{\mu_2} \right) \cdot \|\mathbf{H}\|_F^4 \\ &\quad - \left( \frac{\mu_2}{2} + \frac{2K^2}{\mu_1} \right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 - \left( \frac{8r}{\mu_1} + \frac{r}{L_1} \right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2, \end{aligned}$$

where  $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$ ,  $\mu'_1 = \min\{\mu_1, 2\}$ , and  $C = 18(\beta' + \beta)\alpha r \sigma_1 / \mu_2$ .

**Lemma C.2** (Local Smoothness Property for Low-Rank Structure). Suppose the sample loss function  $\mathcal{L}_n$  satisfies Conditions 4.2 and 4.4. Recall that  $\mathbf{X}^*$  is the unknown rank- $r$  matrix and  $\mathbf{S}^*$  is the unknown  $s$ -sparse matrix. For any matrix  $\mathbf{Z} = [\mathbf{U}; \mathbf{V}] \in \mathbb{R}^{(d_1+d_2) \times r}$  and  $\mathbf{S} \in \mathbb{R}^{d_1 \times d_2}$  with at most  $s'$  nonzero entries satisfying  $s' \leq \tilde{s}$ , we have

$$\begin{aligned} \|\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S})\|_F^2 &\leq \left( 12L_1^2 \|\mathbf{X} - \mathbf{X}^*\|_F^2 + 12(1+K)^2 \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 + \|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2 \right) \cdot \|\mathbf{Z}\|_2^2 \\ &\quad + 12r \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{Z}\|_2^2, \end{aligned}$$

where  $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$ .

*Proof of Lemma B.1.* Recall  $\mathbf{Z}^* = [\mathbf{U}^*; \mathbf{V}^*]$  and  $\mathbf{X}^* = \mathbf{U}^* \mathbf{V}^{*\top}$ , where  $\mathbf{U}^* = \bar{\mathbf{U}}^*(\boldsymbol{\Sigma}^*)^{1/2}$ ,  $\mathbf{V}^* = \bar{\mathbf{V}}^*(\boldsymbol{\Sigma}^*)^{1/2}$ , we have  $\|\mathbf{Z}^*\|_2 = \sqrt{2\sigma_1}$ . According to our initial ball assumption  $\mathbf{Z}^0 \in \mathbb{B}(\sqrt{\sigma_r}/4)$ , there exists an orthogonal matrix  $\mathbf{R} \in \mathbb{R}^{r \times r}$  such that  $\|\mathbf{Z}^0 - \mathbf{Z}^* \mathbf{R}\|_F \leq \sqrt{\sigma_r}/4$ , thus we obtain

$$\sqrt{\sigma_1} \leq \|\mathbf{Z}^*\|_2 - \|\mathbf{Z}^0 - \mathbf{Z}^* \mathbf{R}\|_2 \leq \|\mathbf{Z}^0\|_2 \leq \|\mathbf{Z}^*\|_2 + \|\mathbf{Z}^0 - \mathbf{Z}^* \mathbf{R}\|_F \leq 2\sqrt{\sigma_1}.$$

Recall (3.1) and the definition of  $\mathcal{C}_1, \mathcal{C}_2$  in (3.2), then it is obvious that  $\mathbf{U}^* \in \mathcal{C}_1$  and  $\mathbf{V}^* \in \mathcal{C}_2$ . Consider a fixed iteration stage  $t$ , we denote

$$\begin{aligned} \tilde{\mathbf{U}}^{t+1} &= \mathbf{U}^t - \eta \nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}^t \mathbf{V}^{t\top} + \mathbf{S}^t) - \frac{1}{2} \eta \mathbf{U}^t (\mathbf{U}^{t\top} \mathbf{U}^t - \mathbf{V}^{t\top} \mathbf{V}^t), \\ \tilde{\mathbf{V}}^{t+1} &= \mathbf{V}^t - \eta \nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{U}^t \mathbf{V}^{t\top} + \mathbf{S}^t) - \frac{1}{2} \eta \mathbf{V}^t (\mathbf{V}^{t\top} \mathbf{V}^t - \mathbf{U}^{t\top} \mathbf{U}^t). \end{aligned}$$

Denote  $\tilde{\mathbf{Z}}^{t+1} = [\tilde{\mathbf{U}}^{t+1}; \tilde{\mathbf{V}}^{t+1}]$ , and  $\mathbf{Z}^t = [\mathbf{U}^t; \mathbf{V}^t]$ , for any iteration stage  $t$ , then according to (B.2), we have  $\tilde{\mathbf{Z}}^{t+1} = \mathbf{Z}^t - \eta \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t)$ . Besides, according to Algorithm 1, we obtain

$$\mathbf{U}^{t+1} = \mathcal{P}_{\mathcal{C}_1}(\tilde{\mathbf{U}}^{t+1}) \quad \text{and} \quad \mathbf{V}^{t+1} = \mathcal{P}_{\mathcal{C}_2}(\tilde{\mathbf{V}}^{t+1}).$$

Recall  $\mathbf{Z}^* = [\mathbf{U}^*; \mathbf{V}^*]$ , and  $\mathbf{R}^t = \operatorname{argmin}_{\mathbf{R} \in \mathbb{Q}_r} \|\mathbf{Z}^t - \mathbf{Z}^* \mathbf{R}\|_F$ , for any  $t$ . Denote  $\mathbf{H}^t = \mathbf{Z}^t - \mathbf{Z}^* \mathbf{R}^t$ . Since  $\mathcal{C}_1, \mathcal{C}_2$  are both rotation-invariant constraint sets, and  $\mathbf{U}^* \in \mathcal{C}_1, \mathbf{V}^* \in \mathcal{C}_2$ , we have

$$\begin{aligned} d^2(\mathbf{Z}^{t+1}, \mathbf{Z}^*) &\leq \|\mathbf{Z}^{t+1} - \mathbf{Z}^* \mathbf{R}^t\|_F^2 \\ &\leq \|\mathbf{Z}^t - \eta \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t) - \mathbf{Z}^* \mathbf{R}^t\|_F^2 \\ &= d^2(\mathbf{Z}^t, \mathbf{Z}^*) - 2\eta \langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t), \mathbf{H}^t \rangle + \eta^2 \|\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t)\|_F^2, \end{aligned} \tag{C.1}$$

where the first inequality follows from Definition 4.1, and the second inequality is due to the nonexpansive property of projection  $\mathcal{P}_{\mathcal{C}_i}$  onto  $\mathcal{C}_i$ , where  $i \in \{1, 2\}$ . Therefore, it suffices to lower bound the inner product term  $\langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t), \mathbf{H}^t \rangle$  and upper bound the term  $\|\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t)\|_F^2$ , respectively. According to Algorithm 1, we have  $(\mathbf{U}^t, \mathbf{V}^t)$  satisfies the condition of  $(\mathbf{U}, \mathbf{V})$  in Lemma C.1, and  $\mathbf{S}^t$  has at most  $\gamma\beta$ -fraction nonzero entries per row and column with  $\|\mathbf{S}\|_0 \leq \gamma's \leq \tilde{s}$ . Denote  $\mathbf{X}^t = \mathbf{U}^t \mathbf{V}^{t\top}$ , then according to Lemma C.1, we obtain

$$\begin{aligned} \langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t), \mathbf{H}^t \rangle &\geq \frac{\mu_1}{4} \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + \frac{1}{16} \|\mathbf{U}^{t\top} \mathbf{U}^t - \mathbf{V}^{t\top} \mathbf{V}^t\|_F^2 + \left( \frac{\mu'_1}{20} \sigma_r - C \right) \cdot \|\mathbf{H}^t\|_F^2 \\ &\quad - \left( \frac{L_1 + 1}{8} + \frac{1}{\mu_2} \right) \cdot \|\mathbf{H}^t\|_F^4 - \left( \frac{\mu_2}{2} + \frac{2K^2}{\mu_1} \right) \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 - \left( \frac{8r}{\mu_1} + \frac{r}{L_1} \right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2, \end{aligned}$$

where  $\mu'_1 = \min\{\mu_1, 2\}$ , and  $C = 18(\gamma + 1)\beta\alpha r\sigma_1/\mu_2$ . Besides, according to Lemma C.2, we have

$$\begin{aligned} \|\nabla_{\mathbf{Z}}\tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t)\|_F^2 &\leq \left(12L_1^2\|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + 12(1+K)^2 \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + \|\mathbf{U}^{t\top}\mathbf{U}^t - \mathbf{V}^{t\top}\mathbf{V}^t\|_F^2\right) \cdot \|\mathbf{Z}^t\|_2^2 \\ &\quad + 12r\|\nabla_{\mathbf{X}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{Z}^t\|_2^2. \end{aligned}$$

Note that under the assumption of  $\mathbf{Z}^t \in \mathbb{B}(c_2\sqrt{\sigma_r})$ , where  $c_2 \leq 1/4$ , we have  $\|\mathbf{Z}^t\|_2 \leq \|\mathbf{Z}^*\mathbf{R}^t\|_2 + \|\mathbf{Z}^t - \mathbf{Z}^*\mathbf{R}^t\|_2 \leq 2\sqrt{\sigma_1}$ , since  $\|\mathbf{Z}^*\|_2^2 = 2\sigma_1$ . Thus, if we set the step size  $\eta = c_1/\sigma_1$ , where  $c_1 \leq \min\{1/32, \mu_1/(192L_1^2)\}$ , and we assume  $\beta \leq 1/(c_3\alpha r\kappa)$  with  $c_3$  large enough such that  $c_3 \geq 720(\gamma + 1)\mu_2/\mu'_1$ , we have

$$\begin{aligned} -2\eta\langle\nabla_{\mathbf{Z}}\tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t), \mathbf{H}^t\rangle + \eta^2\|\nabla_{\mathbf{Z}}\tilde{F}_n(\mathbf{Z}^t, \mathbf{S}^t)\|_F^2 &\leq -\frac{\eta\mu_1}{4}\|\mathbf{X}^t - \mathbf{X}^*\|_F^2 - \frac{\eta\mu'_1\sigma_r}{20}\|\mathbf{H}^t\|_F^2 \\ &\quad + \eta\left(\frac{L_1+1}{4} + \frac{2}{\mu_2}\right) \cdot \|\mathbf{H}^t\|_F^4 + C_1\|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + C_2\|\nabla_{\mathbf{X}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2, \end{aligned}$$

where  $C_1 = 48\eta^2(1+K)^2\sigma_1 + \eta(\mu_2 + 4K^2/\mu_1)$ , and  $C_2 = 48\eta^2r\sigma_1 + 2\eta(8r/\mu_1 + r/L_1)$ . Note that according to our assumption,  $\|\mathbf{H}^t\|_F^2 \leq c_2^2\sigma_r$  with  $c_2^2 \leq \mu'_1/[10(L_1 + 1 + 8/\mu_2)]$ , thus by (C.1), we obtain

$$d^2(\mathbf{Z}^{t+1}, \mathbf{Z}^*) \leq \left(1 - \frac{\eta\mu'_1\sigma_r}{40}\right)d^2(\mathbf{Z}^t, \mathbf{Z}^*) - \frac{\eta\mu_1}{4}\|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + C_1\|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + C_2\|\nabla_{\mathbf{X}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2,$$

which completes the proof.  $\square$

## C.2 Proof of Lemma B.2

In order to prove Lemma B.2, we need to utilize the following lemma. Inspired by [56], we present Lemma C.3, which characterizes a nearly non-expansiveness property of the truncation operator  $\mathcal{T}_\theta$ , as long as  $\theta$  is large enough. We provides its proof in Section D.3 for completeness.

**Lemma C.3.** Suppose  $\mathbf{S}^* \in \mathbb{R}^{d_1 \times d_2}$  is the unknown sparse matrix with at most  $\beta$ -fraction nonzero entries per row and column. For any matrix  $\mathbf{S} \in \mathbb{R}^{d_1 \times d_2}$ , we have

$$\|\mathcal{T}_{\gamma\beta}(\mathbf{S}) - \mathbf{S}^*\|_F^2 \leq \left(1 + \sqrt{\frac{2}{\gamma-1}}\right)^2 \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2,$$

where  $\gamma > 1$  is a parameter.

Now we are ready to prove Lemma B.2.

*Proof of Lemma B.2.* Consider a fixed iteration stage  $t$ . For the sparse structure, according to Algorithm 1, we have

$$\mathbf{S}^{t+1} = \mathcal{T}_{\gamma\beta} \circ \mathcal{H}_{\gamma's}(\mathbf{S}^t - \tau\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{U}^t\mathbf{V}^{t\top} + \mathbf{S}^t)).$$

Denote  $\bar{\mathbf{S}}^{t+1} = \mathcal{H}_{\gamma's}(\mathbf{S}^t - \tau\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{U}^t\mathbf{V}^{t\top} + \mathbf{S}^t))$ , then we have  $\mathbf{S}^{t+1} = \mathcal{T}_{\gamma\beta}(\bar{\mathbf{S}}^{t+1})$ . To begin with according to Lemma C.3, we have

$$\|\mathbf{S}^{t+1} - \mathbf{S}^*\|_F^2 = \|\mathcal{T}_{\gamma\beta}(\bar{\mathbf{S}}^{t+1}) - \mathbf{S}^*\|_F^2 \leq \left(1 + \sqrt{\frac{2}{\gamma-1}}\right)^2 \cdot \|\bar{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2. \quad (\text{C.2})$$

Moreover, denote  $\Omega = \Omega^* \cup \Omega^t \cup \Omega^{t+1}$ , where  $\Omega^* = \text{supp}(\mathbf{S}^*)$ ,  $\Omega^t = \text{supp}(\mathbf{S}^t)$  and  $\Omega^{t+1} = \text{supp}(\bar{\mathbf{S}}^{t+1})$ . Obviously, the cardinality of  $\Omega$  satisfies  $\gamma's \leq |\Omega| \leq (2\gamma' + 1)s$ . Based on  $\Omega$ , we define  $\tilde{\mathbf{S}}^{t+1}$  as follows

$$\tilde{\mathbf{S}}^{t+1} = \mathcal{P}_\Omega(\mathbf{S}^t - \tau\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{U}^t\mathbf{V}^{t\top} + \mathbf{S}^t)) = \mathbf{S}^t - \tau\mathcal{P}_\Omega(\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{U}^t\mathbf{V}^{t\top} + \mathbf{S}^t)), \quad (\text{C.3})$$

where  $\mathcal{P}_\Omega$  is the projection operator onto the index set  $\Omega$ . Note that  $\Omega^{t+1} \subseteq \Omega$ , thus we have  $\bar{\mathbf{S}}^{t+1} = \mathcal{H}_{\gamma's}(\tilde{\mathbf{S}}^{t+1})$ . According to Lemma 3.3 in [35], we have

$$\|\bar{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2 \leq \left(1 + \frac{2}{\sqrt{\gamma'-1}}\right) \cdot \|\tilde{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2. \quad (\text{C.4})$$



Therefore, it is sufficient to upper bound  $\|\tilde{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2$ . By (C.3), we have

$$\|\tilde{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2 = \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 - 2\tau \underbrace{\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t), \mathbf{S}^t - \mathbf{S}^* \rangle}_{I_1} + \tau^2 \underbrace{\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t))\|_F^2}_{I_2}, \quad (\text{C.5})$$

where the equality holds because  $\langle \mathcal{P}_\Omega(\mathbf{A}), \mathbf{B} \rangle = \langle \mathbf{A}, \mathcal{P}_\Omega(\mathbf{B}) \rangle$ . In the following discussions, we are going to bound  $I_1$  and  $I_2$  respectively. Consider the term  $I_1$  first, we have

$$\begin{aligned} I_1 &= \underbrace{\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*), \mathbf{S}^t - \mathbf{S}^* \rangle}_{I_{11}} + \underbrace{\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{S}^t - \mathbf{S}^* \rangle}_{I_{12}} \\ &\quad + \underbrace{\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{S}^t - \mathbf{S}^* \rangle}_{I_{13}}. \end{aligned} \quad (\text{C.6})$$

As for the first term  $I_{11}$  in (C.6), according to Lemma B.3, we have

$$I_{11} \geq \frac{\mu_2}{2} \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + \frac{1}{2L_2} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*))\|_F^2. \quad (\text{C.7})$$

Note that we have  $\text{supp}(\mathbf{S}^t - \mathbf{S}^*) \subseteq \Omega^t \cup \Omega^*$ , where  $\Omega^t \cup \Omega^*$  has at most  $(\gamma + 1)\beta$ -fraction nonzero entries per row and column. Denote  $\mathbf{R}^t$  as the optimal rotation with respect to  $\mathbf{Z}^t = [\mathbf{U}^t; \mathbf{V}^t]$ , and  $\mathbf{H}^t = \mathbf{Z}^t - \mathbf{Z}^* \mathbf{R}^t$ . According to Condition 4.4, we obtain the bound of  $I_{12}$  in (C.6)

$$\begin{aligned} |I_{12}| &\leq |\langle \mathbf{X}^t - \mathbf{X}^*, \mathbf{S}^t - \mathbf{S}^* \rangle| + K \|\mathbf{X}^t - \mathbf{X}^*\|_F \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F \\ &\leq \|\mathcal{P}_{\Omega^t \cup \Omega^*}(\mathbf{X}^t - \mathbf{X}^*)\|_F \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F + K \|\mathbf{X}^t - \mathbf{X}^*\|_F \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F \\ &\leq \sqrt{18(\gamma + 1)\beta\alpha r\sigma_1} \|\mathbf{H}^t\|_F \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F + K \|\mathbf{X}^t - \mathbf{X}^*\|_F \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F, \end{aligned} \quad (\text{C.8})$$

where the second inequality holds because  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_F$ , and the last inequality follows from Lemma 14 in [56]. As for the last term  $I_{13}$  in (C.6), we have

$$|I_{13}| \leq \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty} \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_{1,1} \leq \sqrt{(\gamma' + 1)s} \cdot \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty} \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F, \quad (\text{C.9})$$

where the first inequality holds because  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_{\infty, \infty} \cdot \|\mathbf{B}\|_{1,1}$ , and the second inequality follows from the fact that  $\mathbf{S}^t - \mathbf{S}^*$  has at most  $(\gamma' + 1)s$  nonzero entries. Therefore, plugging (C.7), (C.8) and (C.9) into (C.6), we obtain the lower bound of  $I_1$

$$\begin{aligned} I_1 &\geq \frac{\mu_2}{8} \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + \frac{1}{2L_2} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*))\|_F^2 - \frac{2K^2}{\mu_2} \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 \\ &\quad - \frac{36(\gamma + 1)\beta\alpha r\sigma_1}{\mu_2} \|\mathbf{H}^t\|_F^2 - \frac{2(\gamma' + 1)s}{\mu_2} \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2. \end{aligned} \quad (\text{C.10})$$

Next, consider the term  $I_2$  in (C.5). We have

$$\begin{aligned} I_2 &\leq 3\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*))\|_F^2 + 3\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F^2 \\ &\quad + 3\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F^2 \end{aligned} \quad (\text{C.11})$$

As for the second term in (C.11), according to the definition of Frobenius norm, we have

$$\begin{aligned} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F &= \sup_{\|\mathbf{W}\| \leq 1} \langle \mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)), \mathbf{W} \rangle \\ &= \sup_{\|\mathbf{W}\| \leq 1} \langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathcal{P}_\Omega(\mathbf{W}) \rangle \\ &\leq (1 + K) \cdot \|\mathbf{X}^t - \mathbf{X}^*\|_F \cdot \|\mathcal{P}_\Omega(\mathbf{W})\|_F \\ &\leq (1 + K) \cdot \|\mathbf{X}^t - \mathbf{X}^*\|_F, \end{aligned} \quad (\text{C.12})$$

where the second equality holds because  $\langle \mathcal{P}_\Omega(\mathbf{A}), \mathbf{B} \rangle = \langle \mathbf{A}, \mathcal{P}_\Omega(\mathbf{B}) \rangle$ , and the first inequality holds because of Condition 4.4. As for the last term in (C.11), note that  $|\Omega| \leq (2\gamma' + 1)s$ , thus we have

$$\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\|_F^2 \leq (2\gamma' + 1)s \|\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2. \quad (\text{C.13})$$

Therefore, plugging (C.12) and (C.13) into (C.11), we obtain the upper bound of  $I_2$

$$I_2 \leq 3\|\mathcal{P}_\Omega(\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^t) - \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X}^t + \mathbf{S}^*))\|_F^2 + 3(1+K)^2 \cdot \|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + 3(2\gamma' + 1)s\|\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2. \quad (\text{C.14})$$

If we set the step size  $\tau \leq 1/(3L_2)$ , then by plugging (C.10) and (C.14) into (C.5), we have

$$\|\tilde{\mathbf{S}}^{t+1} - \mathbf{S}^*\|_F^2 \leq \left(1 - \frac{\mu_2\tau}{4}\right) \cdot \|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + C_3\|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + C_4\|\mathbf{H}^t\|_F^2 + C_5\|\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2, \quad (\text{C.15})$$

where  $C_3 = 4\tau K^2/\mu_2 + 3\tau^2(1+K)^2$ ,  $C_4 = 72\tau(\gamma+1)\beta\alpha r\sigma_1/\mu_2$  and  $C_5 = 4\tau(\gamma'+1)s/\mu_2 + 3\tau^2(2\gamma'+1)s$ . Thus combining (C.2), (C.4) and (C.15), we obtain

$$\|\mathbf{S}^{t+1} - \mathbf{S}^*\|_F^2 \leq \rho\|\mathbf{S}^t - \mathbf{S}^*\|_F^2 + C(\gamma, \gamma') \cdot (C_3\|\mathbf{X}^t - \mathbf{X}^*\|_F^2 + C_4\|\mathbf{H}^t\|_F^2 + C_5\|\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_{\infty, \infty}^2),$$

which completes the proof.  $\square$

### C.3 Proof of Lemma B.3

In order to proof Lemma B.3, we need to make use of the following lemma, which can be derived following the standard proof of Lipschitz continuous gradient property [40].

**Lemma C.4.** Suppose the sample loss function  $\mathcal{L}_n$  satisfies Conditions 4.3. Given a fixed rank- $r$  matrix  $\mathbf{X} \in \mathbb{R}^{d_1 \times d_2}$ , then for any sparse matrices  $\mathbf{S}_1, \mathbf{S}_2 \in \mathbb{R}^{d_1 \times d_2}$  with cardinality at most  $\tilde{s}$ , we have

$$\begin{aligned} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) &\geq \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2) + \langle \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_1 - \mathbf{S}_2 \rangle \\ &\quad + \frac{1}{2L_2} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2))\|_F^2, \end{aligned}$$

where  $\Omega \subseteq [d_1] \times [d_2]$  is an index set with cardinality at most  $\tilde{s}$  such that  $\text{supp}(\mathbf{S}_1) \subseteq \Omega$  and  $\mathcal{P}_\Omega$  is the projection operator onto  $\Omega$ .

Now we are ready to prove Lemma B.3.

*Proof of Lemma B.3.* Since the sample loss function  $\mathcal{L}_n$  satisfies the restricted strong convexity Condition 4.3, we have

$$\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2) \geq \mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) + \langle \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1), \mathbf{S}_2 - \mathbf{S}_1 \rangle + \frac{\mu_2}{2} \|\mathbf{S}_2 - \mathbf{S}_1\|_F^2. \quad (\text{C.16})$$

According to Lemma C.4, we have

$$\begin{aligned} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) &\geq \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2) + \langle \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_1 - \mathbf{S}_2 \rangle \\ &\quad + \frac{1}{2L_2} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2))\|_F^2. \end{aligned} \quad (\text{C.17})$$

Therefore, combining (C.16) and (C.17), we obtain

$$\begin{aligned} \langle \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_1 - \mathbf{S}_2 \rangle &\geq \frac{\mu_2}{2} \|\mathbf{S}_1 - \mathbf{S}_2\|_F^2 \\ &\quad + \frac{1}{2L_2} \|\mathcal{P}_\Omega(\nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \nabla_{\mathbf{S}}\mathcal{L}_n(\mathbf{X} + \mathbf{S}_2))\|_F^2, \end{aligned}$$

which completes the proof.  $\square$

## D Proofs of Auxiliary Lemmas in Appendix C

To begin with, we introduce the following notations for simplicity. Consider  $\mathbf{Z} \in \mathbb{R}^{(d_1+d_2) \times r}$ , for  $\mathbf{U} \in \mathbb{R}^{d_1 \times r}$  and  $\mathbf{V} \in \mathbb{R}^{d_2 \times r}$ , and  $\mathbf{X} = \mathbf{UV}^\top$ , we let  $\mathbf{Z} = [\mathbf{U}; \mathbf{V}]$ . Let  $\mathbf{R} = \operatorname{argmin}_{\tilde{\mathbf{R}} \in \mathbb{Q}^r} \|\mathbf{Z} - \mathbf{Z}^* \tilde{\mathbf{R}}\|_F$  be the optimal rotation regarding to  $\mathbf{Z}$ , and  $\mathbf{H} = \mathbf{Z} - \mathbf{Z}^* \mathbf{R} = [\mathbf{H}_U; \mathbf{H}_V]$  with  $\mathbf{H}_U \in \mathbb{R}^{d_1 \times r}$  and  $\mathbf{H}_V \in \mathbb{R}^{d_2 \times r}$ .

Besides, we introduce the following projection metrics, which are essential for proving the following lemmas. Denote by  $\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2, \bar{\mathbf{U}}_3$  the left singular matrices of  $\mathbf{X}, \mathbf{U}, \mathbf{H}_U$  respectively. Let  $\tilde{\mathbf{U}}$  be the matrix spanned by the column vectors of  $\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2$  and  $\bar{\mathbf{U}}_3$ , i.e.,

$$\operatorname{col}(\tilde{\mathbf{U}}) = \operatorname{span}\{\bar{\mathbf{U}}_1, \bar{\mathbf{U}}_2, \bar{\mathbf{U}}_3\} = \operatorname{col}(\bar{\mathbf{U}}_1) + \operatorname{col}(\bar{\mathbf{U}}_2) + \operatorname{col}(\bar{\mathbf{U}}_3). \quad (\text{D.1})$$

It is easy to show that  $\tilde{\mathbf{U}}$  is an orthonormal matrix with at most  $3r$  columns. Here, the sum of two subspaces is defined as  $\mathbf{U}_1 + \mathbf{U}_2 = \{\mathbf{u}_1 + \mathbf{u}_2 \mid \mathbf{u}_1 \in \mathbf{U}_1, \mathbf{u}_2 \in \mathbf{U}_2\}$ . Similarly, denote by  $\bar{\mathbf{V}}_1, \bar{\mathbf{V}}_2, \bar{\mathbf{V}}_3$  the right singular matrices of  $\mathbf{X}, \mathbf{V}, \mathbf{H}_V$  respectively. Again, let  $\tilde{\mathbf{V}}$  be the matrix spanned by the column of  $\bar{\mathbf{V}}_1, \bar{\mathbf{V}}_2$  and  $\bar{\mathbf{V}}_3$ , i.e.,

$$\operatorname{col}(\tilde{\mathbf{V}}) = \operatorname{span}\{\bar{\mathbf{V}}_1, \bar{\mathbf{V}}_2, \bar{\mathbf{V}}_3\} = \operatorname{col}(\bar{\mathbf{V}}_1) + \operatorname{col}(\bar{\mathbf{V}}_2) + \operatorname{col}(\bar{\mathbf{V}}_3), \quad (\text{D.2})$$

where the rank of  $\tilde{\mathbf{V}}$  is at most  $3r$ .

### D.1 Proof of Lemma C.1

*Proof.* Recall  $\mathbf{Z} = [\mathbf{U}; \mathbf{V}]$ . We denote  $\tilde{\mathbf{Z}} = [\mathbf{U}; -\mathbf{V}] \in \mathbb{R}^{(d_1+d_2) \times r}$ , then we can rewrite the regularization term  $\|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2$  as  $\|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2$  and its gradient with respect to  $\mathbf{Z}$  as  $\nabla_{\mathbf{Z}}(\|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2) = 4\tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top \mathbf{Z}$ . According to the formula of  $\nabla \tilde{F}_n(\mathbf{Z}, \mathbf{S})$  in (B.2), we have

$$\langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S}), \mathbf{H} \rangle = \underbrace{\langle \nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{UV}^\top + \mathbf{S}), \mathbf{H}_U \rangle + \langle \nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{UV}^\top + \mathbf{S}), \mathbf{H}_V \rangle}_{I_1} + \frac{1}{2} \underbrace{\langle \tilde{\mathbf{Z}} \tilde{\mathbf{Z}}^\top \mathbf{Z}, \mathbf{H} \rangle}_{I_2}, \quad (\text{D.3})$$

where  $\tilde{\mathbf{Z}} = [\mathbf{U}; -\mathbf{V}]$ , and  $\mathbf{H}_U, \mathbf{H}_V$  denote the top  $d_1 \times r$  and bottom  $d_2 \times r$  submatrices of  $\mathbf{H}$  respectively. Note that  $\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{UV}^\top + \mathbf{S}) = \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) \mathbf{V}$ , and  $\nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{UV}^\top + \mathbf{S}) = [\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S})]^\top \mathbf{U}$ . Consider the term  $I_1$  in (D.3) first, we have

$$\begin{aligned} I_1 &= \langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}), \mathbf{UV}^\top - \mathbf{U}^* \mathbf{V}^{*\top} + \mathbf{H}_U \mathbf{H}_V^\top \rangle \\ &= \underbrace{\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top \rangle}_{I_{11}} + \underbrace{\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top \rangle}_{I_{12}} \\ &\quad + \underbrace{\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}), \mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top \rangle}_{I_{13}}. \end{aligned} \quad (\text{D.4})$$

In the following discussions, we are going to bound  $I_{11}, I_{12}$  and  $I_{13}$  respectively. For the first term  $I_{11}$  in (D.4), we have

$$\begin{aligned} |I_{11}| &\leq \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2 \cdot (\|\mathbf{X} - \mathbf{X}^*\|_* + \|\mathbf{H}_U \mathbf{H}_V^\top\|_*) \\ &\leq \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2 \cdot (\sqrt{2r} \|\mathbf{X} - \mathbf{X}^*\|_F + \sqrt{r} \|\mathbf{H}_U \mathbf{H}_V\|_F) \\ &\leq \frac{\mu_1}{16} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{L_1}{16} \|\mathbf{H}\|_F^4 + \left( \frac{8r}{\mu_1} + \frac{r}{L_1} \right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2, \end{aligned} \quad (\text{D.5})$$

where the first inequality holds because of Von Neumann trace inequality, the second inequality is due to  $\mathbf{X} - \mathbf{X}^*$  has rank at most  $2r$  and  $\mathbf{H}_U \mathbf{H}_V^\top$  has rank at most  $r$ , and the last inequality holds because  $\|\mathbf{H}_U \mathbf{H}_V^\top\|_F \leq \|\mathbf{H}_U\| \cdot \|\mathbf{H}_V\|_F \leq \|\mathbf{H}\|_F^2/2$  and  $2ab \leq ta^2 + b^2/t$ , for any  $t > 0$ . As for the second term  $I_{12}$  in (D.4), note that  $\mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top$  has rank at most  $3r$ , thus according to the structural Lipschitz gradient Condition 4.4, we have

$$\begin{aligned} |I_{12}| &\leq |\langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top \rangle| + K \|\mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \\ &\leq |\langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} - \mathbf{X}^* \rangle| + \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{H}_U \mathbf{H}_V^\top\|_F + K \|\mathbf{X} - \mathbf{X}^* + \mathbf{H}_U \mathbf{H}_V^\top\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \\ &\leq |\langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} - \mathbf{X}^* \rangle| + \frac{1+K}{2} \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{H}\|_F^2 + K \|\mathbf{X} - \mathbf{X}^*\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F, \end{aligned} \quad (\text{D.6})$$

where the second inequality follows from triangle inequality and the fact that  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_F$ , and the last inequality is due to triangle inequality and the fact that  $\|\mathbf{H}_U \mathbf{H}_V^\top\|_F \leq \|\mathbf{H}_U\|_F \cdot \|\mathbf{H}_V\|_F \leq \|\mathbf{H}\|_F^2/2$ . Therefore, it suffices to bound the first term  $|\langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} - \mathbf{X}^* \rangle|$ . Denote the support of  $\mathbf{S} - \mathbf{S}^*$  by  $\Omega$ , then according to our assumption,  $\Omega$  has at most  $\beta' + \beta$  fraction nonzero entries per row and column. By Lemma 14 in [56], we further obtain

$$\begin{aligned} |I_{12}| &\leq \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathcal{P}_\Omega(\mathbf{U}\mathbf{V}^\top - \mathbf{U}\mathbf{V}^*)\|_F + \frac{1+K}{2} \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{H}\|_F^2 + K \|\mathbf{X} - \mathbf{X}^*\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \\ &\leq \sqrt{18(\beta' + \beta)\alpha r \sigma_1} \|\mathbf{H}\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F + \frac{1+K}{2} \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{H}\|_F^2 + K \|\mathbf{X} - \mathbf{X}^*\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \\ &\leq \frac{\mu_1}{8} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \left(\frac{\mu_2}{2} + \frac{2K^2}{\mu_1}\right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 + \frac{18(\beta' + \beta)\alpha r \sigma_1}{\mu_2} \|\mathbf{H}\|_F^2 + \frac{(1+K)^2}{4\mu_2} \|\mathbf{H}\|_F^4, \end{aligned} \quad (\text{D.7})$$

where the first inequality holds because  $|\langle \mathbf{A}, \mathcal{P}_\Omega(\mathbf{B}) \rangle| \leq \|\mathcal{P}_\Omega(\mathbf{A})\|_F \cdot \|\mathbf{B}\|_F$ , and the second inequality is due to Lemma 14 in [56], and the last inequality holds because  $2ab \leq ta^2 + b^2/t$ , for any  $t > 0$ . Finally, we consider the last term  $I_{13}$  in (D.4). Recall the orthonormal projection matrices  $\tilde{\mathbf{U}}$  and  $\tilde{\mathbf{V}}$  in (D.1) and (D.2). According to Lemma B.2 in [50], we have

$$\begin{aligned} \langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}), \mathbf{X} - \mathbf{X}^* \rangle &\geq \frac{1}{4L_1} \|\tilde{\mathbf{U}}^\top (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\|_F^2 \\ &\quad + \frac{1}{4L_1} \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\|_F^2 + \frac{\mu_1}{2} \|\mathbf{X} - \mathbf{X}^*\|_F^2. \end{aligned} \quad (\text{D.8})$$

As for the remaining term in  $I_{13}$ , we have

$$\begin{aligned} |\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}), \mathbf{H}_U \mathbf{H}_V^\top \rangle| &= |\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}), \tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{H}_U \mathbf{H}_V^\top \rangle| \\ &\leq \frac{1}{2} \|\tilde{\mathbf{U}}^\top (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\|_F \cdot \|\mathbf{H}\|_F^2 \\ &\leq \frac{1}{2L_1} \|\tilde{\mathbf{U}}^\top (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\|_F^2 + \frac{L_1}{8} \|\mathbf{H}\|_F^4, \end{aligned} \quad (\text{D.9})$$

where the equality is due to the fact that  $\text{col}(\bar{\mathbf{U}}_3) \subseteq \text{col}(\tilde{\mathbf{U}})$ , where  $\bar{\mathbf{U}}_3$  is the left singular matrix of  $\mathbf{H}_U$ , which implies that  $\tilde{\mathbf{U}} \tilde{\mathbf{U}}^\top \mathbf{H}_U = \mathbf{H}_U$ , the first inequality holds because  $|\langle \mathbf{A}, \mathbf{B}\mathbf{C} \rangle| \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\mathbf{C}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_2 \cdot \|\mathbf{C}\|_F$  and  $\|\tilde{\mathbf{U}}\|_2 = 1$ , and the last inequality holds because  $2ab \leq ta^2 + b^2/t$ , for any  $t > 0$ . Similarly, we have

$$\begin{aligned} |\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}), \mathbf{H}_U \mathbf{H}_V^\top \rangle| &\leq \frac{1}{2} \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\|_F \cdot \|\mathbf{H}\|_F^2 \\ &\leq \frac{1}{2L_1} \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\|_F^2 + \frac{L_1}{8} \|\mathbf{H}\|_F^4. \end{aligned} \quad (\text{D.10})$$

Therefore, combining (D.8), (D.9) and (D.10), we obtain the lower bound of  $I_{13}$

$$I_{13} \geq \frac{\mu_1}{2} \|\mathbf{X} - \mathbf{X}^*\|_F^2 - \frac{L_1}{8} \|\mathbf{H}\|_F^4. \quad (\text{D.11})$$

Hence, combining (D.5), (D.7) and (D.11), we further obtain the lower bound of  $I_1$  in (D.3)

$$\begin{aligned} I_1 &\geq \frac{3\mu_1}{8} \|\mathbf{X} - \mathbf{X}^*\|_F^2 - \frac{18(\beta' + \beta)\alpha r \sigma_1}{\mu_2} \|\mathbf{H}\|_F^2 - \left(\frac{L_1}{8} + \frac{(1+K)^2}{4\mu_2}\right) \cdot \|\mathbf{H}\|_F^4 \\ &\quad - \left(\frac{\mu_2}{2} + \frac{2K^2}{\mu_1}\right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 - \left(\frac{8r}{\mu_1} + \frac{r}{L_1}\right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_F^2. \end{aligned} \quad (\text{D.12})$$

Besides, according to Lemma C.1 in [49], we obtain the following lower bound regarding  $I_2$  in (D.3)

$$I_2 \geq \frac{1}{2} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 - \frac{1}{2} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F \cdot \|\mathbf{H}\|_F^2 \geq \frac{1}{4} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 - \frac{1}{4} \|\mathbf{H}\|_F^4. \quad (\text{D.13})$$

Note that  $K \in (0, 1)$ , by plugging (D.12) and (D.13) into (D.3), we have

$$\begin{aligned} \langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S}), \mathbf{H} \rangle &\geq \frac{3\mu_1}{8} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{1}{8} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 - \frac{18(\beta' + \beta)\alpha r \sigma_1}{\mu_2} \|\mathbf{H}\|_F^2 - \left( \frac{L_1 + 1}{8} + \frac{1}{\mu_2} \right) \cdot \|\mathbf{H}\|_F^4 \\ &\quad - \left( \frac{\mu_2}{2} + \frac{2K^2}{\mu_1} \right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 - \left( \frac{8r}{\mu_1} + \frac{r}{L_1} \right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2. \end{aligned} \quad (\text{D.14})$$

Furthermore, denote  $\tilde{\mathbf{Z}}^* = [\mathbf{U}^*; -\mathbf{V}^*]$ , then we obtain the following result

$$\begin{aligned} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 &= \langle \mathbf{Z}\mathbf{Z}^\top - \mathbf{Z}^*\mathbf{Z}^{*\top}, \tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top - \tilde{\mathbf{Z}}^*\tilde{\mathbf{Z}}^{*\top} \rangle + \langle \mathbf{Z}^*\mathbf{Z}^{*\top}, \tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top \rangle + \langle \mathbf{Z}\mathbf{Z}^\top, \tilde{\mathbf{Z}}^*\tilde{\mathbf{Z}}^{*\top} \rangle \\ &\geq \langle \mathbf{Z}\mathbf{Z}^\top - \mathbf{Z}^*\mathbf{Z}^{*\top}, \tilde{\mathbf{Z}}\tilde{\mathbf{Z}}^\top - \tilde{\mathbf{Z}}^*\tilde{\mathbf{Z}}^{*\top} \rangle \\ &= \|\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top}\|_F^2 + \|\mathbf{V}\mathbf{V}^\top - \mathbf{V}^*\mathbf{V}^{*\top}\|_F^2 - 2\|\mathbf{X} - \mathbf{X}^*\|_F^2, \end{aligned} \quad (\text{D.15})$$

where the first equality follows from  $\tilde{\mathbf{Z}}^{*\top} \mathbf{Z}^* = 0$ , and the inequality is due to the fact that  $\langle \mathbf{A}\mathbf{A}^\top, \mathbf{B}\mathbf{B}^\top \rangle = \|\mathbf{A}^\top \mathbf{B}\|_F^2 \geq 0$ . Therefore, by (D.15), we obtain

$$4\|\mathbf{X} - \mathbf{X}^*\|_F^2 + \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 = \|\mathbf{Z}\mathbf{Z}^\top - \mathbf{Z}^*\mathbf{Z}^{*\top}\|_F^2 \geq 4(\sqrt{2} - 1)\sigma_r \|\mathbf{H}\|_F^2, \quad (\text{D.16})$$

where the inequality follows from Lemma 5.4 in [46] and the fact that  $\sigma_r^2(\mathbf{Z}^*) = 2\sigma_r(\mathbf{X}^*) = 2\sigma_r$ . Denote  $\mu'_1 = \min\{\mu_1, 2\}$ , then by plugging (D.16) into (D.14), we obtain

$$\begin{aligned} \langle \nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S}), \mathbf{H} \rangle &\geq \frac{\mu_1}{4} \|\mathbf{X} - \mathbf{X}^*\|_F^2 + \frac{1}{16} \|\tilde{\mathbf{Z}}^\top \mathbf{Z}\|_F^2 + \left( \frac{\mu'_1}{20} \sigma_r - \frac{18(\beta' + \beta)\alpha r \sigma_1}{\mu_2} \right) \cdot \|\mathbf{H}\|_F^2 \\ &\quad - \left( \frac{L_1 + 1}{8} + \frac{1}{\mu_2} \right) \cdot \|\mathbf{H}\|_F^4 - \left( \frac{\mu_2}{2} + \frac{2K^2}{\mu_1} \right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 - \left( \frac{8r}{\mu_1} + \frac{r}{L_1} \right) \cdot \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2, \end{aligned}$$

which completes the proof.  $\square$

## D.2 Proof of Lemma C.2

*Proof.* According to the formula of  $\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S})$  in (B.2), we have

$$\|\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S})\|_F^2 \leq 2\|\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S})\|_F^2 + 2\|\nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S})\|_F^2 + \|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2 \cdot \|\mathbf{Z}\|_2^2, \quad (\text{D.17})$$

where the inequality follows from the fact that  $\|\mathbf{A} + \mathbf{B}\|_F^2 \leq 2\|\mathbf{A}\|_F^2 + 2\|\mathbf{B}\|_F^2$ ,  $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_2 \cdot \|\mathbf{B}\|_F$ , and  $\max\{\|\mathbf{U}\|_2, \|\mathbf{V}\|_2\} \leq \|\mathbf{Z}\|_2$ . Consider the first term  $\|\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S})\|_F^2$ . Denote  $\mathbf{X} = \mathbf{U}\mathbf{V}^\top$ , then we have

$$\begin{aligned} \|\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S})\|_F^2 &\leq 3 \underbrace{\|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\mathbf{V}\|_F^2}_{I_1} \\ &\quad + 3 \underbrace{\|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*))\mathbf{V}\|_F^2}_{I_2} + 3 \underbrace{\|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\mathbf{V}\|_F^2}_{I_3}, \end{aligned} \quad (\text{D.18})$$

where the inequality holds because  $\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U}\mathbf{V}^\top + \mathbf{S}) = \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S})\mathbf{V}$  and  $\|\mathbf{A} + \mathbf{B} + \mathbf{C}\|_F^2 \leq 3(\|\mathbf{A}\|_F^2 + \|\mathbf{B}\|_F^2 + \|\mathbf{C}\|_F^2)$ . In the following discussion, we are going to upper bound  $I_1, I_2$  and  $I_3$  separately. As for  $I_1$ , according to the orthonormal projection matrix  $\tilde{\mathbf{V}}$  defined in (D.2), we have

$$\begin{aligned} I_1 &= \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\tilde{\mathbf{V}}^\top \mathbf{V}\|_F^2 \\ &\leq \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\|_F^2 \cdot \|\tilde{\mathbf{V}}^\top \mathbf{V}\|_2^2 \\ &\leq \|(\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}))\tilde{\mathbf{V}}\|_F^2 \cdot \|\mathbf{V}\|_2^2, \end{aligned} \quad (\text{D.19})$$

where the equality holds because  $\text{col}(\mathbf{V}) \subseteq \text{col}(\tilde{\mathbf{V}})$ , which implies that  $\tilde{\mathbf{V}}\tilde{\mathbf{V}}^\top \mathbf{V} = \mathbf{V}$ , the first inequality is due to the fact that  $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_2$ , and the last inequality holds because  $\|\mathbf{A}\mathbf{B}\|_2 \leq \|\mathbf{A}\|_2 \cdot \|\mathbf{B}\|_2$  and the

fact that  $\tilde{\mathbf{V}}$  is orthonormal. Moreover, consider the second term  $I_2$  in (D.18). According to the definition of Frobenius norm, we have

$$\begin{aligned} \left\| (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)) \mathbf{V} \right\|_F &= \sup_{\|\mathbf{W}\|_F \leq 1} \langle (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)) \mathbf{V}, \mathbf{W} \rangle \\ &\leq \sup_{\|\mathbf{W}\|_F \leq 1} (1 + K) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{W} \mathbf{V}^\top\|_F \\ &\leq (1 + K) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{V}\|_2, \end{aligned} \quad (\text{D.20})$$

where the first inequality follows from the structural Lipschitz gradient Condition 4.4 and the fact that  $|\langle \mathbf{A}, \mathbf{B} \rangle| \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_F$ , and the second one holds because  $\|\mathbf{A} \mathbf{B}\|_F \leq \|\mathbf{A}\|_F \cdot \|\mathbf{B}\|_2$  and  $\|\mathbf{W}\|_F \leq 1$ . Finally, consider the last term  $I_3$  in (D.18), we have

$$I_3 \leq \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{V}\|_F^2 \leq r \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{V}\|_2^2. \quad (\text{D.21})$$

Thus, combining (D.19), (D.20) and (D.21), we obtain

$$\begin{aligned} \|\nabla_{\mathbf{U}} \mathcal{L}_n(\mathbf{U} \mathbf{V}^\top + \mathbf{S})\|_F^2 &\leq 3 \left\| (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S})) \tilde{\mathbf{V}} \right\|_F^2 \cdot \|\mathbf{V}\|_2^2 \\ &\quad + 3(1 + K)^2 \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 \cdot \|\mathbf{V}\|_2^2 + 3r \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{V}\|_2^2. \end{aligned} \quad (\text{D.22})$$

As for the second term  $\|\nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{U} \mathbf{V}^\top + \mathbf{S})\|_F^2$  in (D.17), based on similar techniques, we obtain

$$\begin{aligned} \|\nabla_{\mathbf{V}} \mathcal{L}_n(\mathbf{U} \mathbf{V}^\top + \mathbf{S})\|_F^2 &\leq 3 \left\| \tilde{\mathbf{U}}^\top (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S})) \right\|_F^2 \cdot \|\mathbf{U}\|_2^2 \\ &\quad + 3(1 + K)^2 \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 \cdot \|\mathbf{U}\|_2^2 + 3r \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{U}\|_2^2, \end{aligned} \quad (\text{D.23})$$

where  $\tilde{\mathbf{U}}$  is an orthonormal matrix defined in (D.1). According to Lemma C.1 in [50] and Condition 4.2, we have

$$\begin{aligned} I &= \left\| (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S})) \tilde{\mathbf{V}} \right\|_F^2 + \left\| \tilde{\mathbf{U}}^\top (\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S})) \right\|_F^2 \\ &\leq 4L_1 (\mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \mathcal{L}_n(\mathbf{X} + \mathbf{S})) - \langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}), \mathbf{X}^* - \mathbf{X} \rangle \leq 2L_1^2 \cdot \|\mathbf{X} - \mathbf{X}^*\|_F^2. \end{aligned} \quad (\text{D.24})$$

Therefore, plugging (D.22), (D.23) and (D.24) into (D.17), we obtain

$$\begin{aligned} \|\nabla_{\mathbf{Z}} \tilde{F}_n(\mathbf{Z}, \mathbf{S})\|_F^2 &\leq \left( 12L_1^2 \|\mathbf{X} - \mathbf{X}^*\|_F^2 + 12(1 + K)^2 \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 + \|\mathbf{U}^\top \mathbf{U} - \mathbf{V}^\top \mathbf{V}\|_F^2 \right) \cdot \|\mathbf{Z}\|_2^2 \\ &\quad + 12r \|\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*)\|_2^2 \cdot \|\mathbf{Z}\|_2^2, \end{aligned}$$

where the inequality holds because  $\max\{\|\mathbf{U}\|_2, \|\mathbf{V}\|_2\} \leq \|\mathbf{Z}\|_2$ . Thus, we finish the proof.  $\square$

### D.3 Proof of Lemma C.3

*Proof.* Denote the support of  $\mathbf{S}^*$  and  $\mathcal{T}_{\gamma\beta}(\mathbf{S})$  by  $\Omega^*$  and  $\Omega$  respectively. According to the definition of the truncation operator  $\mathcal{T}_\alpha$ , we have

$$\begin{aligned} \|\mathcal{T}_{\gamma\beta}(\mathbf{S}) - \mathbf{S}^*\|_F^2 &= \|\mathcal{P}_\Omega(\mathcal{T}_{\gamma\beta}(\mathbf{S}) - \mathbf{S}^*)\|_F^2 + \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathcal{T}_{\gamma\beta}(\mathbf{S}) - \mathbf{S}^*)\|_F^2 \\ &= \|\mathcal{P}_\Omega(\mathbf{S} - \mathbf{S}^*)\|_F^2 + \|\mathcal{P}_{\Omega^* \setminus \Omega}(-\mathbf{S}^*)\|_F^2, \end{aligned} \quad (\text{D.25})$$

where the second inequality holds because  $[\mathcal{T}_{\gamma\beta}(\mathbf{S})]_{i,j} = S_{i,j}$  if  $(i, j) \in \Omega$ , and  $[\mathcal{T}_{\gamma\beta}(\mathbf{S})]_{i,j} = 0$  otherwise. For any  $(i, j) \in \Omega^* \setminus \Omega$ , we claim

$$|(\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*)_{i,j}| \leq \max \left\{ \underbrace{|(\mathbf{S} - \mathbf{S}^*)_{i,*}^{(\gamma\beta d_2 - \beta d_2)}|}_{I_1}, \underbrace{|(\mathbf{S} - \mathbf{S}^*)_{*,j}^{(\gamma\beta d_1 - \beta d_1)}|}_{I_2} \right\}, \quad (\text{D.26})$$

where we denote the  $k$ -th largest element in magnitude of  $(\mathbf{S} - \mathbf{S}^*)_{i,*}$  by  $(\mathbf{S} - \mathbf{S}^*)_{i,*}^{(k)}$ , and the  $k$ -th largest element in magnitude of  $(\mathbf{S} - \mathbf{S}^*)_{*,j}$  by  $(\mathbf{S} - \mathbf{S}^*)_{*,j}^{(k)}$ . In the following discussion, we are going to prove claim (D.26) by

contradiction. Suppose  $|\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*|_{i,j} = |S_{i,j}| > \max\{I_1, I_2\}$ , where  $(i, j) \in \Omega^* \setminus \Omega$ . Noticing  $\mathbf{S}^*$  has at most  $\beta$ -fraction nonzero entries per row and column, we have

$$I_1 \geq |\mathbf{S}_{i,*}^{(\gamma\beta d_2 - \beta d_2)}| \quad \text{and} \quad I_2 \geq |\mathbf{S}_{*,j}^{(\gamma\beta d_1 - \beta d_1)}|.$$

Thus we have  $|S_{i,j}| \geq \max\{|\mathbf{S}_{i,*}^{(\gamma\beta d_2 - \beta d_2)}|, |\mathbf{S}_{*,j}^{(\gamma\beta d_1 - \beta d_1)}|\}$ , which contradicts with the fact that  $(i, j) \in \Omega^* \setminus \Omega$ . Therefore, based on (D.26), we obtain

$$\begin{aligned} \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*)\|_F^2 &= \sum_{(i,j) \in \Omega^* \setminus \Omega} |(\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*)_{i,j}|^2 \\ &\leq \sum_{(i,j) \in \Omega^* \setminus \Omega} \frac{\|(\mathbf{S} - \mathbf{S}^*)_{i,*}\|_2^2}{(\gamma - 1)\beta d_2} + \sum_{(i,j) \in \Omega^* \setminus \Omega} \frac{\|(\mathbf{S} - \mathbf{S}^*)_{*,j}\|_2^2}{(\gamma - 1)\beta d_1} \\ &\leq \frac{2}{\gamma - 1} \|\mathbf{S} - \mathbf{S}^*\|_F^2, \end{aligned} \tag{D.27}$$

where the first inequality is due to (D.26), and the second inequality holds because for each row and column of  $\Omega^*$ , it has at most  $\beta$ -fraction nonzero elements. Thus we obtain

$$\begin{aligned} \|\mathcal{P}_{\Omega^* \setminus \Omega}(-\mathbf{S}^*)\|_F^2 &= \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^*) - \mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*)\|_F^2 \\ &\leq (1 + c) \cdot \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^*)\|_F^2 + \left(1 + \frac{1}{c}\right) \cdot \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^* + \mathbf{S}^*)\|_F^2 \\ &\leq (1 + c) \cdot \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^*)\|_F^2 + \frac{c+1}{c} \cdot \frac{2}{\gamma - 1} \|\mathbf{S} - \mathbf{S}^*\|_F^2, \end{aligned} \tag{D.28}$$

where the second inequality holds because  $\|\mathbf{A} + \mathbf{B}\|_F^2 \leq (1 + c) \cdot \|\mathbf{A}\|_F^2 + (1 + 1/c) \cdot \|\mathbf{B}\|_F^2$ , for any  $c > 0$ , and the second inequality is due to (D.27). Therefore, plugging in (D.28) into (D.25), we have

$$\begin{aligned} \|\mathcal{T}_{\gamma\beta}(\mathbf{S}) - \mathbf{S}^*\|_F^2 &\leq \|\mathcal{P}_{\Omega}(\mathbf{S} - \mathbf{S}^*)\|_F^2 + (1 + c) \cdot \|\mathcal{P}_{\Omega^* \setminus \Omega}(\mathbf{S} - \mathbf{S}^*)\|_F^2 + \frac{c+1}{c} \cdot \frac{2}{\gamma - 1} \|\mathbf{S} - \mathbf{S}^*\|_F^2 \\ &\leq \left(1 + c + \frac{2(c+1)}{c(\gamma - 1)}\right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2 \\ &= \left(1 + \frac{2}{\gamma - 1} + 2\sqrt{\frac{2}{\gamma - 1}}\right) \cdot \|\mathbf{S} - \mathbf{S}^*\|_F^2, \end{aligned}$$

where we set  $c = \sqrt{(\gamma - 1)}/2$  in the last step. Thus we complete the proof.  $\square$

## E Proofs of Specific Models

In this section, we provide proofs for specific models. In the following discussions, we let  $d = \max\{d_1, d_2\}$ .

### E.1 Proofs of Robust Matrix Sensing

For matrix sensing, recall that we have the linear measurement operator  $\mathcal{A}$  with each sensing matrix  $\mathbf{A}_i$  sampled independently from  $\Sigma$ -Gaussian ensemble, where  $\text{vec}(\mathbf{A}_i) \sim N(0, \Sigma)$ . In particular, we consider  $\Sigma = \mathbf{I}$  and here  $\text{vec}(\mathbf{A}_i)$  denotes the vectorization of matrix  $\mathbf{A}_i$ . In order to prove the results for matrix sensing, we first lay out several lemmas, which are essential to prove the results for robust matrix sensing. The first lemma is useful to verify the restricted strong convexity and smoothness conditions in Condition 4.2.

**Lemma E.1.** [37] Suppose we have the linear measurement operator  $\mathcal{A}$  with each sensing matrix  $\mathbf{A}_i$  sampled independently from  $\mathbf{I}$ -Gaussian ensemble, then there exists constants  $c_0, c_1$  such that for all  $\Delta \in \mathbb{R}^{d_1 \times d_2}$  with rank at most  $2\tilde{r}$ , it holds with probability at least  $1 - \exp(-c_0 n)$  that

$$\left| \frac{\|\mathcal{A}(\Delta)\|_2^2}{n} - \frac{1}{2} \|\text{vec}(\Delta)\|_2^2 \right| \leq c_1 \frac{\tilde{r}d}{n} \|\Delta\|_F^2. \tag{E.1}$$

The second lemma is useful to verify the restricted strong convexity and smoothness conditions in Condition 4.3.

**Lemma E.2.** [43] For any random matrix  $\mathbf{A} \in \mathbb{R}^{n \times d_1 d_2}$ , which is drawn from the  $\Sigma$ -Gaussian ensemble, and the cardinalities of all vector  $\mathbf{s} \in \mathbb{R}^{d_1 d_2}$  satisfy  $|\mathbf{s}| \leq \tilde{s}$ . If we have sample size  $n \geq c_2 \tilde{s} \log d$ , then the following inequality holds with probability at least  $1 - c_3 \exp(-c_4 n)$

$$c_5 \|\Sigma^{1/2} \mathbf{s}\|_2^2 - c_6 \frac{\log d}{n} \|\mathbf{s}\|_1 \leq \frac{\|\mathbf{A} \mathbf{s}\|_2^2}{n} \leq c_7 \|\Sigma^{1/2} \mathbf{s}\|_2^2 + c_8 \frac{\log d}{n} \|\mathbf{s}\|_1,$$

where  $\{c_i\}_{i=2}^8$  are universal constants.

The next lemma verifies the structural Lipschitz gradient condition in Condition 4.4.

**Lemma E.3.** Consider robust matrix sensing with objective loss function defined in section 4.2. There exist constants  $C_0, C_1$  such that the following inequality holds with probability at least  $1 - \exp(-C_0 d)$

$$\begin{aligned} |\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{X} \rangle - \langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle| &\leq K \|\mathbf{X}\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F, \\ |\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{S} \rangle - \langle \mathbf{X} - \mathbf{X}^*, \mathbf{S} \rangle| &\leq K \|\mathbf{X} - \mathbf{X}^*\|_F \cdot \|\mathbf{S}\|_F, \end{aligned}$$

for all low-rank matrices  $\mathbf{X}, \mathbf{X}^*$  with rank at most  $\tilde{r}$  and all sparse matrices  $\mathbf{S}, \mathbf{S}^*$  with sparsity at most  $\tilde{s}$ , where  $\tilde{r}, \tilde{s}$  are defined in Condition 4.2, and the structural Lipschitz gradient parameter  $K = C_1 \sqrt{(rd + s) \log d/n}$ .

The last lemma verifies the condition in Condition 4.5 for robust matrix sensing.

**Lemma E.4.** Consider robust matrix sensing, suppose each sensing matrix  $\mathbf{A}_i$  is sampled independently from I-Gaussian ensemble and each element of noise vector  $\epsilon$  follows i.i.d. sub-Gaussian distribution with parameter  $\nu$ . Then we have the following inequalities hold with probability at least  $1 - C_2/d$  in terms of spectral norm and infinity norm respectively

$$\left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{A}_i \right\|_2 \leq C_3 \nu \sqrt{\frac{d}{n}} \quad \text{and} \quad \left\| \frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{A}_i \right\|_{\infty, \infty} \leq C_4 \nu \sqrt{\frac{\log d}{n}},$$

where  $C_2, C_3, C_4$  are universal constants.

Now, we are ready to prove Corollary 4.11.

*Proof of Corollary 4.11.* In order to prove Corollary 4.11, we only need to verify the restricted strong convex and smoothness conditions in Conditions 4.2 and 4.3, the structural Lipschitz gradient condition in Condition 4.4, and the condition in Condition 4.5.

Recall that we have the sample loss function for robust matrix sensing as  $\mathcal{L}_n(\mathbf{X} + \mathbf{S}) := \|\mathbf{y} - \mathcal{A}_n(\mathbf{X} + \mathbf{S})\|_2^2 / (2n)$ . Therefore, for all given sparse matrices  $\mathbf{S}$ , we have the following holds for all matrices  $\mathbf{X}_1, \mathbf{X}_2$  with rank at most  $\tilde{r}$

$$\mathcal{L}_n(\mathbf{X}_1 + \mathbf{S}) - \mathcal{L}_n(\mathbf{X}_2 + \mathbf{S}) - \langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}_2 + \mathbf{S}), \mathbf{X}_2 - \mathbf{X}_1 \rangle = \frac{\|\mathcal{A}(\Delta)\|_2^2}{n},$$

where  $\Delta = \mathbf{X}_2 - \mathbf{X}_1$ . According to Lemma E.1, if we have  $n > c'_1 \tilde{r} d$ , where  $c'$  is some constants. Then, with probability at least  $1 - \exp(-c_0 n)$ , we have the restricted strong convexity and smoothness conditions in Condition 4.2 hold with parameter  $\mu_1 = 4/9$  and  $L_1 = 5/9$ . In addition, for all given low-rank matrices  $\mathbf{X}$ , we have the following holds for all matrices  $\mathbf{S}_1, \mathbf{S}_2$  with sparsity at most  $\tilde{s}$

$$\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2) - \langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_2 - \mathbf{S}_1 \rangle = \frac{\|\mathcal{A}(\Delta)\|_2^2}{n},$$

where  $\Delta = \mathbf{S}_2 - \mathbf{S}_1$ . Furthermore, we can obtain  $\|\mathcal{A}(\Delta)\|_2^2 = \|\mathbf{A} \delta\|_2^2$ , where we have  $\mathbf{A} \in \mathbb{R}^{n \times d_1 d_2}$  with each row  $\mathbf{A}_{i*} = \text{vec}(\mathbf{A}_i)$ , and  $\delta = \text{vec}(\Delta)$ . Therefore, according to Lemma E.2, we have

$$c_1 \|\delta\|_2^2 - c_2 \frac{\log d}{n} \|\delta\|_1 \leq \frac{\|\mathbf{A} \delta\|_2^2}{n} \leq c_3 \|\delta\|_2^2 + c_4 \frac{\log d}{n} \|\delta\|_1.$$



Thus provided that  $n > c_5 \tilde{s} \log d$ , with probability at least  $1 - c_3 \exp(-c_4 n)$ , the restricted strong convexity and smoothness conditions in Condition 4.3 hold with parameters  $\mu_2 = 4/9$  and  $L_2 = 5/9$ .

Next, according to Lemma E.3, with probability at least  $1 - \exp(-C_0 d)$ , we can establish the structural Lipschitz gradient condition in Condition 4.5 with parameter  $K = C_1 \sqrt{(rd + s) \log d/n}$ .

Finally, we will verify the condition in Condition 4.5. By the definition of the objective loss function for robust matrix sensing, we have  $\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*) = \sum_{i=1}^n \epsilon_i \mathbf{A}_i/n$  and  $\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*) = \sum_{i=1}^n \epsilon_i \mathbf{A}_i/n$ . Therefore, according to Lemma E.4, with probability at least  $1 - C_2/d$ , we can establish the condition in Condition 4.5 with parameters  $\epsilon_1 = C_3 \nu \sqrt{d/n}$  and  $\epsilon_2 = C_4 \nu \sqrt{\log d/n}$ . This completes the proof.  $\square$

## E.2 Proofs of Robust PCA

Note that since robust PCA under fully observed model is a special case of robust PCA under partially observed model, thus we just lay out the proofs of robust PCA under partially observed model. To prove the results of partially observed robust PCA, we need the following lemmas, which are essential to establish the restricted strong convexity and smoothness conditions in Conditions 4.2 and 4.3. Note that the following lemmas only work for robust PCA under noisy observation model.

**Lemma E.5.** [38] There exist universal constants  $\{c_i\}_{i=1}^4$  such that if the number of observations  $n \geq c_1 r d \log d$ , and the following condition is satisfied for all  $\Delta \in \mathbb{R}^{d_1 \times d_2}$

$$\sqrt{\frac{d_1 d_2}{r}} \frac{\|\Delta\|_{\infty, \infty}}{\|\Delta\|_F} \cdot \frac{\|\Delta\|_*}{\|\Delta\|_F} \leq \frac{1}{c_2} \sqrt{n/(d \log d)}, \quad (\text{E.2})$$

we have, with probability at least  $1 - c_3/d$ , that the following holds

$$\left| \frac{\|\mathcal{A}(\Delta)\|_2}{\sqrt{n}} - \frac{\|\Delta\|_F}{\sqrt{d_1 d_2}} \right| \leq \frac{1}{10} \frac{\|\Delta\|_F}{\sqrt{d_1 d_2}} \left( 1 + \frac{c_4 \sqrt{d_1 d_2} \|\Delta\|_{\infty, \infty}}{\sqrt{n} \|\Delta\|_F} \right).$$

**Lemma E.6.** There exist universal constants  $\{c_i\}_{i=1}^5$  such that as long as  $n \geq c_1 \log d$ , we have with probability at least  $1 - c_2 \exp(-c_3 \log d)$  that

$$\left| \frac{\|\mathcal{A}(\Delta)\|_2}{\sqrt{n}} - \frac{\|\Delta\|_F}{\sqrt{d_1 d_2}} \right| \leq \frac{1}{2} \frac{\|\Delta\|_F}{\sqrt{d_1 d_2}} + \frac{c_5 \|\Delta\|_{\infty, \infty}}{\sqrt{n}} \quad \text{for all } \Delta \in \mathcal{C}(n), \quad (\text{E.3})$$

where we have the set  $\mathcal{C}(n)$  as follows

$$\mathcal{C}(n) = \left\{ \Delta \in \mathbb{R}^{d_1 \times d_2} \mid \frac{\|\Delta\|_{1,1}}{\|\Delta\|_F} \cdot \frac{\|\Delta\|_{\infty, \infty}}{\|\Delta\|_F} \leq c_4 \sqrt{\frac{n}{d_1 d_2 \log d}} \right\}.$$

The next lemma verifies the structural Lipschitz gradient condition in Condition 4.4.

**Lemma E.7.** Consider partially observed robust PCA with objective loss function defined in section 4.2. There exist constants  $C_0, C_1$  such that the following inequality holds with probability at least  $1 - \exp(-C_0 d)$

$$\begin{aligned} |\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{X} \rangle - \langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle| &\leq K \|\mathbf{X}\|_F \cdot \|\mathbf{S} - \mathbf{S}^*\|_F, \\ |\langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}^*) - \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{S} \rangle - \langle \mathbf{X} - \mathbf{X}^*, \mathbf{S} \rangle| &\leq K \|\mathbf{X} - \mathbf{X}^*\|_F \cdot \|\mathbf{S}\|_F, \end{aligned}$$

for all low-rank matrices  $\mathbf{X}, \mathbf{X}^*$  with rank at most  $\tilde{r}$  and all sparse matrices  $\mathbf{S}, \mathbf{S}^*$  with sparsity at most  $\tilde{s}$ , where  $\tilde{r}, \tilde{s}$  are defined in Condition 4.2, and the structural Lipschitz gradient parameter  $K = C_1 \sqrt{(rd + s) \log d/n}$ .

The last lemma verifies the condition in Condition 4.5 for partially observed robust PCA.

**Lemma E.8.** Consider partially observed robust PCA. If  $\mathbf{A}_{jk} = \mathbf{e}_j \mathbf{e}_k^\top$  is uniformly distributed on  $\Omega$ , then for i.i.d. zero mean random variables  $\epsilon_{jk}$  with variance  $\nu^2$ , we have the following inequalities hold with probability at least  $1 - C_2/d$  in terms of spectral norm and infinity norm respectively

$$\left\| \frac{1}{p} \sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk} \right\|_2 \leq C_3 \nu \sqrt{\frac{d \log d}{p}} \quad \text{and} \quad \left\| \frac{1}{p} \sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk} \right\|_{\infty, \infty} \leq C_4 \nu \sqrt{\frac{\log d}{p}},$$

where  $C_2, C_3, C_4$  are universal constants, and  $p = n/(d_1 d_2)$ .

Now, we are ready to prove Corollary 4.15.

*Proof of Corollary 4.15.* To prove Corollary 4.15, we need to verify the restricted strong convexity and smoothness conditions in Conditions 4.2 and 4.3, the structural Lipschitz gradient condition in Condition 4.4, and the condition in Condition 4.5.

In the following discussion, we let  $\mathbf{A}_{jk} = \mathbf{e}_j \mathbf{e}_k^\top$ , where  $\mathbf{e}_i, \mathbf{e}_j$  are basis vectors with  $d_1$  and  $d_2$  dimensions, and we let  $\mathcal{A}$  be the corresponding transformation operator. In addition, let the number of observations to be  $|\Omega| = n$ . Therefore, the objective loss function for robust PCA in 4.2 can be rewritten as

$$\mathcal{L}_n(\mathbf{X} + \mathbf{S}) := \frac{1}{2p} \sum_{(j,k) \in \Omega} (\langle \mathbf{A}_{jk}, \mathbf{X} + \mathbf{S} \rangle - Y_{jk})^2.$$

Therefore, for all given sparse matrices  $\mathbf{S}$ , we have the following holds for all matrices  $\mathbf{X}_1, \mathbf{X}_2$  satisfying incoherence condition with rank at most  $\tilde{r}$

$$\mathcal{L}_n(\mathbf{X}_1 + \mathbf{S}) - \mathcal{L}_n(\mathbf{X}_2 + \mathbf{S}) - \langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}_2 + \mathbf{S}), \mathbf{X}_2 - \mathbf{X}_1 \rangle = \frac{\|\mathcal{A}(\Delta)\|_2^2}{p},$$

where  $\Delta = \mathbf{X}_1 - \mathbf{X}_2$ , and  $p = n/(d_1 d_2)$ . Now, we are ready to prove the restricted strong convexity and smoothness conditions in Condition 4.2.

**Case 1:** If  $\Delta$  not satisfies condition (E.2), then we have

$$\begin{aligned} \|\Delta\|_F^2 &\leq C_0 (\sqrt{d_1 d_2} \|\Delta\|_\infty) \|\Delta\|_* \sqrt{\frac{d \log d}{nr}} \\ &\leq 2C_0 \alpha_1 \sqrt{d_1 d_2} \|\Delta\|_* \sqrt{\frac{d \log d}{nr}} \\ &\leq 2C_0 \alpha_1 \sqrt{2\tilde{r} d_1 d_2} \|\Delta\|_F \sqrt{\frac{d \log d}{nr}}, \end{aligned}$$

where  $\tilde{\alpha} = \alpha r / \sqrt{d_1 d_2}$  due to the incoherence condition of low rank matrices  $\mathbf{X}_1$  and  $\mathbf{X}_2$ , and the last inequality comes from  $\text{rank}(\Delta) \leq 2\tilde{r}$ . Thus, by the definition of  $\tilde{r}$ , we can obtain

$$\|\Delta\|_F^2 \leq C_1 \alpha^2 \sigma_1^2 \frac{r^2 d \log d}{n}. \quad (\text{E.4})$$

**Case 2:** If  $\Delta$  satisfies condition (E.2), then according to Lemma E.5, we have

$$\left| \frac{\|\mathcal{A}(\Delta)\|_2}{\sqrt{p}} - \|\Delta\|_F \right| \leq \frac{\|\Delta\|_F}{10} \left( 1 + \frac{C_2 \sqrt{d_1 d_2} \|\Delta\|_{\infty, \infty}}{\sqrt{n} \|\Delta\|_F} \right).$$

Thus if  $C_2 \sqrt{d_1 d_2} \|\Delta\|_{\infty, \infty} / (\sqrt{n} \|\Delta\|_F) \geq C_3$ , we have

$$\|\Delta\|_F^2 \leq C_4 \frac{\tilde{\alpha}^2}{p}.$$

Otherwise, if  $C_2 \sqrt{d_1 d_2} \|\Delta\|_{\infty, \infty} / (\sqrt{n} \|\Delta\|_F) \leq C_3$ , we have

$$\frac{8}{9} \|\Delta\|_F^2 \leq \frac{\|\mathcal{A}(\Delta)\|_2^2}{p} \leq \frac{10}{9} \|\Delta\|_F^2,$$

which gives us the restricted strong convexity and smoothness conditions in Condition 4.2 with parameters  $\mu_1 = 8/9, L_1 = 10/9$ .

Next, we prove the restricted strong convexity and smoothness conditions in Condition 4.3. For all given low-rank matrices  $\mathbf{X}$ , we have the following holds for all matrices  $\mathbf{S}_1, \mathbf{S}_2$  with at most  $\tilde{s}$  nonzero entries and infinity norm bound  $\alpha_1 / \sqrt{d_1 d_2}$

$$\mathcal{L}_n(\mathbf{X} + \mathbf{S}_1) - \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2) - \langle \nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X} + \mathbf{S}_2), \mathbf{S}_2 - \mathbf{S}_1 \rangle = \frac{\|\mathcal{A}(\Delta)\|_2^2}{p},$$

where  $\mathbf{\Delta} = \mathbf{S}_1 - \mathbf{S}_2$ , and  $p = n/(d_1 d_2)$ .

**Case 1:** If  $\mathbf{\Delta} \notin \mathcal{C}(n)$ , then we can get

$$\|\mathbf{\Delta}\|_F^2 \leq C_5 (\sqrt{d_1 d_2} \|\mathbf{\Delta}\|_{\infty, \infty}) \cdot \|\mathbf{\Delta}\|_{1,1} \sqrt{\frac{\log d}{n}} \leq 2C_5 \alpha_1 \|\mathbf{\Delta}\|_{1,1} \sqrt{\frac{\log d}{n}},$$

where the last inequality is due to the fact that  $\|\mathbf{\Delta}\|_{\infty} = \|\mathbf{S}_1 - \mathbf{S}_2\|_{\infty} \leq 2\alpha_1/\sqrt{d_1 d_2}$ . Therefore, we can obtain

$$\|\mathbf{\Delta}\|_F^2 \leq 2C_5 \sqrt{2\tilde{s}} \alpha_1 \|\mathbf{\Delta}\|_F \sqrt{\frac{\log d}{n}},$$

where the inequality holds because  $\mathbf{\Delta}$  has at most  $2\tilde{s}$  nonzero entries. Therefore, by the definition of  $\tilde{s}$ , we have

$$\|\mathbf{\Delta}\|_F^2 \leq C_6 \alpha_1^2 \frac{s \log d}{n}. \quad (\text{E.5})$$

**Case 2:** If  $\mathbf{\Delta} \in \mathcal{C}(n)$ , we have

$$\left| \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{p}} - \|\mathbf{\Delta}\|_F \right| \leq \frac{1}{2} \|\mathbf{\Delta}\|_F + \frac{c_5 \sqrt{d_1 d_2} \|\mathbf{\Delta}\|_{\infty, \infty}}{\sqrt{n} \|\mathbf{\Delta}\|_F},$$

If  $\sqrt{n} \|\mathbf{\Delta}\|_F \leq C_7 \sqrt{d_1 d_2} \|\mathbf{\Delta}\|_{\infty, \infty}$ , we can obtain  $\|\mathbf{\Delta}\|_F^2 \leq C_7' \alpha_1^2/n$ . Otherwise, if we have  $\sqrt{n} \|\mathbf{\Delta}\|_F \geq C_7 \sqrt{d_1 d_2} \|\mathbf{\Delta}\|_{\infty, \infty}$ , according to Lemma E.6, we obtain

$$\frac{8}{9} \|\mathbf{\Delta}\|_F^2 \leq \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2^2}{p} \leq \frac{10}{9} \|\mathbf{\Delta}\|_F^2,$$

which implies the restricted strong convexity and smoothness conditions in Condition 4.3 hold with parameters  $\mu_2 = 8/9, L_2 = 10/9$ .

Next, according to Lemma E.7, with probability at least  $1 - \exp(-C_0 d)$ , we can establish the structural Lipschitz gradient condition in Condition 4.5 with parameter  $K = C_1 \sqrt{(rd + s) \log d/n}$ .

Finally, we verify the condition in Condition 4.5. By the definition of the objective loss function for robust PCA, we have  $\nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*) = \sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk}/p$  and  $\nabla_{\mathbf{S}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*) = \sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk}/p$ , where  $\epsilon_{jk}$  are i.i.d. Gaussian variables with variance  $\nu^2/(d_1 d_2)$ . Therefore, according to Lemma E.8, with probability at least  $1 - C_8'/d$ , we have  $\|\sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk}/p\|_2^2 \leq C_8 \nu^2 d \log d/n$ . In addition, we have  $\|\sum_{j,k \in \Omega} \epsilon_{jk} \mathbf{A}_{jk}/p\|_{\infty, \infty}^2 \leq C_9 \nu^2 \log d/n$ . Furthermore, we have additional estimation error bounds (E.4) and (E.5) when we derive the restricted strong convexity and smoothness conditions. Therefore, we can establish the condition in Condition 4.5 with parameters  $\epsilon_1^2 = C_8 \max\{\alpha_1^2, \nu^2\} d/n$  and  $\epsilon_2^2 = C_9 \max\{\alpha_1^2, \nu^2\} \log d/n$ . This completes the proof.  $\square$

## F Proofs of Technical Lemmas in Appendix E

### F.1 Proof of Lemma E.3

*Proof.* In order to verify the structural Lipschitz gradient condition, we need to make use of the Bernstein-type inequality for sub-exponential random variables in [47] as well as the corresponding covering arguments for low-rank and sparse structures, respectively.

By the definition of the objective loss function of matrix sensing, we have

$$\langle \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}) - \nabla_{\mathbf{X}} \mathcal{L}_n(\mathbf{X}^* + \mathbf{S}^*), \mathbf{X} \rangle = \frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} \rangle = \frac{1}{n} \sum_{i=1}^n Y_i,$$

where  $Y_i = \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} \rangle$ . Note that  $\langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle, \langle \mathbf{A}_i, \mathbf{X} \rangle$  follow i.i.d. normal distribution  $N(0, \|\mathbf{S} - \mathbf{S}^*\|_F^2/n)$  and  $N(0, \|\mathbf{X}\|_F^2/n)$  respectively. Thus  $Y_i$  follows i.i.d. chi-square distribution which is also sub-exponential. Besides,  $\mathbb{E}(Y_i) = \langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle$ , and we have

$$\|Y_i - \mathbb{E}[Y_i]\|_{\psi_1} \leq 2\|Y_i\|_{\psi_1} \leq 2\|\langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle\|_{\psi_2} \cdot \|\langle \mathbf{A}_i, \mathbf{X} \rangle\|_{\psi_2} \leq 2C^2 \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{X}\|_F = \lambda,$$

where  $C$  is a universal constant. Thus, by applying Proposition 5.16 in [47], for  $Y_i - \mathbb{E}[Y_i]$ , we obtain

$$\mathbb{P}\left\{\left|\sum_{i=1}^n \frac{1}{n}(Y_i - \mathbb{E}[Y_i])\right| \geq t\right\} \leq 2 \exp\left[-c \min\left(\frac{nt^2}{\lambda^2}, \frac{nt}{\lambda}\right)\right].$$

According the covering argument for low-rank matrices Lemma 3.1 in [8] and covering number for sparse matrices in [48], we have

$$\begin{aligned} & \mathbb{P}\left\{\sup_{(\mathbf{S}-\mathbf{S}^*) \in \mathcal{N}_\epsilon^{cs}, \mathbf{X} \in \mathcal{N}_\epsilon^{3r}} \left|\frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} \rangle - \langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle\right| \geq t\right\} \\ & \leq 2|\mathcal{N}_\epsilon^{cs}| |\mathcal{N}_\epsilon^{3r}| \exp\left[-c_1 \min\left(\frac{nt^2}{\lambda^2}, \frac{nt}{\lambda}\right)\right] \\ & \leq 2\left(\frac{9}{\epsilon}\right)^{(d_1+d_2+1)3r} \cdot \left(\frac{c_2 d_1 d_2}{cs\epsilon}\right)^{cs} \cdot \exp\left[-c_1 \min\left(\frac{nt^2}{\lambda^2}, \frac{nt}{\lambda}\right)\right] \\ & \leq \exp\left[c_3(rd \log(1/\epsilon) + s \max\{\log d, \log(1/\epsilon)\}) - c_1 \min\left(\frac{nt^2}{\lambda^2}, \frac{nt}{\lambda}\right)\right] \leq \exp(-c'd), \end{aligned} \quad (\text{F.1})$$

where  $c_1, c_2, c_3$  are constants,  $\lambda = 2C^2$ , and the first inequality follows from union bound, the second inequality is due to the covering arguments, and the last inequality holds by setting  $t = c_4 \sqrt{(rd + s) \log d / \sqrt{n}}$ . Besides, note that for any  $\mathbf{X} \in \mathcal{M}_{3r}$ ,  $\mathbf{S} \in \mathbf{S}^* + \mathcal{M}_{cs}$ , there exists  $\mathbf{X}_1 \in \mathcal{N}_\epsilon^{3r}$ ,  $\mathbf{S}_1 \in \mathbf{S}^* + \mathcal{N}_\epsilon^{cs}$  such that  $\|\mathbf{X} - \mathbf{X}_1\|_F \leq \epsilon$  and  $\|\mathbf{S} - \mathbf{S}_1\|_F \leq \epsilon$ . Thus, we have

$$\begin{aligned} & \left|\frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} \rangle - \frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S}_1 - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X}_1 \rangle\right| \\ & \leq \left|\frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} - \mathbf{X}_1 \rangle\right| + \left|\frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}_1 \rangle \langle \mathbf{A}_i, \mathbf{X}_1 \rangle\right| \\ & \leq \sqrt{L_1 L_2} \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{X} - \mathbf{X}_1\|_F + \sqrt{L_1 L_2} \|\mathbf{S} - \mathbf{S}_1\|_F \cdot \|\mathbf{X}_1\|_F \leq 2\epsilon \sqrt{L_1 L_2}, \end{aligned} \quad (\text{F.2})$$

where the first inequality holds because of triangle inequality, and the second inequality follows from the restricted strong smoothness condition for both low-rank and sparse structures. Similarly, we have

$$|\langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle - \langle \mathbf{S}_1 - \mathbf{S}^*, \mathbf{X}_1 \rangle| \leq \|\mathbf{S} - \mathbf{S}^*\|_F \cdot \|\mathbf{X} - \mathbf{X}_1\|_F + \|\mathbf{S} - \mathbf{S}_1\|_F \cdot \|\mathbf{X}_1\|_F \leq 2\epsilon, \quad (\text{F.3})$$

Therefore, combining (F.1), (F.2) and (F.3), by triangle inequality, we obtain

$$\sup_{(\mathbf{S}-\mathbf{S}^*) \in \mathcal{M}_{cs}, \mathbf{X} \in \mathcal{M}_{3r}} \left|\frac{1}{n} \sum_{i=1}^n \langle \mathbf{A}_i, \mathbf{S} - \mathbf{S}^* \rangle \langle \mathbf{A}_i, \mathbf{X} \rangle - \langle \mathbf{S} - \mathbf{S}^*, \mathbf{X} \rangle\right| \leq t + 2\epsilon \sqrt{L_1 L_2} + 2\epsilon,$$

with probability at least  $1 - \exp(-c'd)$ . We establish the incoherence condition by setting  $\epsilon = t / (2\sqrt{L_1 L_2} + 2)$  in (F.3). By similar techniques, we can prove the second inequality in Lemma E.3. Note that we obtain  $K = C \sqrt{(rd + s) \log d / n}$  in Lemma E.3.  $\square$

## F.2 Proof of Lemma E.4

*Proof.* The first inequality in Lemma E.4 has been established in [37] Lemma 6. We provided the second inequality using Bernstein-type inequality and Union Bound. Recall that, we have

$$\left\|\frac{1}{n} \sum_{i=1}^n \epsilon_i \mathbf{A}_i\right\|_{\infty, \infty} = \max_{j,k} \left|\frac{1}{n} \sum_{i=1}^n \epsilon_i A_{jk}^i\right| = \max_{j,k} \left|\frac{1}{n} \sum_{i=1}^n Z_{jk}^i\right|,$$

where we let  $Z_{jk}^i = \epsilon_i A_{jk}^i$ . Since  $Z_{jk}^i$  are independent centered sub-exponential random variables for  $i = 1, \dots, n$  with  $\max_i \|Z_{jk}^i\|_{\psi_1} \leq 2 \max_i \|\epsilon_i\|_{\psi_2} \cdot \|A_{jk}^i\|_{\psi_2} \leq 2\nu$ , according to Proposition 5.16 in [47], we have

$$\mathbb{P}\left\{\left|\frac{1}{n} \sum_{i=1}^n Z_{jk}^i\right| \geq t\right\} \leq 2 \exp\left(-C' \frac{nt^2}{\nu^2}\right).$$

Thus by union bound, we have

$$\mathbb{P}\left\{\max_{j,k}\left|\frac{1}{n}\sum_{i=1}^n Z_{jk}^i\right|\geq t\right\}\leq 2d_1d_2\exp\left(-C'\frac{nt^2}{\nu^2}\right).$$

Let  $t = C_2\nu\sqrt{\log d/n}$ , we have the following inequality holds with probability at least  $1 - C/d$

$$\left\|\frac{1}{n}\sum_{i=1}^n \epsilon_i \mathbf{A}_i\right\|_{\infty,\infty}\leq C_2\nu\sqrt{\frac{\log d}{n}}.$$

□

Thus, we complete the proof.

### F.3 Proof of Lemma E.6

The proof of this lemma is inspired by the proof of Theorem 1 in [38], and we extended it to the sparse case. In order to prove Lemma E.6, we only need to prove the inequality (E.3) holds with high probability. Specifically, we consider the following event

$$E = \left\{\exists \mathbf{S} \in \mathcal{C}(n) \mid \left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}}\right| \leq \frac{1}{2}\frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}}\right\}.$$

Therefore, we want to establish the probability for event  $E$ , and we need the following lemmas.

**Lemma F.1.** Consider the robust PCA under observation model in section 4.2, for  $\ell = 1, 2, \dots$ , we have

$$\mathbb{P}(E_\ell) \leq \exp(-c_1n\alpha^{2\ell}\mu^2),$$

where we have

$$E_\ell := \left\{\exists \mathbf{S} \in \mathcal{B}'(\alpha^\ell\mu) \mid \left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}}\right| \geq \frac{5}{12}\frac{\alpha^\ell\mu}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}}\right\},$$

and

$$\mathcal{B}'(\alpha^\ell\mu) = \left\{\mathbf{S} \in \mathcal{C}(n, s) \mid \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}} \leq \frac{\alpha^\ell\mu}{\sqrt{d_1d_2}}\right\}.$$

*Proof of Lemma E.6.* The reminder of this proof is to derive the probability of the event  $E$ . In order to establish the probability of the event  $E$ , we make use of the peeling argument of the Frobenius norm  $\|\mathbf{S}\|_F$ . Let  $\mu = c\sqrt{\log d/n}$ , and  $\alpha = 6/5$ . For  $\ell = 1, 2, \dots$ , we define the sets

$$\mathcal{S}_\ell := \left\{\mathbf{S} \in \mathcal{C}(n, s) \mid \frac{\alpha^{\ell-1}\mu}{\sqrt{d_1d_2}} \leq \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}} \leq \frac{\alpha^\ell\mu}{\sqrt{d_1d_2}}\right\}.$$

Therefore, if the event  $E$  holds, there exist a matrix  $\mathbf{S}$  that must belongs to  $\mathcal{S}_\ell$  for some  $\ell = 1, 2, \dots$  such that

$$\left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}}\right| \geq \frac{1}{2}\frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}} \geq \frac{1}{2}\frac{\alpha^{\ell-1}\mu}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}} = \frac{5}{12}\frac{\alpha^\ell\mu}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}},$$

where the last equality is due to the fact that  $\alpha = 6/5$ .

Next, consider following events  $E_\ell$ , for  $\ell = 1, 2, \dots$

$$E_\ell := \left\{\exists \mathbf{S} \in \mathcal{B}'(\alpha^\ell\mu) \mid \left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}}\right| \geq \frac{5}{12}\frac{\alpha^\ell\mu}{\sqrt{d_1d_2}} + \frac{32\|\mathbf{S}\|_{\infty,\infty}}{\sqrt{n}}\right\},$$

where we have the constraint set

$$\mathcal{B}'(\alpha^\ell\mu) = \left\{\mathbf{S} \in \mathcal{C}(n, s) \mid \frac{\|\mathbf{S}\|_F}{\sqrt{d_1d_2}} \leq \frac{\alpha^\ell\mu}{\sqrt{d_1d_2}}\right\}.$$

Since  $\mathbf{S} \in \mathcal{S}_\ell$  implies that  $\mathbf{S} \in \mathcal{B}(\alpha^\ell \mu)$ , we can get  $E \subset \bigcup_{\ell=1}^{\infty} E_\ell$ . Therefore, we only need to upper bound the probability  $\mathbb{P}(\bigcup_{\ell=1}^{\infty} E_\ell)$ . In order to do so, we need upper bound the probability  $\mathbb{P}(E_\ell)$ . According to Lemma F.1, we have  $\mathbb{P}(E_\ell) \leq \exp(-c_1 n \alpha^{2\ell} \mu^2)$ . Therefore, we can obtain

$$\mathbb{P}(E) \leq \mathbb{P}\left(\bigcup_{\ell=1}^{\infty} E_\ell\right) \leq \sum_{\ell=1}^{\infty} \mathbb{P}(E_\ell) \leq \sum_{\ell=1}^{\infty} \exp(-c_1 n \alpha^{2\ell} \mu^2).$$

Thus according to the inequality  $a \leq e^a$ , we can obtain

$$\mathbb{P}(E) \leq \sum_{\ell=1}^{\infty} \exp(-2\ell c_1 n \mu^2 \log \alpha) \leq \frac{\exp(-2c_1 n \mu^2 \log \alpha)}{1 - \exp(-2c_1 n \mu^2 \log \alpha)} = \frac{\exp(-c_2 \log d)}{1 - \exp(-c_2 \log d)},$$

where the last equality comes from the definition  $\mu = c\sqrt{\log d/n}$ , and this implies  $\mathbb{P}(E) \leq c_3 \exp(-c_2 \log d)$ .  $\square$

#### F.4 Proof of Lemma E.7

*Proof.* The proof of this Lemma is similar to the proof of Lemma E.3, using Proposition 5.16 in [47] and covering number argument, with probability at least  $1 - \exp(-c_1 d)$ , we can obtain the restricted Lipschitz gradient condition in Condition 4.4 with parameter  $K = c_2 \sqrt{(rd + s) \log d/n}$ .  $\square$

#### F.5 Proof of Lemma E.8

*Proof.* For the first inequality in Lemma E.8, it has been established in [38] Proposition 1. For the second inequality in Lemma E.8, we use the similar proof as in the proof of Lemma E.4. By proposition 5.16 in [47] and union bound, with probability at least  $1 - C/d$ , we can obtain the required inequality.  $\square$

### G Proof of Auxiliary Lemmas in Appendix F

In order to prove Lemma F.1, we need the following lemmas.

**Lemma G.1.** We have the following holds with probability at least  $1 - C \exp(-C_1 n D^2)$

$$\max_{k=1, \dots, N(D/8)} \left| \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} \right| \leq \frac{D}{8\sqrt{d_1 d_2}} + \frac{32\|\mathbf{S}\|_{\infty, \infty}}{\sqrt{n}}.$$

**Lemma G.2.** We have the following holds

$$\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} \leq \frac{D}{2\sqrt{d_1 d_2}},$$

where we have

$$\mathcal{D}(\delta) := \{\mathbf{\Delta} \in \mathbb{R}^{d_1 \times d_2} \mid \|\mathbf{\Delta}\|_F \leq \delta, \|\mathbf{\Delta}\|_{1,1} \leq 2\rho(D), \|\mathbf{\Delta}\|_0 \leq 2\tilde{s}\},$$

and  $\rho(D) \leq D^2 / (c\sqrt{\log d/n})$ .

*Proof of Lemma F.1.* The proof of this lemma is inspired by the proof of Lemma 3 in [38]. Note that since the definition of the constraint set  $\mathcal{C}(n)$  and  $E$  is invariant to rescaling of  $\mathbf{S}$ , we can assume w.l.o.g. that  $\|\mathbf{S}\|_{\infty, \infty} = 1/\sqrt{d_1 d_2}$ . Therefore, it is equivalent to consider following events

$$E_\ell := \left\{ \exists \mathbf{S} \in \mathcal{B}(\alpha^\ell \mu) \mid \left| \frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} \right| \geq \frac{3\alpha^\ell \mu}{4\sqrt{d_1 d_2}} + \frac{32}{\sqrt{nd_1 d_2}} \right\}$$

where we have the constraint set

$$\mathcal{B}(\alpha^\ell \mu) = \left\{ \mathbf{S} \in \mathcal{C}(n, s) \mid \|\mathbf{S}\|_{\infty, \infty} \leq \frac{1}{\sqrt{d_1 d_2}}, \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} \leq \frac{\alpha^\ell \mu}{\sqrt{d_1 d_2}}, \|\mathbf{S}\|_{1,1} \leq \rho(\alpha^\ell \mu) \right\},$$

where  $\rho(\alpha^\ell \mu) \leq (\alpha^\ell \mu)^2 / (c\sqrt{\log d/n})$ . Define

$$Z_n(\alpha^\ell \mu) := \sup_{\mathbf{S} \in \mathcal{B}(\alpha^\ell \mu)} \left| \frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} \right|.$$

For simplicity, we use  $D$  to denote  $\alpha^\ell \mu$  in the following discussion. Therefore, we just need to prove the following probability bound

$$\mathbb{P}\left(Z_n(D) \geq \frac{3D}{4\sqrt{d_1 d_2}} + \frac{32}{\sqrt{nd_1 d_2}}\right) \leq c_3 \exp(-c_4 n D^2).$$

Suppose  $\mathbf{S}^1, \dots, \mathbf{S}^{N(\delta)}$  are a  $\delta$ -covering of  $\mathcal{B}(D)$  in terms of Frobenius norm. Therefore, for any  $\mathbf{S} \in \mathcal{B}(D)$ , there exist a matrix  $\mathbf{\Delta} \in \mathbb{R}^{d_1 \times d_2}$  and some index  $k \in \{1, \dots, N(\delta)\}$  satisfying  $\mathbf{S} = \mathbf{S}^k + \mathbf{\Delta}$ , where  $\|\mathbf{\Delta}\|_F \leq \delta$ . Thus we can obtain

$$\begin{aligned} \frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} &= \frac{\|\mathcal{A}(\mathbf{S}^k + \mathbf{\Delta})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k + \mathbf{\Delta}\|_F}{\sqrt{d_1 d_2}} \\ &\leq \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} + \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} + \frac{\|\mathbf{\Delta}\|_F}{\sqrt{d_1 d_2}} \\ &\leq \left| \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} \right| + \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} + \frac{\delta}{\sqrt{d_1 d_2}}. \end{aligned}$$

In addition we can get

$$\left| \frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} \right| \leq \left| \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} \right| + \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} + \frac{\delta}{\sqrt{d_1 d_2}}.$$

Therefore, we have

$$Z_n(D) \leq \frac{\delta}{\sqrt{d_1 d_2}} + \max_{k=1, \dots, N(\delta)} \left| \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} \right| + \sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}}, \quad (\text{G.1})$$

where we have  $\mathcal{D}(\delta) := \{\mathbf{\Delta} \in \mathbb{R}^{d_1 \times d_2} \mid \|\mathbf{\Delta}\|_F \leq \delta, \|\mathbf{\Delta}\|_{1,1} \leq 2\rho(D), \|\mathbf{\Delta}\|_0 \leq 2cs\}$ . We establish the high probability bound of (G.1) with  $\delta = D/8$ . First, according to Lemma G.1, we have

$$\max_{k=1, \dots, N(D/8)} \left| \frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}} \right| \leq \frac{D}{8\sqrt{d_1 d_2}} + \frac{32\|\mathbf{S}\|_{\infty, \infty}}{\sqrt{n}}, \quad (\text{G.2})$$

holds with probability at least  $1 - c \exp(-c_1 n D^2)$ .

Next, according to Lemma G.2, we have

$$\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} \leq \frac{D}{2\sqrt{d_1 d_2}}, \quad (\text{G.3})$$

holds with probability at least  $1 - c_2 \exp(-c_3 n D^2)$ .

Therefore, combining (G.2) and (G.3), we can get

$$Z_n(D) \leq \frac{D}{8\sqrt{d_1 d_2}} + \frac{D}{8\sqrt{d_1 d_2}} + \frac{D}{2\sqrt{d_1 d_2}} + \frac{32\|\mathbf{S}\|_{\infty, \infty}}{\sqrt{n}} \leq \frac{3D}{4\sqrt{d_1 d_2}} + \frac{32}{\sqrt{nd_1 d_2}},$$

holds with probability at least  $1 - c_3 \exp(-c_4 n D^2)$ , and the last inequality comes from that  $\|\mathbf{S}\|_{\infty, \infty} \leq 1/\sqrt{d_1 d_2}$ .  $\square$

## H Proofs of Auxiliary Lemmas in Appendix G

### H.1 Proof of Lemma G.1

*Proof.* First, we prove that for a fixed matrix  $\mathbf{S}$ , we have the following inequality holds

$$\mathbb{P}\left(\left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}}\right| \geq \frac{\delta}{\sqrt{d_1 d_2}} + \frac{32\|\mathbf{S}\|_{\infty, \infty}}{\sqrt{n}}\right) \leq C \exp(-C_1 n \delta^2).$$

Since we have

$$\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} = \frac{1}{\sqrt{n}} \sqrt{\sum_{j,k \in \Omega} \langle \mathbf{A}_{jk}, \mathbf{S} \rangle^2} = \frac{1}{\sqrt{n}} \sup_{\|\mathbf{w}\|_2=1} \sum_{j,k \in \Omega} w_i \langle \mathbf{A}_{jk}, \mathbf{S} \rangle,$$

we consider

$$\begin{aligned} \sqrt{d_1 d_2} \frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} &= \frac{\sqrt{d_1 d_2}}{\sqrt{n}} \sqrt{\sum_{j,k \in \Omega} \langle \mathbf{A}_{jk}, \mathbf{S} \rangle^2} = \frac{1}{\sqrt{n}} \sup_{\|\mathbf{w}\|_2=1} \sum_{j,k \in \Omega} w_{jk} \langle \sqrt{d_1 d_2} \mathbf{A}_{jk}, \mathbf{S} \rangle \\ &= \frac{1}{\sqrt{n}} \sup_{\|\mathbf{w}\|_2=1} \sum_{j,k \in \Omega} w_{jk} Y_{jk}, \end{aligned}$$

where we have the random variables  $Y_{jk}$  satisfying  $|Y_{jk}| = |\langle \sqrt{d_1 d_2} \mathbf{A}_{jk}, \mathbf{S} \rangle| \leq \sqrt{d_1 d_2} \|\mathbf{S}\|_{\infty, \infty} = 1$ . Therefore, according to lemma I.1, we have

$$\mathbb{P}\left(\left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \mathbb{E}\left[\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}}\right]\right| \geq \frac{\delta}{\sqrt{d_1 d_2}} + \frac{16}{\sqrt{nd_1 d_2}}\right) \leq C \exp(-C_1 n \delta^2). \quad (\text{H.1})$$

In addition, we have

$$\begin{aligned} \left|\frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}} - \mathbb{E}\left[\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}}\right]\right| &= \left|\sqrt{\mathbb{E}\left[\frac{\|\mathcal{A}(\mathbf{S})\|_2^2}{n}\right]} - \mathbb{E}\left[\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}}\right]\right| \\ &\leq \sqrt{\mathbb{E}\left[\frac{\|\mathcal{A}(\mathbf{S})\|_2^2}{n}\right] - \mathbb{E}\left[\left(\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}}\right)^2\right]} \leq \frac{16}{\sqrt{nd_1 d_2}}. \end{aligned} \quad (\text{H.2})$$

Therefore, combining (H.1) and (H.2), we can obtain

$$\mathbb{P}\left(\left|\frac{\|\mathcal{A}(\mathbf{S})\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}\|_F}{\sqrt{d_1 d_2}}\right| \geq \frac{\delta}{d_1 d_2} + \frac{32}{\sqrt{nd_1 d_2}}\right) \leq C \exp(-C_1 n \delta^2).$$

Next, according to Lemma 4 in [38], there exists a  $\delta$ -covering of  $\mathcal{B}(D)$  such that

$$\log N(\delta) \leq C_3 (\rho(D)/\delta)^2 \log d.$$

Therefore, we can get

$$\begin{aligned} \mathbb{P}\left[\max_{k=1, \dots, N(D/8)} \left|\frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}}\right| \geq \frac{\delta}{\sqrt{d_1 d_2}} + \frac{32}{\sqrt{nd_1 d_2}}\right] &\leq C \exp(-C_1 n \delta^2 + \log N(\delta)) \\ &\leq C \exp(-C_1 n \delta^2 + C_3 (\rho(D)/\delta)^2 \log d). \end{aligned}$$

Since we have  $\delta = D/8$  and  $\rho(D) = C_4 D^2 / \sqrt{\log d/n}$ , we can obtain

$$\mathbb{P}\left[\max_{k=1, \dots, N(D/8)} \left|\frac{\|\mathcal{A}(\mathbf{S}^k)\|_2}{\sqrt{n}} - \frac{\|\mathbf{S}^k\|_F}{\sqrt{d_1 d_2}}\right| \geq \frac{\delta}{\sqrt{d_1 d_2}} + \frac{32\|\mathbf{S}\|_{\infty, \infty}}{\sqrt{n}}\right] \leq C \exp(-C_2 n \delta^2),$$

which complete the proof.  $\square$



## H.2 Proof of Lemma G.2

*Proof.* According to Lemma 5 in [38], we have following results

$$\mathbb{P}\left[\left|\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \sqrt{d_1 d_2} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} - \mathbb{E}\left[\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \sqrt{d_1 d_2} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}}\right]\right| \geq \delta\right] \leq C \exp(-C_1 n \delta^2), \quad (\text{H.3})$$

and

$$\left(\mathbb{E}\left[\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \sqrt{d_1 d_2} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}}\right]\right)^2 \leq 16 \sqrt{d_1 d_2} \|\mathbf{\Delta}\|_{\infty} \mathbb{E}\left[\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{1}{n} \sum_{j,k \in \Omega} \xi_{jk} \langle \mathbf{A}_{jk}, \mathbf{\Delta} \rangle\right] + \delta^2,$$

where  $\xi_{jk}$  are independent Rademacher variables. Furthermore, by the duality between norms, we can obtain

$$\frac{1}{n} \sum_{j,k \in \Omega} \xi_{jk} \langle \mathbf{A}_{jk}, \mathbf{\Delta} \rangle \leq \left\| \frac{1}{n} \sum_{j,k \in \Omega} \xi_{jk} \mathbf{A}_{jk} \right\|_{\infty} \cdot \|\mathbf{\Delta}\|_{1,1} \leq \rho(D) \left\| \frac{1}{n} \sum_{j,k \in \Omega} \xi_{jk} \mathbf{A}_{jk} \right\|_{\infty, \infty},$$

where the last inequality is due to the fact that  $\mathbf{\Delta} \in \mathcal{D}(\delta)$ . Finally, we have

$$\left\| \frac{1}{n} \sum_{j,k \in \Omega} \xi_{jk} \mathbf{A}_{jk} \right\|_{\infty, \infty} \leq C \sqrt{\frac{\log d}{n}}. \quad (\text{H.4})$$

To prove this, we use Hoeffding's inequality and Union Bound. By the definition of  $\mathbf{A}_i$ , we can obtain

$$\left\| \frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{A}_i \right\|_{\infty, \infty} = \max_{j,k} \left| \frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{e}_j^i \mathbf{e}_k^i \right| = \max_{j,k} \left| \frac{1}{n} \sum_{i=1}^n Z_{jk}^i \right|,$$

where we have  $Z_{jk}^i = \xi_i A_{jk}$ . Thus we can get  $|Z_{jk}^i| \leq |\xi_i| = 1$ , and we conclude that  $Z_{jk}^i$  are independent centered sub-Gaussian random variables for  $i = 1, \dots, n$ . Therefore, following the same procedure as in the proof of Lemma E.4, we can obtain inequality (H.4). Therefore, we can obtain

$$\left(\mathbb{E}\left[\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}}\right]\right)^2 \leq C \frac{\|\mathbf{\Delta}\|_{\infty, \infty} \rho(D)}{\sqrt{d_1 d_2}} \sqrt{\frac{\log d}{n}} + \frac{\delta^2}{d_1 d_2} \leq C' \frac{D^2}{d_1 d_2},$$

where the last inequality comes from the definition of  $\rho(D)$ ,  $\delta$  and  $\|\mathbf{\Delta}\|_{\infty, \infty} \leq 2/\sqrt{d_1 d_2}$ . It implies that

$$\mathbb{E}\left[\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}}\right] \leq C'' \frac{D}{\sqrt{d_1 d_2}}. \quad (\text{H.5})$$

Thus combining (H.3) and (H.5), we have

$$\sup_{\mathbf{\Delta} \in \mathcal{D}(\delta)} \frac{\|\mathcal{A}(\mathbf{\Delta})\|_2}{\sqrt{n}} \leq \frac{D}{2\sqrt{d_1 d_2}}$$

holds with probability at least  $1 - C \exp(-C_1 n D^2)$ .  $\square$

## I Other Auxiliary Lemmas

**Lemma I.1.** [31] Consider independent random variables  $Y_1, \dots, Y_n$  such that  $a_i \leq Y_i \leq b_i$  for  $i = 1, \dots, n$ . Let

$$Z := \sup_{\mathbf{t} \in \mathcal{T}} \sum_{i=1}^n t_i Y_i,$$

where  $\mathcal{T}$  is a family of vectors  $\mathbf{t} \in \mathbb{R}^n$  such that  $\sigma = \sup_{\mathbf{t} \in \mathcal{T}} \left(\sum_{i=1}^n t_i^2 (b_i - a_i)^2\right)^{1/2} \leq \infty$ . Then, for any  $r \geq 0$ , we have

$$\mathbb{P}(|Z - m_Z| \geq r) \leq 4 \exp\left(-\frac{r^2}{4\sigma^2}\right),$$

where  $m_Z$  is a median of  $Z$ . Furthermore, we have

$$|\mathbb{E}(Z) - m_Z| \leq 4\sqrt{\pi}\sigma \quad \text{and} \quad \text{Var}(Z) \leq 16\sigma^2.$$