

A. Introduction

In this supplement, we first provide additional experimental results on the proposed estimator with MCP regularization, followed by the details of technical proof for the main results, including proofs of theorems and auxiliary lemmas.

B. Additional Experimental Results

Regarding matrix completion and matrix sensing, we present additional experimental results of the proposed estimator with MCP penalty. Due to the similar properties and parameter settings of these two nonconvex penalties, the MCP penalty and SCAD penalty, the numerical behaviour of the proposed estimator with MCP penalty resembles the one with SCAD penalty, as shown in Figure 2.

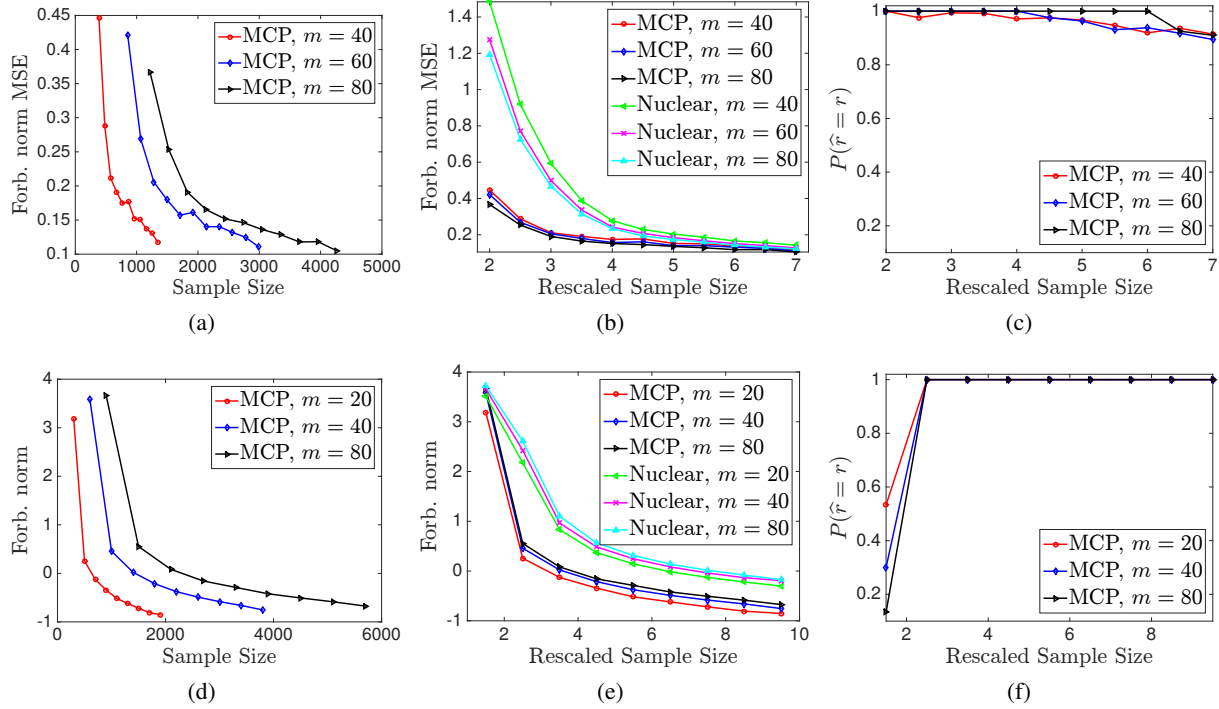


Figure 2. Simulation Results for Matrix Completion and Matrix Sensing with MCP penalty. Accordingly, the size of matrix and the rank are $m \times m$. The results of matrix completion, with rank $r = \lfloor \log^2 m \rfloor$, in Figure 2(a)-2(c) with the rescaled sample size $N = n/(rm \log m)$; while matrix sensing, for rank $r = 10$, is studied in Figure 2(d)-2(f) with rescaled sample size $N = n/(rm)$.

In detail, Figure 2(a)- 2(c) are the results for matrix completion. With the same settings as experiments shown in Figure 1, we have that the estimator with MCP penalty, a particular case of the proposed estimator with nonconvex penalty, behaves in accordance with our theoretical analysis and outperforms the estimator with nuclear norm. For the other example, *i.e.*, matrix sensing, the results in Figure 2(d)- 2(f) manifest the superiority of the estimator with MCP penalty. Particularly, for both examples, we have with high probability, the rank of the underlying matrix is recovered with high probability.

C. Background

For matrix $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$, which is exactly low-rank and has rank r , we have the singular value decomposition (SVD) form of $\Theta^* = \mathbf{U}^* \mathbf{\Gamma}^* \mathbf{V}^{*\top}$, where $\mathbf{U}^* \in \mathbb{R}^{m_1 \times r}$, $\mathbf{V}^* \in \mathbb{R}^{m_2 \times r}$ are matrices consist of left and right singular vectors, and $\mathbf{\Gamma}^* = \text{diag}(\gamma_1^*, \dots, \gamma_r^*) \in \mathbb{R}^{r \times r}$. Based on \mathbf{U}^* , \mathbf{V}^* , we define the following two subspaces of $\mathbb{R}^{m_1 \times m_2}$:

$$\mathcal{F}(\mathbf{U}^*, \mathbf{V}^*) := \{\Delta \mid \text{row}(\Delta) \subseteq \mathbf{V}^* \text{ and } \text{col}(\Delta) \subseteq \mathbf{U}^*\},$$

and

$$\mathcal{F}^\perp(\mathbf{U}^*, \mathbf{V}^*) := \{\Delta \mid \text{row}(\Delta) \perp \mathbf{V}^* \text{ and } \text{col}(\Delta) \perp \mathbf{U}^*\},$$

where $\mathbf{\Delta} \in \mathbb{R}^{m_1 \times m_2}$ is an arbitrary matrix, and $\text{row}(\mathbf{\Delta}) \subseteq \mathbb{R}^{m_2}$, $\text{col}(\mathbf{\Delta}) \subseteq \mathbb{R}^{m_1}$ are the row space and column space of the matrix $\mathbf{\Delta}$, respectively. We will use the shorthand notation of $\mathcal{F}, \mathcal{F}^\perp$, when $(\mathbf{U}^*, \mathbf{V}^*)$ are clear from the context. Define $\Pi_{\mathcal{F}}, \Pi_{\mathcal{F}^\perp}$ as the projection operator onto the subspaces \mathcal{F} and \mathcal{F}^\perp :

$$\begin{aligned}\Pi_{\mathcal{F}}(\mathbf{A}) &= \mathbf{U}^* \mathbf{U}^{*\top} \mathbf{A} \mathbf{V}^* \mathbf{V}^{*\top}, \\ \Pi_{\mathcal{F}^\perp}(\mathbf{A}) &= (\mathbf{I}_{m_1} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{A} (\mathbf{I}_{m_2} - \mathbf{V}^* \mathbf{V}^{*\top}).\end{aligned}\tag{C.1}$$

Thus, for all $\mathbf{\Delta} \in \mathbb{R}^{m_1 \times m_2}$, we have its orthogonal complement $\mathbf{\Delta}''$ with respect to the true low-rank matrix $\mathbf{\Theta}^*$ as follows:

$$\begin{aligned}\mathbf{\Delta}'' &= (\mathbf{I}_{m_1} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{\Delta} (\mathbf{I}_{m_2} - \mathbf{V}^* \mathbf{V}^{*\top}), \\ \mathbf{\Delta}' &= \mathbf{\Delta} - \mathbf{\Delta}'',\end{aligned}\tag{C.2}$$

where $\mathbf{\Delta}'$ is the component which has overlapped row and column space with $\mathbf{\Theta}^*$. (Negahban et al., 2012) gives detailed discussion about the concept of decomposibility and a large class of decomposable norms, among which the decomposability of the nuclear norm and Frobenius norm is relevant to our problem. For low-rank estimation, we have the equality that $\|\mathbf{\Theta}^* + \mathbf{\Delta}'\|_* = \|\mathbf{\Theta}^*\|_* + \|\mathbf{\Delta}''\|_*$ with $\mathbf{\Delta}''$ defined above.

D. Proof of the Main Results

D.1. Proof of Theorem 3.4

We first define $\tilde{\mathcal{L}}_{n,\lambda}(\cdot)$ as follows,

$$\tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}) = \mathcal{L}_n(\mathbf{\Theta}) + \mathcal{Q}_\lambda(\mathbf{\Theta}).$$

Based on the the restrict strongly convexity of \mathcal{L}_n , and the curvature parameter of the non-convex penalty, if $\kappa(\mathfrak{X}) > \zeta_-$, we have the restrict strongly convexity of $\tilde{\mathcal{L}}_{n,\lambda}(\cdot)$, as stated in the following lemma.

Lemma D.1. Under Assumption 3.1, if it is assumed that $\mathbf{\Theta}_1 - \mathbf{\Theta}_2 \in \mathcal{C}$, we have

$$\tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}_2) \geq \tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}_1) + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}_1), \mathbf{\Theta}_2 - \mathbf{\Theta}_1 \rangle + \frac{\kappa(\mathfrak{X}) - \zeta_-}{2} \|\mathbf{\Theta}_2 - \mathbf{\Theta}_1\|_F^2.$$

Proof. Proof is provided in Section F.1. □

In the following, we prove that $\hat{\mathbf{\Delta}} = \hat{\mathbf{\Theta}} - \mathbf{\Theta}^*$ lies in the cone \mathcal{C} , where

$$\mathcal{C} = \{\mathbf{\Delta} \in \mathbb{R}^{m_1 \times m_2} \mid \|\Pi_{\mathcal{F}^\perp}(\mathbf{\Delta})\|_* \leq 5 \|\Pi_{\mathcal{F}}(\mathbf{\Delta})\|_*\}.$$

Lemma D.2. Under Assumption 3.1, the condition $\kappa(\mathfrak{X}) > \zeta_-$, and the regularization parameter $\lambda \geq 2 \|\mathfrak{X}^*(\epsilon)\|_2 / n$, we have

$$\|\Pi_{\mathcal{F}}(\hat{\mathbf{\Theta}} - \mathbf{\Theta}^*)\|_* \leq 5 \|\Pi_{\mathcal{F}^\perp}(\hat{\mathbf{\Theta}} - \mathbf{\Theta}^*)\|_*.$$

Proof. Proof is provided in Section F.2. □

Now we are ready to prove Theorem 3.4.

Proof of Theorem 3.4. According to Lemma D.1, we have

$$\tilde{\mathcal{L}}_{n,\lambda}(\hat{\mathbf{\Theta}}) \geq \tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}^*) + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}^*), \hat{\mathbf{\Theta}} - \mathbf{\Theta}^* \rangle + \frac{\kappa(\mathfrak{X}) - \zeta_-}{2} \|\hat{\mathbf{\Theta}} - \mathbf{\Theta}^*\|_F^2,\tag{D.1}$$

$$\tilde{\mathcal{L}}_{n,\lambda}(\mathbf{\Theta}^*) \geq \tilde{\mathcal{L}}_{n,\lambda}(\hat{\mathbf{\Theta}}) + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\hat{\mathbf{\Theta}}), \mathbf{\Theta}^* - \hat{\mathbf{\Theta}} \rangle + \frac{\kappa(\mathfrak{X}) - \zeta_-}{2} \|\mathbf{\Theta}^* - \hat{\mathbf{\Theta}}\|_F^2.\tag{D.2}$$

Meanwhile, since $\|\cdot\|_*$ is convex, we have

$$\lambda\|\widehat{\Theta}\|_* \geq \lambda\|\Theta^*\|_* + \lambda\langle\widehat{\Theta} - \Theta^*, \mathbf{W}^*\rangle, \quad (\text{D.3})$$

$$\lambda\|\Theta^*\|_* \geq \lambda\|\widehat{\Theta}\|_* + \lambda\langle\Theta^* - \widehat{\Theta}, \mathbf{W}^*\rangle, \quad (\text{D.4})$$

where $\mathbf{W}^* \in \|\Theta^*\|_*$.

Adding (D.1) to (D.4), we have

$$0 \geq \langle\nabla\tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda\mathbf{W}^*, \widehat{\Theta} - \Theta^*\rangle + \langle\nabla\tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda\widehat{\mathbf{W}}, \Theta^* - \widehat{\Theta}\rangle + (\kappa(\mathfrak{X}) - \zeta_-)\|\widehat{\Theta} - \Theta^*\|_F^2.$$

Since $\widehat{\Theta}$ is the solution to the SDP (2.2), $\widehat{\Theta}$ satisfies the optimality condition (variational inequality), for any $\Theta' \in \mathbb{R}^{m_1 \times m_2}$, it holds that

$$\max_{\Theta'} \langle\nabla\tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda\widehat{\mathbf{W}}, \widehat{\Theta} - \Theta'\rangle \leq 0,$$

which implies

$$\langle\nabla\tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda\widehat{\mathbf{W}}, \Theta^* - \widehat{\Theta}\rangle \geq 0.$$

Hence,

$$\begin{aligned} (\kappa(\mathfrak{X}) - \zeta_-)\|\widehat{\Theta} - \Theta^*\|_F^2 &\leq \langle\nabla\tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda\mathbf{W}^*, \Theta^* - \widehat{\Theta}\rangle \\ &\leq \left\langle \Pi_{\mathcal{F}^\perp}(\nabla\tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda\mathbf{W}^*), \Theta^* - \widehat{\Theta} \right\rangle + \left\langle \Pi_{\mathcal{F}}(\nabla\tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda\mathbf{W}^*), \Theta^* - \widehat{\Theta} \right\rangle. \end{aligned} \quad (\text{D.5})$$

Recall that $\gamma^* = \gamma(\Theta^*)$ is the vector of (ordered) singular values of Θ^* . In the following, we decompose (D.5) into three parts with regard to the magnitudes of the singular values of Θ^* .

- (1) $i \in S^c$ that $(\gamma^*)_i = 0$;
- (2) $i \in S_1$ that $(\gamma^*)_i \geq \nu$;
- (3) $i \in S_2$ that $\nu > (\gamma^*)_i > 0$.

Note that $S_1 \cup S_2 = S$.

(1) For $i \in S^c$, it correspond to the projector $\Pi_{\mathcal{F}^\perp}(\cdot)$ since $\gamma(\Pi_{\mathcal{F}^\perp}(\Theta^*)) = (\gamma^*)_{S^c} = \mathbf{0}$.

Based on the regularity condition (iii) in Assumption 3.3 that $q'_\lambda(0) = 0$, we have that $\nabla\mathcal{Q}_\lambda(\Theta^*) = \mathbf{U}^*q'_\lambda(\Gamma^*)\mathbf{V}^{*\top}$ where $\Gamma^* \in \mathbb{R}^{r \times r}$ is the diagonal matrix with $\text{diag}(\Gamma^*) = \gamma^*$, we have

$$\begin{aligned} \Pi_{\mathcal{F}^\perp}(\nabla\mathcal{Q}_\lambda(\Theta^*)) &= (\mathbf{I}_{m_1} - \mathbf{U}^*\mathbf{U}^{*\top})\mathbf{U}^*q'_\lambda(\Gamma^*)\mathbf{V}^{*\top}(\mathbf{I}_{m_2} - \mathbf{V}^*\mathbf{V}^{*\top}) \\ &= (\mathbf{U}^* - \mathbf{U}^*)q'_\lambda(\Gamma^*)(\mathbf{V}^{*\top} - \mathbf{V}^{*\top}) \\ &= \mathbf{0}. \end{aligned}$$

Therefore,

$$\Pi_{\mathcal{F}^\perp}(\nabla\mathcal{Q}_\lambda(\Theta^*)) = \mathbf{0}.$$

Meanwhile, we have

$$\|\Pi_{\mathcal{F}^\perp}(\nabla\mathcal{L}_n(\Theta^*))\|_2 \leq \|\nabla\mathcal{L}_n(\Theta^*)\|_2 = \frac{\|\mathfrak{X}^*(\epsilon)\|_2}{n} \leq \lambda.$$

For $\mathbf{Z}^* = -\lambda^{-1}\Pi_{\mathcal{F}^\perp}(\nabla\mathcal{L}_n(\Theta^*))$, we have $\mathbf{W}^* = \mathbf{U}^*\mathbf{V}^{*\top} + \mathbf{Z}^* \in \partial\|\Theta^*\|_*$ because $\|\mathbf{Z}^*\|_2 \leq 1$ and $\mathbf{Z}^* \in \mathcal{F}^\perp$, which satisfies the condition of \mathbf{W}^* to be subgradient of $\|\Theta^*\|_*$. With this particular choice of \mathbf{W}^* , we have

$$\Pi_{\mathcal{F}^\perp}(\nabla\mathcal{L}_n(\Theta^*) + \lambda\mathbf{W}^*) = \Pi_{\mathcal{F}^\perp}(\nabla\mathcal{L}_n(\Theta^*)) + \lambda\mathbf{Z}^* = \mathbf{0},$$

which implies that

$$\langle \Pi_{\mathcal{F}^\perp} (\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle = \langle \mathbf{0}, \Theta^* - \hat{\Theta} \rangle = 0. \quad (\text{D.6})$$

(2) Consider $i \in S_1$ that $(\gamma^*)_i \geq \nu$. Let $|S_1| = r_1$. Define a subspace of \mathcal{F} associated with S_1 as follows

$$\mathcal{F}_{S_1}(\mathbf{U}^*, \mathbf{V}^*) := \{\Delta \in \mathbb{R}^{m_1 \times m_2} \mid \text{row}(\Delta) \subset \mathbf{V}_{S_1}^* \text{ and } \text{col}(\Delta) \subset \mathbf{U}_{S_1}^*\},$$

where $\mathbf{U}_{S_1}^*$ and $\mathbf{V}_{S_1}^*$ is the matrix with the i^{th} row of \mathbf{U}^* and \mathbf{V}^* where $i \in S_1$.

Recall that $\mathcal{P}_\lambda(\Theta^*) = \mathcal{Q}_\lambda(\Theta^*) + \lambda \|\Theta^*\|_*$. We have

$$\nabla \mathcal{P}_\lambda(\Theta^*) = \nabla \mathcal{Q}_\lambda(\Theta^*) + \lambda(\mathbf{U}^* \mathbf{V}^{*\top} + \mathbf{Z}^*).$$

Projecting $\nabla \mathcal{P}_\lambda(\Theta^*)$ into the subspace \mathcal{F}_{S_1} , we have

$$\begin{aligned} \Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{P}_\lambda(\Theta^*)) &= \Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{Q}_\lambda(\Theta^*) + \lambda \mathbf{U}^* \mathbf{V}^{*\top} + \lambda \mathbf{Z}^*) \\ &= \mathbf{U}_{S_1}^* q'_\lambda(\Gamma_{S_1}^*)(\mathbf{V}_{S_1}^*)^\top + \lambda \mathbf{U}_{S_1}^* (\mathbf{V}_{S_1}^*)^\top \\ &= \mathbf{U}_{S_1}^* (q'_\lambda(\Gamma_{S_1}^*) + \lambda \mathbf{I}_{S_1})(\mathbf{V}_{S_1}^*)^\top, \end{aligned}$$

where $\Gamma_{S_1}^* \in \mathbb{R}^{r_1 \times r_1}$ and $(q'_\lambda(\Gamma_{S_1}^*) + \lambda \mathbf{I}_{S_1})$ is a diagonal matrix that $(q'_\lambda(\Gamma_{S_1}^*) + \lambda \mathbf{I}_{S_1})_{ii} = 0$ for $i \notin S_1$, and for all $i \in S_1$,

$$(q'_\lambda(\Gamma_{S_1}^*) + \lambda \mathbf{I}_{S_1})_{ii} = q'_\lambda((\gamma^*)_i) + \lambda = p'_\lambda((\gamma^*)_i) = 0,$$

where the last equality is because $p_\lambda(\cdot)$ satisfies the regularity condition (i) with $(\gamma^*)_i \geq \nu$ for $i \in S_1$. Thus, we have $q'_\lambda(\mathbf{D}_{S_1}) + \lambda \mathbf{I}_{S_1} = \mathbf{0}$, which indicates that $\Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{P}_\lambda(\Theta^*)) = \mathbf{0}$. Therefore, we have

$$\begin{aligned} \langle \Pi_{\mathcal{F}_{S_1}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle &= \langle \Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*) + \nabla \mathcal{P}_\lambda(\Theta^*)), \Theta^* - \hat{\Theta} \rangle \\ &= \langle \Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*)), \Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta}) \rangle \\ &\leq \|\Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 \cdot \|\Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta})\|_*, \end{aligned}$$

where the last inequality is derived from the Hölder inequality. What remains is to bound $\|\Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta})\|_*$. By the properties of projection on to the subspace \mathcal{F}_{S_1} , we have

$$\|\Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta})\|_* \leq \sqrt{r_1} \|\Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta})\|_F \leq \sqrt{r_1} \|\Theta^* - \hat{\Theta}\|_F,$$

where the second inequality is due to the fact that $\text{rank}(\Pi_{\mathcal{F}_{S_1}}(\Theta^* - \hat{\Theta})) \leq r_1$. Therefore, we have

$$\langle \Pi_{\mathcal{F}_{S_1}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle \leq \sqrt{r_1} \|\Pi_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 \cdot \|\Theta^* - \hat{\Theta}\|_F. \quad (\text{D.7})$$

(3) Finally, consider $i \in S_2$ that $(\gamma^*)_i \leq \nu$. Let $|S_2| = r_2$. Define a subspace of \mathcal{F} associated with S_2 as follows

$$\mathcal{F}_{S_2}(\mathbf{U}^*, \mathbf{V}^*) := \{\Delta \in \mathbb{R}^{m_1 \times m_2} \mid \text{row}(\Delta) \subset \mathbf{V}_{S_2}^* \text{ and } \text{col}(\Delta) \subset \mathbf{U}_{S_2}^*\},$$

where $\mathbf{U}_{S_2}^*$ and $\mathbf{V}_{S_2}^*$ is the matrix with the i^{th} row of \mathbf{U}^* and \mathbf{V}^* where $i \in S_2$. It is obvious that for all $\Delta \in \mathbb{R}^{m_1 \times m_2}$, the following decomposition holds

$$\Pi_{\mathcal{F}}(\Delta) = \Pi_{\mathcal{F}_{S_1}}(\Delta) + \Pi_{\mathcal{F}_{S_2}}(\Delta).$$

In addition, since $\mathbf{U}^*, \mathbf{V}^*$ are unitary matrices, we have

$$\mathcal{F}_{S_1} \subset \mathcal{F}_{S_2}^\perp, \text{ and } \mathcal{F}_{S_2} \subset \mathcal{F}_{S_1}^\perp,$$

where $\mathcal{F}_{S_1}^\perp, \mathcal{F}_{S_2}^\perp$ denote the complementary subspace of \mathcal{F}_{S_1} and \mathcal{F}_{S_2} , respectively. Similar to analysis in (2) on S_1 , we have

$$\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{Q}_\lambda(\Theta^*)) = \mathbf{U}_{S_2}^* q'_\lambda(\mathbf{\Gamma}_{S_2}^*)(\mathbf{V}_{S_2}^*)^\top,$$

where $q'_\lambda(\mathbf{\Gamma}_{S_2}^*)$ is a diagonal matrix that $(q'_\lambda(\mathbf{\Gamma}_{S_2}^*))_{ii} = 0$ for $i \notin S_2$, and for all $i \in S_2$, $(q'_\lambda(\mathbf{\Gamma}_{S_2}^*))_{ii} = q'_\lambda((\gamma^*)_i) \leq \lambda$, since $(\gamma^*)_i \leq \nu$ and $q_\lambda(\cdot)$ satisfies the regularity condition (iv). Therefore

$$\|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{Q}_\lambda(\Theta^*))\|_2 = \max_{i \in S_2} (q'_\lambda(\mathbf{\Gamma}_{S_2}^*))_{ii} \leq \lambda. \quad (\text{D.8})$$

Meanwhile, we have

$$\|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\lambda \mathbf{W}^*)\|_2 \leq \|\mathbf{\Pi}_{\mathcal{F}}(\lambda \mathbf{U}^* \mathbf{V}^{*\top})\|_2 = \lambda, \quad (\text{D.9})$$

where the first inequality is due the fact that $\mathcal{F}_{S_2} \in \mathcal{F}$, and last equality comes from the fact that $\|\mathbf{U}^* \mathbf{V}^{*\top}\|_2 = 1$. Therefore, we have

$$\|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\lambda \mathbf{W}^*)\|_2 \leq \lambda. \quad (\text{D.10})$$

In addition, we have the fact that $\|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 \leq \|\nabla \mathcal{L}_n(\Theta^*)\|_2 \leq \lambda$, which indicates that

$$\begin{aligned} \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle &= \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{L}_n(\Theta^*) + \nabla \mathcal{Q}_\lambda(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle \\ &= \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{L}_n(\Theta^*)), \Theta^* - \hat{\Theta} \rangle + \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{Q}_\lambda(\Theta^*)), \Theta^* - \hat{\Theta} \rangle + \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle \\ &\leq \left[\|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 + \|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \mathcal{Q}_\lambda(\Theta^*))\|_2 + \|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\lambda \mathbf{W}^*)\|_2 \right] \|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\Theta^* - \hat{\Theta})\|_*, \end{aligned}$$

where the last inequality is due to Hölder's inequality. Since we have obtained the bound for each term, as in (D.8), (D.9), (D.10), we have

$$\begin{aligned} \langle \mathbf{\Pi}_{\mathcal{F}_{S_2}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*), \Theta^* - \hat{\Theta} \rangle &\leq 3\lambda \|\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\Theta^* - \hat{\Theta})\|_* \\ &\leq 3\lambda \sqrt{r_2} \|\Theta^* - \hat{\Theta}\|_F, \end{aligned} \quad (\text{D.11})$$

where the last inequality utilizes the fact that $\text{rank}(\mathbf{\Pi}_{\mathcal{F}_{S_2}}(\Theta^* - \hat{\Theta})) \leq r_2$.

Adding (D.6), (D.7), and (D.11), we have

$$\begin{aligned} (\kappa(\mathfrak{X}) - \zeta_-) \|\hat{\Theta} - \Theta^*\|_F^2 &\leq \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) + \lambda \mathbf{W}^*, \Theta^* - \hat{\Theta} \rangle \\ &\leq \sqrt{r_1} \|\mathbf{\Pi}_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 \cdot \|\Theta^* - \hat{\Theta}\|_F + 3\lambda \sqrt{r_2} \|\Theta^* - \hat{\Theta}\|_F, \end{aligned}$$

which indicate that

$$\|\hat{\Theta} - \Theta^*\|_F \leq \frac{\sqrt{r_1}}{\kappa(\mathfrak{X}) - \zeta_-} \|\mathbf{\Pi}_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 + \frac{3\lambda \sqrt{r_2}}{\kappa(\mathfrak{X}) - \zeta_-}.$$

This completes the proof. \square

D.2. Proof of Theorem 3.5

Before presenting the proof of Theorem 3.5, we need the following lemma.

Lemma D.3 (Deterministic Bound). Suppose $\Theta^* \in \mathbb{R}^{m_1 \times m_2}$ has rank r , $\mathfrak{X}(\cdot)$ satisfies RSC with respect to \mathcal{C} . Then the error bound between the oracle estimator $\hat{\Theta}_O$ and true Θ^* satisfies

$$\|\hat{\Theta}_O - \Theta^*\|_F \leq \frac{2\sqrt{r} \|\mathbf{\Pi}_{\mathcal{F}}(\nabla \mathcal{L}_n(\Theta^*))\|_2}{\kappa(\mathfrak{X})}, \quad (\text{D.12})$$

Proof. Proof is provided in Section F.3. \square

Proof of Theorem 3.5. Suppose $\widehat{\mathbf{W}} \in \partial \|\widehat{\boldsymbol{\Theta}}\|_*$, since $\widehat{\boldsymbol{\Theta}}$ is the solution to the SDP (2.2), the variational inequality yields

$$\max_{\boldsymbol{\Theta}' } \langle \widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}', \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\boldsymbol{\Theta}}) + \lambda \widehat{\mathbf{W}} \rangle \leq 0. \quad (\text{D.13})$$

In the following, we will show that there exists some $\widehat{\mathbf{W}}_O \in \partial \|\widehat{\boldsymbol{\Theta}}_O\|_*$ such that, for all $\boldsymbol{\Theta}' \in \mathbb{R}^{m_1 \times m_2}$,

$$\max_{\boldsymbol{\Theta}' } \langle \widehat{\boldsymbol{\Theta}}_O - \boldsymbol{\Theta}', \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\boldsymbol{\Theta}}_O) + \lambda \widehat{\mathbf{W}}_O \rangle \leq 0. \quad (\text{D.14})$$

Recall that $\tilde{\mathcal{L}}_{n,\lambda}(\boldsymbol{\Theta}) = \mathcal{L}_n(\boldsymbol{\Theta}) + \mathcal{Q}_\lambda(\boldsymbol{\Theta})$. By projecting the components of the inner product of the LHS in (D.14) into two complementary spaces \mathcal{F} and \mathcal{F}^\perp , we have the following decomposition

$$\begin{aligned} & \langle \widehat{\boldsymbol{\Theta}}_O - \boldsymbol{\Theta}', \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\boldsymbol{\Theta}}_O) + \lambda \widehat{\mathbf{W}}_O \rangle \\ &= \underbrace{\langle \boldsymbol{\Pi}_{\mathcal{F}}(\widehat{\boldsymbol{\Theta}}_O - \boldsymbol{\Theta}'), \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\boldsymbol{\Theta}}_O) + \lambda \widehat{\mathbf{W}}_O \rangle}_{I_1} + \underbrace{\langle \boldsymbol{\Pi}_{\mathcal{F}^\perp}(\widehat{\boldsymbol{\Theta}}_O - \boldsymbol{\Theta}'), \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\boldsymbol{\Theta}}_O) + \lambda \widehat{\mathbf{W}}_O \rangle}_{I_2}. \end{aligned} \quad (\text{D.15})$$

Analysis of Term I_1 . Let $\boldsymbol{\gamma}^* = \boldsymbol{\gamma}(\boldsymbol{\Theta}^*)$, $\widehat{\boldsymbol{\gamma}}_O = \boldsymbol{\gamma}(\widehat{\boldsymbol{\Theta}}_O)$ be the vector of (ordered) singular values of $\boldsymbol{\Theta}^*$ and $\widehat{\boldsymbol{\Theta}}_O$, respectively. By the perturbation bounds for singular values, the Weyl's inequality (Weyl, 1912), we have that

$$\max_i |(\boldsymbol{\gamma}^*)_i - (\widehat{\boldsymbol{\gamma}}_O)_i| \leq \|\boldsymbol{\Theta}^* - \widehat{\boldsymbol{\Theta}}_O\|_2 \leq \|\boldsymbol{\Theta}^* - \widehat{\boldsymbol{\Theta}}_O\|_F.$$

Since Lemma D.3 provides the Frobenius norm on the estimation error $\boldsymbol{\Theta}^* - \widehat{\boldsymbol{\Theta}}_O$, we obtain that

$$\max_i |(\boldsymbol{\gamma}^*)_i - (\widehat{\boldsymbol{\gamma}}_O)_i| \leq \frac{2\sqrt{r}}{n\kappa(\mathfrak{X})} \|\boldsymbol{\mathfrak{X}}^*(\boldsymbol{\epsilon})\|_2.$$

If it is assumed that $S = \text{supp}(\boldsymbol{\sigma}^*)$, we have $|S| = r$. The triangle inequality yields that

$$\begin{aligned} \min_{i \in S} |(\widehat{\boldsymbol{\gamma}}_O)_i| &= \min_{i \in S} |(\widehat{\boldsymbol{\gamma}}_O)_i - (\boldsymbol{\gamma}^*)_i + (\boldsymbol{\gamma}^*)_i| \geq -\max_{i \in S} |(\widehat{\boldsymbol{\gamma}}_O - \boldsymbol{\gamma}^*)_i| + \min_{i \in S} |(\boldsymbol{\gamma}^*)_i| \\ &\geq -\frac{2\sqrt{r}}{n\kappa(\mathfrak{X})} \|\boldsymbol{\mathfrak{X}}^*(\boldsymbol{\epsilon})\|_2 + \nu + \frac{2\sqrt{r}}{n\kappa(\mathfrak{X})} \|\boldsymbol{\mathfrak{X}}^*(\boldsymbol{\epsilon})\|_2 \\ &= \nu, \end{aligned}$$

where the inequality on the second line is derived based on the condition that $\min_{i \in S} |(\boldsymbol{\gamma}^*)_i| \geq \nu + 2n^{-1}\sqrt{r}\|\boldsymbol{\mathfrak{X}}^*(\boldsymbol{\epsilon})\|_*/\kappa(\mathfrak{X})$. Based on the definition of oracle estimator (3.2), $\widehat{\boldsymbol{\Theta}}_O \in \mathcal{F}$, which implies $\text{rank}(\widehat{\boldsymbol{\Theta}}_O) = r$. Therefore, we have

$$(\widehat{\boldsymbol{\gamma}}_O)_1 \geq (\widehat{\boldsymbol{\gamma}}_O)_2 \geq \dots \geq (\widehat{\boldsymbol{\gamma}}_O)_r \geq \nu > 0 = (\widehat{\boldsymbol{\gamma}}_O)_{r+1} = (\widehat{\boldsymbol{\gamma}}_O)_m = 0. \quad (\text{D.16})$$

By the definition of Oracle estimator, we have $\widehat{\boldsymbol{\Theta}}_O = \mathbf{U}^* \widehat{\boldsymbol{\Gamma}}_O \mathbf{V}^{*\top}$, where $\widehat{\boldsymbol{\Gamma}}_O$ is the diagonal matrix with $\text{diag}(\widehat{\boldsymbol{\Gamma}}_O) = \widehat{\boldsymbol{\gamma}}_O$. Since $\mathcal{P}_\lambda(\boldsymbol{\Theta}) = \mathcal{Q}_\lambda(\boldsymbol{\Theta}) + \lambda \|\boldsymbol{\Theta}\|_*$, we have

$$\begin{aligned} \boldsymbol{\Pi}_{\mathcal{F}}(\nabla \mathcal{P}_\lambda(\widehat{\boldsymbol{\Theta}}_O)) &= \boldsymbol{\Pi}_{\mathcal{F}}(\nabla \mathcal{Q}_\lambda(\widehat{\boldsymbol{\Theta}}_O) + \lambda \partial \|\widehat{\boldsymbol{\Theta}}_O\|_*) \\ &= \boldsymbol{\Pi}_{\mathcal{F}}(\mathbf{U}^* q'_\lambda(\widehat{\boldsymbol{\Gamma}}_O) \mathbf{V}^{*\top} + \lambda \mathbf{U}^* \mathbf{V}^{*\top} + \lambda \widehat{\mathbf{Z}}_O) \\ &= \mathbf{U}^* \left(q'_\lambda((\widehat{\boldsymbol{\Gamma}}_O)_S) + \lambda \mathbf{I}_r \right) \mathbf{V}^{*\top}, \end{aligned} \quad (\text{D.17})$$

where $\widehat{\mathbf{Z}}_O \in \mathcal{F}^\perp$, $\|\widehat{\mathbf{Z}}_O\|_2 \leq 1$, and $(\widehat{\boldsymbol{\Gamma}}_O)_S \in \mathbb{R}^{r \times r}$ is a diagonal matrix with $\text{diag}((\widehat{\boldsymbol{\Gamma}}_O)_S) = (\widehat{\boldsymbol{\gamma}}_O)_S$. The first equality in (D.17) is based on the definition of $\nabla \mathcal{Q}_\lambda(\cdot)$ and $\partial \|\cdot\|_*$, while the second is to simply project each component into the subspace \mathcal{F} . Since $p_\lambda(t) = q_\lambda(t) + \lambda|t|$, we have $p'_\lambda(t) = q'_\lambda(t) + \lambda t$ for all $t > 0$. Consider the diagonal matrix $q'_\lambda((\widehat{\boldsymbol{\Gamma}}_O)_S) + \lambda \mathbf{I}_r$, we have the i^{th} ($i \in S$) element on the diagonal that

$$\left(q'_\lambda((\widehat{\boldsymbol{\Gamma}}_O)_S) + \lambda \mathbf{I}_r \right)_{ii} = q'_\lambda((\widehat{\boldsymbol{\gamma}}_O)_i) + \lambda = p'_\lambda((\widehat{\boldsymbol{\gamma}}_O)_i).$$

Since $p_\lambda(\cdot)$ satisfies the regularity condition (ii), that $p'_\lambda(t) = 0$ for all $t \geq \nu$, we have $p'_\lambda((\widehat{\gamma}_O)_i) = 0$ for $i \in S$, in light of the fact that $(\widehat{\gamma}_O)_i \geq \nu > 0$. Therefore, the diagonal matrix $q'_\lambda((\widehat{\Gamma}_O)_S) + \lambda \mathbf{I}_r = \mathbf{0}$, substituting which into (D.17) yields

$$\mathbf{\Pi}_{\mathcal{F}}(\nabla \mathcal{P}_\lambda(\widehat{\Theta}_O)) = \mathbf{0}. \quad (\text{D.18})$$

Since $\widehat{\Theta}_O$ is a minimizer of (3.2) over \mathcal{F} , we have the following optimality condition that for all $\Theta' \in \mathbb{R}^{m_1 \times m_2}$,

$$\max_{\Theta'} \langle \mathbf{\Pi}_{\mathcal{F}}(\widehat{\Theta}_O - \Theta'), \nabla \mathcal{L}_n(\widehat{\Theta}_O) \rangle \leq 0. \quad (\text{D.19})$$

Substitute (D.18) and (D.19) into item I_1 , we have for all $\widehat{\mathbf{W}}_O \in \partial \|\widehat{\Theta}_O\|_*$,

$$\begin{aligned} & \max_{\Theta'} \langle \mathbf{\Pi}_{\mathcal{F}}(\widehat{\Theta}_O - \Theta'), \nabla \widetilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O \rangle \\ &= \max_{\Theta'} \langle \mathbf{\Pi}_{\mathcal{F}}(\widehat{\Theta}_O - \Theta'), \nabla \mathcal{L}_n(\widehat{\Theta}_O) \rangle + \max_{\Theta'} \langle \mathbf{\Pi}_{\mathcal{F}}(\widehat{\Theta}_O - \Theta'), \mathbf{\Pi}_{\mathcal{F}}(\nabla \mathcal{P}_\lambda(\widehat{\Theta}_O)) \rangle \\ &\leq 0. \end{aligned} \quad (\text{D.20})$$

Analysis of Term I_2 . By definition of $\nabla \mathcal{Q}_\lambda(\Theta)$, and the condition that $q'_\lambda(\cdot)$ satisfies the regularity condition (iii) in Assumption 3.3, we have the SVD of $\nabla \mathcal{Q}_\lambda(\Theta_O)$ as $\nabla \mathcal{Q}_\lambda(\widehat{\Theta}_O) = \mathbf{U}^* q'_\lambda(\widehat{\Gamma}_O) \mathbf{V}^{*\top}$, where $\widehat{\Gamma}_O \in \mathbb{R}^{r \times r}$ is a diagonal matrix. Projecting $\nabla \mathcal{Q}_\lambda(\widehat{\Theta}_O)$ into \mathcal{F}^\perp yields that

$$\begin{aligned} \mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{Q}_\lambda(\widehat{\Theta}_O)) &= (\mathbf{I}_{m_1} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{U}^* q'_\lambda(\widehat{\Gamma}_O) \mathbf{V}^{*\top} (\mathbf{I}_{m_1} - \mathbf{V}^* \mathbf{V}^{*\top}) \\ &= (\mathbf{U}^* - \mathbf{U}^*) q'_\lambda(\widehat{\Gamma}_O)_{S^c} (\mathbf{V}^{*\top} - \mathbf{V}^{*\top}) \\ &= \mathbf{0}. \end{aligned}$$

Thus,

$$\mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{Q}_\lambda(\widehat{\Theta}_O)) = \mathbf{0}. \quad (\text{D.21})$$

Therefore,

$$I_2 = \langle \mathbf{\Pi}_{\mathcal{F}^\perp}(-\Theta'), \mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{L}_n(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O) \rangle.$$

Moreover, the triangle inequality yields

$$\begin{aligned} \|\nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_2 &\leq \|\nabla \mathcal{L}_n(\Theta^*)\|_2 + \|\nabla \mathcal{L}_n(\Theta^*) - \nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_2 \\ &\leq \|\nabla \mathcal{L}_n(\Theta^*)\|_2 + \|\nabla \mathcal{L}_n(\Theta^*) - \nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_F \\ &\leq \|\nabla \mathcal{L}_n(\Theta^*)\|_2 + \rho(\mathfrak{X}) \|\Theta^* - \widehat{\Theta}_O\|_F, \end{aligned} \quad (\text{D.22})$$

where the second inequality comes from the fact that $\|\nabla \mathcal{L}_n(\Theta^*) - \nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_2 \leq \|\nabla \mathcal{L}_n(\Theta^*) - \nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_F$, while the last inequality is obtained by the restricted strong smoothness (Assumption 3.2), which is equivalent to

$$\|\nabla \mathcal{L}_n(\Theta) - \nabla \mathcal{L}_n(\Theta + \widehat{\Delta}_O)\|_F \leq \rho(\mathfrak{X}) \|\widehat{\Delta}_O\|_F,$$

over the restricted set \mathcal{C} ; since $\mathbf{\Pi}_{\mathcal{F}^\perp}(\widehat{\Delta}_O) = \mathbf{0}$, it is evident that $\widehat{\Delta}_O \in \mathcal{C}$.

Substitute (D.12) of Lemma D.3 into (D.22), we have

$$\left\| \mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{L}_n(\widehat{\Theta}_O)) \right\|_2 \leq \|\nabla \mathcal{L}_n(\widehat{\Theta}_O)\|_2 \leq \|\nabla \mathcal{L}_n(\Theta^*)\|_2 + \frac{2\sqrt{r}\rho(\mathfrak{X})}{n\kappa(\mathfrak{X})} \|\mathfrak{X}^*(\epsilon)\|_2 \leq \lambda,$$

where the last inequality follows from the choice of λ .

By setting $\widehat{\mathbf{Z}}_O = -\lambda^{-1} \mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{L}_n(\widehat{\Theta}_O))$, such that $\widehat{\mathbf{W}}_O = \mathbf{U}^* \mathbf{V}^{*\top} + \widehat{\mathbf{Z}}_O \in \partial \|\widehat{\Theta}_O\|_*$ since $\widehat{\mathbf{Z}}_O$ satisfies the condition $\widehat{\mathbf{Z}}_O \in \mathcal{F}^\perp$, $\|\widehat{\mathbf{Z}}_O\|_2 \leq 1$, we have

$$\mathbf{\Pi}_{\mathcal{F}^\perp}(\nabla \mathcal{L}_n(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O) = \mathbf{0},$$

which implies that

$$I_2 = \langle \Pi_{\mathcal{F}^\perp}(-\Theta'), \mathbf{0} \rangle = 0. \quad (\text{D.23})$$

Substitute (D.20) and (D.23) into (D.15), we obtain (D.14) that

$$\max_{\Theta'} \langle \widehat{\Theta}_O - \Theta', \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O \rangle \leq 0.$$

Now we are going to prove that $\widehat{\Theta}_O = \Theta^*$.

Applying Lemma D.1, we have

$$\tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) \geq \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O), \widehat{\Theta} - \widehat{\Theta}_O \rangle + \frac{\kappa(\mathfrak{X}) - \zeta_-}{2} \|\widehat{\Theta}_O - \widehat{\Theta}\|_F^2, \quad (\text{D.24})$$

$$\tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) \geq \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}), \widehat{\Theta}_O - \widehat{\Theta} \rangle + \frac{\kappa(\mathfrak{X}) - \zeta_-}{2} \|\widehat{\Theta}_O - \widehat{\Theta}\|_F^2. \quad (\text{D.25})$$

On the other hand, because of the convexity of nuclear norm $\|\cdot\|_*$, we obtain

$$\lambda \|\widehat{\Theta}\|_* \geq \lambda \|\widehat{\Theta}_O\|_* + \lambda \langle \widehat{\Theta} - \widehat{\Theta}_O, \widehat{\mathbf{W}}_O \rangle, \quad (\text{D.26})$$

$$\lambda \|\widehat{\Theta}_O\|_* \geq \lambda \|\widehat{\Theta}\|_* + \lambda \langle \widehat{\Theta}_O - \widehat{\Theta}, \widehat{\mathbf{W}}_O \rangle. \quad (\text{D.27})$$

Add (D.24) to (D.27), we obtain

$$0 \geq \underbrace{\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta}_O - \widehat{\Theta} \rangle}_{I_3} + \underbrace{\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta} - \widehat{\Theta}_O \rangle}_{I_4} + (\kappa(\mathfrak{X}) - \zeta_-) \|\widehat{\Theta}_O - \widehat{\Theta}\|_F^2. \quad (\text{D.28})$$

Analysis of Term I_3 . By (D.13), we have

$$\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta} - \widehat{\Theta}_O \rangle \leq \max_{\Theta'} \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta} - \Theta' \rangle \leq 0. \quad (\text{D.29})$$

Therefore $I_3 \geq 0$.

Analysis of Term I_4 . By (D.14), we have

$$\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta}_O - \widehat{\Theta} \rangle \leq \max_{\Theta'} \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}_O) + \lambda \widehat{\mathbf{W}}_O, \widehat{\Theta}_O - \Theta' \rangle \leq 0. \quad (\text{D.30})$$

Therefore $I_4 \geq 0$. Substituting (D.29) and (D.30) into (D.28) yields that

$$(\kappa(\mathfrak{X}) - \zeta_-) \|\widehat{\Theta}_O - \widehat{\Theta}\|_F^2 \leq 0,$$

which holds if and only if

$$\widehat{\Theta}_O = \widehat{\Theta}, \quad (\text{D.31})$$

because $\kappa(\mathfrak{X}) > \zeta_-$.

By Lemma D.3, we obtain the error bound

$$\|\widehat{\Theta} - \Theta^*\|_F = \|\widehat{\Theta}_O - \Theta^*\|_F \leq \frac{2\sqrt{r} \|\Pi_{\mathcal{F}}(\nabla \mathcal{L}_n(\Theta^*))\|_2}{\kappa(\mathfrak{X})},$$

which completes the proof. \square

E. Proof of the Results for Specific Examples

In this section, we provide the detailed proofs for corollaries of specific examples presented in Section 3.2. We will first establish the RSC condition for both examples, followed by proofs of the corollaries and more results on oracle property respecting two specific examples of matrix completion.

Particularly, the proofs include the following components: (i) establish the RSC condition, obtaining $\kappa(\mathfrak{X})$ by which Assumption 3.1 holds with high probability; (ii) estimate $\|\nabla \mathcal{L}_n(\Theta^*)\|_2$ for the choice of the regularity parameter λ ; (iii) establish the RSS condition, obtaining $\rho(\mathfrak{X})$ by which Assumption 3.2 holds with high probability.

E.1. Matrix Completion

As shown in (Candès & Recht, 2012) with various examples, it is insufficient to recover the low-rank matrix, since it is infeasible to recover overly “spiky” matrices which have very few large entries. Some existing work (Candès & Recht, 2012) imposes stringent matrix incoherence conditions to preclude such matrices; these assumptions are relaxed in more recent work (Negahban & Wainwright, 2012; Gunasekar et al., 2014) by restricting the spikiness ratio, which is defined as follows:

$$\alpha_{\text{sp}}(\Theta) = \frac{\sqrt{m_1 m_2} \|\Theta\|_{\infty}}{\|\Theta\|_F}.$$

Assumption E.1. There exists a known α^* , such that

$$\|\Theta^*\|_{\infty} = \frac{\alpha_{\text{sp}}(\Theta^*) \|\Theta^*\|_F}{\sqrt{m_1 m_2}} \leq \alpha^*.$$

For the example of matrix completion, we have the following matrix concentration inequality, which follows from Proof of Corollary 1 in (Negahban & Wainwright, 2012).

Proposition E.2. Let \mathbf{X}_i uniformly distributed on \mathcal{X} , and $\{\xi_k\}_{k=1}^n$ be a finite sequence of independent Gaussian variables with variance σ^2 . There exist constants C_1, C_2 that with probability at least $1 - C_2/M$, we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{X}_i \right\|_2 \leq C_1 \sigma \sqrt{\frac{M \log M}{m_1 m_2 n}}.$$

Furthermore, the following Lemma plays a key role in obtaining faster rate for estimator with nonconvex penalties. Particularly, the following Lemma will provide an upper bound on $\|\Pi_{\mathcal{F}}(\nabla \mathcal{L}_n(\Theta^*))\|_2$.

Lemma E.3. If ξ_i is Gaussian noise with variance σ^2 . \mathcal{S} is a r -dimensional subspace. It holds with probability at least $1 - C_2/M$,

$$\left\| \Pi_{\mathcal{S}} \left(\frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{X}_i \right) \right\|_2 \leq C_1 \sigma \sqrt{\frac{r \log M}{m_1 m_2 n}},$$

where C_1, C_2 are universal constants.

Proof. Proof is provided in Section F.4. □

In addition, we have the following Lemma (Theorem 1 in (Negahban & Wainwright, 2012)), which plays central role in establishing the RSC condition.

Lemma E.4. There are universal constants, $k_1, k_2, C_1, \dots, C_5$, such that as long as $n > C_2 M \log M$, if the following condition is satisfied that

$$\sqrt{m_1 m_2} \frac{\|\Delta\|_{\infty}}{\|\Delta\|_F} \frac{\|\Delta\|_*}{\|\Delta\|_F} \leq \frac{\sqrt{rn}}{k_1 r_1 \sqrt{\log M} + k_2 \sqrt{r_2 M \log M}}, \quad (\text{E.1})$$

we have

$$\left| \frac{\|\tilde{\mathbf{x}}_n(\Delta)\|_2}{\sqrt{n}} - \frac{\|\Delta\|_F}{\sqrt{m_1 m_2}} \right| \leq \frac{7}{8} \frac{\|\Delta\|_F}{\sqrt{m_1 m_2}} \left[1 + \frac{C_1 \alpha_{\text{sp}}(\Delta)}{\sqrt{n}} \right], \quad (\text{E.2})$$

with probability greater than $1 - C_3 \exp(-C_4 M \log M)$.

Proof of Corollary 3.6. With regard to the example of matrix completion, we consider a partially observed setting, *i.e.*, only the entries over the subset \mathcal{X} . A uniform sampling model is assumed that

$$\forall (i, j) \in \mathcal{X}, i \sim \text{uniform}([m_1]), j \sim \text{uniform}([m_2]).$$

Recall that $\hat{\Delta} = \hat{\Theta} - \Theta^*$. In this proof, we consider two cases, depending on if the condition in (E.1) holds or not.

1. The condition in (E.1) does not hold.
2. The condition in (E.1) does hold.

CASE 1. If the condition in (E.1) is violated, it implies that

$$\begin{aligned} \|\widehat{\Delta}\|_F^2 &\leq \sqrt{m_1 m_2} \|\widehat{\Delta}\|_\infty \cdot \|\widehat{\Delta}\|_* \frac{k_1 r_1 \sqrt{\log M} + k_2 \sqrt{r_2 M \log M}}{\sqrt{rn}} \\ &\leq \sqrt{m_1 m_2} (2\alpha^*) (\|\widehat{\Delta}'\|_* + \|\widehat{\Delta}''\|_*) \frac{k_1 r_1 \sqrt{\log M} + k_2 \sqrt{r_2 M \log M}}{\sqrt{rn}} \\ &\leq 12\alpha^* \sqrt{r m_1 m_2} \|\widehat{\Delta}'\|_F \frac{k_1 r_1 \sqrt{\log M} + k_2 \sqrt{r_2 M \log M}}{\sqrt{rn}}, \end{aligned}$$

where $\widehat{\Delta}' = \Pi_{\mathcal{F}}(\widehat{\Delta})$ and $\widehat{\Delta}'' = \Pi_{\mathcal{F}^\perp}(\widehat{\Delta})$, the second inequality follows from $\|\widehat{\Delta}\|_\infty \leq \|\widehat{\Theta}\|_\infty + \|\Theta^*\|_\infty \leq 2\alpha^*$, and the decomposability of nuclear norm that $\|\widehat{\Delta}\|_* = \|\widehat{\Delta}'\|_* + \|\widehat{\Delta}''\|_*$; while the third inequality is based on the cone condition $\|\widehat{\Delta}'\|_* \leq 5\|\widehat{\Delta}''\|_*$ and $\|\widehat{\Delta}'\|_* \leq \sqrt{r}\|\widehat{\Delta}'\|_F$.

Moreover, since $\|\widehat{\Delta}'\|_F \leq \|\widehat{\Delta}\|_F$, we obtain that

$$\frac{1}{\sqrt{m_1 m_2}} \|\widehat{\Delta}\|_F \leq 12\alpha^* \left(k_1 r_1 \sqrt{\frac{\log M}{n}} + k_1 \sqrt{\frac{r_2 M \log M}{n}} \right). \quad (\text{E.3})$$

CASE 2. The condition in (E.1) is satisfied.

As implied by (E.2), we have

$$\frac{\|\mathfrak{X}_n(\Delta)\|_2}{\sqrt{n}} \geq \frac{1}{8} \frac{\|\Delta\|_F}{\sqrt{m_1 m_2}} \left[1 - \frac{C'_1 \alpha_{\text{sp}}(\Delta)}{\sqrt{n}} \right],$$

If $C'_1 \alpha_{\text{sp}}(\widehat{\Delta})/\sqrt{n} > 1/2$, we have

$$\|\widehat{\Delta}\|_F \leq 2C_2 \sqrt{m_1 m_2} \frac{\|\widehat{\Delta}\|_\infty}{\sqrt{n}} \leq 4C_2 \alpha^* \sqrt{\frac{m_1 m_2}{n}}. \quad (\text{E.4})$$

If $C'_1 \alpha_{\text{sp}}(\widehat{\Delta})/\sqrt{n} \leq 1/2$, we have

$$\frac{\|\mathfrak{X}_n(\widehat{\Delta})\|_2^2}{n} \geq \frac{C_6^2}{m_1 m_2} \|\widehat{\Delta}\|_F^2. \quad (\text{E.5})$$

In order to establish the RSC condition, we need to show that (E.5) is equivalent to Assumption 3.1.

$$\begin{aligned} &\mathcal{L}_n(\Theta^* + \widehat{\Delta}) - \mathcal{L}_n(\Theta^*) - \langle \nabla \mathcal{L}_n(\Theta^*), \widehat{\Delta} \rangle \\ &= \frac{1}{2n} \sum_{i=1}^n (\langle \Theta^* + \widehat{\Delta}, \mathbf{X}_i \rangle - y_i)^2 + \frac{1}{2n} \sum_{i=1}^n (\langle \Theta^*, \mathbf{X}_i \rangle - y_i)^2 - \frac{1}{n} \sum_{i=1}^n (\langle \Theta^*, \mathbf{X}_i \rangle - y_i) \langle \mathbf{X}_i, \widehat{\Delta} \rangle \\ &= \frac{\|\mathfrak{X}_n(\widehat{\Delta})\|_2^2}{n}. \end{aligned}$$

Thus, we have that (E.5) establishes the RSC condition, and $\kappa(\mathfrak{X}) = C_6^2/(m_1 m_2)$.

After establishing the RSC condition, what remains is to upper bound $n^{-1} \|\mathfrak{X}^*(\epsilon)\|_2$ and $n^{-1} \|\Pi_{\mathcal{F}_{S_1}}(\mathfrak{X}^*(\epsilon))\|_2$. By Proposition E.2, we have that with high probability,

$$\frac{1}{n} \|\mathfrak{X}^*(\epsilon)\|_2 \leq C_6 \sigma \sqrt{\frac{M \log M}{m_1 m_2 n}}; \quad (\text{E.6})$$

By Lemma E.3, we have that with high probability,

$$\frac{1}{n} \|\Pi_{\mathcal{F}_{S_1}}(\mathfrak{X}^*(\epsilon))\|_2 \leq C_7 \sigma \sqrt{\frac{r_1 \log M}{m_1 m_2 n}}. \quad (\text{E.7})$$

Substituting (E.6) and (E.7) into Theorem 3.4, we have that there exist positive constants C'_1, C'_2 such that

$$\frac{1}{\sqrt{m_1 m_2}} \|\widehat{\Theta} - \Theta^*\|_F \leq C'_1 \sigma r_1 \sqrt{\frac{\log M}{n}} + C'_2 \sigma \sqrt{\frac{r_2 M \log M}{n}}. \quad (\text{E.8})$$

Putting pieces (E.3), (E.4), and (E.8) together, we have

$$\frac{1}{\sqrt{m_1 m_2}} \|\widehat{\Theta} - \Theta^*\|_F \leq \max\{\alpha^*, \sigma\} \left[C_3 r_1 \sqrt{\frac{\log M}{n}} + C_4 \sqrt{\frac{r_2 M \log M}{n}} \right],$$

which completes the proof. \square

Corollary E.5. Under the conditions of Theorem 3.5, suppose \mathbf{X}_i uniformly distributed on \mathcal{X} . These exists positive constants C_1, \dots, C_4 , for any $t > 0$, if $\kappa(\mathfrak{X}) = C_1/(m_1 m_2) > \zeta_-$ and γ^* satisfies

$$\min_{i \in S} |(\gamma^*)_i| \geq \nu + C_2 \sigma \sqrt{r m_1 m_2} \sqrt{\frac{M \log M}{n}},$$

where $S = \text{supp}(\sigma^*)$, for estimator in (2.2) with regularization parameter

$$\lambda \geq C_3 (1 + \sqrt{r}) \sigma \sqrt{\frac{M \log M}{n m_1 m_2}},$$

we have that with high probability, $\widehat{\Theta} = \widehat{\Theta}_O$, which yields that $\text{rank}(\widehat{\Theta}) = \text{rank}(\widehat{\Theta}_O) = \text{rank}(\Theta^*) = r$. In addition, we have

$$\frac{1}{\sqrt{m_1 m_2}} \|\widehat{\Theta} - \Theta^*\|_F \leq C_4 r \sigma \sqrt{\frac{\log M}{n}}. \quad (\text{E.9})$$

Proof of Corollary E.5. As shown in the proof of Corollary 3.6, we have $\kappa(\mathfrak{X}) = C_1/(m_1 m_2)$, together with (E.6) and (E.7), in order to prove Corollary E.5, according to Theorem 3.5, what remains is to obtain $\rho(\mathfrak{X})$ in Assumption 3.2. It can be shown that Assumption 3.2 is equivalent as

$$\frac{\rho(\mathfrak{X})}{2} \|\widehat{\Delta}\|_F^2 \geq \frac{1}{n} \|\mathfrak{X}(\widehat{\Delta})\|_2^2.$$

We consider the following cases depending on if (E.1) holds or not.

CASE 1. If the condition in (E.1) is violated,

$$\frac{1}{n} \|\mathfrak{X}(\widehat{\Delta})\|_F^2 \leq \|\widehat{\Delta}\|_\infty^2 \leq \|\widehat{\Delta}\|_F^2,$$

which implies that $\rho(\mathfrak{X}) = 1$.

CASE 2. The condition in (E.1) is satisfied. As implied by Lemma E.4, when $n \geq C_5^2 \alpha^* \geq C_5^2 \alpha_{\text{sp}}(\widehat{\Delta})$, we have that with high probability, the following holds:

$$\frac{C_6}{m_1 m_2} \|\widehat{\Delta}\|_F^2 \geq \frac{1}{n} \|\mathfrak{X}(\widehat{\Delta})\|_2^2.$$

Thus, $\rho(\mathfrak{X}) = C_6/(m_1 m_2)$, which completes the proof. \square

E.2. Matrix Sensing With Dependent Sampling

In this subsection, we provide the proof for the results on matrix sensing. In particular, we will first establish the RSC condition for the application of matrix sensing, followed by the proof on faster convergence rate and more results on the oracle property.

In order to establish the RSC condition, we need the following lemma (Proposition 1 in (Negahban & Wainwright, 2011)).

Lemma E.6. Consider the sampling operator of Σ -ensemble, it holds with probability at least $1 - 2 \exp(-n/32)$ that

$$\frac{\|\mathfrak{X}(\Delta)\|_2}{\sqrt{n}} \geq \frac{1}{4} \|\sqrt{\Sigma} \text{vec}(\Delta)\|_2 - 12\pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \|\Delta\|_*.$$

In addition, we need the upper bound of $n^{-1} \|\mathfrak{X}^*(\epsilon)\|_2$, as stated in the following Proposition (Lemma 6, (Negahban & Wainwright, 2011)).

Proposition E.7. With high probability, there are universal constants C_1, C_2 and C_3 such that

$$\mathbb{P} \left[\frac{\|\mathfrak{X}^*(\epsilon)\|_2}{n} \geq C_1 \sigma \pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \right] \leq C_2 \exp(-C_3(m_1 + m_2)),$$

where $\pi(\Sigma)^2 = \sup_{\|u\|_2=1, \|v\|_2=1} \text{Var}(u^\top \Sigma v)$.

Proof of Corollary 3.8. To begin with, we need to establish the RSC condition as in Assumption 3.1. According to Lemma E.6, we have that

$$\frac{\|\mathfrak{X}(\hat{\Delta})\|_2}{\sqrt{n}} \geq \frac{\sqrt{\lambda_{\min}(\Sigma)}}{4} \|\hat{\Delta}\|_F - 12\pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \|\hat{\Delta}\|_*.$$

By the decomposability of nuclear norm, we have that

$$\|\hat{\Delta}\|_* = \|\hat{\Delta}'\|_* + \|\hat{\Delta}''\|_* \leq 6\|\hat{\Delta}'\|_* = 6\sqrt{r}\|\hat{\Delta}'\|_F \leq 6\sqrt{r}\|\hat{\Delta}\|_F, \quad (\text{E.10})$$

where $\hat{\Delta}' = \Pi_{\mathcal{F}}(\hat{\Delta})$ and $\hat{\Delta}'' = \Pi_{\mathcal{F}^\perp}(\hat{\Delta})$.

By substituting (E.10) into Proposition E.6, we have that

$$\begin{aligned} \frac{\|\mathfrak{X}(\hat{\Delta})\|_2}{\sqrt{n}} &\geq \frac{\sqrt{\lambda_{\min}(\Sigma)}}{4} \|\hat{\Delta}\|_F - 72\sqrt{r}\pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \|\hat{\Delta}\|_F \\ &= \left[\frac{\sqrt{\lambda_{\min}(\Sigma)}}{4} - 72\sqrt{r}\pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \right] \|\hat{\Delta}\|_F. \end{aligned}$$

Thus, for $n > C_1 r \pi^2(\Sigma) m_1 m_2 / \lambda_{\min}(\Sigma)$ where C_1 is sufficiently large such that

$$72\sqrt{r}\pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}} \right) \leq \frac{\lambda_{\min}(\Sigma)}{8},$$

we have

$$\frac{\|\mathfrak{X}(\hat{\Delta})\|_2}{\sqrt{n}} \geq \frac{\sqrt{\lambda_{\min}(\Sigma)}}{8} \|\hat{\Delta}\|_F,$$

which implies that

$$\frac{\|\mathfrak{X}(\hat{\Delta})\|_2^2}{n} \geq \frac{\lambda_{\min}(\Sigma)}{64} \|\hat{\Delta}\|_F^2.$$

Therefore, $\kappa(\mathfrak{X}) = \lambda_{\min}(\Sigma)/32$ such that the following holds,

$$\frac{\|\mathfrak{X}(\hat{\Delta})\|_2^2}{n} \geq \frac{\kappa(\mathfrak{X})}{2} \|\hat{\Delta}\|_F^2,$$

which establishes the RSC condition for matrix sensing.

On the other hand, we have

$$\|\mathbf{\Pi}_{\mathcal{F}_{S_1}}(\nabla \mathcal{L}_n(\Theta^*))\|_2 = \|\mathbf{U}_{S_1}^* \mathbf{U}_{S_1}^{*\top} \nabla \mathcal{L}_n(\Theta^*) \mathbf{V}_{S_1}^* \mathbf{V}_{S_1}^{*\top}\|_2 = \|\mathbf{U}_{S_1}^{*\top} \nabla \mathcal{L}_n(\Theta^*) \mathbf{V}_{S_1}^*\|_2,$$

where the second inequality follows from the property of left and right singular vectors $\mathbf{U}_{S_1}^*, \mathbf{V}_{S_1}^*$.

It is worth noting that $\mathbf{U}_{S_1}^{*\top} \nabla \mathcal{L}_n(\Theta^*) \mathbf{V}_{S_1}^* \in \mathbb{R}^{r_1 \times r_1}$. By Proposition E.7, we have that

$$\begin{aligned} \|\mathbf{U}^{*\top} \nabla \mathcal{L}_n(\Theta^*) \mathbf{V}^*\|_2 &\leq 2C_0 \sigma \pi(\Sigma) \sqrt{\frac{M}{n}}, \\ \|\mathbf{U}_{S_1}^{*\top} \nabla \mathcal{L}_n(\Theta^*) \mathbf{V}_{S_1}^*\|_2 &\leq 2C_0 \sigma \pi(\Sigma) \sqrt{\frac{r_1}{n}}, \end{aligned} \quad (\text{E.11})$$

which hold with probability at least $1 - C_1 \exp(-C_2 r_1)$.

The upper bound is obtained directed from Theorem 3.4 and (E.11). Thus, we complete the proof. \square

Corollary E.8. Under the condition of Theorem 3.5, for some universal constants C_1, \dots, C_6 if $\kappa(\mathfrak{X}) = C_1 \lambda_{\min}(\Sigma) > \zeta$ and γ^* satisfies

$$\min_{i \in S} |(\gamma^*)_i| \geq \nu + C_2 \sigma \pi(\Sigma) \frac{\sqrt{r}(\sqrt{m_1} + \sqrt{m_2})}{\sqrt{n} \lambda_{\min}(\Sigma)},$$

where $S = \text{supp}(\gamma^*)$, for estimator in (2.2) with regularization parameter

$$\lambda \geq C_3 \left(1 + \frac{\sqrt{r} \lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)}\right) \sigma \pi(\Sigma) \left(\sqrt{\frac{m_1}{n}} + \sqrt{\frac{m_2}{n}}\right),$$

we have that $\widehat{\Theta} = \widehat{\Theta}_O$, which yields that $\text{rank}(\widehat{\Theta}) = \text{rank}(\widehat{\Theta}_O) = \text{rank}(\Theta^*) = r$, with probability at least $1 - C_4 \exp(-C_5(m_1 + m_2))$. In addition, we have

$$\|\widehat{\Theta} - \Theta^*\|_F \leq \frac{C_6 r \pi(\Sigma)}{\sqrt{n} \lambda_{\min}(\Sigma)}. \quad (\text{E.12})$$

Proof of Corollary E.8. The proof follows from the proof of Corollary 3.8 and Theorem 3.5. As shown in the proof of Corollary 3.8, we have $\kappa(\mathfrak{X}) = C_1 \lambda_{\min}(\Sigma)$, together with (E.11), in order to prove Corollary E.8, according to Theorem 3.5, what remains is to obtain $\rho(\mathfrak{X})$ in Assumption 3.2, respecting the example of matrix sensing.

According to Assumption 3.2, we have that $\rho(\mathfrak{X}) = \lambda_{\max}(\mathbf{H}_n)$, where \mathbf{H}_n is the Hessian matrix of $\mathcal{L}_n(\cdot)$. Based on the definition of $\mathcal{L}_n(\cdot)$, we have

$$\mathbf{H}_n = n^{-1} \sum_{i=1}^n \text{vec}(\mathbf{X}_i) \text{vec}(\mathbf{X}_i)^\top.$$

Thus $\mathbb{E}[\mathbf{H}_n] = \Sigma$. By concentration, we have that when n is sufficiently large, with high probability, $\lambda_{\max}(\mathbf{H}_n) \leq 2\lambda_{\max}(\Sigma)$, which is equivalent to $\rho(\mathfrak{X}) \leq 2\lambda_{\max}(\Sigma)$, holding with high probability, where n is sufficiently large. This completes the proof. \square

F. Proof of Auxiliary Lemmas

F.1. Proof of Lemma D.1

Proof. By the restricted strong convexity assumption (Assumption 3.1), we have

$$\mathcal{L}_n(\Theta_2) \geq \mathcal{L}_n(\Theta_1) + \langle \nabla \mathcal{L}_n(\Theta_1), \Theta_2 - \Theta_1 \rangle + \frac{\kappa(\mathfrak{X})}{2} \|\Theta_2 - \Theta_1\|_F^2. \quad (\text{F.1})$$

In the following, we will show the strong smoothness of $\mathcal{Q}_\lambda(\cdot)$, based on the regularity condition (ii), which imposes constraint on the level of nonconvexity of $q_\lambda(\cdot)$. Assume $\gamma_1 = \gamma(\Theta_1), \gamma_2 = \gamma(\Theta_2)$ are the vectors of singular values

of Θ_1, Θ_2 , respectively, and the singular values in γ_1, γ_2 are nonincreasing. For Θ_1, Θ_2 , we have the following singular value decompositions:

$$\begin{aligned}\Theta_1 &= \mathbf{U}_1 \mathbf{\Gamma}_1 \mathbf{V}_1^\top, \\ \Theta_2 &= \mathbf{U}_2 \mathbf{\Gamma}_2 \mathbf{V}_2^\top,\end{aligned}$$

where $\mathbf{\Gamma}_1, \mathbf{\Gamma}_2 \in \mathbb{R}^{m \times m}$ are diagonal matrix with $\mathbf{\Gamma}_1 = \text{diag}(\gamma_1), \mathbf{\Gamma}_2 = \text{diag}(\gamma_2)$. For each pair of singular values of Θ_1, Θ_2 : $((\gamma_1)_i, (\gamma_2)_i)$ where $i = 1, 2, \dots, m$, we have

$$-\zeta_- ((\gamma_1)_i - (\gamma_2)_i)^2 \leq [q'_\lambda((\gamma_1)_i) - q'_\lambda((\gamma_2)_i)]((\gamma_1)_i - (\gamma_2)_i),$$

which is equivalent to

$$\langle (-q'_\lambda(\mathbf{\Gamma}_1)) - (-q'_\lambda(\mathbf{\Gamma}_2)), \mathbf{\Gamma}_1 - \mathbf{\Gamma}_2 \rangle \leq \zeta_- \|\mathbf{\Gamma}_1 - \mathbf{\Gamma}_2\|_F^2,$$

which yields

$$\langle (-\nabla \mathcal{Q}_\lambda(\Theta_1)) - (-\nabla \mathcal{Q}_\lambda(\Theta_2)), \Theta_1 - \Theta_2 \rangle \leq \zeta_- \|\Theta_1 - \Theta_2\|_F^2. \quad (\text{F.2})$$

Since (F.2) is the definition of strongly smoothness of $-\mathcal{Q}(\cdot)$, it can be show to be equivalent to the following inequality that

$$\mathcal{Q}_\lambda(\Theta_2) \geq \mathcal{Q}_\lambda(\Theta_1) + \langle \nabla \mathcal{Q}(\Theta_1), \Theta_2 - \Theta_1 \rangle - \frac{\zeta_-}{2} \|\Theta_2 - \Theta_1\|_F^2. \quad (\text{F.3})$$

Adding up (F.1) and (F.3), we complete the proof. \square

F.2. Proof of Lemma D.2

Proof. By Lemma D.1, we have that

$$\tilde{\mathcal{L}}_{n,\lambda}(\hat{\Theta}) + \lambda \|\hat{\Theta}\|_* - \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*) - \lambda \|\Theta^*\|_* \geq \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \hat{\Theta} - \Theta^* \rangle + \lambda \|\hat{\Theta}\|_* - \lambda \|\Theta^*\|_*. \quad (\text{F.4})$$

For the first term on the RHS in (F.4), we have the following lower bound

$$\begin{aligned}\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \hat{\Theta} - \Theta^* \rangle &= \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \Pi_{\mathcal{F}}(\hat{\Theta} - \Theta^*) \rangle + \langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \Pi_{\mathcal{F}^\perp}(\hat{\Theta} - \Theta^*) \rangle \\ &\geq - \underbrace{\|\Pi_{\mathcal{F}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*))\|_2}_{I_1} \|\Pi_{\mathcal{F}}(\hat{\Theta} - \Theta^*)\|_* \\ &\quad - \underbrace{\|\Pi_{\mathcal{F}^\perp}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*))\|_2}_{I_2} \|\Pi_{\mathcal{F}^\perp}(\hat{\Theta} - \Theta^*)\|_*,\end{aligned} \quad (\text{F.5})$$

where the last inequality follows from Hölder's inequality.

Analysis of term I_1 . It can be shown that $\nabla \mathcal{L}_n(\Theta^*) = -\tilde{\mathfrak{X}}^*(\epsilon)/n$. Based on the condition that $\lambda > 2n^{-1} \|\tilde{\mathfrak{X}}^*(\epsilon)\|_2$, we have that

$$\|\nabla \mathcal{L}_n(\Theta^*)\|_2 \leq \lambda/2. \quad (\text{F.6})$$

Moreover, by condition (iv) in Assumption 3.3 and (F.6), we obtain that

$$\|\Pi_{\mathcal{F}}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*))\|_2 = \|\Pi_{\mathcal{F}}(\nabla \mathcal{L}_n(\Theta^*) + \mathcal{Q}_\lambda(\Theta^*))\|_2 \leq 3\lambda/2.$$

Analysis of term I_2 . Since $\Pi_{\mathcal{F}^\perp}(\Theta^*) = \mathbf{0}$, we have that

$$\|\Pi_{\mathcal{F}^\perp}(\nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*))\|_2 = \|\Pi_{\mathcal{F}^\perp}(\nabla \mathcal{L}_n(\Theta^*))\|_2 \leq \lambda/2. \quad (\text{F.7})$$

Putting pieces (F.6) and (F.7) into (F.5), we obtain

$$\langle \nabla \tilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \hat{\Theta} - \Theta^* \rangle \geq -3\lambda/2 \|\Pi_{\mathcal{F}}(\hat{\Theta} - \Theta^*)\|_* - \lambda/2 \|\Pi_{\mathcal{F}^\perp}(\hat{\Theta} - \Theta^*)\|_*. \quad (\text{F.8})$$

Meanwhile, we have the lower bound on $\lambda\|\widehat{\Theta}\|_* - \lambda\|\Theta\|_*$ that

$$\begin{aligned}\lambda\|\widehat{\Theta}\|_* - \lambda\|\Theta\|_* &= \lambda\|\Pi_{\mathcal{F}}(\widehat{\Theta})\|_* + \lambda\|\Pi_{\mathcal{F}^\perp}(\widehat{\Theta})\|_* - \lambda\|\Theta\|_* \\ &\geq -\lambda\|\Pi_{\mathcal{F}}(\widehat{\Theta} - \Theta^*)\|_* + \lambda\|\Pi_{\mathcal{F}^\perp}(\widehat{\Theta} - \Theta^*)\|_*\end{aligned}\quad (\text{F.9})$$

Adding (F.8) and (F.9) yields that

$$\langle \nabla \widetilde{\mathcal{L}}_{n,\lambda}(\Theta^*), \widehat{\Theta} - \Theta^* \rangle + \lambda\|\widehat{\Theta}\|_* - \lambda\|\Theta\|_* = -5\lambda/2\|\Pi_{\mathcal{F}}(\widehat{\Theta} - \Theta^*)\|_* + \lambda/2\|\Pi_{\mathcal{F}^\perp}(\widehat{\Theta} - \Theta^*)\|_*.\quad (\text{F.10})$$

Due to the fact that $\widehat{\Theta}$ is the global minimizer of (2.2), provided the condition that $\kappa(\mathfrak{X}) > \zeta_-$, we have

$$\widetilde{\mathcal{L}}_{n,\lambda}(\widehat{\Theta}) + \lambda\|\widehat{\Theta}\|_* - \widetilde{\mathcal{L}}_{n,\lambda}(\Theta) - \lambda\|\Theta^*\|_* \leq 0.\quad (\text{F.11})$$

Substituting (F.10) and (F.11) into (F.4), since $\lambda > 0$, we have that

$$\|\Pi_{\mathcal{F}^\perp}(\widehat{\Theta} - \Theta^*)\|_* \leq 5\|\Pi_{\mathcal{F}}(\widehat{\Theta} - \Theta^*)\|_*,$$

which completes the proof. \square

F.3. Proof of Lemma D.3

Proof. $\widehat{\Delta}_O = \widehat{\Theta}_O - \Theta^*$. According to observation model (2.1) and definition of $\mathfrak{X}(\cdot)$, we have

$$\begin{aligned}\mathcal{L}_n(\widehat{\Theta}_O) - \mathcal{L}_n(\Theta^*) &= \frac{1}{2n} \sum_{i=1}^n (y_i - \mathfrak{X}_i(\Theta^* + \widehat{\Delta}_O))^2 - \frac{1}{2n} \sum_{i=1}^n (y_i - \mathfrak{X}_i(\Theta^*))^2 \\ &= \frac{1}{2n} \sum_{i=1}^n (\epsilon_i - \mathfrak{X}_i(\widehat{\Delta}_O))^2 - \frac{1}{2n} \sum_{i=1}^n \epsilon_i^2 \\ &= \frac{1}{2n} \|\mathfrak{X}(\widehat{\Delta}_O)\|_2^2 - \frac{1}{n} \langle \mathfrak{X}^*(\epsilon), \widehat{\Delta}_O \rangle,\end{aligned}$$

where $\mathfrak{X}^*(\epsilon) = \sum_{i=1}^n \epsilon_i \mathbf{X}_i$ is the adjoint of the operator \mathfrak{X} . Because the oracle estimator $\widehat{\Theta}_O$ minimizes $\mathcal{L}_n(\cdot)$ over the subspace \mathcal{F} , while $\Theta^* \in \mathcal{F}$, we have $\mathcal{L}_n(\widehat{\Theta}_O) - \mathcal{L}_n(\Theta^*) \leq 0$, which yields

$$\frac{1}{2n} \|\mathfrak{X}(\widehat{\Delta}_O)\|_2^2 \leq \frac{1}{n} \langle \mathfrak{X}^*(\epsilon), \widehat{\Delta}_O \rangle.\quad (\text{F.12})$$

On the other hand, recall that by the RSC condition (Assumption 3.1), we have

$$\mathcal{L}_n(\Theta + \Delta) \geq \mathcal{L}_n(\Theta) + \langle \nabla \mathcal{L}_n(\Theta), \Delta \rangle + \kappa(\mathfrak{X})/2 \|\Delta\|_F^2,$$

which implies that

$$\frac{1}{2n} \|\mathfrak{X}(\widehat{\Delta}_O)\|_2^2 - \frac{1}{n} \langle \mathfrak{X}^*(\epsilon), \widehat{\Delta}_O \rangle - \langle \nabla \mathcal{L}_n(\Theta^*), \Delta \rangle = \frac{1}{2n} \|\mathfrak{X}(\widehat{\Delta}_O)\|_2^2 \geq \frac{\kappa(\mathfrak{X})}{2} \|\widehat{\Delta}_O\|_F^2.\quad (\text{F.13})$$

Substituting (F.13) into (F.12), we have

$$\frac{\kappa(\mathfrak{X})}{2} \|\widehat{\Delta}_O\|_F^2 \leq \frac{1}{2n} \|\mathfrak{X}(\widehat{\Delta}_O)\|_2^2 \leq \frac{1}{n} \langle \mathfrak{X}^*(\epsilon), \widehat{\Delta}_O \rangle.\quad (\text{F.14})$$

Therefore,

$$\|\widehat{\Delta}_O\|_F^2 \leq \frac{2\langle \Pi_{\mathcal{F}}(\mathfrak{X}^*(\epsilon)), \widehat{\Delta}_O \rangle}{n\kappa(\mathfrak{X})} \leq \frac{2\|\Pi_{\mathcal{F}}(\mathfrak{X}^*(\epsilon))\|_2 \|\widehat{\Delta}_O\|_*}{n\kappa(\mathfrak{X})},$$

where the last inequality is due to Hölder inequality. Moreover, since the rank Δ_O is r , we have the fact that $\|\widehat{\Delta}_O\|_* \leq \sqrt{r} \|\widehat{\Delta}_O\|_F$, which indicates that

$$\|\widehat{\Delta}_O\|_F^2 \leq \frac{2\sqrt{r} \|\Pi_{\mathcal{F}}(\mathfrak{X}^*(\epsilon))\|_2 \cdot \|\widehat{\Delta}_O\|_F}{n\kappa(\mathfrak{X})}.$$

Therefore, we have the following deterministic error bound

$$\|\widehat{\Delta}_O\|_F \leq \frac{2\sqrt{r}\|\Pi_{\mathcal{F}}(\mathfrak{X}^*(\epsilon))\|_2}{n\kappa(\mathfrak{X})} = \frac{2\sqrt{r}\|\Pi_{\mathcal{F}}(\nabla\mathcal{L}_n(\Theta^*))\|_2}{\kappa(\mathfrak{X})},$$

where the last equality results from the fact that $\nabla\mathcal{L}_n(\Theta^*) = -\mathfrak{X}^*(\epsilon)/n$.

Thus, we complete the proof. \square

F.4. Proof of Lemma E.3

In order to prove Lemma E.3, we need the Ahlswede-Winter Matrix Bound. To begin with, we introduce the definition of $\|\cdot\|_{\psi_1}$ and $\|\cdot\|_{\psi_2}$, followed by some established results on $\|\cdot\|_{\psi_1}$ and $\|\cdot\|_{\psi_2}$.

The sub-Gaussian norm of X , denoted by $\|X\|_{\psi_2}$, is defined as follows

$$\|X\|_{\psi_2} = \sup_{p \geq 1} p^{-1/2}(\mathbb{E}|X|^p)^{1/p}.$$

It is known that if $\mathbb{E}[X] = 0$, then $\mathbb{E}[\exp(tX)] \leq \exp(Ct^2\|X\|_{\psi_2}^2)$ for all $t \in \mathbb{R}$.

The sub-Exponential norm of X , denoted by $\|X\|_{\psi_1}$, is defined as follows

$$\|X\|_{\psi_1} = \sup_{p \geq 1} p^{-1}(\mathbb{E}|X|^p)^{1/p}.$$

By (Vershynin, 2010), we have the following Lemma.

Lemma F.1. For Z_1 and Z_2 being two sub-Gaussian random variables, $Z_1 Z_2$ is a sub-exponential random variable with

$$\|Z_1 Z_2\|_{\psi_1} \leq C \max\{\|Z_1\|_{\psi_2}^2, \|Z_2\|_{\psi_2}^2\},$$

where $C > 0$ is an absolute constant.

Theorem F.2 (Ahlswede-Winter Matrix Bound). (Negahban & Wainwright, 2012) Let $\mathbf{Z}_1, \dots, \mathbf{Z}_n$ be random matrices of size $m_1 \times m_2$. Let $\|\mathbf{Z}_i\|_{\psi_1} \leq K$ for all i such that $\|\mathbf{Z}_i\|_{\psi_1}$ is upper bounded by K . Furthermore, we have $\delta_i^2 = \max\{\|\mathbb{E}[\mathbf{Z}_i^\top \mathbf{Z}_i]\|_2, \|\mathbb{E}[\mathbf{Z}_i \mathbf{Z}_i^\top]\|_2\}$, and $\delta^2 = \sum_{i=1}^n \delta_i^2$. Then we have

$$\mathbb{P}\left(\left\|\sum_{i=1}^n \mathbf{Z}_i\right\|_2 \geq t\right) \leq m_1 m_2 \max\left\{\exp\left(-\frac{t^2}{4\delta^2}\right), \exp\left(-\frac{t}{2K}\right)\right\}.$$

Now we are ready to prove Lemma E.3.

Proof of Lemma E.3. Since \mathbf{U}^* and \mathbf{V}^* are singular vectors, for $\mathcal{S} = \mathcal{F}(\mathbf{U}^*, \mathbf{V}^*)$, we have

$$\begin{aligned} \frac{1}{n}\left\|\Pi_{\mathcal{S}}\left(\sum_{i=1}^n \xi_i \mathbf{X}_i\right)\right\|_2 &= \frac{1}{n}\left\|\mathbf{U}^* \mathbf{U}^{*\top} \left(\sum_{i=1}^n \xi_i \mathbf{X}_i\right) \mathbf{V}^* \mathbf{V}^{*\top}\right\|_2 \\ &= \frac{1}{n}\left\|\mathbf{U}^{*\top} \left(\sum_{i=1}^n \xi_i \mathbf{X}_i\right) \mathbf{V}^*\right\|_2. \end{aligned}$$

Recall that $\mathbf{X}_i = \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top$. Let $\mathbf{Y}_i = \epsilon_i \mathbf{X}_i = \epsilon_i \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top$. We have $\|\mathbf{Y}_i\|_{\psi_1} \leq C\sigma^2$. Let $\mathbf{Z}_i = \mathbf{U}^{*\top} \mathbf{Y}_i \mathbf{V}^* \in \mathbb{R}^{r \times r}$. We have

$$\|\mathbf{Z}_i\|_{\psi_1} = \|\mathbf{U}^{*\top} \mathbf{Y}_i \mathbf{V}^*\|_{\psi_1}.$$

Based on the definition of \mathbf{Y}_i , we have that $\|\mathbf{Z}_i\|_{\psi_1} < C\sigma$. By applying Theorem F.1, we have

$$\|\mathbf{Z}_i\|_{\psi_1} \leq C'\sigma^2.$$

Thus, $K = C'\sigma^2$.

Furthermore, we have

$$\begin{aligned}\mathbb{E}[\mathbf{Z}_i \mathbf{Z}_i^\top] &= \mathbb{E}[\mathbf{U}^{*\top} \mathbf{Y}_i \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{Y}_i^\top \mathbf{U}^*] = \mathbb{E}[\epsilon_i^2 \mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^*] \\ &= \sigma^2 \mathbb{E}[\mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^*]\end{aligned}$$

Based on the definition of spectral norm, we have

$$\begin{aligned}\|\mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^*\|_2 &= \max_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^* \mathbf{a} \\ &= \max_{\|\mathbf{b}\|_2=1} \mathbf{b}^\top \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{b},\end{aligned}$$

where the second equality follows by setting $\mathbf{b} = \mathbf{U}^* \mathbf{a} \in \mathbb{R}^{m_1}$. In addition, we have

$$\mathbf{b}^\top \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{b} = \mathbf{b}_{j(i)}^\top \mathbf{v}_k^* \mathbf{v}_k^{*\top} \mathbf{b}_{j(i)} = \mathbf{b}_{j(i)}^2 \|\mathbf{v}_k^*\|_2^2,$$

where \mathbf{v}_k^* is the k -th row of \mathbf{V}^* . Thus

$$\begin{aligned}\|\mathbb{E}[\mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^*]\|_2 &= \left\| \frac{1}{m_1 m_2} \sum_{j=1}^{m_1} \sum_{k=2}^{m_2} \mathbf{U}^{*\top} \mathbf{e}_j \mathbf{e}_k^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_k \mathbf{e}_j^\top \mathbf{U}^* \right\|_2 \\ &= \frac{1}{m_1 m_2} \max_{\|\mathbf{a}\|_2=1} \mathbf{a}^\top \sum_{j=1}^{m_1} \sum_{k=2}^{m_2} \mathbf{U}^{*\top} \mathbf{e}_j \mathbf{e}_k^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_k \mathbf{e}_j^\top \mathbf{U}^* \mathbf{a} \\ &= \frac{1}{m_1 m_2} \max_{\|\mathbf{b}\|_2=1} \sum_{j=1}^{m_1} \sum_{k=2}^{m_2} b_j^2 \|\mathbf{v}_k^*\|_2^2.\end{aligned}$$

Since $\sum_{j=1}^{m_1} b_j^2 = 1$ and $\sum_{k=2}^{m_2} \|\mathbf{v}_k^*\|_2^2 = \|\mathbf{V}^*\|_F^2 = r$, we obtain that

$$\|\mathbb{E}[\mathbf{U}^{*\top} \mathbf{e}_{j(i)} \mathbf{e}_{k(i)}^\top \mathbf{V}^* \mathbf{V}^{*\top} \mathbf{e}_{k(i)} \mathbf{e}_{j(i)}^\top \mathbf{U}^*]\|_2 = \frac{r}{m_1 m_2}.$$

Therefore, we have

$$\|\mathbb{E}[\mathbf{Z}_i \mathbf{Z}_i^\top]\|_2 = \frac{\sigma^2 r}{m_1 m_2},$$

and the same result also applies to $\|\mathbb{E}[\mathbf{Z}_i^\top \mathbf{Z}_i]\|_2$.

By applying Theorem F.2, we obtain that

$$\mathbb{P}\left(\left\|\sum_{i=1}^n \xi_i \mathbf{Z}_i\right\|_2 \geq t\right) \leq m_1 m_2 \max\left\{\exp\left(-\frac{m_1 m_2 t^2}{4n\sigma^2 r}\right), \exp\left(-\frac{t}{2\sigma^2}\right)\right\}.$$

Thus, with probability at least $1 - C_2 M^{-1}$, we have

$$\left\|\sum_{i=1}^n \xi_i \mathbf{Z}_i\right\|_2 \leq C_1 \sigma \sqrt{\frac{nr \log M}{m_1 m_2}}$$

where $M = \max(m_1, m_2)$. It immediately implies that

$$\left\|\frac{1}{n} \sum_{i=1}^n \xi_i \mathbf{Z}_i\right\|_2 \leq C_1 \sigma \sqrt{\frac{r \log M}{m_1 m_2 n}}, \quad (\text{F.15})$$

which completes the proof. \square