

# Achieving Oracle Property for Low-rank Covariance Matrix Estimation from Quadratic Measurements

Anonymous CPAL submission

The covariance matrix is a fundamental statistical quantity that characterizes the linear relationships among multiple random variables. The most common approach in estimating a covariance matrix is based on measurements assuming access to all the variables. However, this assumption is not practical for a dynamic and resource-constrained environment. Compressive covariance sketching offers a viable alternative to estimate the covariance matrix based on compressed measurements, which largely reduces the storage and computational requirements. In this paper, we introduce a novel method for low-rank covariance matrix estimation from quadratic measurements based on non-convex regularization and propose a proximal gradient homotopy method to solve the proposed estimator. The proposed method is theoretically proven to attain superior statistical rates to estimators based on the convex nuclear norm, achieving the oracle rate. The theoretical advancements are corroborated by extensive numerical experiments, highlighting the efficacy and robustness of our method.

## 1. Introduction

The covariance matrix, which quantifies linear correlations among multiple random variables, is a cornerstone of statistical analysis [1–4]. Despite its significance across numerous data and signal analysis applications [5–7], estimating the covariance matrix is challenging since it is not directly observable and can only be estimated from samples. The most naive approach to estimating the covariance matrix is the sample covariance matrix  $\Sigma_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$ , where  $\{\mathbf{x}_t \in \mathbb{R}^d\}_{t=1}^T$  denote  $T$  samples [8–10]. In many situations, directly measuring  $\mathbf{x}_t$  can be too costly, particularly when there are limitations on computational and storage resources. In these cases, a practical solution is to compress the measurements  $\mathbf{x}_t$  into a lower-dimensional representation. The covariance can then be estimated using these compressed measurements.

The Compressive Covariance Sensing (or Compressive Covariance Sketching) (CCS), adapts the principles of compressive sensing for covariance matrix estimation, providing a viable approach to inferring the covariance matrix through compressed measurements [11–16]. In this paper, we study a quadratic measurement model [15, 16] as follows:

$$y_i = \frac{1}{T} \sum_{t=1}^T y_{it} + \eta_i, \quad \text{where } y_{it} = (\mathbf{a}_i^\top \mathbf{x}_t)^2,$$

where  $\mathbf{a}_i$  is a random sensing vector. It is easy to verify that  $y_i = \mathbf{a}_i^\top \Sigma_T \mathbf{a}_i + \eta_i$ . CCS is particularly advantageous when dealing with high-dimensional data, as it significantly reduces both storage and computational requirements [17]. The primary challenge in the CCS problem lies in its inherent ill-posedness: the number of available measurements is insufficient relative to the degrees of freedom of the covariance matrix [15]. To address this issue, exploiting structural information of the covariance matrix—where a small number of components capture most of the data variability—has proven highly effective, dramatically reducing the dimensionality of the estimation problem.

This paper studies the situation where the covariance matrix exhibits a low-rank property [18, 19]. This structural property is inspired by the advancements in phase retrieval [20–25], a field dedicated to recovering rank-one covariance matrices from quadratic measurements. Building on this

foundation, our proposed recovery algorithm extends the PhaseLift [21] methodology to handle low-rank matrices effectively. A direct approach is through rank minimization to achieve accurate recovery of a low-rank matrix. However, this approach is computationally intractable, leading to the development of various surrogate loss functions, which include the nuclear norm [26], Schatten- $p$  norm [27], and weighted nuclear norm [28], among others. The nuclear norm approach has been praised for its close convex approximation to matrix rank, facilitating global optimality [29, 30]. A more advanced variant, weighted nuclear norm, assigns different weights to different singular values, enhancing estimator stability and accuracy over conventional nuclear norm penalty [31]. Despite their effectiveness, these  $\ell_1$ -homotopy penalties introduce a non-negligible bias into estimators, which inevitably compromises accuracy. In contrast, non-convex penalties, such as the Smoothly Clipped Absolute Deviation (SCAD) penalty [32] and the Minimax Concave Penalty (MCP) [33], have demonstrated promising empirical results [34–36].

In this paper, we study the CCS problem by assuming the covariance matrix has a low-rank structure. We propose a novel method for low-rank CCS by incorporating a non-convex penalty. We propose a proximal gradient homotopy method to solve the proposed estimator. Our theoretical analysis confirms that the proposed estimator achieves superior statistical convergence rates compared to traditional nuclear norm-penalized estimators. Under reasonable assumptions, our approach demonstrates the oracle property, ensuring precise recovery of the true rank of the covariance matrix along with improved convergence rates. The efficacy of our proposed method is further substantiated by extensive simulation studies, showcasing its superior performance relative to existing methods. Our main contributions are summarized as follows:

- We propose a new approach for estimating low-rank covariance matrices from a quadratic measurement model using a non-convex penalty, which aids in obtaining unbiased estimates. A proximal gradient homotopy algorithm solves the proposed estimator.
- We rigorously demonstrate the statistical properties of the approximate solution generated by our proposed algorithm. Under weak assumptions, we prove that the proposed estimator achieves the oracle statistical rate in the Frobenius norm.
- Our theoretical findings are corroborated through extensive numerical experiments on synthetic datasets, showcasing the effectiveness and robustness of our proposed method.

## 2. Compressive Covariance Sketching

This section presents a unified framework for low-rank covariance matrix sketching. Subsequently, we formulate our estimator with the non-convex penalty and provide several practical illustrations of the covariance matrix sketching model.

### 2.1. The Quadratic Measurement Model

We consider the quadratic measurement model defined as follows:

$$y_i = \mathbf{a}_i^\top \Sigma_T \mathbf{a}_i + \eta_i, \quad i = 1, \dots, m,$$

where  $\{y_i\}_{i=1}^m$  are the observed measurements,  $\Sigma_T$  denotes the sample covariance matrix,  $\{\mathbf{a}_i\}_{i=1}^m$  are the sensing vectors,  $\{\eta_i\}_{i=1}^m$  represents the noise components, and  $m$  is the total number of measurements. We assume that the noise term  $\{\eta_i\}_{i=1}^m$  are independent and sampled from a sub-Gaussian distribution with zero mean and variance proxy  $\delta^2$ . To facilitate further analysis, we decompose the sample covariance matrix  $\Sigma_T$  into the true covariance matrix  $\Sigma^*$  and a bias term  $\mathbf{E}$ , such that  $\Sigma_T = \Sigma^* + \mathbf{E}$ . For clarity, we define the sensing matrix  $\mathbf{A}_i := \mathbf{a}_i \mathbf{a}_i^\top$ . Observing that  $\mathbf{a}_i^\top \Sigma_T \mathbf{a}_i = \langle \mathbf{A}_i, \Sigma_T \rangle$ , we rewrite the measurement model as  $y_i = \langle \mathbf{A}_i, \Sigma_T \rangle + \eta_i$ . Let  $\mathbf{y} := [y_1, \dots, y_m]^\top$  be the measurement vector and  $\boldsymbol{\eta} := [\eta_1, \dots, \eta_m]^\top$  the noise vector. We introduce the linear operator  $\mathcal{A}(\cdot) : \mathbb{R}^{d \times d} \mapsto \mathbb{R}^m$ , defined by  $\mathcal{A}(\Sigma_T) = [\langle \mathbf{A}_1, \Sigma_T \rangle, \dots, \langle \mathbf{A}_m, \Sigma_T \rangle]$  for any  $\Sigma_T \in \mathbb{R}_+^{d \times d}$ . With these definitions, the measurement model can be compactly expressed as

84  $\mathbf{y} = \mathcal{A}(\Sigma_T) + \boldsymbol{\eta}$ . Notably, the quadratic measurement model—previously explored by [15, 21]—  
 85 offers a significantly simpler implementation and greater computational efficiency than models that  
 86 utilize full-rank measurement matrices with independently and identically distributed entries.

## 87 2.2. The Proposed Estimator

88 We denote the matrix to be estimated as  $\Sigma \in \mathbb{R}^{d \times d}$  and propose the following regularized least  
 89 squares type optimization estimator

$$\min_{\Sigma \succeq 0} \left\{ \frac{1}{2m} \|\mathbf{y} - \mathcal{A}(\Sigma)\|_2^2 + \sum_{i=1}^d p_\lambda(\sigma_i(\Sigma)) \right\}, \quad (1)$$

90 where  $\sigma_i(\Sigma)$  denotes the  $i$ -th largest eigenvalue of  $\Sigma$ , and  $p_\lambda(\cdot)$  is a non-convex penalty function  
 91 governed by a regularization parameter  $\lambda > 0$ . As detailed in Assumption 1, we impose specific  
 92 conditions on it.

93 **Assumption 1.** *The function  $p_\lambda(t)$  defined on  $[0, +\infty)$  can be decomposed as  $p_\lambda(t) = \lambda t + q_\lambda(t)$ , and the  
 94 functions  $p_\lambda(t)$  and  $q_\lambda(t)$  satisfy the following conditions:*

- 95 (a) *For all  $t_1 \geq t_2 \geq 0$ , the derivative satisfies  $0 \leq p'_\lambda(t_1) \leq p'_\lambda(t_2) \leq \lambda$  and  $\lim_{t \rightarrow 0} p'_\lambda(t) = \lambda$ ;*
- 96 (b)  *$q'_\lambda(t)$  is monotonic and Lipschitz continuous. Specifically, for  $t_1 \geq t$ , there exists a constant  $\zeta^- \geq 0$   
 97 such that  $q'_\lambda(t_1) - q'_\lambda(t) \geq -\zeta^-(t_1 - t)$ .*
- 98 (c)  *$p_\lambda(t)$  is non-decreasing with  $p_\lambda(0) = 0$  and is differentiable almost everywhere on  $(0, +\infty)$ ;*
- 99 (d) *There exists an  $\alpha > 0$  such that  $p'_\lambda(t) = 0$  for  $t \geq \alpha$ ;*

100 According to [32], a well-designed penalty function should induce an estimator exhibiting three  
 101 desirable properties: unbiasedness, sparsity, and continuity. These properties align with the con-  
 102 ditions outlined in Assumption 1. Several functions satisfying Assumption 1 have been explored.  
 103 Below, we present a few representative examples:

- 104 • **Smooth clipped absolute deviation penalty:** This penalty, due to [32], is given by

$$p'_\lambda(t) := \begin{cases} \lambda, & \text{for } 0 < t \leq \lambda, \\ \frac{b\lambda - t}{b-1}, & \text{for } \lambda \leq t \leq b\lambda, \\ 0, & \text{for } t \geq b\lambda, \end{cases}$$

105 where  $b > 2$  is an additional tuning parameter. The authors in [32] proposed  $b = 3.7$   
 106 through a Bayesian rationale, applicable when the variable dimension is less than 100.

- 107 • **Minimax concave penalty regularizer:** This penalty, due to [33], is defined as follows:

$$p_\lambda(t) := \text{sign}(t) \lambda \cdot \int_0^{|t|} \left(1 - \frac{z}{\lambda b}\right)_+ dz$$

108 for some  $b > 0$ .

## 109 2.3. Applications

110 Low-rank CCS serves as a foundational technique in a multitude of applications across diverse do-  
 111 mains. These applications range from traffic data monitoring [37] and array signal processing [38]  
 112 to collaborative filtering [39, 40] and metric learning [41, 42]. In these fields, the primary focus is  
 113 often on extracting second-order or higher-order statistics—such as the power spectrum and cyclic  
 114 power spectrum—rather than directly reconstructing random signals [43–45]. Here, we briefly re-  
 115 view two representative examples: PhaseLift and spectrum estimation.

- 116 • **PhaseLift:** PhaseLift is a pivotal method for addressing the phase retrieval problem, which  
 117 aims to recover a complex vector  $\mathbf{x} \in \mathbb{C}^d$  from magnitude-only measurements. Specifically,

---

**Algorithm 1: The Projection Proximal Gradient Algorithm**


---

**Input:**  $\lambda, \varepsilon, \mathbf{y}, \{\mathbf{a}_i\}_{i=1}^m$   
**1 Initialize**  $\eta, \delta, \Sigma_0 = \mathbf{0}, \phi_k = 0$ ;  
**2 for**  $k = 0, 1, 2, \dots, K - 1$  **do**  
    $\lambda_{k+1} = \eta \lambda_k, \varepsilon_{k+1} = \frac{\lambda_k}{4}$ ;  
    $(\Sigma_{k+1}, \phi_{k+1}) \leftarrow \text{ProxGrad}(\lambda_{k+1}, \varepsilon_{k+1}, \Sigma_k, \phi_k)$ ;  
    $\mathbf{U}_{k+1} \Psi_{k+1} \mathbf{U}_{k+1}^\top \leftarrow \text{EigDec}(\Sigma_{k+1})$ ;  
    $\Psi_{k+1} = \max(0, \Psi_{k+1})$ ;  $\Sigma_{k+1} = \mathbf{U}_{k+1} \Psi_{k+1} \mathbf{U}_{k+1}^\top$ ;  
**7 end**  
**Output:**  $\Sigma_K$

---

phase retrieval involves reconstructing  $\mathbf{x}$  using  $m$  intensity measurements of the form  $b_i = |\langle \mathbf{x}, \mathbf{z}_i \rangle|^2$ . PhaseLift reformulates this problem by lifting it into a high-dimensional space, transitioning from recovering the vector  $\mathbf{x}$  to recovering the matrix  $\mathbf{X} = \mathbf{x}\mathbf{x}^\text{H}$ . Here,  $\mathbf{x}^\text{H}$  denotes the conjugate transpose of  $\mathbf{x}$ , and  $\mathbf{X}$  is a rank-1 complex Hermitian matrix. This lifting technique effectively transforms the original nonlinear problem into a tractable linear problem through convex optimization. For a comprehensive discussion on the equivalence conditions between the original and lifted formulations, we refer readers to [21].

• **Spectrum Estimation:** Spectrum estimation is an essential technique for analyzing the frequency characteristics of signals within stochastic processes. Its primary objective is to accurately recover the power spectrum (or power spectral density) of a signal, which is intrinsically linked to its autocorrelation function. Specifically, the task of spectrum sensing involves estimating the power spectrum of a wide-sense stationary analog signal  $x(t)$ . The system under consideration employs a sampling architecture composed of  $M$  parallel branches. Each  $i$ -th branch modulates the input signal  $x(t)$  with a potentially complex-valued periodic waveform  $p_i(t)$  that has a period of  $NT$ . This modulation is followed by an integrate-and-dump circuit operating with the same period  $NT$ , effectively reducing the sampling rate to  $\frac{1}{N}$  of the Nyquist rate. The output of the  $i$ -th branch at the  $k$ -th sampling interval is given by

$$y_i[k] = \frac{1}{NT} \int_{kNT}^{(k+1)NT} p_i(t)x(t)dt.$$

Leveraging this relationship, the goal is to reconstruct the power spectrum of  $x(t)$  using the collected samples  $\{y_i[k]\}_{i,k}$  [12]. This approach facilitates efficient spectrum estimation by exploiting the multi-branch sampling architecture and the properties of wide-sense stationary signals.

### 3. Optimization Algorithm

We propose solving problem (1) via a projection proximal gradient algorithm. Define  $f(\Sigma) = \frac{1}{2m} \|\mathbf{y} - \mathcal{A}(\Sigma)\|_2^2$ . Given Assumption 1, the penalty term  $\sum_{i=1}^d p_\lambda(\sigma_i(\Sigma))$  can be expressed as a combination of  $\lambda \|\Sigma\|_*$  and  $Q_\lambda(\Sigma)$ , where  $Q_\lambda(\Sigma) = \sum_{i=1}^d q_\lambda(\sigma_i(\Sigma))$  is the concave component. Define  $\tilde{f}(\Sigma) = f(\Sigma) + Q_\lambda(\Sigma)$ . Consequently, the optimization problem is reformulated as follows:

$$\min_{\Sigma} \quad \tilde{f}(\Sigma) + \lambda \|\Sigma\|_* + \mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\Sigma), \quad (2)$$

where  $\mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\Sigma)$  is the indicator function for the semidefinite cone, taking value zero if  $\Sigma \succeq \mathbf{0}$  and infy otherwise.

To derive the estimate  $\hat{\Sigma}$  from problem (2), we develop a projection proximal gradient method with backtracking line search [46, 47]. Starting from  $\Sigma_0 = \mathbf{0}$  for simplicity, we generate a sequence  $\{\Sigma_k\}_{k \geq 0}$  via the proximal gradient update described below.

---

**Algorithm 2:**  $(\Sigma_j, \phi) \leftarrow \text{ProxGrad}(\lambda, \varepsilon, \Sigma, \phi)$

---

**Input:**  $\lambda, \varepsilon, \Sigma, \phi$   
**1 Initialize**  $j = 0, \Sigma_j = \Sigma$ ;  
**2 while**  $\omega(\tilde{\Sigma}_j) > \varepsilon$  **do**  
**3**      $j = j + 1$ ;  
**4**      $\Sigma_j \leftarrow \arg \min_{\Sigma} \tilde{g}(\Sigma \mid \Sigma_{j-1})$ ;  
**5**     **while**  $\langle \nabla \tilde{f}(\Sigma_{j-1}), \Sigma_j - \Sigma_{j-1} \rangle + \frac{\phi}{2} \|\Sigma - \Sigma_j\|_F^2 < 0$  **do**  
**6**          $\Sigma_j \leftarrow \arg \min_{\Sigma} \tilde{g}(\Sigma \mid \Sigma_{j-1})$ ;  
**7**          $\phi = 2\phi$ ;  
**8**     **end**  
**9**      $\phi = \frac{\phi}{2}$ ;  
**10 end**  
**Output:**  $(\Sigma_j, \phi)$

---

At the  $(k + 1)$ -th iteration, we construct an isotropic quadratic approximation of  $\tilde{f}(\Sigma)$  at a given matrix  $\Sigma_k$  as follows:

$$\tilde{g}(\Sigma \mid \Sigma_k) = \tilde{f}(\Sigma_k) + \langle \nabla \tilde{f}(\Sigma_k), \Sigma - \Sigma_k \rangle + \frac{\phi}{2} \|\Sigma - \Sigma_k\|_F^2 + \lambda \|\Sigma\|_* + \mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\Sigma), \quad (3)$$

where  $\phi$  is an appropriately chosen quadratic coefficient. The next iterate  $\Sigma_{k+1}$  is then obtained by solving the optimization problem:

$$\Sigma_{k+1} = \arg \min_{\Sigma} \tilde{g}(\Sigma \mid \Sigma_k)$$

Suppose  $\hat{\Sigma}$  is an exact minimizer. It satisfies the optimality condition:

$$\langle \hat{\Sigma} - \Sigma, \nabla \tilde{f}(\hat{\Sigma}) + \lambda \widehat{\mathbf{W}} + \widehat{\mathbf{E}} \rangle \leq 0,$$

where  $\widehat{\mathbf{W}} \in \partial \|\hat{\Sigma}\|_*$  and  $\widehat{\mathbf{E}} \in \partial \mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\hat{\Sigma})$  are sub-gradients of the nuclear norm and the indicator function, respectively. However, obtaining the exact solution  $\hat{\Sigma}$  can be computationally intensive. Therefore, we define an approximate solution to (3).

**Definition 2.** An approximate solution  $\tilde{\Sigma}$  to (3) is said to be  $\varepsilon$ -optimal if its sub-optimal measure  $\omega(\tilde{\Sigma})$  satisfies  $\omega(\tilde{\Sigma}) \leq \varepsilon$ , where

$$\begin{aligned} \omega(\tilde{\Sigma}) &= \min_{\tilde{\mathbf{W}} \in \partial \|\tilde{\Sigma}\|_*, \tilde{\mathbf{E}} \in \partial \mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\tilde{\Sigma})} \max_{\Sigma'} \left\{ \frac{\langle \tilde{\Sigma} - \Sigma', (\nabla \tilde{f}(\tilde{\Sigma}) + \lambda \tilde{\mathbf{W}} + \tilde{\mathbf{E}}) \rangle}{\|\tilde{\Sigma} - \Sigma'\|_*} \right\} \\ &= \min_{\tilde{\mathbf{W}} \in \partial \|\tilde{\Sigma}\|_*, \tilde{\mathbf{E}} \in \partial \mathbb{I}_{\{\Sigma \succeq \mathbf{0}\}}(\tilde{\Sigma})} \left\{ \|\nabla \tilde{f}(\tilde{\Sigma}) + \lambda \tilde{\mathbf{W}} + \tilde{\mathbf{E}}\|_2 \right\}, \end{aligned}$$

and the second equality follows from the duality between the nuclear norm  $\|\cdot\|_*$  and the spectral norm  $\|\cdot\|_2$ .

## 4. Main Theories

In this section, we first introduce some technical backgrounds and assumptions. Following these, we establish the statistical convergence rate of our proposed covariance estimator and the iteration complexity of the proposed algorithm. The proofs are provided in the supplementary material.

## 4.1. Backgrounds and Assumptions

We consider a symmetric positive semi-definite matrix  $\Sigma^* \in \mathbb{R}^{d \times d}$  whose non-zero eigenvalues are collected in the vector  $\sigma^* \in \mathbb{R}^r$ . The eigen-decomposition of  $\Sigma^*$  is given by  $\Sigma^* = \mathbf{U}^* \mathbf{\Gamma}^* \mathbf{U}^{*\top}$ , where  $\mathbf{U}^* \in \mathbb{R}^{d \times r}$  is a matrix whose columns are orthonormal eigenvectors corresponding to the non-zero eigenvalues in  $\sigma^*$ , and  $\mathbf{\Gamma}^* = \text{diag}(\sigma^*) \in \mathbb{R}^{r \times r}$  is a diagonal matrix containing these eigenvalues. Based on this decomposition, we define two complementary subspaces  $\mathcal{S}$  and  $\mathcal{S}^\perp$  as follows:

$$\begin{aligned}\mathcal{T}(\mathbf{U}^*) &:= \{\Delta \mid \text{row}(\Delta) \subseteq \mathbf{U}^* \text{ and } \text{col}(\Delta) \subseteq \mathbf{U}^*\}, \\ \mathcal{T}^\perp(\mathbf{U}^*) &:= \{\Delta \mid \text{row}(\Delta) \perp \mathbf{U}^* \text{ and } \text{col}(\Delta) \perp \mathbf{U}^*\}.\end{aligned}$$

Corresponding to these subspaces, we define two projection operators  $\mathcal{P}(\cdot)$  and  $\mathcal{P}^\perp(\cdot)$  that project matrices onto  $\mathcal{T}$  and  $\mathcal{T}^\perp$ , respectively. For any matrix  $\Delta \in \mathbb{R}^{d \times d}$ , these projections are explicitly given by:

$$\mathcal{P}(\Delta) = \mathbf{U}^* \mathbf{U}^{*\top} \Delta \mathbf{U}^* \mathbf{U}^{*\top}, \mathcal{P}^\perp(\Delta) = (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \Delta (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}),$$

where  $\mathbf{I}$  denotes the identity matrix of the appropriate dimension. We then define a subset  $\mathcal{B}$ , which forms a cone representing a restricted set of directions:

$$\mathcal{B} = \{\Delta \in \mathbb{R}^{d \times d} \mid \|\mathcal{P}^\perp(\Delta)\|_* \leq 5 \|\mathcal{P}(\Delta)\|_*\}.$$

**Assumption 3** (Restricted Strong Convexity). *For the empirical loss function  $f(\cdot)$ , there exists some  $\rho^- > 0$  such that, for all  $\Sigma_2 - \Sigma_1 \in \mathcal{B}$ ,*

$$f(\Sigma_2) \geq f(\Sigma_1) + \langle \nabla f(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle + \frac{\rho^-}{2} \|\Sigma_2 - \Sigma_1\|_{\text{F}}^2.$$

**Assumption 4** (Restricted Strong Smoothness). *For the empirical loss function  $f(\cdot)$ , there exists some  $\infty > \rho^+ \geq \rho^- > 0$  such that, for all  $\Sigma_2 - \Sigma_1 \in \mathcal{B}$ ,*

$$f(\Sigma_1) + \langle \nabla f(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle + \frac{\rho^+}{2} \|\Sigma_2 - \Sigma_1\|_{\text{F}}^2 \geq f(\Sigma_2).$$

Previous research has conclusively demonstrated that random Gaussian measurements satisfy the conditions of restricted strong convexity and restricted strong smoothness [48].

## 4.2. Statistical Properties

We now present our main theorem, which establishes a deterministic error bound for the estimator in the quadratic measurement model.

**Theorem 5.** *Suppose that Assumptions 1, 3, and 4 hold. If  $\lambda \geq 2 \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma))\|_2 / m$  and  $\rho^- > \zeta^-$ , the optimal solution  $\hat{\Sigma}$  satisfies:*

$$\|\hat{\Sigma} - \Sigma^*\|_{\text{F}} \leq \frac{\vartheta \sqrt{|\mathcal{S}_1|}}{\rho^- - \zeta^-} + \frac{3\lambda \sqrt{|\mathcal{S}_2|}}{\rho^- - \zeta^-},$$

where  $|\mathcal{S}_1|$  and  $|\mathcal{S}_2|$  denote the cardinality of the sets  $\mathcal{S}_1 = \{i \mid \sigma_i^* \geq \alpha\}$  and  $\mathcal{S}_2 = \{i \mid \alpha > \sigma_i^* > 0\}$  respectively,  $\vartheta = \|\mathcal{P}_{\mathcal{S}_1}(\nabla f(\Sigma^*))\|_2$ , and  $\mathcal{P}_{\mathcal{S}_1}$  is the projection onto the subspace of  $\mathcal{T}$  associated with  $\mathcal{S}_1$ .

**Remark 6.** *The oracle rate refers to the statistical convergence rate of the oracle estimator, which knows the true rank subspace  $\mathcal{T}(\mathbf{U}^*)$  in advance. The oracle estimator  $\hat{\Sigma}^O$  is defined as*

$$\hat{\Sigma}^O = \arg \min_{\Sigma: \Sigma \in \mathcal{T}(\mathbf{U}^*)} f(\Sigma).$$

The upper bound on the Frobenius norm of the estimation error provided by Theorem 5 consists of two distinct components, each corresponding to different magnitudes of the eigenvalues of the true covariance matrix  $\Sigma^*$ . Specifically, the set  $\mathcal{S}_1$  corresponds to the indices of the larger eigenvalues, while  $\mathcal{S}_2$  includes those of the smaller eigenvalues. By selecting an appropriate value for  $\zeta^-$ , such as  $\zeta^- = \rho^-/2$ , we demonstrate that the proposed estimator achieves faster convergence than the nuclear norm-based estimator. Next, we give the explicit statistical rate of convergence.

198 **Corollary 7** (Oracle Property). Suppose that Assumptions 1, 3, and 4 hold. If  $\rho^- > \zeta^-$ ,  $\lambda$  is selected such  
 199 that

$$\lambda \geq \frac{2 \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2}{m} + \frac{2\sqrt{r}\rho^+ \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2}{m\rho^-}$$

200 and if  $\boldsymbol{\sigma}^*$  satisfies that,

$$\min_{i \in \mathcal{S}} |\sigma_i^*| \geq \alpha + \frac{2\sqrt{r} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2}{m\rho^-},$$

201 where  $\mathcal{S} = \text{supp}(\boldsymbol{\sigma}^*)$ , then the optimal solution  $\hat{\boldsymbol{\Sigma}}$  satisfies:

$$\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{\text{F}} \leq \frac{2 \|\mathcal{P}(\nabla f(\boldsymbol{\Sigma}^*))\|_2 \sqrt{r}}{\rho^-}.$$

202 As established in 7, by appropriately selecting the regularization parameter  $\lambda$  and ensuring that  
 203 the smallest non-zero eigenvalue has a sufficiently large, the proposed estimator in 1 can closely  
 204 approximate the oracle estimator, even in the absence of direct knowledge of the subspace  $\mathcal{T}$ . Con-  
 205 sequently, the estimator in 1 is able to recover the rank of the true covariance matrix  $\boldsymbol{\Sigma}^*$ . It indicates  
 206 the CCS algorithm exhibits the oracle property, enabling near-optimal recovery of the covariance  
 207 matrix, which mirrors the performance of an "oracle" that has full knowledge of the true underlying  
 208 data structure. In practical terms, this means that CCS can produce covariance estimates that are  
 209 almost as accurate as those obtained from full data sampling, even in situations where the number  
 210 of measurements is significantly reduced.

211 Next, we give the explicit statistical rate of convergence under the sub-Gaussian design. First, we  
 212 make the essential assumption about the sensing vector.

213 **Assumption 8.** We assume that the sensing vectors  $\mathbf{a}_i$ 's ( $1 \leq i \leq m$ ) consist of independently and iden-  
 214 tically distributed (i.i.d.) sub-Gaussian random variables. Each vector  $\mathbf{a}_i$  contains  $d$  elements, denoted as  
 215  $(\mathbf{a}_i)_j$  ( $1 \leq j \leq d$ ), which satisfy the following conditions:

$$\mathbb{E}[(\mathbf{a}_i)_j] = 0, \mathbb{E}[(\mathbf{a}_i)_j^2] = 1, \text{ and } \mathbb{E}[(\mathbf{a}_i)_j^4] > 1.$$

216 Assumption 8 indicates that each component of the sensing vector exhibits heavy-tailed behavior.  
 217 As a result, the sensing matrix  $\mathbf{A}_i$  follows a sub-exponential distribution. Let  $\text{vec}(\mathbf{A}_i) \in \mathbb{R}^{d^2}$  rep-  
 218 resent the vectorization of  $\mathbf{A}_i$ . Under Assumption 8,  $\text{vec}(\mathbf{A}_i)$  is distributed according to a sub-  
 219 exponential distribution with mean 0 and variance proxy  $\boldsymbol{\Theta}$ .

220 **Corollary 9.** Let  $\mathbf{x}$  be a sub-Gaussian random vector with zero mean and covariance  $\boldsymbol{\Sigma}^*$ . Consider a collection  
 221 of  $T$  i.i.d. samples  $\{\mathbf{x}_i\}_{i=1}^T$  drawn from  $\mathbf{x}$ . Suppose that Assumptions 1, 3, 4 and 8 hold. There exist universal  
 222 constants  $C_1, C_2, \dots, C_{11}$  such that if  $\lambda \geq 2C_1\delta\varpi\sqrt{d/m}$ , where  $\varpi = \sqrt{\sup_{\|\mathbf{u}\|_2=1, \|\mathbf{v}\|_2=1} \mathbf{u}^\top \mathbf{a}_i \mathbf{a}_i^\top \mathbf{v}}$  and  
 223  $\rho^- = C_2\lambda_{\min}(\boldsymbol{\Theta}) > \zeta^-$ , it holds with probability at least  $1 - C_3 \exp(-2dC_4)$  that

$$\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{\text{F}} \leq \frac{\delta\varpi}{\sqrt{mT}\lambda_{\min}(\boldsymbol{\Theta})} [C_5 r_1 + C_6 \sqrt{r_2 d}].$$

224 Further, if  $\boldsymbol{\sigma}^*$  satisfies that

$$\min_{i \in \mathcal{S}} |\sigma_i^*| \geq \nu + \frac{C_7\delta\varpi\sqrt{rd/m}}{\lambda_{\min}(\boldsymbol{\Theta})},$$

225 where  $\mathcal{S} = \text{supp}(\boldsymbol{\sigma}^*)$ . And if  $\lambda \geq C_8\delta\varpi\sqrt{d/m} \left(1 + \frac{\sqrt{r}\lambda_{\max}(\boldsymbol{\Theta})}{\lambda_{\min}(\boldsymbol{\Theta})}\right)$ , it holds with probability at least  $1 -$   
 226  $C_9 \exp(-C_{10}d)$  that

$$\|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{\text{F}} \leq \frac{C_{11}r\delta\varpi}{\sqrt{mT}\lambda_{\min}(\boldsymbol{\Theta})}.$$

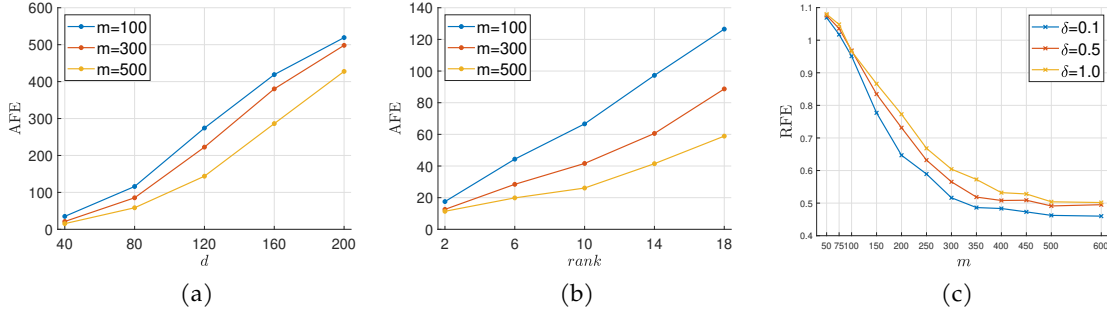


Figure 1: Absolute Frobenius norm error of the proposed estimator changing  $d$  and  $rank$  and Relative Frobenius norm error changing the noise variance proxy  $\delta$ .

## 5. Numerical Experiments

In this section, we perform several numerical experiments to evaluate the performance of the CCS with proposed regularization techniques. All experiments are implemented in MATLAB 2023b, run on an Intel i7 – 10700 2.90 GHz  $\times 8$  with 16 GB of RAM, and all the results reported represent averages over 100 Monte Carlo simulations.

**Data Generation and Hyperparameter Setup:** The synthetic datasets are generated as follows. We generate the true covariance matrix as  $\Sigma^* = \mathbf{L}\mathbf{L}^\top$ , where  $\mathbf{L}$  is a  $d \times r$  matrix with independent Gaussian components. The sensing vector  $\{\mathbf{a}_i\}_{i=1}^m$  are obtained by drawing i.i.d. samples from a sub-Gaussian  $(\mathbf{0}, \mathbf{Y})$  distribution. The noise  $\{\eta_i\}_{i=1}^m$  obeys the sub-Gaussian distribution with mean  $\mathbf{0}$  and a variance proxy of  $\delta^2$ . For all methods, the tuning parameter  $\lambda$  is determined through the application of five-fold cross-validation. Unless otherwise specified, we employ the SCAD as the non-convex penalty function, setting the parameter  $b$  to 3.7, and the variance proxy  $\delta$  is 0.1.

**Measurement Criteria:** To evaluate the performance of the proposed estimator, we utilize several metrics based on the Frobenius norm. Let  $\hat{\Sigma}$  denote the generic covariance matrix obtained by various algorithms and estimators,  $\Sigma^*$  represent the true covariance matrix. We define the Absolute Frobenius Error (AFE) as  $\|\hat{\Sigma} - \Sigma^*\|_F$ . Additionally, we introduce the Mean squared Frobenius Error (MFE) as  $\frac{\|\hat{\Sigma} - \Sigma^*\|_F^2}{d^2}$ , where  $d$  is the dimension of the covariance matrix. Similarly, the Relative Frobenius Error (RFE) is defined as  $\frac{\|\hat{\Sigma} - \Sigma^*\|_F}{\|\Sigma^*\|_F}$ .

Figure 1 shows the performance of our estimator under different parameter settings. Figure 1a shows the variation of the AFE as the dimension  $d$  of the covariance matrix increases, with a fixed rank of 10. Three curves, in different colors, correspond to different sample numbers  $m = \{100, 300, 500\}$ . From the plot, we observe that the estimation error increases with the dimension. A larger sample number is required for convergence as the dimension grows. Figure 1b presents the results when the matrix dimension  $d = 60$ , with three lines corresponding to different sample numbers. The results display that the augmentation of the rank of the true covariance matrix leads to a rise in the estimation error. For figure 1c, the dimension is fixed at  $d = 60$ . The experiment is implemented by changing the variance proxy  $\delta$  of the noise component, which shows higher noise intensity adds up to estimation error.

Figures 2a and 2b illustrate experiments with synthetic data by varying the number of the samples while maintaining a fixed matrix rank of 10. Figure 2a shows the mean squared Frobenius norm error of our proposed non-convex estimator. The estimation error converges to the same value as the dimension increases. Figure 2b compares the performance of the non-convex method to the convex method (Nuclear Norm). The non-convex method demonstrates a faster rate of error reduction as the sample number increases.



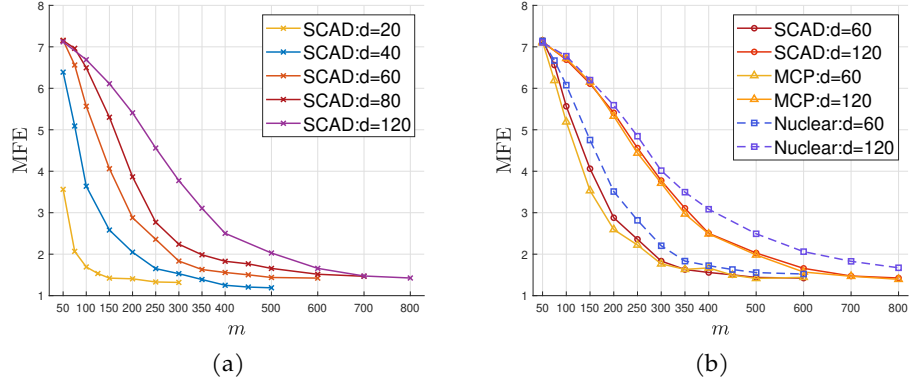


Figure 2: Mean Squared Frobenius norm error of the proposed estimator and the comparison between convex and nonconvex methods.

Table 1: The estimation error of non-convex and convex methods  $\pm$  standard variance under 100 Monte Carlo Simulations.

		$d = 40$				$d = 80$			
		$m = 100$	$m = 200$	$m = 300$	$m = 500$	$m = 100$	$m = 300$	$m = 500$	$m = 600$
SCAD	AFE	76.2997 $\pm$ 13.0608	57.2471 $\pm$ 14.4645	49.4734 $\pm$ 15.8468	43.6211 $\pm$ 18.3060	203.8819 $\pm$ 14.3887	119.7443 $\pm$ 31.1828	102.9800 $\pm$ 36.4549	98.5062 $\pm$ 46.9165
	RFE	0.7281 $\pm$ 0.1261	0.5541 $\pm$ 0.1024	0.4796 $\pm$ 0.1516	0.4087 $\pm$ 0.1627	0.9846 $\pm$ 0.0275	0.5686 $\pm$ 0.1460	0.4883 $\pm$ 0.1725	0.46850 $\pm$ 0.2120
	MFE	3.6385 $\pm$ 1.1104	2.0483 $\pm$ 1.1008	1.5298 $\pm$ 1.2044	1.1893 $\pm$ 1.4276	6.4950 $\pm$ 0.7281	2.2404 $\pm$ 1.2968	1.6570 $\pm$ 1.4817	1.5162 $\pm$ 1.2715
MCP	AFE	72.3491 $\pm$ 16.2992	52.5658 $\pm$ 18.5107	43.5744 $\pm$ 16.5545	40.0282 $\pm$ 18.4621	200.6326 $\pm$ 13.0460	120.2025 $\pm$ 33.7358	101.3771 $\pm$ 37.7358	96.7045 $\pm$ 40.9204
	RFE	0.6823 $\pm$ 0.1459	0.5050 $\pm$ 0.1631	0.4164 $\pm$ 0.1694	0.4730 $\pm$ 0.1763	0.9483 $\pm$ 0.0377	0.5705 $\pm$ 0.1579	0.4851 $\pm$ 0.1927	0.4634 $\pm$ 0.2031
	MFE	3.2715 $\pm$ 1.2436	1.7271 $\pm$ 1.3924	1.1867 $\pm$ 1.4667	1.0001 $\pm$ 1.7633	6.2896 $\pm$ 0.8240	2.2576 $\pm$ 1.3099	1.5992 $\pm$ 1.8947	1.4612 $\pm$ 1.4458
Nuclear	AFE	79.6284 $\pm$ 12.5129	58.1602 $\pm$ 14.0585	50.1646 $\pm$ 19.4004	47.1060 $\pm$ 19.2026	209.4597 $\pm$ 12.3611	122.0572 $\pm$ 36.6691	103.4399 $\pm$ 47.9068	106.7587 $\pm$ 65.1436
	RFE	0.7729 $\pm$ 0.1113	0.5656 $\pm$ 0.1397	0.4824 $\pm$ 0.1879	0.4458 $\pm$ 0.2601	0.9966 $\pm$ 0.0542	0.6408 $\pm$ 0.1837	0.4465 $\pm$ 0.2279	0.5132 $\pm$ 0.3205
	MFE	3.9629 $\pm$ 1.2688	2.1141 $\pm$ 1.0384	1.5728 $\pm$ 1.5795	1.3869 $\pm$ 2.1682	7.2163 $\pm$ 0.7761	2.4620 $\pm$ 1.3350	1.8642 $\pm$ 2.0745	1.78080 $\pm$ 3.3428
		$d = 120$				$d = 160$			
		$m = 100$	$m = 300$	$m = 500$	$m = 700$	$m = 100$	$m = 300$	$m = 500$	$m = 800$
SCAD	AFE	290.3662 $\pm$ 13.0428	233.1275 $\pm$ 34.3399	160.9030 $\pm$ 39.5538	145.6902 $\pm$ 44.8508	416.1755 $\pm$ 13.1164	371.2838 $\pm$ 20.1159	270.8453 $\pm$ 40.8805	183.1849 $\pm$ 68.6043
	RFE	0.9958 $\pm$ 0.0215	0.7469 $\pm$ 0.1097	0.54620 $\pm$ 0.1885	0.4508 $\pm$ 0.1429	1.0001 $\pm$ 0.0122	0.8987 $\pm$ 0.0435	0.6989 $\pm$ 0.0990	0.4396 $\pm$ 0.1666
	MFE	6.6894 $\pm$ 0.5615	3.7742 $\pm$ 1.0867	2.0283 $\pm$ 1.7619	1.4740 $\pm$ 0.9685	6.7657 $\pm$ 0.4275	5.4722 $\pm$ 0.5948	3.0043 $\pm$ 0.8789	1.3108 $\pm$ 1.2793
MCP	AFE	291.0940 $\pm$ 11.3085	231.1940 $\pm$ 34.4239	168.9255 $\pm$ 45.8193	145.3658 $\pm$ 52.7231	417.0620 $\pm$ 13.0938	379.877326 $\pm$ 4.551	273.8305 $\pm$ 41.9579	184.9737 $\pm$ 72.2468
	RFE	0.9915 $\pm$ 0.0234	0.7325 $\pm$ 0.1089	0.5424 $\pm$ 0.1502	0.4648 $\pm$ 0.2616	0.9966 $\pm$ 0.0121	0.9033 $\pm$ 0.0543	0.7293 $\pm$ 0.1200	0.4698 $\pm$ 0.1680
	MFE	6.7208 $\pm$ 0.4928	3.7119 $\pm$ 1.0580	1.9817 $\pm$ 1.0836	1.4674 $\pm$ 1.4120	6.7946 $\pm$ 0.4241	5.6370 $\pm$ 0.7717	3.0960 $\pm$ 0.9044	1.4850 $\pm$ 1.3038
Nuclear	AFE	292.3386 $\pm$ 8.9425	240.4273 $\pm$ 35.7671	189.4290 $\pm$ 50.5312	162.2422 $\pm$ 64.2535	416.5137 $\pm$ 15.3655	374.2606 $\pm$ 24.9383	282.5598 $\pm$ 58.3104	192.6122 $\pm$ 94.5904
	RFE	1.0010 $\pm$ 0.0174	0.7728 $\pm$ 0.1072	0.6043 $\pm$ 0.2912	0.5190 $\pm$ 0.2086	1.0014 $\pm$ 0.0130	0.9004 $\pm$ 0.0548	0.6810 $\pm$ 0.1414	0.46160 $\pm$ 0.1530
	MFE	6.7747 $\pm$ 0.3910	4.0143 $\pm$ 1.2636	2.4919 $\pm$ 3.6096	1.8280 $\pm$ 1.6649	6.7767 $\pm$ 0.4999	5.4715 $\pm$ 0.7326	3.1188 $\pm$ 1.3402	1.4492 $\pm$ 1.1249

Table 1 presents the performance of our proposed methods with different penalties with different dimensions of the covariance matrix. The statistical results of SCAD, MCP, and Nuclear Norm methods are shown in the table with different estimation standards and the standard variance. For each experiment, the rank of the covariance matrix is set to 10, and the variance proxy  $\delta = 0.1$ . The dimension of the matrix to be estimated changes in the range of  $d = 40, 80, 120, 160$ . And different sample numbers for convergence of each dimension can be observed from the results in the table.

## 6. Conclusion and Discussion

In this paper, we investigate compressive covariance sensing and introduce a quadratic measurement model tailored for estimating high-dimensional low-rank covariance matrices. This model effectively addresses the inherent challenges of limited measurements, as well as the constraints on storage and computational resources typically in high-dimensional settings. We propose a novel non-convex penalization approach and employ a proximal gradient homotopy algorithm to solve the resulting estimator. Notably, our proposed estimator achieves superior statistical convergence rates compared to traditional nuclear norm penalized estimators and attains the oracle property. These advancements are corroborated through comprehensive theoretical analysis and empirical experiments. As a limitation, we note that our study is the absence of performance assessments on real-world datasets, such as those encountered in spectrum analysis. Future work will focus on applying our estimator to publicly available real-world datasets to further validate its effectiveness in practical scenarios.

## References

- [1] Wei-Liem Loh. *Estimating covariance matrices*. Stanford University, 1988.
- [2] Mohsen Pourahmadi. *High-dimensional covariance estimation: with high-dimensional data*, volume 882. John Wiley & Sons, 2013.
- [3] Arup Bose and Monika Bhattacharjee. *Large covariance and autocovariance matrices*. Chapman and Hall/CRC, 2018.
- [4] Aygul Zagidullina. *High-dimensional covariance matrix estimation: an introduction to random matrix theory*. Springer Nature, 2021.
- [5] Ted H Szatrowski. Necessary and sufficient conditions for explicit solutions in the multivariate normal estimation problem for patterned means and covariances. *The Annals of Statistics*, pages 802–810, 1980.
- [6] Hsien-sen Hung and Mostafa Kaveh. On the statistical sufficiency of the coherently averaged covariance matrix for the estimation of the parameters of wideband sources. In *ICASSP'87. IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 12, pages 33–36. IEEE, 1987.
- [7] Jianfeng Yao, Shurong Zheng, and ZD Bai. *Sample covariance matrices and high-dimensional data analysis*. Cambridge UP, New York, 2015.
- [8] Emile Richard, Pierre-André Savalle, and Nicolas Vayatis. Estimation of simultaneously sparse and low rank matrices. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 51–58, 2012.
- [9] Azer P Shikhaliev, Lee C Potter, and Yuejie Chi. Low-rank structured covariance matrix estimation. *IEEE Signal Processing Letters*, 26(5):700–704, 2019.
- [10] Quan Wei and Ziping Zhao. Large covariance matrix estimation with oracle statistical rate via majorization-minimization. *IEEE Transactions on Signal Processing*, 2023.
- [11] Dyonisius Dony Ariananda and Geert Leus. Wideband power spectrum sensing using sub-Nyquist sampling. In *2011 IEEE 12th International Workshop on Signal Processing Advances in Wireless Communications*, pages 101–105. IEEE, 2011.
- [12] Dyonisius Dony Ariananda and Geert Leus. Compressive Wideband Power Spectrum Estimation. *IEEE Transactions on Signal Processing*, 60(9):4775–4789, 2012.
- [13] Dyonisius Dony Ariananda, Daniel Romero, and Geert Leus. Compressive angular and frequency periodogram reconstruction for multiband signals. In *2013 5th IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, pages 440–443, 2013.
- [14] Daniel Romero, Dyonisius Dony Ariananda, Zhi Tian, and Geert Leus. Compressive covariance sensing: Structure-based compressive sensing beyond sparsity. *IEEE Signal Processing Magazine*, 33(1):78–93, 2015.
- [15] Yuxin Chen, Yuejie Chi, and Andrea J. Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *IEEE Transactions on Information Theory*, 61(7):4034–4059, 2015.
- [16] Kai Zhong, Prateek Jain, and Inderjit S Dhillon. Efficient matrix sensing using rank-1 gaussian measurements. In *Algorithmic Learning Theory: 26th International Conference, ALT 2015*, pages 3–18. Springer, 2015.
- [17] Xixian Chen, Michael R Lyu, and Irwin King. Toward efficient and accurate covariance matrix estimation on compressed data. In *International Conference on Machine Learning*, pages 767–776. PMLR, 2017.

- [18] Angelika Rohde and Alexandre B. Tsybakov. Estimation of high-dimensional low-rank matrices. *The Annals of Statistics*, 39(2):887–930, 2011.
- [19] Karim Lounici. High-dimensional covariance matrix estimation with missing observations. *Bernoulli*, 20(3):1029–1058, 2014.
- [20] Kishore Jaganathan, Samet Oymak, and Babak Hassibi. Recovery of sparse 1-d signals from the magnitudes of their fourier transform. In *2012 IEEE International Symposium on Information Theory Proceedings*, pages 1473–1477. IEEE, 2012.
- [21] Emmanuel J Candes, Thomas Strohmer, and Vladislav Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.
- [22] Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. Phase retrieval using alternating minimization. *Advances in Neural Information Processing Systems*, 26, 2013.
- [23] Amir Beck and Yonina C Eldar. Sparsity constrained nonlinear optimization: Optimality conditions and algorithms. *SIAM Journal on Optimization*, 23(3):1480–1509, 2013.
- [24] Emmanuel J Candès and Xiaodong Li. Solving quadratic equations via phaselift when there are about as many equations as unknowns. *Foundations of Computational Mathematics*, 5(14):1017–1026, 2014.
- [25] Irene Waldspurger, Alexandre d’Aspremont, and Stéphane Mallat. Phase recovery, maxcut and complex semidefinite programming. *Mathematical Programming*, 149:47–81, 2015.
- [26] Emmanuel J Candès and Terence Tao. The power of convex relaxation: Near-optimal matrix completion. *IEEE Transactions on Information Theory*, 56(5):2053–2080, 2010.
- [27] Feiping Nie, Hua Wang, Xiao Cai, Heng Huang, and Chris Ding. Robust matrix completion via joint Schatten p-norm and lp-norm minimization. In *2012 IEEE 12th International Conference on Data Mining*, pages 566–574. IEEE, 2012.
- [28] Kun Chen, Hongbo Dong, and Kung-Sik Chan. Reduced rank regression via adaptive nuclear norm penalization. *Biometrika*, 100(4):901–920, 2013.
- [29] Sahand Negahban and Martin J Wainwright. Estimation of (near) low-rank matrices with noise and high-dimensional scaling. *The Annals of Statistics*, 39(2):1069–1097, 2011.
- [30] Vladimir Koltchinskii, Karim Lounici, and Alexandre B Tsybakov. Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *The Annals of Statistics*, 39(5):2302, 2011.
- [31] Yan Huang, Guisheng Liao, Yijian Xiang, Lei Zhang, Jie Li, and Arye Nehorai. Low-rank approximation via generalized reweighted iterative nuclear and Frobenius norms. *IEEE Transactions on Image Processing*, 29:2244–2257, 2020.
- [32] Jianqing Fan and Runze Li. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96(456):1348–1360, 2001.
- [33] Cun-Hui Zhang. Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics*, 38(2):894–942, 2010.
- [34] Dehua Liu, Tengfei Zhou, Hui Qian, Congfu Xu, and Zhihua Zhang. A nearly unbiased matrix completion approach. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2013*, pages 210–225. Springer, 2013.
- [35] Zhaoran Wang, Han Liu, and Tong Zhang. Optimal computational and statistical rates of convergence for sparse nonconvex learning problems. *The Annals of Statistics*, 42(6):2164, 2014.

- [36] Huan Gui, Jiawei Han, and Quanquan Gu. Towards faster rates and oracle property for low-rank matrix estimation. In *International Conference on Machine Learning*, pages 2300–2309. PMLR, 2016.
- [37] Hu Shao, William HK Lam, Agachai Sumalee, Anthony Chen, and Martin L Hazelton. Estimation of mean and covariance of peak hour origin–destination demands from day-to-day traffic counts. *Transportation Research Part B: Methodological*, 68:52–75, 2014.
- [38] Björn Ottersten, Peter Stoica, and Richard Roy. Covariance matching estimation techniques for array signal processing applications. *Digital Signal Processing*, 8(3):185–210, 1998.
- [39] Kai Yu, Shenghuo Zhu, John Lafferty, and Yihong Gong. Fast nonparametric matrix factorization for large-scale collaborative filtering. In *Proceedings of the 32nd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 211–218, 2009.
- [40] Yu Zhang, Bin Cao, and Dit Yan Yeung. Multi-domain collaborative filtering. In *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence*, page 725, 2010.
- [41] Ravishankar Sivalingam, Vassilios Morellas, Daniel Boley, and Nikolaos Papanikolopoulos. Metric learning for semi-supervised clustering of region covariance descriptors. In *2009 Third ACM/IEEE International Conference on Distributed Smart Cameras*, pages 1–8. IEEE, 2009.
- [42] Aurélien Bellet, Amaury Habrard, and Marc Sebban. *Metric learning*. Morgan & Claypool Publishers, 2015.
- [43] William A Gardner. Exploitation of spectral redundancy in cyclostationary signals. *IEEE Signal Processing Magazine*, 8(2):14–36, 1991.
- [44] Kyouwoong Kim, Ihsan A Akbar, Kyung K Bae, Jung-Sun Um, Chad M Spooner, and Jeffrey H Reed. Cyclostationary approaches to signal detection and classification in cognitive radio. In *2007 2nd IEEE International Symposium on New Frontiers in Dynamic Spectrum Access Networks*, pages 212–215. IEEE, 2007.
- [45] Georgios B Giannakis. Cyclostationary signal analysis. In *Digital Signal Processing Fundamentals*, pages 433–464. CRC press, 2017.
- [46] Benjamin Rolfs, Bala Rajaratnam, Dominique Guillot, Ian Wong, and Arian Maleki. Iterative thresholding algorithm for sparse inverse covariance estimation. *Advances in Neural Information Processing Systems*, 25, 2012.
- [47] Wenfu Xia, Ziping Zhao, and Ying Sun. C-ISTA: Iterative shrinkage-thresholding algorithm for sparse covariance matrix estimation. In *IEEE Statistical Signal Processing Workshop*, pages 215–219, 2023.
- [48] Sahand N Negahban, Pradeep Ravikumar, Martin J Wainwright, and Bin Yu. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. 2012.
- [49] Qiang Sun, Kean Ming Tan, Han Liu, and Tong Zhang. Graphical nonconvex optimization via an adaptive convex relaxation. In *International Conference on Machine Learning*, pages 4810–4817. PMLR, 2018.

## A. Appendices

### A.1. Technical Lemmata

**Lemma 10.** Under Assumption 3, and assuming that  $\Sigma_1 - \Sigma_2 \in \mathcal{B}$ , we establish

$$\tilde{f}(\Sigma_2) \geq \tilde{f}(\Sigma_1) + \langle \nabla \tilde{f}(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle + \frac{\rho^- - \zeta^-}{2} \|\Sigma_2 - \Sigma_1\|_F^2.$$

*Proof.* Considering the eigen-decompositions of covariance matrices  $\Sigma_1$  and  $\Sigma_2$  given by

$$\Sigma_1 = \mathbf{U}_1 \mathbf{\Lambda}_1 \mathbf{U}_1^\top, \Sigma_2 = \mathbf{U}_2 \mathbf{\Lambda}_2 \mathbf{U}_2^\top,$$

where  $\mathbf{\Lambda}_1 = \text{diag}(\sigma_1)$  and  $\mathbf{\Lambda}_2 = \text{diag}(\sigma_2)$  are diagonal matrices containing eigenvalues  $\sigma_1$  and  $\sigma_2$  of  $\Sigma_1$  and  $\Sigma_2$ , respectively.

Under (b) of Assumption 1, for each component  $i = 1, 2, \dots, d$  in  $\sigma_1$  and  $\sigma_2$ , we have

$$q'_\lambda((\sigma_1)_i) - q'_\lambda((\sigma_2)_i) \geq -\zeta^-((\sigma_1)_i - (\sigma_2)_i),$$

which leads to

$$\langle (-\nabla Q_\lambda(\Sigma_1)) - (-\nabla Q_\lambda(\Sigma_2)), \Sigma_1 - \Sigma_2 \rangle \leq \zeta^- \|\Sigma_1 - \Sigma_2\|_F^2.$$

This inequality illustrates the strong smoothness of  $-Q(\cdot)$  in the region  $\mathcal{B}$ . As a consequence, we establish

$$Q_\lambda(\Sigma_1) + \langle \nabla Q_\lambda(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle - \frac{\zeta^-}{2} \|\Sigma_1 - \Sigma_2\|_F^2 \leq Q_\lambda(\Sigma_2). \quad (4)$$

Additionally, by Assumption 3, we obtain

$$f(\Sigma_2) \geq f(\Sigma_1) + \langle \nabla f(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle + \frac{\rho^-}{2} \|\Sigma_2 - \Sigma_1\|_F^2. \quad (5)$$

Given that  $\tilde{f}(\Sigma)$  is the sum of  $f(\Sigma)$  and  $Q_\lambda(\Sigma)$ , we can combine 4 and 5 to achieve

$$\tilde{f}(\Sigma_2) \geq \tilde{f}(\Sigma_1) + \langle \nabla \tilde{f}(\Sigma_1), \Sigma_2 - \Sigma_1 \rangle + \frac{\rho^- - \zeta^-}{2} \|\Sigma_2 - \Sigma_1\|_F^2.$$

This result demonstrates the restricted strong convexity of the objective function  $\tilde{f}(\Sigma)$ .  $\square$

**Lemma 11.** Suppose Assumption 3 holds. If  $\rho^- > \zeta^-$ , and  $\lambda \geq \frac{2}{m} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*))\|_2$ , we have

$$\left\| \mathcal{P}^\perp(\hat{\Sigma} - \Sigma^*) \right\|_* \leq 5 \left\| \mathcal{P}(\hat{\Sigma} - \Sigma^*) \right\|_*.$$

*Proof.* Given Lemma 10, we obtain

$$\underbrace{\left( \tilde{f}(\hat{\Sigma}) + \lambda \|\hat{\Sigma}\|_* \right) - \left( \tilde{f}(\Sigma^*) + \lambda \|\Sigma^*\|_* \right)}_{\text{I}} \geq \underbrace{\langle \nabla \tilde{f}(\Sigma^*), \hat{\Sigma} - \Sigma^* \rangle}_{\text{II}} + \underbrace{\lambda \|\hat{\Sigma}\|_* - \lambda \|\Sigma^*\|_*}_{\text{III}}. \quad (6)$$

In the following, we analyze the term II. Applying Hölder's inequality, we have

$$\begin{aligned} & \langle \nabla \tilde{f}(\Sigma^*), \hat{\Sigma} - \Sigma^* \rangle \\ &= \langle \nabla \tilde{f}(\Sigma^*), \mathcal{P}(\hat{\Sigma} - \Sigma^*) \rangle + \langle \nabla \tilde{f}(\Sigma^*), \mathcal{P}^\perp(\hat{\Sigma} - \Sigma^*) \rangle \\ &\geq - \underbrace{\left\| \mathcal{P}(\nabla \tilde{f}(\Sigma^*)) \right\|_2}_{\text{II}_1} \cdot \left\| \mathcal{P}(\hat{\Sigma} - \Sigma^*) \right\|_* - \underbrace{\left\| \mathcal{P}^\perp(\nabla \tilde{f}(\Sigma^*)) \right\|_2}_{\text{II}_2} \cdot \left\| \mathcal{P}^\perp(\hat{\Sigma} - \Sigma^*) \right\|_*. \end{aligned} \quad (7)$$

It remains to bound the terms I, II<sub>1</sub>, II<sub>2</sub> and III, respectively.

421 For the term I, since  $\widehat{\Sigma}$  represents the optimal solution to the problem (1). This leads to

$$\left( \tilde{f}(\widehat{\Sigma}) + \lambda \|\widehat{\Sigma}\|_* \right) - \left( \tilde{f}(\Sigma^*) + \lambda \|\Sigma^*\|_* \right) \leq 0. \quad (8)$$

422 For the term  $\text{II}_1$ , considering the choice of  $\lambda$ , we have

$$\|\nabla f(\Sigma^*)\|_2 = \left\| \frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*)) \right\|_2 \leq \frac{\lambda}{2}.$$

423 Based on (a) in Assumption 1 and  $\lambda \geq \frac{2}{m} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*))\|_2$ , we have  $\|\nabla Q_\lambda(\Sigma^*)\|_2 \leq \lambda$ , which  
424 leads to

$$\left\| \mathcal{P}(\nabla \tilde{f}(\Sigma^*)) \right\|_2 = \left\| \mathcal{P}(\nabla f(\Sigma^*) + \nabla Q_\lambda(\Sigma^*)) \right\|_2 \leq \frac{3\lambda}{2}. \quad (9)$$

425 For term  $\text{II}_2$ , we first note from the definition of  $\mathcal{P}^\perp(\cdot)$  that  $\mathcal{P}^\perp(\Sigma^*) = \mathbf{0}$ . Based on (c) in Assumption  
426 (1), it follows that  $Q_\lambda(\mathcal{P}^\perp(\Sigma^*)) = 0$ . This leads to

$$\left\| \mathcal{P}^\perp(\nabla \tilde{f}(\Sigma^*)) \right\|_2 = \left\| \mathcal{P}^\perp(\nabla f(\Sigma^*)) \right\|_2 \leq \frac{\lambda}{2}. \quad (10)$$

427 Then we substitute (9) and (10) into (7), we have

$$\langle \nabla \tilde{f}(\Sigma^*), \widehat{\Sigma} - \Sigma^* \rangle \geq -\frac{3\lambda}{2} \left\| \mathcal{P}(\widehat{\Sigma} - \Sigma^*) \right\|_* - \frac{\lambda}{2} \left\| \mathcal{P}^\perp(\widehat{\Sigma} - \Sigma^*) \right\|_*. \quad (11)$$

428 For the term III, projecting onto the complementary subspaces  $\mathcal{T}$  and  $\mathcal{T}^\perp$  yields the following de-  
429 composition

$$\begin{aligned} \lambda \|\widehat{\Sigma}\|_* - \lambda \|\Sigma^*\|_* &= \lambda \|\mathcal{P}(\widehat{\Sigma})\|_* - \lambda \|\mathcal{P}(\Sigma^*)\|_* + \lambda \|\mathcal{P}^\perp(\widehat{\Sigma})\|_* - \lambda \|\mathcal{P}^\perp(\Sigma^*)\|_* \\ &\geq -\lambda \left\| \mathcal{P}(\widehat{\Sigma} - \Sigma^*) \right\|_* + \lambda \left\| \mathcal{P}^\perp(\widehat{\Sigma} - \Sigma^*) \right\|_*. \end{aligned} \quad (12)$$

430 Based on the inequalities in formulas (8), (11), and (12), when substituted into formula (6), we  
431 obtain

$$0 \geq -\frac{3\lambda}{2} \left\| \mathcal{P}(\widehat{\Sigma} - \Sigma^*) \right\|_* - \frac{\lambda}{2} \left\| \mathcal{P}^\perp(\widehat{\Sigma} - \Sigma^*) \right\|_* - \lambda \left\| \mathcal{P}(\widehat{\Sigma} - \Sigma^*) \right\|_* + \lambda \left\| \mathcal{P}^\perp(\widehat{\Sigma} - \Sigma^*) \right\|_*.$$

432 From this inequality, we have

$$\left\| \mathcal{P}^\perp(\widehat{\Sigma} - \Sigma^*) \right\|_* \leq 5 \left\| \mathcal{P}(\widehat{\Sigma} - \Sigma^*) \right\|_*,$$

433 which implies that  $\widehat{\Delta} = \widehat{\Sigma} - \Sigma^*$  lies in the cone  $\mathcal{B}$ . □

## 434 A.2. Proof of Theorem 5

435 *Proof.* Applying Lemma 10, we establish

$$\begin{aligned} \tilde{f}(\widehat{\Sigma}) &\geq \tilde{f}(\Sigma^*) + \langle \nabla \tilde{f}(\Sigma^*), \widehat{\Sigma} - \Sigma^* \rangle + \frac{\rho^- - \zeta^-}{2} \left\| \widehat{\Sigma} - \Sigma^* \right\|_{\text{F}}^2, \\ \tilde{f}(\Sigma^*) &\geq \tilde{f}(\widehat{\Sigma}) + \langle \nabla \tilde{f}(\widehat{\Sigma}), \Sigma^* - \widehat{\Sigma} \rangle + \frac{\rho^- - \zeta^-}{2} \left\| \Sigma^* - \widehat{\Sigma} \right\|_{\text{F}}^2. \end{aligned}$$

436 By summing these two inequalities, we obtain

$$(\rho^- - \zeta^-) \left\| \widehat{\Sigma} - \Sigma^* \right\|_{\text{F}}^2 - \langle \nabla \tilde{f}(\widehat{\Sigma}) - \nabla \tilde{f}(\Sigma^*), \widehat{\Sigma} - \Sigma^* \rangle \leq 0. \quad (13)$$

437 Considering the convexity of the nuclear norm  $\|\cdot\|_*$ , we deduce

$$\langle \lambda \widehat{\mathbf{W}} - \lambda \mathbf{W}^*, \widehat{\Sigma} - \Sigma^* \rangle \geq 0, \quad (14)$$

438 where  $\widehat{\mathbf{W}} \in \partial \|\widehat{\boldsymbol{\Sigma}}\|_*$  and  $\mathbf{W}^* \in \partial \|\boldsymbol{\Sigma}^*\|_*$ .

439 Since  $\widehat{\boldsymbol{\Sigma}}$  is the optimal solution to (1), it satisfies the optimality condition:

$$\langle \nabla \tilde{f}(\widehat{\boldsymbol{\Sigma}}) + \lambda \widehat{\mathbf{W}} + \widehat{\boldsymbol{\Xi}}, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \leq 0, \quad (15)$$

440 where  $\widehat{\boldsymbol{\Xi}} \in \partial \mathbb{I}_{\{\boldsymbol{\Sigma} \succeq 0\}}(\widehat{\boldsymbol{\Sigma}})$ .

441 Given the definition of the indicator function  $\mathbb{I}_{\{\boldsymbol{\Sigma} \succeq 0\}}(\boldsymbol{\Sigma})$ , it follows that:

$$\langle \widehat{\boldsymbol{\Xi}}, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle = 0. \quad (16)$$

442 By combining the inequalities (13), (14), (15) and (16), we obtain

$$\begin{aligned} (\rho^- - \zeta^-) \|\widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_{\text{F}}^2 &\leq \langle \nabla \tilde{f}(\widehat{\boldsymbol{\Sigma}}) - \nabla \tilde{f}(\boldsymbol{\Sigma}^*), \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \nabla \tilde{f}(\widehat{\boldsymbol{\Sigma}}) + \lambda \widehat{\mathbf{W}} + \widehat{\boldsymbol{\Xi}}, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle - \langle \widehat{\boldsymbol{\Xi}}, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &\quad - \langle \lambda \widehat{\mathbf{W}} - \lambda \mathbf{W}^*, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle - \langle \nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &\leq - \langle \nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*, \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle. \end{aligned} \quad (17)$$

443 Consider the partitioning of the eigenvalues of  $\boldsymbol{\Sigma}^*$  into three distinct sets:

- 444 • For indices  $i \in \mathcal{S}_1$ , the eigenvalues satisfy  $\sigma_i^* \geq \alpha$ .
- 445 • For indices  $i \in \mathcal{S}_2$ , we have  $\alpha > \sigma_i^* > 0$ .
- 446 • For indices  $i \in \overline{\mathcal{S}}$ , we have  $\sigma_i^* = 0$ , with  $\mathcal{S} = \mathcal{S}_1 \cup \mathcal{S}_2$ .

447 Based on this partition, we define two subspace of  $\mathcal{T}$  associated with  $\mathcal{S}_1$  and  $\mathcal{S}_2$  as follows:

$$\mathcal{T}_{\mathcal{S}_1}(\mathbf{U}^*) := \{ \boldsymbol{\Delta} \in \mathbb{R}^{d \times d} \mid \text{row}(\boldsymbol{\Delta}) \subset \mathbf{U}_{\mathcal{S}_1}^* \text{ and } \text{col}(\boldsymbol{\Delta}) \subset \mathbf{U}_{\mathcal{S}_1}^* \},$$

448

$$\mathcal{T}_{\mathcal{S}_2}(\mathbf{U}^*) := \{ \boldsymbol{\Delta} \in \mathbb{R}^{d \times d} \mid \text{row}(\boldsymbol{\Delta}) \subset \mathbf{U}_{\mathcal{S}_2}^* \text{ and } \text{col}(\boldsymbol{\Delta}) \subset \mathbf{U}_{\mathcal{S}_2}^* \},$$

449 where  $\mathbf{U}_{\mathcal{S}_1}^*$  is the sub-matrix of  $\mathbf{U}^*$  formed by the rows indexed by  $i \in \mathcal{S}_1$ , and similarly,  $\mathbf{U}_{\mathcal{S}_2}^*$  is the  
 450 sub-matrix formed by the rows indexed by  $i \in \mathcal{S}_2$ . Corresponding to these subspaces, we define  
 451 two projection operators  $\mathcal{P}_{\mathcal{S}_1}(\cdot)$  and  $\mathcal{P}_{\mathcal{S}_2}(\cdot)$ , which project matrices onto the subspace  $\mathcal{T}_{\mathcal{S}_1}$  and  $\mathcal{T}_{\mathcal{S}_2}$ ,  
 452 respectively. In the following, we decompose the RHS in (17) with respect to these three distinct  
 453 components.

454 For  $i \in \overline{\mathcal{S}}$ , we proceed under (a) in Assumption 1, which ensures that  $q'_\lambda(0) = 0$ . This property  
 455 enables us to express the gradient as follows:

$$\nabla Q_\lambda(\boldsymbol{\Sigma}^*) = \mathbf{U}^* q'_\lambda(\boldsymbol{\Gamma}^*) \mathbf{U}^{*\top},$$

456 where  $\boldsymbol{\Gamma}^* \in \mathbb{R}^{r \times r}$  is the diagonal matrix with  $\text{diag}(\boldsymbol{\Gamma}^*) = \boldsymbol{\sigma}^*$ . By projecting  $\nabla Q_\lambda(\boldsymbol{\Sigma}^*)$  onto  $\mathcal{T}^\perp$ , we  
 457 obtain

$$\begin{aligned} \mathcal{P}^\perp(\nabla Q_\lambda(\boldsymbol{\Sigma}^*)) &= (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{U}^* q'_\lambda(\boldsymbol{\Gamma}^*) \mathbf{U}^{*\top} (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \\ &= (\mathbf{U}^* - \mathbf{U}^*) q'_\lambda(\boldsymbol{\Gamma}^*) (\mathbf{U}^{*\top} - \mathbf{U}^{*\top}) \\ &= \mathbf{0}. \end{aligned} \quad (18)$$

458 Subsequently, we decompose  $\mathbf{W}^* \in \partial \|\boldsymbol{\Sigma}^*\|_*$  as:

$$\mathbf{W}^* = \mathcal{P}(\mathbf{W}^*) + \mathcal{P}^\perp(\mathbf{W}^*) = \mathbf{U}^* \mathbf{U}^{*\top} - \lambda^{-1} \mathcal{P}^\perp(\nabla f(\boldsymbol{\Sigma}^*)), \quad (19)$$

459 which follows from the norm condition:

$$\|\lambda^{-1}\mathcal{P}^\perp(\nabla f(\boldsymbol{\Sigma}^*))\|_2 \leq \lambda^{-1}\|\nabla f(\boldsymbol{\Sigma}^*)\|_2 = \lambda^{-1}m^{-1}\|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2 \leq 1.$$

460 Combining the results from (18) and (19), we obtain

$$\begin{aligned} \langle \mathcal{P}^\perp(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle &= \langle \mathcal{P}^\perp(\nabla f(\boldsymbol{\Sigma}^*)) + \lambda \mathcal{P}^\perp(\mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathbf{0}, \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= 0. \end{aligned} \quad (20)$$

461 For  $i \in \mathcal{S}_1$ , under (d) in Assumption 1, since  $\sigma_i^* \geq \alpha$ , we have  $q'_\lambda(\sigma_i^*) = -\lambda$ . Consequently, we can  
462 project  $\nabla Q_\lambda(\boldsymbol{\Sigma}^*)$  onto  $\mathcal{T}_{\mathcal{S}_1}$  as follows:

$$\mathcal{P}_{\mathcal{S}_1}(\nabla Q_\lambda(\boldsymbol{\Sigma}^*)) = \mathbf{U}_{\mathcal{S}_1}^* q'_\lambda(\boldsymbol{\Gamma}_{\mathcal{S}_1}^*) \mathbf{U}_{\mathcal{S}_1}^{*\top} = -\lambda \mathbf{U}_{\mathcal{S}_1}^* \mathbf{U}_{\mathcal{S}_1}^{*\top},$$

463 where  $\boldsymbol{\Gamma}_{\mathcal{S}_1}^* \in \mathbb{R}^{r \times r}$ . Then, we have

$$\begin{aligned} &\langle \mathcal{P}_{\mathcal{S}_1}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathcal{P}_{\mathcal{S}_1}(\nabla f(\boldsymbol{\Sigma}^*)) + \mathcal{P}_{\mathcal{S}_1}(\lambda \mathbf{W}^*) + \mathcal{P}_{\mathcal{S}_1}(\nabla Q_\lambda(\boldsymbol{\Sigma}^*)), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathcal{P}_{\mathcal{S}_1}(\nabla f(\boldsymbol{\Sigma}^*)) + \mathcal{P}_{\mathcal{S}_1}(\lambda \mathbf{U}^* \mathbf{U}^{*\top}) - \lambda \mathbf{U}_{\mathcal{S}_1}^* \mathbf{U}_{\mathcal{S}_1}^{*\top}, \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathcal{P}_{\mathcal{S}_1}(\nabla f(\boldsymbol{\Sigma}^*)), \mathcal{P}_{\mathcal{S}_1}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*) \rangle. \end{aligned} \quad (21)$$

464 By applying Hölder's inequality and define  $\vartheta = \|\mathcal{P}_{\mathcal{S}_1}(\nabla f(\boldsymbol{\Sigma}^*))\|_2$ , we update (21) to

$$\langle \mathcal{P}_{\mathcal{S}_1}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \geq -\vartheta \|\mathcal{P}_{\mathcal{S}_1}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*)\|_*. \quad (22)$$

465 Since  $\text{rank}(\mathcal{P}_{\mathcal{S}_1}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*)) \leq |\mathcal{S}_1|$ , we further refine (22) as

$$\langle \mathcal{P}_{\mathcal{S}_1}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \geq -\vartheta \sqrt{|\mathcal{S}_1|} \|\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*\|_F.$$

466 For  $i \in \mathcal{S}_2$ , applying Hölder's inequality, we derive

$$\begin{aligned} &\langle \mathcal{P}_{\mathcal{S}_2}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathcal{P}_{\mathcal{S}_2}(\nabla f(\boldsymbol{\Sigma}^*) + \nabla Q_\lambda(\boldsymbol{\Sigma})) + \lambda \mathbf{W}^*, \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &= \langle \mathcal{P}_{\mathcal{S}_2}(\nabla f(\boldsymbol{\Sigma}^*)), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle + \langle \mathcal{P}_{\mathcal{S}_2}(\nabla Q_\lambda(\boldsymbol{\Sigma})), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle + \langle \mathcal{P}_{\mathcal{S}_2}(\lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \rangle \\ &\geq - \left[ \underbrace{\|\mathcal{P}_{\mathcal{S}_2}(\nabla f(\boldsymbol{\Sigma}^*))\|_2}_\text{I} + \underbrace{\|\mathcal{P}_{\mathcal{S}_2}(\nabla Q_\lambda(\boldsymbol{\Sigma}))\|_2}_\text{II} + \underbrace{\|\mathcal{P}_{\mathcal{S}_2}(\lambda \mathbf{W}^*)\|_2}_\text{III} \right] \cdot \|\mathcal{P}_{\mathcal{S}_2}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*)\|_2. \end{aligned} \quad (23)$$

467 It remains to analysis the terms I, II and III in (23) individually.

468 For term I, we have

$$\|\mathcal{P}_{\mathcal{S}_2}(\nabla f(\boldsymbol{\Sigma}^*))\|_2 \leq \|\nabla f(\boldsymbol{\Sigma}^*)\|_2 \leq \lambda. \quad (24)$$

469 For term II, we have

$$\|\mathcal{P}_{\mathcal{S}_2}(\nabla Q_\lambda(\boldsymbol{\Sigma}^*))\|_2 = \left\| \mathbf{U}_{\mathcal{S}_1}^* q'_\lambda(\boldsymbol{\Gamma}_{\mathcal{S}_1}^*) \mathbf{U}_{\mathcal{S}_2}^{*\top} \right\|_2.$$

470 Additionally, under (a) in Assumption 1, for any  $i \in \mathcal{S}_2$ ,  $\sigma_i^*$  satisfies  $\alpha > \sigma_i^* > 0$ , we have  $|q'_\lambda(\sigma_i^*)| \leq$   
471  $\lambda$ . This leads to

$$\|\mathcal{P}_{\mathcal{S}_2}(\nabla Q_\lambda(\boldsymbol{\Sigma}))\|_2 \leq \lambda. \quad (25)$$



472 For term III, we have

$$\|\mathcal{P}_{\mathcal{S}_2}(\lambda \mathbf{W}^*)\|_2 \leq \|\mathcal{P}_{\mathcal{S}_2}(\lambda \mathbf{U}^* \mathbf{U}^{*\top})\|_2 = \lambda. \quad (26)$$

473 By substituting (24), (25) and (26) into (23), we have

$$\left\langle \mathcal{P}_{\mathcal{S}_2}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\rangle \geq -3\lambda \left\| \mathcal{P}_{\mathcal{S}_2}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*) \right\|_2.$$

474 Since  $\text{rank}(\mathcal{P}_{\mathcal{S}_2}(\hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^*)) \leq |\mathcal{S}_2|$ , we have

$$\left\langle \mathcal{P}_{\mathcal{S}_2}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\rangle \geq -3\lambda \sqrt{|\mathcal{S}_2|} \left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F. \quad (27)$$

475 By substituting (20), (22) and (27) into (17), we have

$$\begin{aligned} & (\rho^- - \zeta^-) \left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F^2 \\ & \leq - \left\langle \mathcal{P}_{\mathcal{S}_1}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\rangle - \left\langle \mathcal{P}_{\mathcal{S}_2}(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\rangle \\ & \quad - \left\langle \mathcal{P}^\perp(\nabla \tilde{f}(\boldsymbol{\Sigma}^*) + \lambda \mathbf{W}^*), \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\rangle \\ & \leq \vartheta \sqrt{|\mathcal{S}_1|} \left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F + 3\lambda \sqrt{|\mathcal{S}_2|} \left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F. \end{aligned}$$

476 Consequently, we get

$$\left\| \hat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F \leq \frac{\vartheta \sqrt{|\mathcal{S}_1|}}{\rho^- - \zeta^-} + \frac{3\lambda \sqrt{|\mathcal{S}_2|}}{\rho^- - \zeta^-}.$$

477

□

### 478 A.3. Proof of Corollary 7

479 *Proof.* According to the definitions of  $f(\cdot)$ , we have

$$\begin{aligned} f(\hat{\boldsymbol{\Sigma}}^O) - f(\boldsymbol{\Sigma}^*) &= \frac{1}{2m} \left\| \mathbf{y} - \mathcal{A}(\hat{\boldsymbol{\Sigma}}^O) \right\|_2^2 - \frac{1}{2m} \left\| \mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*) \right\|_2^2 \\ &= \frac{1}{2m} \left\| \boldsymbol{\eta} - \mathcal{A}(\hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^*) \right\|_2^2 - \frac{1}{2m} \left\| \boldsymbol{\eta} \right\|_2^2 \\ &= \frac{1}{2m} \left\| \mathcal{A}(\hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^*) \right\|_2^2 - \frac{1}{m} \left\langle \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*)), \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\rangle \\ &= \frac{1}{2m} \left\| \mathcal{A}(\hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^*) \right\|_2^2 + \left\langle \nabla f(\boldsymbol{\Sigma}^*), \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\rangle \leq 0, \end{aligned} \quad (28)$$

480 where the last inequality applies due to  $\hat{\boldsymbol{\Sigma}}^O$  is the oracle estimator.

481 Next, applying the assumption 3, we have

$$f(\hat{\boldsymbol{\Sigma}}^O) \geq f(\boldsymbol{\Sigma}^*) + \left\langle \nabla f(\boldsymbol{\Sigma}^*), \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\rangle + \frac{\rho^-}{2} \left\| \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\|_F^2. \quad (29)$$

482 By combining inequalities (28) and (29), we obtain

$$\begin{aligned} \frac{\rho^-}{2} \left\| \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\|_F^2 &\leq f(\hat{\boldsymbol{\Sigma}}^O) - f(\boldsymbol{\Sigma}^*) - \left\langle \nabla f(\boldsymbol{\Sigma}^*), \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\rangle \\ &= \frac{1}{2m} \left\| \mathcal{A}(\hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^*) \right\|_2^2 \\ &\leq - \left\langle \nabla f(\boldsymbol{\Sigma}^*), \hat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\rangle. \end{aligned}$$

483 This inequality can be further refined using Hölder's inequality as follows

$$\left\| \hat{\Sigma}^O - \Sigma^* \right\|_F^2 \leq \frac{2 \left\langle \mathcal{P}(\nabla f(\Sigma^*)), \hat{\Sigma}^O - \Sigma^* \right\rangle}{\rho^-} \leq \frac{2 \|\mathcal{P}(\nabla f(\Sigma^*))\|_2 \cdot \left\| \hat{\Sigma}^O - \Sigma^* \right\|_*}{\rho^-}. \quad (30)$$

484 Given that  $\text{rank}(\hat{\Sigma}^O - \Sigma^*) = r$ , we have  $\left\| \hat{\Sigma}^O - \Sigma^* \right\|_* \leq \sqrt{r} \left\| \hat{\Sigma}^O - \Sigma^* \right\|_F$ . Using this in (30) gives  
 485 our desired result

$$\left\| \hat{\Sigma}^O - \Sigma^* \right\|_F^2 \leq \frac{2\sqrt{r} \|\mathcal{P}(\nabla f(\Sigma^*))\|_2 \cdot \left\| \hat{\Sigma}^O - \Sigma^* \right\|_F}{\rho^-}.$$

486 Thus, we have

$$\left\| \hat{\Sigma}^O - \Sigma^* \right\|_F \leq \frac{2\sqrt{r} \|\mathcal{P}(\nabla f(\Sigma^*))\|_2}{\rho^-}.$$

487 Next, we will prove that  $\hat{\Sigma}^O = \hat{\Sigma}$  under our assumptions. From Lemma 10, the following inequalities holds  
 488

$$\begin{aligned} \tilde{f}(\hat{\Sigma}) &\geq \tilde{f}(\hat{\Sigma}^O) + \left\langle \nabla \tilde{f}(\hat{\Sigma}^O), \hat{\Sigma} - \hat{\Sigma}^O \right\rangle + \frac{\rho^- - \zeta^-}{2} \left\| \hat{\Sigma} - \hat{\Sigma}^O \right\|_F^2, \\ \tilde{f}(\hat{\Sigma}^O) &\geq \tilde{f}(\hat{\Sigma}) + \left\langle \nabla \tilde{f}(\hat{\Sigma}), \hat{\Sigma}^O - \hat{\Sigma} \right\rangle + \frac{\rho^- - \zeta^-}{2} \left\| \hat{\Sigma}^O - \hat{\Sigma} \right\|_F^2. \end{aligned}$$

489 By adding the above two inequality, we obtain

$$\begin{aligned} &(\rho^- - \zeta^-) \left\| \hat{\Sigma} - \hat{\Sigma}^O \right\|_F^2 \\ &\leq \left\langle \nabla \tilde{f}(\hat{\Sigma}) - \nabla \tilde{f}(\hat{\Sigma}^O), \hat{\Sigma} - \hat{\Sigma}^O \right\rangle \\ &= \left\langle \nabla \tilde{f}(\hat{\Sigma}) + \lambda \widehat{\mathbf{W}}, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle - \left\langle \nabla \tilde{f}(\hat{\Sigma}^O) + \lambda \widehat{\mathbf{W}}^O, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle - \left\langle \lambda \widehat{\mathbf{W}} - \lambda \widehat{\mathbf{W}}^O, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle \\ &= \underbrace{\left\langle \nabla \tilde{f}(\hat{\Sigma}) + \lambda \widehat{\mathbf{W}}, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle}_I - \underbrace{\left\langle \nabla \tilde{f}(\hat{\Sigma}^O) + \lambda \widehat{\mathbf{W}}^O, \mathcal{P}(\hat{\Sigma}^O - \hat{\Sigma}) \right\rangle}_{III} \\ &\quad - \underbrace{\left\langle \nabla \tilde{f}(\hat{\Sigma}^O) + \lambda \widehat{\mathbf{W}}^O, \mathcal{P}^\perp(\hat{\Sigma}^O - \hat{\Sigma}) \right\rangle}_{III} - \underbrace{\left\langle \lambda \widehat{\mathbf{W}} - \lambda \widehat{\mathbf{W}}^O, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle}_{IV}. \end{aligned}$$

490 Then we analysis the terms I, II, III and IV, respectively.

491 Since  $\hat{\Sigma}$  is the optimal solution of (1), it satisfies the optimality condition. Therefore, we have

$$\left\langle \nabla \tilde{f}(\hat{\Sigma}) + \lambda \widehat{\mathbf{W}} + \hat{\Xi}, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle \leq 0, \quad (31)$$

492 where  $\hat{\Xi} \in \partial \mathbb{I}_{\{\Sigma \succeq 0\}}(\hat{\Sigma})$ . According to the definition of  $\mathbb{I}_{\{\Sigma \succeq 0\}}(\Sigma)$ , we have

$$\left\langle \hat{\Xi}, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle = 0. \quad (32)$$

493 By subtracting (32) from (31), we obtain

$$\left\langle \nabla \tilde{f}(\hat{\Sigma}) + \lambda \widehat{\mathbf{W}}, \hat{\Sigma} - \hat{\Sigma}^O \right\rangle \leq 0. \quad (33)$$

494 Based on the definition of the oracle estimator, the eigen-decompositions of  $\hat{\Sigma}^O$  is  $\hat{\Sigma}^O = \mathbf{U}^* \hat{\Gamma}^O \mathbf{U}^{*\top}$ .

495 Next, we decompose  $\widehat{\mathbf{W}}^O \in \partial \left\| \hat{\Sigma}^O \right\|_*$  as

$$\widehat{\mathbf{W}}^O = \mathcal{P}(\widehat{\mathbf{W}}^O) + \mathcal{P}^\perp(\widehat{\mathbf{W}}^O) = \mathbf{U}^* \mathbf{U}^{*\top} + \hat{\mathbf{Z}}^O,$$

496 where  $\widehat{\mathbf{Z}}^O \in \mathcal{T}^\perp$  and  $\|\widehat{\mathbf{Z}}^O\|_2 \leq 1$ . Since  $\sum_{i=1}^d p_\lambda(\sigma_i(\boldsymbol{\Sigma})) = \lambda \|\boldsymbol{\Sigma}\|_* + Q_\lambda(\boldsymbol{\Sigma})$ , we have

$$\begin{aligned}
\mathcal{P}\left(\nabla\left(\sum_{i=1}^d p_\lambda\left(\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right)\right)\right) &= \mathcal{P}\left(\lambda \widehat{\mathbf{W}}^O + \nabla Q_\lambda\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right) \\
&= \mathcal{P}\left(\lambda\left(\mathbf{U}^* \mathbf{U}^{*\top} + \widehat{\mathbf{Z}}^O\right) + \nabla Q_\lambda\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right) \\
&= \mathcal{P}\left(\lambda \mathbf{U}^* \mathbf{U}^{*\top} + \nabla Q_\lambda\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right) \\
&= \lambda \mathbf{U}^* \mathbf{U}^{*\top} + \mathbf{U}^* q'_\lambda\left(\left(\widehat{\boldsymbol{\Gamma}}^O\right)_S\right) \mathbf{U}^{*\top} \\
&= \mathbf{U}^* \left(\lambda \mathbf{I} + q'_\lambda\left(\left(\widehat{\boldsymbol{\Gamma}}^O\right)_S\right)\right) \mathbf{U}^{*\top}, \tag{34}
\end{aligned}$$

497 and the diagonal matrix  $\left(\lambda \mathbf{I} + q'_\lambda\left(\left(\widehat{\boldsymbol{\Gamma}}^O\right)_S\right)\right)$  satisfies

$$\left[\lambda \mathbf{I} + q'_\lambda\left(\left(\widehat{\boldsymbol{\Gamma}}^O\right)_S\right)\right]_{ii} = \lambda + q'_\lambda\left(\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right) = p'_\lambda\left(\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right), \tag{35}$$

498 where  $i \in \mathcal{S}$ . According to the perturbation bounds for eigenvalues, we have

$$\begin{aligned}
\min_{i \in \mathcal{S}} \left| \sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right) \right| &= \min_{i \in \mathcal{S}} \left| \sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right) - \sigma_i^* + \sigma_i^* \right| \\
&\geq -\max_{i \in \mathcal{S}} \left| \sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right) - \sigma_i^* \right| + \min_{i \in \mathcal{S}} |\sigma_i^*| \\
&\geq -\left\| \widehat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\|_2 + \min_{i \in \mathcal{S}} |\sigma_i^*| \\
&\geq -\left\| \widehat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\|_F + \min_{i \in \mathcal{S}} |\sigma_i^*| \\
&\geq -\frac{2\sqrt{r} \|\mathcal{P}(\nabla f(\boldsymbol{\Sigma}^*))\|_2}{\rho^-} + \min_{i \in \mathcal{S}} |\sigma_i^*|. \tag{36}
\end{aligned}$$

499 Since  $\boldsymbol{\sigma}^*$  satisfies  $\min_{i \in \mathcal{S}} |\sigma_i^*| \geq \alpha + \frac{2\sqrt{r} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2}{m\rho^-} = \alpha + \frac{2\sqrt{r} \|\mathcal{P}(\nabla f(\boldsymbol{\Sigma}^*))\|_2}{\rho^-}$ , (36) can be updated  
500 to

$$\min_{i \in \mathcal{S}} \left| \sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right) \right| \geq \alpha.$$

501 Given  $\widehat{\boldsymbol{\Sigma}}^O \in \mathcal{T}$ , we have  $\text{rank}\left(\widehat{\boldsymbol{\Sigma}}^O\right) = r$ , then  $\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)$  satisfies

$$\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right) \begin{cases} \geq \alpha & i \in \mathcal{S}, \\ = 0 & i \notin \mathcal{S}. \end{cases}$$

502 Under (d) in Assumption 1, for  $i \in \mathcal{S}$ , we have

$$p'_\lambda\left(\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right) = 0. \tag{37}$$

503 Then combining (34), (35) and (37), we get

$$\mathcal{P}\left(\nabla\left(\sum_{i=1}^d p_\lambda\left(\sigma_i\left(\widehat{\boldsymbol{\Sigma}}^O\right)\right)\right)\right) = \mathbf{0}.$$

504 Since  $\widehat{\boldsymbol{\Sigma}}^O$  is the oracle estimator, it satisfies the optimality condition. Therefore, we have the follow-  
505 ing inequality

$$\left\langle \nabla f\left(\widehat{\boldsymbol{\Sigma}}^O\right) + \widehat{\boldsymbol{\Xi}}^O, \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right\rangle \leq 0,$$

506 where  $\widehat{\boldsymbol{\Xi}}^O \in \partial \mathbb{I}_{\{\boldsymbol{\Sigma} \succeq \mathbf{0}\}}\left(\widehat{\boldsymbol{\Sigma}}^O\right)$ .

507 According to the definition of  $\mathbb{I}_{\{\Sigma \succeq 0\}}(\Sigma)$ , we have

$$\langle \widehat{\Sigma}^O, \widehat{\Sigma}^O - \widehat{\Sigma} \rangle = 0.$$

508 Then we have that

$$\langle \nabla f(\widehat{\Sigma}^O), \mathcal{P}(\widehat{\Sigma}^O - \widehat{\Sigma}) \rangle \leq \langle \nabla f(\widehat{\Sigma}^O), \widehat{\Sigma}^O - \widehat{\Sigma} \rangle \leq 0.$$

509 Therefore, we have

$$\begin{aligned} & \langle \nabla \tilde{f}(\widehat{\Sigma}^O) + \lambda \widehat{\mathbf{W}}^O, \mathcal{P}(\widehat{\Sigma}^O - \widehat{\Sigma}) \rangle \\ &= \langle \nabla f(\widehat{\Sigma}^O), \mathcal{P}(\widehat{\Sigma}^O - \widehat{\Sigma}) \rangle + \left\langle \mathcal{P} \left( \nabla \left( \sum_{i=1}^d p_{\lambda}(\sigma_i(\widehat{\Sigma}^O)) \right) \right), \mathcal{P}(\widehat{\Sigma}^O - \widehat{\Sigma}) \right\rangle \\ &\leq 0. \end{aligned} \tag{38}$$

510 By projecting  $\nabla Q_{\lambda}(\widehat{\Sigma}^O)$  into  $\mathcal{T}^{\perp}$ , we have

$$\begin{aligned} \mathcal{P}^{\perp}(\nabla Q_{\lambda}(\widehat{\Sigma}^O)) &= (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \mathbf{U}^* q'_{\lambda}(\widehat{\mathbf{\Gamma}}^O) \mathbf{U}^{*\top} (\mathbf{I} - \mathbf{U}^* \mathbf{U}^{*\top}) \\ &= (\mathbf{U}^* - \mathbf{U}^*) q'_{\lambda}(\widehat{\mathbf{\Gamma}}^O) (\mathbf{U}^* - \mathbf{U}^*) \\ &= \mathbf{0}. \end{aligned}$$

511 Under Assumption 4, it holds that

$$f(\Sigma^*) + \langle \nabla f(\Sigma^*), \widehat{\Sigma}^O - \Sigma^* \rangle + \frac{\rho^+}{2} \|\widehat{\Sigma}^O - \Sigma^*\|_{\mathbf{F}}^2 \geq f(\widehat{\Sigma}^O).$$

512

$$f(\widehat{\Sigma}^O) + \langle \nabla f(\widehat{\Sigma}^O), \Sigma^* - \widehat{\Sigma}^O \rangle + \frac{\rho^+}{2} \|\Sigma^* - \widehat{\Sigma}^O\|_{\mathbf{F}}^2 \geq f(\Sigma^*).$$

513 Adding the above two inequalities up, we get

$$\begin{aligned} \rho^+ \|\widehat{\Sigma}^O - \Sigma^*\|_{\mathbf{F}}^2 &\geq \langle \nabla f(\widehat{\Sigma}^O) - \nabla f(\Sigma^*), \widehat{\Sigma}^O - \Sigma^* \rangle \\ &\geq \|\nabla f(\widehat{\Sigma}^O) - \nabla f(\Sigma^*)\|_{\mathbf{F}} \cdot \|\widehat{\Sigma}^O - \Sigma^*\|_{\mathbf{F}}, \end{aligned}$$

514 which can be simplified as

$$\|\nabla f(\widehat{\Sigma}^O) - \nabla f(\Sigma^*)\|_{\mathbf{F}} \leq \rho^+ \|\widehat{\Sigma}^O - \Sigma^*\|_{\mathbf{F}}.$$

515 Thus, we have

$$\begin{aligned} \|\mathcal{P}^{\perp}(\nabla f(\widehat{\Sigma}^O))\|_2 &\leq \|\nabla f(\widehat{\Sigma}^O)\|_2 \\ &\leq \|\nabla f(\widehat{\Sigma}^O) - \nabla f(\Sigma^*)\|_2 + \|\nabla f(\Sigma^*)\|_2 \\ &\leq \|\nabla f(\widehat{\Sigma}^O) - \nabla f(\Sigma^*)\|_{\mathbf{F}} + \|\nabla f(\Sigma^*)\|_2 \\ &\leq \rho^+ \|\widehat{\Sigma}^O - \Sigma^*\|_{\mathbf{F}} + \|\nabla f(\Sigma^*)\|_2 \\ &\leq \frac{2\rho^+ \sqrt{r} \|\mathcal{P}(\nabla f(\Sigma^*))\|_2}{\rho^-} + \|\nabla f(\Sigma^*)\|_2 \\ &\leq \frac{2\rho^+ \sqrt{r} \|\nabla f(\Sigma^*)\|_2}{\rho^-} + \|\nabla f(\Sigma^*)\|_2 \\ &= \frac{2\rho^+ \sqrt{r} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma))\|_2}{m\rho^-} + \frac{\|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma))\|_2}{m}. \end{aligned}$$

Under the condition that  $\lambda \geq 2m^{-1} \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2 + 2m^{-1} \sqrt{r} \rho^+ \|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2 / \rho^-$ , the upper bound of  $\left\| \mathcal{P}^\perp \left( \nabla f \left( \widehat{\boldsymbol{\Sigma}}^O \right) \right) \right\|_2$  satisfies

$$\left\| \mathcal{P}^\perp \left( \nabla f \left( \widehat{\boldsymbol{\Sigma}}^O \right) \right) \right\|_2 \leq \lambda.$$

So we can decompose  $\widehat{\mathbf{W}}^O \in \partial \left\| \widehat{\boldsymbol{\Sigma}}^O \right\|_*$  as

$$\widehat{\mathbf{W}}^O = \mathcal{P}(\widehat{\mathbf{W}}^O) + \mathcal{P}^\perp(\widehat{\mathbf{W}}^O) = \mathbf{U}^* \mathbf{U}^{*\top} - \lambda^{-1} \mathcal{P}^\perp \left( \nabla f \left( \widehat{\boldsymbol{\Sigma}}^O \right) \right).$$

Then we get

$$\begin{aligned} & \left\langle \nabla \tilde{f} \left( \widehat{\boldsymbol{\Sigma}}^O \right) + \lambda \widehat{\mathbf{W}}^O, \mathcal{P}^\perp \left( \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right) \right\rangle \\ &= \left\langle \mathcal{P}^\perp \left( \nabla f \left( \widehat{\boldsymbol{\Sigma}}^O \right) + \lambda \widehat{\mathbf{W}}^O \right), \mathcal{P}^\perp \left( \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right) \right\rangle + \left\langle \mathcal{P}^\perp \left( \nabla Q_\lambda \left( \widehat{\boldsymbol{\Sigma}}^O \right) \right), \mathcal{P}^\perp \left( \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right) \right\rangle \\ &= \left\langle \mathbf{0}, \mathcal{P}^\perp \left( \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right) \right\rangle + \left\langle \mathbf{0}, \mathcal{P}^\perp \left( \widehat{\boldsymbol{\Sigma}}^O - \widehat{\boldsymbol{\Sigma}} \right) \right\rangle \\ &= 0. \end{aligned} \tag{39}$$

Additionally, since  $\|\cdot\|_*$  is convex,  $\widehat{\mathbf{W}} \in \partial \left\| \widehat{\boldsymbol{\Sigma}} \right\|_*$  and  $\widehat{\mathbf{W}}^O \in \partial \left\| \boldsymbol{\Sigma}^* \right\|_*$ , we have

$$\left\langle \lambda \widehat{\mathbf{W}} - \lambda \widehat{\mathbf{W}}^O, \widehat{\boldsymbol{\Sigma}} - \widehat{\boldsymbol{\Sigma}}^O \right\rangle \geq 0. \tag{40}$$

Therefore, by adding up (33), (38), (39) and (40), we obtain

$$(\rho^- - \zeta^-) \left\| \widehat{\boldsymbol{\Sigma}} - \widehat{\boldsymbol{\Sigma}}^O \right\|_F^2 \leq 0.$$

Since  $\rho^- > \zeta^-$ , it follows that  $\widehat{\boldsymbol{\Sigma}}^O = \widehat{\boldsymbol{\Sigma}}$ . Thus, we have

$$\left\| \widehat{\boldsymbol{\Sigma}} - \boldsymbol{\Sigma}^* \right\|_F = \left\| \widehat{\boldsymbol{\Sigma}}^O - \boldsymbol{\Sigma}^* \right\|_F \leq \frac{2\sqrt{r} \|\mathcal{P}(\nabla f(\boldsymbol{\Sigma}^*))\|_2}{\rho^-}$$

□

## A.4. Proof of Corollary 9

First, we give three essential Lemmas.

**Lemma 12** (Proposition 1 in [29]). *Consider the sampling operator of  $\boldsymbol{\Theta}$ -ensemble, it holds with probability at least  $1 - 2 \exp(-m/32)$  that*

$$\frac{\|\mathcal{A}(\boldsymbol{\Delta})\|_2}{\sqrt{m}} \geq \frac{1}{4} \left\| \sqrt{\boldsymbol{\Theta}} \text{vec}(\boldsymbol{\Delta}) \right\|_2 - 24\pi(\boldsymbol{\Theta}) \sqrt{d/m} \|\boldsymbol{\Delta}\|_*.$$

**Lemma 13** (Lemma 6 in [29]). *There're universal constants  $C_0, C_1$  and  $C_2$  such that*

$$\mathbb{P} \left[ \frac{\|\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}))\|_2}{m} \geq 2C_0 \delta \varpi \sqrt{d/m} \right] \leq C_1 \exp(-2dC_2).$$

**Lemma 14** (Lemma D.1 in [49]). *Let  $\mathbf{x}$  be a sub-Gaussian random vector with zero mean and covariance  $\boldsymbol{\Sigma}^*$  and  $\{\mathbf{x}_i\}_{i=1}^T$  be a collection of i.i.d. samples from  $\mathbf{x}$ . There exists some constants  $c_1, c_2$ , and  $z_0$  such that for all  $z$  with  $0 < z < z_0$ , the sample covariance matrix  $\boldsymbol{\Sigma}_T$  satisfies the following tail bound*

$$\mathbb{P} \left( |\Sigma_{ij}^* - (\Sigma_T)_{ij}| > z \right) \leq c_1 \exp(-c_2 T z^2).$$

Then, we give the proof of Corollary 9.

533 *Proof.* By projecting  $\|\hat{\Delta}\|_*$  into  $\mathcal{T}$  and  $\mathcal{T}^\perp$ , we have

$$\|\hat{\Delta}\|_* = \|\mathcal{P}(\hat{\Delta})\|_* + \|\mathcal{P}^\perp(\hat{\Delta})\|_* \leq 6 \|\mathcal{P}(\hat{\Delta})\|_* = 6\sqrt{r} \|\mathcal{P}(\hat{\Delta})\|_F \leq 6\sqrt{r} \|\hat{\Delta}\|_F.$$

534 Then according to Lemma 12, we have

$$\begin{aligned} \frac{\|\mathcal{A}(\hat{\Delta})\|_2}{\sqrt{m}} &\geq \frac{1}{4} \|\sqrt{\Theta} \text{vec}(\hat{\Delta})\|_2 - 24\pi(\Theta) \sqrt{d/m} \|\hat{\Delta}\|_* \\ &\geq \frac{\sqrt{\lambda_{\min}(\Theta)}}{4} \|\hat{\Delta}\|_F - 144\sqrt{r}\pi(\Theta) \sqrt{d/m} \|\hat{\Delta}\|_F. \end{aligned}$$

535 Then for  $m \geq \frac{C_1 r \pi^2(\Theta) d^2}{\lambda_{\min}(\Theta)}$  where  $C_1$  is sufficiently large such that

$$144\sqrt{r}\pi(\Theta) \sqrt{d/m} \leq \frac{\sqrt{\lambda_{\min}(\Theta)}}{8},$$

536 we have

$$\frac{\|\mathcal{A}(\hat{\Delta})\|_2}{\sqrt{m}} \geq \frac{\sqrt{\lambda_{\min}(\Theta)}}{8} \|\hat{\Delta}\|_F,$$

537 which implies that

$$\frac{\|\mathcal{A}(\hat{\Delta})\|_2^2}{m} \geq \frac{\lambda_{\min}(\Theta)}{64} \|\hat{\Delta}\|_F^2,$$

538 for  $\rho^- = C_2 \lambda_{\min}(\Theta)$ , the above inequality can be updated to

$$\frac{\|\mathcal{A}(\hat{\Delta})\|_2^2}{m} \geq \frac{\rho^-}{64C_2} \|\hat{\Delta}\|_F^2,$$

539 when we set  $C_2 = \frac{1}{32}$ , the Restricted Strong Convexity of  $\mathcal{A}(\cdot)$  is established. Additionally, we have

$$\|\mathcal{P}_{S_1}(\nabla f(\Sigma))\|_2 = \|\mathbf{U}_{S_1}^{*\top} \nabla f(\Sigma) \mathbf{U}_{S_1}^*\|_2.$$

540 According to Lemma 13, since  $\mathbf{U}_{S_1}^{*\top} \nabla f(\Sigma) \mathbf{U}_{S_1}^* \in \mathbb{R}^{d \times d}$ , combining with Lemma 14, we have

$$\begin{aligned} \|\mathbf{U}^{*\top} \nabla f(\Sigma) \mathbf{U}^*\|_2 &\leq 2C_0 \delta \varpi \sqrt{d/(mT)}, \\ \|\mathbf{U}_{S_1}^{*\top} \nabla f(\Sigma) \mathbf{U}_{S_1}^*\|_2 &\leq 2C_0 \delta \varpi \sqrt{|S_1|/(mT)}, \end{aligned}$$

541 which holds with probability at least  $1 - C_1 \exp(-2|S_1|C_2)$ . Based on Theorem 5, since  $\vartheta =$

542  $\|\mathcal{P}_{S_1}(\nabla \tilde{f}(\Sigma^*))\|_2$ ,  $\lambda = \|\nabla f(\Sigma^*)\|_2$  and  $\rho^- = C_2 \lambda_{\min}(\Theta)$

$$\begin{aligned} \|\hat{\Sigma} - \Sigma^*\|_F &\leq \frac{2\vartheta \sqrt{|S_1|}}{\rho^-} + \frac{6\lambda \sqrt{|S_2|}}{\rho^-} \\ &\leq \frac{1}{C_2 \lambda_{\min}(\Theta)} \left( 4C_0 \delta \varpi \sqrt{|S_1|/(mT)} \cdot \sqrt{|S_1|} + 12C_0 \delta \varpi \sqrt{d/(mT)} \cdot \sqrt{|S_2|} \right) \\ &\leq \frac{\delta \varpi}{\lambda_{\min}(\Theta) \sqrt{mT}} \left( 4C_0 |S_1|/C_2 + 12C_0 \sqrt{d}/C_2 \right), \end{aligned}$$

543 which holds with probability at least  $1 - C_1 \exp(-2|S_1|C_2)$ . Therefore, by choosing appropriate

544 value of  $C_0, C_1, C_2, C_5$  and  $C_6$ , we get

$$\|\hat{\Sigma} - \Sigma^*\|_F \leq \frac{\delta \varpi}{\sqrt{mT} \lambda_{\min}(\Theta)} \left[ C_1 |S_1| + C_2 \sqrt{|S_2|d} \right],$$

545 which holds with probability at least  $1 - C_5 \exp(-2dC_6)$ . Since  $\rho^- = C_1 \lambda_{\min}(\Theta)$  and Lemma 13,  
 546 according to Corollary 7, we have

$$\begin{aligned} \min_{i \in \mathcal{S}} |\sigma_i^*| &\geq \alpha + \frac{2\sqrt{r} \|\sum_{i=1}^m \eta_i \mathbf{A}_i\|_2}{m\rho^-} \\ &\geq \alpha + \frac{C_2 \delta \varpi \sqrt{rd/m}}{\lambda_{\min}(\Theta)} \\ &\geq \alpha + \frac{C_2 m^{-1} \|\sum_{i=1}^m \eta_i \mathbf{A}_i\|_2 \sqrt{r}}{2\rho^-}, \end{aligned}$$

547 which holds with probability at least  $1 - C_1 \exp(-2dC_2)$ . With suitable value of  $C_2$ , the above in-  
 548 equality is equal to

$$\min_{i \in \mathcal{S}} |\sigma_i^*| \geq \alpha + \frac{2\sqrt{r} \|\sum_{i=1}^m \eta_i \mathbf{A}_i\|_2}{m\rho^-},$$

549 which holds with probability at least  $1 - C_1 \exp(-2dC_2)$  and is the requirement of  $\sigma^*$  in Corollary  
 550 7. According to the definition of  $f(\cdot)$ , we have

$$\mathbf{H}_m = \frac{\sum_{i=1}^m \text{vec}(\mathbf{A}_i) \text{vec}(\mathbf{A}_i)^\top}{m},$$

551 where  $\mathbf{H}_m$  is the Hessian matrix of  $f(\cdot)$ . Then we have  $\mathbb{E}(\mathbf{H}_m) = \Theta$ .

552 By concentration, when  $m$  is sufficiently large, with high probability,  $\lambda_{\max}(\mathbf{H}_m) \leq 2\lambda_{\max}(\Theta)$ . Based  
 553 on 4, we could obtain that  $\rho^+ = \lambda_{\max}(\mathbf{H}_m)$ . Then we have that when  $m$  is sufficiently large, with  
 554 high probability,  $\rho^+ \leq 2\lambda_{\max}(\Theta)$ . Additionally, given  $\rho^- = C_1 \lambda_{\min}(\Theta)$  and Lemma 13, we have

$$\begin{aligned} \lambda &\geq C_3 \delta \varpi \sqrt{d/m} \left( 1 + \frac{\sqrt{r} \lambda_{\max}(\Theta)}{\lambda_{\min}(\Theta)} \right) \\ &\geq C_3 \delta \varpi \sqrt{d/m} \left( 1 + \frac{\sqrt{r} \rho^+ C_1}{2\rho^-} \right) \\ &\geq \frac{C_3}{2C_1} m^{-1} \left\| \sum_{i=1}^m \eta_i \mathbf{A}_i \right\|_2 \left( 1 + \frac{\sqrt{r} \rho^+ C_1}{2\rho^-} \right), \end{aligned}$$

555 which holds with probability at least  $1 - C_1 \exp(-2dC_2)$ . With suitable value of  $C_1$  and  $C_3$ , the above  
 556 inequality is equal to

$$\lambda \geq 2m^{-1} \left\| \sum_{i=1}^m \eta_i \mathbf{A}_i \right\|_2 + 2m^{-1} \sqrt{r} \rho^+ \left\| \sum_{i=1}^m \eta_i \mathbf{A}_i \right\|_2 / \rho^-,$$

557 which holds with probability at least  $1 - C_1 \exp(-2dC_2)$  and is the requirement of  $\lambda$  in Corollary 7.  
 558 Since the requirements of  $\sigma^*$  and  $\lambda$  are reached with probability at least  $1 - C_1 \exp(-2dC_2)$ , accord-  
 559 ing to Corollary 7, we have that  $\hat{\Sigma}$  is the optimal solution of the estimator,  $\text{rank}(\hat{\Sigma}) = \text{rank}(\Sigma^*) = r$ ,  
 560 with probability at least  $1 - C_4 \exp(-C_5 d)$  by setting the proper values of  $C_4$  and  $C_5$ .

561 In addition, we have

$$\begin{aligned} \left\| \hat{\Sigma} - \Sigma^* \right\|_F &\leq \frac{2\sqrt{r} \left\| \mathcal{P}_{S_1}(\nabla \tilde{f}(\Sigma^*)) \right\|_2}{\rho^-} \\ &\leq \frac{2\sqrt{r} \left\| \mathcal{P}_{S_1}(\nabla \tilde{f}(\Sigma^*)) \right\|_2}{C_1 \lambda_{\min}(\Theta)} \\ &\leq \frac{4\sqrt{r} C_0 \delta \varpi \sqrt{r/(mT)}}{C_1 \lambda_{\min}(\Theta)} \\ &\leq \frac{4C_0 r \delta \varpi}{C_1 \sqrt{mT} \lambda_{\min}(\Theta)} \\ &\leq \frac{C_6 r \delta \varpi}{\sqrt{mT} \lambda_{\min}(\Theta)}, \end{aligned}$$

562 which holds with probability at least  $1 - C_4 \exp(-C_5 d)$  with suitable value of  $C_1, C_2, C_4, C_5$  and  
563  $C_6$ .  $\square$