

Large Sparse Covariance Matrix Sensing from Quadratic Measurements

Wenbin Wang

School of Information Science and Technology
ShanghaiTech University
Shanghai, China
wangwb2023@shanghaitech.edu.cn

Ziping Zhao

School of Information Science and Technology
ShanghaiTech University
Shanghai, China
zipingzhao@shanghaitech.edu.cn

Abstract—Covariance matrix, characterizing the correlation property among multiple random variables, plays a pivotal role in many data analytics areas. A common approach is to directly estimate the covariance matrix from samples, yet it faces substantial challenges when data changes rapidly and the acquisition devices have limited processing power and storage capacity. To address these issues, a classical solution is covariance matrix sensing, which involves compressing signals and estimating the covariance matrix from the sensing measurements. In this paper, we consider a quadratic measurement model and adopt the sparsity assumption, which is prevalent in the high dimensional setting. We propose a regularized least squares type estimator, which can guarantee the positive definiteness and sparse property of the covariance matrices. To obtain the estimator, a multistage convex relaxation algorithm based on majorization-minimization (MM) is developed. We clearly establish the statistical performance of all the sequential approximate solutions produced by the MM-based algorithm, and prove after sufficient iterations the final estimator can actually attain the oracle statistical rate. The theoretical findings are supported by numerical simulations.

Index Terms—Covariance matrix sensing, quadratic measurements, majorization-minimization, positive definiteness, sparsity, non-convex statistical optimization.

I. INTRODUCTION

The covariance matrix, explaining the variability and relationship between different variables, plays a central role in statistical inference and information processing [1]–[3]. A common method for covariance matrix estimation is directly based on the samples. However, this approach often encounters significant challenges due to the rapidly changing data, and the limited processing capability and storage capacity of sampling equipment. Thus, it is desirable to consider recovering covariance matrices from a single data stream and a minimal number of storage measurements. Covariance matrix sensing, a.k.a., covariance matrix sketching, is a classical solution to address these challenging issues. It involves two key processes: sensing the signals, and estimating the covariance matrix from potentially noisy measurements [4].

In this paper, we consider a quadratic covariance matrix sensing model. Denote the ground truth covariance matrix as $\Sigma \in \mathbb{R}^{d \times d}$. We obtain m measurements as follows:

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

where $\{y_i\}_{i=1}^m$ denotes a sequence of measurements, $\{\mathbf{a}_i \in \mathbb{R}^d\}_{i=1}^m$ are the sensing vectors, and $\{\eta_i\}_{i=1}^m$ represent the noise. This model has applications in many fields such as wireless communications [5], optical imaging [6], environmental monitoring [7], etc. In wireless communications, the power spectrum characterizes the distribution of power of a signal across different frequencies, providing crucial insights into the signal properties [8], [9]. Consequently, estimating the power spectrum is vital. Power spectrum estimation enables the reconstruction of signals from far fewer samples than traditionally required by the Nyquist-Shannon sampling theorem, used to infer the spectral characteristics of a signal [10]. This technique focuses on capturing the spectral properties of random signals based on energy observations over time [5]. In this approach, sensing vectors $\{\mathbf{a}_i \in \mathbb{R}^d\}_{i=1}^m$ are employed to interact with the signal \mathbf{x}_t at discrete times, where the vector \mathbf{a}_i could be implemented using techniques such as random demodulators [11]. The average energy of the signal over N observations is calculated as follows:

$$y_i = \frac{1}{N} \sum_{t=1}^N (\mathbf{a}_i^\top \mathbf{x}_t)^2 = \mathbf{a}_i^\top \Sigma_N \mathbf{a}_i, \quad \text{for } i = 1, \dots, m, \quad (2)$$

where $\Sigma_N = \frac{1}{N} \sum_{t=1}^N \mathbf{x}_t \mathbf{x}_t^\top$ symbolizes the sample covariance matrix. The power spectrum can then be inferred by analyzing the structure of the covariance matrix, typically through spectral analysis techniques like Fourier transforms. This sampling model ties the observed energy directly to the power spectrum, providing a robust way to estimate the spectral distribution of the signal based purely on energy measurements, making it particularly suitable for high-frequency applications where phase measurements may be difficult to obtain [5]. In high-dimensional data processing, covariance sketching emerges as a pivotal application of the quadratic model, particularly advantageous in settings where the data dimensions and volume surpass the available computational resources [4], [12]. This technique fundamentally entails the estimation of covariance matrices through the use of compressed or "sketched" data representations. Specifically, the quadratic nature of the model is employed to transform and reduce the dimensionality of data streams or samples, which are then used to approximate the covariance matrix efficiently. For

instance, considering a sequence of high-dimensional vectors $\mathbf{x}_t \in \mathbb{R}^d$, the process involves selecting a random sketching vector \mathbf{a}_{i_t} at each time t and computing the quadratic sketch $y_i = (\mathbf{a}_{i_t}^\top \mathbf{x}_t)^2$. These sketches are subsequently aggregated and normalized to form an estimate of the covariance matrix $\Sigma = \mathbb{E}(\mathbf{x}_t \mathbf{x}_t^\top)$, with adjustments for error terms denoted by η_i . Such methodologies not only facilitate the processing of data with limited memory and computational capabilities but also enhance the ability to handle real-time data streams efficiently, thus exemplifying the practical applications of quadratic models in modern statistical data analysis [13]. Overall, each of these applications necessitates the recovery of covariance matrices from a minimal set of quadratic measurements (1).

The covariance matrix sensing technique is inspired by the success of compressed sensing [14], which asserts it is possible to sense and compress signals simultaneously without lose information. To achieve it effectively, exploiting low-dimensional structures inherent in high-dimensional data is a common approach [15], with sparsity being one of the widely adopted assumption [16], i.e., most off-diagonal elements are almost zero, significantly decreasing the number of parameters that need estimation. Under approximately sparse assumption, [8] and [17] investigated the recovery of second-order statistics for nearly sparse cyclostationary signals through random linear measurements, employing ℓ_1 -minimization but without performance guarantees. Another pioneering work by Dasarthy et al. [4] have suggested estimating a nearly sparse covariance matrix via measurements formatted as $\mathbf{Y} = \mathbf{A}\Sigma\mathbf{A}^\top$, where \mathbf{A} represents a sketching matrix based on expander graphs.

The ℓ_1 penalty has been utilized in covariance matrix sensing. [18] derived a fundamental guarantee on how many samples are sufficient to approximate the ground truth by ℓ_1 regularization and showed that covariance estimation from compressive measurements can be highly robust under sparse assumption. However, it is widely acknowledge that it introduces a non-negligible bias into the estimator. To mitigate this bias, alternatives like the non-convex smoothly clipped absolute deviation (SCAD) penalty [19] and minimax concave penalty (MCP) [20] have been suggested. The benefits of non-convex penalties can also be found in many other topics such as sparse covariance estimation [21], low-rank matrix recovery [22], high-dimensional graphical models [23], graph trend filtering [24], etc.

Based on these insights, we proposed to estimate covariance matrix from a small number of measurements using the non-convex penalty. And we adopts a sub-Gaussian error assumption, different from the error term in [18] which is assumed to be bounded in ℓ_1 or ℓ_2 norm. Analyzing the statistical optimality of oracle-style convergence is a complex task due to the non-convexity of the penalty function. To address this, we designed a multistage convex relaxation algorithm using a majorization-minimization (MM) approach that solves a series of convex subproblems. This ensures that the approximate local optimum shares the same optimal statistical properties

as the unattainable global optimum. Our findings also confirm that the estimator produced by the MM-based algorithm reaches the oracle statistical rate under minimal assumptions. Theoretical analyses are supported by numerical experiments, which demonstrate the method's effectiveness.

Detailed proofs of the results in this paper are available at: <https://www.ncvxopt.com/pubs/WangZhao-CovSensing.pdf>.

II. PROBLEM FORMULATION

We denote the matrix to be estimated as $\Sigma \in \mathbb{R}^{d \times d}$. In general, it is NP-hard to estimate the covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ from $m < \frac{d(d+1)}{2}$ measurements, unless appropriate structural assumptions, such as sparsity, are assumed to be known as a prior. In this paper, we concentrate on the sparse covariance matrices estimation. Additionally, we outline some standard assumptions on sensing vectors and noise in Assumption 1.

Assumption 1. We operate under the assumption that $\{\eta_i\}_{i=1}^m$ and $\{\mathbf{a}_i\}_{i=1}^m$ consist of i.i.d sub-Gaussian random variables. Specifically, η_i ($1 \leq i \leq m$) is drawn from a sub-Gaussian distribution characterized by a mean of 0 and a variance proxy σ^2 . Similarly, \mathbf{a}_i 's ($1 \leq i \leq m$) are considered i.i.d copies of $\mathbf{z} = [z_1, \dots, z_d]^\top$, where each z_i is i.i.d drawn from a distribution satisfying

$$\mathbb{E}(z_i) = 0, \mathbb{E}(z_i^2) = 1, \text{ and } \mathbb{E}(z_i^4) > 1. \quad (3)$$

We propose to estimate the sparse covariance matrices from quadratic measurements using the non-convex penalty. The idea is simple yet powerful. For notational simplicity, let $\mathbf{A}_i := \mathbf{a}_i \mathbf{a}_i^\top$ represent the equivalent sensing matrix, $\mathbf{Y} := [y_1, \dots, y_m]^\top$ express the equivalent measurements vector, and define the linear operator $\mathcal{A}(\Sigma) : \mathbb{R}^{d \times d} \mapsto \mathbb{R}^m$ that maps a matrix $\Sigma \in \mathbb{R}^{d \times d}$ to $\{\langle \Sigma, \mathbf{A}_i \rangle\}_{i=1}^m$. Thus, the measures \mathbf{Y} obeys $y_i := \langle \mathbf{A}_i, \Sigma \rangle + \eta_i$ and we can express the estimator as

$$\min_{\Sigma \succ 0} \left\{ \frac{1}{2m} \|\mathbf{Y} - \mathcal{A}(\Sigma)\|_2^2 - \tau \log \det \Sigma + \sum_{j,k} p_\lambda(|\Sigma_{jk}|) \right\}. \quad (4)$$

In (4), the first term aims to minimize the least squares errors. The second term, a log-determinant barrier function, ensures positive definiteness with a barrier parameter $\tau \geq 0$. The third term p_λ represents a non-convex penalty function governed by a regularization parameter $\lambda > 0$. We impose certain restrictions on it, as illustrated in Assumption 2.

Assumption 2. The function $p_\lambda(t)$ defined on $[0, +\infty)$ satisfies:

- $p_\lambda(t)$ is non-decreasing on $[0, +\infty)$ with $p_\lambda(0) = 0$ and is differentiable almost everywhere on $(0, +\infty)$;
- $0 \leq p'_\lambda(t_1) \leq p'_\lambda(t_2) \leq \lambda$ for all $t_1 \geq t_2 \geq 0$ and $\lim_{t \rightarrow 0} p'_\lambda(t) = \lambda$;
- There exists an $\alpha > 0$ such that $p'_\lambda(t) = 0$ for $t \geq \alpha\lambda$;
- There exists some $c \in (0, \alpha)$ such that $p'_\lambda(c\lambda) \geq \frac{\lambda}{2}$.

Under Assumption 2, the initial trio of conditions guarantee sparsity and unbiasedness. The final criterion is invariably satisfied on account of $p'_\lambda(0) = \lambda$ and $p'_\lambda(\alpha\lambda) = 0$ which is listed for the convenience of subsequent theoretical analysis. Typical examples of $p_\lambda(\cdot)$ in Assumption (2) include SCAD [19] and MCP [19].

III. OPTIMIZATION ALGORITHM

We now construct an MM-based algorithm to solve (4). The multistage convex relaxation is essentially a sequential optimization framework [25]. At each stage, we take weighted ℓ_1 -norm as surrogate function of $\sum_{j,k} p_\lambda(|\Sigma_{jk}|)$. For notational simplicity, we define $f(\Sigma) = \frac{1}{2m} \|\mathbf{Y} - \mathcal{A}(\Sigma)\|_2^2 - \tau \log \det \Sigma$. Consequently, we consider to solve a sequence of convex relaxation subproblems as follows:

$$\min_{\Sigma \succ 0} \left\{ f(\Sigma) + \sum p'_\lambda \left(\left| \widehat{\Sigma}_{jk}^{(q-1)} \right| \right) |\Sigma_{jk}| \right\}, \quad 1 \leq q \leq Q, \quad (5)$$

where $\widehat{\Sigma}^{(q)}$ is the optimal solution to the q -th subproblem.

Each subproblem corresponds to a weighted ℓ_1 -penalized covariance estimation problem, which can generally be reformulated as follows:

$$\min_{\Sigma \succ 0} \{ f(\Sigma) + \|\Lambda \odot \Sigma\|_1 \}, \quad (6)$$

where Λ is the regularized parameter matrix with $\Lambda_{jk} = p'_\lambda \left(\left| \widehat{\Sigma}_{jk} \right| \right) \in [0, \lambda]$. According to the Karush-Kuhn-Tucker (KKT) conditions, the unique sparse global optimal solution $\widehat{\Sigma}$ for (6) satisfies the first-order optimal condition:

$$\nabla f(\widehat{\Sigma}) + \Lambda \odot \widehat{\Sigma} = 0, \quad \text{with } \widehat{\Sigma} \in \partial \left\| \widehat{\Sigma} \right\|_1,$$

in which $\nabla f(\Sigma) = -\frac{1}{m} \mathcal{A}^*(\mathbf{Y} - \mathcal{A}(\Sigma)) - \tau \Sigma^{-1}$. Obviously, it is impossible to find an exact solution of problem (6) which has no analytical solution. In practice, we settle for the second-best solution defined in 3 to (6), and terminate the iterations when the approximate KKT condition holds.

Definition 3. Given a tolerance level ε , $\widetilde{\Sigma}^{(q)}$ qualifies as an ε -optimal solution to (6) if $\omega_{\Lambda^{(q-1)}} \left(\widetilde{\Sigma}^{(q)} \right) \leq \varepsilon$, where

$$\omega_{\Lambda^{(q-1)}} \left(\widetilde{\Sigma} \right) = \min_{\Xi \in \partial \left\| \widetilde{\Sigma} \right\|_1} \left\| \nabla f(\widetilde{\Sigma}) + \Lambda^{(q-1)} \odot \Xi \right\|_{\max},$$

where $\Lambda_{jk}^{(q-1)} = p'_\lambda \left(\left| \widehat{\Sigma}_{jk}^{(q-1)} \right| \right)$.

The MM algorithm is described in Algorithm 1, where we set $\Sigma^{\{0\}} = \mathbf{I}$ as a trivial start. To obtain $\widetilde{\Sigma}^{(q)}$ by solving (6), we applying a proximal Newton method with backtracking line search [26], [27]. We start with $\Sigma_0^{(q)} = \widetilde{\Sigma}^{(q-1)}$ and stop the proximal Newton algorithm when $\omega_{\Lambda^{(q-1)}} \left(\widetilde{\Sigma} \right)$ researches the prefixed optimization error ε , which is given in Definition on this page.

Algorithm 1: The MM-Based Multistage Convex Relaxation Algorithm for Solving (4).

Input: $\mathcal{S}, \tau, \lambda$;
1 Initialize $\widetilde{\Sigma}^{(0)} = \mathbf{I}$
2 for $q = 1, 2, \dots, Q$ **do**
3 $\Lambda_{jk}^{(q-1)} = p'_\lambda \left(\left| \widetilde{\Sigma}_{jk}^{(q-1)} \right| \right)$;
4 obtain $\widetilde{\Sigma}^{(q)}$ by solving (6);
5 $q = q + 1$;
6 end
Output: $\widetilde{\Sigma}^{(Q)}$

IV. STATISTICAL THEORIES

In this section, we first introduce some important definitions and technical assumptions. Following this foundational setup, we establish the statistical convergence rate of our proposed covariance estimator.

A. Assumptions

Definition 4 (Sparse Eigenvalue). Given any positive integer s , we define the largest and smallest s -sparse eigenvalues of $\nabla^2 f(\Sigma)$ as

$$\rho_s^+ = \sup \left\{ \mathbf{v}^\top \nabla^2 f(\Sigma) \mathbf{v} \mid \|\mathbf{v}\|_2^2 = 1, \|\mathbf{v}\|_0 \leq s \right\};$$

$$\rho_s^- = \inf \left\{ \mathbf{v}^\top \nabla^2 f(\Sigma) \mathbf{v} \mid \|\mathbf{v}\|_2^2 = 1, \|\mathbf{v}\|_0 \leq s \right\}.$$

Moreover, we define $\kappa_s = \rho_s^+ / \rho_s^-$ as the s -sparse condition number.

The sparse eigenvalue condition is widely studied in the high-dimensional statistics. And it is closely related to the restricted strongly convex, restricted smoothness and restricted eigenvalues properties [28]. Such condition have been employed by [29]–[31].

Definition 5. Define a local region of Σ^* by

$$\mathcal{B}(\Sigma^*, r) = \{ \Sigma \succ 0 \mid \|\Sigma - \Sigma^*\|_F \leq r \}.$$

In the following, we will show that $f(\Sigma)$ satisfies local strong convexity and local strong smoothness over a sparse domain. And in our analysis, we set the radius r as $\frac{\rho_{2s^*+2\widetilde{s}}^-}{4\tau\kappa}$.

Assumption 6. Given $\Sigma \in \mathcal{B}(\Sigma^*, r)$ for a generic constant r , there exists a generic constant C_1 such that $\nabla^2 f(\Sigma)$ satisfies the sparse eigenvalue properties with parameters $\rho_{2s^*+2\widetilde{s}}^-$ and $\rho_{2s^*+2\widetilde{s}}^+$ satisfying

$$0 < \rho_{2s^*+2\widetilde{s}}^- < \rho_{2s^*+2\widetilde{s}}^+ < +\infty$$

with $\widetilde{s} \geq C_1 \kappa_{2s^*+2\widetilde{s}}^2 s^*$ and $\kappa_{2s^*+2\widetilde{s}} = \rho_{2s^*+2\widetilde{s}}^+ / \rho_{2s^*+2\widetilde{s}}^-$.

Assumption 6 shows that $\nabla^2 f(\Sigma)$ has bounded largest and non-zero smallest sparse eigenvalues, under the condition that Σ is adequately sparse and proximate to Σ^* . This employment of similar conditions has become a common practice

in tackling high-dimensional problems [32]–[35]. Under our design, we confirm the existence of sparse eigenvalue. We refer readers to [36]–[42] for further details.

Assumption 7 (Local Restricted Hessian Smoothness). *Referring \tilde{s} introduced in Assumption 6. Constant $L_{2s^*+2\tilde{s}}$ and r exist such that for any $\Sigma, \Sigma' \in \mathcal{B}(\Sigma^*, r)$ with $\|\Sigma_{\tilde{S}}\|_0 \leq \tilde{s}$ and $\|\Sigma'_{\tilde{S}}\|_0 \leq \tilde{s}$, the following inequality holds:*

$$\sup_{\mathbf{u} \in \Delta, \|\mathbf{u}\|_F=1} \|\mathbf{u}\|_{\nabla^2 f(\Sigma') - \nabla^2 f(\Sigma)}^2 \leq L_{2s^*+2\tilde{s}} \|\Sigma - \Sigma'\|_F^2,$$

where $\Delta = \{\mathbf{u} \mid \text{supp}(\mathbf{u}) \subseteq (\text{supp}(\Sigma) \cup \text{supp}(\Sigma'))\}$ and

$$\begin{aligned} & \|\mathbf{u}\|_{\nabla^2 f(\Sigma') - \nabla^2 f(\Sigma)}^2 \\ &= \text{vec}^\top(\mathbf{u}) (\nabla^2 f(\Sigma') - \nabla^2 f(\Sigma)) \text{vec}(\mathbf{u}). \end{aligned}$$

Assumption 7 ensures that $\nabla^2 f(\Sigma)$ exhibits Lipschitz continuity in the vicinity of Σ^* across a sparse domain.

In the following, we specify some mild conditions for the true covariance matrix. We represent the true covariance matrix as Σ^* , and identify the support set of Σ^* as $\mathcal{S}^* = \{(j, k) \mid \Sigma_{jk}^* \neq 0\}$, with s^* denoting its size, that is, $s^* = |\mathcal{S}^*|$.

Assumption 8. *Given the true covariance matrix Σ^* , there exist generic constants α, c such that*

$$\|\Sigma_{\mathcal{S}^*}^*\|_{\min} = \min_{(j,k) \in \mathcal{S}^*} |\Sigma_{jk}^*| \geq (\alpha + c) \lambda \geq \lambda,$$

where α and c being constants introduced in Assumption (2).

Assumption 8 guarantees the minimum signal strength condition, which is frequently utilized in the study of non-convex penalized regression problems [19], [20], [43]. In our design case, the tuning parameter λ can be taken to be the order of $\sqrt{\frac{\log d}{m}}$ that could not be too large when the measurements m increases, which guarantees that the estimators are close enough to the true model parameter. Taking the signal strength into account, we can obtain the oracle rate of convergence.

Assumption 9. *For the true covariance matrix Σ^* , there exists some constant $\kappa \geq 1$ such that*

$$0 < \frac{1}{\kappa} \leq \lambda_{\min}(\Sigma^*) \leq \lambda_{\max}(\Sigma^*) \leq \kappa < \infty.$$

The condition regarding the smallest and largest eigenvalues of the true covariance matrix in Assumption 9 is standard within the existing research on the sparse covariance matrix estimation problem [44].

Assumption 10. *At every stage of addressing the convex relaxed subproblem for all $q \geq 1$, we set*

$$\varepsilon = \frac{C_3}{\sqrt{m}} \leq \frac{\lambda}{8},$$

where C_3 is a predetermined small constant.

Assumption 10 ensures that the solution $\tilde{\Sigma}^{(q)}$ obtained at each stage, for all $q \geq 1$, achieves adequate precision. This level of accuracy is essential for the convergence analysis of multistage convex relaxation.

B. Statistical Guarantees and Consequences

Now we introduce the primary result, illustrating the contraction property of the solution path $\{\tilde{\Sigma}^{(q)}\}_{q \geq 1}$. For simplicity, we denote the following symbols. For functionals $f(n)$ and $g(n)$, $f(n) \gtrsim g(n)$ if $f(n) \geq cg(n)$, $f(n) \lesssim g(n)$ if $f(n) \leq Cg(n)$, and $f(n) \asymp g(n)$ if $cg(n) \leq f(n) \leq Cg(n)$ for some constant c and C . $\mathcal{O}_p(\cdot)$ is used to denote bounded in probability.

Theorem 11 (Contraction Property). *Under the Assumptions 1, 6, 8 and 9, the ε -optimal solution $\tilde{\Sigma}^{(q)}$ adheres to the following bound:*

$$\begin{aligned} \|\tilde{\Sigma}^{(q)} - \Sigma^*\|_F &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \underbrace{\|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F}_{\text{oracle rate}} \\ &\quad + \frac{1}{\rho_{2s^*+2\tilde{s}}} \underbrace{\varepsilon \sqrt{s^*}}_{\text{optimization error}} + \delta \underbrace{\|\tilde{\Sigma}^{(q-1)} - \Sigma^*\|_F}_{\text{contraction}}, \end{aligned} \quad (7)$$

for $1 \leq q \leq Q$, where $\delta \in (0, 1)$ is the contraction factor.

Remark 12. The oracle rate describes the statistical convergence rate for an oracle estimator, which has prior knowledge the true support set \mathcal{S}^* . The oracle estimator $\hat{\Sigma}^O$ is defined as

$$\hat{\Sigma}^O = \arg \min_{\Sigma: \Sigma_{\mathcal{S}^*} = 0} f(\Sigma).$$

According to the definition, it follows quite straightforwardly that $\hat{\Sigma}^O$ fulfills the condition $\|\hat{\Sigma}^O - \Sigma^*\|_F \lesssim \|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F$, indicating a close approximation to Σ^* within the confines of the true support set's gradient norms.

Theorem 11 reveals the estimation discrepancy between the ε -optimal solution $\tilde{\Sigma}^{(q)}$ and the actual parameter Σ^* is bounded by three components: the oracle rate, the optimization error, and a contraction term. Following this, we detail the precise statistical convergence rate applicable within our framework.

Corollary 13. *Suppose that Assumptions 1, 6, 8, 9 and 10 hold. If $\lambda \asymp \sqrt{\frac{\log d}{m}}$, $\tau \lesssim \sqrt{\frac{1}{m}} \|(\Sigma^*)^{-1}\|_{\max}^{-1}$, the ε -optimal solution $\tilde{\Sigma}^{(1)}$ satisfies*

$$\|\tilde{\Sigma}^{(1)} - \Sigma^*\|_F \lesssim \sqrt{\frac{s^* \log d}{m}}$$

with high probability.

Corollary 13 naturally follows from Theorem 11 when $q = 1$. Furthermore, the contraction property, as discussed in Theorem 11, is induced by the MM-based multistage convex relaxation algorithm. To attain the oracle rate, the optimization error must be strategically selected so that $\varepsilon \leq \frac{\|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F}{\sqrt{s^*}}$ and the parameter Q should be sufficiently large. This leads us to the subsequent conclusion.

Corollary 14. Suppose that Assumptions 1, 6, 7, 8, 9, and 10 hold. If $\lambda \asymp \sqrt{\frac{\log d}{m}}$, $\tau \lesssim \sqrt{\frac{1}{m}} \|(\Sigma^*)^{-1}\|_{\max}^{-1}$, and $Q \gtrsim \log(\lambda\sqrt{m}) \gtrsim \log \log d$, then the ε -optimal solution $\tilde{\Sigma}^{(Q)}$ satisfies

$$\|\tilde{\Sigma}^{(Q)} - \Sigma^*\|_F = \mathcal{O}_p\left(\sqrt{\frac{s^*}{m}}\right).$$

Corollary 14 is a direct consequence of Theorem 11, which suggests that under minimal assumptions, solving no more than approximately $\log \log d$ convex problems is sufficient to reach the oracle rate $\sqrt{\frac{s^*}{m}}$.

V. NUMERICAL EXPERIMENTS

We compare our proposal with the existing estimator, focusing specifically on the estimator that incorporates an ℓ_1 penalty. We conduct a series of Monte Carlo simulations utilizing matrices of dimensions 50×50 . The tuning parameters λ and τ were selected via five-fold cross-validation. Initially, positive symmetric sparse covariance matrices are generated utilizing the “sprandsym” function in MATLAB, defined by a sparsity value k , which represents the quantity of non-zero elements within the matrix. The non-convex penalty function chosen for these experiments is the MCP, with a constant setting of $b = 2$ across all trials.

For each scenario characterized by a specific (m, k) pair, we execute 10 independent iterations. The outcomes are depicted in Figure 1, which illustrates the mean relative error of the estimated covariance matrices as a function of m across varying levels of sparsity. Each measurement includes additive noise, achieved by generating η_i from a normal distribution $\sigma \cdot \mathcal{N}(0, 1)$ with $\sigma = 10^{-4}$. From the simulation results, it is evident that the estimator with MCP performs better than the ℓ_1 penalty-based estimator in terms of Frobenius norm error. This supports our theoretical findings that non-convex penalty can reduce the covariance estimation error.

Fig. 2 illustrates the performance of oracle rate in estimating a covariance matrix from quadratic measurement model. The y -axis represents the deviation between an estimated covariance matrix and the true covariance matrix, measured using the Frobenius norm. The x -axis is a scaled measure related to the support size of true covariance matrix and the measurement complexity, indicated by $\sqrt{\frac{s^*}{m}}$. Two lines are plotted to show the estimation error for two different dimensions ($d = \{50, 60\}$) under the sparsity level of $k = 80$. As the measurement size increase, the error between the estimated and true covariance matrix decreases. The black line ($d = 60$) shows a slightly higher error trend than the blue line ($d = 50$), suggesting that higher-dimensional covariance tends to have larger deviations from the true matrix. This visual representation provides insight into how the quality of covariance matrix estimation varies with the oracle rate.

VI. CONCLUSION AND DISCUSSION

This paper introduces a new sample strategy named covariance matrix sensing, including sensing and compressing

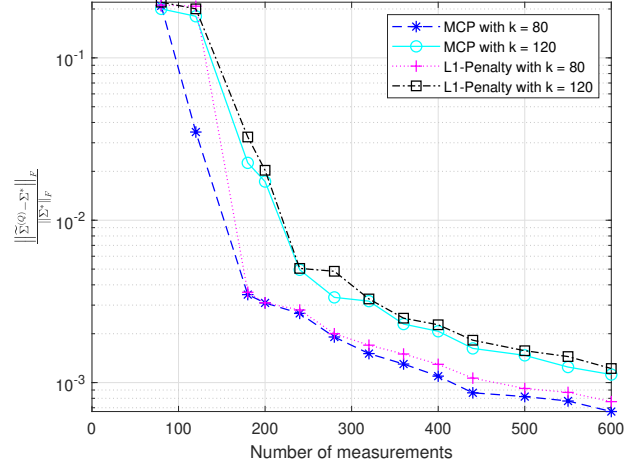


Fig. 1. The mean relative error of the estimated covariance matrices for different sparsity level with noise.

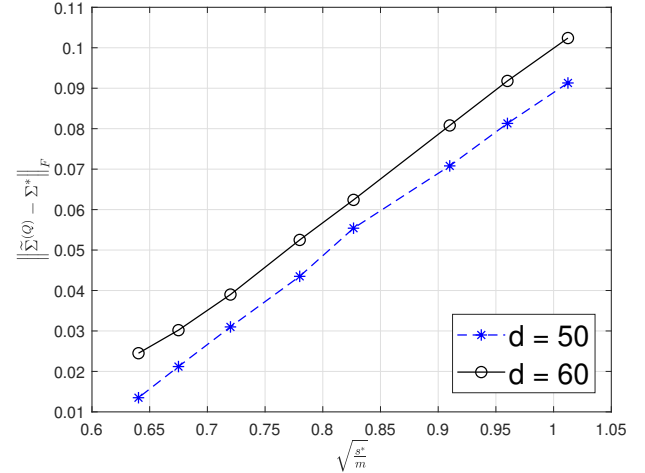


Fig. 2. The oracle rate for different measurements and different dimension of Σ^* .

signals and estimating the covariance matrix from a small number of measurements. Specifically, we estimate large sparse covariance matrices from quadratic samplings through a non-convex penalty, detailing both its theoretical results and empirical validations. We demonstrate that our estimators achieve superior statistical convergence rates when contrasted with current methodologies.

And all theoretical derivations and the suggested algorithm can be seamlessly applied to the bilinear measurement model where $y_i = \mathbf{a}_i^\top \Sigma \mathbf{b}_i$, with \mathbf{a}_i and \mathbf{b}_i being independently produced sensing vectors. This demonstrates that our outcomes are equally applicable to this asymmetric sensing model too.

REFERENCES

- [1] M. Abt and W. J. Welch, "Fisher information and maximum-likelihood estimation of covariance parameters in Gaussian stochastic processes," *Canadian Journal of Statistics*, vol. 26, no. 1, pp. 127–137, 1998.
- [2] M. Pourahmadi, *High-dimensional covariance estimation: with high-dimensional data*. John Wiley & Sons, 2013, vol. 882.
- [3] Z. Bao, X. Ding, J. Wang, and K. Wang, "Statistical inference for principal components of spiked covariance matrices," *The Annals of Statistics*, vol. 50, no. 2, pp. 1144–1169, 2022.
- [4] G. Dasarathy, P. Shah, B. N. Bhaskar, and R. Nowak, "Covariance sketching," in *Proceedings of 2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2012, pp. 1026–1033.
- [5] R. C. Daniels and R. W. Heath, "60 GHz wireless communications: Emerging requirements and design recommendations," *IEEE Vehicular Technology Magazine*, vol. 2, no. 3, pp. 41–50, 2007.
- [6] M. Raymer, M. Beck, and D. McAlister, "Complex wave-field reconstruction using phase-space tomography," *Physical Review Letters*, vol. 72, no. 8, p. 1137, 1994.
- [7] M. Behrendt, "Uncertainty modelling in power spectrum estimation of environmental processes," 2022.
- [8] G. Leus and Z. Tian, "Recovering second-order statistics from compressive measurements," in *Proceedings of 2011 4th IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP)*. IEEE, 2011, pp. 337–340.
- [9] D. D. Ariananda and G. Leus, "Compressive wideband power spectrum estimation," *IEEE Transactions on signal processing*, vol. 60, no. 9, pp. 4775–4789, 2012.
- [10] —, "Compressive Wideband Power Spectrum Estimation," *IEEE Transactions on Signal Processing*, vol. 60, no. 9, pp. 4775–4789, 2012.
- [11] J. A. Tropp, J. N. Laska, M. F. Duarte, J. K. Romberg, and R. G. Baraniuk, "Beyond nyquist: Efficient sampling of sparse bandlimited signals," *IEEE Transactions on Information Theory*, vol. 56, no. 1, pp. 520–544, 2009.
- [12] S. Muthukrishnan *et al.*, "Data streams: Algorithms and applications," *Foundations and Trends® in Theoretical Computer Science*, vol. 1, no. 2, pp. 117–236, 2005.
- [13] G. Dasarathy, P. Shah, B. N. Bhaskar, and R. D. Nowak, "Sketching sparse matrices, covariances, and graphs via tensor products," *IEEE Transactions on Information Theory*, vol. 61, no. 3, pp. 1373–1388, 2015.
- [14] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [15] J. Wright and Y. Ma, *High-dimensional data analysis with low-dimensional models: Principles, computation, and applications*. Cambridge University Press, 2022.
- [16] W. B. Wu and M. Pourahmadi, "Nonparametric estimation of large covariance matrices of longitudinal data," *Biometrika*, vol. 90, no. 4, pp. 831–844, 2003.
- [17] Z. Tian, Y. Tefesse, and B. M. Sadler, "Cyclic feature detection with sub-Nyquist sampling for wideband spectrum sensing," *IEEE Journal of Selected Topics in Signal Processing*, vol. 6, no. 1, pp. 58–69, 2011.
- [18] Y. Chen, Y. Chi, and A. J. Goldsmith, "Exact and Stable Covariance Estimation From Quadratic Sampling via Convex Programming," *IEEE Transactions on Information Theory*, vol. 61, no. 7, pp. 4034–4059, 2015.
- [19] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [20] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *Annals of Statistics*, vol. 38, pp. 894–942, 2010.
- [21] Q. Wei and Z. Zhao, "Large Covariance Matrix Estimation With Oracle Statistical Rate via Majorization-Minimization," *IEEE Transactions on Signal Processing*, vol. 71, pp. 3328–3342, 2023.
- [22] H. Gui, J. Han, and Q. Gu, "Towards faster rates and oracle property for low-rank matrix estimation," in *Proceedings of International Conference on Machine Learning*. PMLR, 2016, pp. 2300–2309.
- [23] Q. Sun, K. M. Tan, H. Liu, and T. Zhang, "Graphical nonconvex optimization via an adaptive convex relaxation," in *Proceedings of International Conference on Machine Learning*. PMLR, 2018, pp. 4810–4817.
- [24] R. Varma, H. Lee, J. Kovačević, and Y. Chi, "Vector-valued graph trend filtering with non-convex penalties," *IEEE Transactions on Signal and Information Processing over Networks*, vol. 6, pp. 48–62, 2019.
- [25] T. Zhang, "Analysis of multi-stage convex relaxation for sparse regularization," *Journal of Machine Learning Research*, vol. 11, no. 3, 2010.
- [26] X. Li, L. Yang, J. Ge, J. Haupt, T. Zhang, and T. Zhao, "On quadratic convergence of DC proximal Newton algorithm in nonconvex sparse learning," *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [27] C.-J. Hsieh, M. A. Sustik, I. S. Dhillon, P. Ravikumar *et al.*, "QUIC: quadratic approximation for sparse inverse covariance estimation," *Journal of Machine Learning Research*, vol. 15, no. 83, pp. 2911–2947, 2014.
- [28] S. Zhou, "Restricted eigenvalue conditions on subgaussian random matrices," *arXiv preprint arXiv:0912.4045*, 2009.
- [29] P. J. Bickel, Y. Ritov, and A. B. Tsybakov, "Simultaneous analysis of Lasso and Dantzig selector," *The Annals of Statistics*, vol. 37, no. 4, pp. 1705 – 1732, 2009.
- [30] P.-L. Loh and M. J. Wainwright, "Regularized m-estimators with non-convexity: statistical and algorithmic theory for local optima," *Journal of Machine Learning Research*, vol. 16, no. 1, pp. 559–616, jan 2015.
- [31] S. N. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu, "A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers," *Statistical Science*, vol. 27, no. 4, pp. 538 – 557, 2012.
- [32] L. Xiao and T. Zhang, "A proximal-gradient homotopy method for the sparse least-squares problem," *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1062–1091, 2013.
- [33] Z. Wang, H. Liu, and T. Zhang, "Optimal computational and statistical rates of convergence for sparse nonconvex learning problems," *Annals of Statistics*, vol. 42, no. 6, p. 2164, 2014.
- [34] X. Li, J. Haupt, R. Arora, H. Liu, M. Hong, and T. Zhao, "A first order free lunch for sqrt-lasso," *arXiv preprint arXiv:1605.07950*, 2016.
- [35] X. Li, T. Zhao, R. Arora, H. Liu, and J. Haupt, "Stochastic variance reduced optimization for nonconvex sparse learning," in *Proceedings of International Conference on Machine Learning*. PMLR, 2016, pp. 917–925.
- [36] V. Sivakumar, A. Banerjee, and P. Ravikumar, "Beyond sub-gaussian measurements: High-dimensional structured estimation with sub-exponential designs," *Advances in Neural Information Processing Systems*, vol. 28, pp. 2197–2205, 2015.
- [37] V. Koltchinskii and S. Mendelson, "Bounding the Smallest Singular Value of a Random Matrix Without Concentration," *International Mathematics Research Notices*, vol. 2015, no. 23, pp. 12 991–13 008, 03 2015.
- [38] R. I. Oliveira, "The lower tail of random quadratic forms with applications to ordinary least squares," *Probability Theory and Related Fields*, vol. 166, no. 3, pp. 1175–1194, Dec 2016.
- [39] S. Mendelson and G. Paouris, "On generic chaining and the smallest singular value of random matrices with heavy tails," *Journal of Functional Analysis*, vol. 262, no. 9, pp. 3775–3811, 2012.
- [40] A. K. Kuchibhotla and A. Chakraborty, "Moving beyond sub-Gaussianity in high-dimensional statistics: Applications in covariance estimation and linear regression," *Information and Inference: A Journal of the IMA*, vol. 11, no. 4, pp. 1389–1456, 2022.
- [41] M. Genzel and C. Kipp, "Generic error bounds for the generalized lasso with sub-exponential data," *Sampling Theory, Signal Processing, and Data Analysis*, vol. 20, no. 2, p. 15, 2022.
- [42] S. van de Geer and A. Muro, "On higher order isotropy conditions and lower bounds for sparse quadratic forms," *Electronic Journal of Statistics*, vol. 8, no. 2, pp. 3031 – 3061, 2014.
- [43] J. Fan, H. Liu, Q. Sun, and T. Zhang, "I-LAMM for sparse learning: Simultaneous control of algorithmic complexity and statistical error," *Annals of Statistics*, vol. 46, no. 2, p. 814, 2018.
- [44] A. J. Rothman, "Positive definite estimators of large covariance matrices," *Biometrika*, vol. 99, no. 3, pp. 733–740, 2012.
- [45] P. Bühlmann and S. Van De Geer, *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.

APPENDIX A
PROOF OF STATISTICAL THEORY

In this appendix, we first provide essential lemmata and then provide the proofs of all the statistical theoretical results in Section IV-B. In Subsection A-A, we begin by the fact that the existence of SE within our design framework. Subsequently, we give the definition of Localized Restricted Eigenvalues (LRE), and then demonstrate the relationship between SE and LRE. Followed by a proposition that connects the SE condition to the localized version of the sparse strong convexity/sparse strong smoothness. Then we introduce Lemma 20, which sets bounds on the estimation error for the general problem ??, followed by Lemma 21 that outlines the estimation error limited for approximate solutions derived using the MM-based algorithm. Subsection A-B elaborates on Theorem 11, leveraging Lemma 20 and 21 to illustrate the solution path's contraction characteristic. Subsection A-C focuses on key concentration inequalities, which are vital for deducing the explicit statistical convergence rate. In Subsections A-D and A-E, the document presents the proofs for Corollary 13 and Corollary 14, respectively.

A. Technical Lemmata

Lemma 15. Let $\mathbf{X} \in \mathbb{R}^{n \times p}$ have independent isotropic sub-exponential rows. Let $\hat{\mathbf{A}} \subseteq S^{p-1}$, $0 < \xi < 1$, and c is a constant that depends on the sub-exponential norm $K = \sup_{\mathbf{u} \in \hat{\mathbf{A}}} \|\langle \mathbf{X}, \mathbf{u} \rangle\|_{\psi_1}$. Let $\omega_e(\hat{\mathbf{A}})$ denote the exponential width of the set. Then for some $\nu > 0$ with probability at least $(1 - \exp(-\nu^2/2))$,

$$\inf_{\mathbf{u} \in \hat{\mathbf{A}}} \|\mathbf{X}\mathbf{u}\|_2 \geq c\xi (1 - \xi^2)^2 \sqrt{n} - 4\omega_e(\hat{\mathbf{A}}) - \xi\nu.$$

And if $\{V_j\}$ is \mathcal{F}_0 -measurable we get for all $\|\mathbf{u}\|_2 = 1$,

$$\sup \|\mathbf{X}\mathbf{u}\|_2 \leq 2^{1-1/h_0} \sqrt{2h\mu_h} + 2^{1-1/h_0} hK,$$

where the definition of related letters involved is consistent with [36], [42].

Remark 16. Consider the matrix $\tilde{\mathbf{A}}$ constructed as follows:

$$\tilde{\mathbf{A}} = \begin{bmatrix} \text{vec}(\mathbf{a}_1 \mathbf{a}_1^\top) \\ \vdots \\ \text{vec}(\mathbf{a}_m \mathbf{a}_m^\top) \end{bmatrix},$$

where $\tilde{\mathbf{A}} \in \mathbb{R}^{m \times d^2}$ and its rows are independent isotropic sub-exponential vectors. Under the designed moment conditions, such a matrix $\tilde{\mathbf{A}}$ satisfies the sparse eigenvalue condition with high probability. Specifically, it is straightforward to derive a non-trivial asymptotic lower and upper bound for $\tilde{\mathbf{A}}$ with exponential probability. Moreover, the column of $\tilde{\mathbf{A}}$ have ℓ_2 norms in the order of \sqrt{m} , which holds with overwhelming probability when $\tilde{\mathbf{A}}$ is under our specified design. For more

details on the proof, we refer readers to [36]–[42] for a comprehensive review.

Definition 17. We denote the local ℓ_1 cone as

$$\mathcal{C}(s, \vartheta, r) = \{\mathbf{u}, \boldsymbol{\Sigma} \mid \mathcal{S}^* \subseteq \mathcal{J}, |\mathcal{J}| \leq s\} \cup \{\mathbf{u}, \boldsymbol{\Sigma} \mid \|\mathbf{u}_{\mathcal{J}^c}\|_1 \leq \vartheta \|\mathbf{u}_{\mathcal{J}}\|_1, \|\boldsymbol{\Sigma} - \boldsymbol{\Sigma}^*\|_F \leq r\}.$$

Then we define the largest and smallest LRE as

$$\psi_{s, \vartheta, r}^+ = \sup_{\mathbf{u}, \boldsymbol{\Sigma}} \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\boldsymbol{\Sigma}) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid (\mathbf{u}, \boldsymbol{\Sigma}) \in \mathcal{C}(s, \vartheta, r) \right\},$$

$$\psi_{s, \vartheta, r}^- = \inf_{\mathbf{u}, \boldsymbol{\Sigma}} \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\boldsymbol{\Sigma}) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid (\mathbf{u}, \boldsymbol{\Sigma}) \in \mathcal{C}(s, \vartheta, r) \right\}.$$

The following proposition demonstrates the relationships between SE and LRE. The proof can be found in [45], thus is omitted here.

Proposition 18. Given any $\boldsymbol{\Sigma}, \boldsymbol{\Sigma}' \in \mathcal{C}(s, \vartheta, r) \cap \mathcal{B}(\boldsymbol{\Sigma}^*, r)$, we have

$$c_1 \psi_{s, \vartheta, r}^- \leq \rho_s^- \leq c_2 \psi_{s, \vartheta, r}^-,$$

$$c_3 \psi_{s, \vartheta, r}^+ \leq \rho_s^+ \leq c_4 \psi_{s, \vartheta, r}^+,$$

where c_1, c_2, c_3 and c_4 are constant.

Lemma 19. Let $D_f(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) = f(\boldsymbol{\Sigma}_1) - f(\boldsymbol{\Sigma}_2) - \langle \nabla f(\boldsymbol{\Sigma}_2), \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2 \rangle$ and define the symmetrized Bregman divergence for the loss function $f(\cdot)$ as

$$D_f^s(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) = D_f(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) + D_f(\boldsymbol{\Sigma}_2, \boldsymbol{\Sigma}_1) = \langle \nabla f(\boldsymbol{\Sigma}_1) - \nabla f(\boldsymbol{\Sigma}_2), \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2 \rangle.$$

For any $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathcal{B}(\boldsymbol{\Sigma}^*, r)$ such that $\|(\boldsymbol{\Sigma}_1)_{\overline{\mathcal{S}^*}}\|_0 \leq \tilde{s}$ and $\|(\boldsymbol{\Sigma}_2)_{\overline{\mathcal{S}^*}}\|_0 \leq \tilde{s}$, we have

$$\frac{1}{2} \rho_{2s^*+2\tilde{s}}^- \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_2^2 \leq D_f(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) \leq \frac{1}{2} \rho_{2s^*+2\tilde{s}}^+ \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_2^2,$$

$$\rho_{2s^*+2\tilde{s}}^- \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_2^2 \leq D_f^s(\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2) \leq \rho_{2s^*+2\tilde{s}}^+ \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_2^2.$$

Proof: By the mean value theorem, there exists a $\theta \in [0, 1]$ such that

$$\tilde{\boldsymbol{\Sigma}} = \theta \boldsymbol{\Sigma}_1 + (1 - \theta) \boldsymbol{\Sigma}_2 \in \mathcal{B}(\boldsymbol{\Sigma}^*, r),$$

$$\|\tilde{\boldsymbol{\Sigma}}\|_0 \leq 2\tilde{s} \text{ and}$$

$$\langle \nabla f(\boldsymbol{\Sigma}_1) - \nabla f(\boldsymbol{\Sigma}_2), \boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2 \rangle = \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}_2\|_{\nabla^2 f(\tilde{\boldsymbol{\Sigma}})}^2.$$

By the definition of the sparse eigenvalue, we obtain the desired result. ■

Lemma 20. Suppose that Assumption 9 hold. Consider the general problem in (6). Assume that there exists a set \mathcal{E} such that

$$\mathcal{S}^* \subseteq \mathcal{E}, |\mathcal{E}| \leq 2s^*, \text{ and } \|\boldsymbol{\Lambda}_{\overline{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2}.$$

Then we have $\lambda \geq 2(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon)$, and the ε -optimal solution $\tilde{\Sigma}$ satisfies

$$\begin{aligned} \|\tilde{\Sigma} - \Sigma^*\|_F &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}^-} \left(\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \varepsilon\sqrt{|\mathcal{E}|} + \|\Lambda_{S^*}\|_F \right) \\ &\leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}^-} \lambda\sqrt{s^*}. \end{aligned}$$

Proof: We know that $\|\nabla f(\Sigma^*)\|_{\max} \leq \frac{\lambda}{4}$ and $\varepsilon \leq \frac{\lambda}{8}$ based on assumptions 8 and 10, hence $\lambda \geq 2(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon)$. Then we construct an intermediate estimator $\tilde{\Sigma}^* = \Sigma^* + \theta(\tilde{\Sigma} - \Sigma^*)$, where θ is taken such that $\|\tilde{\Sigma}^* - \Sigma^*\|_F = r$, if $\|\tilde{\Sigma} - \Sigma^*\|_F > r$; $\theta = 1$ otherwise. By the definition of $\tilde{\Sigma}^*$, we see that $\|\tilde{\Sigma}^* - \Sigma^*\|_F \leq r$. Using Lemma 26, we know that the approximate solution falls in the ℓ_1 -cone. From the construction of $\tilde{\Sigma}^*$, we know $\tilde{\Sigma}^* - \Sigma^* = \theta(\tilde{\Sigma} - \Sigma^*)$. Thus, we have

$$\|(\tilde{\Sigma}^* - \Sigma^*)_{\mathcal{E}}\|_1 \leq 5\|(\tilde{\Sigma}^* - \Sigma^*)_{\bar{\mathcal{E}}}\|_1.$$

Combining the inequality above with the assumption $|\mathcal{E}| \leq 2s^*$ results that $\tilde{\Sigma}^*$ falls in the local ℓ_1 -cone. Combing the Lemma 19, Definition 17 and Proposition 18, it implies the localized restricted strong convexity, i.e.

$$\rho_{2s^*+2\tilde{s}}^- \|\tilde{\Sigma}^* - \Sigma^*\|_F^2 \leq D_f^s(\tilde{\Sigma}^*, \Sigma^*). \quad (8)$$

We use Lemma 27 to bound the right hand side of the above inequality such as

$$\begin{aligned} D_f^s(\tilde{\Sigma}^*, \Sigma^*) &\leq \theta D_f^s(\tilde{\Sigma}, \Sigma^*) \\ &= \theta \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \tilde{\Sigma} - \Sigma^* \rangle. \end{aligned} \quad (9)$$

Plugging (9) back into (8) yields

$$\begin{aligned} \rho_{2s^*+2\tilde{s}}^- \|\tilde{\Sigma}^* - \Sigma^*\|_F^2 &\leq \theta \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \tilde{\Delta} \rangle \\ &= \theta \underbrace{\langle \nabla f(\tilde{\Sigma}) + \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \rangle}_{\text{I}} \\ &\quad - \theta \underbrace{\langle \nabla f(\Sigma^*), \tilde{\Delta} \rangle}_{\text{II}} - \theta \underbrace{\langle \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \rangle}_{\text{III}}, \end{aligned} \quad (10)$$

where $\tilde{\Delta} = \tilde{\Sigma} - \Sigma^*$, $\tilde{\Xi}^{(t+1)} \in \partial \|\Sigma^{(t+1)}\|_1$. Then we establish limits for terms I, II and III, respectively.

Define $\Omega = \nabla f(\tilde{\Sigma}) + \Lambda \odot \tilde{\Xi}$. For term I, partitioning the support of Ω and $\tilde{\Delta}$ into \mathcal{E} and its complement $\bar{\mathcal{E}}$, and then applying the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{I} &= \langle \Omega_{\mathcal{E}}, \tilde{\Delta}_{\mathcal{E}} \rangle + \langle \Omega_{\bar{\mathcal{E}}}, \tilde{\Delta}_{\bar{\mathcal{E}}} \rangle \\ &\leq \|\Omega_{\mathcal{E}}\|_F \|\tilde{\Delta}_{\mathcal{E}}\|_F + \|\Omega_{\bar{\mathcal{E}}}\|_{\max} \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1 \\ &\leq \sqrt{|\mathcal{E}|} \|\Omega_{\mathcal{E}}\|_{\max} \|\tilde{\Delta}_{\mathcal{E}}\|_F + \|\Omega\|_{\max} \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1 \\ &\stackrel{(i)}{\leq} \varepsilon\sqrt{|\mathcal{E}|} \|\tilde{\Delta}_{\mathcal{E}}\|_F + \varepsilon \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1, \end{aligned}$$

where (i) is from $\|\Omega\|_{\max} \leq \varepsilon$ by Definition 3. For term II, dividing the support of $\nabla f(\Sigma^*)$ and $\tilde{\Delta}$ into \mathcal{E} and $\bar{\mathcal{E}}$. Applying the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{II} &= \langle (\nabla f(\Sigma^*))_{\mathcal{E}}, \tilde{\Delta}_{\mathcal{E}} \rangle + \langle (\nabla f(\Sigma^*))_{\bar{\mathcal{E}}}, \tilde{\Delta}_{\bar{\mathcal{E}}} \rangle \\ &\geq -\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F \|\tilde{\Delta}_{\mathcal{E}}\|_F - \|(\nabla f(\Sigma^*))_{\bar{\mathcal{E}}}\|_{\max} \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1 \end{aligned}$$

For term III, dividing the support of $\Lambda \odot \tilde{\Xi}$ and $\tilde{\Delta}$ into S^* and \bar{S}^* , and then applying the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{III} &= \langle (\Lambda \odot \tilde{\Xi})_{S^*}, \tilde{\Delta}_{S^*} \rangle + \langle (\Lambda \odot \tilde{\Xi})_{\bar{S}^*}, \tilde{\Delta}_{\bar{S}^*} \rangle \\ &\stackrel{(i)}{=} \langle (\Lambda \odot \tilde{\Xi})_{S^*}, \tilde{\Delta}_{S^*} \rangle + \langle \Lambda_{\bar{S}^*}, |\tilde{\Delta}_{\bar{S}^*}| \rangle \\ &\geq -\|\Lambda_{S^*}\|_F \|\tilde{\Delta}_{S^*}\|_F + \langle \Lambda_{\bar{\mathcal{E}}}, |\tilde{\Delta}_{\bar{\mathcal{E}}}| \rangle \\ &\stackrel{(ii)}{\geq} -\|\Lambda_{S^*}\|_F \|\tilde{\Delta}_{\mathcal{E}}\|_F + \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1, \end{aligned}$$

where (i) is from

$$\langle (\Lambda \odot \tilde{\Xi})_{\bar{S}^*}, \tilde{\Delta}_{\bar{S}^*} \rangle = \langle \Lambda_{\bar{S}^*}, |\tilde{\Delta}_{\bar{S}^*}| \rangle = \langle \Lambda_{\bar{S}^*}, |\tilde{\Delta}_{\bar{S}^*}| \rangle,$$

and (ii) is from

$$\begin{aligned} \langle \Lambda_{\bar{\mathcal{E}}}, |\tilde{\Delta}_{\bar{\mathcal{E}}}| \rangle &= \sum_{(j,k) \in \bar{\mathcal{E}}} \Lambda_{jk} |\tilde{\Delta}_{jk}| \geq \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \sum_{(j,k) \in \bar{\mathcal{E}}} |\tilde{\Delta}_{jk}| \\ &= \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1, \end{aligned}$$

and

$$\|\tilde{\Delta}_{S^*}\|_F \leq \|\tilde{\Delta}_{\mathcal{E}}\|_F.$$

Combining the above results into 10 yields

$$\begin{aligned} &\rho_{2s^*+2\tilde{s}}^- \|\tilde{\Sigma}^* - \Sigma^*\|_F^2 \\ &\leq \left(\|\Lambda_{S^*}\|_F + \|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \varepsilon\sqrt{|\mathcal{E}|} \right) \times \theta \|\tilde{\Delta}_{\mathcal{E}}\|_F \\ &\quad - \theta \left(\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} - (\|(\nabla f(\Sigma^*))_{\bar{\mathcal{E}}}\|_{\max} + \varepsilon) \right) \|\tilde{\Delta}_{\bar{\mathcal{E}}}\|_1 \\ &\stackrel{(i)}{\leq} \left(\|\Lambda_{S^*}\|_F + \|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \varepsilon\sqrt{|\mathcal{E}|} \right) \\ &\quad \times \theta \left\| (\tilde{\Sigma} - \Sigma^*)_{\mathcal{E}} \right\|_F \\ &\leq \left((\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_{\max} + \varepsilon) \sqrt{|\mathcal{E}|} + \|\Lambda_{S^*}\|_{\max} \sqrt{|S^*|} \right) \\ &\quad \times \theta \left\| (\tilde{\Sigma} - \Sigma^*)_{\mathcal{E}} \right\|_F \\ &\leq \left((\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_{\max} + \varepsilon) \sqrt{2s^*} + \lambda\sqrt{s^*} \right) \\ &\quad \times \theta \left\| (\tilde{\Sigma} - \Sigma^*)_{\mathcal{E}} \right\|_F \\ &\leq \frac{2 + \sqrt{2}}{2} \lambda\sqrt{s^*} \times \theta \underbrace{\left\| (\tilde{\Sigma} - \Sigma^*)_{\mathcal{E}} \right\|_F}_{\text{IV}}, \end{aligned}$$

where (i) is due to the fact that $\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2} \geq \|\nabla f(\Sigma^*)\|_{\max} + \varepsilon$.

For **IV**, we have $\theta \left\| \left(\tilde{\Sigma} - \Sigma^* \right)_{\mathcal{E}} \right\|_{\text{F}} = \left\| \left(\tilde{\Sigma}^* - \Sigma^* \right)_{\mathcal{E}} \right\|_{\text{F}}$. Thus, we obtain

$$\rho_{2s^*+2\tilde{s}} \left\| \tilde{\Sigma}^* - \Sigma^* \right\|_{\text{F}} \leq \frac{2+\sqrt{2}}{2} \lambda \sqrt{s^*},$$

which is a contraction with the construction of $\tilde{\Sigma}^*$. This indicates that $\tilde{\Sigma}^* = \tilde{\Sigma}$. Therefore, the desired bound hold for $\tilde{\Sigma}$. ■

Lemma 21. Suppose that Assumptions 1 and 9 hold. Define the set $\mathcal{E}^{(q)}$ by

$$\mathcal{E}^{(q)} = \mathcal{S}^* \cup \mathcal{S}^{(q)}, \text{ with } \mathcal{S}^{(q)} = \left\{ (j, k) \mid \Lambda_{jk}^{(q-1)} \leq p'_\lambda(u) \right\},$$

where $u = c\lambda$ and $c = \frac{2+\sqrt{2}}{2\rho_{2s^*+2\tilde{s}}}$ is the same given in the Assumption 2. Then we have $\lambda \geq 2(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon)$, and for $q \geq 1$, $|\mathcal{E}^{(q)}| \leq 2s^*$, $\left\| \Lambda_{\mathcal{E}^{(q)}}^{(q-1)} \right\|_{\min} \geq \frac{\lambda}{2}$, and

$$\begin{aligned} \left\| \tilde{\Sigma}^{(q)} - \Sigma^* \right\|_{\text{F}} &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\left\| (\nabla f(\Sigma^*))_{\mathcal{E}^{(q)}} \right\|_{\text{F}} + \varepsilon \sqrt{|\mathcal{E}^{(q)}|} \right) \\ &\quad + \frac{1}{\rho_{2s^*+2\tilde{s}}} \left\| \Lambda_{\mathcal{S}^*}^{(q-1)} \right\|_{\text{F}} \\ &\leq \frac{2+\sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}. \end{aligned}$$

Proof: We first prove $|\mathcal{E}^{(q)}| \leq 2s^*$ holds by induction. For $q = 1$, we have $\Lambda_{jk}^{(0)} = \lambda \geq p'_\lambda(u)$ for $j \neq k$ and thus $|\mathcal{S}^{(1)}| \leq s^*$ and $\mathcal{E}^{(1)} = \mathcal{S}^* \cup \mathcal{S}^{(1)}$, which implies $|\mathcal{E}^{(1)}| \leq 2s^*$ holds. Assume $|\mathcal{E}^{(q)}| \leq 2s^*$ holds at $q-1$, i.e., $|\mathcal{E}^{(q-1)}| \leq 2s^*$ holds for some $q \geq 2$. Next, we will prove $|\mathcal{E}^{(q)}| \leq 2s^*$ holds at q . For any $(j, k) \in \mathcal{S}^{(q)}$, we obtain $\left| \tilde{\Sigma}_{jk}^{(q-1)} \right| \geq u$ and further have

$$\begin{aligned} \sqrt{|\mathcal{S}^{(q)} \setminus \mathcal{S}^*|} &\leq \sqrt{\sum_{(j,k) \in \mathcal{S}^{(q)} \setminus \mathcal{S}^*} \left(u^{-1} \tilde{\Sigma}_{jk}^{(q-1)} \right)^2} \\ &= u^{-1} \left\| \tilde{\Sigma}_{\mathcal{S}^{(q)} \setminus \mathcal{S}^*}^{(q-1)} \right\|_{\text{F}} \\ &= u^{-1} \left\| \left(\tilde{\Sigma}^{(q-1)} - \Sigma^* \right)_{\mathcal{S}^{(q)} \setminus \mathcal{S}^*} \right\|_{\text{F}} \\ &\leq u^{-1} \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}}. \end{aligned} \quad (11)$$

For any $(j, k) \in \overline{\mathcal{S}^{(q-1)}}$, we have $\Lambda_{jk}^{(q-2)} = p'_\lambda \left(\tilde{\Sigma}_{jk}^{(q-2)} \right) \geq p'_\lambda(u) \geq \frac{\lambda}{2}$, which implies

$$\left\| \Lambda_{\mathcal{E}^{(q-1)}}^{(q-2)} \right\|_{\min} \geq \left\| \Lambda_{\mathcal{S}^{(q-1)}}^{(q-1)} \right\|_{\min} \geq p'_\lambda(u) \geq \frac{\lambda}{2}.$$

One also has $|\mathcal{E}^{(q-1)}| \leq 2s^*$ and $\mathcal{S}^* \subseteq \mathcal{E}^{(q-1)}$. Applying Lemma 20 with $\tilde{\Sigma} = \tilde{\Sigma}^{(q-1)}$, $\mathcal{E} = \mathcal{E}^{(q-1)}$, and $\Lambda_{\mathcal{S}^*} = \Lambda_{\mathcal{S}^*}^{(q-2)}$ yields

$$\left\| \tilde{\Sigma}^{(q-2)} - \Sigma^* \right\|_{\text{F}} \leq \frac{2+\sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}.$$

Substituting the above result into the inequality 11 yields

$$\sqrt{|\mathcal{S}^{(q)} \setminus \mathcal{S}^*|} \leq \frac{2+\sqrt{2}}{2u\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*} = \sqrt{s^*}.$$

Thus, we have

$$|\mathcal{E}^{(q)}| = |\mathcal{S}^* \cup (\mathcal{S}^{(q)} \setminus \mathcal{S}^*)| = |\mathcal{S}^*| + |\mathcal{S}^{(q)} \setminus \mathcal{S}^*| \leq 2s^*,$$

completing the induction.

Then by the definition of $\mathcal{E}^{(q)}$ and $\mathcal{S}^{(q)}$, we have

$$\left\| \Lambda_{\mathcal{E}^{(q)}}^{(q-1)} \right\|_{\min} \geq \left\| \Lambda_{\mathcal{S}^{(q)}}^{(q)} \right\|_{\min} \geq p'_\lambda(u) \geq \frac{\lambda}{2}.$$

Applying Lemma 20 with $\tilde{\Sigma} = \tilde{\Sigma}^{(q)}$, $\mathcal{E} = \mathcal{E}^{(q)}$, and $\Lambda_{\mathcal{S}^*} = \Lambda_{\mathcal{S}^*}^{(q-1)}$, the ε -optimal solution $\tilde{\Sigma}^{(q)}$ satisfies

$$\begin{aligned} \left\| \tilde{\Sigma}^{(q)} - \Sigma^* \right\|_{\text{F}} &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\left\| (\nabla f(\Sigma^*))_{\mathcal{E}^{(q)}} \right\|_{\text{F}} + \varepsilon \sqrt{|\mathcal{E}^{(q)}|} + \left\| \Lambda_{\mathcal{S}^*}^{(q-1)} \right\|_{\text{F}} \right) \\ &\leq \frac{2+\sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}. \end{aligned}$$

■

B. Proof of Theorem 11

Proof: Based on Lemma 21, we have

$$\begin{aligned} \left\| \tilde{\Sigma}^{(q)} - \Sigma^* \right\|_{\text{F}} &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \underbrace{\left(\left\| (\nabla f(\Sigma^*))_{\mathcal{E}^{(q)}} \right\|_{\text{F}} + \varepsilon \sqrt{|\mathcal{E}^{(q)}|} \right)}_{\text{V}} \\ &\quad + \frac{1}{\rho_{2s^*+2\tilde{s}}} \underbrace{\left\| \Lambda_{\mathcal{S}^*}^{(q-1)} \right\|_{\text{F}}}_{\text{VI}}. \end{aligned} \quad (12)$$

Then, we proceed to establish bounds for the term **V** and **VI**, respectively.

For term **V**, dividing the support set into \mathcal{S}^* and $\mathcal{E}^{(q)} \setminus \mathcal{S}^*$, we obtain

$$\begin{aligned} \text{V} &\leq \left\| (\nabla f(\Sigma^*))_{\mathcal{S}^*} \right\|_{\text{F}} + \varepsilon \sqrt{s^*} \\ &\quad + (\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon) \sqrt{|\mathcal{E}^{(q)} \setminus \mathcal{S}^*|} \\ &\leq \left\| (\nabla f(\Sigma^*))_{\mathcal{S}^*} \right\|_{\text{F}} + \varepsilon \sqrt{s^*} + \frac{\lambda}{2u} \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}}, \end{aligned} \quad (13)$$

where the second inequality is due to

$$\sqrt{|\mathcal{E}^{(q)} \setminus \mathcal{S}^*|} = \sqrt{|\mathcal{S}^{(q)} \setminus \mathcal{S}^*|} \leq u^{-1} \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}}, \quad (14)$$

which follows from the inequality (11).

By Assumption 2 and 8, for any Σ , if $|\Sigma_{jk} - \Sigma_{jk}^*| \geq u$, then $p'_\lambda(\Sigma_{jk}) \leq \lambda \leq \lambda u^{-1} |\Sigma_{jk} - \Sigma_{jk}^*|$; otherwise,

$p'_\lambda(\Sigma_{jk}) \leq p'_\lambda(|\Sigma_{jk}^*| - u) = 0$. Therefore, for term **VI**, we have

$$\mathbf{VI} \leq \lambda u^{-1} \left\| \tilde{\Sigma}_{S^*}^{(q-1)} - \Sigma_{S^*}^* \right\|_{\text{F}} \leq \lambda u^{-1} \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}}. \quad (15)$$

Substituting the above results into (12) yields

$$\begin{aligned} & \left\| \tilde{\Sigma}^{(q)} - \Sigma^* \right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{s^*+2\tilde{s}}} \left(\|(\nabla f(\Sigma^*))_{S^*}\|_{\text{F}} + \varepsilon \sqrt{s^*} \right) + \delta \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}}, \end{aligned} \quad (16)$$

where $\delta = \frac{3\lambda}{2u\rho_{s^*+2\tilde{s}}} = \frac{3}{2+\sqrt{2}} \in (0, 1)$. ■

C. Concentration Inequality

Lemma 22. *Under Assumptions 8, there exists some constant c_1 , the following result holds*

$$\mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} \leq \lambda \right) \geq 1 - \frac{c_1}{d}.$$

Proof: We start by bounding

$$\mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} > t \right).$$

First of all, using the union bound, we obtain

$$\begin{aligned} & \mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} > t \right) \\ & = \mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \eta_i \mathbf{a}_i \mathbf{a}_i^\top \right\|_{\max} > t \right) \\ & \leq \sum_{l=1}^{d^2} \mathbb{P} \left(\frac{1}{m} \left| \tilde{\mathbf{A}}_{\cdot l}^\top \cdot \boldsymbol{\eta} \right| > t \right), \end{aligned} \quad (17)$$

where $\boldsymbol{\eta} = [\eta_1, \dots, \eta_m]^\top$. Let $\varrho_l = \tilde{\mathbf{A}}_{\cdot l}^\top \cdot \boldsymbol{\eta}$. Since η_i is sub-Gaussian $(0, \sigma^2)$ for $i = 1, \dots, m$, we obtain

$$\begin{aligned} & \mathbb{E}(\exp\{t_0 \varrho_l\} + \exp\{-t_0 \varrho_l\}) \\ & \leq 2 \exp \left\{ m^{-2} \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2^2 \sigma_\eta^2 t_0^2 / 2 \right\}, \end{aligned} \quad (18)$$

which implies

$$\mathbb{P}(|\varrho_l| \geq t) \exp\{t_0 t\} \leq 2 \exp \left\{ m^{-2} \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2^2 \sigma_\eta^2 t_0^2 / 2 \right\}.$$

Taking $t_0 = t \left(m^{-2} \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2^2 \sigma_\eta^2 \right)^{-1}$ yields that

$$\mathbb{P}(|\varrho_l| \geq t) \leq 2 \exp \left\{ -\frac{t^2}{2\sigma_\eta^2 \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2^2 / m^2} \right\}.$$

Plugging it into (17) results

$$\begin{aligned} & \mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} > t \right) \\ & \stackrel{(i)}{\leq} 2d^2 \exp \left\{ -\frac{mt^2}{2\sigma_\eta^2 \max_l \left\{ \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2^2 / m \right\}} \right\} \\ & = c_1 \exp \{-c_2 m t^2 + 2 \log d\}, \end{aligned}$$

where (i) is from Lemma 15 that the column of $\tilde{\mathbf{A}}$ are normalized such that $\max_l \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|_2 \leq \sqrt{m}$, and c_1, c_2 are constants. Then taking $\lambda = \sqrt{\frac{3 \log d}{c_3 m}} \asymp \sqrt{\frac{\log d}{m}}$, we obtain

$$\begin{aligned} & \mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} \leq \lambda \right) \\ & \geq 1 - c_1 \exp(-c_2 m \lambda^2 + 2 \log d) \\ & = 1 - \frac{c_1}{d}. \end{aligned}$$

■

Lemma 23. *Under the same conditions in Lemma 22, the following result hold*

$$\left\| \left(\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right)_{S^*} \right\|_{\text{F}} = O_p \left(\sqrt{\frac{s^*}{m}} \right).$$

Proof: For simplicity, we define

$$\nabla f_1(\Sigma^*) = -\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top.$$

Similar to Lemma 22 analysis, for any M such that $0 < M \sqrt{\frac{1}{m}}$, we obtain

$$\begin{aligned} & \mathbb{P} \left(\|(-\nabla f_1(\Sigma^*))_{S^*}\|_{\max} > M \sqrt{\frac{1}{m}} \right) \\ & \leq c_1 s^* \exp(-c_2 M) \\ & = c_1 \exp(-c_2 M + \log s^*). \end{aligned}$$

By taking M such that $\sqrt{\frac{2 \log s^*}{c_3}} < M$ and $M \rightarrow \infty$ in the above inequality obtains

$$\lim_{M \rightarrow \infty} \sup_m \mathbb{P} \left(\|(-\nabla f_1(\Sigma^*))_{S^*}\|_{\max} > M \sqrt{\frac{1}{m}} \right) = 0.$$

The proof is completed by applying

$$\begin{aligned} & \left\| \left(\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right)_{S^*} \right\|_{\text{F}} \\ & \leq \sqrt{s^*} \left\| \left(\frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right)_{S^*} \right\|_{\max}. \end{aligned}$$

■

D. Proof of Corollary 13

Proof: One has

$$\|\nabla f(\Sigma^*)\|_{\max} \leq \left\| \frac{1}{m} \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} + \tau \|(\Sigma^*)^{-1}\|_{\max}. \quad (19)$$

If τ satisfy $\tau \lesssim \sqrt{\frac{1}{m}} \|(\Sigma^*)^{-1}\|_{\max}^{-1}$, then by Lemma 22, $\lambda \geq 2(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon)$ holds w.h.p.

Applying Lemma 21 with $q = 1$, we obtain

$$\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}. \quad (20)$$

If Assumption 8 holds, that is, $\lambda \asymp \sqrt{\frac{\log d}{m}}$, then $\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \lesssim \sqrt{\frac{s^* \log d}{m}}$ w.h.p. ■

E. Proof of Corollary 14

Proof: One has

$$\|\nabla f(\Sigma^*)\|_{\max} \leq \left\| \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right\|_{\max} + \tau \|(\Sigma^*)^{-1}\|_{\max}.$$

If Assumptions 8 and 10 hold, and τ satisfy

$$\tau \lesssim \sqrt{\frac{1}{m}} \|(\Sigma^*)^{-1}\|_{\max}^{-1},$$

then by Lemma 22, $\lambda \geq 2(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon)$ holds w.h.p.

Applying Theorem 11, we obtain

$$\begin{aligned} & \left\| \tilde{\Sigma}^{(q)} - \Sigma^* \right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{s^*+2\tilde{s}}} \left(\|\nabla f(\Sigma^*)\|_{\text{F}} + \varepsilon \sqrt{s^*} \right) + \delta \left\| \tilde{\Sigma}^{(q-1)} - \Sigma^* \right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{s^*+2\tilde{s}}} \left(\|\nabla f(\Sigma^*)\|_{\text{F}} + \varepsilon \sqrt{s^*} \right) + \delta^{k-1} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{s^*+2\tilde{s}}} \left(\|\nabla f(\Sigma^*)\|_{\text{F}} + \varepsilon \sqrt{s^*} \right) + \delta^{k-1} \frac{2 + \sqrt{2}}{2\rho_{s^*+2\tilde{s}}} \lambda \sqrt{s^*}, \end{aligned} \quad (21)$$

where the last inequality is due to $\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \leq \frac{2 + \sqrt{2}}{2\rho_{s^*+2\tilde{s}}} \lambda \sqrt{s^*}$, which follows from Lemma 21 with $q = 1$.

One has

$$\begin{aligned} & \|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_{\text{F}} \\ & = \left\| \left(- \sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top - \tau (\Sigma^*)^{-1} \right) \right\|_{\mathcal{S}^* \text{F}} \\ & \leq \left\| \left(\sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right) \right\|_{\mathcal{S}^* \text{F}} + \tau \|(\Sigma^*)^{-1}\|_{\mathcal{S}^* \text{F}} \\ & \leq \left\| \left(\sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right) \right\|_{\mathcal{S}^* \text{F}} + \tau \|(\Sigma^*)^{-1}\|_{\text{F}} \end{aligned} \quad (22)$$

By Lemma 23, $\left\| \left(\sum_{i=1}^m \mathbf{a}_i (y_i - \mathbf{a}_i^\top \Sigma^* \mathbf{a}_i) \mathbf{a}_i^\top \right) \right\|_{\mathcal{S}^* \text{F}} \lesssim \sqrt{\frac{s^*}{m}}$ holds w.h.p. ■

APPENDIX B FURTHER INTERMEDIATE RESULTS

Lemma 24. Given $\omega_{\Lambda^{(q-1)}}(\tilde{\Sigma}^{(q)}) \leq \frac{\lambda}{8}$, we have that for all $t \geq 1$ at the $(q+1)$ -th stage,

$$\omega_{\Lambda^{(q)}}(\Sigma_t) \leq \frac{\lambda}{4}$$

and

$$F_{\Lambda^{(q)}}(\Sigma_t) \leq F_{\Lambda^{(q)}}(\Sigma^*) + \frac{\lambda}{4} \|\Sigma_t - \Sigma^*\|_1.$$

Proof: Note that at the $(q+1)$ -th stage, $\Sigma_0 = \tilde{\Sigma}^{(q)}$. Then we have

$$\begin{aligned} \omega_{\Lambda^{(q)}}(\Sigma_t) &= \min_{\Xi \in \partial \|\Sigma_0\|_1} \left\| \nabla f(\Sigma_0) + \Lambda^{(q)} \odot \Xi \right\|_{\max} \\ &\leq \min_{(i)} \left\| \nabla f(\Sigma_0) + \Lambda^{(q-1)} \odot \Xi \right\|_{\max} \\ &\quad + \left\| \left(\Lambda^{(q)} - \Lambda^{(q-1)} \right) \odot \Xi \right\|_{\max} \\ &\leq \omega_{\Lambda^{(q-1)}}(\Sigma_0) + \left\| \Lambda^{(q)} - \Lambda^{(q-1)} \right\|_{\max} \\ &\stackrel{(ii)}{\leq} \frac{\lambda}{8} + \frac{\lambda}{8} \leq \frac{\lambda}{4}, \end{aligned}$$

where (i) is from triangle inequality, (ii) is from the definition of the approximate KKT condition and Ξ , and (iii) is from $\omega_{\Lambda^{(q-1)}}(\Sigma_0) = \omega_{\Lambda^{(q-1)}}(\tilde{\Sigma}^{(q)}) \leq \frac{\lambda}{8}$ and $\left\| \Lambda^{(q)} - \Lambda^{(q-1)} \right\|_{\max} \leq \frac{\lambda}{8}$.

For some $\Xi_t = \arg \min_{\Xi \in \partial \|\Sigma_t\|_1} \left\| \nabla f(\Sigma_t) + \Lambda^{(q)} \odot \Xi \right\|_{\max}$, we have

$$\begin{aligned} & F_{\Lambda^{(q)}}(\Sigma^*) \\ & \stackrel{(i)}{\geq} F_{\Lambda^{(q)}}(\Sigma_t) - \left\langle \nabla f(\Sigma_t) + \Lambda^{(q)} \odot \Xi_t, \Sigma_t - \Sigma^* \right\rangle \\ & \geq F_{\Lambda^{(q)}}(\Sigma_t) - \left\| \nabla f(\Sigma_t) + \Lambda^{(q)} \odot \Xi_t \right\|_{\max} \|\Sigma_t - \Sigma^*\|_1 \\ & \stackrel{(ii)}{\geq} F_{\Lambda^{(q)}}(\Sigma_t) - \frac{\lambda}{4} \|\Sigma_t - \Sigma^*\|_1, \end{aligned}$$

where (i) is from the convexity of F and (ii) is from the fact that for all $t \geq 0$, $\left\| \nabla f(\Sigma_t) + \Lambda^{(q)} \odot \Xi_t \right\|_{\max} \leq \frac{\lambda}{4}$. ■

Lemma 25. Suppose $\left\| \Sigma_{\tilde{S}}^t \right\|_0 \leq \tilde{s}$ and $\omega_{\Lambda^{(q)}}(\Sigma_t) \leq \frac{\lambda}{4}$. Then there exists a generic constant c_1 such that

$$\left\| \Sigma_t - \Sigma^* \right\|_F \leq \frac{c_1 \lambda \sqrt{s^*}}{\rho_{2s^*+2\tilde{s}}}.$$

Proof: It is similar to the proof of Lemma 21, we omit it for simplicity. For more details, we refer readers to [43]. ■

APPENDIX C

PRELIMINARY LEMMAS

Lemma 26. Consider a set \mathcal{E} such that $\mathcal{S}^* \subseteq \mathcal{E}$, if $\left\| \nabla f(\Sigma^*) \right\|_{\max} + \varepsilon \leq \left\| \Lambda_{\mathcal{E}} \right\|_{\min}$, we have

$$\left\| \left(\tilde{\Sigma} - \Sigma^* \right)_{\bar{\mathcal{E}}} \right\|_1 \leq 5 \left\| \left(\tilde{\Sigma} - \Sigma^* \right)_{\mathcal{E}} \right\|_1.$$

Proof: Define $\Omega = \nabla f(\tilde{\Sigma}) + \Lambda \odot \tilde{\Xi}$ and $\tilde{\Delta} = \tilde{\Sigma} - \Sigma^*$. By the mean value theorem, there exist a $\theta \in [0, 1]$, such that

$$\nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*) = \text{mat} \left(H(\theta) \text{vec}(\tilde{\Delta}) \right),$$

where $H(\theta) = \nabla^2 f(\theta \Sigma^* + (1-\theta)\tilde{\Sigma})$. Then we have

$$\begin{aligned} & \left\langle \nabla f(\tilde{\Sigma}) + \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \right\rangle \\ &= \left\langle \nabla f(\Sigma^*) + \text{mat} \left(H(\theta) \text{vec}(\tilde{\Delta}) \right) + \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \right\rangle \\ &\leq \left\| \Omega \right\|_{\max} \left\| \tilde{\Delta} \right\|_1. \end{aligned}$$

Using the fact $\left\| \tilde{\Delta} \right\|_{H(\theta)}^2 \geq 0$, we have

$$0 \leq \left\| \Omega \right\|_{\max} \left\| \tilde{\Delta} \right\|_1 - \underbrace{\left\langle \nabla f(\Sigma^*), \tilde{\Delta} \right\rangle}_{\text{I}} - \underbrace{\left\langle \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \right\rangle}_{\text{II}}. \quad (23)$$

For term I, separating the support of $\nabla f(\Sigma^*)$ and $\tilde{\Delta}$ to \mathcal{E} and $\bar{\mathcal{E}}$, and then using the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{I} &= \left\langle (\nabla f(\Sigma^*))_{\mathcal{E}}, \tilde{\Delta}_{\mathcal{E}} \right\rangle + \left\langle (\nabla f(\Sigma^*))_{\bar{\mathcal{E}}}, \tilde{\Delta}_{\bar{\mathcal{E}}} \right\rangle \\ &\geq -\left\| (\nabla f(\Sigma^*))_{\mathcal{E}} \right\|_F \left\| \tilde{\Delta}_{\mathcal{E}} \right\|_F - \left\| (\nabla f(\Sigma^*))_{\bar{\mathcal{E}}} \right\|_{\max} \left\| \tilde{\Delta}_{\bar{\mathcal{E}}} \right\|_1 \\ &\geq -\left\| (\nabla f(\Sigma^*))_{\mathcal{E}} \right\|_F \left\| \tilde{\Delta} \right\|_F - \left\| (\nabla f(\Sigma^*)) \right\|_{\max} \left\| \tilde{\Delta} \right\|_1. \end{aligned} \quad (24)$$

For term II, separating the support of $\Lambda \odot \tilde{\Xi}$ and $\tilde{\Delta}$ to \mathcal{E} and $\bar{\mathcal{E}}$, and then using the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{II} &= \left\langle (\Lambda \odot \tilde{\Xi})_{\mathcal{E}}, \tilde{\Delta}_{\mathcal{E}} \right\rangle + \left\langle (\Lambda \odot \tilde{\Xi})_{\bar{\mathcal{E}}}, \tilde{\Delta}_{\bar{\mathcal{E}}} \right\rangle \\ &\geq -\left\| \Lambda_{\mathcal{E}} \right\|_{\max} \left\| \tilde{\Delta}_{\mathcal{E}} \right\|_1 + \left\| \Lambda_{\bar{\mathcal{E}}} \right\|_{\min} \left\| \tilde{\Delta}_{\bar{\mathcal{E}}} \right\|_1. \end{aligned} \quad (25)$$

Plugging the above results into 23 yields

$$\begin{aligned} \left\| \tilde{\Delta}_{\bar{\mathcal{E}}} \right\|_1 &\leq \frac{\lambda + \left\| (\nabla f(\Sigma^*)) \right\|_{\max} + \omega_{\Lambda}(\tilde{\Sigma})}{\left\| \Lambda_{\bar{\mathcal{E}}} \right\|_{\min} - \left(\left\| (\nabla f(\Sigma^*)) \right\|_{\max} + \omega_{\Lambda}(\tilde{\Sigma}) \right)} \left\| \tilde{\Delta}_{\mathcal{E}} \right\|_1 \\ &\stackrel{(i)}{\leq} 5 \left\| \tilde{\Delta}_{\mathcal{E}} \right\|_1, \end{aligned}$$

where (i) is from $\left\| \Lambda_{\bar{\mathcal{E}}} \right\|_{\min} \geq \frac{\lambda}{2}$, $\left\| (\nabla f(\Sigma^*)) \right\|_{\max} \leq \frac{\lambda}{4}$ and $\omega_{\Lambda}(\tilde{\Sigma}) \leq \varepsilon$. ■

Lemma 27. Recall $D_f(\Sigma_1, \Sigma_2)$ and $D_f^s(\Sigma_1, \Sigma_2)$ defined in Lemma 19. For $\Sigma(\theta) = \Sigma^* + \theta(\Sigma - \Sigma^*)$ with $\theta \in (0, 1]$, we have

$$D_f^s(\Sigma(\theta), \Sigma^*) \leq \theta D_f^s(\Sigma, \Sigma^*).$$

Proof: Let

$$\begin{aligned} \varphi(\theta) &= D_f(\Sigma(\theta), \Sigma^*) \\ &= f(\Sigma(\theta)) - f(\Sigma^*) - \langle \nabla f(\Sigma^*), \Sigma(\theta) - \Sigma^* \rangle. \end{aligned}$$

Since the derivative of $f(\Sigma(\theta))$ with respect to θ is $\langle \nabla f(\Sigma(\theta)), \Sigma - \Sigma^* \rangle$, it follows that

$$\varphi'(\theta) = \langle \nabla f(\Sigma(\theta)) - \nabla f(\Sigma^*), \Sigma - \Sigma^* \rangle.$$

Therefore, the symmetric Bregman divergence $D_f^s(\Sigma(\theta), \Sigma^*)$ can be written as

$$\begin{aligned} D_f^s(\Sigma(\theta), \Sigma^*) &= \langle \nabla f(\Sigma(\theta)) - \nabla f(\Sigma^*), \theta(\Sigma - \Sigma^*) \rangle \\ &= \theta \varphi'(\theta) \end{aligned}$$

for $0 < \theta \leq 1$. Plugging $\theta = 1$ in the equation above, we have $\varphi'(1) = D_f^s(\Sigma, \Sigma^*)$ as a special case. If we assume that $\varphi(\theta)$ is convex, then $\varphi'(\theta)$ is non-decreasing and thus

$$D_f^s(\Sigma(\theta), \Sigma^*) = \theta \varphi'(\theta) \leq \theta \varphi'(1) = \theta D_f^s(\Sigma, \Sigma^*).$$

It remains to show the convexity of $\varphi(\theta)$ with respect to θ , i.e.

$$\Sigma(\alpha_1 \theta_1 + \alpha_2 \theta_2) \leq \alpha_1 \Sigma(\theta_1) + \alpha_2 \Sigma(\theta_2),$$

for $\theta_1, \theta_2 \in (0, 1]$, $\alpha_1, \alpha_2 \geq 0$ and $\alpha_1 + \alpha_2 = 1$; or equivalently, the convexity of $f(\Sigma(\theta))$ and $\langle \nabla f(\Sigma^*), \Sigma^* - \Sigma(\theta) \rangle$, respectively.

For $\forall \alpha_1, \alpha_2 \geq 0$ such that $\alpha_1 + \alpha_2 = 1$, and $\theta_1, \theta_2 \in (0, 1]$, we have $\Sigma(\alpha_1 \theta_1 + \alpha_2 \theta_2) = \alpha_1 \Sigma(\theta_1) + \alpha_2 \Sigma(\theta_2)$. By the bi-linearity property of the inner product function $\langle \cdot, \cdot \rangle$, and using the linearity property of $\Sigma(\cdot)$, we have the following equality hold

$$\begin{aligned} & -\langle \nabla f(\Sigma^*), \Sigma(\alpha_1 \theta_1 + \alpha_2 \theta_2) - \Sigma^* \rangle \\ &= -\alpha_1 \langle \nabla f(\Sigma^*), \Sigma(\theta_1) - \Sigma^* \rangle - \alpha_2 \langle \nabla f(\Sigma^*), \Sigma(\theta_2) - \Sigma^* \rangle. \end{aligned} \quad (26)$$

On the other side, by the convexity of the loss function $f(\cdot)$, we obtain

$$\begin{aligned} & f(\Sigma(\alpha_1 \theta_1 + \alpha_2 \theta_2)) \\ &= f(\alpha_1 \Sigma(\theta_1) + \alpha_2 \Sigma(\theta_2)) \\ &\leq \alpha_1 f(\Sigma(\theta_1)) + \alpha_2 f(\Sigma(\theta_2)). \end{aligned} \quad (27)$$

By adding 26 and 27 together and using the definition of function $\Sigma(\cdot)$, we obtain

$$\Sigma(\alpha_1 \theta_1 + \alpha_2 \theta_2) \leq \alpha_1 \Sigma(\theta_1) + \alpha_2 \Sigma(\theta_2),$$

which indicates $\Sigma(\theta)$ is a convex function. Thus we complete our proof. ■

Lemma 28. *Under the same condition of Lemma 21, we have the following basic inequality*

$$\left\langle \nabla f\left(\tilde{\Sigma}^{(1)}\right) - \nabla f\left(\Sigma^*\right), \tilde{\Sigma}^{(1)} - \Sigma^* \right\rangle \leq \frac{c_1 \lambda^2 s^*}{\rho_{2s^*+2\tilde{s}}^-}.$$

Proof: Applying Lemma 21 with $q = 1$, we obtain

$$\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}^-} \lambda \sqrt{s^*}. \quad (28)$$

On the other side, applying Lemma 26 yields that

$$\begin{aligned} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 &\leq \left\| \left(\tilde{\Sigma}^{(1)} - \Sigma^* \right)_{\overline{\mathcal{E}\{1\}}} \right\|_1 + \left\| \left(\tilde{\Sigma}^{(1)} - \Sigma^* \right)_{\mathcal{E}^{(1)}} \right\|_1 \\ &\leq 6 \left\| \left(\tilde{\Sigma}^{(1)} - \Sigma^* \right)_{\mathcal{E}^{(1)}} \right\|_1, \end{aligned}$$

where $\mathcal{E}^{(1)}$ can be taken as \mathcal{S}^* . This, combined with 28, results

$$\begin{aligned} \left\| \left(\tilde{\Sigma}^{(1)} - \Sigma^* \right)_{\mathcal{S}^*} \right\|_1 &\leq \sqrt{s^*} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\text{F}} \\ &\leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}^-} \lambda s^*. \end{aligned}$$

Therefore, we obtain

$$\begin{aligned} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 &\leq 6 \left\| \left(\tilde{\Sigma}^{(1)} - \Sigma^* \right)_{\mathcal{S}^*} \right\|_1 \\ &\leq \frac{3(2 + \sqrt{2})}{\rho_{2s^*+2\tilde{s}}^-} \lambda s^*. \end{aligned}$$

Because $\tilde{\Sigma}^{(1)}$ is a ε -optimal solution, we have

$$\begin{aligned} &\left\langle \nabla f\left(\tilde{\Sigma}^{(1)}\right) - \nabla f\left(\Sigma^*\right), \tilde{\Sigma}^{(1)} - \Sigma^* \right\rangle \\ &\leq \left\| \nabla f\left(\tilde{\Sigma}^{(1)}\right) + \Lambda \odot \tilde{\Xi}^{(1)} - \Lambda \odot \tilde{\Xi}^{(1)} - \nabla f\left(\Sigma^*\right) \right\|_{\max} \times \\ &\quad \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 \\ &\leq \left(1 + \frac{1}{4} \right) \lambda \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 \leq \frac{c_1 \lambda^2 s^*}{\rho_{2s^*+2\tilde{s}}^-}, \end{aligned}$$

for some constant c_1 . ■