

Optimal Compressive Covariance Sensing via Quadratic Sampling based on Nonconvex Learning

Wenbin Wang, *Student Member, IEEE* and Ziping Zhao, *Member, IEEE*

Abstract—Covariance matrix, capturing the degree of linear correlation between multiple variables, is a fundamental statistical quantity in diverse data analysis applications. Traditional method for estimation the covariance matrix typically assume full access to all available measurements. However, it becomes impractical when data evolves rapidly or the data acquisition devices have limited processing power and storage. To address these challenges, compressive covariance sensing (CCS) has emerged as a promising approach, enabling covariance matrix estimation from a reduced number of measurements, often significantly smaller than the dimension of the variables. In this paper, we focus on the quadratic (or rank-one) measurement model for CCS under the assumption of a sparse covariance structure, which minimizes memory requirements and computational complexity during the sampling process. We propose a least-squares estimator for the covariance matrix, regularized with positive-definiteness and non-convex sparsity-inducing penalties. To efficiently compute this estimator, we develop a multistage convex relaxation algorithm based on the majorization-minimization (MM) framework. This algorithm provides strong computational and statistical guarantees by leveraging local restricted strong convexity and Hessian smoothness. Our proposed algorithm achieves quadratic convergence to approximate solutions at each stage of convex relaxation. We further rigorously establish the statistical performance of all sequential approximate solutions generated by the MM-based algorithm and demonstrate that the optimal solution reached after sufficient iterations possesses oracle statistical properties in the Frobenius norm. Comprehensive numerical simulations are included to validate our theoretical results.

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

I. INTRODUCTION

The covariance matrix is a fundamental statistical tool used to quantify the degree of linear correlation among multiple random variables [1]–[3]. It plays a pivotal role in multivariate data processing across science and engineering, including component analysis [4], factor analysis [5], beamformer design [6], portfolio optimization [7], among others. Despite its significance, the covariance matrix is not directly observable and must be estimated from empirical data, thus requiring the development of efficient estimation techniques.

Consider n independent observations $\{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^n$ drawn from a zero-mean random vector \mathbf{x} . The sample covariance matrix (SCM), $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$, is a commonly adopted covariance estimator. However, in high-dimensional settings where the number of variables significantly exceeds the number of observations, the SCM becomes unreliable [8], [9]. To address this, it is crucial to exploit structural nature of the

covariance matrix, with sparsity being one of the most favored assumptions. Sparsity suggests that many covariance entries are zero, thereby substantially reducing the number of parameters to be estimated. Thresholding techniques [10]–[12], which leverage sparse assumption, are widely used to adjust the SCM by eliminating elements below a specific threshold, thereby improving estimation accuracy. The statistical properties of these thresholding covariance estimators, including minimax lower bounds [13] and convergence rates [14], have been extensively investigated, validating both their theoretical foundations and practical applicability. Thresholding techniques have been theoretically validated to yield estimators that are asymptotically positive definite [10], [11]. However, practical applications often require that estimators be positive definite even with a finite number of samples. This property is critical in supervised learning contexts, where methods like quadratic discriminant analysis and covariance regularized regression depend on nonnegative definite covariance matrices for the convexity of their optimization landscapes. To ensure both sparsity and positive definiteness in the estimators, Rothman [15] proposed integrating a log-determinant barrier function into the soft-thresholding covariance matrix estimation problem.

Despite covariance estimation based on full measurements has been extensively studied in the literature, it faces great challenge in the scenario characterized by limited sampling capabilities and rapidly changing data environments, where the high dimensionality of the data renders it either impractical to obtain full data samples. A viable alternative approach is to estimate covariance matrix based on compressed measurements, a technique known as compressive covariance sensing or compressed covariance sketching (CCS). CCS offers significant advantages by enabling the recovery of covariance matrices from data streams while substantially reducing computation and storage overhead. Compressive sensing for covariance matrices generalizes the classical compressive sensing problem for signals [16], [17], focusing on extracting second-order statistics rather than reconstructing random vectors.

CCS is an advanced extension of the foundational principles of compressed sensing [16], [17]. In compressed sensing, the presumption of signal sparsity underpins many of the theoretical and practical developments [16], [18]. This assumption has been effectively adapted for CCS to focus on the covariance structure. [19] proposed methods for estimating sparse covariance from quadratic sampling via convex programming, providing guarantee on performance. Unlike the traditional Restricted Isometry Property ($\text{RIP-}\ell_2/\ell_2$) which does not directly apply to CCS [20], a novel mixed-norm $\text{RIP-}\ell_2/\ell_1$ has

been introduced. This criterion uses the Frobenius norm and the ℓ_1 norm to assess input and output strengths, respectively, ensuring that the dimensionality reduction projection retains signal integrity and enables perfect recovery from minimal samples. A similar criterion, $\text{RIP-}\ell_1/\ell_1$, is also applicable to sparse covariance recovery [20]. Unlike previous studies, this paper focuses on the sparse eigenvalue condition, a related but weaker concept compared to the restricted isometry property [21], [22]. Other works have explored the CCS problem under distributed sparsity assumption [23], and further extensions have considered other covariance structures such as banded, Toeplitz, circulant matrices [24].

On the other hand, many practical applications involve random processes where directly reconstructing the random signal is not meaningful. In such case, the focus often shifts towards extracting second-order or higher-order statistics, such as the power spectrum, which contains valuable features and provide reliable information about the process. [25] developed robust compressive techniques for wideband spectrum sensing that leverage the sparsity of two-dimensional cyclic spectra in communications signals. A concurrent work [26] employed a non-negative least squares approach to estimate the signal power spectrum, based on the assumption of sparsity.

A. Contributions

This paper investigates the methodology of covariance matrix sensing, a prevalent technique for estimating high-dimensional sparse covariance matrices. Our primary focus is on the estimation of the covariance matrix utilizing a quadratic measurement model that integrates both a log-determinant penalty and a nonconvex penalty. The introduction of the nonconvex penalty aids in obtaining unbiased estimates; however, it simultaneously introduces considerable difficulties in solving the optimization problem and in performing theoretical analyses. The three questions of interest then are (a) is it feasible to accurately reconstruct the original high-dimensional covariance matrix Σ from the compressed measurements, (b) can a global optimum be achieved, and (c) how can the statistical properties of local optima be characterized. To answer the first question, we examine the sparse eigenvalue, which, to our knowledge, represents the weakest condition documented in the literature for ensuring exact recovery. For the second question, we utilize multiple Monte Carlo simulations to approximate the global optimum. In response to the third question, we propose a multistage convex relaxation technique within the majorization-minimization (MM) framework. This approach guarantees that the approximate local optimum not only converges but also possesses favorable statistical properties that approximate those of the global optimum.

B. Organization

The rest structure of this paper is organized as follows: In Section II, we establish the foundational background knowledge necessary for understanding the problem of covariance matrix sensing. Section III then transitions to the design of the algorithm, wherein we present the development of the proximal Newton method within the MM framework. Following

this, Section IV provides an asymptotic statistical analysis that elucidates the reconstruction performance of the model outlined in (5), specifically addressing the convergence rate of the estimator and the iteration complexity of the algorithm. The findings from numerical experiments are subsequently detailed in Section V. The paper concludes with a discussion in Section VII. The appendices include proofs for all theoretical claims made throughout the paper.

C. Notation

This subsection provides a concise overview of the essential notations employed throughout the paper. Scalar values are denoted by conventional lower-case or upper-case letters, while boldface lower-case letters signify vectors and upper-case letters represent matrices. The symbols $\mathbf{0}$ and \mathbf{I} are used to denote the all-zero matrix and the identity matrix, respectively, with dimensions that are contextually appropriate. The set of all $m \times n$ matrices with real entries is represented by $\mathbb{R}^{m \times n}$.

Let $\mathbf{X} \in \mathbb{R}^{m \times n}$ be a matrix. The elements of the matrix are denoted by X_{ij} and $[\mathbf{X}]_{ij}$, both referring to the entry in the (i, j) -th position. We denote $\mathbf{X} \succeq \mathbf{0}$ (resp. $\mathbf{X} \succ \mathbf{0}$) to indicate that \mathbf{X} is positive semidefinite (resp. definite). The transpose, inverse, and determinant of the matrix \mathbf{X} are represented by \mathbf{X}^\top , \mathbf{X}^{-1} , and $\det(\mathbf{X})$, respectively. The smallest and largest eigenvalues of \mathbf{X} are denoted by $\lambda_{\min}(\mathbf{X})$ and $\lambda_{\max}(\mathbf{X})$. The p -norm of \mathbf{X} is defined as $\|\mathbf{X}\|_p = \left(\sum_{i=1}^m \sum_{j=1}^n |X_{ij}|^p \right)^{\frac{1}{p}}$ for a real $p > 0$. For instance, the Frobenius norm and spectral norm are indicated by $\|\mathbf{X}\|_F$ and $\|\mathbf{X}\|_2$, respectively, while $\|\mathbf{X}\|_1$ represents the sum of the absolute values of all entries of \mathbf{X} . Specifically, the maximum-absolute-value norm of \mathbf{X} is expressed as $\|\mathbf{X}\|_{\max}$, and the minimum-absolute-value norm as $\|\mathbf{X}\|_{\min}$. The notation $\mathbf{X}_{\cdot j}$ ($\mathbf{X}_{k \cdot}$) is used to denote the j -th column (k -th row) of \mathbf{X} . The vectorization of \mathbf{X} , achieved by stacking its columns, is denoted as $\text{vec}(\mathbf{X})$, while $\text{mat}(\mathbf{X})$ represents the inverse operation. Furthermore, the Kronecker product and Hadamard product (also referred to as the entry-wise product) of matrices \mathbf{X} and \mathbf{Y} are denoted by $\mathbf{X} \otimes \mathbf{Y}$ and $\mathbf{X} \odot \mathbf{Y}$, respectively. The Euclidean inner product is defined as $\langle \mathbf{X}, \mathbf{Y} \rangle := \text{tr}(\mathbf{X}\mathbf{Y}^\top)$, where $\text{tr}(\cdot)$ denotes the trace of a matrix. For a real-valued function $f(\mathbf{X})$, the gradient $\nabla f(\mathbf{X})$ is a $d \times d$ matrix with the (i, j) -th element given by $\frac{\partial}{\partial X_{ij}} f(\mathbf{X})$, denoted by $\nabla_{ij} f(\mathbf{X})$, while $\nabla^2 f(\mathbf{X})$ represents the $d^2 \times d^2$ Hessian matrix. Additionally, we define $\|\mathbf{X}\|_{\mathbf{H}}^2 := \text{vec}^\top(\mathbf{X}) \mathbf{H} \text{vec}(\mathbf{X})$ for the purpose of simplified representation.

For an index set \mathcal{E} , its cardinality is represented as $|\mathcal{E}|$, while its complement is denoted by $\bar{\mathcal{E}}$. The notation $\mathbf{X}_{\mathcal{E}}$ refers to the matrix comprising entries X_{ij} for $(i, j) \in \mathcal{E}$ and zero elsewhere. The function $\text{sign}(x)$ is defined as $\text{sign} = x/|x|$ when $x \neq 0$ and is equal to zero otherwise. For the functionals $f(n)$ and $g(n)$, we express $f(n) \gtrsim g(n)$ 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 positive constants c and C . Additionally, $\mathcal{O}_p(\cdot)$ is utilized to indicate boundedness in probability.

II. PRELIMINARIES

In this section, we introduce a comprehensive framework for covariance matrix sensing, delineate its fundamental assumptions, and examine the sparse eigenvalue associated with the design sensing matrix. Subsequently, we provide several practical illustrations of the covariance matrix sensing model. Lastly, we formulate our estimator with nonconvex penalty.

A. The Measurement Model

We consider a conventional framework for indirect measurements, referred to as the quadratic measurement (or rank-one measurement) model, represented mathematically as

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

Here $\{y_i\}_{i=1}^m$ denotes the measurements, \mathbf{S} is the SCM, $\{\mathbf{a}_i \in \mathbb{R}^d\}_{i=1}^m$ signifies the sensing vectors, $\{\eta_i\}_{i=1}^m$ indicates the noise components, and m represents the total number of measurements. The error term η_i ($1 \leq i \leq m$) is presumed to be sampled from a sub-exponential distribution with mean 0 and variance proxy σ^2 . We further analyze the decomposition of the measurements as follows:

$$y_i = \mathbf{a}_i^\top \mathbf{S} \mathbf{a}_i + \eta_i = \mathbf{a}_i^\top (\mathbf{\Sigma}^* + \mathbf{E}) \mathbf{a}_i + \eta_i,$$

where $\mathbf{\Sigma}^*$ is the true covariance matrix, and \mathbf{E} is the bias term. For the sake of clarity, we define $\mathbf{A}_i := \mathbf{a}_i \mathbf{a}_i^\top$ as the corresponding sensing matrix. Consequently, the sampling process can be articulated as $\{\langle \mathbf{A}_i, \mathbf{S} \rangle\}_{i=1}^m$. Thus, the measures $\{y_i\}_{i=1}^m$ can be expressed as $y_i := \langle \mathbf{A}_i, \mathbf{S} \rangle + \eta_i$. Next, we define the vector of measurements as $\mathbf{y} := [y_1, \dots, y_m]^\top$ and the vector of noise as $\boldsymbol{\eta} := [\eta_1, \dots, \eta_m]^\top$. Additionally, we introduce the linear operator $\mathcal{A}(\mathbf{S}) : \mathbb{R}^{d \times d} \mapsto \mathbb{R}^m$, which maps a matrix $\mathbf{S} \in \mathbb{R}^{d \times d}$ to the set $\{\langle \mathbf{A}_i, \mathbf{S} \rangle\}_{i=1}^m$. Consequently, it can be reformulated as $\mathbf{y} = \mathcal{A}(\mathbf{S}) + \boldsymbol{\eta}$. It is noteworthy that the measurement model in (1) has been considered before by [19], [20]. It is essential to emphasize that the measurement model delineated in (1) is considerably more straightforward to implement and entails a reduced computational cost in comparison to full-rank measurement matrices characterized by independently and identically distributed entries.

B. The Design Matrix $\tilde{\mathbf{A}}$

We first present several fundamental assumptions regarding the sensing vectors \mathbf{a}_i , as delineated in Assumption 1.

Assumption 1. We assume that the sensing vectors \mathbf{a}_i 's ($1 \leq i \leq m$) comprise independently and identically distributed (i.i.d.) sub-Gaussian random variables. Each vector \mathbf{a}_i contains d elements, denoted as $(\mathbf{a}_i)_j$ ($1 \leq j \leq d$), which are characterized by the following conditions:

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

This indicates a significant heavy-tailed behavior in these components.

Rewriting Model (1), we equivalently have $y_i = (\mathbf{a}_i \otimes \mathbf{a}_i)^\top \text{vec}(\mathbf{S}) + \eta_i$. Define the design matrix $\tilde{\mathbf{A}}$ as follows:

$$\tilde{\mathbf{A}} = [(\mathbf{a}_1 \otimes \mathbf{a}_1) \quad \dots \quad (\mathbf{a}_m \otimes \mathbf{a}_m)]^\top \in \mathbb{R}^{m \times d^2}.$$

The Assumption 1 asserts that each row of the design matrix $\tilde{\mathbf{A}}$ is composed of i.i.d. sub-exponential random variables. Existing research suggests that to adequately recover a sparse covariance matrix $\mathbf{\Sigma}$ from an underdetermined system of equations, the design matrix $\tilde{\mathbf{A}}$ must exhibit specific favorable properties, specifically adhering to designated incoherence conditions. Various concepts of incoherence have been proposed within the sparse reconstruction literature, such as the Uniform Uncertainty Principle (UUP). Bickel et al. [27] introduced the Restricted Eigenvalue (RE) condition, demonstrating that it is one of the weakest and most general conditions imposed on sparse reconstruction problems in the literature, serving as a relaxation of the UUP. Here, we present one version of the RE condition. First, we outline a formulation of the RE condition as described by Bickel et al. [27].

Definition 2 (Restricted Eigenvalue). For some integer $0 < s_0 < d$ and a positive constant c_0 , the $\text{RE}(s_0, c_0, \mathbf{X})$ for the matrix $\mathbf{X} \in \mathbb{R}^{d \times d}$ necessitates that the following inequality holds:

$$\forall \mathbf{u} \neq 0, \quad \min_{J \subset \{1, \dots, d\}, |J| \leq s_0} \min_{\|\mathbf{u}_J\|_1 \leq c_0 \|\mathbf{u}_{J^c}\|_1} \frac{\|\mathbf{X} \mathbf{u}\|_2}{\|\mathbf{u}_J\|_2} > 0, \quad (3)$$

where \mathbf{v}_J denotes the subvector of $\mathbf{v} \in \mathbb{R}^d$ restricted to the subset J of $\{1, \dots, d\}$.

In the following, we introduce the Sparse Eigenvalues (SE) condition as delineated by Fan et al. [28], which exhibits a significant correlation with the RE properties. For the detailed relationships between RE and SE, please refer to Appendix A.

Definition 3 (Sparse Eigenvalue). For a positive integer $0 < s_0 < d$, the localized sparse eigenvalues are defined as

$$\rho_{s_0, r}^+ = \sup \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\mathbf{\Sigma}) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid \|\mathbf{u}\|_0 \leq s_0, \|\mathbf{\Sigma} - \mathbf{\Sigma}^*\|_F \leq r \right\};$$

$$\rho_{s_0, r}^- = \inf \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\mathbf{\Sigma}) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid \|\mathbf{u}\|_0 \leq s_0, \|\mathbf{\Sigma} - \mathbf{\Sigma}^*\|_F \leq r \right\}.$$

The SE condition has been extensively studied within the domain of high-dimensional sparse recovery or compressed sensing [29]–[33]. This condition serves as a criterion to ensure recovery for anisotropic measurements. For additional conditions closely related to the RE/SE condition, please refer to Van de Geer et al. [34].

The SE condition with parameter s_0 has been demonstrated to be valid for random Gaussian measurements, specifically for a design matrix consisting of $m = \mathcal{O}(s_0 \log d)$ independent instances of a d -dimensional Gaussian random vector \mathbf{x} with covariance matrix $\mathbf{\Sigma}$ [32]. This holds under the assumption (3) is satisfied for the square root of $\mathbf{\Sigma}$. Further analysis by [35] extends these insights by establishing that the exponential width of any set does not exceed $\sqrt{\log d}$ times the Gaussian width of the set. This connection allows for the derivation of Gaussian width-based results in scenarios involving sub-exponential distributions through the application of generic chaining. Therefore, when our design matrix $\tilde{\mathbf{A}}$ satisfies $m = \mathcal{O}(s_0 \log^{\frac{3}{2}} d)$, the SE condition is met with overwhelming probability at least $1 - \exp(-\epsilon^2/2)$ for some $\epsilon > 0$. We present the following assumption:

Assumption 4. *There exists a universal constant c_1 such that for $\tilde{s} \geq c_1 s^*$, the SE property is guaranteed with parameters $\rho_{2s^*+2\tilde{s},r}^-$ and $\rho_{2s^*+2\tilde{s},r}^+$ satisfying*

$$0 < \rho_{2s^*+2\tilde{s},r}^- < \rho_{2s^*+2\tilde{s},r}^+ < +\infty$$

with probability at least $1 - \exp(-\epsilon^2/2)$ for some $\epsilon > 0$.

Assumption 4 asserts that the sparse eigenvalues of the Hessian matrix $\nabla^2 f(\Sigma)$ are lower and upper bounded when Σ is adequate sparsity and in close proximity to Σ^* with high probability (w.h.p). In other words, $\nabla^2 f(\Sigma)$ possesses both bounded maximum and non-zero minimum sparse eigenvalues over a cone.

C. Practical Examples

Practical applications of this quadratic measurement model (1) are extensive and varied. One prominent example is its use in high-frequency communications systems, where noncoherent energy measurements are preferred over phase measurements. In these systems, the model is employed to estimate the energy spectral density of signals, where the focus lies solely on spectral distribution rather than signal recovery [36]. While naturally occurring signals possess infinite duration, the data available is typically limited to finite discrete-time signals. The practical implementation frequently employs random demodulators [37] to capture energy measurements via the sensing vector \mathbf{a}_i , with average energy over T observations given by:

$$y_i = \frac{1}{T} \sum_{t=1}^T |\mathbf{a}_i^\top \mathbf{x}_t|^2 + \eta_i \quad (4)$$

$$= \mathbf{a}_i^\top \mathbf{S} \mathbf{a}_i + \eta_i, \text{ for } i = 1, \dots, m,$$

where $\mathbf{S} = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top$ represents the sample covariance matrix¹ over T time periods. Furthermore, this model plays a crucial role in incoherent imaging within optical systems, offering a precise mathematical framework for describing the response of the imaging process to light intensity. In this scenario, the intensity y_i at each detector is modeled as the ensemble average of the squared projections of the scene state onto the detector's sensitivity pattern. This quadratic formulation is instrumental in determining the energy absorbed by each detector, aiding in the reconstruction of high-quality images from incoherent light measurements. Additional applications encountered in real-time financial [38], IoT sensor data [39] or direction of arrival estimation [40] also leverage quadratic measurement model to captures covariance information on the fly.

D. Proposed Estimator

A successful method for regression estimation is to minimize the least squares errors while adding regularization terms to ensure the low-dimensional structure of the estimator. We

¹In this case, the matrix should be interpreted as the autocorrelation matrix. Nevertheless, assuming the mean is zero, it is equivalent to the covariance matrix. Therefore, the aforementioned trick can be utilized. In Section VI-B, we will discuss the approach to take when the mean is not zero.

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

$$\min_{\Sigma \succ 0} \left\{ \frac{1}{2m} \|\mathbf{y} - \mathcal{A}(\Sigma)\|_2^2 - \tau \log \det \Sigma + \sum_{i,j} p_\lambda(|\Sigma_{ij}|) \right\}. \quad (5)$$

In (5), the first term aims to minimize the empirical errors; the log-determinant barrier function in the second term ensures positive definiteness with a barrier parameter $\tau \geq 0$; p_λ in the third term represents a non-convex penalty function governed by a regularization parameter $\lambda > 0$. We impose certain restrictions on it, as illustrated in Assumption 5.

Assumption 5. *The function $p_\lambda(t) : \mathbb{R} \rightarrow \mathbb{R}$ satisfies:*

- $p_\lambda(t)$ is symmetric around zero with $p_\lambda(0) = 0$, nondecreasing on the nonnegative, differentiable almost everywhere on $(0, +\infty)$, and subdifferentiable at $t = 0$;
- $0 \leq p'_\lambda(t_1) \leq p'_\lambda(t_2) \leq \lambda$ 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 $t \geq \alpha\lambda$;

It is suggested that a well-designed penalty function should yield an estimator with three properties: unbiasedness, sparsity, continuity [41], which coincides with the three conditions in Assumption 5. A variety of functions have been explored through past research on nonconvex regularization, of which we present a few representative examples here:

- Smooth clipped absolute deviation penalty: This penalty, due to [41], is given by

$$p'_\lambda(t) := \begin{cases} \lambda, & \text{for } 0 < t \leq \lambda, \\ \frac{b\lambda - t}{b-1}, & \text{for } \lambda \leq t \leq b\lambda, \\ 0, & \text{for } t \geq b\lambda, \end{cases} \quad (6)$$

where $b > 2$ is an additional tuning parameter. The authors in [41] proposed $b = 3.7$ through a Bayesian rationale, applicable when the variable dimension is less than 100.

- Minimax concave penalty regularizer: This penalty, due to [42], is defined as follows:

$$p_\lambda(t) := \text{sign}(t) \lambda \cdot \int_0^{|t|} \left(1 - \frac{z}{\lambda b}\right)_+ dz \quad (7)$$

for some $b > 0$.

Note that the capped ℓ_1 -penalty, takes the form $p_\lambda(t) := \lambda \cdot \min\{|t|, b\lambda\}$, where $b > 0$ is a fixed parameter. The authors in [22] point out it is a simpler but less smooth version of the smooth clipped absolute deviation penalty.

III. OPTIMIZATION ALGORITHM

In this section, we begin by revisiting the MM framework, which will be pivotal for the forthcoming theoretical analysis. Subsequently, we present an MM-based multistage convex relaxation algorithm to solve (5), along with explicit pseudocode. Additionally, we detail the specific solution method utilized at each stage of the MM-based algorithm — proximal Newton homotopy algorithm.

Algorithm 1: The MM-Based Multistage Convex Relaxation Algorithm for Solving (5).

Input: $\{y_i, \mathbf{a}_i\}_{i=1}^m, \tau, \lambda;$

```

1 Initialize  $\tilde{\Sigma}^{(0)} = \mathbf{I}$ 
2 for  $k = 1, 2, \dots, K$  do
3    $\Lambda_{ij}^{(k-1)} = p'_\lambda \left( \left| \tilde{\Sigma}_{ij}^{(k-1)} \right| \right);$ 
4    $\tilde{\Sigma}^{(k)} = \arg \min F(\Sigma);$ 
5    $k = k + 1;$ 
6 end
Output:  $\tilde{\Sigma}^{(K)}$ 

```

A. The MM Framework and Multistage Convex Relaxation Algorithm

The MM algorithm framework is an iterative process that encompasses two main steps: the Majorization step and the Minimization step. Suppose we want to minimize a real-valued function $F(\mathbf{x})$, below is the basic procedure of this algorithm framework:

- **Majorization Step:** In this step, the goal is to select or construct a surrogate function $\bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)})$ which equals the value of the original objective function $F(\mathbf{x})$ at the current iteration point $\mathbf{x}^{(k-1)}$, but is an upper bound for $F(\mathbf{x})$ at all other points. Specifically, $\bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)}) \geq F(\mathbf{x})$ for all \mathbf{x} and $\bar{F}(\mathbf{x}^{(k-1)} | \mathbf{x}^{(k-1)}) = F(\mathbf{x}^{(k-1)})$.
- **Minimization Step:** After the surrogate function is determined, the next step is to find the point $\mathbf{x}^{(k)}$ that minimizes $\bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)})$, such that $\mathbf{x}^{(k)} \in \arg \min \bar{F}(\mathbf{x} | \mathbf{x}^{(k-1)})$.

These two steps are executed in an alternating fashion, with each iteration aimed at reducing or at least not increasing the value of the objective function $F(\mathbf{x})$, ensuring $F(\mathbf{x}^{(k)}) \leq F(\mathbf{x}^{(k-1)})$. Instead of minimizing $F(\mathbf{x})$ directly, the algorithm focuses on sequentially solving a series of simple optimization problem. The process begins with a feasible initial point $\mathbf{x}^{(0)}$ and continues iteratively through $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots$ until a specified convergence criterion is satisfied.

We now introduce the multistage convex relaxation algorithm based on the MM approach to find the stationary solution of (5). The multistage convex relaxation algorithm is described in Algorithm 1, where we set $\Sigma^{(0)} = \mathbf{I}$ as a trivial start. This algorithm, built on a sequential optimization framework, draws theoretical support from the analysis provided by Zhang et al. [22]. To simplify, we define

$$f(\Sigma) = \frac{1}{2m} \|\mathbf{y} - \mathcal{A}(\Sigma)\|_2^2 - \tau \log \det \Sigma.$$

Throughout the stages of the algorithm, a weighted ℓ_1 -norm serves as the surrogate function for $\sum_{i,j} p_\lambda(\Sigma_{ij})$. Specifically, for each $1 \leq k \leq K$, we minimize the subsequent convex relaxation subproblems sequentially:

$$F(\Sigma) = f(\Sigma) + \sum_{i,j} p'_\lambda \left(\left| \tilde{\Sigma}_{ij}^{(k-1)} \right| \right) |\Sigma_{ij}|, \quad (8)$$

where $\tilde{\Sigma}^{(k)}$ represents the optimal solution to the k -th subproblem. Each subproblem can further be reformulated as follows:

$$\min_{\Sigma \succ 0} \{f(\Sigma) + \|\Lambda \odot \Sigma\|_1\}, \quad (9)$$

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

$$\nabla f(\hat{\Sigma}) + \Lambda \odot \hat{\Sigma} = 0,$$

where $\hat{\Sigma} \in \partial \|\hat{\Sigma}\|_1$ and the gradient

$$\nabla f(\Sigma) = -\frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma)) - \tau \Sigma^{-1}.$$

Here, $\mathcal{A}^*(\cdot)$ is the conjugate operator of $\mathcal{A}(\cdot)$. Given that the problem (9) lacks an analytical solution, we aim for a suboptimal solution defined under a prespecified tolerance level ε , and terminate the iterations when the approximate KKT condition are satisfied.

Definition 6. For a specified tolerance level ε , the solution $\tilde{\Sigma}^{(k)}$ is deemed ε -optimal for the k -th subproblem 9 if the condition $\omega_\Lambda \left(\tilde{\Sigma}^{(k)} \right) \leq \varepsilon$ is satisfied, where

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

In order to guarantee that the solution $\tilde{\Sigma}^{(k)}$ derived at each stage attains sufficient precision for all $k \geq 1$, we present the following assumption, which is essential for the convergence analysis of multistage convex relaxation.

Assumption 7. For each iteration of the convex relaxation process (9) for all $k \geq 1$, we set the accuracy parameter ε as follows:

$$\varepsilon = \frac{c_2}{\sqrt{mn}} \leq \frac{\lambda}{8},$$

where c_2 is a predetermined small constant.

B. Proximal Newton Algorithm

At each stage, the subproblem is obviously convex, enabling the application of a variety of solution techniques. However, exact analytical solutions are unattainable, iterative algorithms are required. Among these, first-order methods such as the proximal gradient algorithm and second-order method like the proximal Newton algorithm are particularly noteworthy. In this study, we adopt the proximal Newton method to enhance computational efficiency. The proximal Newton algorithm operates by solving a proximal subproblem in each iteration, which involves updating the current estimate using both gradient and Hessian information. Specifically, the algorithm computes a Newton direction to serve as the descent direction, followed by determining an appropriate step size that

ensures a sufficient decrease in the objective function, thereby facilitating convergence.

For the sake of notational simplicity, we denote the iteration index within the k -th stage as t and omit the stage index k . We build a quadratic approximation by considering the second-order Taylor expansion of $f(\Sigma)$:

$$\begin{aligned} \bar{f}(\Sigma | \Sigma_t) &= f(\Sigma_t) + \langle \nabla f(\Sigma_t), \Sigma - \Sigma_t \rangle \\ &\quad + \frac{1}{2} \|\Sigma - \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2. \end{aligned} \quad (10)$$

We then solve

$$\Sigma_{t+\frac{1}{2}} = \underset{\Sigma \succ 0}{\operatorname{argmin}} \bar{F}(\Sigma | \Sigma_t, \Lambda), \quad (11)$$

where

$$\bar{F}(\Sigma | \Sigma_t, \Lambda) = \bar{f}(\Sigma | \Sigma_t) + \|\Lambda \odot \Sigma\|_1.$$

Denote

$$\Delta \Sigma_t = \Sigma_{t+\frac{1}{2}} - \Sigma_t$$

as the Newton direction for the function $F(\Sigma)$. Particularly, we adopt coordinate descent algorithms combined with active set strategy to get the Newton direction, which can improve the time complexity. The detailed computation will be introduced in the Appendix B.

An additional backtracking line searching procedure is included to guarantee the descent of the objective value. We try to find a step size $\beta \in (0, 1]$ such that the Armijo condition [43] holds. Specially, we start with a fixed constant $\mu \in (0.5, 1)$ and update $\beta = \mu^q$ from $q = 0$ with a constant decrease rate until we find the smallest $q \in \mathbb{N}$ such that

$$F(\Sigma_t + \beta \Delta \Sigma_t) \leq F(\Sigma_t) + \alpha \beta \delta_t,$$

where $\alpha \in (0, 0.5)$ and

$$\begin{aligned} \delta_t &= \langle \nabla f(\Sigma_t), \Delta \Sigma_t \rangle - \|\Lambda \odot \Sigma_t\|_1 \\ &\quad + \|\Lambda \odot (\Sigma_t + \Delta \Sigma_t)\|_1. \end{aligned} \quad (12)$$

Then Σ_{t+1} is set as $\Sigma_{t+1} = \Sigma_t + \beta \Delta \Sigma_t$. The whole proximal Newton algorithm is summarized in Algorithm 2.

IV. STATISTICAL AND COMPUTATIONAL THEORIES

In this section, we first introduce some technical assumptions and important definitions. 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.

A. Assumptions

We present several assumptions pertinent to the true covariance matrix, denoted as Σ^* . We define the support set of Σ^* as $\mathcal{S}^* = \{(j, k) \mid \Sigma_{jk}^* \neq 0\}$, with s^* representing its size, i.e., $s^* = |\mathcal{S}^*|$. These assumptions serve as essential parameters for the examination of the sparsity and structural characteristics of the true covariance matrix within the framework of our proposed statistical model.

Algorithm 2: Proximal Newton Algorithm (ProxNewton) With Back-tracking Line Search for Solving (9).

Input: $\tilde{\Sigma}^{(k-1)}, \Lambda^{(k-1)}, \varepsilon$;
1 Initialize $t = 0, \Sigma_t^{(k)} = \tilde{\Sigma}^{(k-1)}, \mu = 0.8, \alpha = 0.3$
2 repeat
3 $\Sigma_{t+\frac{1}{2}}^{(k)} = \underset{\Sigma \succ 0}{\operatorname{argmin}} \bar{F}(\Sigma | \Sigma_t^{(k)}, \Lambda^{(k-1)})$
4 $\Delta \Sigma_t^{(k)} = \Sigma_{t+\frac{1}{2}}^{(k)} - \Sigma_t^{(k)}$
5 $\delta_t = \langle \nabla f(\Sigma_t^{(k)}), \Delta \Sigma_t^{(k)} \rangle - \|\Lambda^{(k-1)} \odot \Sigma_t^{(k)}\|_1$
6 $\quad + \|\Lambda^{(k-1)} \odot (\Sigma_t^{(k)} + \Delta \Sigma_t^{(k)})\|_1$
7 $\beta = 1, q = 0$
8 repeat
9 $\beta = \mu^q$
10 $q = q + 1$
11 if $\Sigma_t^{(k)} + \beta \Delta \Sigma_t^{(k)} \preceq 0$ **then**
12 \quad continue;
13 **end**
14 until $F(\Sigma_t^{(k)} + \beta \Delta \Sigma_t^{(k)}) \leq F(\Sigma_t^{(k)}) + \alpha \beta \delta_t$;
15 $\Sigma_{t+1}^{(k)} = \Sigma_t^{(k)} + \beta \Delta \Sigma_t^{(k)}$
16 $t = t + 1$
17 until $\omega_{\Lambda^{(k-1)}}(\Sigma_t^{(k)}) \leq \varepsilon$;
Output: $\tilde{\Sigma}^{(k)} = \Sigma_t^{(k)}$

Assumption 8. For the ground truth covariance matrix Σ^* , there exists finite upper and positive lower bounds on its eigenvalue. Specifically, there exists a constant $\kappa \geq 1$ satisfies

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

Assumption 8 is a common premise in the literature concerning of the estimation of sparse covariance matrices [9], [15], [44], [45]. This assumption offers several advantages, and readers are encouraged to consult [15] for further details.

Assumption 9. Given the true covariance matrix Σ^* , there exist universal constants α, c_3 such that

$$\|\Sigma_{\mathcal{S}^*}^*\|_{\min} = \min_{(i,j) \in \mathcal{S}^*} |\Sigma_{ij}^*| \geq (\alpha + c_3) \lambda,$$

where α is a constant introduced in Assumption 5, and $c_3 \in (0, \alpha)$ satisfies $p'_\lambda(c_3 \lambda) \geq \frac{\lambda}{2}$.

Additionally, it is essential to select λ appropriately to ensure that the regularization is sufficiently large to eliminate irrelevant coordinates, thereby yielding the solution is adequately sparse.

Assumption 10. There exists a generic constant c_4 such that

$$\lambda = c_4 \sqrt{\frac{\log d}{m}} \geq 4(\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon).$$

Assumption 9 is referred to the minimum signal strength condition, commonly applied in the analysis of non-convex penalized regression problems [28], [41], [42]. Assumption 10 sets the tuning parameter λ in the order of $\sqrt{\frac{\log d}{m}}$. It prevents

λ from becoming overly large as the number of measurements m increases, thereby ensuring that the estimators maintain a close correspondence with the true model parameters. Taking these assumptions into account, we can obtain the oracle rate of the convergence.

Definition 11. Define a local region around Σ^* as

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

In our analysis, we set the radius r as $\frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa}$. Under Assumption 10, we prove that the solution produced by our algorithm will fall within the local region $\mathcal{B}(\Sigma^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa})$ after a finite number of iterations.

B. Statistical Guarantees and Consequences

The following will present the main result, which demonstrates the contraction property of the solution path $\left\{\tilde{\Sigma}^{(k)}\right\}_{k \geq 1}$.

Theorem 12 (Contraction Property). *Consider the estimator in 5 and suppose the Assumptions 1, 4, 8, and 9 hold. Then with probability exceeding $1 - \exp(-\epsilon^2/2)$ for some $\epsilon > 0$, the ε -optimal solution $\tilde{\Sigma}^{(k)}$ is bounded by:*

$$\begin{aligned} \left\|\tilde{\Sigma}^{(k)} - \Sigma^*\right\|_F &\leq \frac{1}{\rho_{2s^*+2\tilde{s}}^-} \left(\underbrace{\|(\nabla f(\Sigma^*))_{S^*}\|_F}_{\text{oracle rate}} + \underbrace{\varepsilon\sqrt{s^*}}_{\text{optimization error}} \right) \\ &\quad + \underbrace{\delta \left\|\tilde{\Sigma}^{(k-1)} - \Sigma^*\right\|_F}_{\text{contraction}}, \end{aligned} \quad (13)$$

for $1 \leq k \leq K$, where $\delta \in (0, 1)$ is the contraction factor, provided that $m = \mathcal{O}\left((s^* + \tilde{s}) \log^{\frac{3}{2}} d\right)$.

Remark 13. The oracle estimator $\hat{\Sigma}^O$ is defined with prior knowledge of the true support set S^* , and is given by

$$\begin{aligned} \hat{\Sigma}^O &= \arg \min_{\Sigma} f(\Sigma) \\ \text{s.t.} \quad &\Sigma_{\bar{S}^*} = 0. \end{aligned}$$

The statistical convergence rate of this oracle estimator, commonly referred to as the oracle rate, provides a benchmark for evaluating the theoretical upper limit of estimator performance when the true support set is known. Specifically, the distance between $\hat{\Sigma}^O$ and Σ^* is characterized by $\left\|\hat{\Sigma}^O - \Sigma^*\right\|_F \lesssim \|(\nabla f(\Sigma^*))_{S^*}\|_F$. This inequality indicates that the oracle estimator can theoretically approximate the true parameter matrix Σ^* very closely when the true support sets are known.

Theorem 12 elaborates the estimation discrepancy between the ε -optimal solution $\tilde{\Sigma}^{(k)}$ and the ground truth Σ^* is constrained by three primary factors: the oracle rate, the optimization error, and a contraction term. The oracle rate represents the error bound achievable under ideal conditions using the optimal strategy; the optimization error reflects the performance gap between the actual optimization algorithm

and the optimal strategy; the contraction term indicates the rate at which the error converges to its minimum value as the number of iterations increases. The ensuing results specify the exact statistical convergence rate.

Corollary 14. *Let \mathbf{x} be a sub-Gaussian random vector with zero mean and covariance Σ^* and $\{\mathbf{x}_i\}_{i=1}^n$ be a collection of independent and identically distributed (i.i.d) samples from \mathbf{x} . Suppose that Assumptions 1, 4, 7, 8, and 9 hold, that is, $\lambda \asymp \sqrt{\frac{\log d}{mn}}$ and $\varepsilon \lesssim \sqrt{\frac{1}{mn}}$. If $\tau \lesssim \sqrt{\frac{1}{mn}} \|(\Sigma^*)^{-1}\|_{\max}^{-1}$, then the ε -optimal solution $\tilde{\Sigma}^{(1)}$ satisfies*

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

with high probability (w.h.p.).

Corollary 15. *Let \mathbf{x} be a sub-Gaussian random vector with zero mean and covariance Σ^* and $\{\mathbf{x}_i\}_{i=1}^n$ be a collection of independent and identically distributed (i.i.d) samples from \mathbf{x} . Suppose that Assumptions 1, 4, 7, 8, and 9 hold, that is, $\lambda \asymp \sqrt{\frac{\log d}{mn}}$ and $\varepsilon \lesssim \sqrt{\frac{1}{mn}}$. If $\tau \lesssim \sqrt{\frac{1}{mn}} \|(\Sigma^*)^{-1}\|_{\max}^{-1}$, and $K \geq \log(\lambda\sqrt{mn}) \gtrsim \log \log d$, then the ε -optimal solution $\tilde{\Sigma}^{(K)}$ satisfies*

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

Corollary 14 and Corollary 15 are direct consequence of Theorem 12. The former addresses the scenario where $k = 1$ and describes a contraction property resulting from the MM-based multistage convex relaxation algorithm. It demonstrates that to achieve the oracle rate, the optimization error ε must be chosen such that $\varepsilon \leq \frac{\|(\nabla f(\Sigma^*))_{S^*}\|_F}{\sqrt{s^*}}$, and the parameter K must be sufficiently large. The latter implies that, with minimal assumptions, solving no more than approximately $\log \log d$ convex problems is enough to achieve the oracle rate $\sqrt{\frac{s^*}{mn}}$.

C. Computational Theory

Lemma 16. *Under Assumptions 4 and 8, $f(\Sigma)$ is $\left(\rho_{2s^*+2\tilde{s}}^- + \frac{16\tau^3}{\kappa^2(4\tau+\rho_{2s^*+2\tilde{s}}^-)^2}\right)$ -strongly convex and $\left(\rho_{s^*+2\tilde{s}}^+ + \frac{16\tau^3\kappa^2}{(4\tau-\rho_{2s^*+2\tilde{s}}^-)^2}\right)$ -smooth in the region $\mathcal{B}\left(\Sigma^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa}\right)$, where*

- κ is defined in Assumption 8
- $\rho_{2s^*+2\tilde{s}}^-, \rho_{2s^*+2\tilde{s}}^+$ is defined in Assumption 4
- τ is the positive tuning hyper-parameter.

Lemma 16 ensures that within the specified region, the function exhibits strong convexity and smoothness properties, which are essential for stability and convergence analysis. Next, we present the explicit iteration complexity of the proposed algorithm. We begin by characterizing the convergence for the first stage.

Theorem 17. For $k = 1$, suppose Assumption 1, 4, 7, 8, and 9 hold. After a sufficient number of iterations $T < \infty$, for all $t \geq T$, the following conditions are satisfied: $\|\Sigma_t^{(1)} - \Sigma^*\|_F \leq r$, $\left\|(\Sigma_t^{(1)})_{\bar{S}^*}\right\|_0 \leq \tilde{s}$ and

$$\left\|\Sigma_{t+1}^{(1)} - \hat{\Sigma}^{(1)}\right\|_F \leq \frac{\tau\kappa^3}{\rho_{2s^*+2\tilde{s}}^-} \left\|\Sigma_t^{(1)} - \hat{\Sigma}^{(1)}\right\|_F^2.$$

Moreover, we need at most $T + \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon}\right)$ iterations to terminate the proximal Newton algorithm.

Theorem 18. For $k \geq 2$, suppose 1, 4, 7, 8, and 9 hold. For all iterations t , we have $\|\Sigma_t^{(k)} - \Sigma^*\|_F \leq r$, $\left\|(\Sigma_t^{(k)})_{\bar{S}^*}\right\|_0 \leq \tilde{s}$, which guarantee $\beta = 1$, and

$$\left\|\Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)}\right\|_F \leq \frac{\tau\kappa^3}{\rho_{2s^*+2\tilde{s}}^-} \left\|\Sigma_t^{(k)} - \hat{\Sigma}^{(k)}\right\|_F^2.$$

Moreover, the proximal Newton algorithm requires at most $\log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon}\right)$ iterations to terminate.

Theorem 17 demonstrates that the solution enters the ball $\mathcal{B}\left(\Sigma^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa}\right)$ following the initial stage. To satisfy the approximate KKT conditions, the number of iterations is capped at $\mathcal{O}\left(T + \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon}\right)\right)$. The resulting solution $\hat{\Sigma}^{(1)}$ exhibits notable qualities, laying the foundation for improved computational efficiency in the subsequent stages of our proximal Newton algorithm. Theorem 18 states that the sparsity of the solution is preserved throughout the iterations, and the solution remains within $\mathcal{B}\left(\Sigma^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa}\right)$. Due to the sparsity of the solution, the algorithm exhibits quadratic convergence. Moreover, the algorithm achieves a logarithmic number of iterations, $\mathcal{O}\left(\log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon}\right)\right)$, to satisfy the approximate KKT condition.

V. NUMERICAL EXPERIMENTS

In this section, we examine the practical recovery performance of proposed estimator and provide numerical results for the proposed algorithm for sparse covariance matrix sensing. Section V-A describe the generation of synthetic data and the measurement criteria used in the paper. In Section V-B and Section V-C, the statistical theoretical results of Section IV-B and the computational theoretical results of Section IV-C are demonstrated, respectively. Section V-D compares the relative estimation error of our proposed estimator and existing ℓ_1 estimator. The non-convex penalty function chosen for these experiments is the MCP, with a constant setting of $b = 2$ across all trials. The selection of the tuning parameters λ and τ is determined through the application of five-fold cross-validation. All methods are implemented in MATLAB, and run on an Intel i7-10700 2.90 GHz $\times 8$ with 16 GB of RAM. All results are averaged on 100 Monte Carlo realizations.

A. Synthetic Datasets and Measurement Criteria

Synthetic data for this section is generated by the following methods. For different fixed d , n independent data points are generated by drawing i.i.d. samples from a $\mathcal{N}(0, \Sigma^*)$ distribution. In this paper, we consider the following three typical covariance models and two randomly generated sparse covariance models.

- Banded matrix:

$$\Sigma_{ij}^* = \begin{cases} 1 - \frac{|i-j|}{10}, & |i-j| \leq 10, \\ 0, & \text{otherwise.} \end{cases}$$

- Block matrix: The indices $1, 2, \dots, d$ are partitioned into 10 ordered groups of equal size with $\Sigma_{ij}^* = \begin{cases} 1, & i = j, \\ 0.5, & i \text{ and } j (i \neq j) \text{ belong to the same group,} \\ 0, & \text{otherwise.} \end{cases}$
- Toeplitz matrix: $\Sigma_{ij}^* = 0.75^{|i-j|}$.
- Sprandsym matrix: Σ^* is generated utilizing the “sprandsym” built-in function in MATLAB with a sparsity value s^* .
- Probability matrix: Σ^* is generated with off-diagonal entries drawn i.i.d. from a uniform $(-1, 1)$ distribution. These entries are set to zero with a fixed probability \mathbf{P} . Finally, a multiple of the identity was added to the resulting matrix.

The first two models are sparse, while the third one is approximately sparse. The fourth model is positive, symmetric, and sparse when s^* is small. The fifth model is positive definite, well-conditioned, and can be sparse when \mathbf{P} is appropriate (in this case, either $\mathbf{P} = 0.97$ or $\mathbf{P} = 0.85$ to simulate a very sparse and a somewhat sparse model).

The estimation performance is measured by the absolute error and relative error under both the Frobenius norm and the spectral norm, respectively. Specifically, the Frobenius absolute error (FAE) is defined as $\|\hat{\Sigma} - \Sigma^*\|_F$, where $\hat{\Sigma}^2$ is the estimated covariance matrix and Σ^* is the true value. Similarly, the Frobenius relative error (FRE) is defined as $\frac{\|\hat{\Sigma} - \Sigma^*\|_F}{\|\Sigma^*\|_F}$. The selection performance is examined by the false positive rate (FPR) and the true positive rate (TPR), defined as, respectively

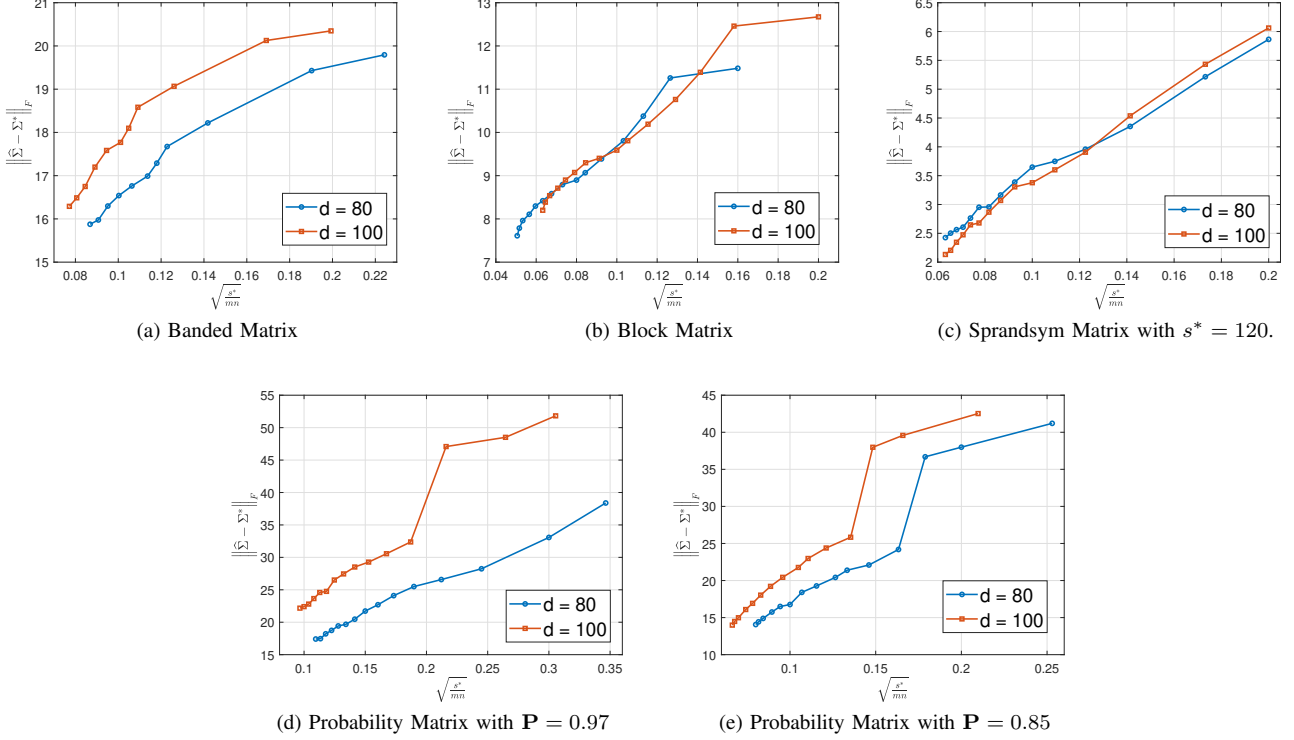
$$\text{FPR} = \frac{\#\{(i, j) : \hat{\Sigma}_{ij} \neq 0 \ \& \ \Sigma_{ij}^* = 0\}}{\#\{(i, j) : \Sigma_{ij}^* = 0\}},$$

$$\text{TPR} = \frac{\#\{(i, j) : \hat{\Sigma}_{ij} \neq 0 \ \& \ \Sigma_{ij}^* \neq 0\}}{\#\{(i, j) : \Sigma_{ij}^* \neq 0\}}.$$

B. Demonstration of Oracle Rate

The oracle rate derived for proposed estimator and algorithm in Section IV-B is shown to be dependent on the sparsity level s^* , the number of measurements m , and number of samples n .

²In the whole Section V and Appendix E, we use $\hat{\Sigma}$ to denote a generic covariance estimator, which can be the estimation results obtained by various algorithms and estimators.



*The Toeplitz model is not reported due to the fact that it is not sparse but approximately sparse.

Fig. 1. The oracle rate for different dimension of different covariance model.

To demonstrate this theoretical result, we apply our proposed algorithm to synthetic datasets, as describe in Section V-A. Fig. 1 illustrates the estimation performance across various matrix configurations differentiated by dimensions d , sparsity level s^* , and measurement count m . The x -axis is $\sqrt{\frac{s^*}{mn}}$, and the y -axis represents FAE. Two distinct lines depicts the FAE for dimensions $d = \{80, 100\}$. We observe that an increase in measurement size m and sample count n corresponds to a reduction in FAE, visually demonstrating how the estimation quality of the covariance matrix varies with the oracle rate. According to Corollary 15, the Frobenius norm difference $\|\hat{\Sigma} - \Sigma^*\|_F$ converges to the statistical errors as the stage index k increases, following the order of $\mathcal{O}\left(\sqrt{\frac{s^*}{mn}}\right)$. We can see that the estimation errors grow approximately linearly with the theoretical rate, which validates our theoretical guarantee.

C. Demonstration of Computational Rate

Fig. 2 illustrates a detailed analysis of the convergence rates observed during distinct phases of the computational process. In the initial phase, the algorithm display a sub-linear convergence, for that the optimization problem is characterized by a lack of strong convexity. As the computation progresses and the solution approaches the sparse region, a notable transformation occurs in the dynamics of convergence. In this region, the characteristics of the feasible solutions change markedly—they become increasingly sparse. Concurrently, the objective function undergoes a transformation into

an essentially low-dimensional form. This new form is not only sparse but also exhibits properties of strong convexity and smoothness, which are crucial for faster convergence. Consequently, our algorithm exhibits a linear convergence rate for $k \geq 2$. Interestingly, this property extends even to the case of $k = 1$, wherein the convergence rate remains sub-linear until the algorithm transitions into the contraction region, at which point it achieves a quadratic rate of convergence.

D. Simulation Experiment Comparisons

We conduct a series of Monte Carlo trials to evaluate the performance of our proposed covariance matrix sensing algorithm using “Sprandsym matrix model” with dimension $d = 100$ and $n = 50$. The efficacy of the algorithm is assessed under varying sparsity levels and noise conditions, and comparisons are made against existing estimators employing different penalization techniques.

Initially, we examine the FRE of the estimated covariance matrices over various sparsity levels, represented by s^* values of 120, 180, 240 and 300. The results, obtained in the absence of noise and while sensing the true covariance Σ^* with sample count n fixed at 50, are depicted in Figure 3a. This figure illustrates the inverse relationship between FRE and the number of measurements m , where an increase in m correlates with enhanced estimation accuracy. Concurrently, within the same number of measurements, a lower s^* value, which corresponds to higher sparsity, results in a reduced FRE, highlighting the influence of sparsity on estimation precision.

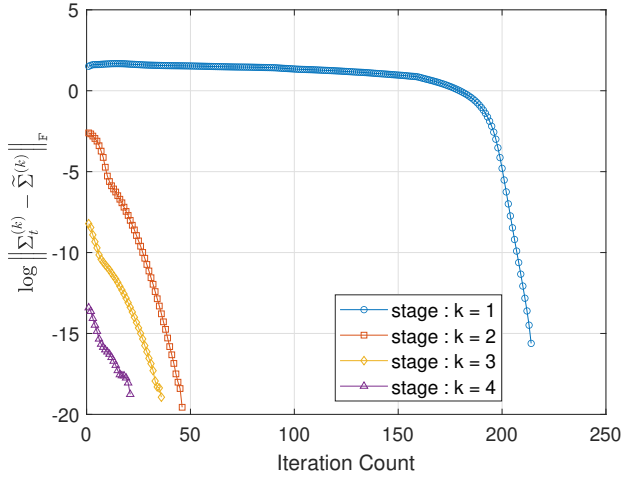


Fig. 2. Computational rate of convergence in each stage for simulation experiment. The x -axis is the iteration count t within the k -th subproblem. A distinct phase transition in the algorithmic convergence rate from sub-linear to super-linear is observable as the iterations progress into the sparse region.

To reflect more realistic conditions, additive noise is introduced into each measurement, where the noise η_i is generated from a sub-exponential distributed random variable scaled by γ , specifically $\gamma \cdot \mathcal{M}(0, 1)$. The influence of this noise, particularly at a level of $\gamma = 10^{-1}$, on the average relative error is illustrated in Figure 3b. These results demonstrate the resilience of our estimator against external noise disturbances, with the trend of FRE across various measurements and sparsity levels remaining consistent even in noisy environments. Concurrently, we examine the variations in FRE relative to sample number as depicted in Figure 3c, while maintaining a constant measurement count of $m = 300$.

Morover, we explore the performance of our proposed algorithm under different noise intensities for the sparsity level $s^* = 300$. The impact of these noise intensities on the accuracy of the estimator is showcased in Figure 4, offering a comprehensive view of how noise intensity affects the reliability of the estimation process. The figure indicates that at lower numbers of measurements, the noise intensity marginally affects the accuracy of covariance matrix recovery. However, as the number of measurements increases, the influence of noise diminishes, leading to more precise recovery outcomes. This trend underscores the importance of sufficient data in mitigating the effects of noise and enhancing the reliability of the estimation process.

Lastly, we conduct a comparative analysis between our proposed covariance matrix sensing estimator and an existing estimator that incorporates an ℓ_1 penalty. Each test scenario incorporated a consistent level of additive noise ($\gamma = 10^{-1}$), with the comparative results depicted in Figure 5. The comparison reveals that our MCP-based estimator consistently outperforms the ℓ_1 -penalty-based estimator in terms of the Frobenius norm error. This supports our theoretical proposition that using a nonconvex penalty can significantly reduce the error in covariance matrix sensing. For rigorous statistical validation of these results, we have detailed the average

metrics and corresponding standard errors (in parentheses) in the Appendix E.

VI. DISCUSSION

A. Bilinear Rank-one Measurement Model

Our research findings provides evidence that covariance matrices can be effectively reconstructed from a limited set of quadratic measurements, contingent upon the sparsity of the underlying structure. This reconstruction process is particularly advantageous in terms of storage efficiency when the sensing vectors are i.i.d. from sub-Gaussian distribution, and when the number of measurements exceeds the basic sampling theoretical limit. Additionally, our analysis naturally extends to the bilinear rank-one measurement model, where measurements take the form $y_i = \mathbf{a}_i^\top \Sigma \mathbf{b}_i$, with the sensing vectors \mathbf{a}_i and \mathbf{b}_i being independently generated. This extension substantiates our results within an asymmetric sensing framework, which has significant implications in various fields, including communication [46], imaging [47], and machine learning [48].

The bilinear model introduces additional complexity due to the interaction between two independent sensing vectors, \mathbf{a}_i and \mathbf{b}_i . Despite this, it is straightforward to demonstrate that this configuration meets the Assumption 3, a critical requirement for successful recovery. The derivation process closely parallels that of the quadratic model, with necessary modifications to accommodate the bilinear structure. Importantly, we also establish that within the context of the bilinear model, the Hessian matrix retains its symmetric properties, thereby enabling our proposed algorithm to sustain both efficiency and effectiveness. As a result, the bilinear measurement model inherits the oracle property, further enhancing the robustness and extensive applicability of our approach.

B. Covariance Sketching for variable with non-zero mean

In the preceding discussion, we consider the scenario where the mean signal vector is assumed to be zero. However, this assumption frequently does not hold true in numerous practical contexts. In case where the mean vector is not known, one might first utilize robust methodologies to estimate the mean vector, and subsequently apply this estimation to perform a robust calculation of the covariance. This bifurcated approach to estimating the mean and covariance not only introduces additional tuning parameters but also increases both statistical variability and computational complexity.

To streamline this process, we propose an end-to-end methodology. We recommend a pairwise difference approach that obviates the necessity for direct mean estimation. Specifically, we define $N := n(n-1)/2$ and consider the pairwise differences

$$\{\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_N\} = \{\mathbf{x}_1 - \mathbf{x}_2, \mathbf{x}_2 - \mathbf{x}_3, \dots, \mathbf{x}_{n-1} - \mathbf{x}_n\},$$

which are identically distributed with zero mean. The covariance of these differences is characterized by $\text{cov}(\tilde{\mathbf{x}}) = 2\Sigma$. Furthermore, the sample covariance matrix can be expressed as a U-statistic. Normally, the sample covariance is calculated as $\Sigma_{\text{sam}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^\top$ where $\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i$. However, using the pairwise differences, it can be reformulated

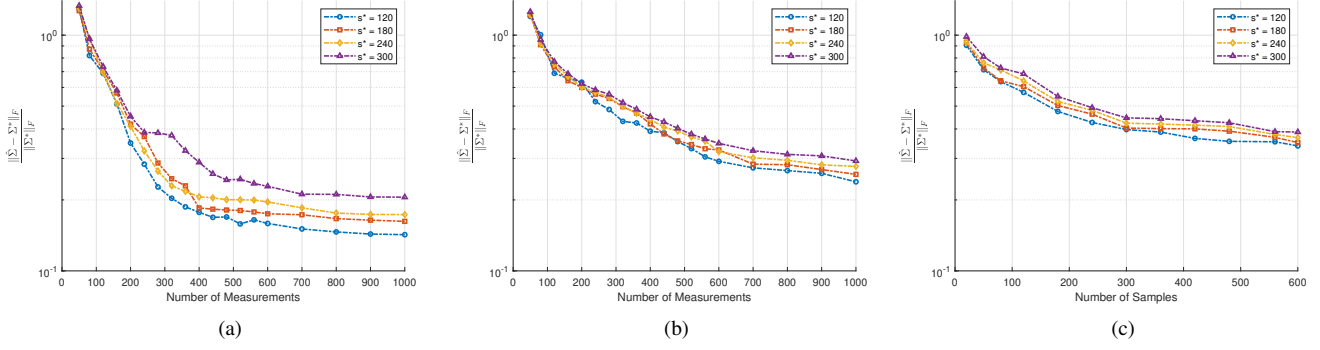


Fig. 3. The FRE of the estimated covariance matrices is examined in two contexts: (a) the true covariance under varying levels of sparsity in the absence of noise; (b) the sample covariance under different sparsity levels when subjected to a noise parameter of $\gamma = 0.1$ and $n = 50$; (c) the sample covariance under different sparsity levels when subjected to a noise parameter of $\gamma = 0.1$ and $m = 300$;

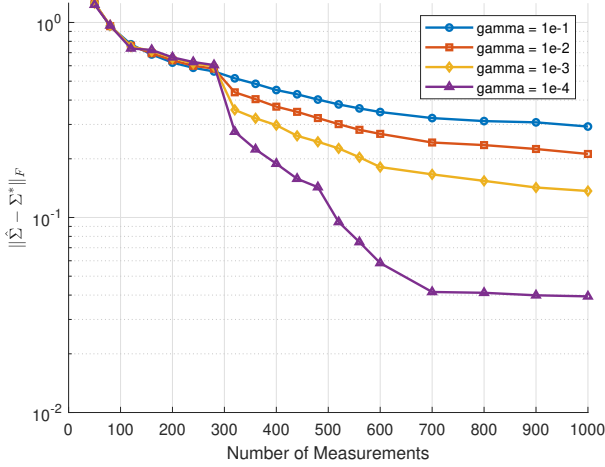


Fig. 4. The FRE of the estimated covariance matrices for different noise levels when $s^* = 300$.

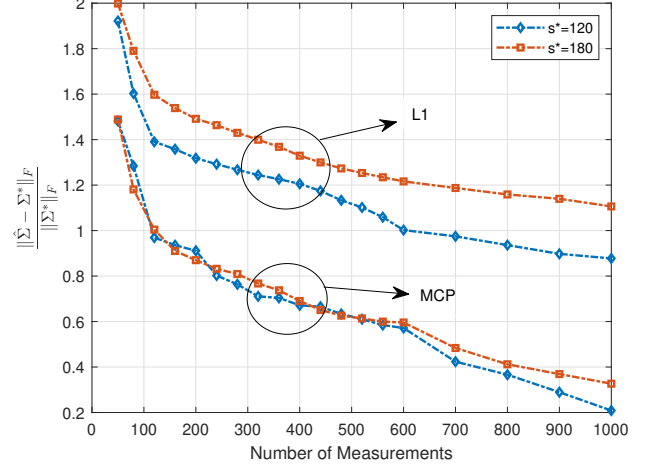


Fig. 5. The FRE of the estimated covariance matrices for different sparsity level with noise (ℓ_1 v.s. MCP).

equivalently as $\Sigma_{\text{sam}} = \frac{1}{N} \sum_{i=1}^N \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top / 2$, thereby simplifying the computation by avoiding the direct estimation of the mean. This reformulation provides a more streamlined approach to covariance estimation, reducing both the computational load and the potential statistical errors associated with separate mean and covariance estimations.

VII. CONCLUSION

In this paper, we present a comprehensive study on covariance matrix sensing, focusing on the high-dimensional sparse covariance matrices estimation through the quadratic measurements model. This model is particularly suited for scenarios with constrained processing capabilities and limited memory, such as real-time data acquisition devices. It encompasses various sampling strategies, primarily those that capture magnitude or energy samples. Our findings demonstrate that sparse covariance matrices can be effectively and accurately reconstructed using a minimal set of quadratic measurements with bottommost storage requirements. We provide rigorous theoretical backing and empirical evidence to support these claims. Notably, our proposed estimators exhibit superior

statistical convergence rates compared to existing methods, highlighting their practical utility and efficiency in real-world applications.

REFERENCES

- [1] M. Pourahmadi, *High-dimensional covariance estimation: with high-dimensional data*. John Wiley & Sons, 2013, vol. 882.
- [2] A. Bose and M. Bhattacharjee, *Large covariance and autocovariance matrices*. Chapman and Hall/CRC, 2018.
- [3] H. Tsukuma and T. Kubokawa, *Shrinkage estimation for mean and covariance matrices*. Springer, 2020.
- [4] H. Abdi and L. J. Williams, "Principal component analysis," *Wiley Interdisciplinary Reviews: Computational Statistics*, vol. 2, no. 4, pp. 433–459, 2010.
- [5] R. L. Gorsuch, *Factor analysis: Classic edition*. Routledge, 2014.
- [6] R. A. Monzingo and T. W. Miller, *Introduction to adaptive arrays*. Scitech publishing, 2004.
- [7] L. K. Chan, J. Karceski, and J. Lakonishok, "On portfolio optimization: Forecasting covariances and choosing the risk model," *The review of Financial studies*, vol. 12, no. 5, pp. 937–974, 1999.
- [8] J. Fan, Y. Liao, and H. Liu, "An overview of the estimation of large covariance and precision matrices," *The Econometrics Journal*, vol. 19, no. 1, pp. C1–C32, 2016.
- [9] C. Lam and J. Fan, "Sparsistency and rates of convergence in large covariance matrix estimation," *Annals of Statistics*, vol. 37, no. 6B, p. 4254, 2009.
- [10] P. J. Bickel and E. Levina, "Covariance regularization by thresholding," *The Annals of Statistics*, vol. 36, no. 6, pp. 2577 – 2604, 2008.
- [11] N. El Karoui, "Operator norm consistent estimation of large-dimensional sparse covariance matrices," 2008.
- [12] A. J. Rothman, E. Levina, and J. Zhu, "Generalized thresholding of large covariance matrices," *Journal of the American Statistical Association*, vol. 104, no. 485, pp. 177–186, 2009.
- [13] T. T. Cai and H. H. Zhou, "Minimax estimation of large covariance matrices under ℓ_1 -norm," *Statistica Sinica*, pp. 1319–1349, 2012.
- [14] C. T. Tony and Z. H. H., "Optimal rates of convergence for sparse covariance matrix estimation," *The Annals of Statistics*, vol. 40, no. 5, pp. 2389–2420, 2012.
- [15] A. J. Rothman, "Positive definite estimators of large covariance matrices," *Biometrika*, vol. 99, no. 3, pp. 733–740, 2012.
- [16] D. L. Donoho, "Compressed sensing," *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.
- [17] E. J. Candes and T. Tao, "Near-optimal signal recovery from random projections: Universal encoding strategies?" *IEEE transactions on Information Theory*, vol. 52, no. 12, pp. 5406–5425, 2006.
- [18] Y. C. Eldar and G. Kutyniok, *Compressed sensing: theory and applications*. Cambridge university press, 2012.
- [19] 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.
- [20] E. J. Candes, T. Strohmer, and V. Voroninski, "Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming," *Communications on Pure and Applied Mathematics*, vol. 66, no. 8, pp. 1241–1274, 2013.
- [21] S. A. van de Geer and P. Bühlmann, "On the conditions used to prove oracle results for the Lasso," *Electronic Journal of Statistics*, vol. 3, no. none, pp. 1360 – 1392, 2009.
- [22] T. Zhang, "Analysis of multi-stage convex relaxation for sparse regularization," *Journal of Machine Learning Research*, vol. 11, no. 3, 2010.
- [23] G. Dasarthy, 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.
- [24] D. Romero, D. D. Ariananda, Z. Tian, and G. Leus, "Compressive Covariance Sensing: Structure-based compressive sensing beyond sparsity," *IEEE Signal Processing Magazine*, vol. 33, no. 1, pp. 78–93, 2016.
- [25] Z. Tian, Y. Tafesse, 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.
- [26] M. A. Lexa, M. E. Davies, J. S. Thompson, and J. Nikolic, "Compressive power spectral density estimation," in *2011 IEEE international conference on acoustics, speech and signal processing (ICASSP)*. IEEE, 2011, pp. 3884–3887.
- [27] P. J. Bickel, Y. Ritov, and A. B. Tsybakov, "Simultaneous analysis of lasso and dantzig selector," 2009.
- [28] 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.
- [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] G. Raskutti, M. J. Wainwright, and B. Yu, "Restricted eigenvalue properties for correlated gaussian designs," *The Journal of Machine Learning Research*, vol. 11, pp. 2241–2259, 2010.
- [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] P. Bühlmann and S. Van De Geer, *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media, 2011.
- [35] 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.
- [36] D. D. Ariananda and G. Leus, "Compressive wideband power spectrum estimation," *IEEE Transactions on Signal Processing*, vol. 60, no. 9, pp. 4775–4789, 2012.
- [37] 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.
- [38] O. Ledoit and M. Wolf, "Improved estimation of the covariance matrix of stock returns with an application to portfolio selection," *Journal of Empirical Finance*, vol. 10, no. 5, pp. 603–621, 2003.
- [39] R. Krishnamurthi, A. Kumar, D. Gopinathan, A. Nayyar, and B. Qureshi, "An overview of iot sensor data processing, fusion, and analysis techniques," *Sensors*, vol. 20, no. 21, p. 6076, 2020.
- [40] N. González-Prelcic and M. E. Domínguez-Jiménez, "Circular sparse rulers based on co-prime sampling for compressive power spectrum estimation," in *Proceedings of 2014 IEEE Global Communications Conference*. IEEE, 2014, pp. 3044–3050.
- [41] 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.
- [42] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *Annals of Statistics*, vol. 38, pp. 894–942, 2010.
- [43] L. Armijo, "Minimization of functions having Lipschitz continuous first partial derivatives," *Pacific Journal of Mathematics*, vol. 16, no. 1, pp. 1–3, 1966.
- [44] H. Liu, L. Wang, and T. Zhao, "Sparse covariance matrix estimation with eigenvalue constraints," *Journal of Computational and Graphical Statistics*, vol. 23, no. 2, pp. 439–459, 2014.
- [45] 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.
- [46] X. Wang and H. V. Poor, "Blind equalization and multiuser detection in dispersive cdma channels," *IEEE Transactions on Communications*, vol. 46, no. 1, pp. 91–103, 1998.
- [47] P. Campisi and K. Egiazarian, *Blind image deconvolution: theory and applications*. CRC press, 2017.
- [48] H. Bristow, A. Eriksson, and S. Lucey, "Fast convolutional sparse coding," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 391–398.
- [49] I. S. Dhillon, C.-J. Hsieh, M. A. Sustik, and P. Ravikumar, "Sparse inverse covariance matrix estimation using quadratic approximation," in *Symposium on Machine Learning in Speech and Language Processing*, 2011.
- [50] J. Friedman, T. Hastie, and R. Tibshirani, "Regularization paths for generalized linear models via coordinate descent," *Journal of statistical software*, vol. 33, no. 1, p. 1, 2010.
- [51] T. T. Wu and K. L. Lange, "Coordinate descent algorithms for lasso penalized regression," *The Annals of Applied Statistics*, vol. 2, pp. 224–244, 2008.
- [52] Y. Ning, T. Zhao, and A. Liu, "A likelihood ratio framework for high-dimensional semiparametric regression," *Annals of Statistics*, vol. 45, no. 6, pp. 2299–2327, Dec. 2017.

- [53] R. Tibshirani, J. Bien, J. Friedman, T. Hastie, N. Simon, J. Taylor, and R. J. Tibshirani, "Strong Rules for Discarding Predictors in Lasso-Type Problems," *Journal of the Royal Statistical Society Series B: Statistical Methodology*, vol. 74, no. 2, pp. 245–266, 11 2011.
- [54] T. Zhao, H. Liu, and T. Zhang, "Pathwise coordinate optimization for sparse learning: Algorithm and theory," *The Annals of Statistics*, vol. 46, no. 1, pp. 180 – 218, 2018.
- [55] 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.
- [56] 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.
- [57] 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.
- [58] I. E.-H. Yen, C.-J. Hsieh, P. K. Ravikumar, and I. S. Dhillon, "Constant Nullspace Strong Convexity and Fast Convergence of Proximal Methods under High-Dimensional Settings," in *Advances in Neural Information Processing Systems*, Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K. Weinberger, Eds., vol. 27. Curran Associates, Inc., 2014.

APPENDIX A

RESTRICTED EIGENVALUE & SPARSE EIGENVALUE

We first introduce an equivalent definition of definition 3.

Definition 19. For three positive integers s_0, ϑ, r , the largest and smallest localized restricted eigenvalues are defined as

$$\psi_{s_0, \vartheta, r}^+ = \sup \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\Sigma) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid (\mathbf{u}, \Sigma) \in \mathcal{I}(s_0, \vartheta, r) \right\};$$

$$\psi_{s_0, \vartheta, r}^- = \inf \left\{ \frac{\mathbf{u}^\top \nabla^2 f(\Sigma) \mathbf{u}}{\mathbf{u}^\top \mathbf{u}} \mid (\mathbf{u}, \Sigma) \in \mathcal{I}(s_0, \vartheta, r) \right\},$$

where $\mathcal{I}(s_0, \vartheta, r)$ is defined as

$$\left\{ (\mathbf{u}, \Sigma) \mid \begin{array}{l} \mathcal{S}^* \subseteq \mathcal{J}, |\mathcal{J}| \leq s_0, \\ \|\mathbf{u}_{\mathcal{J}}\|_1 \leq \vartheta \|\mathbf{u}_{\mathcal{J}^c}\|_1, \|\Sigma - \Sigma^*\|_F \leq r \end{array} \right\},$$

which represents a local ℓ_1 cone.

We now proceed to introduce the relationships between SE and LRE.

Proposition 20. For any $\Sigma \in \mathcal{I}(s, \vartheta, r) \cap \mathcal{B}(\Sigma^*, r)$, the following inequalities hold:

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

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

where c_1, c_2, c_3 and c_4 are constant.

For a detailed proof of Proposition 20, readers are referred to the foundational work in [34], which provides comprehensive mathematical derivations and discussions. This reference is omitted here for brevity.

APPENDIX B

COMPUTING THE NEWTON DIRECTION WITH ACTIVE SET STRATEGY

In this section, we elaborate on the computation of the Newton Direction within the context of our optimization framework. To simplify the exposition, we omit the stage index

k and concentrate on a specific Newton iteration. The gradient and Hessian for $f(\Sigma)$ are given by

$$\nabla f(\Sigma) = -\frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma)) - \tau \Sigma^{-1},$$

and

$$\nabla^2 f(\Sigma) = \frac{1}{m} \sum_{i=1}^m (\mathbf{A}_i \otimes \mathbf{A}_i) + \tau (\Sigma^{-1} \otimes \Sigma^{-1}).$$

The direct computation of the Hessian matrix is computationally intensive, exhibiting a complexity of $\mathcal{O}(md^4)$, which render it impractical for high-dimensional data. However, the inherent symmetric structure of the Hessian permits certain computational optimizations. By utilizing a coordinate descent approach, it is possible to update all variables in $\mathcal{O}(md^3)$ time. This methodology has demonstrated efficacy in addressing Lasso-type problems [49]–[51] with theoretical justifications [52].

In order to precisely articulate the issue, we reformulate (10) as follows:

$$\begin{aligned} \bar{f}(\Delta) &= f(\Sigma_t) + \text{tr} \left((\mathbf{Q} - \tau \mathbf{W})^\top \Delta \right) \\ &\quad + \frac{1}{2} \left(\frac{1}{m} \sum_{i=1}^m \text{tr}(\mathbf{A}_i \Delta \mathbf{A}_i \Delta) + \tau \text{tr}(\mathbf{W} \Delta \mathbf{W} \Delta) \right), \end{aligned} \quad (14)$$

where $\mathbf{Q} = -\frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma))$ and $\mathbf{W} = \Sigma^{-1}$.

In the process of refining the Newton direction, we represent the existing approximation as \mathbf{D} and the revised direction as \mathbf{D}' . When implementing a coordinate descent update for the variable Σ_{ij} (assuming $i < j$), the update is structured to preserve symmetry as follows:

$$\mathbf{D}' = \mathbf{D} + \mu (e_i e_j^\top + e_j e_i^\top),$$

where μ is the scalar update magnitude, and e_i, e_j are the standard basis vectors. The resolution of the one-dimensional problem associated with the equation for (11) is expressed as follows:

$$\arg \min_{\mu} \bar{f}(\mathbf{D} + \mu (e_i e_j^\top + e_j e_i^\top)) + 2\Lambda_{ij} |\Sigma_{ij} + D_{ij} + \mu|.$$

After substituting $\mathbf{D}' = \mathbf{D} + \mu (e_i e_j^\top + e_j e_i^\top)$ for Δ in (14), the function \bar{f} is refined to focus solely on terms dependent on μ , neglecting others. The contributions to the function from three terms are:

- The linear term $\text{tr}((\mathbf{Q} - \tau \mathbf{W})^\top \Delta)$ contributes $2\mu (Q_{ij} - \tau W_{ij})$,
- The penalty term gives $2\Lambda_{ij} |\Sigma_{ij} + D_{ij} + \mu|$,
- The quadratic term yields

$$\begin{aligned} &\frac{1}{m} \sum_{l=1}^m 4\mu (\mathbf{A}_l)_{\cdot i}^\top \mathbf{D} (\mathbf{A}_l)_{\cdot j} + 2\mu^2 \tau (W_{ij}^2 + W_{ii} W_{jj}) \\ &+ \frac{1}{m} \sum_{l=1}^m 2\mu^2 \left((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii} (\mathbf{A}_l)_{jj} \right) + 4\mu \tau \mathbf{W}_{\cdot i}^\top \mathbf{D} \mathbf{W}_{\cdot j}. \end{aligned}$$

Combining all contributions yields a quadratic function of μ , expressed as (15) and the closed-form solution for μ is given

by

$$\mu = -e + \mathcal{T}(e - g/h, \Lambda_{ij}/h), \quad (16)$$

where $\mathcal{T}(u, v) = \text{sign}(u) \max\{|u| - v, 0\}$ is the soft-thresholding function, and

$$\begin{aligned} h &= \tau (W_{ij}^2 + W_{ii}W_{jj}) + \frac{\sum_{l=1}^m ((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii}(\mathbf{A}_l)_{jj})}{m}, \\ g &= Q_{ij} - \tau W_{ij} + \tau \mathbf{W}_{\cdot i}^\top \mathbf{D} \mathbf{W}_{\cdot j} + \frac{\sum_{l=1}^m (\mathbf{A}_l)_{\cdot i}^\top \mathbf{D} (\mathbf{A}_l)_{\cdot j}}{m}, \\ e &= \Sigma_{ij} + D_{ij}. \end{aligned}$$

When $i = j$, the substitution of $\mathbf{D}' = \mathbf{D} + \mu \mathbf{e}_i \mathbf{e}_i^\top$ yields the update formula for D_{ii} as derived from (16). The components of the update formula are defined as follows:

$$\begin{aligned} h &= \tau D_{ii}^2 + \frac{\sum_{l=1}^m (\mathbf{A}_l)_{ii}^2}{m}, \\ g &= Q_{ii} - \tau W_{ii} + \tau \mathbf{W}_{\cdot i}^\top \mathbf{D} \mathbf{W}_{\cdot i} + \frac{\sum_{l=1}^m (\mathbf{A}_l)_{\cdot i}^\top \mathbf{D} (\mathbf{A}_l)_{\cdot i}}{m}, \\ e &= \Sigma_{ii} + D_{ii}. \end{aligned}$$

The primary computational challenge is encountered in the evaluation of the expressions $\mathbf{W}_{\cdot i}^\top \mathbf{D} \mathbf{W}_{\cdot j}$ and $\sum_{l=1}^m (\mathbf{A}_l)_{\cdot i}^\top \mathbf{D} (\mathbf{A}_l)_{\cdot j}$. To alleviate the computational burden, which typically necessitates $\mathcal{O}(md^2)$ time, we employ intermediary $d \times d$ matrices defined as $\mathbf{U}_l = \mathbf{D} \mathbf{A}_l$, $\mathbf{V} = \mathbf{D} \mathbf{W}$. This strategy facilitates the computation of $\mathbf{W}_{\cdot i}^\top \mathbf{V}_{\cdot j}$ and $\mathbf{W}_{\cdot j}^\top (\mathbf{A}_l)_{\cdot i}$ using $\mathcal{O}(d)$ floating-point operations. Furthermore, the maintenance \mathbf{U}_l and \mathbf{V} requires the adjustment of $\mathcal{O}(d)$ elements. For a comprehensive description of the methodology, please consult Algorithm 3.

A. Updating Only a Subset of Variables

We introduce an algorithm that leverages the active set update strategy to exploit the solution sparsity, thereby accelerating computation. The strategy involves two consecutive nested loops. In the outer loop, all coordinates are categorized into two sets based on a heuristic coordinate gradient thresholding rule (strong rule, [53]): the free set (active coordinates) and the fixed set (inactive coordinates). Subsequently, each iteration of the outer loop initiates an inner loop that performs coordinate optimization exclusively on the active set until convergence, while the inactive coordinates remain fixed at zero. Post each inner loop iteration, the algorithm employs a heuristic rule to reassess and update the active set, aiming to further reduce the objective value. This iterative process continues until the composition of the active set stabilizes across subsequent outer

Algorithm 3: Coordinate Descent Algorithms in Conjunction with the active set strategy.

Input: $\nabla f(\Sigma)$, Λ , Σ , $\{\mathbf{a}_l\}_{l=1}^m$
1 Initialize $\mathbf{W} = \Sigma^{-1}$, $\mathbf{A}_l = \mathbf{a}_l \mathbf{a}_l^\top$, $\mathbf{D} = \mathbf{0}$, $\mathbf{U}_l = \mathbf{0}$, $\mathbf{V} = \mathbf{0}$
2 while not converged do
 // Partition the variables into fixed and free sets
3 $\mathcal{C}_{\text{fixed}} = \{(i, j) \mid |\nabla f_{ij}(\Sigma)| \leq \Lambda_{ij} \text{ and } \Sigma_{ij} = 0\}$
4 $\mathcal{C}_{\text{free}} = \{(i, j) \mid |\nabla f_{ij}(\Sigma)| \leq \Lambda_{ij} \text{ and } \Sigma_{ij} \neq 0\}$
5 for $(i, j) \in \mathcal{A}_{\text{free}}$ **do**
 6 $h =$
 $\tau (W_{ij}^2 + W_{ii}W_{jj}) + \frac{\sum_{l=1}^m ((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii}(\mathbf{A}_l)_{jj})}{m}$
 7 $g = \nabla_{ij} f(\Sigma) + \tau \mathbf{W}_{\cdot i}^\top \mathbf{V}_{\cdot j} + \frac{\sum_{l=1}^m (\mathbf{A}_l)_{\cdot i}^\top (\mathbf{U}_l)_{\cdot j}}{m}$
 8 $e = \Sigma_{ij} + D_{ij}$
 9 $\mu = -e + \mathcal{T}(e - g/h, \Lambda_{jk}/h)$
 10 $D_{ij} = D_{ij} + \mu$
 11 $(\mathbf{U}_l)_{i \cdot} = (\mathbf{U}_l)_{i \cdot} + \mu (\mathbf{A}_l)_{j \cdot}$
 12 $(\mathbf{U}_l)_{j \cdot} = (\mathbf{U}_l)_{j \cdot} + \mu (\mathbf{A}_l)_{i \cdot}$
 13 $\mathbf{V}_{i \cdot} = \mathbf{V}_{i \cdot} + \mu \mathbf{W}_{j \cdot}$, $\mathbf{V}_{j \cdot} = \mathbf{V}_{j \cdot} + \mu \mathbf{W}_{i \cdot}$
 14 end
15 end
Output: \mathbf{D}

loop iterations. The efficacy of the active set strategy is well-documented in practical applications [50] and is underpinned by strong theoretical foundations [54].

The free set and fixed set are respectively defined as follows:

$$\begin{aligned} \Sigma_{ij} &\in \mathcal{C}_{\text{fixed}} \text{ if } |\nabla f_{ij}(\Sigma)| \leq \Lambda_{ij} \text{ and } \Sigma_{ij} = 0, \\ \Sigma_{ij} &\in \mathcal{C}_{\text{free}} \text{ if } |\nabla f_{ij}(\Sigma)| > \Lambda_{ij} \text{ and } \Sigma_{ij} \neq 0. \end{aligned}$$

Further, we ascertain that a Newton update confined to the variables within the fixed set does not alter any coordinates in that set. For a detailed discussion on this analysis, we refer readers to the work by [50], [55].

Remark 21. The computation of the Newton direction is streamlined when Σ is diagonal, which is encountered in the initial Newton iteration, wherein the Algorithm is started with an identity matrix. This simple scenario has a succinct closed-form optimal solution presented in (17). On this occasion, it only takes $\mathcal{O}(1)$ time to determine each variable, so the time complexity for solving the first Newton direction drops sharply to $\mathcal{O}(md^2)$. The core cause of this occurrence lies in the fact that, in this case, the Hessian matrix is diagonal, indicating that each one-variable sub-problem is detached from

$$\begin{aligned} &\frac{1}{2} \left(\frac{\sum_{l=1}^m ((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii}(\mathbf{A}_l)_{jj})}{m} + \tau (W_{ij}^2 + W_{ii}W_{jj}) \right) \mu^2 + \\ &\left(Q_{ij} - \tau W_{ij} + \tau \mathbf{W}_{\cdot i}^\top \mathbf{D} \mathbf{W}_{\cdot j} + \frac{\sum_{l=1}^m (\mathbf{A}_l)_{\cdot i}^\top \mathbf{D} (\mathbf{A}_l)_{\cdot j}}{m} \right) \mu + 2\Lambda_{ij} |\Sigma_{ij} + D_{ij} + \mu| \end{aligned} \quad (15)$$

$$\hat{D}_{ij} = \begin{cases} \mathcal{T} \left(-\frac{Q_{ij}}{\tau W_{ii} W_{jj} + \sum_{l=1}^m ((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii} (\mathbf{A}_l)_{jj})}, \frac{\Lambda_{ij}}{\tau W_{ii} W_{jj} + \sum_{l=1}^m ((\mathbf{A}_l)_{ij}^2 + (\mathbf{A}_l)_{ii} (\mathbf{A}_l)_{jj})} \right) & \text{if } i \neq j, \\ -\Sigma_{ii} + \mathcal{T} \left(\Sigma_{ii} - \frac{Q_{ii} - \tau W_{ii}}{\tau W_{ii}^2 + \sum_{l=1}^m (\mathbf{A}_l)_{ii}^2}, \frac{\Lambda_{ii}}{\tau W_{ii}^2 + \sum_{l=1}^m (\mathbf{A}_l)_{ii}^2} \right) & \text{if } i = j. \end{cases} \quad (17)$$

the others. This means that updating each variable individually once suffices to achieve the optimal solution.

APPENDIX C PROOF OF STATISTICAL THEORY

In this appendix, we provide the proofs of all the statistical theoretical results in Section IV-B. In Subsection C-A, we begin by Lemma 22 that connects the SE condition to the localized version of the sparse strong convexity/sparse strong smoothness. Subsequently, we introduce several preliminary lemmas that support Lemma 24. Following this, we present Lemma 25, which sets bounds on the estimation error for the general problem 9, followed by Lemma 26 that outlines the estimation error limited for approximate solutions derived using the MM-based algorithm. Subsection C-B elaborates on Theorem 12, leveraging Lemma 25 and 26 to illustrate the solution path's contraction characteristic. Subsection C-C focuses on key concentration inequalities, which are vital for deducing the explicit statistical convergence rate. In Subsections C-D and C-E, the document presents the proofs for Corollary 14 and Corollary 15, respectively.

A. Technical Lemmata

Lemma 22. Let $f(\cdot)$ be a convex differentiable function, and let $G_f(\Sigma_1, \Sigma_2)$ be defined as

$$G_f(\Sigma_1, \Sigma_2) = f(\Sigma_1) - f(\Sigma_2) - \langle \nabla f(\Sigma_2), \Sigma_1 - \Sigma_2 \rangle.$$

Define the symmetrized Bregman divergence for $f(\cdot)$ as

$$\begin{aligned} G_f^s(\Sigma_1, \Sigma_2) &= G_f(\Sigma_1, \Sigma_2) + G_f(\Sigma_2, \Sigma_1) \\ &= \langle \nabla f(\Sigma_1) - \nabla f(\Sigma_2), \Sigma_1 - \Sigma_2 \rangle. \end{aligned}$$

For any $\Sigma_1, \Sigma_2 \in \mathcal{B}(\Sigma^*, r)$ such that

$$\max \{ \|\Sigma_1\|_{\bar{\mathcal{S}}}^-, \|\Sigma_2\|_{\bar{\mathcal{S}}}^- \} \leq \tilde{s},$$

the following holds:

$$\begin{aligned} \frac{1}{2} \rho_{2s^*+2\tilde{s}}^- \|\Sigma_1 - \Sigma_2\|_2^2 &\leq G_f(\Sigma_1, \Sigma_2) \\ &\leq \frac{1}{2} \rho_{2s^*+2\tilde{s}}^+ \|\Sigma_1 - \Sigma_2\|_2^2, \\ \rho_{2s^*+2\tilde{s}}^- \|\Sigma_1 - \Sigma_2\|_2^2 &\leq G_f^s(\Sigma_1, \Sigma_2) \\ &\leq \rho_{2s^*+2\tilde{s}}^+ \|\Sigma_1 - \Sigma_2\|_2^2. \end{aligned}$$

Proof: According to the mean value theorem and the convexity of the local region $\mathcal{B}(\Sigma^*, r)$, there exists parameters $\theta_1, \theta_2, \theta_3 \in [0, 1]$ such that

$$\begin{aligned} \tilde{\Sigma} &= \theta_1 \Sigma_1 + (1 - \theta_1) \Sigma_2 \in \mathcal{B}(\Sigma^*, r), \\ \tilde{\Sigma}' &= \theta_2 \tilde{\Sigma} + (1 - \theta_2) \Sigma_2 \in \mathcal{B}(\Sigma^*, r), \\ \tilde{\Sigma}'' &= \theta_3 \Sigma_1 + (1 - \theta_3) \Sigma_2 \in \mathcal{B}(\Sigma^*, r), \end{aligned}$$

and satisfies $\max \{ \|\tilde{\Sigma}\|_0, \|\tilde{\Sigma}'\|_0 \} \leq 2\tilde{s}$. The function $G_f(\Sigma_1, \Sigma_2)$ and the symmetrized Bregman divergence can be articulated as follows:

$$\begin{aligned} &f(\Sigma_1) - f(\Sigma_2) - \langle \nabla f(\Sigma_2), \Sigma_1 - \Sigma_2 \rangle \\ &= \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma_2), \Sigma_1 - \Sigma_2 \rangle = \|\Sigma_1 - \Sigma_2\|_{\nabla^2 f(\tilde{\Sigma}')}^2, \end{aligned}$$

and

$$\langle \nabla f(\Sigma_1) - \nabla f(\Sigma_2), \Sigma_1 - \Sigma_2 \rangle = \|\Sigma_1 - \Sigma_2\|_{\nabla^2 f(\tilde{\Sigma}'')}^2.$$

By applying the definition of the sparse eigenvalue, we arrive at the desired conclusion. ■

Lemma 23. Consider $\Sigma(\theta) = \Sigma^* + \theta(\Sigma - \Sigma^*)$ for $\theta \in (0, 1]$. Then, we have

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

Proof: Define $\varphi(\theta)$ as follows:

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

The derivative of $f(\Sigma(\theta))$ with respect to θ is $\langle \nabla f(\Sigma(\theta)), \Sigma - \Sigma^* \rangle$. Thus

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

Consequently, the symmetric Bregman divergence $G_f^s(\Sigma(\theta), \Sigma^*)$ is

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

for $0 < \theta \leq 1$. Given $\varphi'(1) = G_f^s(\Sigma, \Sigma^*)$, and assuming the convexity of $\varphi(\theta)$ based on the convexity of $f(\cdot)$ and the linearity of $\Sigma(\theta)$, $\varphi'(\theta)$ is non-decreasing. Thus

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

completing the proof. For a more detailed proof of the convexity of $\varphi(\theta)$, see [28], [56]. ■

Lemma 24. Consider a set \mathcal{E} such that $\mathcal{S}^* \subseteq \mathcal{E}$. Given the condition $\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon \leq \|\Lambda_{\mathcal{E}}\|_{\min}$, we have

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

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

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

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

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

Given $\|\tilde{\Delta}\|_{H(\theta)}^2 \geq 0$, it follows that:

$$0 \leq \|\Omega\|_{\max} \|\tilde{\Delta}\|_1 - \underbrace{\langle \nabla f(\Sigma^*), \tilde{\Delta} \rangle}_{\text{I}} - \underbrace{\langle \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \rangle}_{\text{II}}. \quad (18)$$

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

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

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

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

Plugging the above results into 18 yields

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

where (i) is from $\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2}$, Assumption 10 and $\omega_{\Lambda}(\tilde{\Sigma}) \leq \varepsilon$. ■

Lemma 25. Suppose that Assumption 7, 8, 9 and 10 hold. Consider the general problem in (9). Let there exists a set \mathcal{E} such that

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

Then the ε -optimal solution $\tilde{\Sigma}$ satisfies

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

Proof: Given the assumption 10, it follows that $\lambda \geq 2(\|(\nabla f(\Sigma^*))_{\max} + \varepsilon)$. We introduce an intermediate estimator $\tilde{\Sigma}^* = \Sigma^* + \theta(\tilde{\Sigma} - \Sigma^*)$, where θ is chosen to ensure $\|\tilde{\Sigma}^* - \Sigma^*\|_{\text{F}} = r$ if $\|\tilde{\Sigma} - \Sigma^*\|_{\text{F}} > r$; otherwise $\theta = 1$. This construction guarantees that $\|\tilde{\Sigma}^* - \Sigma^*\|_{\text{F}} \leq r$. By leveraging Lemma 24, we deduce that the approximate solution is contained within the ℓ_1 -cone. Given the structure

of $\tilde{\Sigma}^*$, we find that $\tilde{\Sigma}^* - \Sigma^* = \theta(\tilde{\Sigma} - \Sigma^*)$. This leads to the following inequality

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

By integrating the aforementioned inequality with the premise that $|\mathcal{E}| \leq 2s^*$, it follows that $\tilde{\Sigma}^*$ resides within the local ℓ_1 -cone. Furthermore, by synthesizing Lemma 22, Definition 19 and Proposition 20, one can deduce the presence of localized restricted strong convexity, i.e.

$$\rho_{2s^*+2\tilde{s}} \|\tilde{\Sigma}^* - \Sigma^*\|_{\text{F}}^2 \leq G_f^s(\tilde{\Sigma}^*, \Sigma^*). \quad (19)$$

We use Lemma 23 to bound the right hand side of the above inequality such as

$$\begin{aligned} G_f^s(\tilde{\Sigma}^*, \Sigma^*) &\leq \theta G_f^s(\tilde{\Sigma}, \Sigma^*) \\ &= \theta \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \tilde{\Sigma} - \Sigma^* \rangle. \end{aligned} \quad (20)$$

Plugging (20) back into (19) yields

$$\begin{aligned} \rho_{2s^*+2\tilde{s}} \|\tilde{\Sigma}^* - \Sigma^*\|_{\text{F}}^2 &\leq \theta \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \tilde{\Delta} \rangle \\ &= \theta \underbrace{\langle \nabla f(\tilde{\Sigma}) + \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \rangle}_{\text{I}} \\ &\quad - \theta \underbrace{\langle \nabla f(\Sigma^*), \tilde{\Delta} \rangle}_{\text{II}} - \theta \underbrace{\langle \Lambda \odot \tilde{\Xi}, \tilde{\Delta} \rangle}_{\text{III}}, \end{aligned} \quad (21)$$

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**, we partition the support of Ω and $\tilde{\Delta}$ into the set \mathcal{E} and its complement $\bar{\mathcal{E}}$. By employing the matrix Hölder's inequality, we derive

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

where (i) is from $\|\Omega\|_{\max} \leq \varepsilon$ by Definition 6. For term **II**, we divide the support of $\nabla f(\Sigma^*)$ and $\tilde{\Delta}$ into the set \mathcal{E} and $\bar{\mathcal{E}}$. Applying the matrix Hölder's inequality, we obtain

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

For term **III**, we separate the support of $\Lambda \odot \tilde{\Xi}$ and $\tilde{\Delta}$ into the set \mathcal{S}^* and $\bar{\mathcal{S}}^*$. And then applying the matrix Hölder's

inequality, we obtain

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

where (i) is from

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

and (ii) is from

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

and

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

Combining the above results into (21) yields

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

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

For IV, we have $\theta \left\| \left(\tilde{\Sigma} - \Sigma^* \right)_{\mathcal{E}} \right\|_F = \left\| \left(\tilde{\Sigma}^* - \Sigma^* \right)_{\mathcal{E}} \right\|_F$. Consequently, we derive

$$\rho_{2s^*+2\tilde{s}}^- \|\tilde{\Sigma}^* - \Sigma^*\|_F \leq \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*},$$

which serves as a contraction in relation to the construction of $\tilde{\Sigma}^*$. This leads to the conclusion that $\tilde{\Sigma}^* = \tilde{\Sigma}$. Therefore, the desired bound is satisfied for $\tilde{\Sigma}$. ■

Lemma 26. Suppose that Assumptions 1 and 8 hold. Consider the problem in (9). Define the set $\mathcal{E}^{(k)} = S^* \cup \mathcal{S}^{(k)}$, where $\mathcal{S}^{(k)} = \left\{ (i, j) \mid \Lambda_{ij}^{(k-1)} \leq p'_\lambda(u) \right\}$ with $u = c\lambda$ and $c = \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}^-} \in (0, \alpha)$, such that $p'_\lambda(c\lambda) \geq \frac{\lambda}{2}$. This can always hold due to Assumption 5. Then for $k \geq 1$, we establish that:

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

Proof: We first show that $|\mathcal{E}^{(k)}| \leq 2s^*$ for all k by induction on k .

Base Case ($k = 1$):

For $k = 1$, we have $\Lambda_{ij}^{(0)} = \lambda$ for $i \neq j$, where λ is a threshold parameter satisfying $\lambda \geq p'_\lambda(u)$. It follows that $|\mathcal{S}^{(1)}| \leq s^*$. And we see $\mathcal{E}^{(1)} = S^* \cup \mathcal{S}^{(1)}$, which implies $|\mathcal{E}^{(1)}| \leq 2s^*$ holds.

Inductive Step ($k \geq 2$):

Assume for some $k \geq 2$, $|\mathcal{E}^{(k-1)}| \leq 2s^*$ holds true. We aim to show that $|\mathcal{E}^{(k)}| \leq 2s^*$. For any $(i, j) \in \mathcal{S}^{(k)}$, we obtain $|\tilde{\Sigma}_{ij}^{(k-1)}| \geq u$ and further have

$$\begin{aligned}
\sqrt{|\mathcal{S}^{(k)} \setminus \mathcal{S}^*|} &\leq \sqrt{\sum_{(i,j) \in \mathcal{S}^{(k)} \setminus \mathcal{S}^*} \left(u^{-1} \tilde{\Sigma}_{ij}^{(k-1)} \right)^2} \\
&= u^{-1} \left\| \tilde{\Sigma}_{\mathcal{S}^{(k)} \setminus \mathcal{S}^*}^{(k-1)} \right\|_F \\
&= u^{-1} \left\| \left(\tilde{\Sigma}^{(k-1)} - \Sigma^* \right)_{\mathcal{S}^{(k)} \setminus \mathcal{S}^*} \right\|_F \\
&\leq u^{-1} \left\| \tilde{\Sigma}^{(k-1)} - \Sigma^* \right\|_F.
\end{aligned} \tag{22}$$

For any $(i, j) \in \overline{\mathcal{S}^{(k-1)}}$, it follows that $\Lambda_{ij}^{(k-2)} = p'_\lambda(\tilde{\Sigma}_{ij}^{(k-2)}) \geq p'_\lambda(u) \geq \frac{\lambda}{2}$. This chain of inequalities indicates that

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

Furthermore, it is established that $|\mathcal{E}^{(k-1)}| \leq 2s^*$ and $S^* \subseteq \mathcal{E}^{(k-1)}$. By applying Lemma 25 with $\tilde{\Sigma} = \tilde{\Sigma}^{(k-1)}$, $\mathcal{E} = \mathcal{E}^{(k-1)}$, and $\Lambda_{S^*} = \Lambda_{S^*}^{(k-2)}$ yields

$$\left\| \tilde{\Sigma}^{(k-1)} - \Sigma^* \right\|_F \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}^-} \lambda \sqrt{s^*}.$$

Substituting this result into the inequality 22 yields

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

Thus, we have

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

completing the induction process.

Then according to the definition of $\mathcal{E}^{(k)}$ and $\mathcal{S}^{(k)}$, it follows that

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

By employing Lemma 25 with $\tilde{\Sigma} = \tilde{\Sigma}^{(k)}$, $\mathcal{E} = \mathcal{E}^{(k)}$, and $\Lambda_{\mathcal{S}^*} = \Lambda_{\mathcal{S}^*}^{(k-1)}$, the ε -optimal solution $\tilde{\Sigma}^{(k)}$ to 9 satisfies

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

B. Proof of Theorem 12

Proof: Based on Lemma 26, we have

$$\begin{aligned} \left\| \tilde{\Sigma}^{(k)} - \Sigma^* \right\|_{\text{F}} & \leq \underbrace{\frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\left\| (\nabla f(\Sigma^*))_{\mathcal{E}^{(k)}} \right\|_{\text{F}} + \varepsilon \sqrt{|\mathcal{E}^{(k)}|} \right)}_{\text{I}} \\ & \quad + \underbrace{\frac{1}{\rho_{2s^*+2\tilde{s}}} \left\| \Lambda_{\mathcal{S}^*}^{(k-1)} \right\|_{\text{F}}}_{\text{II}}. \end{aligned} \quad (23)$$

Then, we proceed to establish bounds for the term **I** and **II**, respectively.

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

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

where the second inequality is due to

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

which follows from the inequality (22).

By Assumption 5 and 9, for any Σ , if $|\Sigma_{ij} - \Sigma_{ij}^*| \geq u$, then $p'_\lambda(\Sigma_{ij}) \leq \lambda \leq \lambda u^{-1} |\Sigma_{ij} - \Sigma_{ij}^*|$; otherwise, $p'_\lambda(\Sigma_{ij}) \leq p'_\lambda(|\Sigma_{ij}^*| - u) = 0$. Therefore, for term **II**, we have

$$\text{II} \leq \lambda u^{-1} \left\| \tilde{\Sigma}_{\mathcal{S}^*}^{(k-1)} - \Sigma_{\mathcal{S}^*}^* \right\|_{\text{F}} \leq \lambda u^{-1} \left\| \tilde{\Sigma}^{(k-1)} - \Sigma^* \right\|_{\text{F}}.$$

Substituting the above results into (23) yields

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

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

C. Concentration Inequality

Lemma 27 (Hanson-Wright Inequality). *Let $\mathbf{x} = (x_1, \dots, x_n) \in \mathbb{R}^n$ be a random vector with independent components x_i which satisfy $\mathbb{E}(x_i) = 0$ and $\|x_i\|_{\psi_2} \leq L$. Let \mathbf{B} be an $n \times n$ matrix. Then for any $t > 0$,*

$$\begin{aligned} & \mathbb{P} \left\{ \left| \mathbf{x}^\top \mathbf{B} \mathbf{x} - \mathbb{E}(\mathbf{x}^\top \mathbf{B} \mathbf{x}) \right| > t \right\} \\ & \leq 2 \exp \left(-c' \min \left(\frac{t^2}{L^4 \|\mathbf{B}\|_{\text{F}}^2}, \frac{t}{L^2 \|\mathbf{B}\|} \right) \right). \end{aligned} \quad (24)$$

Lemma 28 (Lemma D.1 in [56]). *Let \mathbf{x} be a sub-Gaussian random vector with zero mean and covariance Σ^* and $\{\mathbf{x}_i\}_{i=1}^n$ be a collection of i.i.d. samples from \mathbf{x} . There exists some constant c_1, c_2 , and t_0 such that for all t with $0 < t < t_0$, the sample covariance matrix \mathbf{S} satisfies the following tail bound*

$$\mathbb{P}(|\Sigma_{ij}^* - S_{ij}| > t) \leq c_1 \exp(-c_2 n t^2).$$

Lemma 29. *Under Assumptions 9 and the same conditions in Lemma 28, the columns of $\tilde{\mathbf{A}}$ are normalized such that $\max_l \|\tilde{\mathbf{A}}_{\cdot l}\|_2 \leq \sqrt{m}$, there exists some constant c_1 such that*

$$\mathbb{P} \left(\left\| \frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*)) \right\|_{\max} \leq \lambda \right) \geq 1 - \frac{c_1}{d}.$$

Proof: We begin by establishing a bound on the probability that $\left\| \frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*)) \right\|_{\max} > t$ using a union bound approach.

First, consider the bound $\mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \eta'_i \mathbf{a}_i \mathbf{a}_i^\top \right\|_{\max} > t \right)$, where $\eta'_i = \mathbf{a}_i^\top (\Sigma_N - \Sigma^*) \mathbf{a}_i + \eta_i$. By applying the Hanson-Wright inequality (24), we establish

$$\begin{aligned} & \mathbb{P}(\mathbf{a}_i^\top (\Sigma_N - \Sigma^*) \mathbf{a}_i > t) \\ & = \mathbb{P}(|\langle \mathbf{a}_i \mathbf{a}_i^\top, (\Sigma_N - \Sigma^*) \rangle| > t) \\ & \leq 2 \exp \left(-c' \min \left(\frac{t^2}{L^4 \|\Sigma_N - \Sigma^*\|_{\text{F}}^2}, \frac{t}{L^2 \|\Sigma_N - \Sigma^*\|} \right) \right), \end{aligned}$$

where c' is a positive constant, and L is a scaling factor. This indicates that $\langle \mathbf{a}_i \mathbf{a}_i^\top, \Sigma_N - \Sigma^* \rangle$ is a sub-exponential random variable. Additionally, given that η_i follows a sub-exponential distribution, η'_i is also sub-exponential random variable. By summing over individual probabilities, we obtain

$$\begin{aligned} & \mathbb{P} \left(\left\| \frac{1}{m} \mathcal{A}^*(\mathbf{y} - \mathcal{A}(\Sigma^*)) \right\|_{\max} > t \right) \\ & = \mathbb{P} \left(\left\| \frac{1}{m} \sum_{i=1}^m \eta'_i \mathbf{a}_i \mathbf{a}_i^\top \right\|_{\max} > t \right) \\ & \leq \sum_{l=1}^d \mathbb{P} \left(\frac{1}{m} \left| \tilde{\mathbf{A}}_{\cdot l}^\top \cdot \boldsymbol{\eta}' \right| > t \right), \end{aligned} \quad (25)$$

Define $\zeta_l = \tilde{\mathbf{A}}_{\cdot l}^\top \cdot \boldsymbol{\eta}'$. Given that η'_i is sub-exponential $(0, \xi^2)$ for $i = 1, \dots, m$, we apply Lemma 28 and employ concentration inequalities:

$$\mathbb{E}(\exp\{t_0 \zeta_l\} + \exp\{-t_0 \zeta_l\}) \leq 2 \exp \left\{ nm^{-2} \left\| \tilde{\mathbf{A}}_{\cdot l} \right\|^2 \xi^2 t_0^2 / 2 \right\},$$

which implies

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

Selecting $t_0 = t \left(nm^{-2} \left\|\tilde{\mathbf{A}}_{\cdot l}\right\|_2^2 \xi^2\right)^{-1}$ yields that

$$\mathbb{P}(|\zeta_l| \geq t) \leq 2 \exp\left\{-nt^2 \cdot \left(2\xi^2 \left\|\tilde{\mathbf{A}}_{\cdot l}\right\|_2^2/m^2\right)^{-1}\right\}.$$

Plugging it into (25) results

$$\begin{aligned} & \mathbb{P}\left(\left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))\right\|_{\max} > t\right) \\ & \stackrel{(i)}{\leq} 2d^2 \exp\left\{-nmt^2 \cdot \left(2\xi^2 \max_l \left\{\left\|\tilde{\mathbf{A}}_{\cdot l}\right\|_2^2/m\right\}\right)^{-1}\right\} \\ & = c_1 \exp\{-c_2 nmt^2 + 2\log d\}, \end{aligned}$$

where (i) is from 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

$\lambda = \sqrt{\frac{3\log d}{c_3 mn}} \asymp \sqrt{\frac{\log d}{mn}}$, we obtain

$$\begin{aligned} & \mathbb{P}\left(\left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))\right\|_{\max} \leq \lambda\right) \\ & \geq 1 - c_1 \exp(-c_2 nm\lambda^2 + 2\log d) \\ & = 1 - \frac{c_1}{d}. \end{aligned}$$

Lemma 30. *Under the same conditions as in Lemma 29, the following result hold*

$$\left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\text{F}} = \mathcal{O}_p\left(\sqrt{\frac{s^*}{mn}}\right).$$

Proof: Following the analysis presented in Lemma 29 analysis, consider any M such that $0 < M\sqrt{\frac{1}{mn}}$, we have

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

Setting M such that $\sqrt{\frac{2\log s^*}{c_3}} < M$ and letting $M \rightarrow \infty$, we obtain

$$\lim_{M \rightarrow \infty} \sup_m \mathbb{P}\left(\left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\max} > M\sqrt{\frac{1}{mn}}\right) = 0.$$

The proof is completed by applying

$$\begin{aligned} & \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\text{F}} \\ & \leq \sqrt{s^*} \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\max}. \end{aligned}$$

D. Proof of Corollary 14

Proof: We begin by bound the gradient norm of $f(\boldsymbol{\Sigma})$ at $\boldsymbol{\Sigma}^*$:

$$\begin{aligned} \|\nabla f(\boldsymbol{\Sigma}^*)\|_{\max} & \leq \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))\right\|_{\max} \\ & \quad + \tau \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\max}. \end{aligned}$$

Assuming $\lambda \asymp \sqrt{\frac{\log d}{mn}}$ and $\tau \lesssim \sqrt{\frac{1}{mn}} \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\max}^{-1}$, and invoking Lemma 29, we establish that $\lambda \geq 2(\|\nabla f(\boldsymbol{\Sigma}^*)\|_{\max} + \varepsilon)$ holds w.h.p.

Applying Lemma 26 with $k = 1$, we obtain

$$\left\|\tilde{\boldsymbol{\Sigma}}^{(1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}. \quad (26)$$

Under the condition $\lambda \asymp \sqrt{\frac{\log d}{mn}}$, the Frobenius norm difference is asymptotically bounded as $\left\|\tilde{\boldsymbol{\Sigma}}^{(1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \lesssim \sqrt{\frac{s^* \log d}{mn}}$ w.h.p. ■

E. Proof of Corollary 15

Proof: We begin by bound the gradient norm of $f(\boldsymbol{\Sigma})$ at $\boldsymbol{\Sigma}^*$:

$$\begin{aligned} \|\nabla f(\boldsymbol{\Sigma}^*)\|_{\max} & \leq \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))\right\|_{\max} \\ & \quad + \tau \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\max}. \end{aligned}$$

Under Assumptions 7 and 9 condition, and with parameters λ and τ scaled as $\lambda \asymp \sqrt{\frac{\log d}{mn}}$ and $\tau \lesssim \sqrt{\frac{1}{mn}} \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\max}^{-1}$, then by Lemma 29, $\lambda \geq 2(\|\nabla f(\boldsymbol{\Sigma}^*)\|_{\max} + \varepsilon)$ holds w.h.p.

Next, applying Theorem 12, we derive

$$\begin{aligned} & \left\|\tilde{\boldsymbol{\Sigma}}^{(1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\|(\nabla f(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\|_{\text{F}} + \varepsilon \sqrt{s^*}\right) + \delta \left\|\tilde{\boldsymbol{\Sigma}}^{(k-1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\|(\nabla f(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\|_{\text{F}} + \varepsilon \sqrt{s^*}\right) + \delta^{k-1} \left\|\tilde{\boldsymbol{\Sigma}}^{(1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \\ & \leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left(\|(\nabla f(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\|_{\text{F}} + \varepsilon \sqrt{s^*}\right) + \delta^{k-1} \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}, \end{aligned} \quad (27)$$

where the last inequality is due to $\left\|\tilde{\boldsymbol{\Sigma}}^{(1)} - \boldsymbol{\Sigma}^*\right\|_{\text{F}} \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}$, which follows from Lemma 26 with $k = 1$.

Finally, considering the norm of the projected gradient within the support set \mathcal{S}^* , we obtain

$$\begin{aligned} & \|(\nabla f(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\|_{\text{F}} \\ & = \left\|\left(-\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*)) - \tau(\boldsymbol{\Sigma}^*)^{-1}\right)_{\mathcal{S}^*}\right\|_{\text{F}} \\ & \leq \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\text{F}} + \tau \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\mathcal{S}^*} \\ & \leq \left\|\frac{1}{m}\mathcal{A}^*(\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\right\|_{\text{F}} + \tau \left\|(\boldsymbol{\Sigma}^*)^{-1}\right\|_{\text{F}} \end{aligned} \quad (28)$$

By Lemma 30, $\left\| \frac{1}{m} \mathcal{A}^* (\mathbf{y} - \mathcal{A}(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*} \right\|_F = \mathcal{O}_p \left(\sqrt{\frac{s^*}{mn}} \right)$.
If $\tau \lesssim \sqrt{\frac{1}{mn}} \left\| (\boldsymbol{\Sigma}^*)^{-1} \right\|_{\max}^{-1}$, then $\|(\nabla f(\boldsymbol{\Sigma}^*))_{\mathcal{S}^*}\|_F = \mathcal{O}_p \left(\sqrt{\frac{s^*}{mn}} \right)$.

Furthermore, if $K \geq 1 + \frac{\log(\lambda\sqrt{mn})}{\log \delta^{-1}} \gtrsim \log(\lambda\sqrt{mn}) \gtrsim \log \log d$, we derive

$$\delta^{K-1} \lambda \sqrt{s^*} \leq \frac{1}{\lambda \sqrt{mn}} \lambda \sqrt{s^*} \leq \sqrt{\frac{s^*}{mn}}.$$

By integrating all findings with the condition $\varepsilon \lesssim \sqrt{\frac{1}{mn}}$, we conclude that $\left\| \tilde{\boldsymbol{\Sigma}}^{(K)} - \boldsymbol{\Sigma}^* \right\|_F = \mathcal{O}_p \left(\sqrt{\frac{s^*}{mn}} \right)$. ■

APPENDIX D

PROOF OF COMPUTATIONAL THEORY 18

In this appendix, we first show $f(\boldsymbol{\Sigma})$ exhibits properties of the local strong convexity, local strong smoothness and local restricted Hessian smoothness. Then we demonstrate the sparsity of the solution is preserved within a vicinity of the true model parameter $\boldsymbol{\Sigma}^*$, provided that the starting point is a sparse solution. Following this, it is established that proximal Newton updates achieve a quadratic rate of convergence towards a local minimizer at each step within each stage, assuming the initialization occurs within a suitably refined sparse region. Subsequently, Lemmas 38 and 35 confirm that each subproblem experiences a substantial reduction. Furthermore, we outline the necessary number of iterations required at each convex relaxation stage to satisfy the approximate KKT conditions.

A. Proof of Lemma 16

Proof: We first verify the convexity of the set $\mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}}{4\tau\kappa} \right)$. For any $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}}{4\tau\kappa} \right)$, the linearity of the space implies that any convex combination of these points also resides within the set due to the properties of the norm and the definition of the set.

To establish the sparse strong convexity of $f(\boldsymbol{\Sigma})$ within this set, consider two arbitrary points $\boldsymbol{\Sigma}_3, \boldsymbol{\Sigma}_4$ and convex combination $\boldsymbol{\Sigma}_\rho = \rho \boldsymbol{\Sigma}_3 + (1-\rho) \boldsymbol{\Sigma}_4$ for $\rho \in [0, 1]$. The function $f(\boldsymbol{\Sigma})$ can be expanded using a Taylor series approximation around $\boldsymbol{\Sigma}_4$:

$$\begin{aligned} f(\boldsymbol{\Sigma}_3) &= f(\boldsymbol{\Sigma}_4) + \langle \nabla f(\boldsymbol{\Sigma}_4), \boldsymbol{\Sigma}_3 - \boldsymbol{\Sigma}_4 \rangle \\ &\quad + \frac{1}{2} \text{vec}^\top(\boldsymbol{\Sigma}_3 - \boldsymbol{\Sigma}_4) \nabla^2 f(\boldsymbol{\Sigma}_\rho) \text{vec}(\boldsymbol{\Sigma}_3 - \boldsymbol{\Sigma}_4). \end{aligned}$$

The next step is to estimate the lower bound of $\lambda_{\min}(\nabla^2 f(\boldsymbol{\Sigma}_\rho))$. Using the fact that $\lambda_{\min}(\mathbf{X}^{-1}) = \|\mathbf{X}\|_2^{-1}$

and $\|\mathbf{X} \otimes \mathbf{Y}\|_2 = \|\mathbf{X}\|_2 \|\mathbf{Y}\|_2$ for any matrices \mathbf{X}, \mathbf{Y} , we deduce

$$\begin{aligned} &\lambda_{\min}(\nabla^2 f(\boldsymbol{\Sigma}_\rho)) \\ &= \lambda_{\min} \left(\sum_{i=1}^m \mathbf{A}_i \otimes \mathbf{A}_i^\top + \tau \boldsymbol{\Sigma}_\rho^{-1} \otimes \boldsymbol{\Sigma}_\rho^{-1} \right) \\ &\geq \lambda_{\min} \left(\sum_{i=1}^m \mathbf{A}_i \otimes \mathbf{A}_i^\top \right) + \tau \lambda_{\min}(\boldsymbol{\Sigma}_\rho^{-1} \otimes \boldsymbol{\Sigma}_\rho^{-1}) \\ &\stackrel{(i)}{\geq} \rho_{2s^*+2\tilde{s}}^- + \tau \|\boldsymbol{\Sigma}_\rho\|_2^{-2}, \end{aligned}$$

where (i) is from Assumption 4. One has

$$\begin{aligned} \|\boldsymbol{\Sigma}_\rho\|_2 &= \|\boldsymbol{\Sigma}^* - \boldsymbol{\Sigma}^* + \rho \boldsymbol{\Sigma}_1 + (1-\rho) \boldsymbol{\Sigma}_2\|_2 \\ &\leq \|\boldsymbol{\Sigma}^*\|_2 + \rho \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}^*\|_2 + (1-\rho) \|\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}^*\|_2 \\ &\leq \|\boldsymbol{\Sigma}^*\|_2 + \rho \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}^*\|_F + (1-\rho) \|\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}^*\|_F \\ &\stackrel{(i)}{\leq} \kappa + \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa} \leq \kappa + \frac{\kappa \rho_{2s^*+2\tilde{s}}^-}{4\tau}, \end{aligned}$$

where (i) holds due to $\|\boldsymbol{\Sigma}^*\|_2 \leq \kappa$ and $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa} \right)$; the last inequality holds due to $\kappa \geq 1$. Thus, $f(\boldsymbol{\Sigma})$ is shown to be strongly convex over $\mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa} \right)$ with a convexity parameter of $\left(\rho_{2s^*+2\tilde{s}}^- + \frac{16\tau^3}{\kappa^2(4\tau + \rho_{2s^*+2\tilde{s}}^-)^2} \right)$.

Next, we establish the strong smoothness of $f(\boldsymbol{\Sigma})$, which involves setting an upper bound for $\lambda_{\max}(\nabla^2 f(\boldsymbol{\Sigma}_\rho))$. By the fact that $\lambda_{\max}(\mathbf{X}^{-1}) = \lambda_{\min}^{-1}(\mathbf{X})$ and $\lambda_{\min}(\mathbf{X} \otimes \mathbf{Y})$ for any matrices \mathbf{X}, \mathbf{Y} . Therefore,

$$\begin{aligned} &\lambda_{\max}(\nabla^2 f(\boldsymbol{\Sigma}_\rho)) \\ &= \lambda_{\max} \left(\sum_{i=1}^m \mathbf{A}_i \otimes \mathbf{A}_i^\top + \tau \boldsymbol{\Sigma}_\rho^{-1} \otimes \boldsymbol{\Sigma}_\rho^{-1} \right) \\ &\leq \lambda_{\max} \left(\sum_{i=1}^m \mathbf{A}_i \otimes \mathbf{A}_i^\top \right) + \tau \lambda_{\max}(\boldsymbol{\Sigma}_\rho^{-1} \otimes \boldsymbol{\Sigma}_\rho^{-1}) \\ &\leq \rho_{2s^*+2\tilde{s}}^+ + \tau \lambda_{\min}^{-2}(\boldsymbol{\Sigma}_\rho). \end{aligned}$$

One has

$$\begin{aligned} &\lambda_{\min}(\boldsymbol{\Sigma}_\rho) \\ &= \lambda_{\min}(\boldsymbol{\Sigma}^* - \boldsymbol{\Sigma}^* + \rho \boldsymbol{\Sigma}_1 + (1-\rho) \boldsymbol{\Sigma}_2) \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}^*) + \rho \lambda_{\min}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}^*) + (1-\rho) \lambda_{\min}(\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}^*) \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}^*) - \rho \lambda_{\max}(\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}^*) - (1-\rho) \lambda_{\max}(\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}^*) \\ &\geq \lambda_{\min}(\boldsymbol{\Sigma}^*) - \rho \|\boldsymbol{\Sigma}_1 - \boldsymbol{\Sigma}^*\|_F - (1-\rho) \|\boldsymbol{\Sigma}_2 - \boldsymbol{\Sigma}^*\|_F \\ &\geq \frac{1}{\kappa} - \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa}, \end{aligned}$$

utilizing the bounds given by $\lambda_{\min}(\boldsymbol{\Sigma}^*) \geq \frac{1}{\kappa}$ and $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa} \right)$. Thus, $f(\boldsymbol{\Sigma})$ is proven to be $\left(\rho_{2s^*+2\tilde{s}}^+ + \frac{16\tau^3\kappa^2}{(4\tau - \rho_{2s^*+2\tilde{s}}^-)^2} \right)$ -smooth over $\mathcal{B} \left(\boldsymbol{\Sigma}^*, \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa} \right)$. ■

B. Preliminary Lemmas

Lemma 31 (Local Restricted Hessian Smoothness). *Recalled \tilde{s} introduced in Assumption 4. We assert the existence of constant $L_{2s^*+2\tilde{s}}$ and r 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 is satisfied:*

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

where $\Omega = \{\mathbf{U} \mid \text{supp}(\mathbf{U}) \subseteq (\text{supp}(\Sigma) \cup \text{supp}(\Sigma'))\}$.

Proof: To establish the Lipschitz continuity of $\nabla^2 f(\Sigma)$, consider the third derivative $\nabla^3 f(\Sigma)$, which is expressed as $-2\tau\Sigma^{-1} \otimes \Sigma^{-1} \otimes \Sigma^{-1}$. Employing a proof technique analogous to that utilized in Lemma 16, we have

$$\lambda_{\max}(|\nabla^3 f(\Sigma)|) \leq 2\tau\lambda_{\min}^{-3}(\Sigma) \leq 2\tau\kappa^3.$$

This finding allows us to conveniently define the Lipschitz constant for the Hessian matrix, represented as $L_{2s^*+2\tilde{s}}$, to be $2\tau\kappa^3$. ■

Lemma 32. *Under the same condition of Lemma 26, we have the following basic inequality*

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

Proof: Applying Lemma 26 with $k = 1$, we obtain

$$\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_F \leq \frac{2 + \sqrt{2}}{2\rho_{2s^*+2\tilde{s}}} \lambda \sqrt{s^*}. \quad (29)$$

Moreover, by invoking Lemma 24 yields that

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

where $\mathcal{E}^{(1)}$ can be taken as \mathcal{S}^* . Combining these results with 29, we conclude

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

Therefore, we obtain

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

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

$$\begin{aligned} &\left\langle \nabla f(\tilde{\Sigma}^{(1)}) - \nabla f(\Sigma^*), \tilde{\Sigma}^{(1)} - \Sigma^* \right\rangle \\ &\leq \left\| \nabla f(\tilde{\Sigma}^{(1)}) + \Lambda \odot \tilde{\Xi}^{(1)} \right\|_{\max} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 \\ &\quad + \left\| \Lambda \odot \tilde{\Xi}^{(1)} + \nabla f(\Sigma^*) \right\|_{\max} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 \\ &\leq c_0 \lambda \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_1 \leq \frac{c_1 \lambda^2 s^*}{\rho_{2s^*+2\tilde{s}}}, \end{aligned}$$

for some constants c_0 and c_1 . ■

Lemma 33. *Given $\omega_{\Lambda^{(k-1)}}(\tilde{\Sigma}^{(k)}) \leq \frac{\lambda}{8}$, we have that for all $t \geq 1$ at the $(k+1)$ -th stage,*

$$\omega_{\Lambda^{(k)}}(\Sigma_t^{(k+1)}) \leq \frac{\lambda}{4}$$

and

$$F_{\Lambda^{(k)}}(\Sigma_t^{(k+1)}) \leq F_{\Lambda^{(k)}}(\Sigma^*) + \frac{\lambda}{4} \left\| \Sigma_t^{(k)} - \Sigma^* \right\|_1.$$

Proof: We begin by establishing the foundational inequality at the $(k+1)$ -th stage. Note that $\Sigma_0^{(k+1)} = \tilde{\Sigma}^{(k)}$. By definition, we have

$$\omega_{\Lambda^{(k)}}(\Sigma_0) = \min_{\Xi \in \partial \|\Sigma_0\|_1} \left\| \nabla f(\Sigma_0) + \Lambda^{(k)} \odot \Xi \right\|_{\max}.$$

Using the triangle inequality, we can expand the above expression as follows:

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

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

Continuing, for any $t \geq 0$, consider

$$\Xi_t = \arg \min_{\Xi \in \partial \|\Sigma_t\|_1} \left\| \nabla f(\Sigma_t) + \Lambda^{(k)} \odot \Xi \right\|_{\max}.$$

By the convexity of F , we have

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

where (ii) is from the fact that $\left\| \nabla f(\Sigma_t) + \Lambda^{(k)} \odot \Xi_t \right\|_{\max} \leq \frac{\lambda}{4}$. ■

Lemma 34. *Suppose $\|(\Sigma_t)_{\tilde{s}}\|_0 \leq \tilde{s}$ and $\omega_{\Lambda^{(k)}}(\Sigma_t) \leq \frac{\lambda}{4}$. 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: The proof is analogous to that of Lemma 26. Therefore, we omit it for brevity. For more details, we refer readers to [28]. ■

Lemma 35. *Recall $\Sigma_{t+\frac{1}{2}}$ defined in (11) and δ_t as in (12). Denote $\Delta \Sigma_t = \Sigma_{t+\frac{1}{2}} - \tilde{\Sigma}_t$. Then we have*

$$\delta_t \leq -\|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2.$$

Proof: We note $\Delta \Sigma_t$ is the solution for $\bar{F}(\Sigma \mid \Sigma_t, \Lambda^{(k-1)})$. For any $\eta_t \in (0, 1]$, we have formula

$$\begin{aligned}
& \eta_t \langle \nabla f(\Sigma_t), \Delta \Sigma_t \rangle + \frac{\eta_t^2}{2} \|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2 + \left\| \Lambda^{(k-1)} \odot (\Sigma_t + \eta_t \Delta \Sigma_t) \right\|_1 \\
& \geq \langle \nabla f(\Sigma_t), \Delta \Sigma_t \rangle + \frac{1}{2} \|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2 + \left\| \Lambda^{(k-1)} \odot (\Sigma_t + \Delta \Sigma_t) \right\|_1.
\end{aligned} \tag{30}$$

$$\begin{aligned}
& \eta_t \langle \nabla f(\Sigma_t), \Delta \Sigma_t \rangle + \frac{\eta_t^2}{2} \|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2 + \eta_t \left\| \Lambda^{(k-1)} \odot (\Sigma_t + \eta_t \Delta \Sigma_t) \right\|_1 + (1 - \eta_t) \left\| \Lambda^{(k-1)} \odot \Sigma_t \right\|_1 \\
& \geq \langle \nabla f(\Sigma_t), \Delta \Sigma_t \rangle + \frac{1}{2} \|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2 + \left\| \Lambda^{(k-1)} \odot (\Sigma_t + \Delta \Sigma_t) \right\|_1.
\end{aligned} \tag{31}$$

(30). By the convexity of ℓ_1 -norm, we have formula (31). Rearranging the terms, canceling the $(1 - \eta_t)$ factor from both sides and let $\eta_t \rightarrow 1$, we obtain the desired inequality $\delta_t \leq -\|\Delta \Sigma_t\|_{\nabla^2 f(\Sigma_t)}^2$. ■

C. Proof of Theorem 18

Lemma 36 (Sparsity Preserving Lemma). *Suppose that Assumptions 8 and 9 hold. Given $\Sigma_t^{(k)} \in \mathcal{B}(\Sigma^*, r)$ and $\left\| \left(\Sigma_t^{(k)} \right)_{\bar{\mathcal{S}}} \right\|_0 \leq \tilde{s}$, there exists a generic constant c_1 such that*

$$\left\| \left(\Sigma_{t+1}^{(k)} \right)_{\bar{\mathcal{S}}} \right\|_0 \leq \tilde{s} \quad \text{and} \quad \left\| \Sigma_{t+1}^{(k)} - \Sigma^* \right\|_F \leq \frac{c_1 \lambda \sqrt{s^*}}{\rho_{2s^*+2\tilde{s}}^-}.$$

Proof: To simplify notation, we omit the stage index (k) . Considering Σ_{t+1} as the proximal Newton update, the system satisfies the following equation:

$$\text{mat}(\nabla^2 f(\Sigma_t) \text{vec}(\Sigma_{t+1} - \Sigma_t)) + \nabla f(\Sigma_t) + \Lambda \odot \Xi_{t+1} = \mathbf{0},$$

where $\Xi_{t+1} \in \partial \|\Sigma_{t+1}\|_1$.

From Lemma 25, under identical conditions, it holds that $\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2}$ and $|\mathcal{E}| \leq 2s^*$ for some set $\mathcal{E} \supseteq \mathcal{S}^*$. In the following, we decompose

$$\text{mat}(\nabla^2 f(\Sigma_t) \text{vec}(\Sigma_{t+1} - \Sigma_t)) + \nabla f(\Sigma_t)$$

into four parts:

- $\mathbf{V}_1 := \text{mat}(\nabla^2 f(\Sigma_t) \text{vec}(\Sigma_{t+1} - \Sigma^*))$,
- $\mathbf{V}_2 := \text{mat}(\nabla^2 f(\Sigma_t) \text{vec}(\Sigma^* - \Sigma_t))$,
- $\mathbf{V}_3 := \nabla f(\Sigma_t) - \nabla f(\Sigma^*)$,
- $\mathbf{V}_4 := \nabla f(\Sigma^*)$.

(1) For \mathbf{V}_1 , we define set

$$\mathcal{H}_1 = \left\{ (i, j) \in \bar{\mathcal{E}} \mid \left| (\mathbf{V}_1)_{ij} \right| \geq \frac{\lambda}{4} \right\}.$$

Referring to Lemma 33, we establish that

$$F_{\Lambda}(\Sigma_{t+1}) \leq F_{\Lambda}(\Sigma^*) + \frac{\lambda}{4} \|\Sigma_{t+1} - \Sigma^*\|_1.$$

This relationship leads to the following inequality:

$$\begin{aligned}
& f(\Sigma_{t+1}) - f(\Sigma^*) \\
& \leq \lambda (\|\Sigma^*\|_1 - \|\Sigma_{t+1}\|_1) + \frac{\lambda}{4} \|\Sigma_{t+1} - \Sigma^*\|_1 \\
& = \lambda (\|\Sigma_{\mathcal{E}}^*\|_1 - \|(\Sigma_{t+1})_{\mathcal{E}}\|_1 - \|(\Sigma_{t+1})_{\bar{\mathcal{E}}}\|_1) + \frac{\lambda}{4} \|\Sigma_{t+1} - \Sigma^*\|_1 \\
& \leq \frac{5\lambda}{4} \|(\Sigma_{t+1})_{\mathcal{E}} - \Sigma_{\mathcal{E}}^*\|_1 - \frac{3\lambda}{4} \|(\Sigma_{t+1})_{\bar{\mathcal{E}}} - \Sigma_{\bar{\mathcal{E}}}^*\|_1,
\end{aligned} \tag{32}$$

where the equality holds since $\Sigma_{\bar{\mathcal{E}}}^* = \mathbf{0}$. On the other hand, we have

$$\begin{aligned}
& f(\Sigma_{t+1}) - f(\Sigma^*) \\
& \stackrel{(i)}{\geq} \langle \nabla f(\Sigma^*), \Sigma_{t+1} - \Sigma^* \rangle \\
& \geq -\|\nabla f(\Sigma^*)\|_{\max} \|\Sigma_{t+1} - \Sigma^*\|_1 \\
& \stackrel{(ii)}{\geq} -\frac{\lambda}{4} \|\Sigma_{t+1} - \Sigma^*\|_1 \\
& = -\frac{\lambda}{4} \|(\Sigma_{t+1})_{\mathcal{E}} - \Sigma_{\mathcal{E}}^*\|_1 - \frac{\lambda}{4} \|(\Sigma_{t+1})_{\bar{\mathcal{E}}} - \Sigma_{\bar{\mathcal{E}}}^*\|_1,
\end{aligned} \tag{33}$$

where (i) is from the convexity of f and (ii) is from Assumption 9. Combining (32) and (33), we conclude

$$\|(\Sigma_{t+1})_{\bar{\mathcal{E}}} - \Sigma_{\bar{\mathcal{E}}}^*\|_1 \leq 3 \|(\Sigma_{t+1})_{\mathcal{E}} - \Sigma_{\mathcal{E}}^*\|_1.$$

Consider a subset $\mathcal{S}' \subset \mathcal{H}_1$ with $|\mathcal{S}'| = s' \leq \tilde{s}$. Select a vector $\mathbf{v} \in \mathbb{R}^{d^2}$ such that $\|\mathbf{v}\|_{\max} = 1$ and $\|\mathbf{v}\|_0 = s'$, satisfying the condition $s' \frac{\lambda}{4} \leq \mathbf{v}^\top \nabla^2 f(\Sigma_t) \text{vec}(\Sigma_{t+1} - \Sigma^*)$. Then, we have

$$\begin{aligned}
s' \frac{\lambda}{4} & \leq \mathbf{v}^\top \nabla^2 f(\Sigma_t) \text{vec}(\Sigma_{t+1} - \Sigma^*) \\
& \leq \left\| \mathbf{v}^\top (\nabla^2 f(\Sigma_t))^{\frac{1}{2}} \right\|_2 \left\| (\nabla^2 f(\Sigma_t))^{\frac{1}{2}} \text{vec}(\Sigma_{t+1} - \Sigma^*) \right\|_2 \\
& \stackrel{(i)}{\leq} c_1 \sqrt{\rho_{2s^*+2\tilde{s}}^+ \rho_{s'}^+} \|\mathbf{v}\|_2 \|\Sigma^* - \Sigma_{t+1}\|_F \\
& \stackrel{(ii)}{\leq} c_1 \sqrt{s' \rho_{2s^*+2\tilde{s}}^+ \rho_{s'}^+} \|\Sigma^* - \Sigma_{t+1}\|_F \\
& \stackrel{(iii)}{\leq} \frac{c_1 \sqrt{s' \rho_{2s^*+2\tilde{s}}^+ \rho_{s'}^+}}{\rho_{2s^*+2\tilde{s}}^-} \lambda \sqrt{s^*},
\end{aligned}$$

where (i) is from SE condition, (ii) is from the definition of \mathbf{v} , and (iii) is from Lemma 34. From above inequalities, it follows that

$$s' \leq \frac{c_1 \rho_{2s^*+2\tilde{s}}^+ \rho_{s'}^+ s^*}{(\rho_{2s^*+2\tilde{s}}^-)^2} \leq c_1 \varpi_{2s^*+2\tilde{s}}^2 s^*,$$

where $\varpi_{2s^*+2\tilde{s}} = \rho_{2s^*+2\tilde{s}}^+ / \rho_{2s^*+2\tilde{s}}^-$ is the condition number. Given that $s' = |\mathcal{S}'|$ achieves the maximum possible value such that $s' \leq \tilde{s}$ for any subset \mathcal{S}' of \mathcal{H}_1 . We conclude $\mathcal{S}' = \mathcal{H}_1$ to attain the maximum and thus $|\mathcal{H}_1| = s' \leq c_1 \varpi_{2s^*+2\tilde{s}}^2 s^*$.

(2) For \mathbf{V}_2 , we define set

$$\mathcal{H}_2 = \left\{ (i, j) \in \bar{\mathcal{E}} \mid \left| (\mathbf{V}_2)_{ij} \right| \geq \frac{\lambda}{4} \right\}.$$

Similar to the argument of \mathcal{H}_1 , we have $|\mathcal{H}_2| \leq c_2 \varpi_{2s^*+2\tilde{s}}^2 s^*$.

(3) For \mathbf{V}_3 , we define set

$$\mathcal{H}_3 = \left\{ (i, j) \in \bar{\mathcal{E}} \mid |(\mathbf{V}_3)_{ij}| \geq \frac{\lambda}{4} \right\}.$$

Consider a vector $\mathbf{V} \in \mathbb{R}^{d \times d}$ such that $\|\mathbf{V}\|_{\max} = 1$, $\|\mathbf{V}\|_0 = |\mathcal{H}_3|$, and

$$\begin{aligned} & \langle \mathbf{V}, \nabla f(\Sigma_t) - \nabla f(\Sigma^*) \rangle \\ &= \sum_{(i', j') \in \mathcal{H}_3} \mathbf{V}_{i'j'} (\nabla f(\Sigma_t) - \nabla f(\Sigma^*))_{i'j'} \\ &= \sum_{(i', j') \in \mathcal{H}_3} (\nabla f(\Sigma_t) - \nabla f(\Sigma^*))_{ij} \\ &\geq \frac{\lambda}{4}. \end{aligned} \quad (36)$$

Then we have

$$\begin{aligned} & \langle \mathbf{V}, \nabla f(\Sigma_t) - \nabla f(\Sigma^*) \rangle \\ &\leq \|\mathbf{V}\|_F \|\nabla f(\Sigma_t) - \nabla f(\Sigma^*)\|_F \\ &\stackrel{(i)}{\leq} \sqrt{|\mathcal{H}_3|} \|\nabla f(\Sigma_t) - \nabla f(\Sigma^*)\|_F \\ &\stackrel{(ii)}{\leq} \rho_{2s^*+2\tilde{s}}^+ \sqrt{|\mathcal{H}_3|} \|\Sigma_t - \Sigma^*\|_F, \end{aligned} \quad (37)$$

where (i) is from the definition of \mathbf{V} , and (ii) is from the mean value theorem and analogous argument for \mathcal{H}_3 .

Combining (36) and (37), we have

$$\begin{aligned} \lambda |\mathcal{H}_3| &\leq 4\rho_{2s^*+2\tilde{s}}^+ \sqrt{|\mathcal{H}_3|} \|\Sigma_t - \Sigma^*\|_F \\ &\leq 8\lambda \varpi_{2s^*+2\tilde{s}} \sqrt{3s^* |\mathcal{H}_3|}, \end{aligned}$$

where (i) is from Lemma 34. This implies $|\mathcal{H}_3| \leq c_3 \varpi_{2s^*+2\tilde{s}}^2 s^*$.

(4) For \mathbf{V}_4 , we define set

$$\mathcal{H}_4 = \left\{ (i, j) \in \bar{\mathcal{E}} \mid |(\mathbf{V}_4)_{ij}| \geq \frac{\lambda}{4} \right\}.$$

By assumption 9, we have

$$\begin{aligned} 0 \leq |\mathcal{H}_4| &\leq \sum_{(i,j) \in \bar{\mathcal{E}}} \frac{4}{\lambda} |(\nabla f(\Sigma^*))_{ij}| \cdot \mathbb{I} \left(|(\nabla f(\Sigma^*))_{ij}| > \frac{\lambda}{4} \right) \\ &= \sum_{(i,j) \in \bar{\mathcal{E}}} \frac{4}{\lambda} |(\nabla f(\Sigma^*))_{ij}| \cdot 0 = 0. \end{aligned} \quad (38)$$

Combining the results for Set $\mathcal{H}_1 \sim \mathcal{H}_4$, we have that there exists some constant c_0 such that

$$\left\| \left(\Sigma_{t+\frac{1}{2}} \right)_{\bar{\mathcal{S}}^*} \right\|_0 \leq c_0 \kappa_{2s^*+2\tilde{s}}^2 s^* \leq \tilde{s}.$$

■

Lemma 37. Suppose Assumptions 8 and 9 hold. If $\Sigma^{(t)} \in \mathcal{B}(\Sigma^*, r)$ and $\left\| \Sigma_{\bar{\mathcal{S}}^*}^{(t+1)} \right\|_0 \leq \tilde{s}$, then for each stage $k \geq 2$, we have

$$\left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \leq \frac{\tau \kappa^3}{\rho_{2s^*+2\tilde{s}}} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2.$$

Proof: First, we reformulate the prox-Newton update (11) as formula (34). By the Lemma 36, we have $\left\| \left(\Sigma_{t+1}^{(k)} \right)_{\bar{\mathcal{S}}^*} \right\|_0 \leq \tilde{s}$. And by the KKT condition, we have formula (35). By the strictly non-expansive property of the proximal operator, we obtain formula (39).

Note that both $\|\Sigma_{t+1}\|_0 \leq \tilde{s}$ and $\left\| \hat{\Sigma}^{(k)} \right\|_0 \leq \tilde{s}$. On the other hand, from the SE properties, we have

$$\left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_{\nabla^2 f(\hat{\Sigma}^{(k)})}^2 \geq \rho_{2s^*+2\tilde{s}}^- \left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2. \quad (40)$$

Combining (39) and (40), we have formula (41). Then the proof is finished. ■

In the subsequent analysis, it is essential to utilize the property that the iterates $\Sigma_t \in \mathcal{B}(\hat{\Sigma}^{(k)}, 2r)$ rather than $\Sigma_t \in \mathcal{B}(\Sigma^*, r)$ for the convergence analysis of the proximal Newton method. This property is valid due to the simultaneous inclusion of $\Sigma_t \in \mathcal{B}(\Sigma^*, r)$ and $\hat{\Sigma}^{(k)} \in \mathcal{B}(\Sigma^*, r)$. Consequently, it follows that $\Sigma_t \in \mathcal{B}(\hat{\Sigma}^{(k)}, 2r)$, where $2r = \frac{\rho_{2s^*+2\tilde{s}}^-}{2\tau\kappa}$ represents the radius of the quadratic convergence region for the proximal Newton algorithm.

Lemma 38. Suppose that Assumptions 1, 4, 7, 8, and 7 hold. If $\Sigma_t \in \mathcal{B}(\hat{\Sigma}^{(k)}, 2r)$ and $\|(\Sigma_t)_{\bar{\mathcal{S}}^*}\|_0 \leq \tilde{s}$ at each stage $k \geq 2$ with $\alpha = 0.3$, then $\eta_t = 1$. Further we have

$$F(\Sigma_{t+1}^{(k)}) \leq F(\Sigma_t^{(k)}) + \frac{1}{4} \delta_t.$$

$$\Sigma_{t+1}^{(k)} = \operatorname{argmin}_{\Sigma} \left\langle \nabla f(\Sigma_t^{(k)}), \Sigma - \Sigma_t^{(k)} \right\rangle + \frac{1}{2} \left\| \Sigma - \Sigma_t^{(k)} \right\|_{\nabla^2 f(\Sigma_t^{(k)})}^2 + \left\| \Lambda^{(k-1)} \odot \Sigma \right\|_1 \quad (34)$$

$$\hat{\Sigma}^{(k)} = \operatorname{argmin}_{\Sigma} \left\langle \nabla f(\hat{\Sigma}^{(k)}), \Sigma - \hat{\Sigma}^{(k)} \right\rangle + \frac{1}{2} \left\| \Sigma - \hat{\Sigma}^{(k)} \right\|_{\nabla^2 f(\Sigma_t^{(k)})}^2 + \left\| \Lambda^{(k-1)} \odot \hat{\Sigma}^{(k)} \right\|_1 \quad (35)$$

$$\begin{aligned}
\left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_{\nabla^2 f(\hat{\Sigma}^{(k)})}^2 &\leq \text{vec}^\top \left(\Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right) \left[\nabla^2 f \left(\Sigma_t^{(k)} \right) \text{vec} \left(\Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right) + \left(\nabla f \left(\hat{\Sigma}^{(k)} \right) - \nabla f \left(\Sigma_t^{(k)} \right) \right) \right] \\
&\leq \left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \left\| \nabla^2 f \left(\Sigma_t^{(k)} \right) \text{vec} \left(\Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right) + \left(\nabla f \left(\hat{\Sigma}^{(k)} \right) - \nabla f \left(\Sigma_t^{(k)} \right) \right) \right\|_2.
\end{aligned} \tag{39}$$

$$\begin{aligned}
&\left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \\
&\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \left\| \nabla^2 f \left(\Sigma_t^{(k)} \right) \text{vec} \left(\Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right) + \left(\nabla f \left(\hat{\Sigma}^{(k)} \right) - \nabla f \left(\Sigma_t^{(k)} \right) \right) \right\|_2 \\
&= \frac{1}{\rho_{2s^*+2\tilde{s}}} \left\| \int_0^1 \left[\nabla^2 f \left(\Sigma_t^{(k)} + \theta \left(\hat{\Sigma}^{(k)} - \Sigma_t^{(k)} \right) \right) - \nabla^2 f \left(\Sigma_t^{(k)} \right) \right] \cdot \text{vec} \left(\hat{\Sigma}^{(k)} - \Sigma_t^{(k)} \right) d\theta \right\|_2 \\
&\leq \frac{1}{\rho_{2s^*+2\tilde{s}}} \int_0^1 \left\| \left[\nabla^2 f \left(\Sigma_t^{(k)} + \theta \left(\hat{\Sigma}^{(k)} - \Sigma_t^{(k)} \right) \right) - \nabla^2 f \left(\Sigma_t^{(k)} \right) \right] \cdot \text{vec} \left(\hat{\Sigma}^{(k)} - \Sigma_t^{(k)} \right) \right\|_2 d\theta \\
&\stackrel{(i)}{\leq} \frac{\tau \kappa^3}{\rho_{2s^*+2\tilde{s}}} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2,
\end{aligned} \tag{41}$$

where (i) is from the local restricted Hessian smoothness of f .

Proof: Recall $\Sigma_{t+\frac{1}{2}}$ defined in (11) and denote $\Delta \Sigma_t = \Sigma_{t+\frac{1}{2}} - \Sigma_t$. Then we have

$$\begin{aligned}
&\left\| \Delta \Sigma_t^{(k)} \right\|_F^2 \\
&\stackrel{(i)}{\leq} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F + \left\| \Sigma_{t+\frac{1}{2}}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \\
&\stackrel{(ii)}{\leq} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F + \frac{\tau \kappa^3}{\rho_{2s^*+2\tilde{s}}} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2 \\
&\stackrel{(iii)}{\leq} \frac{3}{2} \left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F,
\end{aligned} \tag{42}$$

where (i) is from triangle inequality, (ii) is from Lemma 37, and (iii) is from $\left\| \Sigma_t^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \leq R \leq \frac{\rho_{2s^*+2\tilde{s}}}{2\tau \kappa^3}$.

By Lemma 36, we have $\|(\Sigma_t)_{\bar{S}^*}\|_0 \leq 2\tilde{s}$. Then by expanding F , we have formula (43). The $\eta_t = 1$ can be demonstrated analogously from the analysis in [55] and [57], thus we omit it. ■

Lemma 39. Suppose that Assumptions 1, 4, 7, 8, and 9 hold. To achieve the approximate KKT condition $\omega(\Sigma_t) \leq \varepsilon$ for any $\varepsilon > 0$ at each stage $k \geq 2$, the number of iteration for proximal Newton updates is at most

$$\log \log (3\rho_{2s^*+2\tilde{s}}^+/\varepsilon).$$

Proof: Based on the solution Σ_{t-1} obtained from the $(t-1)$ -th iteration, we consider the optimal solution at t -th iteration. By the KKT condition, we have

$$\text{mat}(\nabla^2 f(\Sigma_{t-1}) \text{vec}(\Sigma_t - \Sigma_{t-1})) + \nabla f(\Sigma_{t-1}) + \Lambda \odot \Xi_t = \mathbf{0},$$

where $\Xi_t \in \partial \|\Sigma_t\|_1$. Given any $\mathbf{V} \in \mathbb{R}^{d \times d}$ satisfying $\|\mathbf{V}\|_2 \leq \|\mathbf{V}\|_1 = 1$ and $\|\mathbf{V}\|_0 \leq 2s^* + 2\tilde{s}$, we derive formula (44). By taking the supremum of the L.H.S of (44) with respect to \mathbf{V} , we have

$$\|\nabla f(\Sigma_t) + \Lambda \odot \Xi_t\|_{\max} \leq 2\rho_{2s^*+2\tilde{s}}^+ \|\Sigma_t - \Sigma_{t-1}\|_F. \tag{45}$$

Then from Lemma (37), we have

$$\begin{aligned}
&\left\| \Sigma_{t+1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \\
&\leq (\tau \kappa^3 / \rho_{2s^*+2\tilde{s}}^-)^{1+2+4+\dots+2^{t-1}} \left\| \Sigma_0^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^{2^t} \\
&\leq \left((\tau \kappa^3 / \rho_{2s^*+2\tilde{s}}^-) \left\| \Sigma_0^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2 \right)^{2^t}.
\end{aligned}$$

By (44) and (42) by taking $\Delta \Sigma_{t-1} = \Sigma_t - \Sigma_{t-1}$, we obtain

$$\begin{aligned}
&\omega_{\Lambda^{(k-1)}}(\Sigma_t^{(k)}) \\
&\leq 2\rho_{2s^*+2\tilde{s}}^+ \left\| \Sigma_t^{(k)} - \Sigma_{t-1}^{(k)} \right\|_F \\
&\leq 3\rho_{2s^*+2\tilde{s}}^+ \left\| \Sigma_{t-1}^{(k)} - \hat{\Sigma}^{(k)} \right\|_F \\
&\leq 3\rho_{2s^*+2\tilde{s}}^+ \left((\tau \kappa^3 / \rho_{2s^*+2\tilde{s}}^-) \left\| \Sigma_0^{(k)} - \hat{\Sigma}^{(k)} \right\|_F^2 \right)^{2^t}.
\end{aligned}$$

$$\begin{aligned}
& F\left(\Sigma_t^{(k)} + \Delta \Sigma_t^{(k)}\right) - F\left(\Sigma_t^{(k)}\right) \\
&= f\left(\Sigma_t^{(k)} + \Delta \Sigma_t^{(k)}\right) - f\left(\Sigma_t^{(k)}\right) + \left\| \Lambda^{(k-1)} \odot \left(\Sigma_t^{(k)} + \Delta \Sigma_t^{(k)}\right) \right\|_1 - \left\| \Lambda^{(k-1)} \odot \Sigma_t^{(k)} \right\|_1 \\
&\stackrel{(i)}{\leq} \left\langle \nabla f\left(\Sigma_t^{(k)}\right), \Delta \Sigma_t^{(k)} \right\rangle + \frac{1}{2} \left\| \Delta \Sigma_t^{(k)} \right\|_{\nabla^2 f\left(\Sigma_t^{(k)}\right)}^2 + \frac{\tau \kappa^3}{3} \left\| \Delta \Sigma_t^{(k)} \right\|_F^3 + \left\| \Lambda^{(k-1)} \odot \left(\Sigma_t^{(k)} + \Delta \Sigma_t^{(k)}\right) \right\|_1 - \left\| \Lambda^{(k-1)} \odot \Sigma_t^{(k)} \right\|_1 \\
&\stackrel{(ii)}{\leq} \delta_t - \frac{1}{2} \delta_t + \frac{\tau \kappa^3}{3} \left\| \Delta \Sigma_t^{(k)} \right\|_F^3 \stackrel{(iii)}{\leq} \frac{1}{2} \delta_t + \frac{\tau \kappa^3}{3 \rho_{2s^*+2\tilde{s}}} \left\| \Delta \Sigma_t^{(k)} \right\|_{\nabla^2 f(\Sigma)}^2 \left\| \Delta \Sigma_t^{(k)} \right\|_F \\
&\stackrel{(iv)}{\leq} \left(\frac{1}{2} - \frac{\tau \kappa^3}{3 \rho_{2s^*+2\tilde{s}}} \left\| \Delta \Sigma_t^{(k)} \right\|_F \right) \delta_t \stackrel{(v)}{\leq} \frac{1}{4} \delta_t,
\end{aligned} \tag{43}$$

where (i) is from the restricted Hessian smooth condition, (ii) and (iv) are from Lemma 35, (iii) is from the same argument of (40), and (v) is from (42), $\delta_t < 0$, and $\left\| \Sigma_t^{(k)} - \widehat{\Sigma}^{(k)} \right\|_F \leq r \leq \frac{\rho_{2s^*+2\tilde{s}}}{2\tau\kappa^3}$.

$$\begin{aligned}
& \langle \nabla f(\Sigma_t) + \Lambda \odot \Xi_t, \mathbf{V} \rangle \\
&= \langle \nabla f(\Sigma_t) - \nabla f(\Sigma_{t-1}) - \text{mat}(\nabla^2 f(\Sigma_{t-1}) \text{vec}(\Sigma_t - \Sigma_{t-1})), \mathbf{V} \rangle \\
&= \langle \nabla f(\Sigma_t) - \nabla f(\Sigma_{t-1}), \mathbf{V} \rangle - (\nabla^2 f(\Sigma_{t-1}) \text{vec}(\Sigma_t - \Sigma_{t-1}))^\top \text{vec}(\mathbf{V}) \\
&\stackrel{(i)}{\leq} \left\| \left(\nabla^2 f(\check{\Sigma}) \right)^{\frac{1}{2}} \text{vec}(\Sigma_t - \Sigma_{t-1}) \right\|_2 \cdot \left\| \text{vec}^\top(\mathbf{V}) \left(\nabla^2 f(\check{\Sigma}) \right)^{\frac{1}{2}} \right\|_2 \\
&\quad + \left\| \left(\nabla^2 f(\Sigma_{t-1}) \right)^{\frac{1}{2}} \text{vec}(\Sigma_t - \Sigma_{t-1}) \right\|_2 \cdot \left\| \text{vec}^\top(\mathbf{V}) \left(\nabla^2 f(\Sigma_{t-1}) \right)^{\frac{1}{2}} \right\|_2 \\
&\stackrel{(ii)}{\leq} 2\rho_{2s^*+2\tilde{s}}^+ \left\| \Sigma_t - \Sigma_{t-1} \right\|_F,
\end{aligned} \tag{44}$$

where (i) is from mean value theorem with some $\check{\Sigma} = (1 - \theta) \Sigma_{t-1} + \theta \Sigma_t$ for some $\theta \in [0, 1]$ and Cauchy-Schwarz inequality, and (ii) is from the SE properties.

By requiring the R.H.S equal to ε we obtain

$$\begin{aligned}
t &= \log \left(\log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right) / \log \left(\frac{\rho_{2s^*+2\tilde{s}}^-}{\tau \kappa^3 \left\| \Sigma_0^{(k)} - \widehat{\Sigma}^{(k)} \right\|_F} \right) \right) \\
&= \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right) - \log \log \left(\frac{\rho_{2s^*+2\tilde{s}}^-}{\tau \kappa^3 \left\| \Sigma_0^{(k)} - \widehat{\Sigma}^{(k)} \right\|_F} \right) \\
&\stackrel{(i)}{\leq} \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right) - \log \log 4 \\
&\leq \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right),
\end{aligned}$$

where (i) is from the fact that $\left\| \Sigma_0^{(k)} - \widehat{\Sigma}^{(k)} \right\|_F \leq r = \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau\kappa^3}$. ■

D. Proof of Theorem 17

Lemma 40. Suppose that Assumptions 1, 4, 7, 8, and 9 hold. After sufficiently many iterations $T < \infty$, the following results

hold for all $t \geq T$:

$$\left\| (\Sigma_t)_{\overline{\mathcal{S}^*}} \right\|_0 \leq \tilde{s} \text{ and } F_{\Lambda^{(0)}}\left(\Sigma_t^{(1)}\right) \leq F_{\Lambda^{(0)}}\left(\Sigma^*\right) + \frac{15\lambda^2 s^*}{4\rho_{2s^*+2\tilde{s}}^-}.$$

Proof: Building upon the framework established in Lemma 35, Lemma 38, it can be concluded that the objective $F_{\Lambda^{(0)}}$ exhibits a sufficient decrease in each iteration of the proximal Newton step, as also discussed in [58]. Consequently, there exists a constant T such that for all $t \geq T$, we have

$$F_{\Lambda^{(0)}}\left(\Sigma_t^{(1)}\right) \leq F_{\Lambda^{(0)}}\left(\Sigma^*\right) + \frac{\lambda}{4} \left\| \Sigma_t^{(1)} - \Sigma^* \right\|_1,$$

where $\left\| \Sigma_t^{(1)} - \Sigma^* \right\|_1 \leq \frac{c_1 \lambda \sqrt{s^*}}{\rho_{2s^*+2\tilde{s}}^-}$ derived from a similar argument presented in the proof of Lemma 25. The subsequent analysis follows a pattern analogous to that of Lemma 36. Let $S_n = \left\{ (i, j) \mid \left| \nabla f(\check{\Sigma})_{ij} \right| = \lambda \right\}$ and consider the set $\left\{ (i, j) \mid (\check{\Sigma})_{ij} \neq 0 \right\} \subseteq S_n$. It is required to show $|S_n| \leq s^* + \tilde{s}$. To demonstrate this, we decompose S_n into two distinct parts:

- $S_n^1 : \left\{ (i, j) \notin \overline{\mathcal{S}^*} \mid \left| \left(\nabla f(\check{\Sigma}) - \nabla f(\Sigma^*) \right)_{ij} \right| \geq \frac{\lambda}{2} \right\},$
- $S_n^2 : \left\{ (i, j) \notin \overline{\mathcal{S}^*} \mid \left| \nabla f(\Sigma^*)_{ij} \right| > \frac{\lambda}{2} \right\}.$

We have $S_n \subseteq \mathcal{S}^* \cup S_n^1 \cup S_n^2$ and establish upper bounds for each.

(1) For S_n^1 , consider S' with maximum size $s' = |S'| \leq \tilde{s}$ such that $S' \subseteq S_n^1$. Then there exists a matrix \mathbf{V} such that $\|\mathbf{V}\|_{\max} = 1$ and $\|\mathbf{V}\|_0 = s'$, satisfying

$$\frac{\lambda s'}{2} \leq \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \mathbf{V} \rangle.$$

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

$$\begin{aligned} & \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*) \\ &= \text{mat} \left(\nabla^2 f \left(\theta \tilde{\Sigma} + (1 - \theta) \Sigma^* \right) \text{vec} \left(\tilde{\Sigma} - \Sigma^* \right) \right). \end{aligned}$$

For simplicity, we define $\mathbf{H} = \nabla^2 f \left(\rho \tilde{\Sigma} + (1 - \rho) \Sigma^* \right)$. Applying the Cauchy-Schwartz inequality, we have

$$\begin{aligned} \frac{\lambda s'}{2} &\leq \text{vec}^\top(\mathbf{V}) \mathbf{H} \text{vec}(\tilde{\Sigma} - \Sigma^*) \\ &\leq \underbrace{\left\| \text{vec}^\top(\mathbf{V}) (\mathbf{H})^{\frac{1}{2}} \right\|_2}_{\text{I}} \underbrace{\left\| (\mathbf{H})^{\frac{1}{2}} \text{vec}(\tilde{\Sigma} - \Sigma^*) \right\|_2}_{\text{II}}. \end{aligned} \quad (46)$$

We now establish upper bounds for terms **I** and **II** respectively. Let $\tilde{\Sigma}, \Sigma^* \in \mathcal{B}(\Sigma^*, r)$, any convex combination of $\tilde{\Sigma}$ and Σ^* also falls in $\mathcal{B}(\Sigma^*, r)$. The localized sparse eigenvalue condition can be applied to \mathbf{H} .

i) For term **I**, based on Definition (3) and Assumption 4, we have

$$\begin{aligned} \left\| \text{vec}^\top(\mathbf{U}) (\mathbf{H})^{\frac{1}{2}} \right\|_2 &\leq \sqrt{\rho_{2s^*+2\tilde{s}}^+} \|\mathbf{U}\|_F^2 \\ &\leq \sqrt{\rho_{2s^*+2\tilde{s}}^+} (\|\mathbf{U}\|_1 \|\mathbf{U}\|_{\max})^{\frac{1}{2}} \\ &\leq \sqrt{\rho_{2s^*+2\tilde{s}}^+} \sqrt{s'}. \end{aligned}$$

ii) For **II**, using Lemma 32, we have

$$\begin{aligned} & \left\| (\mathbf{H})^{\frac{1}{2}} \text{vec}(\tilde{\Sigma} - \Sigma^*) \right\|_2 \\ &\leq \langle \nabla f(\tilde{\Sigma}) - \nabla f(\Sigma^*), \tilde{\Sigma} - \Sigma^* \rangle \\ &\leq \frac{c_1 \lambda^2 s^*}{\rho_{2s^*+2\tilde{s}}^-}. \end{aligned}$$

By substituting the bounds for **I** and **II** back into (46), we obtain

$$\frac{\lambda s'}{2} \leq \sqrt{\rho_{2s^*+2\tilde{s}}^+} \sqrt{s'} \times \frac{c_1 \lambda \sqrt{s^*}}{\rho_{2s^*+2\tilde{s}}^-}.$$

Multiplying both sides of this inequality by $(\frac{\lambda}{2})^{\frac{1}{2}}$ and squaring gives us

$$s' \leq \frac{4\rho_{2s^*+2\tilde{s}}^+ c_2 s^*}{(\rho_{2s^*+2\tilde{s}}^-)^2} < \tilde{s}, \quad (47)$$

where the last inequality follows from our assumption. Since $s' = |S'|$ achieves the maximum possible value such that $s' \leq \tilde{s}$ for any subset S' of S_n^1 and 47 shows that $s' < \tilde{s}$, we conclude that $S' = S_n^1$, and thus

$$|S_n^1| = s' \leq \left\lfloor \frac{4\rho_{2s^*+2\tilde{s}}^+ c_2 s^*}{(\rho_{2s^*+2\tilde{s}}^-)^2} \right\rfloor < \tilde{s}.$$

(2) For S_n^2 , Given $\|\nabla f(\Sigma^*)\|_{\max} + \varepsilon \leq \frac{\lambda}{4}$, it follows that $S_n^2 = \emptyset$ and thus $|S_n^2| = 0$.

Combining all the results, we have $\left\| \Sigma_{\mathcal{S}^*}^{(t)} \right\|_0 \leq \tilde{s}$. ■

Lemma 41. Suppose that Assumptions 4, 8 and 9 hold. If $\left\| \left(\Sigma_t^{(1)} \right)_{\mathcal{S}^*} \right\|_0 \leq \tilde{s}$, and $F_{\Lambda^{(0)}} \left(\Sigma_t^{(1)} \right) \leq F_{\Lambda^{(0)}} \left(\Sigma^* \right) + \frac{15\lambda^2 s^*}{4\rho_{2s^*+2\tilde{s}}^-}$, we have

$$\left\| \Sigma_t^{(1)} - \Sigma^* \right\|_2 \leq \frac{c_1 \lambda \sqrt{s^*}}{\rho_{2s^*+2\tilde{s}}^-}$$

for some constant c_1 and

$$\left\| \Sigma_{t+1}^{(1)} - \hat{\Sigma}^{(1)} \right\|_2 \leq \frac{\tau \kappa^3}{\rho_{2s^*+2\tilde{s}}^-} \left\| \Sigma_t^{(1)} - \hat{\Sigma}^{(1)} \right\|_2^2.$$

Proof: The estimation error is derived analogously from [28] and [57], thus we omit it here. The claim of quadratic convergence follows directly from Lemma (37) given sparse solutions. ■

Lemma 42. Suppose that Assumptions 4, 8, and 9 hold. If $\left\| \left(\Sigma_T^{(1)} \right)_{\mathcal{S}^*} \right\|_0 \leq \tilde{s}$, and $F_{\Lambda^{(0)}} \left(\Sigma_T^{(1)} \right) \leq F_{\Lambda^{(1)}} \left(\Sigma^* \right) + \frac{15\lambda^2 s^*}{4\rho_{2s^*+2\tilde{s}}^-}$ at some iteration T , we need at most

$$T_1 \leq \log \log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right)$$

extra iterations of the proximal Newton updates such that $\omega_{\Lambda^{(0)}} \left(\Sigma_{T+T_1}^{(1)} \right) \leq \frac{\lambda}{8}$.

Proof: The maximum of the number of iterations for the proximal Newton update is determined by integrating Lemma (40) and Lemma (39). Note that

$$T_1 \leq \log \frac{\log \left(\frac{3\rho_{2s^*+2\tilde{s}}^+}{\varepsilon} \right)}{\log \left(\frac{\rho_{2s^*+2\tilde{s}}^-}{\tau \kappa^3 \left\| \Sigma^{(T+1)} - \hat{\Sigma}^{(1)} \right\|_F} \right)}.$$

Then we obtain the result from $\left\| \Sigma_{T+1}^{(1)} - \hat{\Sigma}^{(1)} \right\|_F \leq r = \frac{\rho_{2s^*+2\tilde{s}}^-}{4\tau \kappa}$. ■

APPENDIX E SUPPLEMENTARY TABLES

TABLE I
QUANTITATIVE COMPARISON AMONG DIFFERENT PENALTIES FOR THE BANDED MATRIX

	ℓ_1				MCP					ℓ_1				MCP			
	$m = 100$	$m = 500$	$m = 1500$	$m = 3000$	$m = 100$	$m = 500$	$m = 1500$	$m = 3000$		$m = 100$	$m = 500$	$m = 1500$	$m = 3000$	$m = 100$	$m = 500$	$m = 1500$	$m = 3000$
	$d = 80, n = 50$									$d = 80, n = 100$							
FAE	27.7825 (0.2616)	24.1565 (0.1648)	22.1561 (0.1246)	19.9998 (0.1051)	24.3152 (0.2551)	20.1542 (0.1566)	19.4295 (0.1047)	18.6415 (0.1015)	26.5923 (0.2436)	23.8546 (0.1534)	21.4165 (0.1132)	19.4620 (0.1052)	23.7515 (0.2435)	19.5612 (0.1530)	18.5164 (0.1106)	17.8962 (0.0945)	
FRE	1.3876 (0.2546)	1.2065 (0.1591)	1.1066 (0.1198)	0.9989 (0.9895)	1.2144 (0.2357)	1.0066 (0.1296)	0.9704 (0.1210)	0.9311 (0.0890)	1.3282 (0.1973)	1.1914 (0.1275)	1.0697 (0.1016)	0.9720 (0.0999)	1.1863 (0.1773)	0.9770 (0.1176)	0.9248 (0.1012)	0.8938 (0.1001)	
FPR	0.1615 (0.0702)	0.1528 (0.0765)	0.1137 (0.0686)	0.1053 (0.0552)	0.1256 (0.0699)	0.1038 (0.0438)	0.0681 (0.0423)	0.0311 (0.0217)	0.1572 (0.0425)	0.1315 (0.0407)	0.1152 (0.0319)	0.0873 (0.0287)	0.1137 (0.0366)	0.0852 (0.0287)	0.0489 (0.0209)	0.0241 (0.0174)	
TPR	0.8072 (0.0368)	0.8491 (0.0382)	0.8774 (0.0397)	0.8955 (0.0399)	0.8413 (0.0397)	0.8618 (0.0407)	0.8957 (0.0412)	0.9214 (0.0420)	0.8117 (0.0371)	0.8546 (0.0392)	0.8810 (0.0401)	0.9016 (0.0408)	0.8455 (0.0400)	0.8821 (0.0412)	0.9016 (0.0418)	0.9425 (0.0431)	
$d = 100, n = 50$																	
FAE	28.1276 (0.2492)	26.5312 (0.2151)	23.1562 (0.1813)	20.1532 (0.1513)	24.8612 (0.2263)	21.0017 (0.2052)	20.1562 (0.1673)	18.8964 (0.1233)	27.8613 (0.2161)	24.2361 (0.2111)	22.2613 (0.1562)	19.9131 (0.1162)	24.2315 (0.2023)	21.1156 (0.1762)	19.8216 (0.1230)	18.1566 (0.1055)	
FRE	1.1012 (0.2435)	1.0387 (0.2051)	0.9066 (0.1736)	0.7890 (0.1435)	0.9734 (0.2315)	0.8223 (0.2184)	0.7892 (0.1763)	0.7398 (0.1256)	1.0908 (0.2251)	0.9489 (0.2056)	0.8716 (0.1813)	0.7796 (0.1562)	0.9487 (0.2130)	0.8267 (0.1504)	0.7761 (0.1261)	0.7109 (0.1023)	
FPR	0.1511 (0.0851)	0.1341 (0.0793)	0.1056 (0.0712)	0.0915 (0.0663)	0.1105 (0.0714)	0.1004 (0.0627)	0.0557 (0.0487)	0.0303 (0.0226)	0.1462 (0.0487)	0.1227 (0.0415)	0.1005 (0.0326)	0.0905 (0.0291)	0.1126 (0.0392)	0.0756 (0.0311)	0.0448 (0.0265)	0.0131 (0.0199)	
TPR	0.8157 (0.0337)	0.8581 (0.0391)	0.9167 (0.0401)	0.9294 (0.0413)	0.8780 (0.0402)	0.8860 (0.0412)	0.9070 (0.0422)	0.9421 (0.0430)	0.8307 (0.0380)	0.8706 (0.0401)	0.8913 (0.0409)	0.9212 (0.0411)	0.8605 (0.0416)	0.8915 (0.0423)	0.9172 (0.0427)	0.9512 (0.0436)	