

Large Covariance Matrix Estimation for Groups of Highly Correlated Variables via Nonconvex Optimization

Shanshan Zou

*School of Information Science and Technology
ShanghaiTech University
Shanghai, China
zoushsh2022@shanghaitech.edu.cn*

Ziping Zhao

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

Abstract—This paper addresses the problem of covariance matrix estimation in scenarios where the underlying variables can be divided into groups, with variables within each group being highly correlated. The resulting covariance matrix exhibits both sparse structures and approximate low-rank properties due to the presence of highly correlated blocks. By appropriately permuting the variables, this covariance matrix can be approximately transformed to a block diagonal form. In this work, we investigate the estimation of covariance matrices under this group structure in high dimensions. We propose a least squares based covariance estimation framework with a trace norm and nonconvex sparsity regularizer to promote both low-rankness and sparsity. Additionally, we introduce a spectral constraint to ensure the positive semi-definiteness of the covariance matrix, even in finite samples, while allowing for the incorporation of prior spectral information. To solve this nonconvex statistical estimation problem, we develop an algorithm based on the majorization-minimization (MM) framework, which iteratively solves a sequence of convex subproblems via the alternating direction method of multipliers (ADMM). We provide theoretical guarantees, proving that the proposed algorithm converges to an estimator that achieves the oracle statistical rate under mild technical assumptions. The theoretical results are supported by numerical experiments.

Index Terms—Covariance matrix estimation, joint sparse and low-rank, block diagonal, correlated clusters, nonconvex statistical optimization.

I. INTRODUCTION

Covariance matrix estimation is a fundamental problem in multivariate data analysis, with applications spanning a variety of fields such as statistics [1], biology [2], finance [3], [4], signal processing [5], [6], and machine learning [7]. The sample covariance matrix is commonly used due to its simplicity, but it is known to perform poorly in high dimensions, where the number of variables exceeds or is comparable to the sample size [8]. To improve estimation accuracy in such regimes, various structural assumptions have been explored, including sparsity [9]–[11], low-rank [12]–[14], sparse plus low-rank [15], [16], banded [17], Toeplitz [18], and Hankel [19] structures.

In high dimensional settings, one of the most widely adopted structural assumptions for covariance matrix estimation is sparsity. Common methods for sparse covariance estimation include regularized techniques such as thresholding [9], [10], banding [20], [21], and tapering [17], [22]. However, these methods do not inherently guarantee a positive semi-definite estimate in finite samples. Therefore, ensuring the positive semi-definiteness of the estimated matrix becomes a key consideration. In [23], a positive definite covariance estimator was introduced by incorporating a logarithmic barrier function into the soft-thresholding estimator [10]. Nevertheless, this approach requires the smallest eigenvalue of the true covariance matrix to be bounded away from zero; otherwise, the influence of the perturbation from the log-determinant barrier becomes difficult to control. To address this limitation, subsequent works [24], [25] proposed using a positive definite constraint, which eliminates the need for the log-determinant

barrier penalty. Furthermore, this constraint can be generalized to a more general spectral constraint [26], [27], allowing the incorporation of prior spectral information about the true covariance matrix.

In this paper, we focus on high-dimensional covariance matrix estimation where the underlying variables can be divided into groups of highly correlated variables, forming correlated clusters—a structure commonly encountered in real-world applications. For example, in biology, genes can be grouped into pathways, with stronger connections within a pathway than between pathways [28]. Similarly, in financial markets, the stocks can be grouped by asset classes [29]. Various methods, such as clustering based on prior information [29], k -means clustering [30], and hierarchical clustering [31], can be used to detect this group (or cluster) information. By permuting the variables based on this group information, the resulting covariance matrix becomes block diagonal, a special case of sparsity structure. Additionally, the covariance matrix exhibits approximate low-rankness, as the blocks formed by highly correlated variables tend to be approximately low-rank.

For estimating such block diagonal covariance matrices for groups of highly correlated variables, traditional sparse covariance estimation methods are not suited since they cannot well capture the block diagonal structures. This highlights the need for specialized algorithms designed to handle this problem. In [32], [33], convex optimization approaches based on penalized least squares were examined, where a trace penalty [13] is incorporated into the soft-thresholding sparse covariance estimator to promote low-rankness. Additionally, a positive semi-definite constraint is imposed to ensure the covariance matrix remains positive semi-definite. In [32], generalized forward-backward and incremental proximal descent methods were employed to efficiently solve the optimization problem, while the alternating direction method of multipliers (ADMM) was used in [33]. In [33], it was proved that the estimator achieves a statistical rate of $O_p\left(\sqrt{\frac{s^* \log d}{n}}\right)$, where n is the number of samples, d is the dimension, and s^* represents the number of non-zero elements. This rate aligns with the minimax rate for sparse covariance matrix estimation [21], [23], [24].

In this paper, we propose a nonconvex formulation for covariance matrix estimation in scenarios where variables are grouped into highly correlated clusters. It is well-known that the convex ℓ_1 -norm penalty used in soft-thresholding estimators introduces a non-negligible estimation bias [34]. To address this, we propose using a nonconvex penalty to estimate the covariance matrix. Additionally, we incorporate a generic spectral constraint, which offers greater flexibility when prior information—such as the minimum and maximum eigenvalues of the true covariance matrix is available. To solve this estimation problem, we develop an algorithm based on majorization-

minimization (MM). The algorithm handles the nonconvex problem by iteratively solving a sequence of convex subproblems, each of which is solved using the ADMM method. We rigorously establish the statistical properties of the proposed MM based algorithm and prove that the resulting estimator achieves the oracle statistical rate of $O_p\left(\sqrt{\frac{s^*}{n}}\right)$ in the Frobenius norm, outperforming the rate of $O_p\left(\sqrt{\frac{s^* \log d}{n}}\right)$ in [33]. Numerical experiment results support the theoretical findings and validate the superiority of our algorithm over existing methods.

II. PROBLEM FORMULATION

Given a collection of observations $\{\mathbf{x}_i\}_{i=1}^n$ from a zero-mean random variable $\mathbf{x} \in \mathbb{R}^d$, the sample covariance matrix (SCM) is defined as

$$\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top.$$

To estimate a block diagonal covariance for groups of highly correlated variables, we propose the following least squares covariance fitting formulation with joint low-rank and sparse regularization:

$$\underset{\Sigma \in \mathcal{E}}{\text{minimize}} \quad \frac{1}{2} \|\Sigma - \mathbf{S}\|_F^2 + \tau \text{tr}(\Sigma) + \sum_{i \neq j} p_\lambda(|\Sigma_{ij}|), \quad (1)$$

where $\|\cdot\|_F$ is the Frobenius norm, $\text{tr}(\cdot)$ is the trace with tuning parameter $\tau > 0$, and $p_\lambda : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ is a nonconvex penalty function with tuning parameter $\lambda > 0$. \mathcal{E} is the spectral constraint, defined as

$$\mathcal{E} = \{\Sigma \mid \beta \mathbf{I} \preceq \Sigma \preceq \alpha \mathbf{I}\}.$$

where $\alpha \geq 0$ and $\beta \geq 0$. The spectral constraint is used to further promote the flexibility of the positive semi-definite constraint. It enables the estimated covariance matrix to achieve positive semi-definiteness, and avoid the indefiniteness which may cause trouble in downstream analysis. In addition, for certain downstream analysis of covariance matrices estimation, such as discriminant analysis, we need a good condition number for the estimated covariance matrix. Sometimes, we may want to incorporate the prior knowledge about the smallest and largest eigenvalues into the inference [26], [27].

In [13], the trace norm is proposed to encourage the low-rankness, since the eigenvalues of the positive semi-definite matrix are all greater than or equals to 0. To eliminate the estimation bias brought by ℓ_1 penalty, we propose using a nonconvex penalty to promote sparseness. We consider a class of nonconvex penalty functions p_λ satisfying the following assumptions.

Assumption 1. The function $p_\lambda(t)$ defined on $[0, +\infty)$ satisfies:

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

Prototypical examples of the penalty function p_λ in Assumption 1 include smoothly clipped absolute deviation (SCAD) [34] and minimax concave penalty (MCP) [35].

The problem (1) is generic. When $\tau = 0$, we can deal with the sparse covariance matrix estimation problem [24], [36]. When $\lambda = 0$, our estimator will reduce to low-rank covariance matrix estimation [13]. Compared to [36], which proposed a nonconvex estimation for sparse covariance matrices, our model is more flexible since it can

Algorithm 1: The Multistage Convex Relaxation Algorithm.

Input: $\mathbf{S}, \tau, \lambda$;
1 Initialize $\widehat{\Sigma}^{(0)}$;
2 for $k = 1, 2, \dots, K$ **do**
3 $\Lambda_{ij}^{(k-1)} = p'_\lambda(|\widehat{\Sigma}_{ij}^{(k-1)}|)$;
4 obtain $\widehat{\Sigma}^{(k)}$ by solving (3);
5 $k = k + 1$;
6 end
Output: $\widehat{\Sigma}^{(K)}$.

also solve the sparse covariance matrix estimation problem and the low-rank covariance matrix estimation problem. Also, as mentioned in [24], the influence of the perturbation from the log-determinant barrier becomes difficult to control.

III. OPTIMIZATION ALGORITHM

A. The MM Framework

Consider the minimization of a continuous function $F(\Theta)$. Initialized as $\Theta^{(0)}$, the MM algorithm [37] generates a sequence of feasible points $\{\Theta^{(k)}\}_{k \geq 1}$ by the following induction. At point $\Theta^{(k-1)}$, in the majorization step, we design a surrogate function $\bar{F}(\Theta \mid \Theta^{(k-1)})$ that locally approximates the objective function $F(\Theta)$, satisfying

$$\begin{cases} \bar{F}(\Theta \mid \Theta^{(k-1)}) \geq F(\Theta), \\ \bar{F}(\Theta^{(k-1)} \mid \Theta^{(k-1)}) = F(\Theta^{(k-1)}). \end{cases}$$

Then, in the minimization step, we update $\Theta^{(k)}$ as

$$\Theta^{(k)} \in \arg \min_{\Theta} \left\{ \bar{F}(\Theta \mid \Theta^{(k-1)}) \right\}.$$

B. The MM based Multistage Convex Relaxation Algorithm

The MM framework is employed to solve (1). In each iteration of MM, a weighted ℓ_1 surrogate function is constructed for $\sum_{i \neq j} p_\lambda(|\Sigma_{ij}|)$. This leads to a multistage procedure that solves a sequence of convex relaxation subproblems. Specifically, starting with an initial estimate $\widehat{\Sigma}^{(0)}$, we consider a sequence of convex optimization problems:

$$\underset{\Sigma \in \mathcal{E}}{\text{minimize}} \quad f(\Sigma) + \sum_{i \neq j} p'_\lambda(|\widehat{\Sigma}_{ij}^{(k-1)}|) |\Sigma_{ij}|, \quad 1 \leq k \leq K, \quad (2)$$

where $f(\Sigma) = \frac{1}{2} \|\Sigma - \mathbf{S} + \tau \mathbf{I}\|_F^2$, and $\widehat{\Sigma}^{(k)}$ denotes the optimal solution to the k -th subproblem. The solution to the k -th subproblem $\widehat{\Sigma}^{(k)}$ is given by

$$\underset{\Sigma \in \mathcal{E}}{\text{minimize}} \quad f(\Sigma) + \left\| \Lambda^{(k-1)} \odot \Sigma \right\|_{1, \text{off}}, \quad (3)$$

where $\Lambda_{ij}^{(k-1)} = p'_\lambda(|\widehat{\Sigma}_{ij}^{(k-1)}|)$ and $\|\cdot\|_{1, \text{off}}$ is the off-diagonal ℓ_1 -norm. The MM based multistage convex relaxation algorithm is summarized in Algorithm 1.

C. Solving the Subproblem (3)

We use the ADMM algorithm to solve each of the subproblems (3). Firstly, we introduce an auxiliary variable $\Psi \in \mathbb{R}^{d \times d}$, so that we can equivalently convert (3) into

$$\begin{aligned} & \underset{\Sigma}{\text{minimize}} \quad f(\Sigma) + \left\| \Lambda^{(k-1)} \odot \Psi \right\|_{1, \text{off}} \\ & \text{subject to} \quad \beta \mathbf{I} \preceq \Sigma \preceq \alpha \mathbf{I} \\ & \quad \quad \quad \Psi = \Sigma. \end{aligned}$$

Algorithm 2: ADMM for Solving (3)

Input: \mathbf{S} , $\mathbf{\Lambda}^{(k-1)}$, $\widehat{\mathbf{\Sigma}}^{(k-1)}$, T , τ , ρ , ϵ ;
1 Initialize $\mathbf{\Psi}_0^{(k)} = \widehat{\mathbf{\Sigma}}^{(k-1)}$, $\mathbf{\Gamma}_0^{(k)} = \mathbf{0}$, $t = 0$;
2 repeat
3 $t = t + 1$;
4 $\mathbf{\Pi}_t^{(k)} = \frac{\rho \mathbf{\Psi}_{t-1}^{(k)} - \mathbf{\Gamma}_{t-1}^{(k)} + \mathbf{S} - \tau \mathbf{I}}{\rho + 1}$;
5 $\mathbf{\Sigma}_t^{(k)} = \sum_{i=1}^d \min(\max(\lambda_i, \alpha), \beta) \mathbf{v}_i \mathbf{v}_i^\top$, where λ_i is the i -th eigenvalue of $\mathbf{\Pi}_t^{(k)}$, and \mathbf{v}_i is the corresponding eigenvector;
6 $\mathbf{\Psi}_t^{(k)} = \mathcal{T}_{\mathbf{\Lambda}^{(k-1)}/\rho}(\mathbf{\Sigma}_t^{(k)} + \frac{1}{\rho} \mathbf{\Gamma}_{t-1}^{(k)})$, where $\mathcal{T}_{\mathbf{\Lambda}}$ is the element-wise soft-thresholding operator;
7 $\mathbf{\Gamma}_t^{(k)} = \mathbf{\Gamma}_{t-1}^{(k)} + \rho(\mathbf{\Sigma}_t^{(k)} - \mathbf{\Psi}_t^{(k)})$;
8 until $\|\mathbf{\Sigma}_t^{(k)} - \mathbf{\Psi}_t^{(k)}\|_F < \epsilon$ or $t > T$;
Output: $\widehat{\mathbf{\Sigma}}^{(k)} = \mathbf{\Sigma}_t^{(k)}$.

Then we define the Lagrange multiplier $\mathbf{\Gamma} \in \mathbb{R}^{d \times d}$ associated with the equality constraint $\mathbf{\Psi} = \mathbf{\Sigma}$ to obtain the augmented Lagrangian function as

$$\mathcal{L}_\rho(\mathbf{\Sigma}, \mathbf{\Psi}, \mathbf{\Gamma}) := \frac{1}{2} \|\mathbf{\Sigma} - \mathbf{S} + \tau \mathbf{\Gamma}\|_F^2 + \|\mathbf{\Lambda}^{(k-1)} \odot \mathbf{\Psi}\|_{1,\text{off}} + \langle \mathbf{\Gamma}, \mathbf{\Sigma} - \mathbf{\Psi} \rangle + \frac{\rho}{2} \|\mathbf{\Sigma} - \mathbf{\Psi}\|_F^2,$$

where $\rho > 0$ is the penalty parameter. ADMM procedure consists of iteratively minimizing \mathcal{L}_ρ with respect to $\mathbf{\Sigma}$, minimizing \mathcal{L}_ρ with respect to $\mathbf{\Psi}$, and updating the multiplier $\mathbf{\Gamma}$. The ADMM algorithm is summarized in Algorithm 2.

IV. MAIN THEORY

In this section, we first introduce some necessary assumptions for the statistical analysis, and then establish the statistical convergence rate of the proposed covariance estimator. Due to space limitation, all the proof to the results of the paper is given in [38].

A. Assumptions

We denote the true covariance matrix by $\mathbf{\Sigma}^*$. Let $\mathcal{S}^* = \{(i, j) \mid \Sigma_{ij}^* \neq 0\}$ be the support set of $\mathbf{\Sigma}^*$ and s^* be its cardinality, i.e., $s^* = |\mathcal{S}^*|$. In the following, we impose some mild conditions on the true covariance matrix $\mathbf{\Sigma}^*$.

Assumption 2. For the true covariance matrix $\mathbf{\Sigma}^*$, there exists $\kappa \geq 1$ such that

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

This assumption requires the eigenvalues of true covariance matrix $\mathbf{\Sigma}^*$ to be finite and bounded below from a positive number, which is a standard assumption for sparse matrix estimation [33].

Assumption 3. The true covariance matrix $\mathbf{\Sigma}^*$ satisfies

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

where α and c are constants defined in Assumption 1.

This assumption is referred to as the minimum signal strength condition, which is commonly employed in the analysis of nonconvex penalized regression problems [34], [35].

B. Statistical Analysis

Theorem 4. Suppose that Assumption 2 holds. If $\lambda \geq 2 \|\nabla f(\mathbf{\Sigma}^*)\|_{\max}$, then the optimal solution $\widehat{\mathbf{\Sigma}}^{(k)}$ satisfies the following δ -contraction property:

$$\|\widehat{\mathbf{\Sigma}}^{(k)} - \mathbf{\Sigma}^*\|_F \leq \underbrace{\|(\nabla f(\mathbf{\Sigma}^*))_{\mathcal{S}^*}\|_F}_{\text{oracle rate}} + \delta \underbrace{\|\widehat{\mathbf{\Sigma}}^{(k-1)} - \mathbf{\Sigma}^*\|_F}_{\text{contraction}}$$

for $1 \leq k \leq K$, where $\delta \in (0, 1)$ is the contraction parameter.

Theorem 4 establishes an upper bound on the estimation error between the optimal solution $\widehat{\mathbf{\Sigma}}^{(k)}$ and the true parameter $\mathbf{\Sigma}^*$, comprising two distinct components, the oracle rate and a contraction term. Next, we will derive an explicit statistical rate of convergence for the estimator under the assumption of sub-Gaussian design.

Corollary 5. Let \mathbf{x} be a sub-Gaussian random vector with zero mean and covariance $\mathbf{\Sigma}^*$ 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 and 2 hold. If

$$\lambda \asymp \sqrt{\frac{\log d}{n}}, \quad \tau \lesssim \sqrt{\frac{1}{n}},$$

then the optimal solution $\widehat{\mathbf{\Sigma}}^{(1)}$ satisfies

$$\|\widehat{\mathbf{\Sigma}}^{(1)} - \mathbf{\Sigma}^*\|_F \lesssim \sqrt{\frac{s^* \log d}{n}}$$

with high probability.

Corollary 5 follows from Theorem 4 by setting $k = 1$. Notably, the convergence rate of $\sqrt{\frac{s^* \log d}{n}}$ established in Corollary 5 coincides with the rate obtained by [23].

Corollary 6. Let \mathbf{x} be a sub-Gaussian random vector with zero mean and covariance $\mathbf{\Sigma}^*$ and $\{\mathbf{x}_i\}_{i=1}^n$ be a collection of i.i.d. samples from \mathbf{x} . Suppose that Assumptions 1, 2, and 3 hold. If

$$\lambda \asymp \sqrt{\frac{\log d}{n}}, \quad \tau \lesssim \sqrt{\frac{1}{n}},$$

and $K \gtrsim \log(\lambda\sqrt{n}) \gtrsim \log \log d$, then the optimal solution $\widehat{\mathbf{\Sigma}}^{(K)}$ satisfies

$$\|\widehat{\mathbf{\Sigma}}^{(K)} - \mathbf{\Sigma}^*\|_F = O_p\left(\sqrt{\frac{s^*}{n}}\right).$$

Corollary 6 follows directly from Theorem 4, which reveals that under mild assumptions, it suffices to solve no more than approximately $\log \log d$ convex problems to attain the oracle rate of convergence $\sqrt{\frac{s^*}{n}}$. Notably, this implies that our proposed estimator exhibits a superior statistical rate of convergence in the Frobenius norm compared to the estimator of [23]. We obtain the same oracle statistical rate as [36]. However, compared to [36], we consider a more general problem formulation which secures the positive semi-definiteness of the estimator by enforcing a spectral constraint, and can be naturally extended to more flexible settings. When $\tau = 0$, we can deal with the sparse covariance matrix estimation problem in [36].

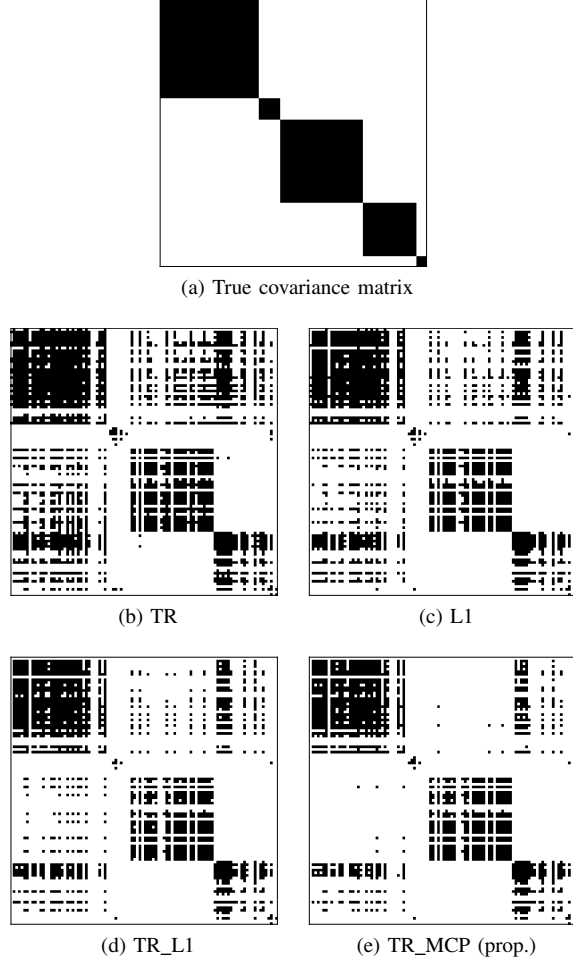
V. NUMERICAL EXPERIMENTS

In this section, we conduct numerical experiments on synthetic datasets to verify the theoretical properties of the proposed estimator.

We generate data in the following way. We draw n vectors $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}^*)$ for a block diagonal covariance matrix $\mathbf{\Sigma}^* \in \mathbb{R}^{d \times d}$. We use r blocks of random sizes and of the form $\mathbf{v}\mathbf{v}^\top$ where the entries

TABLE I: QUANTITATIVE COMPARISON AMONG DIFFERENT METHODS

	SCM	TR	L1	TR_L1	TR_MCP (prop.)
$n = 30, d = 50, r = 5, s = 724$					
$\ \Sigma^{(K)} - \Sigma^*\ _F$	1.0607±0.1742	0.9375±0.1533	0.9337±0.1581	0.9177±0.1360	0.7736±0.1714
FPR	-	0.1935±0.0473	0.1370±0.0375	0.0894±0.0279	0.0531±0.0433
TPR	-	0.9744±0.0241	0.9872±0.0282	0.9877±0.0156	0.9891±0.0196
Time	-	0.0437	0.0391	0.0297	2.7594
$n = 50, d = 80, r = 8, s = 1672$					
$\ \Sigma^{(K)} - \Sigma^*\ _F$	1.5100±0.2521	1.2938±0.2454	1.2916±0.2841	1.2528±0.2450	1.1360±0.2487
FPR	-	0.2034±0.0501	0.1301±0.0357	0.0823±0.0499	0.0528±0.0378
TPR	-	0.9623±0.0133	0.9754±0.0285	0.9848±0.0147	0.9884±0.0159
Time	-	0.0563	0.0531	0.3359	4.3219
$n = 70, d = 100, r = 10, s = 2430$					
$\ \Sigma^{(K)} - \Sigma^*\ _F$	1.7865±0.1463	1.4699±0.1553	1.4561±0.1211	1.4079±0.2078	1.3380±0.2053
FPR	-	0.1938±0.0446	0.1238±0.0586	0.0938±0.0437	0.0473±0.0471
TPR	-	0.9618±0.0143	0.9730±0.0257	0.9801±0.0258	0.9837±0.0260
Time	-	0.0891	0.0594	0.4906	6.3344

Fig. 1: Support of Σ^* , and of the estimates given by TR, L1, TR_L1, and TR_MCP.

of \mathbf{v} are drawn i.i.d. from the uniform distribution on $[-1, 1]$. We compute the SCM as $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top$. We compare the estimation performance of different methods based on the estimation errors in Frobenius norm for the covariance matrix (i.e., $\|\Sigma^{(K)} - \Sigma^*\|_F$), the false positive rate (FPR), the true positive rate (TPR). We conduct experiments in different settings of n, d, r , and s . The regularization parameters for the trace norm and the ℓ_1 norm in other estimators are selected via 4-fold cross-validation.

In Table I, we apply our method TR_MCP with trace norm plus MCP regularization, as well as SCM, low-rank covariance estimation with trace norm regularization (TR), sparse covariance estimation with ℓ_1 norm regularization (L1), and joint low-rank and sparse covariance estimation with trace norm plus ℓ_1 norm regularization (TR_L1) to the SCM. We generate 10 datasets for each setting of n, d, r , and s , and compute the mean and standard error of the above evaluation metrics for each dataset. Notably, our proposed method TR_MCP outputs uniformly better results in terms of estimation errors, FPR, and TPR compared to other methods. We also report the average CPU time (in seconds) in Table I. And TR_L1 is slightly faster than other methods. We do not report the FPR, TPR, and computational time of SCM, since it is a direct method.

We also show in Fig. 1 the support of the true covariance matrix and the supports recovered by TR, L1, TR_L1, and TR_MCP. The output covariance matrix is thresholded in absolute value. The support recovery demonstrates TR_MCP discovers the best block diagonal structure than others.

VI. CONCLUSIONS

In this paper, we have proposed a novel approach for large covariance matrix estimation with groups of highly correlated variables using the nonconvex penalty and presented both its theoretical and empirical results. Our proposed estimators are proven to have better statistical rates of convergence compared to existing approaches. To the best of our knowledge, this is the first work to obtain the oracle statistical rate of convergence for the large covariance matrix estimation problem for groups of highly correlated variables.

REFERENCES

- [1] P. M. Bentler and P. Dudgeon, "Covariance structure analysis: Statistical practice, theory, and directions," *Annu. Rev. Psychol.*, vol. 47, no. 1, pp. 563–592, 1996.
- [2] J. Schäfer and K. Strimmer, "A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics," *Stat. Appl. Genet. Mol. Biol.*, vol. 4, no. 1, 2005.
- [3] Z. Zhao and D. P. Palomar, "Mean-reverting portfolio with budget constraint," *IEEE Trans. Signal Process.*, vol. 66, no. 9, pp. 2342–2357, 2018.
- [4] Z. Zhao, R. Zhou, and D. P. Palomar, "Optimal mean-reverting portfolio with leverage constraint for statistical arbitrage in finance," *IEEE Trans. Signal Process.*, vol. 67, no. 7, pp. 1681–1695, 2019.
- [5] J. Capon, "High-resolution frequency-wavenumber spectrum analysis," *Proc. IEEE*, vol. 57, no. 8, pp. 1408–1418, 1969.
- [6] S. U. Pillai, *Array Signal Processing*. Springer, 2012.
- [7] C. M. Bishop and N. M. Nasrabadi, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [8] I. M. Johnstone, "On the distribution of the largest eigenvalue in principal components analysis," *Ann. Statist.*, vol. 29, no. 2, pp. 295–327, 2001.
- [9] P. J. Bickel and E. Levina, "Covariance regularization by thresholding," *Ann. Stat.*, vol. 36, no. 6, pp. 2577 – 2604, 2008.
- [10] A. J. Rothman, E. Levina, and J. Zhu, "Generalized thresholding of large covariance matrices," *J. Am. Stat. Assoc.*, vol. 104, no. 485, pp. 177–186, 2009.
- [11] J. Fan, Y. Liao, and H. Liu, "An overview of the estimation of large covariance and precision matrices," *Econometrics J.*, vol. 19, no. 1, pp. C1–C32, 2016.
- [12] S. Negahban and M. J. Wainwright, "Estimation of (near) low-rank matrices with noise and high-dimensional scaling," *Ann. Statist.*, vol. 39, no. 2, 2011.
- [13] K. Lounici, "High-dimensional covariance matrix estimation with missing observations," *Bernoulli*, pp. 1029–1058, 2014.
- [14] A. P. Shikhaliyev, L. C. Potter, and Y. Chi, "Low-rank structured covariance matrix estimation," *IEEE Signal Process. Lett.*, vol. 26, no. 5, pp. 700–704, 2019.
- [15] A. Agarwal, S. Negahban, and M. J. Wainwright, "Noisy matrix decomposition via convex relaxation: Optimal rates in high dimensions," *Ann. Statist.*, vol. 40, no. 2, 2012.
- [16] S. Zou and Z. Zhao, "Large covariance matrix estimation based on factor models via nonconvex optimization," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.*, pp. 9656–9660, IEEE, 2024.
- [17] T. T. Cai, C.-H. Zhang, and H. H. Zhou, "Optimal rates of convergence for covariance matrix estimation," *Ann. Statist.*, 2010.
- [18] T. T. Cai, Z. Ren, and H. H. Zhou, "Optimal rates of convergence for estimating Toeplitz covariance matrices," *Probab. Theory Relat. Fields.*, vol. 156, no. 1-2, pp. 101–143, 2013.
- [19] G. Camba-Mendez and G. Kapetanios, "Testing the rank of the Hankel covariance matrix: a statistical approach," *IEEE Trans. Autom. Control.*, vol. 46, no. 2, pp. 331–336, 2001.
- [20] P. J. Bickel and E. Levina, "Regularized estimation of large covariance matrices," *Ann. Statist.*, 2008.
- [21] T. T. Cai and H. H. Zhou, "Minimax estimation of large covariance matrices under ℓ_1 -norm," *Stat. Sin.*, pp. 1319–1349, 2012.
- [22] R. Furrer and T. Bengtsson, "Estimation of high-dimensional prior and posterior covariance matrices in Kalman filter variants," *J. Multivariate Anal.*, vol. 98, no. 2, pp. 227–255, 2007.
- [23] A. J. Rothman, "Positive definite estimators of large covariance matrices," *Biometrika*, vol. 99, no. 3, pp. 733–740, 2012.
- [24] L. Xue, S. Ma, and H. Zou, "Positive-definite ℓ_1 -penalized estimation of large covariance matrices," *J. Am. Stat. Assoc.*, vol. 107, no. 500, pp. 1480–1491, 2012.
- [25] H. Liu, L. Wang, and T. Zhao, "Sparse covariance matrix estimation with eigenvalue constraints," *J. Comput. Graph. Stat.*, vol. 23, no. 2, pp. 439–459, 2014.
- [26] J.-H. Won, J. Lim, S.-J. Kim, and B. Rajaratnam, "Condition-number-regularized covariance estimation," *J. R. Stat. Soc. B*, vol. 75, no. 3, pp. 427–450, 2013.
- [27] Y.-H. Xiao, P.-L. Li, and S. Lu, "Sparse estimation of high-dimensional inverse covariance matrices with explicit eigenvalue constraints," *J. Oper. Res. Soc. China*, vol. 9, pp. 543–568, 2021.
- [28] U. Alon, N. Barkai, D. A. Notterman, K. Gish, S. Ybarra, D. Mack, and A. J. Levine, "Broad patterns of gene expression revealed by clustering analysis of tumor and normal colon tissues probed by oligonucleotide arrays," *Proc. Nat. Acad. Sci.*, vol. 96, no. 12, pp. 6745–6750, 1999.
- [29] L. Žiganić, S. Begušić, and Z. Kostanjčar, "Block-diagonal idiosyncratic covariance estimation in high-dimensional factor models for financial time series," *J. Comput. Sci-neth.*, p. 102348, 2024.
- [30] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A k -means clustering algorithm," *J. R. Stat. Soc.*, vol. 28, no. 1, pp. 100–108, 1979.
- [31] R. Sibson, "SLINK: An optimally efficient algorithm for the single-link cluster method," *Comput. J.*, vol. 16, no. 1, pp. 30–34, 1973.
- [32] E. Richard, P.-A. Savalle, and N. Vayatis, "Estimation of simultaneously sparse and low rank matrices," in *Proc. Int. Conf. Mach. Learn.*, pp. 51–58, 2012.
- [33] S.-L. Zhou, N.-H. Xiu, Z.-Y. Luo, and L.-C. Kong, "Sparse and low-rank covariance matrix estimation," *J. Oper. Res. Soc. China*, vol. 3, no. 2, pp. 231–250, 2015.
- [34] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *J. Am. Stat. Assoc.*, vol. 96, no. 456, pp. 1348–1360, 2001.
- [35] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *Ann. Statist.*, vol. 38, no. 1, pp. 894–942, 2010.
- [36] Q. Wei and Z. Zhao, "Large covariance matrix estimation with oracle statistical rate via majorization-minimization," *IEEE Trans. Signal Process.*, 2023.
- [37] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Trans. Signal Process.*, vol. 65, no. 3, pp. 794–816, 2016.
- [38] <https://www.ncvxo.com/pubs/ZouZhaoBlockDiagonalCovariance.pdf>

APPENDIX

In this appendix, we first provide some auxiliary lemmas, and then provide the proofs of all the statistical theoretical results in Section IV.

A. Auxiliary Lemmas

Lemma 7. *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 constants c_1 , c_2 , and t_0 such that for all t with $0 < t < t_0$, the SCM \mathbf{S} satisfies the following tail bound*

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

Lemma 8. *Under the same conditions in Lemma 7, if taking $\lambda = \sqrt{\frac{3 \log d}{c_2 n}} \asymp \sqrt{\frac{\log d}{n}} < t_0$, then the following result holds*

$$\mathbb{P}(\|\Sigma^* - \mathbf{S}\|_{\max} \leq \lambda) \geq 1 - \frac{c_1}{d}.$$

Proof: Applying Lemma 7 and union bound, for any λ such that $0 < \lambda < t_0$, we obtain

$$\begin{aligned} \mathbb{P}(\|\Sigma^* - \mathbf{S}\|_{\max} > \lambda) &\leq c_1 d^2 \exp(-c_2 n \lambda^2) \\ &= c_1 \exp(-c_2 n \lambda^2 + 2 \log d). \end{aligned}$$

For n sufficiently large such that $n > \frac{3 \log d}{c_2 t_0^2}$, by taking $\lambda = \sqrt{\frac{3 \log d}{c_2 n}} \asymp \sqrt{\frac{\log d}{n}} < t_0$, we obtain

$$\begin{aligned} \mathbb{P}(\|\Sigma^* - \mathbf{S}\|_{\max} \leq \lambda) &\geq 1 - c_1 \exp(-c_2 n \lambda^2 + 2 \log d) \\ &= 1 - \frac{c_1}{d}. \end{aligned}$$

■

Lemma 9. *Under the same conditions in Lemma 7, the following result holds*

$$\|(\Sigma^* - \mathbf{S})_{S^*}\|_F = O_p\left(\sqrt{\frac{s^*}{n}}\right).$$

Proof: Applying Lemma 7 and union bound, for any M such that $0 < M\sqrt{\frac{1}{n}} < t_0$, we obtain

$$\begin{aligned} & \mathbb{P} \left(\|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_{\max} > M\sqrt{\frac{1}{n}} \right) \\ & \leq c_1 s^* \exp(-c_2 M^2) \\ & = c_1 \exp(-c_2 M^2 + \log s^*). \end{aligned}$$

By taking M such that $\sqrt{\frac{2 \log s^*}{c_2}} < M < t_0 \sqrt{n}$ and $M \rightarrow \infty$ in the above inequality obtains

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

The proof is completed by applying $\|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_F \leq \sqrt{s^*} \|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_{\max}$. ■

B. Technical Lemmas

Each subproblem in (2) corresponds to a weighted ℓ_1 penalized covariance estimation problem, which generally can be written into the following form:

$$\underset{\beta \mathbf{I} \preceq \mathbf{\Sigma} \preceq \alpha \mathbf{I}}{\text{minimize}} \quad f(\mathbf{\Sigma}) + \|\mathbf{\Lambda} \odot \mathbf{\Sigma}\|_{1, \text{off}}, \quad (4)$$

where $\mathbf{\Lambda}$ is a $d \times d$ matrix of regularization parameters with $\Lambda_{ij} \in [0, \lambda]$.

Lemma 10. Suppose that Assumption 2 holds. Consider the general problem in (4). Assume that there exists a set \mathcal{E} such that

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

If $\lambda \geq 2 \|\nabla f(\mathbf{\Sigma}^*)\|_{\max}$, then the optimal solution $\hat{\mathbf{\Sigma}}$ satisfies

$$\begin{aligned} \|\hat{\mathbf{\Sigma}} - \mathbf{\Sigma}^*\|_F & \leq \|(\nabla f(\mathbf{\Sigma}^*))_{\mathcal{E}}\|_F + \|\mathbf{\Lambda}_{S^*}\|_F \\ & \leq \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*}. \end{aligned}$$

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

$$\langle \nabla f(\mathbf{\Sigma}) - \nabla f(\mathbf{\Sigma}^*), \mathbf{\Delta} \rangle = \text{vec}^\top(\mathbf{\Delta}) \nabla^2 f(\mathbf{\Sigma}^* + \rho \mathbf{\Delta}) \text{vec}(\mathbf{\Delta}),$$

where $\mathbf{\Delta} = \mathbf{\Sigma} - \mathbf{\Sigma}^*$. One has

$$\begin{aligned} & \text{vec}^\top(\mathbf{\Delta}) \nabla^2 f(\mathbf{\Sigma}^* + \rho \mathbf{\Delta}) \text{vec}(\mathbf{\Delta}) \\ & \geq \lambda_{\min} (\nabla^2 f(\mathbf{\Sigma}^* + \rho \mathbf{\Delta})) \|\mathbf{\Delta}\|_F^2 = \|\mathbf{\Delta}\|_F^2. \end{aligned}$$

Hence, we obtain

$$\|\mathbf{\Delta}\|_F^2 \leq \langle \nabla f(\mathbf{\Sigma}) - \nabla f(\mathbf{\Sigma}^*), \mathbf{\Delta} \rangle. \quad (5)$$

Define the Lagrangian function of (4) as

$$\begin{aligned} \mathcal{L}(\mathbf{\Sigma}, \mathbf{Z}_1, \mathbf{Z}_2) & := f(\mathbf{\Sigma}) + \|\mathbf{\Lambda} \odot \mathbf{\Sigma}\|_{1, \text{off}} \\ & \quad - \langle \mathbf{Z}_1, \mathbf{\Sigma} - \alpha \mathbf{I} \rangle + \langle \mathbf{Z}_2, \mathbf{\Sigma} - \beta \mathbf{I} \rangle \end{aligned}$$

where \mathbf{Z}_1 and \mathbf{Z}_2 are $d \times d$ matrix of dual variables. From convex optimization theory, we know that any optimal solution $\hat{\mathbf{\Sigma}}$ to (4) satisfies the following KKT condition:

$$\begin{aligned} \nabla f(\hat{\mathbf{\Sigma}}) + \mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}} - \hat{\mathbf{Z}}_1 + \hat{\mathbf{Z}}_2 & = \mathbf{0}, \quad \text{with } \hat{\mathbf{\Sigma}} \in \partial \|\hat{\mathbf{\Sigma}}\|_{1, \text{off}}; \\ \langle \hat{\mathbf{Z}}_1, \hat{\mathbf{\Sigma}} - \alpha \mathbf{I} \rangle & = 0; \\ \langle \hat{\mathbf{Z}}_2, \hat{\mathbf{\Sigma}} - \beta \mathbf{I} \rangle & = 0; \\ \beta \mathbf{I} & \succeq \hat{\mathbf{\Sigma}} \succeq \alpha \mathbf{I}; \\ \hat{\mathbf{Z}}_1 & \succeq \mathbf{0}; \\ \hat{\mathbf{Z}}_2 & \succeq \mathbf{0}, \end{aligned} \quad (6)$$

where $\nabla f(\mathbf{\Sigma}) = \mathbf{\Sigma} - \mathbf{S} + \tau \mathbf{I}$.

Let $\hat{\mathbf{\Delta}} = \hat{\mathbf{\Sigma}} - \mathbf{\Sigma}^*$. Applying the inequality (5) with $\mathbf{\Sigma} = \hat{\mathbf{\Sigma}}$ and $\mathbf{\Delta} = \hat{\mathbf{\Delta}}$ yields

$$\begin{aligned} \|\hat{\mathbf{\Delta}}\|_F^2 & \leq \langle \nabla f(\hat{\mathbf{\Sigma}}) - \nabla f(\mathbf{\Sigma}^*), \hat{\mathbf{\Delta}} \rangle \\ & = \underbrace{\langle \nabla f(\hat{\mathbf{\Sigma}}) + \mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}} - \hat{\mathbf{Z}}_1 + \hat{\mathbf{Z}}_2, \hat{\mathbf{\Delta}} \rangle}_{\text{I}} - \underbrace{\langle \nabla f(\mathbf{\Sigma}^*), \hat{\mathbf{\Delta}} \rangle}_{\text{II}} \\ & \quad - \underbrace{\langle \mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}}, \hat{\mathbf{\Delta}} \rangle}_{\text{III}} + \underbrace{\langle \hat{\mathbf{Z}}_1, \hat{\mathbf{\Delta}} \rangle}_{\text{IV}} - \underbrace{\langle \hat{\mathbf{Z}}_2, \hat{\mathbf{\Delta}} \rangle}_{\text{V}}, \end{aligned} \quad (7)$$

where $\hat{\mathbf{\Sigma}} \in \partial \|\hat{\mathbf{\Sigma}}\|_{1, \text{off}}$. It remains to bound terms I, II, III, IV, and V, respectively.

For term I, using the KKT condition in (6), we obtain

$$\text{I} = \langle \nabla f(\hat{\mathbf{\Sigma}}) + \mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}} - \hat{\mathbf{Z}}_1 + \hat{\mathbf{Z}}_2, \hat{\mathbf{\Delta}} \rangle = 0$$

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

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

For term III, separating the support of $\mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}}$ and $\hat{\mathbf{\Delta}}$ to S^* and \bar{S}^* , and then using the matrix Hölder's inequality, we obtain

$$\begin{aligned} \text{III} & = \langle (\mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}})_{S^*}, \hat{\mathbf{\Delta}}_{S^*} \rangle + \langle (\mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}})_{\bar{S}^*}, \hat{\mathbf{\Delta}}_{\bar{S}^*} \rangle \\ & = \langle (\mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}})_{S^*}, \hat{\mathbf{\Delta}}_{S^*} \rangle + \langle \mathbf{\Lambda}_{\bar{S}^*}, |\hat{\mathbf{\Delta}}_{\bar{S}^*}| \rangle \\ & \geq -\|\mathbf{\Lambda}_{S^*}\|_F \|\hat{\mathbf{\Delta}}_{S^*}\|_F + \langle \mathbf{\Lambda}_{\bar{S}^*}, |\hat{\mathbf{\Delta}}_{\bar{S}^*}| \rangle \\ & \geq -\|\mathbf{\Lambda}_{S^*}\|_F \|\hat{\mathbf{\Delta}}\|_F + \|\mathbf{\Lambda}_{\bar{S}^*}\|_{\min} \|\hat{\mathbf{\Delta}}_{\bar{S}^*}\|_1, \end{aligned}$$

where the second equality is due to

$$\langle (\mathbf{\Lambda} \odot \hat{\mathbf{\Sigma}})_{\bar{S}^*}, \hat{\mathbf{\Delta}}_{\bar{S}^*} \rangle = \langle \mathbf{\Lambda}_{\bar{S}^*}, |\hat{\mathbf{\Sigma}}_{\bar{S}^*}| \rangle = \langle \mathbf{\Lambda}_{\bar{S}^*}, |\hat{\mathbf{\Delta}}_{\bar{S}^*}| \rangle,$$

and the second inequality is due to

$$\begin{aligned} \langle \mathbf{\Lambda}_{\bar{S}^*}, |\hat{\mathbf{\Delta}}_{\bar{S}^*}| \rangle & = \sum_{(i,j) \in \bar{\mathcal{E}}} \Lambda_{ij} |\hat{\Delta}_{ij}| \geq \|\mathbf{\Lambda}_{\bar{S}^*}\|_{\min} \sum_{(i,j) \in \bar{\mathcal{E}}} |\hat{\Delta}_{ij}| \\ & = \|\mathbf{\Lambda}_{\bar{S}^*}\|_{\min} \|\hat{\mathbf{\Delta}}_{\bar{S}^*}\|_1. \end{aligned}$$

For term IV, we obtain

$$\begin{aligned} \text{IV} & = \langle \hat{\mathbf{Z}}_1, \hat{\mathbf{\Sigma}} \rangle - \langle \hat{\mathbf{Z}}_1, \mathbf{\Sigma}^* \rangle \\ & = \langle \hat{\mathbf{Z}}_1, \hat{\mathbf{\Sigma}} - \alpha \mathbf{I} \rangle + \langle \hat{\mathbf{Z}}_1, \alpha \mathbf{I} \rangle - \langle \hat{\mathbf{Z}}_1, \mathbf{\Sigma}^* \rangle \\ & = \langle \hat{\mathbf{Z}}_1, \alpha \mathbf{I} - \mathbf{\Sigma}^* \rangle \leq 0. \end{aligned}$$

For term V, we obtain

$$\begin{aligned} V &= \langle \widehat{\mathbf{Z}}_2, \widehat{\Sigma} \rangle - \langle \widehat{\mathbf{Z}}_2, \Sigma^* \rangle \\ &= \langle \widehat{\mathbf{Z}}_2, \widehat{\Sigma} - \beta \mathbf{I} \rangle + \langle \widehat{\mathbf{Z}}_2, \beta \mathbf{I} \rangle - \langle \widehat{\mathbf{Z}}_2, \Sigma^* \rangle \\ &= \langle \widehat{\mathbf{Z}}_2, \beta \mathbf{I} - \Sigma^* \rangle \geq 0. \end{aligned}$$

Substituting the above results into (7) yields

$$\begin{aligned} \|\widehat{\Delta}\|_F^2 &\leq (\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \|\Lambda_{S^*}\|_F) \|\widehat{\Delta}\|_F \\ &\quad + (\|\nabla f(\Sigma^*)\|_{\max} - \|\Lambda_{\bar{\mathcal{E}}}\|_{\min}) \|\widehat{\Delta}\|_1 \\ &\leq (\|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \|\Lambda_{S^*}\|_F) \|\widehat{\Delta}\|_F, \end{aligned} \quad (8)$$

where the second inequality is due to $\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \geq \frac{\lambda}{2} \geq \|\nabla f(\Sigma^*)\|_{\max}$. Dividing by $\|\widehat{\Delta}\|_F$ on both sides of the inequality (8), we have

$$\begin{aligned} \|\widehat{\Delta}\|_F &\leq \|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_F + \|\Lambda_{S^*}\|_F \\ &\leq \|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_{\max} \sqrt{|\mathcal{E}|} + \|\Lambda_{S^*}\|_{\max} \sqrt{|\mathcal{S}^*|} \\ &\leq \|(\nabla f(\Sigma^*))_{\mathcal{E}}\|_{\max} \sqrt{2s^*} + \lambda \sqrt{s^*} \\ &\leq \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*}. \end{aligned}$$

Lemma 11. Suppose that Assumptions 2 hold. Consider the problem in (3). Define the set $\mathcal{E}^{(k)}$ by

$$\mathcal{E}^{(k)} = \mathcal{S}^* \cup \mathcal{S}^{(k)}, \text{ with } \mathcal{S}^{(k)} = \{(i, j) \mid \Lambda_{ij}^{(k-1)} < p'_\lambda(u)\},$$

where $u = c\lambda$ and $c = \frac{2+\sqrt{2}}{2}$ is the same as that given in Assumption 2. If $\lambda \geq 2\|\nabla f(\Sigma^*)\|_{\max}$, then for $k \geq 1$, we have $|\mathcal{E}^{(k)}| \leq 2s^*$, $\|\Lambda_{\mathcal{E}^{(k)}}^{(k-1)}\|_{\min} \geq \frac{\lambda}{2}$, and

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

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

$$\begin{aligned} \sqrt{|\mathcal{S}^{(k)} \setminus \mathcal{S}^*|} &\leq \sqrt{\sum_{(i,j) \in \mathcal{S}^{(k)} \setminus \mathcal{S}^*} (u^{-1} \widehat{\Sigma}_{ij}^{(k-1)})^2} \\ &= u^{-1} \|\widehat{\Sigma}_{\mathcal{S}^{(k)} \setminus \mathcal{S}^*}^{(k-1)}\|_F \\ &= u^{-1} \left\| \left(\widehat{\Sigma}^{(k-1)} - \Sigma^* \right)_{\mathcal{S}^{(k)} \setminus \mathcal{S}^*} \right\|_F \\ &\leq u^{-1} \|\widehat{\Sigma}^{(k-1)} - \Sigma^*\|_F. \end{aligned} \quad (9)$$

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

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

One also has $|\mathcal{E}^{(k-1)}| \leq 2s^*$ and $\mathcal{S}^* \subseteq \mathcal{E}^{(k-1)}$. Applying Lemma 10 with $\widehat{\Sigma} = \widehat{\Sigma}^{(k-1)}$, $\mathcal{E} = \mathcal{E}^{(k-1)}$, and $\Lambda_{S^*} = \Lambda_{\mathcal{S}^*}^{(k-2)}$ yields

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

Substituting the above result into the inequality (9) yields

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

Thus, we have

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

completing the induction.

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

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

Applying Lemma 10 with $\widehat{\Sigma} = \widehat{\Sigma}^{(k)}$, $\mathcal{E} = \mathcal{E}^{(k)}$, and $\Lambda_{S^*} = \Lambda_{\mathcal{S}^*}^{(k-1)}$, the optimal solution $\widehat{\Sigma}^{(k)}$ to (3) satisfies

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

C. Proof of Theorem 4

Proof: By Lemma 11, we have

$$\|\widehat{\Sigma}^{(k)} - \Sigma^*\|_F \leq \underbrace{\|(\nabla f(\Sigma^*))_{\mathcal{E}^{(k)}}\|_F}_I + \underbrace{\|\Lambda_{S^*}^{(k-1)}\|_F}_{II} \quad (10)$$

Next, we bound the terms I and II, respectively.

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

$$\begin{aligned} I &\leq \|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F + \|\nabla f(\Sigma^*)\|_{\max} \sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} \\ &\leq \|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F + \frac{1}{2} \lambda u^{-1} \|\widehat{\Sigma}^{(k-1)} - \Sigma^*\|_F, \end{aligned}$$

where the second inequality is due to

$$\sqrt{|\mathcal{E}^{(k)} \setminus \mathcal{S}^*|} = \sqrt{|\mathcal{S}^{(k)} \setminus \mathcal{S}^*|} \leq u^{-1} \|\widehat{\Sigma}^{(k-1)} - \Sigma^*\|_F,$$

which follows from the inequality (9).

By Assumptions 1 and 3, 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 V, we have

$$II \leq \lambda u^{-1} \|\widehat{\Sigma}_{S^*}^{(k-1)} - \Sigma_{S^*}^*\|_F \leq \lambda u^{-1} \|\widehat{\Sigma}^{(k-1)} - \Sigma^*\|_F.$$

Substituting the above results into (10) yields

$$\|\widehat{\Sigma}^{(k)} - \Sigma^*\|_F \leq \|(\nabla f(\Sigma^*))_{\mathcal{S}^*}\|_F + \delta \|\widehat{\Sigma}^{(k-1)} - \Sigma^*\|_F,$$

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

D. Proof of Corollary 5

Proof: Since $\nabla f(\mathbf{\Sigma}) = \mathbf{\Sigma} - \mathbf{S} + \tau \mathbf{I}$, one has

$$\|\nabla f(\mathbf{\Sigma}^*)\|_{\max} \leq \|\mathbf{\Sigma}^* - \mathbf{S}\|_{\max} + \tau.$$

If λ and τ satisfy

$$\lambda \asymp \sqrt{\frac{\log d}{n}}, \quad \tau \lesssim \sqrt{\frac{1}{n}},$$

then by Lemma 8, $\lambda \geq 2 \|\nabla f(\mathbf{\Sigma}^*)\|_{\max}$ holds with high probability (w.h.p).

Applying Lemma 11 with $k = 1$, we obtain

$$\|\widehat{\mathbf{\Sigma}}^{(1)} - \mathbf{\Sigma}^*\|_F \leq \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*}.$$

If $\lambda \asymp \sqrt{\frac{\log d}{n}}$, then $\|\widehat{\mathbf{\Sigma}}^{(1)} - \mathbf{\Sigma}^*\|_F \lesssim \sqrt{\frac{s^* \log d}{n}}$ w.h.p. \blacksquare

E. Proof of Corollary 6

Proof: One has

$$\|\nabla f(\mathbf{\Sigma}^*)\|_{\max} \leq \|\mathbf{\Sigma}^* - \mathbf{S}\|_{\max} + \tau.$$

If λ , τ , and ε satisfy

$$\lambda \asymp \sqrt{\frac{\log d}{n}}, \quad \tau \lesssim \sqrt{\frac{1}{n}},$$

then by Lemma 8, $\lambda \geq 2 \|\nabla f(\mathbf{\Sigma}^*)\|_{\max}$ holds w.h.p.

Applying Theorem 4, we obtain

$$\begin{aligned} & \|\widehat{\mathbf{\Sigma}}^{(k)} - \mathbf{\Sigma}^*\|_F \\ & \leq \|(\nabla f(\mathbf{\Sigma}^*))_{S^*}\|_F + \delta \|\widehat{\mathbf{\Sigma}}^{(k-1)} - \mathbf{\Sigma}^*\|_F \\ & \leq \frac{1}{1 - \delta} \|(\nabla f(\mathbf{\Sigma}^*))_{S^*}\|_F + \delta^{k-1} \|\widehat{\mathbf{\Sigma}}^{(1)} - \mathbf{\Sigma}^*\|_F \\ & \leq \frac{1}{1 - \delta} \|(\nabla f(\mathbf{\Sigma}^*))_{S^*}\|_F + \delta^{k-1} \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*}, \end{aligned}$$

where the last inequality is due to $\|\widehat{\mathbf{\Sigma}}^{(1)} - \mathbf{\Sigma}^*\|_F \leq \frac{2 + \sqrt{2}}{2} \lambda \sqrt{s^*}$, which follows from Lemma 11 with $k = 1$.

One has

$$\begin{aligned} \|(\nabla f(\mathbf{\Sigma}^*))_{S^*}\|_F &= \|(\mathbf{\Sigma}^* - \mathbf{S} + \tau \mathbf{I})_{S^*}\|_F \\ &\leq \|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_F + \tau \|\mathbf{I}_{S^*}\|_F \\ &\leq \|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_F + \tau \sqrt{s^*}. \end{aligned}$$

By Lemma 9, $\|(\mathbf{\Sigma}^* - \mathbf{S})_{S^*}\|_F = O_p\left(\sqrt{\frac{s^*}{n}}\right)$. If $\tau \lesssim \sqrt{\frac{1}{n}}$, then

$$\|(\nabla f(\mathbf{\Sigma}^*))_{S^*}\|_F = O_p\left(\sqrt{\frac{s^*}{n}}\right).$$

If $K \geq 1 + \frac{\log(\lambda\sqrt{n})}{\log \delta - 1} \gtrsim \log(\lambda\sqrt{n}) \gtrsim \log \log d$, then we have

$$\delta^{K-1} \lambda \sqrt{s^*} \leq \frac{1}{\lambda \sqrt{n}} \lambda \sqrt{s^*} \leq \sqrt{\frac{s^*}{n}},$$

which yields that $\|\widehat{\mathbf{\Sigma}}^{(K)} - \mathbf{\Sigma}^*\|_F = O_p\left(\sqrt{\frac{s^*}{n}}\right)$. \blacksquare