

Robust Large Covariance Matrix Estimation Under Huber Loss

Jialin Yu

Abstract—XXX

Index Terms—XXX

I. INTRODUCTION

XXX The estimation of covariance matrices is a fundamental problem in modern multivariate data analysis. It has broad applications in many fields such as statistics [1], biology [2], finance [3], [4], [5], signal processing [6], [7], [8], [9], [10], machine learning [11], etc. For example, many dimension reduction techniques, including principal component analysis [12] and linear and quadratic discriminant analysis [13], require the estimation of a covariance matrix in advance from a given collection of data points. Other prominent examples include portfolio optimization [14] and beamformer design [15]. When the dimension of the covariance matrix becomes large, the estimation problem is generally challenging. It is well-known that when the dimension is larger than the sample size, the commonly used sample covariance matrix (SCM) is singular, which may cause trouble in many applications. In addition, the number of parameters to be estimated grows quadratically with the dimension of the covariance matrix. Therefore, the problem of large (or high-dimensional) covariance matrix estimation has received considerable attention over the past decade.

In order to estimate large covariance matrices effectively, one of the most popular assumptions is sparsity, i.e., a majority of the off-diagonal elements are nearly zeros, which largely reduces the number of parameters to be estimated. The sparsity assumption is reasonable in real applications [16]. A commonly used method for estimating sparse covariance matrices is called thresholding [17], [18], [19], which is to set small elements in the SCM to zeros. Some statistical properties of the thresholding covariance estimator, including minimax lower bounds [20] and rates of convergence [21], have been established in the literature. The thresholding covariance estimator is proved only to be asymptotically positive definite [17], [18]. However, in practice it is more desirable to require the positive definiteness of the estimator under finite samples. Besides the thresholding methods for general sparse covariance matrix estimation, when covariances have sparsely banded structures, specialized covariance estimation methods include banding [22] and tapering [23], [24], [25] have been proposed. Nevertheless, this paper will focus on the general sparse cases.

Theoretical properties of large covariance estimators discussed in the literature often hinge heavily on the Gaussian or sub-Gaussian¹ assumption ([26]). See, for example, Corollary 8 of [26]. Such an assumption is typically very restrictive

in practice. For example, a recent fMRI study by Eklund et al. (2016) reported that most of the common software packages for fMRI analysis, such as SPM and FSL, can result in inflated false-positive rates up to 70% under 5% nominal levels, and questioned a number of fMRI studies among approximately 40,000 studies according to PubMed. Their results suggested that The principal cause of the invalid cluster inferences is spatial autocorrelation functions that do not follow the assumed Gaussian shape. Eklund et al. (2016) plotted the empirical versus theoretical spatial autocorrelation functions for several datasets. The empirical autocorrelation functions have much heavier tails compared to their theoretical counterparts under the commonly used assumption of a Gaussian random field, which causes the failure of fMRI inferences. Similar phenomenon has also been discovered in genomic studies (Liu et al., 2003; Purdom and Holmes, 2005) and in quantitative finance (Cont, 2001). It is therefore imperative to develop robust inferential procedures that are less sensitive to the distributional assumptions.

(Or we substitute the previous paragraph with the following explanation from Adaptive Huber Regression) The sub-Gaussian tails requirement, albeit being convenient for theoretical analysis, is not realistic in many practical applications since modern data are often collected with low quality. For example, a recent study on functional magnetic resonance imaging (fMRI) (Eklund, Nichols and Knutsson, 2016) shows that the principal cause of invalid fMRI inferences is that the data do not follow the assumed Gaussian shape, which speaks to the need of validating the statistical methods being used in the field of neuroimaging. In a microarray data example considered in Wang, Peng and Li (2015), it is observed that some gene expression levels have heavy tails as their kurtosises are much larger than 3, despite of the normalization methods used. In finance, the power-law nature of the distribution of returns has been validated as a stylized fact (Cont, 2001). Fan et al. (2016) argued that heavy-tailed distribution is a stylized feature for high dimensional data and proposed a shrinkage principle to attenuate the influence of outliers. Standard statistical procedures that are based on the method of least squares often behave poorly in the presence of heavy-tailed data² (Catoni, 2012). It is therefore of ever-increasing interest to develop new statistical methods that are robust against heavy-tailed errors and other potential forms of contamination.)

There are two ways to adjust to heavy-tailed data in existing literature: [27] assumes the polynomial-tail condition and proposes a quadratic loss ℓ_1 penalized robust covariance estimator, and methods summarized in [28] that assume finite fourth moments condition and use robust (Huber) loss in their formulations. However, neither of them achieves the oracle rate

in [26]. It is then imperative to find a formulation that achieves the oracle rate in the presence of heavy-tailedness. Implementing the robust loss with ℓ_1 penalty, however, fails to boost the performance significantly as we will demonstrate via (example or numerical analysis?). Given that the existing literature that achieves the oracle rate i.e.[26] possesses concave penalty, we propose a formulation with Huber loss and concave penalty. Under this formulation, the resulting estimator assumes only the polynomial-tail condition and achieves the oracle rate.

A. Notations

The following notation is adopted. Standard lower-case or upper-case letters stand for scalars and boldface lower-case (upper-case) letters denote vectors (matrices). Both X_{ij} and $[\mathbf{X}]_{ij}$ denote the (i, j) -th entry of the matrix \mathbf{X} . \mathbb{R}_+ denotes the set of non-negative real numbers, $\mathbb{R}^{m \times n}$ denotes the set of real $m \times n$ matrices. $\mathbf{0}$ and $\mathbf{1}$ stand for the all-zero and all-one vector/matrix, respectively. \mathbf{I} stands for the identity matrix. $\mathbf{X} \succ \mathbf{0}$ ($\mathbf{X} \succeq \mathbf{0}$) means \mathbf{X} is positive definite (semidefinite). $\mathbf{x} \geq \mathbf{0}$ denotes each element of \mathbf{x} is non-negative.

Let $\|\mathbf{X}\|_\infty = \max_{k,l} |X_{kl}|$ and $\|\mathbf{X}\|_{\min} = \min_{k,l} |X_{kl}|$. Let $\|\mathbf{X}\|_{1,\text{off}} = \sum_{k \neq l} |X_{kl}|$ denote the sum-absolute-value norm for all entries and for off-diagonals. We write $[d]$ for the set $\{1, 2, \dots, d\}$ and $\lfloor x \rfloor$ for the largest integer not exceeding x . For an index set \mathcal{E} , we use $|\mathcal{E}|$ to denote its cardinality, $\bar{\mathcal{E}}$ to denote its complement. Use $\mathbf{X}_{\mathcal{E}}$ to denote the matrix whose (i, j) -th entry is equal to X_{ij} if $(i, j) \in \mathcal{E}$, and zero otherwise. Let $\mathbf{A} \circ \mathbf{B}$ denote the Hadamard product between matrix \mathbf{A} and \mathbf{B} . Let $\partial f(\cdot)$ denote the subdifferential of a multivariate function f .

Let $\text{sgn}(x)$ denote the sign of variable x , i.e., $\text{sgn}(x) = x/|x|$. For functions $f(n)$ and $g(n)$, we denote $f(n) \leq g(n)$ if $f(n) \leq Cg(n)$, $f(n) \geq g(n)$ if $f(n) \geq cg(n)$ and $f(n) \asymp g(n)$ if $cg(n) \leq f(n) \leq Cg(n)$ for some positive constants c and C .

II. PROBLEM FORMULATION

Given samples \mathbf{x}_i , $i = 1, \dots, n$ from a heavy-tailed distribution with covariance matrix Σ^* , let $N := n(n-1)/2$ and define the paired data

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

which are identically distributed from a random vector \mathbf{y} with mean $\mathbf{0}$ and covariance matrix $\text{cov}(\mathbf{y}) = 2\Sigma^*$. Let

$$\begin{aligned} L_\alpha(\Sigma) &:= \sum_{k,\ell} \frac{1}{N} \sum_{1 \leq i < j \leq n} \rho_\alpha(\Sigma_{k\ell} - (x_{ik} - x_{jk})(x_{il} - x_{jl})/2) \\ &= \sum_{k,\ell} \frac{1}{N} \sum_{m=1}^N \rho_\alpha(\Sigma_{k\ell} - y_{mk}y_{m\ell}/2) \end{aligned}$$

with $\rho_\alpha : \mathbb{R} \rightarrow \mathbb{R}_+$ a Huber loss function defined as

$$\rho_\alpha(x) = \begin{cases} x^2/2 & \text{if } |x| \leq \alpha, \\ \alpha|x| - \alpha^2/2 & \text{if } |x| > \alpha, \end{cases} \quad (1)$$

where $\alpha > 0$ is a robustification parameter. Compared with the squared error loss, large values of x are down-weighted

in the Huber loss, yielding robustness. Generally speaking, minimizing Huber's loss produces a biased estimator of the mean, and parameter α can be chosen to control the bias. In other words, α quantifies the tradeoff between bias and robustness. As observed by Sun, Zhou and Fan (2018), in order to achieve an optimal tradeoff, α should adapt to the sample size, the dimension, and the noise level.

We will consider the following covariance estimation problem

$$\min_{\Sigma \succ \mathbf{0}} \left\{ L_\alpha(\Sigma) - \tau \log \det \Sigma + \sum_{k \neq \ell} p_\lambda(|\Sigma_{k\ell}|) \right\}, \quad (2)$$

where $\tau \log \det$ is a positive-definiteness penalty function with a regularization parameter $\tau > 0$ and $p_\lambda : \mathbb{R}_+ \rightarrow \mathbb{R}_+$ is a non-convex penalty function with a regularization parameter $\lambda > 0$. We consider a class of non-convex penalty functions $p_\lambda(\cdot)$ satisfying the following assumptions.

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

- $p_\lambda(x) = \lambda^2 p(x/\lambda)$ with some function $p(x)$ defined on $[0, +\infty)$.
- $p(x)$ is non-decreasing and concave on $[0, +\infty)$ with $p(0) = 0$ and is differentiable almost everywhere on $(0, +\infty)$;
- $0 \leq p'(x_1) \leq p'(x_2) \leq \lambda$ for all $x_1 \geq x_2 \geq 0$ and $\lim_{x \rightarrow 0} p'(x) = 1$.
- There exists some $\phi_1 > \phi_0 > \sqrt{5}$ such that $p'(\phi_0) > 0$ and $p'(x) = 0, \forall x \geq \phi_1$.

III. ALGORITHM

A. The MM Algorithmic Framework: A Brief Review

Consider the minimization of a continuous function $F(\mathbf{y})$. Initialized as $\mathbf{y}^{(0)}$, the MM algorithm [29], [30] generates a sequence of feasible points $\{\mathbf{y}^{(t)}\}_{t \geq 1}$ by the following induction. At point $\mathbf{y}^{(t-1)}$, in the majorization step, we design a surrogate function $\bar{F}(\mathbf{y} | \mathbf{y}^{(t-1)})$ that locally approximates the objective function $F(\mathbf{y})$, satisfying

$$\begin{cases} \bar{F}(\mathbf{y} | \mathbf{y}^{(t-1)}) \geq F(\mathbf{y}), \\ \bar{F}(\mathbf{y}^{(t-1)} | \mathbf{y}^{(t-1)}) = F(\mathbf{y}^{(t-1)}). \end{cases}$$

Then, in the minimization step, we update $\mathbf{y}^{(t)}$ as

$$\mathbf{y}^{(t)} \in \arg \min_{\mathbf{y}} \left\{ \bar{F}(\mathbf{y} | \mathbf{y}^{(t-1)}) \right\}.$$

B. Robust covariance matrix estimation via MM

We follow the MM framework to solve (2). In each iteration of MM, we find a weighted ℓ_1 surrogate function of $\sum_{i \neq j} p_\lambda(|\Sigma_{ij}|)$. Consequently, we consider a multistage procedure that solves a sequence of convex relaxation subproblems, which is also known as an iteratively re-weighted ℓ_1 algorithm [31] or a difference of convex algorithm [32]. Specifically, starting with an initial estimate $\tilde{\Sigma}^{(0)}$, we consider a sequence of convex optimization problems

$$\min_{\Sigma \succ \mathbf{0}} \left\{ L_\alpha(\Sigma) - \tau \log \det \Sigma + \sum_{k \neq \ell} p'_\lambda(|\widehat{\Sigma}_{kl}^{(t-1)}|) |\Sigma_{k\ell}| \right\}, \quad (3)$$

where $t = 1, 2, \dots$, and $\widehat{\Sigma}^{(t)}$ is the optimal solution to the t -th subproblem.

Each subproblem in (3) corresponds to a weighted ℓ_1 penalized covariance estimation problem, which generally can be written in the following form:

$$\min_{\Sigma > 0} \left\{ L_\alpha(\Sigma) - \tau \log \det \Sigma + \|\Lambda \circ \Sigma\|_{1,\text{off}} \right\} \quad (4)$$

where Λ is a $d \times d$ matrix of regularization parameters with $\Lambda_{ij} \in [0, \lambda]$. By convex optimization theory, any optimal solution $\widehat{\Sigma}$ to (4) satisfies the following first-order optimality condition:

$$\nabla L_\alpha(\widehat{\Sigma}) - \tau \widehat{\Sigma}^{-1} + \Lambda \circ \Xi = \mathbf{0}, \quad \text{with } \Xi \in \partial \|\widehat{\Sigma}\|_{1,\text{off}}.$$

Since analytical solution does not exist for (4), the exact solution $\widehat{\Sigma}$ can never be achieved. In practice, due to optimization error from iterative methods, we define the ϵ -optimal solution to problem (4) as follows:

Definition 2. For a pre-specified tolerance level $\epsilon > 0$, we say $\widetilde{\Sigma}$ is an ϵ -optimal solution to (4) if

$$\min_{\Xi \in \partial \|\widetilde{\Sigma}\|_{1,\text{off}}} \left\| \nabla L_\alpha(\widetilde{\Sigma}) - \tau \widetilde{\Sigma}^{-1} + \Lambda \circ \Xi \right\|_\infty \leq \epsilon.$$

We use $\widetilde{\Sigma}^{(t)}$ to denote an ϵ -optimal solution to the t -th subproblem, which is given by

$$\min_{\Sigma > 0} \left\{ L_\alpha(\Sigma) - \tau \log \det \Sigma + \|\Lambda^{(t-1)} \circ \Sigma\|_{1,\text{off}} \right\}, \quad (5)$$

where $\Lambda_{kl}^{(t-1)} = p'_\lambda(|\widetilde{\Sigma}_{kl}^{(t-1)}|)$ for all $k, l \in [d]$.

IV. THEORETICAL RESULTS

Recall that the underlying true covariance matrix is denoted by Σ^* . Let $\mathcal{S} = \{(k, l) \mid \Sigma_{kl}^* \neq 0\}$ be the support set of Σ^* and s be its cardinality, i.e., $s = |\mathcal{S}|$. In the following, we impose some mild conditions on the true covariance matrix Σ^* and the distribution of the i.i.d. samples \mathbf{x}_i , $i = 1, \dots, n$.

Assumption 3. For the true covariance matrix, assume $\Sigma^* \succ \mathbf{0}$.

Assumption 4. $\mathbf{x}_i \in \mathbb{R}^d$ is a heavy-tailed random variable, i.e. $\mathbb{E}[|x_{ij}|^{4(1+\gamma)}] \leq \sigma_x^{2(1+\gamma)}$ for all $1 \leq j \leq d$ with some positive σ_x .

We assume Assumptions 1, 3, and 4 hold for the rest of the paper. (declare this within every result.)

Definition 5. Given $\mathbb{B}^\infty(r) := \{\Delta \in \mathbb{R}^{d \times d} : \|\Delta\|_\infty \leq r\}$. We define $\mathcal{E}_1(r, \kappa)$ for any $\Sigma \in \Sigma^* + \mathbb{B}^\infty(r)$,

$$\langle \nabla L_\alpha(\Sigma) - \nabla L_\alpha(\Sigma^*), \Sigma - \Sigma^* \rangle \geq \kappa \|\Sigma - \Sigma^*\|_F^2.$$

Proposition 6. Suppose that Assumptions 1 and 3 hold. Let $\tau, \lambda, r > 0$ satisfy

$$\tau \leq \frac{\lambda s^{1/2}}{\|(\Sigma^*)^{-1}\|_F}, \quad r > 5\lambda s^{1/2}.$$

Then, conditioned on the event $\mathcal{E}_1(r, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$, any ϵ -optimal solution to (4) satisfies

$$\left\| \widetilde{\Sigma}^{(1)} - \Sigma^* \right\|_F \leq 5\lambda s^{1/2}$$

If we take $\lambda \asymp \sqrt{\frac{\log d}{n}}$ in Proposition 6, then we can see that the rate of convergence $\sqrt{\frac{s \log d}{n}}$ is the same as that obtained by Wei's paper.

Definition 7. Let $\mathbb{C}(l) := \{\Delta \in \mathbb{R}^{d \times d} : \|\Delta\|_1 \leq l \|\Delta\|_F\}$.

The definition of $\mathbb{C}(l)$ converts a bound for $\left\| \widetilde{\Sigma}^{(T)} - \Sigma^* \right\|_F$ into a bound for $\left\| \widetilde{\Sigma}^{(T)} - \Sigma^* \right\|_1$. Assume that we have bounded $\left\| \widetilde{\Sigma}^{(T)} - \Sigma^* \right\|_F$, i.e. $\left\| \widetilde{\Sigma}^{(T)} - \Sigma^* \right\|_F \leq M$ for some $M > 0$. Then, given that

$$\widetilde{\Sigma}^{(T)} \in \Sigma^* + \mathbb{C}(l), \quad (6)$$

we have $\left\| \widetilde{\Sigma}^{(T)} - \Sigma^* \right\|_1 \leq lM$. To guarantee (6), a trivial choice would be $l \asymp d$, but we will show that $l \asymp s^{1/2}$ is actually enough.

The following Proposition demonstrates the contraction property of the solution path $\left\{ \widetilde{\Sigma}^{(t)} \right\}_{t \geq 1}$.

Proposition 8. Suppose that Assumptions 1 and 3 hold. By Assumption 1, there exists some $\phi_0 > \sqrt{5}$ such that $p'(\phi_0) > 0$. Let

$$\tau \leq \frac{\lambda s^{1/2}}{\|(\Sigma^*)^{-1}\|_F} \quad (7)$$

and choose $c > 0$ so that

$$0.5p'(\phi_0)(c^2 + 1)^{1/2} + 2 = 0.5c\phi_0 \quad (8)$$

Set $l = (2 + \frac{2}{p'(\phi_0)})(c^2 + 1)^{1/2} s^{1/2} + \frac{2}{p'(\phi_0)} s^{1/2}$ and let $r > 0$ satisfy

$$c\phi_0 \lambda s^{1/2} \leq r \quad (9)$$

Under the minimal signal strength condition $\|\Sigma_{\mathcal{S}}^*\|_{\min} \geq \phi_0 \lambda$ and conditioned on event $\mathcal{E}_1(r, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda\}$, any ϵ -optimal solution $\widetilde{\Sigma}^{(t)}$ to (5) satisfies $\widetilde{\Sigma}^{(t)} \in \Sigma^* + \mathbb{C}(l)$, and we have

$$\left\| \widetilde{\Sigma}^{(t)} - \Sigma^* \right\|_F \leq \delta \left\| \widetilde{\Sigma}^{(t-1)} - \Sigma^* \right\|_F + 2 \left\{ \|p'_\lambda(|\Sigma_{\mathcal{S}}^*| - \phi_0 \lambda)\|_F + s^{1/2} \epsilon + \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_F + \tau \|(\Sigma^*)^{-1}\|_F \right\} \quad (10)$$

where $\delta = \sqrt{5}/\phi_0 \in (0, 1)$.

Remark 9. Let

$$r^{\text{ora}} := 2 \left\{ \|p'_\lambda(|\Sigma_{\mathcal{S}}^*| - \phi_0 \lambda)\|_F + \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_F + s^{1/2} \epsilon + \tau \|(\Sigma^*)^{-1}\|_F \right\},$$

then Proposition 2.2 states that

$$\left\| \widetilde{\Sigma}^{(t)} - \Sigma^* \right\|_F \leq \delta \left\| \widetilde{\Sigma}^{(t-1)} - \Sigma^* \right\|_F + r^{\text{ora}},$$

where $\delta = \sqrt{5}/\phi_0 \in (0, 1)$. And by induction we further have

$$\left\| \tilde{\Sigma}^{(t)} - \Sigma^* \right\|_F \leq \delta^{t-1} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_F + (1 - \delta)^{-1} r^{\text{ora}}.$$

In Proposition 8, we have demonstrated the exact choice of l in 6, which satisfies $l \asymp s^{1/2}$. Meanwhile, we have shown that the estimation error between the ϵ -optimal solution $\tilde{\Sigma}^{(t)}$ and the true covariance Σ^* is bounded by two terms, namely, r^{ora} and a contraction term.

Lemma 10. *Suppose that Assumption 4 holds. Define $\epsilon_{kl} = \Sigma_{kl}^* - y_{mk}y_{m\ell}/2$ for all $k, l \in [d]$. Then $\max \left\{ \mathbb{E} \left[\epsilon_{kl}^2 \right], \mathbb{E} \left[\epsilon_{kl}^{2(1+\gamma)} \right] \right\} \leq K$ for all $k, l \in [d]$ with some constant K that depends only on σ_x and γ .*

Note that for the sake of simplicity, in 10, $\mathbb{E} \left[\epsilon_{kl}^2 \right]$ and $\mathbb{E} \left[\epsilon_{kl}^{2(1+\gamma)} \right]$ are bounded with the same constant K . Lemma 10 is a direct implication of Assumption 4. It uses the moment condition on \mathbf{x}_i to achieve the moment conditions on ϵ_{kl} , which will play an important role in the following propositions.

Proposition 11. *Suppose that Assumption 4 holds. Let K be the constant defined in Lemma (10). Assume α satisfies $4K/\alpha^2 < 1/6$. Then, with at least $1 - d^2 \exp(-n/12)$ probability, for all $\Sigma \in \Sigma^* + \mathbb{B}^\infty(\alpha/2)$,*

$$\langle \nabla L_\alpha(\Sigma) - \nabla L_\alpha(\Sigma^*), \Sigma - \Sigma^* \rangle \geq \frac{1}{2} \|\Sigma - \Sigma^*\|_F^2$$

Proposition 11 implies that with proper assumptions on α , $\mathcal{E}_1(\alpha/2, 1/2)$ happens with high probability. For instance, by taking $n \gtrsim (\log d)^{1+\frac{1}{2\gamma}}$, we have $1 - d^2 \exp(-n/12) \geq 1 - 2/d$.

Proposition 12. *Suppose that Assumption 4 holds. Let K be the constant defined in Lemma (10). Assume $\alpha = \sqrt{Kn}/\log d$, then*

$$\|\nabla L_\alpha(\Sigma^*)\|_\infty \leq 8\sqrt{\frac{K \log d}{n}} \quad (11)$$

with at least $1 - 2/d$ probability. Furthermore, if we have $n \geq K^{-1}(\log d)^{1+\frac{1}{2\gamma}}$, then for any $\beta > 0$,

$$\Pr \left\{ \|\nabla L_\alpha(\Sigma^*)\|_F \geq (\beta + 1)\sqrt{\frac{Ks}{n}} \right\} \leq \frac{2}{\beta}. \quad (12)$$

In Proposition 12, 11 indicates that $\{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda\}$ happens with high probability if we take $\epsilon \leq \sqrt{1/n}$ and $\lambda \asymp \sqrt{\log d/n}$. 12 implies that for any $p \in [0, 1)$, there exists β such that $\|\nabla L_\alpha(\Sigma^*)\|_F \leq (\beta + 1)\sqrt{Ks/n}$ holds with at least p probability as long as $n \geq K^{-1}(\log d)^{1+\frac{1}{2\gamma}} \asymp (\log d)^{1+\frac{1}{2\gamma}}$. Weakly speaking, $\|\nabla L_\alpha(\Sigma^*)\|_F = O_P(\sqrt{s/n})$.

Theorem 13. *Suppose that Assumptions 1, 3 and 4 hold. By Assumption 1, there exists some $\phi_1 > \phi_0 > \sqrt{5}$ such that*

$$p'(\phi_0) > 0 \quad \text{and} \quad p'(\phi) = 0, \forall \phi \geq \phi_1$$

Suppose the sample size satisfies $n \gtrsim (\log d)^{1+\frac{1}{2\gamma}}$. Take $\lambda \asymp \sqrt{\log d/n}$, $l \asymp s^{1/2}$ and let $\alpha = \sqrt{Kn}/\log d$, $r = \alpha/2$, $\tau \leq \|(\Sigma^*)^{-1}\|_F^{-1} \cdot \sqrt{Ks/n}$, $\epsilon \leq \sqrt{K/n}$. Then, under the

minimal signal strength condition $\|\Sigma_S^*\|_{\min} \geq (\phi_0 + \phi_1)\lambda$, the multi-stage estimator $\tilde{\Sigma}^{(T)}$ with $T \gtrsim \frac{\log(\log d)}{\log(1/\delta)}$ satisfies the bound

$$\left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F \lesssim \beta \sqrt{\frac{s}{n}}, \quad \left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_I \lesssim \beta \frac{s}{\sqrt{n}}$$

with at least $1 - 2/d - d^2 \exp(-n/12) - 2/\beta$ probability and K to be the constant defined in Lemma (10). This immediately implies the following weaker conclusion

$$\left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F = O_P \left(\sqrt{\frac{s}{n}} \right), \quad \left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_I = O_P \left(\frac{s}{\sqrt{n}} \right)$$

with at least $1 - 2/d$ probability.

Theorem 13 is a direct consequence of the combination of Proposition 8, 11 and 12, which implies that under weak assumptions, we just need to solve no more than approximately $\log \log d$ convex problems to achieve the oracle rate $\sqrt{s/n}$. It is easy to see that the proposed estimator achieves the oracle statistical rate of convergence under weaker assumptions on the distribution of sample data \mathbf{x}_i . In existing literature, Wei's paper proposes a formulation with squared loss that requires sub-Gaussian assumption to achieve the oracle rate. With Huber loss in our formulation, we only need to assume the existence of $4(1 + \gamma)$ -th moment of \mathbf{x}_i . In other words, our estimator is capable of handling heavy-tailed data.

V. NUMERICAL SIMULATIONS

VI. CONCLUSIONS

XXX

APPENDIX A

PROOFS OF STATISTICAL THEORY

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

Lemma 14. *Let \mathcal{E} be a subset of $[d]$ that contains \mathcal{S} . For any $\Sigma \in \mathbb{R}^{d \times d}$ satisfying $\Sigma_{\mathcal{E}} = \mathbf{0}$ and $\epsilon > 0$, provided $\Lambda = (\Lambda)_{kl}$ satisfies $\|\Lambda_{\mathcal{E}}\|_{\min} > \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_\infty + \epsilon$, any ϵ -optimal solution $\tilde{\Sigma}$ to (4) satisfies*

$$\begin{aligned} & \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_1 \\ & \leq (\|\Lambda_{\mathcal{E}}\|_{\min} - \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_\infty - \epsilon)^{-1} \\ & \quad \cdot \left\{ (\|\Lambda_{\mathcal{E}}\|_\infty + \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_\infty + \epsilon) \cdot \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_1 \right. \\ & \quad \left. + \|\tau \Sigma^{-1}\|_F \cdot \left\| \tilde{\Sigma} - \Sigma \right\|_F \right\} \end{aligned}$$

Proof: For any $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1, \text{off}}$, define $U(\Xi) = \nabla L_\alpha(\tilde{\Sigma}) - \tau \tilde{\Sigma}^{-1} + \Lambda \circ \Xi \in \mathbb{R}^{d \times d}$. By convexity of $L_\alpha(\Sigma)$ and $-\log \det \Sigma$:

$$\langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \geq 0 \quad \text{and} \quad \langle \Sigma^{-1} - \tilde{\Sigma}^{-1}, \tilde{\Sigma} - \Sigma \rangle \geq 0.$$

Therefore,

$$\begin{aligned}
& \langle \mathbf{U}(\Xi), \tilde{\Sigma} - \Sigma \rangle \geq \|\mathbf{U}(\Xi)\|_\infty \left\| \tilde{\Sigma} - \Sigma \right\|_1 \\
& = \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle + \langle \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\
& \quad - \langle \tau \Sigma^{-1}, \tilde{\Sigma} - \Sigma \rangle + \langle \tau \Sigma^{-1} - \tau \tilde{\Sigma}^{-1}, \tilde{\Sigma} - \Sigma \rangle \\
& \quad + \langle \Lambda \circ \Xi, \tilde{\Sigma} - \Sigma \rangle \\
& \geq 0 + \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 \\
& \quad + \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
& \quad + \|\tau \Sigma^{-1}\|_F \cdot \left\| \tilde{\Sigma} - \Sigma \right\|_F + 0 + \langle \Lambda \circ \Xi, \tilde{\Sigma} - \Sigma \rangle
\end{aligned}$$

Moreover, we have

$$\begin{aligned}
& \langle \Lambda \circ \Xi, \tilde{\Sigma} - \Sigma \rangle \\
& = \langle (\Lambda \circ \Xi)_{\bar{\mathcal{E}}}, (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \rangle + \langle (\Lambda \circ \Xi)_\mathcal{E}, (\tilde{\Sigma} - \Sigma)_\mathcal{E} \rangle \\
& \geq \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 - \|\Lambda_\mathcal{E}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1
\end{aligned}$$

Together, the last two displays imply

$$\begin{aligned}
& \|\mathbf{U}(\Xi)\|_\infty \left\| \tilde{\Sigma} - \Sigma \right\|_1 \\
& \geq \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 + \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
& \quad + \|\tau \Sigma^{-1}\|_F \cdot \left\| \tilde{\Sigma} - \Sigma \right\|_F + \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
& \quad - \|\Lambda_\mathcal{E}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1
\end{aligned}$$

Since the right-hand side of this inequality does not depend on Ξ , taking the infimum with respect to $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1, \text{off}}$ on both sides to reach

$$\begin{aligned}
& \epsilon \left\| \tilde{\Sigma} - \Sigma \right\|_1 \\
& \geq \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 + \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
& \quad + \|\tau \Sigma^{-1}\|_F \cdot \left\| \tilde{\Sigma} - \Sigma \right\|_F + \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
& \quad - \|\Lambda_\mathcal{E}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1
\end{aligned}$$

Decompose $\left\| \tilde{\Sigma} - \Sigma \right\|_1$ as $\left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 + \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1$, the stated result follows immediately. \blacksquare

Lemma 15. Consider some $\Sigma \in \mathbb{R}^{d \times d}$ satisfying $\Sigma_{\mathcal{S}^c} = 0$, and let $\mathcal{E} \subseteq [d]$ be a subset that contains \mathcal{S} that has cardinality $|\mathcal{E}| = k$. Assume that $\Lambda = (\Lambda)_{kl}$ satisfies $\|\Lambda\|_\infty \leq \lambda$ and $\|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \geq \rho\lambda > 0$ for some $\rho \in (0, 1]$, $\tau \leq \frac{\lambda s^{1/2}}{\|\Sigma^{-1}\|_F}$. Conditioned on event $\{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\rho\lambda\}$, any ϵ -optimal solution $\tilde{\Sigma}$ to (4) satisfies $\tilde{\Sigma} \in \Sigma + \mathcal{C}(l)$, where $l = \left(2 + \frac{2}{\rho}\right) k^{1/2} + \frac{2}{\rho} s^{1/2}$. Moreover, let $r, \kappa > 0$ satisfy

$$r > \kappa^{-1} \left\{ 2\lambda s^{1/2} + 0.5\rho\lambda k^{1/2} \right\}.$$

Then, conditioned on the event $\mathcal{E}_1(r, \kappa) \cap \{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\rho\lambda\}$,

$$\begin{aligned}
\left\| \tilde{\Sigma} - \Sigma \right\|_F & \leq \kappa^{-1} \left\{ \|\Lambda_{\mathcal{S}}\|_F + \|\nabla L_\alpha(\Sigma)\|_\mathcal{E} + k^{1/2}\epsilon + \|\tau \Sigma^{-1}\|_F \right\} \\
& \leq \kappa^{-1} \left\{ 2\lambda s^{1/2} + 0.5\rho\lambda k^{1/2} \right\} < r
\end{aligned}$$

Proof: Conditioned on the stated event, Lemma 14 indicates

$$\left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \leq \left(1 + \frac{2}{\rho}\right) \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 + \frac{2}{\rho} \sqrt{s} \left\| \tilde{\Sigma} - \Sigma \right\|_F.$$

Therefore,

$$\left\| \tilde{\Sigma} - \Sigma \right\|_1 \leq \left(1 + \frac{2}{\rho}\right) \sqrt{k} \left\| \tilde{\Sigma} - \Sigma \right\|_F + \frac{2}{\rho} \sqrt{s} \left\| \tilde{\Sigma} - \Sigma \right\|_F,$$

which implies that $\tilde{\Sigma} \in \Sigma + \mathcal{C}(l)$.

Now we prove the second statement. Define $\eta = \sup\{u \in [0, 1] : (1-u)\Sigma + u\tilde{\Sigma} \in \mathbb{B}(r)\}$, where $\mathbb{B}(r) = \{\Delta \in \mathbb{R}^{d \times d} : \|\Delta\|_F \leq r\}$. Note that $\eta = 1$ if $\tilde{\Sigma} \in \Sigma + \mathbb{B}(r)$ and $\eta \in (0, 1)$ otherwise. Let $\tilde{\Sigma}_\eta := (1-\eta)\Sigma + \eta\tilde{\Sigma}$. Notice that if $\left\| \tilde{\Sigma}_\eta - \Sigma \right\|_F < r$, then $\tilde{\Sigma}_\eta = \tilde{\Sigma}$. By the convexity of Huber loss, we have

$$\begin{aligned}
& \langle \nabla L_\alpha(\tilde{\Sigma}_\eta) - \nabla L_\alpha(\Sigma), \tilde{\Sigma}_\eta - \Sigma \rangle \\
& \leq \eta \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \quad (13)
\end{aligned}$$

Recall that $\mathbb{B}^\infty(r) := \{\Delta \in \mathbb{R}^{d \times d} : \|\Delta\|_\infty \leq r\}$. Since $\tilde{\Sigma}_\eta - \Sigma \in \mathbb{B}(r) \subseteq \mathbb{B}^\infty(r)$, conditioned on event \mathcal{E}_1 , we have

$$\langle \nabla L_\alpha(\tilde{\Sigma}_\eta) - \nabla L_\alpha(\Sigma), \tilde{\Sigma}_\eta - \Sigma \rangle \geq \kappa \left\| \tilde{\Sigma}_\eta - \Sigma \right\|_F^2 \quad (14)$$

Now we upper bound the right-hand side of (13). For any $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1, \text{off}}$, write

$$\begin{aligned}
& \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\
& = \underbrace{\langle \mathbf{U}(\Xi), \tilde{\Sigma} - \Sigma \rangle}_{:= \Pi_1} - \underbrace{\langle \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle}_{:= \Pi_2} \\
& \quad - \underbrace{\langle \Lambda \circ \Xi, \tilde{\Sigma} - \Sigma \rangle}_{:= \Pi_3} + \underbrace{\langle \tau \Sigma^{-1}, \tilde{\Sigma} - \Sigma \rangle}_{:= \Pi_4} - \underbrace{\tau \langle \Sigma^{-1} - \tilde{\Sigma}^{-1}, \tilde{\Sigma} - \Sigma \rangle}_{\geq 0} \quad (15)
\end{aligned}$$

where $\mathbf{U}(\Xi) := \nabla L_\alpha(\tilde{\Sigma}) - \tau \tilde{\Sigma}^{-1} + \Lambda \circ \Xi \in \mathbb{R}^{d \times d}$. We have

$$\begin{aligned}
|\Pi_1| & \leq \|\mathbf{U}(\Xi)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 + \|\mathbf{U}(\Xi)\|_\mathcal{E} \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 \\
|\Pi_2| & \leq \|\nabla L_\alpha(\Sigma)\|_\mathcal{E} \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{E} \right\|_1 + \|\nabla L_\alpha(\Sigma)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 \\
|\Pi_4| & \leq \|\tau \Sigma^{-1}\|_F \cdot \left\| \tilde{\Sigma} - \Sigma \right\|_F
\end{aligned}$$

Turning to Π_3 , decompose $\Lambda \circ \Xi$ and $\tilde{\Sigma} - \Sigma$ according to $\mathcal{S} \cup (\mathcal{E}/\mathcal{S}) \cup \bar{\mathcal{E}}$ to reach

$$\begin{aligned}
\Pi_3 & = \langle (\Lambda \circ \Xi)_\mathcal{S}, (\tilde{\Sigma} - \Sigma)_\mathcal{S} \rangle \\
& \quad + \langle (\Lambda \circ \Xi)_{\mathcal{E}/\mathcal{S}}, (\tilde{\Sigma} - \Sigma)_{\mathcal{E}/\mathcal{S}} \rangle + \langle (\Lambda \circ \Xi)_{\bar{\mathcal{E}}}, (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \rangle
\end{aligned}$$

Since $\Sigma_{\bar{\mathcal{E}}} = \mathbf{0}$ and $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1, \text{off}}$, we have $\langle (\Lambda \circ \Xi)_{\bar{\mathcal{E}}}, (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \rangle = \langle \Lambda_{\bar{\mathcal{E}}}, \tilde{\Sigma}_{\bar{\mathcal{E}}} \rangle = \langle \Lambda_{\bar{\mathcal{E}}}, \left| \tilde{\Sigma} - \Sigma \right|_{\bar{\mathcal{E}}} \rangle$. Also, $\langle (\Lambda \circ \Xi)_{\mathcal{E}/\mathcal{S}}, (\tilde{\Sigma} - \Sigma)_{\mathcal{E}/\mathcal{S}} \rangle = \langle (\Lambda \circ \Xi)_{\mathcal{E}/\mathcal{S}}, \tilde{\Sigma}_{\mathcal{E}/\mathcal{S}} \rangle \geq 0$.

$$\begin{aligned}
\Pi_3 & \geq \|\Lambda_{\bar{\mathcal{E}}}\|_{\min} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{\mathcal{E}}} \right\|_1 - \|\Lambda_{\mathcal{S}}\|_F \left\| (\tilde{\Sigma} - \Sigma)_\mathcal{S} \right\|_F
\end{aligned}$$

Combining (15) with our estimation for Π_1, Π_2, Π_3 and Π_4 , we have

$$\begin{aligned} & \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\ & \leq -\{\|\Lambda_{\tilde{\mathcal{E}}}\|_{\min} - \|\nabla L_\alpha(\Sigma)\|_\infty - \|U(\Xi)\|_\infty\} \left\| (\tilde{\Sigma} - \Sigma)_{\tilde{\mathcal{E}}} \right\|_1 \\ & \quad + \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_{\mathbb{F}} \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_{\mathbb{F}} + \|(U(\Xi))_{\mathcal{E}}\|_{\mathbb{F}} \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_{\mathbb{F}} \\ & \quad + \|\Lambda_{\mathcal{S}}\|_{\mathbb{F}} \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_{\mathbb{F}} + \|\tau \Sigma^{-1}\|_{\mathbb{F}} \left\| \tilde{\Sigma} - \Sigma \right\|_{\mathbb{F}} \end{aligned}$$

Taking the infimum with respect to $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1, \text{off}}$ on both sides, it follows that

$$\begin{aligned} & \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\ & \leq -\{\|\lambda_{\tilde{\mathcal{E}}}\|_{\min} - \|\nabla L_\alpha(\Sigma)\|_\infty - \epsilon\} \left\| (\tilde{\Sigma} - \Sigma)_{\tilde{\mathcal{E}}} \right\|_1 \\ & \quad + \{\|\nabla L_\alpha(\Sigma)\|_{\mathbb{F}} + k^{1/2}\epsilon\} \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_{\mathbb{F}} \\ & \quad + \|\Lambda_{\mathcal{S}}\|_{\mathbb{F}} \left\| (\tilde{\Sigma} - \Sigma)_{\mathcal{E}} \right\|_{\mathbb{F}} + \|\tau \Sigma^{-1}\|_{\mathbb{F}} \left\| \tilde{\Sigma} - \Sigma \right\|_{\mathbb{F}} \end{aligned} \quad (16)$$

With $\tilde{\Sigma}_\eta - \Sigma = \eta(\tilde{\Sigma} - \Sigma)$, it follows from (13), (14) and (16) that conditioned on $\mathcal{E}_1(r, \kappa) \cap \{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\rho\lambda\}$,

$$\begin{aligned} & \kappa \left\| \tilde{\Sigma}_\eta - \Sigma \right\|_{\mathbb{F}}^2 \leq \\ & \left\{ \|\Lambda_{\mathcal{S}}\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_{\mathbb{F}} + k^{1/2}\epsilon + \|\tau \Sigma^{-1}\|_{\mathbb{F}} \right\} \left\| \tilde{\Sigma}_\eta - \Sigma \right\|_{\mathbb{F}} \end{aligned}$$

Therefore,

$$\begin{aligned} & \left\| \tilde{\Sigma}_\eta - \Sigma \right\|_{\mathbb{F}} \\ & \leq \kappa^{-1} \left\{ \|\Lambda_{\mathcal{S}}\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma)_{\mathcal{E}}\|_{\mathbb{F}} + k^{1/2}\epsilon + \|\tau \Sigma^{-1}\|_{\mathbb{F}} \right\} \\ & \leq \kappa^{-1} \{ \lambda s^{1/2} + 0.5\rho\lambda k^{1/2} + \lambda s^{1/2} \} < r \end{aligned} \quad (17)$$

Since $\tilde{\Sigma}_\eta - \Sigma$ falls in the interior of $\mathbb{B}(r)$, we must have $\tilde{\Sigma} - \Sigma = \tilde{\Sigma}_\eta - \Sigma \in \mathbb{B}(r)$. Consequently, (17) also holds for $\tilde{\Sigma} - \Sigma$. ■

A. Proof of Proposition 6

Proof: With the initial estimate $\tilde{\Sigma}^{(0)} = \mathbf{I}$, we have $\Lambda_{kl}^{(0)} = p'_\lambda(0) = \lambda$ for all $k, l \in [d]$. Then the result follows immediately from Lemma 15 with $\Sigma = \Sigma^*$, $\mathcal{E} = \mathcal{S}$ and $\rho = 1$. ■

B. Proof of Proposition 8

Proof: With the initial estimate $\tilde{\Sigma}^{(0)} = \mathbf{0}$, we have $\Lambda_{kl}^{(0)} = p'_\lambda(0) = \lambda$ for all $k, l \in [d]$. Then, applying Lemma 15 with $\Sigma = \Sigma^*$, $\mathcal{E} = \mathcal{S}$ and $\rho = p'(\phi_0)$ we obtain that, conditioned on the event $\mathcal{E}_1(r, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda\}$,

$$\begin{aligned} & \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_{\mathbb{F}} \\ & \leq 2 \left\{ \left\| \Lambda_{\mathcal{S}}^{(0)} \right\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_{\mathbb{F}} + s^{1/2}\epsilon + \tau \|(\Sigma^*)^{-1}\|_{\mathbb{F}} \right\} \\ & \leq 2 \{2 + 0.5p'(\phi_0)\} s^{1/2}\lambda \end{aligned} \quad (18)$$

For $t \geq 1$, define the augmented set

$$\mathcal{E}_t = \mathcal{S} \cup \left\{ (k, l) : \Lambda_{kl}^{(t-1)} < p'(\phi_0)\lambda \right\}$$

which depends on the solution $\tilde{\Sigma}^{(t-1)}$ from the previous step. We claim that the above constructed sets satisfy

$$|\mathcal{E}_t| < (c^2 + 1)s \quad (19)$$

If (19) is true, it follows from Lemma 15 with $\Sigma = \Sigma^*$, $\mathcal{E} = \mathcal{E}_t$, $k = (c^2 + 1)$ and $\rho = p'(\phi_0)$ that conditioned on the event $\mathcal{E}_1(r, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda\}$, we have $\tilde{\Sigma}^{(t)} \in \Sigma^* + \mathcal{C}(l)$ and

$$\begin{aligned} & \left\| \tilde{\Sigma}^{(t)} - \Sigma^* \right\|_{\mathbb{F}} \\ & \leq 2 \left\{ \left\| \Lambda_{\mathcal{S}}^{(t-1)} \right\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma^*)_{\mathcal{E}_t}\|_{\mathbb{F}} + |\mathcal{E}_t|^{1/2}\epsilon + \tau \|(\Sigma^*)^{-1}\|_{\mathbb{F}} \right\} \\ & \leq 2 \left\{ 2\lambda s^{1/2} + 0.5p'(\phi_0)\lambda(c^2 + 1)^{1/2}s^{1/2} \right\} = c\phi_0\lambda s^{1/2} \leq r \end{aligned} \quad (20)$$

where the last two steps follow from (8) and (9). We will prove our claim (19) by induction. For $t = 1$, we have $\Lambda_{kl}^{(0)} = p'_\lambda(0) = \lambda$ for all $k, l \in [d]$, so $\mathcal{E}_1 = \mathcal{S}$ and (19) holds.

Next, assume (19) holds for some $t \geq 1$, from which (20) follows. To bound $|\mathcal{E}_{t+1}|$, note that for any $(k, l) \in \mathcal{E}_{t+1} \setminus \mathcal{S}$, $p'_\lambda(\tilde{\Sigma}_{kl}^{(t)}) = \Lambda_{kl}^{(t)} < p'(\phi_0)\lambda = p'_\lambda(\phi_0\lambda)$. Together with the monotonicity of p'_λ , this implies that $|\tilde{\Sigma}_{kl}^{(t)}| > \phi_0\lambda$. Recalling that $\Sigma_{kl}^* = 0$ for $(k, l) \in \mathcal{E}_{t+1} \setminus \mathcal{S}$, we obtain

$$\begin{aligned} |\mathcal{E}_{t+1} \setminus \mathcal{S}|^{1/2} & < \left\{ \sum_{(k,l) \in \mathcal{E}_{t+1} \setminus \mathcal{S}} \left(|\tilde{\Sigma}_{kl}^{(t)}| / \phi_0\lambda \right)^2 \right\}^{1/2} \\ & = \left\| (\tilde{\Sigma}^{(t+1)} - \Sigma^*)_{\mathcal{E}_{t+1} \setminus \mathcal{S}} \right\|_{\mathbb{F}} / (\phi_0\lambda) \leq cs^{1/2} \end{aligned} \quad (21)$$

The last inequality follows from (20). Therefore $|\mathcal{E}_{t+1}| \leq |\mathcal{S}| + |\mathcal{E}_{t+1} \setminus \mathcal{S}| < (c^2 + 1)^{1/2}s$, which completes the induction step.

Now we start to prove (10). Recall that for each (k, l) , $\Lambda_{kl}^{(t-1)} = p'_\lambda(\tilde{\Sigma}_{kl}^{(t-1)})$. If $|\tilde{\Sigma}_{kl}^{(t-1)} - \Sigma_{kl}^*| \geq \phi_0\lambda$, then $\Lambda_{kl}^{(t-1)} \leq \lambda \leq |\tilde{\Sigma}_{kl}^{(t-1)} - \Sigma_{kl}^*| / \phi_0$; otherwise if $|\tilde{\Sigma}_{kl}^{(t-1)} - \Sigma_{kl}^*| < \phi_0\lambda$, then $\Lambda_{kl}^{(t-1)} = p'_\lambda(\tilde{\Sigma}_{kl}^{(t-1)}) \leq p'_\lambda(|\tilde{\Sigma}_{kl}^{(t-1)} - \Sigma_{kl}^*| + \phi_0\lambda)$ due to the monotonicity of p'_λ . Putting together the pieces, we conclude that

$$\left\| \Lambda_{\mathcal{S}}^{(t-1)} \right\|_{\mathbb{F}} \leq \|p'_\lambda(|\Sigma_{\mathcal{S}}^*| - \phi_0\lambda)\|_{\mathbb{F}} + \phi_0^{-1} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{S}} \right\|_{\mathbb{F}} \quad (22)$$

For the remaining terms that involve \mathcal{E}_t in (20), by the triangle inequality

$$\begin{aligned} & \|\nabla L_\alpha(\Sigma^*)_{\mathcal{E}_t}\|_{\mathbb{F}} + |\mathcal{E}_t|^{1/2}\epsilon \\ & \leq \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_{\mathbb{F}} + s^{1/2}\epsilon + \|\nabla L_\alpha(\Sigma^*)_{\mathcal{E}_t \setminus \mathcal{S}}\|_{\mathbb{F}} + |\mathcal{E}_t \setminus \mathcal{S}|^{1/2}\epsilon \\ & \leq \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_{\mathbb{F}} + s^{1/2}\epsilon \\ & \quad + (\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon) \cdot \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{E}_t \setminus \mathcal{S}} \right\|_{\mathbb{F}} / (\phi_0\lambda) \\ & \leq \|\nabla L_\alpha(\Sigma^*)_{\mathcal{S}}\|_{\mathbb{F}} + s^{1/2}\epsilon + \frac{p'(\phi_0)}{2\phi_0} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{E}_t \setminus \mathcal{S}} \right\|_{\mathbb{F}} \end{aligned} \quad (23)$$

where the second equality follows from (21) and the last inequality follows from $\{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda\}$. Plugging (22) and (23) into (20) yields

$$\begin{aligned}
& \left\| \tilde{\Sigma}^{(t)} - \Sigma^* \right\|_{\mathbb{F}} \\
& \leq 2 \left\{ \|p'_\lambda(|\Sigma_S^*| - \phi_0\lambda)\|_{\mathbb{F}} + \phi_0^{-1} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_S \right\|_{\mathbb{F}} \right. \\
& \quad \left. + \|\nabla L_\alpha(\Sigma^*)_S\|_{\mathbb{F}} + s^{1/2}\epsilon \right. \\
& \quad \left. + \frac{p'(\phi_0)}{2\phi_0} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{E}_t \setminus S} \right\|_{\mathbb{F}} + \tau \|(\Sigma^*)^{-1}\|_{\mathbb{F}} \right\} \\
& \leq 2 \left\{ \|p'_\lambda(|\Sigma_S^*| - \phi_0\lambda)\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma^*)_S\|_{\mathbb{F}} + s^{1/2}\epsilon \right. \\
& \quad \left. + \frac{\sqrt{5}}{2\phi_0} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{E}_t} \right\|_{\mathbb{F}} + \tau \|(\Sigma^*)^{-1}\|_{\mathbb{F}} \right\} \\
& = 2 \left\{ \|p'_\lambda(|\Sigma_S^*| - \phi_0\lambda)\|_{\mathbb{F}} + \|\nabla L_\alpha(\Sigma^*)_S\|_{\mathbb{F}} + s^{1/2}\epsilon \right. \\
& \quad \left. + \tau \|(\Sigma^*)^{-1}\|_{\mathbb{F}} \right\} + \frac{\sqrt{5}}{\phi_0} \left\| (\tilde{\Sigma}^{(t-1)} - \Sigma^*)_{\mathcal{E}_t} \right\|_{\mathbb{F}}
\end{aligned}$$

The second inequality follows because $p'(\phi_0) \leq 1$ and for any $a, b \geq 0$, $\sqrt{a} + \sqrt{b/4} \leq \sqrt{5(a+b)/4}$. The proof of (10) is then completed. \blacksquare

C. Proof of Lemma 10

Proof: Let $\sigma_y := 4\sigma_x$. Given that $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N\}$ are identically distributed and with the c_r -inequality,

$$\begin{aligned}
\mathbb{E} \left[|y_{mj}|^{4(1+\gamma)} \right] &= \mathbb{E} \left[|x_{1j} - x_{2j}|^{4(1+\gamma)} \right] \\
&\leq 2^{4(1+\gamma)-1} \left\{ \mathbb{E} \left[|x_{1j}|^{4(1+\gamma)} \right] + \mathbb{E} \left[|x_{2j}|^{4(1+\gamma)} \right] \right\} \\
&\leq 2^{4(1+\gamma)} \sigma_x^{2(1+\gamma)} = \sigma_y^{2(1+\gamma)}, \quad \forall j \in [d].
\end{aligned}$$

The second inequality follows from Assumption (4). Using the c_r -inequality again, $\mathbb{E} \left[\epsilon_{kl}^{2(1+\gamma)} \right]$ is bounded by a constant K as follows:

$$\begin{aligned}
& \mathbb{E} \left[\epsilon_{kl}^{2(1+\gamma)} \right] \\
& \leq 2^{2(1+\gamma)-1} \left(|\Sigma_{kl}^*|^{2(1+\gamma)} + \mathbb{E} \left[|y_{mk}y_{ml}/2|^{2(1+\gamma)} \right] \right) \\
& = 2^{2(1+\gamma)-1} \left(\mathbb{E} \left[|y_{mk}y_{ml}/2|^{2(1+\gamma)} \right] + (\sigma_y/2)^{2(1+\gamma)} \right) \\
& \leq 2^{2(1+\gamma)-1} \left\{ \left(\sigma_y^{2(1+\gamma)} + 1 \right)^{1+\gamma} / 2^{2(1+\gamma)} + (\sigma_y/2)^{2(1+\gamma)} \right\} \\
& = 2^{-1} \left\{ \left(\sigma_y^{2(1+\gamma)} + 1 \right)^{1+\gamma} + \sigma_y^{2(1+\gamma)} \right\} := K.
\end{aligned}$$

The last inequality uses the fact that $\mathbb{E} \left[y_{mj}^4 \right] \leq \mathbb{E} \left[|y_{mj}|^{4(1+\gamma)} \right] + 1$ for $j \in [d]$. By similar reasoning,

$$\mathbb{E} \left[\epsilon_{kl}^2 \right] = \text{Var} \left(y_{mk}y_{ml}/2 \right) \leq \mathbb{E} \left[y_{mk}^2 y_{ml}^2 \right] / 4 \leq (\sigma_y^{2(1+\gamma)} + 1) / 4 < K.$$

D. Proof of Proposition 11

Proof: For fixed $k, l \in [d]$, $\Sigma_{kl}^* - y_{mk}y_{ml}/2$ are identically distributed for $m \in [N]$. Let $D_{kl} = (1/N) \sum_{m=1}^N \mathbb{1}(|\Sigma_{kl}^* - y_{mk}y_{ml}/2| \leq \alpha/2)$. By Chebyshev's inequality,

$$\mathbb{E}[D_{kl}] = \Pr(|\Sigma_{kl}^* - y_{mk}y_{ml}/2| \leq \alpha/2) \geq 1 - 4K/\alpha^2 > 5/6.$$

The last inequality holds with $4K/\alpha^2 < 1/6$.

Let $h_{kl}(\mathbf{x}_i, \mathbf{x}_j) := \mathbb{1}(|\Sigma_{kl}^* - (x_{ik} - x_{jk})(x_{il} - x_{jl})/2| \leq \alpha/2)$. Then it's easy to see that

$$D_{kl} = (1/N) \sum_{1 \leq i < j \leq n} h_{kl}(\mathbf{x}_i, \mathbf{x}_j)$$

Let $q = \lfloor n/2 \rfloor$ and $\sum_{\mathcal{P}}$ denote the summation over all $n!$ permutations (i_1, \dots, i_n) of $[n] := \{1, \dots, n\}$. Using the same technique used in the Proof of Theorem 3.1 in [28], it can be shown that

$$\begin{aligned}
\Pr(D_{kl} \leq 1/2) &\leq \Pr \left(e^{-nD_{kl}} \geq e^{-n/2} \right) \\
&\leq e^{n/2} \frac{1}{n!} \sum_{\mathcal{P}} \prod_{j=1}^q \mathbb{E} e^{-(n/q) \cdot h_{kl}(\mathbf{x}_{i_{2j-1}}, \mathbf{x}_{i_{2j}})}. \tag{24}
\end{aligned}$$

Considering the fact that $e^{-(n/q)t} \leq 1 - (1 - e^{-n/q})t$ for $0 \leq t \leq 1$,

$$\begin{aligned}
& \mathbb{E} e^{-(n/q) \cdot h_{kl}(\mathbf{x}_{i_{2j-1}}, \mathbf{x}_{i_{2j}})} \\
& \leq 1 - (1 - e^{-n/q}) \mathbb{E} \left[h_{kl}(\mathbf{x}_{i_{2j-1}}, \mathbf{x}_{i_{2j}}) \right] \\
& \leq 1 - (5/6) \cdot (1 - e^{-n/q}). \tag{25}
\end{aligned}$$

With (24) and (25),

$$\Pr(D_{kl} \leq 1/2) \leq e^{n/2} \left[1/6 + (5/6) \cdot e^{-n/q} \right]^q. \tag{26}$$

Given $q = \lfloor n/2 \rfloor$, it is easy to verify that $1/6 + (5/6) \cdot e^{-n/q} \leq e^{-(7/12) \cdot (n/q)}$ for $n \geq 11$. Combining this with (26) yields

$$\Pr(D_{kl} \leq 1/2) \leq e^{-n/12}.$$

With union bound we have

$$\Pr \left[\min_{k,l} D_{kl} < \frac{1}{2} \right] \leq d^2 \exp(-n/12).$$

Let $\mathcal{G}_{kl} = \{m \in [N] : |\Sigma_{kl}^* - y_{mk}y_{ml}/2| \leq \alpha/2\}$. Under the event that $\min_{k,l} D_{kl} \geq 1/2$,

$$\begin{aligned}
& \frac{1}{N} \sum_{m=1}^N \left\{ \rho'_\alpha(\Sigma_{kl} - y_{mk}y_{ml}/2) - \rho'_\alpha(\Sigma_{kl}^* - y_{mk}y_{ml}/2) \right\} \cdot (\Sigma_{kl} - \Sigma_{kl}^*) \\
& \geq \frac{1}{N} \sum_{m \in \mathcal{G}_{kl}} \left\{ \rho'_\alpha(\Sigma_{kl} - y_{mk}y_{ml}/2) - \rho'_\alpha(\Sigma_{kl}^* - y_{mk}y_{ml}/2) \right\} \\
& \quad \cdot (\Sigma_{kl} - \Sigma_{kl}^*)
\end{aligned}$$

$$\blacksquare = \frac{|\mathcal{G}_{kl}|}{N} \cdot (\Sigma_{kl} - \Sigma_{kl}^*)^2 \geq \frac{1}{2} (\Sigma_{kl} - \Sigma_{kl}^*)^2,$$

where the last equality uses $\Sigma \in \Sigma^* + \mathbb{B}^\infty(\alpha/2)$ and the last inequality follows from $|\mathcal{G}_{kl}|/N = D_{kl}$. Therefore,

$$\begin{aligned} & \langle \nabla L_\alpha(\Sigma) - \nabla L_\alpha(\Sigma^*), \Sigma - \Sigma^* \rangle \\ &= \sum_{k,l} \frac{1}{N} \sum_{m=1}^N \{ \rho'_\alpha(\Sigma_{kl} - y_{mk}y_{ml}/2) - \\ & \quad \rho'_\alpha(\Sigma_{kl}^* - y_{mk}y_{ml}/2) \} \cdot (\Sigma_{kl} - \Sigma_{kl}^*) \\ & \geq \frac{1}{2} \|\Sigma - \Sigma^*\|_F^2 \end{aligned}$$

with at least $1 - d^2 \exp(-n/12)$ probability. \blacksquare

We adopt the following notations before the next stage of proof. Recall that $L_\alpha(\Sigma) = \sum_{k,\ell} \frac{1}{N} \sum_{m=1}^N \rho_\alpha(\Sigma_{k\ell} - y_{mk}y_{m\ell}/2)$. Define $\mathbf{B}(\Sigma) := \mathbb{E}[\nabla L_\alpha(\Sigma)]$, and $\mathbf{W}^* := \nabla L_\alpha(\Sigma^*) - \mathbb{E}[\nabla L_\alpha(\Sigma)]$.

Lemma 16. *Recall that K is the constant defined in Lemma 10. $|(\mathbf{B}(\Sigma^*))_{kl}| < \frac{K}{\alpha^{1+2\gamma}}$. Let $\alpha = \sqrt{\frac{Kn}{\log d}}$. If we have $n \geq K^{-1}(\log d)^{1+\frac{1}{2\gamma}}$, then $\|\mathbf{B}(\Sigma^*)\|_\infty \leq \sqrt{Kn/n}$ and $\|\mathbf{B}(\Sigma^*)\|_F \leq \sqrt{Ks/n}$.*

Proof: For fixed $k, l \in [d]$, $\Sigma_{k\ell}^* - y_{mk}y_{m\ell}/2$ are identically distributed for $m \in [N]$. Let $\epsilon_{kl} = \Sigma_{k\ell}^* - y_{mk}y_{m\ell}/2$, then

$$\begin{aligned} |\mathbb{E}[\rho'_\alpha(\epsilon_{kl})]| &= |\mathbb{E}[\epsilon_{kl}I(|\epsilon_{kl}| \leq \alpha) + \alpha \text{sgn}(\epsilon_{kl})I(|\epsilon_{kl}| > \alpha)]| \\ &= |\mathbb{E}[\epsilon_{kl} + (\alpha \text{sgn}(\epsilon_{kl}) - \epsilon_{kl})I(|\epsilon_{kl}| > \alpha)]| \\ &= |\mathbb{E}\{[\epsilon_{kl} - \alpha \text{sgn}(\epsilon_{kl})]I(|\epsilon_{kl}| > \alpha)\}| \\ &\leq |\mathbb{E}\{(|\epsilon_{kl}| - \alpha \text{sgn}(\epsilon_{kl}))I(|\epsilon_{kl}| > \alpha)\}| \\ &\leq \frac{|\mathbb{E}\{(\epsilon_{kl}^{2(1+\gamma)} - \alpha^{2(1+\gamma)})I(|\epsilon_{kl}| > \alpha)\}|}{\alpha^{1+2\gamma}} \\ &< \frac{K}{\alpha^{1+2\gamma}}. \end{aligned}$$

Therefore, for all k, l ,

$$|(\mathbf{B}(\Sigma^*))_{kl}| = \frac{1}{N} \left| \sum_{m=1}^N \mathbb{E}[\rho'_\alpha(\Sigma_{k\ell}^* - y_{mk}y_{m\ell}/2)] \right| < \frac{K}{\alpha^{1+2\gamma}}.$$

If we have $\alpha = \sqrt{Kn/\log d}$ and $n \geq K^{-1}(\log d)^{1+\frac{1}{2\gamma}}$, then

$$\|\mathbf{B}(\Sigma^*)\|_\infty \leq \frac{K}{(Kn/\log d)^{1/2+\gamma}} \leq \sqrt{\frac{K}{n}},$$

and

$$\|\mathbf{B}(\Sigma^*)\|_F \leq \sqrt{s} \cdot \|\mathbf{B}(\Sigma^*)\|_\infty \leq \sqrt{\frac{Ks}{n}}. \quad \blacksquare$$

E. Proof of Proposition 12

Proof: For fixed k, l , denote $\Sigma_{k\ell}^* - (x_{ik} - x_{jk})(x_{il} - x_{jl})/2$ as $\epsilon_{i,j}$. Let $h_{kl}(\mathbf{x}_i, \mathbf{x}_j) := \rho'_\alpha(\epsilon_{i,j})$. Then it's easy to see that

$$(\nabla L_\alpha(\Sigma^*))_{kl} = (1/N) \sum_{1 \leq i < j \leq n} h_{kl}(\mathbf{x}_i, \mathbf{x}_j).$$

Let $q = \lfloor n/2 \rfloor$ and $\sum_{\mathcal{P}}$ denote the summation over all $n!$ permutations (i_1, \dots, i_n) of $[n] := \{1, \dots, n\}$. Using the same

technique used in the Proof of Theorem 3.1 in [28], it can be shown that

$$\begin{aligned} \Pr((\nabla L_\alpha(\Sigma^*))_{kl} \geq y) &\leq \Pr\left(e^{(q/\alpha) \cdot (\nabla L_\alpha(\Sigma^*))_{kl}} \geq e^{(q/\alpha) \cdot y}\right) \\ &\leq e^{-(q/\alpha) \cdot y} \frac{1}{n!} \sum_{\mathcal{P}} \prod_{j=1}^q \mathbb{E} e^{(1/\alpha) \cdot h_{kl}(\mathbf{x}_{i_{2j-1}}, \mathbf{x}_{i_{2j}})}. \end{aligned} \quad (27)$$

Let $\rho_1(\cdot)$ denote the Huber loss function defined in (1) with $\alpha = 1$. Note that $h_{kl}(\mathbf{x}_i, \mathbf{x}_j) = \rho'_\alpha(\epsilon_{i,j}) = \alpha \rho'_1(\epsilon_{i,j}/\alpha)$. In addition, it is easy to verify the inequality that

$$-\log(1 - x + x^2) \leq \rho'_1(x) \leq \log(1 + x + x^2)$$

Therefore,

$$\begin{aligned} \mathbb{E} e^{(1/\alpha) \cdot h_{kl}(\mathbf{x}_{i_{2j-1}}, \mathbf{x}_{i_{2j}})} &\leq 1 + \mathbb{E} [\epsilon_{i_{2j-1}, i_{2j}}/\alpha] + \mathbb{E} [(\epsilon_{i_{2j-1}, i_{2j}}/\alpha)^2] \\ &\leq 1 + K/\alpha^2 \leq e^{K/\alpha^2}, \end{aligned} \quad (28)$$

where the second inequality holds because $\mathbb{E} [\epsilon_{i_{2j-1}, i_{2j}}] = 0$ and $\mathbb{E} [\epsilon_{i_{2j-1}, i_{2j}}^2] \leq K$ by Lemma 10. Combining (27) and (28) yields

$$\Pr((\nabla L_\alpha(\Sigma^*))_{kl} \geq y) \leq e^{-(q/\alpha) \cdot y + qK/\alpha^2}.$$

Similarly, it can be shown that $\Pr((\nabla L_\alpha(\Sigma^*))_{kl} \leq -y) \leq e^{-(q/\alpha) \cdot y + qK/\alpha^2}$. Since $\alpha = \sqrt{Kn/\log d}$, by taking $y = 8\sqrt{K \log d/n}$, we conclude that

$$\begin{aligned} \Pr\left(|(\nabla L_\alpha(\Sigma^*))_{kl}| \geq 8\sqrt{K \log d/n}\right) \\ \leq 2e^{-(q/n) \cdot 7 \log d} \leq 2/d^3. \end{aligned}$$

The last inequality follows from $q = \lfloor n/2 \rfloor$. With the union bound, we have

$$\|\nabla L_\alpha(\Sigma^*)\|_\infty \leq 8\sqrt{\frac{K \log d}{n}}$$

with at least $1 - 2/d$ probability.

Now we prove the second statement. We have

$$\mathbb{E} \|\mathbf{W}_S^*\|_F = \mathbb{E} \sqrt{\sum_{(k,l) \in \mathcal{S}} (W_{kl}^*)^2} \leq \sqrt{\sum_{(k,l) \in \mathcal{S}} \mathbb{E} (W_{kl}^*)^2}. \quad (29)$$

The inequality holds because $f(t_1, t_2, \dots, t_s) = \sqrt{t_1 + t_2 + \dots + t_s}$ is concave. Recall that for fixed k, l , denote $\Sigma_{k\ell}^* - (x_{ik} - x_{jk})(x_{il} - x_{jl})/2$ as $\epsilon_{i,j}$. Note that with Lemma 10, $\mathbb{E} [\epsilon_{i,j}^2] \leq K$. Recall that

$$W_{kl}^* = \frac{1}{N} \sum_{1 \leq i < j \leq n} \{ \rho'_\alpha(\epsilon_{i,j}) - \mathbb{E} \rho'_\alpha(\epsilon_{i,j}) \}.$$

Then

$$\begin{aligned} N^2 \mathbb{E} (W_{kl}^*)^2 \\ &= \sum_{(i,j,i',j') \in \mathcal{G}} \mathbb{E} \{ \{ \rho'_\alpha(\epsilon_{i,j}) - \mathbb{E} \rho'_\alpha(\epsilon_{i,j}) \} \\ & \quad \cdot \{ \rho'_\alpha(\epsilon_{i',j'}) - \mathbb{E} \rho'_\alpha(\epsilon_{i',j'}) \} \}, \end{aligned}$$

where $\mathcal{G} = \{(i, j, i', j') : 1 \leq i < j \leq n, 1 \leq i' < j' \leq n\}$ denote the set of pairs of indices. Let $\mathcal{G}_d = \{(i, j, i', j') \in \mathcal{G} : i, j, i', j' \text{ are not equal to each other}\}$. Notice that $(i, j, i', j') \in \mathcal{G}_d$ indicates that $\epsilon_{i,j}$ and $\epsilon_{i',j'}$ are

independent because $\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_{i'}, \mathbf{x}_{j'}$ are independent samples, and consequently,

$$\sum_{(i,j,i',j') \in \mathcal{G}_d} \mathbb{E} \{ \{ \rho'_\alpha(\epsilon_{i,j}) - \mathbb{E} \rho'_\alpha(\epsilon_{i,j}) \} \cdot \{ \rho'_\alpha(\epsilon_{i',j'}) - \mathbb{E} \rho'_\alpha(\epsilon_{i',j'}) \} \} = 0.$$

Therefore,

$$\begin{aligned} & N^2 \mathbb{E} (W_{kl}^*)^2 \\ &= \sum_{(i,j,i',j') \in \overline{\mathcal{G}_d}} \mathbb{E} \{ \{ \rho'_\alpha(\epsilon_{i,j}) - \mathbb{E} \rho'_\alpha(\epsilon_{i,j}) \} \cdot \{ \rho'_\alpha(\epsilon_{i',j'}) - \mathbb{E} \rho'_\alpha(\epsilon_{i',j'}) \} \}. \end{aligned} \quad (30)$$

A more basic bound yields

$$\begin{aligned} & \mathbb{E} \{ \{ \rho'_\alpha(\epsilon_{i,j}) - \mathbb{E} \rho'_\alpha(\epsilon_{i,j}) \} \cdot \{ \rho'_\alpha(\epsilon_{i',j'}) - \mathbb{E} \rho'_\alpha(\epsilon_{i',j'}) \} \} \\ & \leq [\text{Var} \{ \rho'_\alpha(\epsilon_{i,j}) \} \cdot \text{Var} \{ \rho'_\alpha(\epsilon_{i',j'}) \}]^{1/2} \\ & \leq [\text{Var} \{ \epsilon_{i,j} \} \cdot \text{Var} \{ \epsilon_{i',j'} \}]^{1/2} \leq K \end{aligned} \quad (31)$$

By definition of \mathcal{G}_d , it's easy to verify that $|\overline{\mathcal{G}_d}| = (2n-3)N$. Combining this with (29), (30) and (31), we have

$$\mathbb{E} \|\mathbf{W}_S^*\|_F \leq \sqrt{\frac{(2n-3)Ks}{N}} < \sqrt{\frac{4Ks}{n}} \quad (32)$$

With Markov's inequality, we have

$$\Pr \left\{ \|\mathbf{W}_S^*\|_F \geq \beta \sqrt{\frac{Ks}{n}} \right\} \leq \frac{\mathbb{E} \|\mathbf{W}_S^*\|_F}{\beta \sqrt{Ks/n}} < \frac{2}{\beta}$$

Recall that $\nabla L_\alpha(\Sigma^*) = \mathbf{B}(\Sigma^*) + \mathbf{W}^*$. With Lemma 16, we have $\|\mathbf{B}(\Sigma^*)\|_F \leq \sqrt{Ks/n}$. Combing the this with (32), we have

$$\Pr \left\{ \|\nabla L_\alpha(\Sigma^*)\|_F \geq (\beta+1) \sqrt{\frac{Ks}{n}} \right\} \leq \frac{2}{\beta}. \quad \blacksquare$$

F. Proof of Theorem 13

Proof: Choose $c > 0$ so that

$$0.5p'(\phi_0) (c^2 + 1)^{1/2} + 2 = 0.5c\phi_0$$

Set $l = \left(2 + \frac{2}{p'(\phi_0)}\right) (c^2 + 1)^{1/2} s^{1/2} + \frac{2}{p'(\phi_0)} s^{1/2}$. Then it's easy to verify that given $r = \alpha/2$,

$$c\phi_0 \lambda s^{1/2} \leq r \quad \text{and} \quad 5\lambda s^{1/2} \leq r.$$

Apply Proposition 8 to conclude that, for any positive constant β , conditioned on event $\mathcal{E}_1(r, 1/2) \cap \{ \|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda \} \cap \{ \|\nabla L_\alpha(\Sigma^*)\|_F \leq (\beta+1)\sqrt{Ks/n} \}$, we have $\tilde{\Sigma}^{(t)} \in \Sigma^* + \mathbb{C}(l)$ and

$$\begin{aligned} \left\| \tilde{\Sigma}^{(t)} - \Sigma^* \right\|_F & \leq \delta \left\| \tilde{\Sigma}^{(t-1)} - \Sigma^* \right\|_F + 2 \{ \|p'_\lambda(|\Sigma_S^*| - \phi_0\lambda)\|_F \\ & \quad + s^{1/2}\epsilon + \|\nabla L_\alpha(\Sigma^*)\|_F + \tau \|(\Sigma^*)^{-1}\|_F \} \end{aligned} \quad (33)$$

By $\|\Sigma_S^*\|_{\min} \geq (\phi_0 + \phi_1)\lambda$, we have $\|p'_\lambda(|\Sigma_S^*| - \phi_0\lambda)\|_F = 0$, so

$$\begin{aligned} \left\| \tilde{\Sigma}^{(t)} - \Sigma^* \right\|_F & \leq \delta \left\| \tilde{\Sigma}^{(t-1)} - \Sigma^* \right\|_F \\ & \quad + 2 \{ s^{1/2}\epsilon + \|\nabla L_\alpha(\Sigma^*)\|_F + \tau \|(\Sigma^*)^{-1}\|_F \} \end{aligned}$$

Therefore by Remark 9,

$$\begin{aligned} \left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F & \leq \delta^{T-1} \left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_F \\ & \quad + \frac{2}{1-\delta} \{ s^{1/2}\epsilon + \|\nabla L_\alpha(\Sigma^*)\|_F + \tau \|(\Sigma^*)^{-1}\|_F \} \\ & \leq \sqrt{s/n} + \frac{2}{1-\delta} (\beta+3) \sqrt{Ks/n} \leq \beta \sqrt{s/n}. \end{aligned}$$

The second inequality uses the bound of $\left\| \tilde{\Sigma}^{(1)} - \Sigma^* \right\|_F$ given in Proposition 6. Provided that $\tilde{\Sigma}^{(T)} - \Sigma^* \in \mathbb{C}(l)$ and $l \asymp s^{1/2}$,

$$\left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_1 \leq s^{1/2} \cdot \beta \sqrt{\frac{s}{n}}.$$

Take $n \geq K^{-1}(\log d)^{1+\frac{1}{2\gamma}}$. By Proposition 11, Proposition 12 and the union bound, event $\mathcal{E}_1(r, 1/2) \cap \{ \|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda \} \cap \{ \|\nabla L_\alpha(\Sigma^*)\|_F \leq (\beta+1)\sqrt{Ks/n} \}$ happens with at least $1 - 2/d - d^2 \exp(-n/12) - 2/\beta$ probability.

To prove the weaker conclusion, condition only on the event $\{ \|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5p'(\phi_0)\lambda \}$, which holds with at least $1 - 2/d$ probability. Compute:

$$\begin{aligned} & \lim_{M \rightarrow \infty} \limsup_n \Pr \left\{ \left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F \geq M \sqrt{s/n} \right\} \\ &= \lim_{\beta \rightarrow \infty} \limsup_n \Pr \left\{ \left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F \geq \right. \\ & \quad \left. \sqrt{s/n} + \frac{2}{1-\delta} (\beta+3) \sqrt{s/n} \right\} \\ & \leq \lim_{\beta \rightarrow \infty} \limsup_n d^2 \exp(-n/12) + 2/\beta \\ &= \lim_{\beta \rightarrow \infty} 2/\beta = 0. \end{aligned}$$

That is,

$$\left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_F = O_P(\sqrt{s/n}).$$

Similarly, $\left\| \tilde{\Sigma}^{(T)} - \Sigma^* \right\|_1 = O_P(s/\sqrt{n})$. \blacksquare

REFERENCES

- [1] A. P. Dempster, "Covariance selection," *Biometrics*, vol. 28, no. 1, pp. 157–175, 1972.
- [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] O. Ledoit and M. Wolf, "Improved estimation of the covariance matrix of stock returns with an application to portfolio selection," *J. Empir. Finance.*, vol. 10, no. 5, pp. 603–621, 2003.
- [4] Z. Zhao and D. P. Palomar, "Mean-reverting portfolio with budget constraint," *IEEE Trans. Signal Process.*, vol. 66, no. 9, pp. 2342–2357, 2018.
- [5] 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.

- [6] Y. Chen, A. Wiesel, Y. C. Eldar, and A. O. Hero, "Shrinkage algorithms for MMSE covariance estimation," *IEEE Trans. Signal Process.*, vol. 58, no. 10, pp. 5016–5029, 2010.
- [7] E. Ollila, D. E. Tyler, V. Koivunen, and H. V. Poor, "Complex elliptically symmetric distributions: Survey, new results and applications," *IEEE Trans. Signal Process.*, vol. 60, no. 11, pp. 5597–5625, 2012.
- [8] Y. Sun, P. Babu, and D. P. Palomar, "Robust estimation of structured covariance matrix for heavy-tailed elliptical distributions," *IEEE Trans. Signal Process.*, vol. 64, no. 14, pp. 3576–3590, 2016.
- [9] A. Aubry, A. De Maio, and L. Pallotta, "A geometric approach to covariance matrix estimation and its applications to radar problems," *IEEE Trans. Signal Process.*, vol. 66, no. 4, pp. 907–922, 2018.
- [10] E. Ollila, D. P. Palomar, and F. Pascal, "Shrinking the eigenvalues of M-estimators of covariance matrix," *IEEE Trans. Signal Process.*, vol. 69, pp. 256–269, 2021.
- [11] C. M. Bishop, *Pattern Recognition and Machine Learning*. Springer, 2006.
- [12] I. T. Jolliffe, *Principal Component Analysis*. Springer, 2002.
- [13] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [14] H. Markowitz, "Portfolio selection," *J. Finance.*, vol. 7, no. 1, pp. 77–91, 1952.
- [15] R. A. Monzingo and T. W. Miller, *Introduction to Adaptive Arrays*. SciTech Publishing, 2004.
- [16] W. B. Wu and M. Pourahmadi, "Nonparametric estimation of large covariance matrices of longitudinal data," *Biometrika*, vol. 90, no. 4, pp. 831–844, 2003.
- [17] P. J. Bickel and E. Levina, "Covariance regularization by thresholding," *Ann. Statist.*, vol. 36, no. 6, pp. 2577–2604, 2008.
- [18] N. El Karoui, "Operator norm consistent estimation of large-dimensional sparse covariance matrices," *Ann. Statist.*, vol. 36, no. 6, pp. 2717–2756, 2008.
- [19] 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.
- [20] T. T. Cai and H. H. Zhou, "Minimax estimation of large covariance matrices under ℓ_1 -norm," *Stat. Sin.*, vol. 22, pp. 1319–1378, 2012.
- [21] —, "Optimal rates of convergence for sparse covariance matrix estimation," *Ann. Statist.*, vol. 40, no. 5, pp. 2389–2420, 2012.
- [22] P. J. Bickel and E. Levina, "Regularized estimation of large covariance matrices," *Ann. Statist.*, vol. 36, no. 1, pp. 199–227, 2008.
- [23] T. T. Cai, C.-H. Zhang, and H. H. Zhou, "Optimal rates of convergence for covariance matrix estimation," *Ann. Statist.*, vol. 38, no. 4, pp. 2118–2144, 2010.
- [24] X. Chen, Z. J. Wang, and M. J. McKeown, "Shrinkage-to-tapering estimation of large covariance matrices," *IEEE Trans. Signal Process.*, vol. 60, no. 11, pp. 5640–5656, 2012.
- [25] E. Ollila and A. Breloy, "Regularized tapered sample covariance matrix," *IEEE Trans. Signal Process.*, vol. 70, pp. 2306–2320, 2022.
- [26] Q. Wei and Z. Zhao, "Large covariance matrix estimation with oracle statistical rate," in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process.* IEEE, 2023.
- [27] 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.
- [28] Y. Ke, S. Minsker, Z. Ren, Q. Sun, and W.-X. Zhou, "User-friendly covariance estimation for heavy-tailed distributions," *Statistical Science*, vol. 34, no. 3, pp. 454–471, 2019.
- [29] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *Am. Stat.*, vol. 58, no. 1, pp. 30–37, 2004.
- [30] 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, 2017.
- [31] E. J. Candes, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted ℓ_1 minimization," *J. Fourier. Anal. Appl.*, vol. 14, no. 5, pp. 877–905, 2008.
- [32] H. A. Le Thi, T. P. Dinh, H. M. Le, and X. T. Vo, "DC approximation approaches for sparse optimization," *Eur. J. Oper. Res.*, vol. 244, no. 1, pp. 26–46, 2015.