

Robust Sparse Large Covariance Matrix Estimation Under Huber Loss

XXX

Abstract—XXX

I. INTRODUCTION

[Covariance estimation]

The estimation of covariance matrices is a fundamental problem in modern multivariate data analysis. It has broad applications in many fields such as statistics, biology, finance, signal processing, machine learning, etc. For example, many dimension reduction techniques, including principal component analysis [1] and linear and quadratic discriminant analysis [2], require the estimation of a covariance matrix in advance from a given collection of data points. Theoretical properties of large covariance estimators discussed in the literature often hinge heavily on the Gaussian or sub-Gaussian assumption. Given that data from fields including genomic studies and quantitative finance usually do not follow the assumed Gaussian or sub-Gaussian shape, such an assumption is typically very restrictive in practice. It is therefore imperative to develop robust inferential procedures that are less sensitive to the distributional assumptions.

[why robustness: outliers]

Heavy-tailed distribution¹ is a viable model for data contaminated by outliers that are typically encountered in applications. Due to heavy-tailedness, the probability that some observations are sampled far away from the “true” parameter of the population is nonnegligible. We refer to these outlying data points as stochastic outliers. A procedure that is robust against such outliers, evidenced by its better finite-sample performance than a nonrobust method, is called a tail-robust procedure. The tail robustness is different from the classical notion of robustness that is often characterized by the breakdown point. Nevertheless, as stated in [3], it does not provide any information on the convergence properties of an estimator, such as consistency and efficiency. Tail-robustness is a concept that combines robustness, consistency, and finite-sample error bounds. We will denote tail-robustness as robustness in the following sections.

[why sparsity: high-dimensional]

[Robust and sparse covariance: existing meth²]

[Robust scatter matrix estimation based on M-estimation, but they cannot exactly obtain a covariance matrix (in most cases, up to a scalar)]

1. Cannot exactly obtain a covariance matrix and in most cases, can only obtain a covariance matrix up to a scalar under ellipsoid assumption for the distribution.

¹The distribution of a random variable X is said to be heavy-tailed if the moment generating function of X is infinite for all $t > 0$, that is, $\mathbb{E}e^{tX} = \infty$ for all $t > 0$. Hence we can only assume weaker moment conditions (for instance, $\mathbb{E}X^4 < \infty$) for heavy-tailed distributions.

2. The method can be adapted for sparse covariance estimation where all nonzero covariates are contained in a submatrix, but not applicable to covariance matrices with general sparsity pattern.

[in this paper, we XXX]

Notations. XXX

II. PROBLEM FORMULATION

Given zero-mean samples \mathbf{x}_i , $i = 1, \dots, n$ from a heavy-tailed distribution, define

$$L_\alpha(\Sigma) := \sum_{k,\ell} \frac{1}{n} \sum_{i=1}^n \rho_\alpha(\Sigma_{k\ell} - x_{ik}x_{i\ell})$$

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

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

Further, define

$$\hat{\Sigma} \in \arg \min_{\Sigma} \left\{ L_\alpha(\Sigma) + \lambda \|\Sigma\|_{1,\text{off}} \right\} \quad (1)$$

In this paper, we want to show $\hat{\Sigma}$ achieves the minimax optimal statistical rate for robust sparse covariance estimation.

III. THEORETICAL RESULTS

We denote the underlying true covariance matrix by Σ^* . Let $S = \{(i, j) \mid \Sigma_{ij}^* \neq 0\}$ be the support set of Σ^* and s be its cardinality, i.e., $s = |S|$. 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 1. $\mathbf{x}_i \in \mathbb{R}^d$ is a heavy-tailed random variable with zero mean, i.e. $\mathbb{E}[x_{ij}] = 0$ and $\mathbb{E}[|x_{ij}|^4] \leq \sigma^2$ for all $1 \leq j \leq d$ with some positive σ .

Remark 2. Assumption 1 immediately implies that there exists constant $K > 0$, such that $\mathbb{E}[(\Sigma_{kl}^* - x_{ik}x_{il})^2] \leq K$ for all $k, l \in [d]$.

Lemma 3. Assume $\|\nabla L_\alpha(\hat{\Sigma})\|_\infty < \beta$ holds with some $\beta = O(1)$. Take $\alpha = \sqrt{Kn/\log d}$. If the sample size satisfies $n \gtrsim \log d$, then

$$\|\hat{\Sigma} - \Sigma^*\|_\infty \lesssim \sqrt{\log d/n} + \beta \quad (2)$$

holds with high probability.

Proof: For fixed k, l , let $\hat{\theta} := (\hat{\Sigma})_{kl}$ and define

$$\Psi(\theta) := \frac{1}{n} \sum_{i=1}^n \rho'_\alpha(\theta - x_{ik}x_{il}), \quad \theta \in \mathbb{R}.$$

Note that $|\Psi(\hat{\theta})| = |(\nabla L_\alpha(\hat{\Sigma}))_{kl}| < \beta$ always hold. In addition, it is easy to verify the inequality that

$$-\log(1 - x + x^2) \leq \rho'_1(x) \leq \log(1 + x + x^2) \quad (3)$$

By (3) and the fact that $\alpha^{-1}\rho'_\alpha(t) = \rho'_1(t/\alpha)$,

$$\begin{aligned} \mathbb{E}e^{(n/\alpha) \cdot \Psi(\theta)} &= \prod_{i=1}^n \mathbb{E}e^{\rho'_1((\theta - x_{ik}x_{il})/\alpha)} \\ &\leq \prod_{i=1}^n \mathbb{E} \left\{ 1 + \alpha^{-1}(\theta - x_{ik}x_{il}) + \alpha^{-2}(\theta - x_{ik}x_{il})^2 \right\} \\ &\leq \prod_{i=1}^n \left[1 + \alpha^{-1}(\theta - \Sigma_{kl}^*) + \alpha^{-2} \left\{ (\theta - \Sigma_{kl}^*)^2 + K \right\} \right] \\ &\leq \exp \left[n\alpha^{-1}(\theta - \Sigma_{kl}^*) + n\alpha^{-2} \left\{ (\theta - \Sigma_{kl}^*)^2 + K \right\} \right]. \end{aligned} \quad (4)$$

Similarly, it can be shown that

$$\begin{aligned} \mathbb{E}e^{-(n/\alpha) \cdot \Psi(\theta)} &\leq \exp \left[-n\alpha^{-1}(\theta - \Sigma_{kl}^*) + n\alpha^{-2} \left\{ (\theta - \Sigma_{kl}^*)^2 + K \right\} \right]. \end{aligned} \quad (5)$$

For $\eta \in (0, 1)$, define

$$\begin{aligned} B_-(\theta) &= (\theta - \Sigma_{kl}^*) + \left\{ (\theta - \Sigma_{kl}^*)^2 + K \right\} / \alpha - (\alpha/n) \log \eta \\ B_+(\theta) &= -(\theta - \Sigma_{kl}^*) + \left\{ (\theta - \Sigma_{kl}^*)^2 + K \right\} / \alpha + (\alpha/n) \log \eta \end{aligned}$$

Together, (4), (5) and Markov's inequality imply

$$\begin{aligned} \Pr(\Psi(\theta) > B_-(\theta)) &\leq e^{-nB_-(\theta)/\alpha} \cdot \mathbb{E}e^{(n/\alpha) \cdot \Psi(\theta)} \leq \eta, \\ \text{and } \Pr(\Psi(\theta) < B_+(\theta)) &\leq e^{-nB_+(\theta)/\alpha} \cdot \mathbb{E}e^{-(n/\alpha) \cdot \Psi(\theta)} \leq \eta. \end{aligned}$$

Let θ_+ be the smallest solution of the quadratic equation $B_+(\theta_+) = \beta$, and θ_- be the largest solution of the quadratic equation $B_-(\theta_-) = -\beta$. We need to check that θ_- and θ_+ are well-defined. Let Δ_- and Δ_+ denote the discriminant of $B_-(\theta) = -\beta$ and $B_+(\theta) = \beta$, respectively. Since $\alpha = \sqrt{Kn/\log d}$, $\beta = O(1)$ and by taking $n \gtrsim \log d$, $\eta = 1/d^3$, we have

$$\Delta_- = 1 - (4/\alpha) \cdot (K/\alpha - (\alpha/n) \cdot \log \eta + \beta) > 0,$$

which implies that θ_- is well-defined as a solution to $B_-(\theta) = -\beta$ on $(\Sigma_{kl}^* - \alpha/2, \Sigma_{kl}^*)$. Similarly, θ_+ is also well-defined. Then, with at least $1 - 2\eta$ probability,

$$\Psi(\theta_+) \geq B_+(\theta_+) = \beta \quad \text{and} \quad \Psi(\theta_-) \leq B_-(\theta_-) = -\beta.$$

Recall that $|\Psi(\hat{\theta})| < \beta$ always hold, and given that $\Psi(\theta)$ is nondecreasing, $\Psi(\theta_-) < \Psi(\hat{\theta}) < \Psi(\theta_+)$ immediately implies $\theta_- \leq \hat{\theta} \leq \theta_+$.

Now we estimate θ_- and θ_+ . Notice that by convexity, the following holds for all $\theta \in (\Sigma_{kl}^* - \alpha/2, \Sigma_{kl}^*)$:

$$B_-(\theta) \leq (1/2) \cdot (\theta - \Sigma_{kl}^*) + B_-(\Sigma_{kl}^*),$$

which immediately implies that

$$\theta_- - \Sigma_{kl}^* \geq -2(K/\alpha - (\alpha/n) \log \eta + \beta).$$

It can be seen that assuming $B_+(\theta_+) - \beta = K/\alpha + (\alpha/n) \log \eta - \beta > 0$, we have $\theta_+ \in (\Sigma_{kl}^*, \Sigma_{kl}^* + \alpha/2)$, and similarly

$$\theta_+ - \Sigma_{kl}^* \leq 2(K/\alpha + (\alpha/n) \log \eta - \beta). \quad (6)$$

Otherwise if $B_+(\theta_+) - \beta \leq 0$, then $\theta_+ \leq 0$. Combining this with (6), we have

$$\theta_+ - \Sigma_{kl}^* \leq \max \{2(K/\alpha + (\alpha/n) \log \eta - \beta), 0\}.$$

Therefore, with $\theta_- \leq \hat{\theta} \leq \theta_+$,

$$|\hat{\theta} - \Sigma_{kl}^*| \leq 2(K/\alpha - (\alpha/n) \log \eta + \beta).$$

With $\eta = 1/d^3$ and the union bound, we have that with at least $1 - 2/d$ probability, $\|\hat{\Sigma} - \Sigma^*\|_\infty \lesssim \sqrt{\log d/n} + \beta$. ■

Proposition 4. Let $\tilde{\Sigma}$ denote an ϵ -optimal solution to (1). Then, $\tilde{\Sigma} \in \Sigma + \mathbb{C}(l)$, where $l = 4s^{1/2}$. Further, assume $\|\tilde{\Sigma} - \Sigma^*\|_\infty \leq \alpha/2$. Conditioned on the event $\mathcal{E}_1(\alpha/2, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$,

$$\|\tilde{\Sigma} - \Sigma^*\|_F \leq 3\lambda s^{1/2} \quad \text{and} \quad \|\tilde{\Sigma} - \Sigma^*\|_l \leq 12\lambda s.$$

Proposition 4 gives the deterministic interpretation of Theorem 7. In the following propositions we will analyze the probability of the conditioned event $\mathcal{E}_1(\alpha/2, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$ mentioned in Proposition 4.

Proposition 5. Suppose that Assumption 1 holds. Recall that K is the constant defined in Remark 2. Assume $n \gtrsim \log d$. Then, for any $\kappa \in (0, 1)$ and $C > 0$,

$\langle \nabla L_\alpha(\Sigma) - \nabla L_\alpha(\Sigma^*), \Sigma - \Sigma^* \rangle \geq \min \{\kappa, \kappa/2C\} \|\Sigma - \Sigma^*\|_F^2$ holds uniformly for all $\Sigma \in \Sigma^* + \mathbb{B}^\infty(C\alpha)$ with high probability.

Proof: Let $D_{kl} = (1/n) \sum_{i=1}^n 1(|\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2)$. By Chebyshev's inequality,

$$\mathbb{E}[D_{kl}] = \Pr(|\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2) \geq 1 - 4K/\alpha^2 > (1 + \kappa)/2.$$

The last inequality holds because $4K/\alpha^2 < (1 - \kappa)/2$, which follows from $n \gtrsim \log d$.

For each fixed $k, l \in [d]$, let $X_i = 1(|\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2)$. To invoke Bernstein's inequality, compute

$$\begin{aligned} \text{Var}[X_i] &= \Pr(|\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2) \\ &\quad \cdot (1 - \Pr(|\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2)) \\ &\leq 1/4 \end{aligned}$$

and with $|X_i - \mathbb{E}[X_i]| \leq 1$,

$$\mathbb{E}|X_i - \mathbb{E}[X_i]|^l \leq \mathbb{E}|X_i - \mathbb{E}[X_i]|^2 \cdot 1 \leq 1/4.$$

Therefore, with Bernstein's inequality

$$\begin{aligned} \Pr \left(\left| \sum_{i=1}^n \{X_i - \mathbb{E}[X_i]\} \right| \geq (1 - \kappa)n/2 \right) \\ \leq 2 \cdot \exp \left(-\frac{(1 - \kappa)^2 n^2 / 8}{n/4 + (1 - \kappa)n/2} \right) = 2 \cdot \exp \left(-\frac{(1 - \kappa)^2 n}{6 - 4\kappa} \right) \end{aligned}$$

and

$$\begin{aligned}
& \Pr \{D_{kl} < \kappa\} \\
& \leq \Pr \{|D_{kl} - \mathbb{E}[D_{kl}]\} \geq (1 - \kappa)/2\} \\
& = \Pr \left\{ \left| (1/n) \sum_{i=1}^n \{X_i - \mathbb{E}[X_i]\} \right| \geq (1 - \kappa)/2 \right\} \\
& \leq 2 \cdot \exp \left(-\frac{(1 - \kappa)^2 n}{6 - 4\kappa} \right).
\end{aligned}$$

With union bound we have

$$\Pr \left[\min_{k,l} D_{kl} < \kappa \right] \leq 2d^2 \cdot \exp \left(-\frac{(1 - \kappa)^2 n}{6 - 4\kappa} \right) < 1/d,$$

where the last inequality follows from $n \gtrsim \log d$. Let $\mathcal{G}_{kl} := \{i \in [n] : |\Sigma_{kl}^* - x_{ik}x_{il}| \leq \alpha/2\}$. Under the event that $\min_{k,l} D_{kl} \geq \kappa$,

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \{\rho'_\alpha(\Sigma_{kl} - x_{ik}x_{il}) - \rho'_\alpha(\Sigma_{kl}^* - x_{ik}x_{il})\} \cdot (\Sigma_{kl} - \Sigma_{kl}^*) \\
& \geq \frac{1}{n} \sum_{i \in \mathcal{G}_{kl}} \{\rho'_\alpha(\Sigma_{kl} - x_{ik}x_{il}) - \rho'_\alpha(\Sigma_{kl}^* - x_{ik}x_{il})\} \cdot (\Sigma_{kl} - \Sigma_{kl}^*) \\
& \geq \frac{1}{n} \sum_{i \in \mathcal{G}_{kl}} \min\{|\Sigma_{kl} - \Sigma_{kl}^*|, \alpha/2\} \cdot |\Sigma_{kl} - \Sigma_{kl}^*| \\
& \geq \frac{1}{n} \sum_{i \in \mathcal{G}_{kl}} \min\{1, 1/2C\} (\Sigma_{kl} - \Sigma_{kl}^*)^2 \\
& \geq \kappa \min\{1, 1/2C\} (\Sigma_{kl} - \Sigma_{kl}^*)^2
\end{aligned}$$

The second last inequality holds since $\Sigma \in \Sigma^* + \mathbb{B}^\infty(C\alpha)$ implies $\alpha/2 \geq |\Sigma_{kl} - \Sigma_{kl}^*|/2C$, and the last inequality follows from $|\mathcal{G}_{kl}|/n = D_{kl}$. Therefore

$$\begin{aligned}
& \langle \nabla L_\alpha(\Sigma) - \nabla L_\alpha(\Sigma^*), \Sigma - \Sigma^* \rangle \\
& = \sum_{k,l} \frac{1}{n} \sum_{i=1}^n \{\rho'_\alpha(\Sigma_{kl} - x_{ik}x_{il}) - \rho'_\alpha(\Sigma_{kl}^* - x_{ik}x_{il})\} \cdot (\Sigma_{kl} - \Sigma_{kl}^*) \\
& \geq \kappa \cdot \min\{1, 1/2C\} \cdot \|\Sigma - \Sigma^*\|_F^2
\end{aligned}$$

with at least $1 - 1/d$ probability. \blacksquare

Proposition 5 implies that for any $\kappa \in (0, 1)$ and $C > 0$, with $n \gtrsim \log d$, event $\mathcal{E}_1(C, \min\{\kappa, \kappa/2C\})$ happens with high probability.

Proposition 6. Suppose that Assumption 1 holds. Let K be the constant defined in Remark 2. Assume $\alpha = \sqrt{Kn/\log d}$, then

$$\|\nabla L_\alpha(\Sigma^*)\|_\infty \leq 8\sqrt{\frac{K \log d}{n}} \quad (7)$$

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

In Proposition 6, (7) indicates that $\{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$ happens with high probability if we take $\lambda \asymp \sqrt{\log d/n}$ and $\epsilon \lesssim \sqrt{\log d/n}$.

Theorem 7. (minimax-optimal rate) Suppose that Assumption 1 holds. Take $\lambda \asymp \sqrt{\log d/n}$ and let $\alpha = \sqrt{Kn/\log d}$, $\epsilon \lesssim \sqrt{\log d/n}$. If the sample size satisfies $n \gtrsim \log d$, then

$$\|\hat{\Sigma} - \Sigma^*\|_F \lesssim \sqrt{\frac{s \log d}{n}} \quad \text{and} \quad \|\hat{\Sigma} - \Sigma^*\|_1 \lesssim s \sqrt{\frac{\log d}{n}}$$

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

Proof: The proof combines Proposition 4 with Lemma 3, Proposition 5 and Proposition 6. To invoke Proposition 4, we first notice that given $\|\nabla L_\alpha(\hat{\Sigma}) + \lambda \Xi\|_\infty \leq \epsilon$ for some $\Xi \in \partial \|\hat{\Sigma}\|_{1,\text{off}}$, we must have $\|\nabla L_\alpha(\hat{\Sigma})\|_\infty < 2\lambda + \epsilon$ always hold. Lemma 3 indicates that

$$\|\hat{\Sigma} - \Sigma^*\|_\infty \lesssim \sqrt{\log d/n} + 2\lambda + \epsilon \lesssim \sqrt{\log d/n} \leq \alpha/2$$

where the last inequality hold with $n \gtrsim \log d$.

By Proposition 6, $\|\nabla L_\alpha(\Sigma^*)\|_\infty \leq 8\sqrt{K \log d/n}$. With $\epsilon \lesssim \sqrt{\log d/n}$ and $\lambda \asymp \sqrt{\log d/n}$, event $\{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$ happens with at least $1 - 2/d$ probability. Still, with $n \gtrsim \log d$, Proposition 5 indicates that $\mathcal{E}_1(\alpha/2, 1/2)$ happens with high probability. With union bound, event $\mathcal{E}_1(\alpha/2, 1/2) \cap \{\|\nabla L_\alpha(\Sigma^*)\|_\infty + \epsilon \leq 0.5\lambda\}$ holds with high probability. Under this event and by Proposition 4,

$$\|\hat{\Sigma} - \Sigma^*\|_F \leq 3\lambda s^{1/2} \quad \text{and} \quad \|\hat{\Sigma} - \Sigma^*\|_1 \leq 12\lambda s. \quad \blacksquare$$

IV. NUMERICAL SIMULATION

XX

V. CONCLUSION

XX

APPENDIX

Lemma 8. For any $\Sigma \in \mathbb{R}^{d \times d}$ satisfying $\Sigma_{\bar{S}} = \mathbf{0}$ and $\epsilon > 0$, provided $\lambda > \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty + \epsilon$, any ϵ -optimal solution $\tilde{\Sigma}$ to (1) satisfies

$$\begin{aligned}
& \|(\tilde{\Sigma} - \Sigma)_{\bar{S}}\|_1 \\
& \leq (\lambda - \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty - \epsilon)^{-1} \\
& \quad \cdot (\lambda + \|\nabla L_\alpha(\Sigma)_S\|_\infty + \epsilon) \cdot \|(\tilde{\Sigma} - \Sigma)_S\|_1
\end{aligned}$$

Proof: For any $\Xi \in \partial \|\tilde{\Sigma}\|_{1,\text{off}}$, define $U(\Xi) = \nabla L_\alpha(\tilde{\Sigma}) + \lambda \Xi \in \mathbb{R}^{d \times d}$. By convexity of $L_\alpha(\Sigma)$:

$$\langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \geq 0.$$

Therefore,

$$\begin{aligned}
& \|U(\Xi)\|_\infty \|\tilde{\Sigma} - \Sigma\|_1 \geq \langle U(\Xi), \tilde{\Sigma} - \Sigma \rangle \\
& = \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle + \langle \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\
& \quad + \langle \lambda \Xi, \tilde{\Sigma} - \Sigma \rangle \\
& \geq 0 - \|\nabla L_\alpha(\Sigma)_S\|_\infty \|(\tilde{\Sigma} - \Sigma)_S\|_1 \\
& \quad - \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty \|(\tilde{\Sigma} - \Sigma)_{\bar{S}}\|_1 + \langle \lambda \Xi, \tilde{\Sigma} - \Sigma \rangle
\end{aligned}$$

Moreover, we have

$$\begin{aligned} & \langle \lambda \Xi, \tilde{\Sigma} - \Sigma \rangle \\ &= \lambda \langle \Xi_{\bar{S}}, (\tilde{\Sigma} - \Sigma)_{\bar{S}} \rangle + \lambda \langle \Xi_S, (\tilde{\Sigma} - \Sigma)_S \rangle \\ &\geq \lambda \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 - \lambda \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1 \end{aligned}$$

Together, the last two displays imply

$$\begin{aligned} & \|U(\Xi)\|_\infty \left\| \tilde{\Sigma} - \Sigma \right\|_1 \\ &\geq -\|\nabla L_\alpha(\Sigma)_S\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1 - \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \\ &\quad + \lambda \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 - \lambda \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1 \end{aligned}$$

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

$$\begin{aligned} & \epsilon \left\| \tilde{\Sigma} - \Sigma \right\|_1 \\ &\geq -\|\nabla L_\alpha(\Sigma)_S\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1 - \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \\ &\quad + \lambda \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 - \lambda \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1 \end{aligned}$$

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

Lemma 9. *Conditioned on event $\{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\lambda\}$, any ϵ -optimal solution $\tilde{\Sigma}$ to (1) satisfies $\tilde{\Sigma} \in \Sigma + \mathbb{C}(l)$, where $l = 4s^{1/2}$. Moreover, assume $\tilde{\Sigma} \in \Sigma + \mathbb{B}^\infty(C\alpha)$. Then, conditioned on the event $\mathcal{E}_1(C\alpha, \kappa) \cap \{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\lambda\}$,*

$$\begin{aligned} \left\| \tilde{\Sigma} - \Sigma \right\|_F &\leq \kappa^{-1} \left\{ \lambda s^{1/2} + \|\nabla L_\alpha(\Sigma)_S\|_F + s^{1/2}\epsilon \right\} \\ &\leq 1.5\kappa^{-1}\lambda s^{1/2}. \end{aligned}$$

Proof: Conditioned on the stated event, Lemma 8 indicates

$$\left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \leq 3 \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_1.$$

Therefore,

$$\left\| \tilde{\Sigma} - \Sigma \right\|_1 \leq 4s^{1/2} \left\| \tilde{\Sigma} - \Sigma \right\|_F,$$

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

Now we prove the second statement. Since $\tilde{\Sigma} - \Sigma \in \mathbb{B}^\infty(C\alpha)$, conditioned on event $\mathcal{E}_1(C\alpha, \kappa)$, we have

$$\langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \geq \kappa \left\| \tilde{\Sigma} - \Sigma \right\|_F^2 \quad (8)$$

Now we upper bound the right-hand side of (8). For any $\Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1,\text{off}}$, write

$$\begin{aligned} & \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\ &= \underbrace{\langle U(\Xi), \tilde{\Sigma} - \Sigma \rangle}_{:=\Pi_1} - \underbrace{\langle \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle}_{:=\Pi_2} - \underbrace{\langle \lambda \Xi, \tilde{\Sigma} - \Sigma \rangle}_{:=\Pi_3} \end{aligned} \quad (9)$$

where $U(\Xi) := \nabla L_\alpha(\tilde{\Sigma}) + \lambda \Xi \in \mathbb{R}^{d \times d}$. We have

$$\begin{aligned} |\Pi_1| &\leq \|U(\Xi)\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 + \|(U(\Xi))_S\|_F \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \\ |\Pi_2| &\leq \|\nabla L_\alpha(\Sigma)_S\|_F \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \\ &\quad + \|\nabla L_\alpha(\Sigma)_{\bar{S}}\|_\infty \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \end{aligned}$$

Turning to Π_3 , decompose $\lambda \Xi$ and $\tilde{\Sigma} - \Sigma$ according to $S \cup \bar{S}$ to reach

$$\begin{aligned} \Pi_3 &= \langle (\lambda \Xi)_S, (\tilde{\Sigma} - \Sigma)_S \rangle + \langle (\lambda \Xi)_{\bar{S}}, (\tilde{\Sigma} - \Sigma)_{\bar{S}} \rangle \\ \text{Since } \Sigma_{\bar{S}} &= \mathbf{0} \text{ and } \Xi \in \partial \left\| \tilde{\Sigma} \right\|_{1,\text{off}}, \text{ we have } \langle (\lambda \Xi)_{\bar{S}}, (\tilde{\Sigma} - \Sigma)_{\bar{S}} \rangle \\ &= \langle (\lambda \Xi)_{\bar{S}}, \tilde{\Sigma}_{\bar{S}} \rangle = \lambda \left\| \tilde{\Sigma}_{\bar{S}} \right\|_1 = \lambda \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1. \end{aligned}$$

$$\Pi_3 \geq \lambda \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 - \lambda s^{1/2} \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F$$

Combining (9) with our estimation for Π_1, Π_2 and Π_3 , we have

$$\begin{aligned} & \leq -\{\lambda - \|\nabla L_\alpha(\Sigma)\|_\infty - \|U(\Xi)\|_\infty\} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \\ & \quad + \|\nabla L_\alpha(\Sigma)_S\|_F \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F + \|(U(\Xi))_S\|_F \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \\ & \quad + \lambda s^{1/2} \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \end{aligned}$$

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

$$\begin{aligned} & \langle \nabla L_\alpha(\tilde{\Sigma}) - \nabla L_\alpha(\Sigma), \tilde{\Sigma} - \Sigma \rangle \\ & \leq -\{\lambda - \|\nabla L_\alpha(\Sigma)\|_\infty - \epsilon\} \left\| (\tilde{\Sigma} - \Sigma)_{\bar{S}} \right\|_1 \\ & \quad + \{\|\nabla L_\alpha(\Sigma)_S\|_F + s^{1/2}\epsilon\} \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \\ & \quad + \lambda s^{1/2} \left\| (\tilde{\Sigma} - \Sigma)_S \right\|_F \end{aligned} \quad (10)$$

It follows from $\tilde{\Sigma} \in \Sigma + \mathbb{B}^\infty(C\alpha)$, (8) and (10) that conditioned on $\mathcal{E}_1(C\alpha, \kappa) \cap \{\|\nabla L_\alpha(\Sigma)\|_\infty + \epsilon \leq 0.5\lambda\}$,

$$\begin{aligned} \kappa \left\| \tilde{\Sigma} - \Sigma \right\|_F^2 &\leq \\ & \left\{ \lambda s^{1/2} + \|\nabla L_\alpha(\Sigma)_S\|_F + s^{1/2}\epsilon \right\} \left\| \tilde{\Sigma} - \Sigma \right\|_F \end{aligned}$$

Therefore,

$$\begin{aligned} & \left\| \tilde{\Sigma} - \Sigma \right\|_F \\ & \leq \kappa^{-1} \left\{ \lambda s^{1/2} + \|\nabla L_\alpha(\Sigma)_S\|_F + s^{1/2}\epsilon \right\} \\ & \leq \kappa^{-1} \{ \lambda s^{1/2} + 0.5\lambda s^{1/2} \} = 1.5\kappa^{-1}\lambda s^{1/2} \end{aligned} \quad (11)$$

■

A. Proof of Proposition 4

Proof: $\left\| \tilde{\Sigma} - \Sigma^* \right\|_F \leq 3\lambda s^{1/2}$ follows immediately from Lemma 9 with $\Sigma = \Sigma^*$ and $C = \kappa = 1/2$. Combining this with $\tilde{\Sigma} \in \Sigma^* + \mathbb{C}(l)$, where $l = 4s^{1/2}$, yields $\left\| \tilde{\Sigma} - \Sigma^* \right\|_1 \leq 12\lambda s$. ■

B. Proof of Proposition 5

We adopt the following notations for the next stage of proof. Recall that $L_\alpha(\Sigma) = \sum_{k,\ell} \frac{1}{n} \sum_{i=1}^n \rho_\alpha(\Sigma_{k\ell} - x_{ik}x_{i\ell})$. Define $\mathbf{B}^* := \mathbb{E}[\nabla L_\alpha(\Sigma^*)]$, and $\mathbf{W}^* := \nabla L_\alpha(\Sigma^*) - \mathbb{E}[\nabla L_\alpha(\Sigma^*)]$.

Lemma 10. *Recall that K is the constant defined in Remark 2. We have $|(\mathbf{B}^*)_{kl}| = |\mathbb{E}[\rho'_\alpha(\epsilon_{kl})]| < \frac{K}{\alpha}$ for all $k, l \in [d]$.*

Proof: For fixed $k, l \in [d]$, let $\epsilon_{kl} := \Sigma_{k\ell}^* - x_{ik}x_{i\ell}$, then

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

Therefore, for all k, l

$$|(\mathbf{B}^*)_{kl}| = \frac{1}{n} \left| \sum_{i=1}^n \mathbb{E}[\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})] \right| < \frac{K}{\alpha}.$$

■

C. Proof of Proposition 6

Proof: $W_{kl}^* = \frac{1}{n} \sum_{i=1}^n \{\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell}) - \mathbb{E}[\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})]\}$. Given that $|\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})| \leq \alpha$, for all $m \geq 2$:

$$\begin{aligned} &\mathbb{E}[\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})]^m \\ &\leq \alpha^{m-2} \cdot \text{Var}[\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})] \\ &\leq \alpha^{m-2} \cdot \text{Var}[\Sigma_{k\ell}^* - x_{ik}x_{i\ell}] \\ &\leq \alpha^{m-2} K \leq \alpha^{m-2} K \cdot m!/2 \end{aligned}$$

The second inequality follows given $\rho'_\alpha(\cdot)$ is 1-Lipschitz. With Bernstein's inequality,

$$\Pr \left(\left| \sum_{i=1}^n \{\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell}) - \mathbb{E}[\rho'_\alpha(\Sigma_{k\ell}^* - x_{ik}x_{i\ell})]\} \right| \geq 7\sqrt{Kn \log d} \right)$$

$$\begin{aligned} &\leq 2 \cdot \exp \left(-\frac{(7\sqrt{Kn \log d})^2/2}{Kn + \alpha \cdot 7\sqrt{Kn \log d}} \right) \\ &= 2 \cdot \exp \left(-\frac{49 \log d}{16} \right) < \frac{2}{d^3} \end{aligned}$$

Recall that $\nabla L_\alpha(\Sigma^*) = \mathbf{B}^* + \mathbf{W}^*$. With Lemma 10, we have $\|\mathbf{B}^*\|_\infty < \frac{K}{\alpha} \leq \sqrt{K \log d/n}$. Combing the two parts together and 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. ■

REFERENCES

- [1] I. T. Jolliffe, *Principal Component Analysis*. Springer, 2002.
- [2] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2009.
- [3] 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.