

Robust Scatter Matrix Estimation for High Dimensional Distributions With Heavy Tail

Junwei Lu¹, Fang Han, and Han Liu

Abstract—This paper studies large scatter matrix estimation for heavy tailed distributions. The contributions of this paper are twofold. First, we propose and advocate to use a new distribution family, the pair-elliptical, for modeling the high dimensional data. The pair-elliptical is more flexible and easier to check the goodness of fit compared to the elliptical. Secondly, built on the pair-elliptical family, we advocate using quantile-based statistics for estimating the scatter matrix. For this, we provide a family of quantile-based statistics. They outperform the existing ones for better balancing the efficiency and robustness. In particular, we show that the propose estimators have comparable performance to the moment-based counterparts under the Gaussian assumption. The method is also tuning-free compared to Catoni's M-estimator for covariance matrix estimation. We further apply the method to conduct a variety of statistical methods. The corresponding theoretical properties as well as numerical performances are provided.

Index Terms—Heavy-tailed distribution, pair-elliptical distribution, quantile-based statistics, scatter matrix.

I. INTRODUCTION

LARGE covariance matrix estimation is a core problem in multivariate statistics. Pearson's sample covariance matrix is widely used for estimation and proves to enjoy certain optimality under the subgaussian assumption [1]–[5]. However, this assumption is not realistic in many real applications where data are heavy-tailed [6]–[8].

To handle heavy-tailed data, rank-based statistics are proposed. Compared to Pearson's sample covariance, rank-based estimators achieve extra efficiency via exploiting the dataset's geometric structures. Such structures, like symmetry, are naturally involved in the data generating scheme and allow for both efficient and robust inference. Conducting rank-based covariance matrix estimation includes two steps. The first

step is to estimate the (latent) correlation matrix. For this, [9]–[14], and [15] exploit Spearman's rho and Kendall's tau estimators. They work under the nonparanormal or the transelliptical distribution family. The second step is to estimate marginal variances. For this, [16], [17], and [18] exploit Catoni's M-estimator [19]. However, Catoni's estimator requires to tune parameters. Moreover, it is sensitive to outliers and accordingly is not a robust estimator.

In this paper, we strengthen the results in the literature in two directions. First, we propose and advocate to use a new distribution family, the pair-elliptical. The pair-elliptical family is strictly larger and requires less symmetry structure than the elliptical. We provide detailed studies on the relation between the pair-elliptical and several heavy tailed distribution families, including the nonparanormal, elliptical, and transelliptical. Moreover, it is easier to test the goodness of fit for the pair-elliptical. For conducting such a test, we combine the existing results in low dimensions [20]–[24] with the familywise error rate controlling techniques including the Bonferroni's correction, the Holm's step-down procedure [25], and the higher criticism method [26], [27].

Secondly, built on the pair-elliptical family, we propose a new set of quantile-based statistics for estimating scatter/covariance matrices¹. We also provide the theoretical properties of the proposed methods. In particular, we show that the proposed quantile-based methods outperform the existing ones for better balancing the robustness and efficiency. As applications, we exploit the proposed estimators for conducting several high dimensional statistical methods, and show the advantages of using the quantile-based statistics both theoretically and empirically.

A. Other Related Works

The quantile-based statistics, such as the median absolute deviation [29] and the Q_n estimators [30], [31], have been used in estimating marginal standard deviations. Their properties in parameter estimation and robustness to outliers are further studied in low dimensions [32]. Moreover, these estimators have been generalized to estimate the dispersions between random variables [33]–[36].

¹The scatter matrix is any matrix proportional to the covariance matrix. See [28] for more details.

Manuscript received May 18, 2016; revised August 24, 2018; accepted January 15, 2020. Date of publication June 11, 2021; date of current version July 14, 2021. (Corresponding author: Junwei Lu.)

Junwei Lu is with the Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115 USA (e-mail: junweilu@hsph.harvard.edu).

Fang Han is with the Department of Statistics, University of Washington, Seattle, WA 98195 USA (e-mail: fanghan@uw.edu).

Han Liu is with the Department of Electrical Engineering and Computer Science, Northwestern University, Evanston, IL 60208 USA (e-mail: hanliu@northwestern.edu).

Communicated by C. Caramanis, Associate Editor for Machine Learning. Color versions of one or more figures in this article are available at <https://doi.org/10.1109/TIT.2021.3088381>.

Digital Object Identifier 10.1109/TIT.2021.3088381

Given these results, we mainly make three contributions: (i) Methodologically, we propose new quantile-based scatter matrix estimators that generalize the existing MAD and Q_n estimators for better balancing the efficiency and robustness. (ii) Theoretically, we provide more insights on the quantile-based methods. They confirm that the quantile-based estimators are also good alternatives to the prevailing moment-based estimators in high dimensions. (iii) We propose a projection method for overcoming the lack of positive semidefiniteness, which is typical in the robust scatter matrix estimation. This approach maintains the efficiency as well as robustness to data contamination, while the prevailing SVD decomposition approach [36] cannot.

Of note, the effectiveness of quantile-base methods is being realized in other fields in high dimensional statistics. For example, [37], [38], and [39] provide analysis on the penalized quantile regression and show that it can handle the case that the noise term is very heavy-tailed. Our method, although very different from theirs, shares similar properties.

B. Notation System

Let $\mathbf{M} = [M_{jk}] \in \mathbb{R}^{d \times d}$ be a matrix and $\mathbf{v} = (v_1, \dots, v_d)^T \in \mathbb{R}^d$ be a vector. We denote the subvector of \mathbf{v} by \mathbf{v}_I whose entries are indexed by a set $I \subset \{1, \dots, d\}$. We define $\mathbf{M}_{I,J}$ as the submatrix of \mathbf{M} whose rows and columns are indexed by I and J . For $0 < q < \infty$, we define the ℓ_0 , ℓ_q , and ℓ_∞ vector (pseudo-)norms to be $\|\mathbf{v}\|_0 := \sum_{j=1}^d I(v_j \neq 0)$, $\|\mathbf{v}\|_q := (\sum_{j=1}^d |v_j|^q)^{1/q}$, and $\|\mathbf{v}\|_\infty := \max_{1 \leq j \leq d} |v_j|$. Here $I(\cdot)$ represents the indicator function. For a matrix \mathbf{M} , we define the matrix ℓ_q , ℓ_{\max} , and Frobenius ℓ_F -norms of \mathbf{M} as $\|\mathbf{M}\|_q := \max_{\|\mathbf{v}\|_q=1} \|\mathbf{M}\mathbf{v}\|_q$, $\|\mathbf{M}\|_{\max} := \max_{j,k} |M_{jk}|$, and $\|\mathbf{M}\|_F := (\sum_{j,k} |M_{jk}|^2)^{1/2}$. For any matrix $\mathbf{M} \in \mathbb{R}^{d \times d}$, we define $\text{diag}(\mathbf{M})$ as the diagonal matrix with the same diagonal entries as \mathbf{M} , and $\mathbf{I}_d \in \mathbb{R}^{d \times d}$ to be the d by d identity matrix. Let $\lambda_j(\mathbf{M})$ and $\mathbf{u}_j(\mathbf{M})$ represent the j -th largest eigenvalue and the corresponding eigenvector of \mathbf{M} , and $\langle \mathbf{M}_1, \mathbf{M}_2 \rangle := \text{Tr}(\mathbf{M}_1^T \mathbf{M}_2)$ be the inner product of \mathbf{M}_1 and \mathbf{M}_2 . For any two random vectors \mathbf{X} and \mathbf{Y} , we write $\mathbf{X} \stackrel{D}{=} \mathbf{Y}$ if and only if \mathbf{X} and \mathbf{Y} are identically distributed. Throughout the paper, we let c, C be two generic absolute constants, whose values may vary at different locations.

C. Paper Organization

The rest of the paper is organized as follows. In the next section we provide the theoretical evaluation of the impacts of heavy tails on the moment-based estimators. This motivates our work. Section III proposes the pair-elliptical family, and reveals the connection among the Gaussian, elliptical, non-paranormal, transelliptical, and pair-elliptical. In Section IV, we introduce the generalized MAD and Q_n estimators for estimating scatter/covariance matrices. Section V provides the theoretical results. Section VI discusses parameter selection. In Section VII, we apply the proposed estimators to conduct multiple multivariate methods. We put experiments on synthetic and real data in Section VIII, more discussions in Section IX, and technical proofs in the appendix.

II. IMPACTS OF HEAVY TAILS ON MOMENT-BASED ESTIMATORS

This section illustrates the motivation of quantile-based estimators. In particular, we show how moment-based estimators fail for heavy tailed data. These estimators include the sample mean and sample covariance matrix. Such estimators are known to be efficient under stringent moment assumptions [5]. However, their performance drops down when such assumptions are violated [9], [40].

We characterize the heavy tailedness by the L_p norm. In detail, for any random variable $X \in \mathbb{R}$ and integer $p \geq 1$, we define the L_p norm of X as

$$\|X\|_{L_p} := (\mathbb{E}|X|^p)^{1/p}.$$

The random variable X is heavy tailed if there exists some $p > 0$ such that

$$\|X\|_{L_q} \leq K \leq \infty \text{ for } q \leq p, \text{ and } \|X\|_{L_{p+1}} = \infty.$$

The heavy tailedness of X is measured by how large p could be such that the p -th moment exists.

In the following we first provide an upper bound of the sample mean. It illustrates the “optimal rate but sub-optimal scaling” phenomenon.

Theorem II.1: Suppose $\mathbf{X} = (X_1, \dots, X_d)^T \in \mathbb{R}^d$ is a random vector with the population mean $\boldsymbol{\mu}$. Assume \mathbf{X} satisfies $\|X_j\|_{L_p} \leq K$, where we assume $d = O(n^\gamma)$ and $p = 2 + 2\gamma + \delta$. Letting $\bar{\boldsymbol{\mu}}$ be the sample mean of n independent observations of \mathbf{X} , we then have, with probability no smaller than $1 - 2d^{-2.5} - (\log d)^{p/2} n^{-\delta/2}$,

$$\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_\infty \leq 12K \cdot \sqrt{\frac{\log d}{n}}.$$

Theorem II.1 shows that, for preserving the $O_P(\sqrt{\log d/n})$ rate of convergence, p determines how large the dimension d can be compared to n . For example, when at most $(4 + \epsilon)$ -th moment exists for \mathbf{X} for some $\epsilon > 0$, the sample mean attains the optimal rate $O_P(\sqrt{\log d/n})$ under the suboptimal scaling $d = O(n)$.

The results in Theorem II.1 cannot be improved without adding more assumptions. Via a worst case analysis, the next theorem characterizes the sharpness of Theorem II.1.

Theorem II.2: For any fixed constant C , $p = 2 + 2\gamma$ with $\gamma > 0$, and $d = n^{\gamma+\delta_0}$ for some $\delta_0 > 0$, there exists certain random vector \mathbf{X} , satisfying

$$\|X_i\|_{L_q} < K, \text{ for some absolute constant } K > 0$$

and all $q \leq p$, such that, with probability tending to 1, we have

$$\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_\infty \geq \sqrt{\frac{C \log d}{n}}.$$

Theorems II.1 and Theorem II.2 together illustrate the constraints of applying moment-based estimators to study heavy tailed distributions. This motivates us to consider alternative methods that are more efficient in handling heavy tailedness.

III. PAIR-ELLIPTICAL DISTRIBUTION

In this section, we introduce the pair-elliptical distribution family. We first briefly review several existing distribution families: Gaussian, elliptical, nonparanormal, and transelliptical. Then we elaborate the relations between the pair-elliptical and aforementioned families.

A. Multivariate Distribution Families

We start by first introducing the elliptical distribution. The elliptical family contains symmetric but possibly very heavy tailed distributions.

Definition III.1 (Elliptical Distribution, [41]): A d -dimensional random vector \mathbf{X} is said to follow an elliptical distribution if and only if there exists a vector $\boldsymbol{\mu} \in \mathbb{R}^d$, a nonnegative random variable $\xi \in \mathbb{R}$, a matrix $\mathbf{A} \in \mathbb{R}^{d \times q}$ ($q \leq d$) of rank q , a random vector $\mathbf{U} \in \mathbb{R}^q$ uniformly distributed in q -dimension sphere \mathbb{S}^{q-1} and independent from ξ , such that

$$\mathbf{X} \stackrel{D}{=} \boldsymbol{\mu} + \xi \mathbf{A} \mathbf{U}.$$

In this case, we represent $\mathbf{X} \sim EC_d(\boldsymbol{\mu}, \mathbf{S}, \xi)$, where $\mathbf{S} := \mathbf{A} \mathbf{A}^T$ is of rank q .

Remark III.2: An equivalent definition of the elliptical distribution is: Any random vector $\mathbf{X} \sim EC_d(\boldsymbol{\mu}, \mathbf{S}, \xi)$ is elliptically distributed if and only if the characteristic function of \mathbf{X} is of the form $\exp(it^T \boldsymbol{\mu}) \phi(t^T \mathbf{S} t)$, where i is the imaginary number satisfying $i^2 = -1$, ϕ is a properly defined characteristic function, and there exists a one to one map between ξ and ϕ . In this case, we represent $\mathbf{X} \sim EC_d(\boldsymbol{\mu}, \mathbf{S}, \phi)$. Moreover, when the elliptical distribution is absolutely continuous, the density function is of the form $g((\mathbf{x} - \boldsymbol{\mu})^T \mathbf{S}^{-1} (\mathbf{x} - \boldsymbol{\mu}))$ for some nonnegative function $g(\cdot)$. In this case, we represent $\mathbf{X} \sim EC_d(\boldsymbol{\mu}, \mathbf{S}, g)$.

Although elliptical distributions have been extensively explored in modeling many real world data, including financial [42]–[45] and imaging data [46], [47], it can be quite restrictive due to the symmetry constraint [48]. One way to handle asymmetric data is to exploit the copula technique. This results to the transelliptical family (meta-elliptical family) proposed and discussed in [49] and [8]. Below we give the formal definition of the transelliptical distribution in [8].

Definition III.3 (Transelliptical Distribution, [8]): A continuous random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ follows a transelliptical distribution, denoted by $\mathbf{X} \sim TE_d(\boldsymbol{\Sigma}^0, \xi; f_1, \dots, f_d)$, if there exist univariate strictly increasing functions f_1, \dots, f_d such that

$$(f_1(X_1), \dots, f_d(X_d))^T \sim EC_d(\mathbf{0}, \boldsymbol{\Sigma}^0, \xi),$$

where $\text{diag}(\boldsymbol{\Sigma}^0) = \mathbf{I}_d$ and $\mathbb{P}(\xi = 0) = 0$. (III.1)

In particular, when

$$(f_1(X_1), \dots, f_d(X_d))^T \sim N_d(\mathbf{0}, \boldsymbol{\Sigma}^0), \text{ where } \text{diag}(\boldsymbol{\Sigma}^0) = \mathbf{I}_d,$$

\mathbf{X} follows a nonparanormal distribution [9], [50]. Here $\boldsymbol{\Sigma}^0$ is called the latent generalized correlation matrix.

B. Pair-Elliptical Distribution

In this section we propose a new distribution family, the pair-elliptical. Compared to the elliptical and transelliptical, the pair-elliptical distribution is of more interest to us. Specifically, it balances the modeling flexibility and interpretability in covariance/scatter matrices estimation.

Definition III.4: A continuous random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ is said to follow a pair-elliptical distribution, denoted by $\mathbf{X} \sim PE_d(\boldsymbol{\mu}, \mathbf{S}, \xi)$, if and only if any pair of random variables $(X_j, X_k)^T$ of \mathbf{X} is elliptically distributed. In other words, we have

$$(X_j, X_k)^T \sim EC_2(\boldsymbol{\mu}_{\{j,k\}}, \mathbf{S}_{\{j,k\},\{j,k\}}, \xi)$$

for all $j \neq k \in \{1, \dots, d\}$ and $\mathbb{P}(\xi = 0) = 0$.

As a special example, a distribution is said to be pair-normal, written as $PN_d(\boldsymbol{\mu}, \mathbf{S})$, if any pairs of \mathbf{X} is bivariate Gaussian distributed.

It is obvious that the pair-elliptical family contains the elliptical distribution family. Moreover, the elliptical is a strict subfamily of the pair-elliptical by considering the following example.

Example III.5: Let $f(X_1, X_2, X_3)$ be the density function of a three dimensional standard Gaussian distribution with median $\mathbf{0}$ and covariance matrix \mathbf{I}_3 , and $\mathbf{X} = (X_1, X_2, X_3)^T$ be a 3-dimensional random vector with the density function

$$g(X_1, X_2, X_3) = \begin{cases} 2f(X_1, X_2, X_3), & \text{if } X_1 X_2 X_3 \geq 0, \\ 0, & \text{otherwise.} \end{cases} \quad (\text{III.2})$$

The distribution in Example III.5 with density in (III.2) is bivariate Gaussian distributed for any pairwise marginal distributions, and therefore belongs to the pair-elliptical family. On the other hand, this distribution is marginally Gaussian distributed but not multivariate Gaussian distributed, and accordingly cannot be elliptically distributed.

Example III.5 also shows that the pair-elliptical distribution can be asymmetric. Moreover, the pair-elliptical distribution has a naturally defined scatter matrix \mathbf{S} , which is proportional to the covariance matrix $\boldsymbol{\Sigma}$ when $\mathbb{E}\xi^2$ exists. This makes the pair-elliptical compatible with many multivariate methods such as principal component analysis and linear discriminant analysis.

The rest of this section focuses on characterizing the relations among the Gaussian, elliptical², transelliptical, nonparanormal, pair-elliptical, and pair-normal families. Recall that in this paper we are only interested in the continuous distributions with density existing. It is obvious that the Gaussian family is a strict subfamily of the elliptical, and the elliptical is also a strict subfamily of both the transelliptical and the pair-elliptical. The next proposition shows that the only intersection between the elliptical and the nonparanormal is the Gaussian.

Proposition III.6 ([14]): If a random vector is both nonparanormally and elliptically distributed, it must follow a Gaussian distribution.

²In the rest of this section we only focus on the continuous elliptical distributions with $\mathbb{P}(\xi = 0) = 0$. And we are only interested in those whose covariance matrix is not the identity.

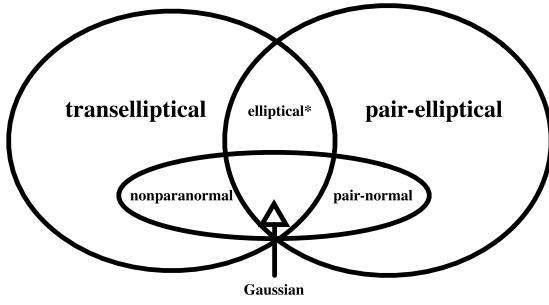


Fig. 1. The Venn diagram illustrating the relations of the Gaussian, elliptical, nonparanormal, transelliptical, pair-normal, and pair-elliptical families. Here “elliptical*” represents the continuous elliptical family with $\mathbb{P}(\xi = 0) = 0$.

In the next proposition, we show that the only intersection between the transelliptical and the pair-elliptical is the elliptical.

Proposition III.7: If a random vector is both transelliptically and pair-elliptically distributed, it must follow an elliptical distribution.

We defer the proof of Proposition III.7 to the appendix. In the end, let’s consider the relation among the pair-normal and all the other distribution families. By definition, the pair-normal is a strict subfamily of the pair-elliptical. On the other hand, the next proposition shows that any random scaled version of the pair-normal is pair-elliptically distributed.

Proposition III.8: Let $\mathbf{Y} \sim PN_d(\boldsymbol{\mu}, \mathbf{S})$ follow a pair-normal distribution. Then for any nonnegative random variable ξ with $\mathbb{P}(\xi = 0) = 0$ and independent of \mathbf{Y} , we have $\mathbf{X} = \boldsymbol{\mu}' + \xi \mathbf{Y}$ follows a pair-elliptical distribution.

In the end, we have the following proposition, which characterizes pair-normal’s connections to the elliptical and nonparanormal distributions.

Proposition III.9: For the pair-normal, elliptical, and nonparanormal distributions, we have

- (i) The only intersection between the pair-normal and elliptical is the Gaussian;
- (ii) The only intersection between the pair-normal and nonparanormal is the Gaussian.

In conclusion, the Venn diagram in Figure 1 summarizes the relation among the Gaussian, elliptical, nonparanormal, transelliptical, pair-normal, and pair-elliptical. From the figure, we can see that the Gaussian distribution locates in the central area whose covariance can be well estimated by the sample covariance matrix. The transelliptical covers the left-hand side of the diagram, where we advocate using rank-based estimators to estimate the covariance matrix. The pair-elliptical covers a new regime on the right-hand side of the diagram, where we will introduce the quantile-based estimators for estimating the covariance matrix.

C. Goodness of Fit Test of the Pair-Elliptical

This section proposes a goodness of fit test of the pair-elliptical. The pair-elliptical family has its advantages here: Both the transelliptical and elliptical require global geometric constraints over all covariates; In comparison, the

pair-elliptical only requires a local pairwise symmetry structure, which could be more easily checked. In this section, we combine the test of elliptical symmetry proposed in [24] with the step-down procedure in [25] for performing the pair-elliptical goodness of fit test.

Specifically, we propose a test of pair-elliptical:

$$H_0 : \text{The data are pair-elliptically distributed.} \quad (\text{III.3})$$

The proposed test is in two steps: In the first step, we test the pairwise elliptical symmetry; In the second step, we use the Holm’s step-down procedure to control the family-wise error.

In the first step, we apply the statistic proposed in [24] for testing pairwise elliptical symmetry. Let $\bar{\mathbf{Z}}$ and $\hat{\mathbf{S}}$ be the sample mean and sample covariance of $\{\mathbf{Z}_i\}_{i=1}^n$. We standardize the data by letting $\mathbf{Y}_i := \hat{\mathbf{S}}^{-1/2}(\mathbf{Z}_i - \bar{\mathbf{Z}})$ and $t(\mathbf{Z}_i) := \sqrt{2}\tilde{\mathbf{Y}}_i/w_i$, where $\tilde{\mathbf{Y}}_i := (\mathbf{Y}_{i1} + \mathbf{Y}_{i2})/2$ and $w_i^2 := \sum_{j=1}^2 (\mathbf{Y}_{ij} - \tilde{\mathbf{Y}}_i)^2$ for $i = 1, \dots, n$. Under H_0 , we have $t(\mathbf{Z}_i) \xrightarrow{D} t_1$ for $i = 1, \dots, n$, where t_1 is the t distribution with degree of freedom 1. To study the goodness-of-fit of the t -distribution, we define $M := \lfloor \sqrt{n} \rfloor$, where $\lfloor \cdot \rfloor$ represents the integer part of a real number and $E := n/M$. Let T_ℓ be the $\ell/M \times 100\%$ quantile of the t_1 distribution for $\ell = 0, \dots, M$, where $T_0 := -\infty$ and $T_M := +\infty$. We also denote the observed frequency $O_\ell := |\{t(\mathbf{Z}_i) : T_{\ell-1} < t(\mathbf{Z}_i) \leq T_\ell\}|$ for $1 \leq \ell \leq M$. [24] consider the following Pearson’s chi-squared test statistic:

$$Z(\{\mathbf{Z}_i\}) := \sum_{\ell=1}^M \frac{(O_\ell - E)^2}{E}.$$

By its nature, $Z(\{\mathbf{Z}_i\})$ is asymptotically chi-squared distributed with degrees of freedom $M - 1$.

In the second step, we screen the data to find whether there is any pair $\{X_j, X_k\}$ that does not follow an elliptical distribution. Considering the following null hypothesis for any $1 \leq j, k \leq d$:

$$H_{jk} : \{\mathbf{Z}_i\} \text{ are elliptically distributed,} \quad (\text{III.4})$$

we use the Holm’s step-down procedure [25] to control the family-wise error rate. Denote the p-values of $Z(\{(X_{ij}, X_{ik})\})$ as π_{jk} and let m_{jk} be the rank statistic of π_{jk} such that

$$m_{jk} := |\{\pi_{j'k'} \mid \pi_{j'k'} < \pi_{jk}, 1 \leq j', k' \leq d\}|.$$

The Holm’s adjusted p-values are defined as

$$\pi_{jk}^H := \max \{1 - (1 - \pi_{j'k'})^{t_{j'k'}} \mid m_{j'k'} \leq m_{jk}, j', k' \}, \quad (\text{III.5})$$

where $t_{jk} = 1 - \frac{2(m_{jk}-1)}{d(d-1)}$. Applying the adjusted p-values, we reject H_{jk} if π_{jk}^H is smaller than the level of significance α . Let $\omega_0 := \{(j, k) \mid H_{jk} \text{ in (III.4) is true, } 1 \leq j \neq k \leq d\}$, [25] shows that we can control the family-wise error rate as

$$\mathbb{P}_{\omega_0} \left(\pi_{jk}^H \leq \alpha \text{ for some } (j, k) \in \omega_0 \right) \leq \alpha.$$

Under the setting of goodness of fit test and H_0 in (III.3), we have $\omega_0 = \{(j, k) \mid 1 \leq j \neq k \leq d\}$ and therefore

$$\mathbb{P}_{H_0}(\pi_{jk}^H \leq \alpha \text{ for some } j \neq k) \leq \alpha.$$

IV. QUANTILE-BASED SCATTER MATRIX ESTIMATION

This section introduces the quantile-based scatter matrix estimators. To this end, we first briefly review the existing quantile-based estimators, including MAD and Q_n proposed in [29] and [30]. Secondly, we generalize these two estimators for estimating scatter matrices. Thirdly, we introduce the projection idea for constructing a positive semidefinite scatter matrix estimator.

A. Robust Scale Estimation

This section briefly reviews the robust estimators of the marginal standard deviation. To this end, we first define the population and sample versions of the quantile function. For any random variable $Z \in \mathbb{R}$ and fixed value $r \in [0, 1]$, let $Q(Z; r)$ represent the $r \times 100\%$ quantile of Z :

$$Q(Z; r) := \inf\{z : r \leq \mathbb{P}(Z \leq z)\}.$$

The $r \times 100\%$ quantile is said to be unique if and only if there exists one and only one $z \in \mathbb{R}$ such that $\mathbb{P}(Z \leq z) = r$. Letting $Z^{(1)} \leq \dots \leq Z^{(n)}$ be the ordered statistics of i.i.d. data $Z_1, \dots, Z_n \stackrel{D}{=} Z$, we define the empirical version of $Q(Z; r)$ as

$$\hat{Q}(\{Z_i\}; r) := Z^{(k^*)}, \text{ where } k^* := \arg \min_{i \in \{1, \dots, n\}} \left\{ \frac{i}{n} \geq r \right\}. \quad (\text{IV.1})$$

We then introduce the standard deviation estimators based on the quantiles. These include the median absolute deviation (MAD) and Q_n estimators. MAD is defined as follows:

$$\text{MAD estimator} : c^{\text{MAD}} \cdot \hat{Q}\left(\left|Z_i - \hat{Q}\left(\{Z_i\}; \frac{1}{2}\right)\right|; \frac{1}{2}\right),$$

where c^{MAD} is the constant making its population counterpart the standard deviation. In general, c^{MAD} is different for different distributions. MAD is robust to outliers with $1/2$ breakdown points [32], but has relatively low efficiency compared to the sample standard deviation under the Gaussian model.

To improve the efficiency while preserving the robustness, the Q_n estimator is proposed:

$$Q_n \text{ estimator} : c^{Q_n} \cdot \hat{Q}\left(\left\{|Z_i - Z_{i'}|\right\}_{i < i'}; 1/4\right),$$

where c^{Q_n} is another constant comparable to c^{MAD} . Q_n is known to be more efficient than MAD.

B. Robust Scatter Matrix Estimators

In this section, we propose our approach for estimating the scatter matrix. This is via generalizing the aforementioned quantile-based scale estimators.

Assume $\mathbf{X}_1, \dots, \mathbf{X}_n$ are n independent observations of a d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ with

$\mathbf{X}_i = (X_{i1}, \dots, X_{id})^T$. We first propose to estimate the marginal standard deviations by generalizing the MAD and Q_n estimators. Specifically, we define generalized median absolute deviation (gMAD) as follows: For the j -th entry of \mathbf{X} , the population and sample versions of gMAD are:

$$\begin{aligned} \text{gMAD} : \sigma^{\text{M}}(X_j; r) &:= Q(|X_j - Q(X_j; 1/2)|; r), \\ \hat{\sigma}^{\text{M}}(X_j; r) &:= \hat{Q}(|X_{ij} - \hat{Q}(\{X_{ij}\}_{i=1}^n; 1/2)|; r). \end{aligned} \quad (\text{IV.2})$$

Here the median is replaced by the $r \times 100\%$ quantile³. We then define the generalized Q_n estimator (gQNE) using the same idea. The population and sample versions of the gQNE for the j -th entry of \mathbf{X} are:

$$\begin{aligned} \text{gQNE} : \sigma^{\text{Q}}(X_j; r) &:= Q(|X_j - \tilde{X}_j|; r), \\ \hat{\sigma}^{\text{Q}}(X_j; r) &:= \hat{Q}(\{|X_{ij} - X_{i'j}|\}_{i < i'}; r), \end{aligned} \quad (\text{IV.3})$$

where $\tilde{\mathbf{X}} := (\tilde{X}_1, \dots, \tilde{X}_d)^T$ is an independent copy of \mathbf{X} . It is easy to check that, when setting $r = 1/2$ and $r = 1/4$ in (IV.2) and (IV.3), we recover the median absolute deviation (MAD) and Q_n estimators. This explains why we call them the generalized MAD and Q_n estimators. Of note, for any $j \in \{1, \dots, d\}$, we have $\text{median}(X_j - \tilde{X}_j) = 0$. Therefore, gQNE is a generalization to the gMAD estimator without requiring estimating the medians.

For estimating the scatter matrix, besides estimating the marginal scales, we also need to estimate the dispersion between any two random variables. For this, we follow the idea in [33]. We first remind that

$$\text{Cov}(X, Y) = \frac{1}{4} [\{\sigma(X + Y)\}^2 - \{\sigma(X - Y)\}^2],$$

where for any random variable Z , $\sigma(Z)$ represents the population standard deviation of Z . We then define the robust estimators of the dispersion between X and Y based on gMAD and gQNE as follows:

$$\begin{aligned} \sigma^{\text{M}}(X, Y; r) &:= \frac{1}{4} [\{\sigma^{\text{M}}(X + Y; r)\}^2 - \{\sigma^{\text{M}}(X - Y; r)\}^2]; \\ \sigma^{\text{Q}}(X, Y; r) &:= \frac{1}{4} [\{\sigma^{\text{Q}}(X + Y; r)\}^2 - \{\sigma^{\text{Q}}(X - Y; r)\}^2]. \end{aligned}$$

Let $\hat{\sigma}^{\text{M}}(X, Y)$ and $\hat{\sigma}^{\text{Q}}(X, Y; r)$ be the corresponding empirical versions. For any d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)^T$, we then define the d by d robust gMAD and gQNE scatter matrices $\mathbf{R}^{\text{M};r} = [\mathbf{R}_{jk}^{\text{M};r}]$ and $\mathbf{R}^{\text{Q};r} = [\mathbf{R}_{jk}^{\text{Q};r}]$ as follows: For any $j \in \{1, \dots, d\}$ and $k < j$, we write

$$\begin{aligned} \mathbf{R}_{jj}^{\text{M};r} &= (\sigma^{\text{M}}(X_j; r))^2, \quad \mathbf{R}_{jk}^{\text{M};r} = \mathbf{R}_{kj}^{\text{M};r} = \sigma^{\text{M}}(X_j, X_k; r); \\ \mathbf{R}_{jj}^{\text{Q};r} &= (\sigma^{\text{Q}}(X_j; r))^2, \quad \mathbf{R}_{jk}^{\text{Q};r} = \mathbf{R}_{kj}^{\text{Q};r} = \sigma^{\text{Q}}(X_j, X_k; r). \end{aligned}$$

In the later section we will show that $\mathbf{R}^{\text{M};r}$ and $\mathbf{R}^{\text{Q};r}$ are indeed scatter matrices under the pair-elliptical family. Let $\hat{\mathbf{R}}^{\text{M};r}$ and $\hat{\mathbf{R}}^{\text{Q};r}$ be the empirical versions of $\mathbf{R}^{\text{M};r}$ and $\mathbf{R}^{\text{Q};r}$ via replacing $\sigma^{\text{M}}(\cdot)$ and $\sigma^{\text{Q}}(\cdot)$ by $\hat{\sigma}^{\text{M}}(\cdot)$ and $\hat{\sigma}^{\text{Q}}(\cdot)$. $\hat{\mathbf{R}}^{\text{M};r}$ and $\hat{\mathbf{R}}^{\text{Q};r}$ are the proposed robust scatter matrix estimators.

There are two remarks. First, we do not discuss how to select r in this section, which will be studied in more details

³Later we will show that using the r -th quantile instead of the median can potentially increase the efficiency of the estimator, in the cost of losing some robustness though.

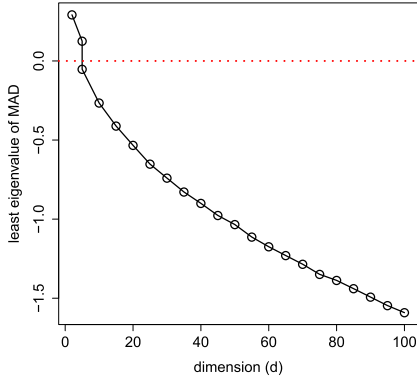


Fig. 2. The plot of the averaged least eigenvalues of a MAD scatter matrix (i.e., $\hat{\mathbf{R}}^{M;1/2}$) against the dimension d ranging from 2 to 200. Here the $n = 50$ observations are coming from the standard Gaussian distribution with dimension d , and the simulations are conducted with 100 repetitions.

in Section VI. Secondly, we note that $\hat{\mathbf{R}}^{M;r}$ and $\hat{\mathbf{R}}^{Q;r}$ are both symmetric matrices by definition. However, they are not necessarily positive semidefinite. We will discuss this issue in the next section.

C. Projection Method

In this section we introduce the projection idea to overcome the lack of positive semidefiniteness (PSD) in robust covariant matrix estimation. It is known that when the dimension is close to or higher than the sample size, the robust covariance matrix estimator can be non-PSD [28]. To illustrate this, Figure 2 shows the averaged least eigenvalue of the MAD scatter matrix estimator under the standard multivariate Gaussian model, with the sample size n fixed to be 50 and the dimension d increasing from 2 to 200.

The lack of PSD can cause problems for many high dimensional multivariate methods. To solve it, we propose a general projection method. In detail, for arbitrary non-PSD matrix estimator $\hat{\mathbf{R}}$, we consider the projection of $\hat{\mathbf{R}}$ to the positive semidefinite matrix cone:

$$\tilde{\mathbf{R}} = \arg \min_{\mathbf{M} \succeq \mathbf{0}} \|\mathbf{M} - \hat{\mathbf{R}}\|, \quad (\text{IV.4})$$

where $\mathbf{M} \succeq \mathbf{0}$ represents that \mathbf{M} is PSD and $\|\cdot\|$ is a certain matrix norm of interest. For any given norm $\|\cdot\|$, a computationally efficient algorithm to solve (IV.4) is given in Supplementary Material Section K.

Due to reasons that will become clearer later, we are interested in the projection with regard to the matrix element wise supremum norm $\|\cdot\|_{\max}$ in (IV.4). Of note, $\tilde{\mathbf{R}}$ and $\hat{\mathbf{R}}$ have the same breakdown point because $\tilde{\mathbf{R}}$ is independent of the data conditioning on $\hat{\mathbf{R}}$. Moreover, we have the following property about $\tilde{\mathbf{R}}$.

Lemma IV.1: Let $\tilde{\mathbf{R}}$ be the solution to (IV.4) with certain matrix norm $\|\cdot\|$ of interest. We have, for any $t \geq 0$ and $\mathbf{M} \in \mathbb{R}^{d \times d}$ with $\mathbf{M} \succeq \mathbf{0}$,

$$\mathbb{P}(\|\tilde{\mathbf{R}} - \mathbf{R}\| \geq t) \leq \mathbb{P}(\|\hat{\mathbf{R}} - \mathbf{R}\| \geq \frac{t}{2}).$$

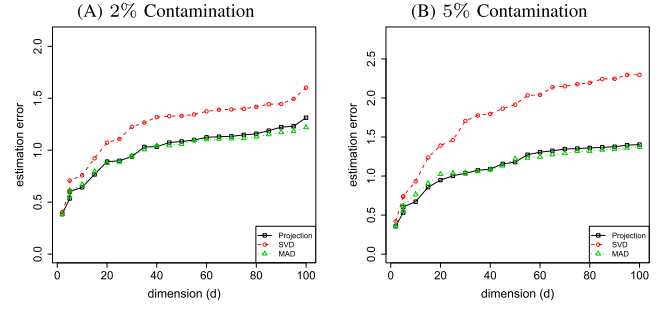


Fig. 3. The plot of the averaged estimation errors of the MAD (“MAD”) and PSD scatter matrix estimators using the projection and SVD decomposition ideas (denoted by “Projection” and “SVD”). The distances are calculated based on the $\|\cdot\|_{\max}$ norm and are plotted against the dimension d ranging from 2 to 100. Here the sample size is 50 and observations are coming from a standard Gaussian distributed data with dimension d , and 2% and 5% data points are randomly chosen and replaced by $+N(3, 3)$ or $-N(3, 3)$. The results are obtained based on 200 repetitions.

Of note, [36] propose an alternative approach to solve the non-PSD problem. Their method exploits the SVD decomposition of any given non-PSD matrix. However, Maronna’s method is not a robust procedure and is sensitive to outliers. More specifically, Figure 3 shows the averaged distance between the population scatter matrix and three different scatter matrix estimators: the possibly non-PSD MAD estimator (denoted by “MAD”), the PSD estimator calculated by using Maronna’s SVD decomposition idea (denoted by “SVD”), and the PSD estimator calculated by our projection idea with regard to the $\|\cdot\|_{\max}$ norm (denoted by “Projection”). Figures 3 (A) and (B) illustrate the results regarding a standard Gaussian distributed data (i.e., the data follow a $N_d(\mathbf{0}, \mathbf{I}_d)$ distribution) with 2% and 5% points being randomly chosen and replaced by $+N(3, 3)$ or $-N(3, 3)$. It shows that the PSD estimator obtained by projection is as insensitive as the MAD estimator (and their estimation accuracy is very close). On the other hand, Maronna’s method is very sensitive to such data contamination.

V. THEORETICAL RESULTS

This section provides the theoretical results of the proposed quantile-based gQNE and gMAD scatter matrix estimators. The section is divided into two parts: In the first part, under the pair-elliptical family, we characterize the relations among the population gQNE, gMAD statistics and Pearson’s covariance matrix; In the second part, we provide the theoretical analysis for gQNE and gMAD estimators.

A. Quantile-Based Estimators Under the Pair-Elliptical

In this section we show that the population gMAD and gQNE statistics, $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$, are scatter matrices of \mathbf{X} when \mathbf{X} is pair-elliptically distributed.

We first focus on gMAD. The next theorem characterizes a sufficient condition under which $\mathbf{R}^{M;r}$ is proportional to the covariance matrix. It also quantifies the scale constant $c^{M;r}$ that connects $\mathbf{R}^{M;r}$ to the covariance matrix.

Theorem V.1: Suppose that $\mathbf{X} = (X_1, \dots, X_d)^T$ is a d -dimensional random vector with the covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$. Then there exists some constant $c^{M;r}$ such that

$$\mathbf{R}^{M;r} = c^{M;r} \Sigma,$$

if for any $j \neq k \in \{1, \dots, d\}$,

$$\begin{aligned} \sqrt{c^{M;r}} &= Q\left(\left|\frac{X_j - Q(X_j, 1/2)}{\sigma(X_j)}\right|; r\right) \\ &= Q\left(\left|\frac{X_j + X_k - Q(X_j + X_k, 1/2)}{\sigma(X_j + X_k)}\right|; r\right) \\ &= Q\left(\left|\frac{X_j - X_k - Q(X_j - X_k, 1/2)}{\sigma(X_j - X_k)}\right|; r\right), \end{aligned} \quad (\text{V.1})$$

and the above quantiles are all unique.

We then study gQNE. The next theorem gives a sufficient condition under which $\mathbf{R}^{Q;r}$ is proportional to the covariance matrix and again quantifies the scale constant $c^{Q;r}$.

Theorem V.2: Suppose that $\mathbf{X} = (X_1, \dots, X_d)^T$ is a d -dimensional random vector with the covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$. Let $\tilde{\mathbf{X}}$ be an independent copy of \mathbf{X} and $\mathbf{Z} = (Z_1, \dots, Z_d)^T := \mathbf{X} - \tilde{\mathbf{X}}$. Then there exists some constant $c^{Q;r}$ such that

$$\mathbf{R}^{Q;r} = c^{Q;r} \Sigma,$$

if for any $j \neq k \in \{1, \dots, d\}$,

$$\begin{aligned} \sqrt{c^{Q;r}/2} &= Q\left(\left|\frac{Z_j}{\sigma(Z_j)}\right|; r\right) = Q\left(\left|\frac{Z_j + Z_k}{\sigma(Z_j + Z_k)}\right|; r\right) \\ &= Q\left(\left|\frac{Z_j - Z_k}{\sigma(Z_j - Z_k)}\right|; r\right), \end{aligned} \quad (\text{V.2})$$

and the above quantiles are all unique.

For any random variable $X \in \mathbb{R}$, Y is said to be the normalized version of X if $Y = (X - Q(X, 1/2))/\sigma(X)$. Accordingly, we have that (V.1) holds if the normalized versions of X_j , $X_j + X_k$, and $X_j - X_k$ are all identically distributed, and (V.2) holds if the normalized versions of Z_j , $Z_j + Z_k$, and $Z_j - Z_k$ are all identically distributed.

The next theorem shows that (V.1) and (V.2) hold under the pair-elliptical family.

Theorem V.3: For any pair-elliptically distributed random vector $\mathbf{X} \sim PE_d(\mu, \mathbf{S}, \xi)$, we have both $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$ are proportional to \mathbf{S} . In particular, when $\mathbb{E}\xi^2 < \infty$, we have both $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$ are proportional to the covariance matrix $\text{Cov}(\mathbf{X})$ and

$$c^{M;r} = \left(Q\left(X_0; \frac{1+r}{2}\right)\right)^2 \quad \text{and} \quad c^{Q;r} = 2\left(Q\left(Z_0; \frac{1+r}{2}\right)\right)^2, \quad (\text{V.3})$$

where X_0 and Z_0 are the normalized versions of X_1 and Z_1 .

Remark V.4: Theorem V.3 shows that, under the pair-elliptical family, $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$ are both proportional to $\text{Cov}(\mathbf{X})$ when the covariance exists. Of note, by Theorems V.1 and V.2, $\mathbf{R}^{M;r}$ or $\mathbf{R}^{Q;r}$ is proportional to $\text{Cov}(\mathbf{X})$ as long as (V.1) or (V.2) holds, and therefore can be applied to study potentially much larger family than the pair-elliptical.

B. Theoretical Properties of gMAD and gQNE

This section studies the estimation accuracy for the proposed scatter matrix estimators $\hat{\mathbf{R}}^{M;r}$ and $\hat{\mathbf{R}}^{Q;r}$. We show that the proposed methods are capable of handling heavy-tailed distributions, and shed light towards robust alternatives to many multivariate methods in high dimensions.

Before proceeding to the main results, we first introduce some extra notation. For any random vector $\mathbf{X} = (X_1, \dots, X_d)^T$ and any $j \neq k \in \{1, \dots, d\}$, we denote $F_{1;j}$, $\bar{F}_{1;j}$, $F_{2;j,k}$, $\bar{F}_{2;j,k}$, $F_{3;j,k}$, and $\bar{F}_{3;j,k}$ to be the distribution functions of X_j , $|X_j - Q(X_j, 1/2)|$, $X_j + X_k$, $|X_j + X_k - Q(X_j + X_k, 1/2)|$, $X_j - X_k$, and $|X_j - X_k - Q(X_j - X_k, 1/2)|$. We suppose that, for some constants κ_1 and η_1 that might scale with n , the following assumption holds:

$$\begin{aligned} (\text{A1}). \quad & \min_{\{j, |y - Q(F_{1;j}, 1/2)| < \kappa_1\}} \frac{d}{dy} F_{1;j}(y) \geq \eta_1, \\ & \min_{\{j, |y - Q(\bar{F}_{1;j}, r)| < \kappa_1\}} \frac{d}{dy} \bar{F}_{1;j}(y) \geq \eta_1, \\ & \min_{\{j \neq k, |y - Q(F_{2;j,k}, 1/2)| < \kappa_1\}} \frac{d}{dy} F_{2;j,k}(y) \geq \eta_1, \\ & \min_{\{j \neq k, |y - Q(\bar{F}_{2;j,k}, r)| < \kappa_1\}} \frac{d}{dy} \bar{F}_{2;j,k}(y) \geq \eta_1, \\ & \min_{\{j \neq k, |y - Q(F_{3;j,k}, 1/2)| < \kappa_1\}} \frac{d}{dy} F_{3;j,k}(y) \geq \eta_1, \\ & \min_{\{j \neq k, |y - Q(\bar{F}_{3;j,k}, r)| < \kappa_1\}} \frac{d}{dy} \bar{F}_{3;j,k}(y) \geq \eta_1, \end{aligned}$$

where for the random variable X with distribution function F , we denote $Q(F; r) := Q(X; r)$. Assumption (A1) requires that the density functions do not degenerate around median or the r -th quantiles. It is easy to check that Assumption (A1) is satisfied when we choose $\eta_1^{-1} = O(\sqrt{\|\Sigma\|_{\max}})$ for Gaussian distribution. Based on Assumption (A1), we have the following theorem, characterizing the estimation accuracy of the gMAD estimator.

Theorem V.5 (gMAD Concentration): Suppose that Assumption (A1) holds and κ_1 is lower bounded by a positive absolute constant. Then we have, for n large enough, with probability no smaller than $1 - 24\alpha^2$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{M;r} - \mathbf{R}^{M;r}\|_{\max} &\leq \\ &\max \left\{ \frac{6}{\eta_1^2} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{4\sqrt{\|\mathbf{R}^{M;r}\|_{\max}}}{\eta_1} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}. \end{aligned}$$

In particular, when \mathbf{X} is pair-elliptically distributed with the covariance matrix Σ existing, we have, with probability no smaller than $1 - 24\alpha^2$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{M;r} - c^{M;r} \Sigma\|_{\max} &\leq \\ &\max \left\{ \frac{6}{\eta_1^2} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{4\sqrt{\|c^{M;r} \Sigma\|_{\max}}}{\eta_1} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}. \end{aligned}$$

Theorem V.5 shows that, when $\kappa_1, \eta_1, \|\Sigma\|_{\max}$, and $c^{M;r}$ are upper and lower bounded by positive absolute constants, the convergence rate of $\hat{\mathbf{R}}^{M;r}$ with regard to the $\|\cdot\|_{\max}$ is $O_P(\sqrt{\log d/n})$. This is comparable to the existing results under subgaussian settings (See, for example, Theorem 1 in [4] and the discussion therein).

We then proceed to quantify the estimation accuracy of the gQNE estimator $\hat{\mathbf{R}}^{Q;r}$. Let $\tilde{\mathbf{X}} = (\tilde{X}_1, \dots, \tilde{X}_d)^T$ be an independent copy of \mathbf{X} . For any $j \neq k \in \{1, \dots, d\}$, let $G_{1;j}$, $G_{2;j,k}$, and $G_{3;j,k}$ be the distribution functions of $|X_j - \tilde{X}_j|$, $|X_j + X_k - (\tilde{X}_j + \tilde{X}_k)|$, and $|X_j - X_k - (\tilde{X}_j - \tilde{X}_k)|$. Suppose that for some constants κ_2 and η_2 that might scale with n , the following assumption holds:

$$(A2). \quad \begin{aligned} \min_{\{j, |y - Q(G_{1;j}; r)| < \kappa_2\}} \frac{d}{dy} G_{1;j}(y) &\geq \eta_2, \\ \min_{\{j \neq k, |y - Q(G_{2;j,k}; r)| < \kappa_2\}} \frac{d}{dy} G_{2;j,k}(y) &\geq \eta_2, \\ \min_{\{j \neq k, |y - Q(G_{3;j,k}; r)| < \kappa_2\}} \frac{d}{dy} G_{3;j,k}(y) &\geq \eta_2. \end{aligned}$$

Provided that Assumption (A2) holds, we have the following theorem. It gives the rate of convergence for $\hat{\mathbf{R}}^{Q;r}$ with regard to the element-wise supremum norm.

Theorem V.6 (gQNE concentration): Suppose that Assumption (A2) holds and κ_2 is lower bounded by a positive absolute constant. Then we have, for n large enough, with probability no smaller than $1 - 8\alpha$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} &\leq \\ \max \left\{ \frac{2}{\eta_2^2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{2\sqrt{\|\mathbf{R}^{Q;r}\|_{\max}}}{\eta_2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}. \end{aligned}$$

In particular, when \mathbf{X} is pair-elliptically distributed with the covariance matrix Σ existing, we have, with probability no smaller than $1 - 8\alpha$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{Q;r} - c^{Q;r} \Sigma\|_{\max} &\leq \\ \max \left\{ \frac{2}{\eta_2^2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{2\sqrt{\|c^{Q;r} \Sigma\|_{\max}}}{\eta_2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}. \end{aligned}$$

Similar to Theorem V.5, when $\kappa_2, \eta_2, \|\Sigma\|_{\max}$, and $c^{Q;r}$ are upper and lower bounded by positive absolute constants, the convergence rate of $\hat{\mathbf{R}}^{Q;r}$ is $O_P(\sqrt{\log d/n})$. Theorems V.5 and V.6 imply that, under the pair-elliptical family, the quantile-based estimators $\hat{\mathbf{R}}^{M;r}$ and $\hat{\mathbf{R}}^{Q;r}$ can be good alternatives to the sample covariance matrix.

Remark V.7: Consider the Gaussian distribution with the diagonal values of Σ lower bounded by an absolute constant. Then, for any fixed $r \in (0, 1)$ and lower bounded $\kappa_i, i = 1, 2$, Assumption (A1) and (A2) are satisfied with $\eta_1^{-1}, \eta_2^{-1} = O(\sqrt{\|\Sigma\|_{\max}})$. This implies that

$$\|\hat{\mathbf{R}}^{M;r} - c^{M;r} \Sigma\|_{\max} = O_P\left(\|\Sigma\|_{\max} \sqrt{\log d/n}\right)$$

$$\text{and } \|\hat{\mathbf{R}}^{Q;r} - c^{Q;r} \Sigma\|_{\max} = O_P\left(\|\Sigma\|_{\max} \sqrt{\log d/n}\right).$$

Let $\tilde{\mathbf{R}}^{M;r}$ and $\tilde{\mathbf{R}}^{Q;r}$ be the solutions to (IV.4). According to Lemma IV.1, we can also establish the concentration for $\tilde{\mathbf{R}}^{M;r}$ and $\tilde{\mathbf{R}}^{Q;r}$.

Corollary V.8: Under Assumptions (A1) and (A2), we have with probability no smaller than $1 - 24\alpha^2$,

$$\begin{aligned} \|\tilde{\mathbf{R}}^{M;r} - \mathbf{R}^{M;r}\|_{\max} &\leq \\ \max \left\{ \frac{3}{\eta_1^2} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{2\sqrt{\|\mathbf{R}^{M;r}\|_{\max}}}{\eta_1} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}; \end{aligned}$$

and with probability no smaller than $1 - 8\alpha$,

$$\begin{aligned} \|\tilde{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} &\leq \\ \max \left\{ \frac{1}{\eta_2^2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \right. \\ &\quad \left. \frac{\sqrt{\|\mathbf{R}^{Q;r}\|_{\max}}}{\eta_2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}. \end{aligned}$$

VI. SELECTION OF THE PARAMETER r

Theorems V.5 and V.6 show that the estimation accuracy of gMAD and gQNE estimators depends on the selection of the parameter r . In particular, the estimation error in estimating $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$ is related to η_1, η_2 and $c^{M;r}, c^{Q;r}$. On the other hand, r determines the breakdown points of $\mathbf{R}^{M;r}$ and $\mathbf{R}^{Q;r}$. Accordingly, the parameter r reflects the tradeoff between efficiency and robustness.

This section focuses on selecting the parameter r . The idea is to explore the parameter r that makes the corresponding estimator attain the highest statistical efficiency, given that the breakdown point is less than a predetermined critical value. Using Theorems V.5 and V.6, we have $\|\hat{\mathbf{R}}^{M;r}/c^{M;r} - \Sigma\|_{\max}$ and $\|\hat{\mathbf{R}}^{Q;r}/c^{Q;r} - \Sigma\|_{\max}$ are small when $\eta_1 \sqrt{c^{M;r}}$ and $\eta_2 \sqrt{c^{Q;r}}$ are large. Therefore, we aim at finding a parameter r such that the first derivatives of $\{\bar{F}_{1;j}, \bar{F}_{2;j,k}, \bar{F}_{3;j,k}\}$ or $\{G_{1;j}, G_{2;j,k}, G_{3;j,k}\}$ in a small interval around r times $\sqrt{c^{M;r}}$ or $\sqrt{c^{Q;r}}$ is the highest.

To this end, we separately estimate the derivatives and the scale parameters $\sqrt{c^{M;r}}$ and $\sqrt{c^{Q;r}}$. First, we estimate the derivatives of $\{\bar{F}_{1;j}, \bar{F}_{2;j,k}, \bar{F}_{3;j,k}\}$ or $\{G_{1;j}, G_{2;j,k}, G_{3;j,k}\}$ using the kernel density estimator [51]. For example for calculating the derivate of $\bar{F}_{1;j}$, we propose to use the data points:

$$|X_{1j} - \hat{Q}(\{X_{ij}\}_{i=1}^n, 1/2)|, \dots, |X_{nj} - \hat{Q}(\{X_{ij}\}_{i=1}^n, 1/2)|.$$

After obtaining the density estimators $\{\hat{f}_{1;j}, \hat{f}_{2;j,k}, \hat{f}_{3;j,k}\}$ or $\{\hat{g}_{1;j}, \hat{g}_{2;j,k}, \hat{g}_{3;j,k}\}$, we calculate the estimators $\hat{c}^{M;r}$ and $\hat{c}^{Q;r}$ of $c^{M;r}$ and $c^{Q;r}$ by comparing the scale of the standard deviation and its robust alternative (either gMAD or gQNE) for any chosen entry. We denote the empirical cumulative densities to be $\{\hat{F}_{1;j}, \hat{F}_{2;j,k}, \hat{F}_{3;j,k}\}$ and $\{\hat{G}_{1;j}, \hat{G}_{2;j,k}, \hat{G}_{3;j,k}\}$. We also denote their inverse function as $\{\hat{F}_{1;j}^{-1}, \hat{F}_{2;j,k}^{-1}, \hat{F}_{3;j,k}^{-1}\}$ and $\{\hat{G}_{1;j}^{-1}, \hat{G}_{2;j,k}^{-1}, \hat{G}_{3;j,k}^{-1}\}$.

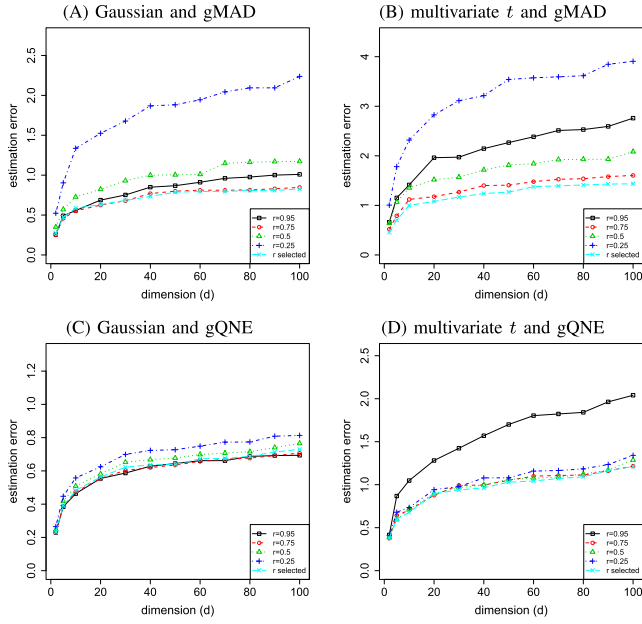


Fig. 4. The plot of the averaged estimation accuracy for the gMAD and gQNE estimators (with $r = 0.25$ to 0.95 and a selected r using the procedure described in Section VI). The distances are calculated based on the $\|\cdot\|_{\max}$ norm and are plotted against the dimension d ranging from 2 to 100. Here the $n = 50$ observations are coming from standard Gaussian or multivariate t distributed data with dimension d . The results are obtained based on 200 repetitions.

In the end, we define the statistics

$$q^{M;r} = \sqrt{c^{M;r}} \min_{j,k} \left\{ \hat{f}_{1;j}(\hat{F}_{1;j}^{-1}(r)), \hat{f}_{2;j,k}(\hat{F}_{2;j,k}^{-1}(r)), \hat{f}_{3;j,k}(\hat{F}_{3;j,k}^{-1}(r)) \right\},$$

$$q^{Q;r} = \sqrt{c^{Q;r}} \min_{j,k} \left\{ \hat{g}_{1;j}(\hat{G}_{1;j}^{-1}(r)), \hat{g}_{2;j,k}(\hat{G}_{2;j,k}^{-1}(r)), \hat{g}_{3;j,k}(\hat{G}_{3;j,k}^{-1}(r)) \right\}.$$

The estimators \hat{r}^M and \hat{r}^Q are then obtained as:

$$\hat{r}^M := \text{Truncate} \left(\arg \max_{r \in [0,1]} q^{M;r}; \delta_1, \delta_2 \right),$$

$$\hat{r}^Q := \text{Truncate} \left(\arg \max_{r \in [0,1]} q^{Q;r}; \delta_1, \delta_2 \right), \quad (\text{VI.1})$$

where δ_1, δ_2 are two pre-defined constants and for any $v \in [0, 1]$, we write

$$\text{Truncate}(v; \delta_1, \delta_2) := \begin{cases} \delta_1, & \text{if } v \leq \delta_1, \\ v, & \text{if } \delta_1 < v < \delta_2, \\ \delta_2, & \text{if } v \geq \delta_2. \end{cases}$$

Here we control the range of $r \in [\delta_1, \delta_2]$. This is for making the procedure robust and stable. Based on the empirical results, we recommend setting $\delta_1 = 0.25$ and $\delta_2 = 0.75$. In practice, we can also compare different values of $q^{M;r}$ and $q^{Q;r}$, and then choose one that achieve the best balance between estimation accuracy and robustness.

To illustrate the power of the proposed selection procedure, let's consider the following cases: Each time we randomly draw $n = 50$ i.i.d. standard Gaussian or multivariate t (with degrees of freedom 5) distributed data points with dimension d ranging from 2 to 100. We explore the gMAD and gQNE estimators with different quantile parameter r and scale all the scatter matrices such that their population versions are the

covariance matrix. We repeat the experiments for 200 times. Figure 4 illustrates the averaged estimation errors with regard to the $\|\cdot\|_{\max}$ norm. It shows that the procedure using the selected parameter r has the averagely smallest estimation error.

VII. APPLICATIONS

We now turn to various consequences of Theorems V.5 and V.6 in conducting multiple multivariate statistical methods in high dimensions. We mainly focus on the methods based on the gQNE estimator, while noting that the analysis for the methods based on the gMAD estimator is similar.

A. Sparse Covariance Matrix Estimation

We begin with estimating the sparse covariance matrix in high dimensions. Suppose we have

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} \mathbf{X} \in \mathbb{R}^d \text{ with covariance matrix } \Sigma.$$

Our target is to estimate the covariance matrix Σ itself. Pearson's sample covariance matrix performs poorly in high dimensions [52], [53]. One common remedy is to assume that the covariance matrix is sparse. This motivates different regularization procedures [2], [3], [54]–[57].

We focus on the method in [56] to illustrate the power of quantile-based statistics. This method directly produces a positive definite estimator and does not require any extra structure for the covariance matrix except for sparsity. The model we focus on is:

$$\mathcal{M}^{\text{COV-Q}}(\Sigma; s) := \{ \mathbf{X} \in \mathcal{M}^Q(\Sigma) : \lambda_d(\Sigma) > 0 \\ \text{and } \text{card}(\{(j, k) : \Sigma_{jk} \neq 0, j \neq k\}) \leq s \}.$$

Motivated by the above model, we solve the following optimization equation for obtaining the estimator $\hat{\Sigma}$:

$$\hat{\Sigma} = \arg \min_{\mathbf{M} = \mathbf{M}^T} \frac{1}{2} \|\hat{\Sigma} - \mathbf{M}\|_F^2 + \lambda \sum_{j \neq k} |\mathbf{M}_{j,k}|, \text{ s.t. } \lambda_d(\mathbf{M}) \geq \epsilon, \quad (\text{VII.1})$$

where $\hat{\Sigma}$ can be any positive semidefinite matrix based on the quantile-based scatter matrix estimator for approximating the covariance matrix, λ is a tuning parameter inducing sparsity, and ϵ is a small enough constant. Here we focus on using the gQNE estimator for generating the positive semidefinite matrix $\hat{\Sigma}$. For estimating the covariance instead of the scatter matrix, we need an extra efficient estimator of the scale parameter $c^{Q;r}$. Here $c^{Q;r}$ is described in Theorem V.2. When the exact distribution of the pair-elliptical is known, the value $c^{Q;r}$ could be theoretically calculated using (V.3). On the other hand, when the exact distribution is unknown, we note that estimating $c^{Q;r}$ is equivalent to estimating the marginal standard deviation for at least one entry X_j of \mathbf{X} . The proposed procedure then is in four steps:

1. Calculate the positive semidefinite matrix $\tilde{\mathbf{R}}^{Q;r}$:

$$\tilde{\mathbf{R}}^{Q;r} := \arg \min_{\mathbf{R} \succeq 0} \|\mathbf{R} - \hat{\mathbf{R}}^{Q;r}\|_{\max}. \quad (\text{VII.2})$$

Equation (VII.2) can be solved by a matrix maximum norm projection algorithm. See Supplementary Material Section K for details.

2. Choose the j -th entry with the lowest empirical fourth centered moments and use the sample standard deviation $\hat{\sigma}_j$ to estimate $\sigma(X_j)$.
3. Estimate $c^{Q;r}$ by $\hat{\mathbf{R}}_{jj}^{Q;r}/\hat{\sigma}_j^2$, and estimate Σ by

$$\hat{\Sigma}^{Q;r} = \frac{\hat{\sigma}_j^2}{\hat{\mathbf{R}}_{jj}^{Q;r}} \cdot \hat{\mathbf{R}}^{Q;r}. \quad (\text{VII.3})$$

4. Produce the final sparse matrix estimator $\tilde{\Sigma}^{Q;r}$ by plugging $\hat{\Sigma}^{Q;r}$ into (VII.1):

$$\tilde{\Sigma}^{Q;r} = \arg \min_{\mathbf{M}=\mathbf{M}^T} \frac{1}{2} \|\hat{\Sigma}^{Q;r} - \mathbf{M}\|_F^2 + \lambda \sum_{j \neq k} |\mathbf{M}_{j,k}|,$$

s.t. $\lambda_d(\mathbf{M}) \geq \epsilon$, where we prefix ϵ to be 10^{-5} as recommended in [56].

We then provide the statistical properties of the quantile-based estimator $\tilde{\Sigma}^{Q;r}$ as follows.

Corollary VII.1: Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are n independent observations of $\mathbf{X} \in \mathcal{M}^{\text{COV-Q}}(\Sigma; s)$ and the assumptions in Theorem 4.2 hold. We denote

$$\zeta := \max \left\{ \frac{4}{\eta_2^2} \left(\sqrt{\frac{4 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \frac{4 \sqrt{\|c^{Q;r} \Sigma\|_{\max}}}{\eta_2} \left(\sqrt{\frac{4 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\},$$

and

$$\zeta_j := \max \left\{ \frac{2}{\eta_2^2} \left(\sqrt{\frac{\log(1/\alpha)}{n}} + \frac{1}{n} \right)^2, \frac{2 \sqrt{c^{Q;r} \Sigma_{jj}}}{\eta_2} \left(\sqrt{\frac{\log(1/\alpha)}{n}} + \frac{1}{n} \right) \right\}.$$

- (i) We have, with probability no smaller than $1 - 8\alpha$,

$$\|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \leq \zeta. \quad (\text{VII.4})$$

- (ii) If $\mathbb{E}X_j^4 \leq K$ for some absolute constant K , then there exists an absolute constant c_1 only depending on K , such that when n is large enough, with probability no smaller than $1 - n^{-2\xi} - 12\alpha$,

$$\|\hat{\Sigma}^{Q;r} - \Sigma\|_{\max} \leq \left(\frac{c_1 n^{-1/2+\xi}}{\hat{\mathbf{R}}_{jj}^{Q;r} - \zeta_j} + \frac{1}{c^{Q;r}} \cdot \frac{\zeta_j}{\hat{\mathbf{R}}_{jj}^{Q;r} - \zeta_j} \right) \cdot (\|\mathbf{R}^{Q;r}\|_{\max} + \zeta) + \frac{\zeta}{c^{Q;r}}. \quad (\text{VII.5})$$

- (iii) If λ takes the same value as the right-hand side of (VII.5), then with probability no smaller than $1 - n^{-2\xi} - 12\alpha$, we have

$$\|\tilde{\Sigma}^{Q;r} - \Sigma\|_F \leq 5\lambda(s+d)^{1/2}.$$

In the following, for notation simplicity, we denote

$$\psi(r, j, \xi, \alpha) := \left(\frac{c_1 n^{-1/2+\xi}}{\hat{\mathbf{R}}_{jj}^{Q;r} - \zeta_j} + \frac{1}{c^{Q;r}} \cdot \frac{\zeta_j}{\hat{\mathbf{R}}_{jj}^{Q;r} - \zeta_j} \right) \cdot (\|\mathbf{R}^{Q;r}\|_{\max} + \zeta) + \frac{\zeta}{c^{Q;r}}. \quad (\text{VII.6})$$

Corollary VII.1 shows that, when $\eta_1, c^{Q;r}, \|\Sigma\|_{\max}$, and $\min_j \Sigma_{jj}$ are upper and lower bounded by positive absolute constants, via picking $\lambda > \psi(r, j, \xi, \alpha) = O(\sqrt{\log d/n})$, $\tilde{\Sigma}^{Q;r}$ approximates Σ in the rate

$$\|\tilde{\Sigma}^{Q;r} - \Sigma\|_F = O_P\left(\sqrt{\frac{(s+d) \log d}{n}}\right),$$

which is the minimax optimal rate shown in [58] and the parametric rate obtained in [56].

Remark VII.2: Following the steps introduced in this section, we could similarly calculate $\hat{\mathbf{R}}^{Q;r}$, $\hat{\Sigma}^{Q;r}$, and $\tilde{\Sigma}^{Q;r}$ for the gMAD estimator $\hat{\mathbf{R}}^{Q;r}$.

Remark VII.3: For handling heavy tailed data, the quantile-based and rank-based methods are intrinsically different. The rank-based estimator first calculates the correlation matrix, then estimate the d marginal variances via employing Catoni's M-estimator. However, for the quantile-based estimator, as (VII.3) shows, we only need to estimate one marginal variance. Moreover, the tuning in Catoni M-estimator is unnecessary.

B. Inverse Covariance Matrix Estimation

This section studies estimating inverse covariance matrix. Suppose $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ are i.i.d. copies of $\mathbf{X} \in \mathbb{R}^d$ with covariance matrix Σ . We are interested in estimating the precision matrix $\Theta := \Sigma^{-1}$. Precision matrix is closely related to graphical models and has been widely studied in the literature: [59] propose a neighborhood pursuit method; [60]–[62], and [63] apply penalized likelihood methods to obtain the estimators; [64] and [40] propose graphical Dantzig selector and CLIME estimator; [10] and [9], [14] generalize the Gaussian graphical model to the nonparanormal distribution, and [14] further generalize it to the transelliptical distribution.

In this section, we plug the covariance estimator $\hat{\Sigma}^{Q;r}$ in (VII.3) into the formulation of CLIME proposed by [40] and obtain the precision matrix estimator:

$$\begin{aligned} \hat{\Theta}^{Q;r} &= \arg \max_{\Theta \in \mathbb{R}^{d \times d}} \sum_{j,k} |\Theta_{jk}|, \\ \text{subject to } \|\hat{\Sigma}^{Q;r} \Theta - \mathbf{I}_d\|_{\max} &\leq \lambda. \end{aligned} \quad (\text{VII.7})$$

The CLIME estimator in (VII.7) can be reformulated as d linear programming [40]. Let \mathbf{A} be a symmetric matrix. For $q \in [0, 1]$, we define $\|\mathbf{A}\|_{q,\infty} := \max_i \sum_j |\mathbf{A}_{ij}|^q$. Considering the set

$$\mathcal{S}_d(q, s, M) := \{\Theta : \|\Theta\|_{1,\infty} \leq M \text{ and } \|\Theta\|_{q,\infty} \leq s\},$$

the following result gives the estimation accuracy of the precision matrix.

Corollary VII.4: Assume that $\Theta = \Sigma^{-1} \in \mathcal{S}_d(q, s, M)$ for some $q \in [0, 1]$ and the tuning parameter in (VII.7) satisfies

$$\lambda \geq \|\Theta\|_{1,\infty} \psi(r, j, \xi, \alpha),$$

where $\psi(r, j, \xi, \alpha)$ is defined in (VII.6). Then there exist constants C_1, C_2 such that

$$\begin{aligned} \|\hat{\Theta}^{Q;r} - \Theta\|_2 &\leq C_1 s \lambda^{1-q}, \|\hat{\Theta}^{Q;r} - \Theta\|_{\max} \leq \|\Theta\|_{1,\infty} \lambda \\ \text{and } \|\hat{\Theta}^{Q;r} - \Theta\|_F^2 &\leq C_2 s d \lambda^{2-q}. \end{aligned}$$

If $\eta_1, c^{Q;r}, \|\Sigma\|_{\max}$, and $\min_j \Sigma_{jj}$ are all upper and lower bounded by absolute constants, we choose $\lambda = O(\sqrt{\log d/n})$, and Corollary VII.4 gives us the rate

$$\|\hat{\Theta}^{Q;r} - \Theta\|_2 = O_P\left(s\left(\frac{\log d}{n}\right)^{\frac{1-q}{2}}\right),$$

$$\|\hat{\Theta}^{Q;r} - \Theta\|_{\max} = O_P\left(\sqrt{\frac{\log d}{n}}\right),$$

and

$$\frac{1}{d}\|\hat{\Theta}^{Q;r} - \Theta\|_F^2 \leq O_P\left(s\left(\frac{\log d}{n}\right)^{\frac{2-q}{2}}\right).$$

According to the minimax estimation rate of precision matrix estimators established in [65], the estimation rates in terms of all the above three matrix norms are optimal.

C. Sparse Principal Component Analysis

In this section we consider conducting principal component analysis (PCA) and sparse PCA. Recall that in (sparse) principal component analysis, we have

$$\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n \stackrel{i.i.d.}{\sim} \mathbf{X} \in \mathbb{R}^d \text{ with covariance matrix } \Sigma,$$

and our target is to estimate the eigenspace spanned by the m leading eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_m$ of Σ . The conventional PCA uses the leading eigenvectors of the sample covariance matrix for estimation. In high dimensions where d can be much larger than n , a sparsity constraint on $\mathbf{u}_1, \dots, \mathbf{u}_m$ is sometimes recommended [66], motivating a series of methods referred to as sparse PCA.

In this section, let $\mathbf{U}_m := (\mathbf{u}_1, \dots, \mathbf{u}_m) \in \mathbb{R}^{d \times m}$ represent the combination of eigenvectors of interest. We aim to estimate the projection matrix $\Pi_m := \mathbf{U}_m \mathbf{U}_m^T$. The model we focus on is:

$$\mathcal{M}^{\text{PCA-Q}}(\Sigma; s, m) := \left\{ \mathbf{X} \in \mathcal{M}^Q(\Sigma) : \sum_{j=1}^d I\left(\sum_{k=1}^m \mathbf{u}_{kj}^2 \neq 0\right) \leq s, \lambda_m(\Sigma) - \lambda_{m+1}(\Sigma) > 0 \right\},$$

where \mathbf{u}_{kj} is the j -th entry of \mathbf{u}_k and $\lambda_m(\Sigma)$ is the m -th largest eigenvalue of Σ . Motivated from the above model, via exploiting the gQNE estimator, we propose the quantile-based (sparse) PCA estimators (Q-PCA) as the optimum to the following Fantope projection problem [67]:

$$\hat{\Pi}_m^{Q;r} = \arg \max_{\Pi \in \mathbb{R}^{d \times d}} \text{Tr}(\Pi^T \hat{\mathbf{R}}^{Q;r}) - \lambda \|\Pi\|_{1,1},$$

subject to $0 \preceq \Pi \preceq \mathbf{I}_d$ and $\text{Tr}(\Pi) = m$, (VII.8)

where $\|\Pi\|_{1,1} = \sum_{1 \leq j, k \leq d} |\Pi_{jk}|$ and $\mathbf{A} \preceq \mathbf{B}$ implies $\mathbf{B} - \mathbf{A}$ is a positive semidefinite matrix. Intrinsically $\hat{\Pi}_m^{Q;r}$ is the estimator of the eigenspace spanned by the m leading eigenvectors of $\mathbf{R}^{Q;r}$. Because the scatter matrix shares the same eigenspace with the covariance matrix, under the model $\mathcal{M}^{\text{PCA-Q}}(\Sigma; s, m)$, $\hat{\Pi}_m^{Q;r}$ is also an estimator of Π_m . We then have the following corollary, stating that the Q-PCA estimator achieves the parametric rate of convergence in estimating Π_m under the model $\mathcal{M}^{\text{PCA-Q}}(\Sigma; s, m)$.

Corollary VII.5: Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_n$ are n independent observations of $\mathbf{X} \in \mathcal{M}^{\text{PCA-Q}}(\Sigma; s, m)$ and the

assumptions in Theorem V.6 hold. If the tuning parameter in (VII.8) satisfies $\lambda \geq \psi(r, j, \xi, \alpha)$ (reminding that $\psi(r, j, \xi, \alpha)$ is defined in (VII.6)), we have

$$\|\hat{\Pi}_m^{Q;r} - \Pi_m\|_F \leq \frac{4s\lambda}{\lambda_m(\Sigma) - \lambda_{m+1}(\Sigma)}.$$

Corollary VII.5 shows that, under appropriate conditions, the convergence rate of the Q-PCA estimator is $O_P(s\sqrt{\log d/n})$, which is the parametric rate in [67].

D. Discriminant Analysis

In this section, we consider the linear discriminant analysis for conducting high dimensional classification [12], [68]–[72]. Let data points $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$ be independently drawn from a joint distribution of (\mathbf{X}, Y) , where $\mathbf{X} \in \mathbb{R}^d$ and $Y \in \{1, 2\}$ is the binary label. We denote $I_1 = \{i : Y_i = 1\}$, $I_2 = \{i : Y_i = 2\}$, and $n_1 = |I_1|, n_2 = |I_2|$. Define $\pi = \mathbb{P}(Y = 1)$, $\mu_1 = \mathbb{E}(\mathbf{X} | Y = 1)$, $\mu_2 = \mathbb{E}(\mathbf{X} | Y = 2)$, $\Sigma = \text{Cov}(\mathbf{X} | Y = 1) = \text{Cov}(\mathbf{X} | Y = 2)$. If the classifier is defined as $h(\mathbf{x}) = I(f(\mathbf{x}) < c) + 1$ for some function f and constant c , we measure the quality of classification by employing the Rayleigh quotient [17]:

$$\text{Rq}(f) = \frac{\text{Var}\{\mathbb{E}[f(\mathbf{X}) | Y]\}}{\text{Var}\{f(\mathbf{X}) - \mathbb{E}[f(\mathbf{X}) | Y]\}}.$$

For the linear functions $f(\mathbf{x}) = \beta^T \mathbf{x} + c$, the Rayleigh quotient has the formulation

$$\text{Rq}(\beta) = \pi(1 - \pi) \frac{[\beta^T(\mu_1 - \mu_2)]^2}{\beta^T \Sigma \beta}.$$

The Rayleigh quotient is minimized when $\beta = \beta^* = \Sigma^{-1}(\mu_1 - \mu_2)$. When $\mathbf{X} | (Y = 1)$ and $\mathbf{X} | (Y = 2)$ are multivariate Gaussian distribution, it matches the Fisher's linear discriminant rule $h_F(\mathbf{x}) = I(\mathbf{x}^T \beta^* + c^*) + 1$, where $c^* = (\mu_1 + \mu_2)^T \beta^* / 2$. In order to estimate β^* and c^* , we apply the estimator proposed in [71]:

$$\hat{\beta}^{Q;r} = \arg \max_{\beta \in \mathbb{R}^d} \|\beta\|_1,$$

$$\text{subject to } \|\hat{\Sigma}^{Q;r} \beta - (\hat{\mu}_1 - \hat{\mu}_2)\|_{\max} \leq \lambda, \quad (\text{VII.9})$$

where $\hat{\mu}_1, \hat{\mu}_2$ are some estimators of μ_1, μ_2 .

In the following, we suggest two kinds of estimators of the mean vectors. For simplicity, we only describe the estimator for μ_1 and similar procedures can also be applied to estimate μ_2 . Suppose that $\mathbf{X}_1, \dots, \mathbf{X}_{n_1}$ correspond to $Y = 1$. The first estimator for μ_1 is the sample median $\hat{\mu}_M = (\hat{\mu}_{M1}, \dots, \hat{\mu}_{Md})^T$ where

$$\hat{\mu}_{Mj} = \text{median}(\{\mathbf{X}_{1j}, \dots, \mathbf{X}_{n_1 j}\}).$$

Due to the Hoeffding's inequality, we can derive that there exists some constant c_m that $\mathbb{P}(|\hat{\mu}_{Mj} - \mu_j| > c_m \sqrt{\log(1/\delta)/n}) \leq \delta$ for any $j = 1, \dots, d$.

[19] proposes an alternative M-estimator for μ_1 . Considering a strictly increasing function h such that $-\log(1 - y + y^2/2) \leq h(y) \leq \log(1 + y + y^2/2)$ and some $\delta \in (0, 1)$ satisfying $n > 2 \log(1/\delta)$, $v \geq \max\{\sigma_1^2, \dots, \sigma_d^2\}$, we define

$$\alpha_\delta^2 = \frac{2 \log(1/\delta)}{n_1 \left(v + \frac{2v \log(1/\delta)}{n - 2 \log(1/\delta)}\right)}.$$

The estimator $\hat{\mu}_C = (\hat{\mu}_{C1}, \dots, \hat{\mu}_{Cd})^T$ is obtained by solving the equation:

$$\sum_{i=1}^{n_1} h(\alpha_\delta(\mathbf{x}_{ij} - \hat{\mu}_{Cj})) = 0.$$

It is shown in [19] that there exists some constant C such that with probability at least $1 - (n \vee d)^{-1}$, $\mathbb{P}(\|\hat{\mu}_{Cj} - \mu_{Cj}\|^2 \geq \frac{2v \log(1/\delta)}{n-2 \log(1/\delta)}) \leq \delta$ for all $j = 1, \dots, d$. By union bound, for both estimators we have

$$\|\hat{\mu}_M - \mu_1\|_\infty \vee \|\hat{\mu}_C - \mu_1\|_\infty = O_P\left(\sqrt{\frac{\log d}{n}}\right). \quad (\text{VII.10})$$

Therefore, we have the following result on how well the Rayleigh quotient of estimated classifier can approximate the optimal one.

Corollary VII.6: We assume that $n_1 \asymp n_2$, $\Delta_d = (\mu_1 - \mu_2)^T \beta^* \geq M$ for some constant $M > 0$. If we choose the tuning parameter in (VII.9) that $\lambda \geq \|\hat{\mu}_1 - \hat{\mu}_2 - \mu_1 + \mu_2\|_\infty + \psi(r, j, \xi, \alpha) \|\beta^*\|_1$, we have

$$\frac{\text{Rq}(\hat{\beta}^{Q;r})}{\text{Rq}(\beta^*)} \geq 1 - 2M^{-1} \|\beta^*\|_1^2 \|\psi(r, j, \xi, \alpha) - 20M^{-1} \lambda \|\beta^*\|_1.$$

According to the rate in (VII.10), if $\eta_1, c^{Q;r}, \|\Sigma\|_{\max}$, and $\min_j \Sigma_{jj}$ are universally bounded from above and below by positive absolute constants and $\|\beta^*\|_1$ is also upper bounded by a constant independent of d and n , we choose the tuning parameter $\lambda = O(\sqrt{\log d/n})$ and obtain the rate through Corollary VII.6 as

$$\frac{\text{Rq}(\hat{\beta}^{Q;r})}{\text{Rq}(\beta^*)} \geq 1 - O\left(\sqrt{\frac{\log d}{n}}\right).$$

The above result matches the parametric rate in [71].

VIII. EXPERIMENTS

This section compares the empirical performance of quantile-based estimators to these based on the Pearson, Kendall's tau, and Spearman's rho covariance estimators for both synthetic and real data. Since for synthetic data analysis, conducting inverse covariance estimation and sparse PCA has reflected the quality of matrix estimation very well, these two methods are the main focus for synthetic data. For the real data, on the other hand, we aim to apply our methods to the classification of genomic data.

A. Synthetic Data

This section focuses on conducting inverse covariance matrix estimation and sparse PCA. Let $\{\hat{\sigma}_1, \dots, \hat{\sigma}_d\}$ be the sample standard deviations of \mathbf{X} . Let $\hat{\mathbf{S}}^P$, $\hat{\mathbf{S}}^K$ and $\hat{\mathbf{S}}^S$ be the Pearson's, Kendall's tau and Spearman's rho sample correlation matrices. We compare $\hat{\Sigma}^{Q,r}$ and $\hat{\Sigma}^{M,r}$ introduced in Section VII-A⁴ to three covariance matrix estimators:

- 1) Pearson's sample covariance matrix: $\hat{\Sigma}_{ij}^P = \hat{\sigma}_i \hat{\sigma}_j \hat{\mathbf{S}}_{ij}^P$, for $i, j = 1, \dots, d$;

⁴Here the parameter r is set as in Section VI with δ_1 and δ_2 to be 0.25 and 0.75.

TABLE I

COMPARISON OF THE GMAD, gQNE, PEARSON'S, KENDALL'S TAU AND SPEARMAN'S RHO COVARIANCE ESTIMATORS COMBINED WITH THE CLIME AND SPARSE PCA (SPCA) ALGORITHMS. THE ESTIMATION ERRORS OF CLIME ESTIMATORS ARE IN TERMS OF FROBENIUS NORMS. THE ESTIMATION ERRORS OF SPCA ESTIMATORS ARE IN TERMS OF $\|\hat{\Pi}_1 - \Pi_1\|_F$. THE ERRORS ARE AVERAGED OVER 100 REPETITIONS WITH STANDARD DEVIATIONS IN THE PARENTHESES

	d	Scheme 1			Scheme 2		
		100	200	500	100	200	500
CLIME	gMAD	5.47 (0.22)	6.29 (0.28)	8.82 (0.60)	7.73 (0.64)	9.06 (0.80)	11.09 (0.92)
	gQNE	4.75 (0.20)	5.42 (0.19)	7.53 (0.22)	7.58 (0.68)	8.76 (0.79)	10.95 (1.03)
	Pearson	4.75 (0.22)	5.38 (0.22)	7.39 (0.26)	7.77 (0.82)	8.95 (1.05)	11.37 (1.45)
	Kendall	5.38 (0.23)	6.35 (0.25)	8.84 (0.29)	8.47 (0.55)	10.01 (0.62)	13.34 (0.66)
	Spearman	4.86 (0.20)	5.49 (0.20)	7.60 (0.30)	7.99 (0.64)	9.17 (0.78)	11.33 (1.04)
SPCA	gMAD	0.35 (0.23)	0.41 (0.27)	0.48 (0.26)	0.48 (0.25)	0.53 (0.25)	0.60 (0.27)
	gQNE	0.19 (0.16)	0.25 (0.24)	0.34 (0.26)	0.37 (0.26)	0.43 (0.26)	0.47 (0.29)
	Pearson	0.17 (0.14)	0.20 (0.20)	0.31 (0.28)	0.78 (0.27)	0.80 (0.23)	0.91 (0.16)
	Kendall	0.20 (0.15)	0.23 (0.20)	0.32 (0.28)	0.45 (0.30)	0.53 (0.28)	0.58 (0.30)
	Spearman	0.21 (0.16)	0.25 (0.23)	0.34 (0.28)	0.46 (0.32)	0.54 (0.28)	0.57 (0.29)

- 2) Kendall's tau covariance matrix: $\hat{\Sigma}_{ij}^K = \hat{\sigma}_i \hat{\sigma}_j \hat{\mathbf{S}}_{ij}^K$, for $i, j = 1, \dots, d$;
- 3) Spearman's rho covariance matrix: $\hat{\Sigma}_{ij}^S = \hat{\sigma}_i \hat{\sigma}_j \hat{\mathbf{S}}_{ij}^S$, for $i, j = 1, \dots, d$.

Moreover, given a covariance matrix Σ , we consider the following three schemes for data generation:

- Scheme 1: $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. copies of $\mathbf{X} \sim N_d(\mathbf{0}, \Sigma)$;
- Scheme 2: $\mathbf{X}_1, \dots, \mathbf{X}_n$ are i.i.d. copies of \mathbf{X} following the multivariate t -distribution with the degree of freedom 5 and $\text{Cov}(\mathbf{X}) = \Sigma$

In the following, we plug the five covariance matrix estimators into the inverse covariance estimation procedure discussed in Section VII-B and the sparse PCA procedure in Section VII-C.

1) *Inverse Covariance Matrix Estimation:* We consider the numerical performance using the setting in [40]. Let $\Omega = (\omega_{ij})_{1 \leq i, j \leq d} := \Sigma^{-1}$ represent the inverse covariance matrix. We consider $\Omega = \mathbf{B} + \delta \mathbf{I}_d$. Each off-diagonal entry of \mathbf{B} is independent and of the value 0.5 with probability 0.1 and the value 0 with probability 0.9. The value δ is chosen such that the condition number of Ω is d . The diagonal of Σ is renormalized to be ones.

We generate the data under the three schemes with dimensions $d = 90, 120, 200$, sample size $n = 100$, and repetition time 100. We measure the estimation error by Frobenius norm $\|\hat{\Sigma} - \Sigma\|_F$. The numerical results are presented in Table I.

2) *Sparse PCA:* We use the setting in [13] to investigate the numerical performance of the sparse principal component

analysis. We consider the spike covariance $\Sigma = \lambda_1 v_1 v_1^T + \lambda_2 v_2 v_2^T + \mathbf{I}_d$, where $\lambda_1 = 5$ and $\lambda_2 = 2$, $v_{1j} = 1/\sqrt{10}$ for $j = 1 \leq j \leq 10$, $v_{1j} = 0$ for $j > 10$ and $v_{2j} = 1/\sqrt{10}$ for $11 \leq j \leq 20$, $v_{2j} = 0$ otherwise. The data sample size is $n = 50$ for $d = 90, 120, 200$ with repetition 100 times. We measure the difference between the true projection matrix $\Pi_1 = v_1 v_1^T$ and its estimator $\hat{\Pi}_1$ through (VII.8) by the Frobenius norm $\|\hat{\Pi}_1 - \Pi_1\|_F$. The results are given in Table I.

3) *Summary of Synthetic Data Results:* From the numerical results summarized in Table I, the gMAD and gQNE covariance estimators perform better than other estimators for both CLIME and sparse PCA under Scheme 2 and Scheme 3. This implies that gMAD and gQNE have better performance in studying heavy-tailed distributions and contaminated data. This is due to the advantage that only one variance is required to be estimated in (VII.3), while the other covariance estimators demand variance estimators for d covariates. Moreover, gQNE outperforms gMAD estimator. This matches the assertion that gQNE is more efficient than gMAD. Under the Scheme 1, since the synthetic data follow Gaussian distribution, the sample covariance estimator is efficient and outperforms the gMAD and gQNE covariance estimators. However, the performance is still comparable even under this scheme.

B. Genomic Data Analysis

In this section, we apply the method introduced in Section VII-D to the large-scale microarray dataset, GPL96 [73]. The dataset contains 20,263 probes with 8,124 samples belonging to 309 tissue types. The probes correspond to 12,679 genes and we average the probes from same genes. The tissue types include prostate tumor (148 samples), B cell lymphoma (213 samples), and many others. The target of our real data analysis is to compare the performance of classifying the prostate tumor from the B cell lymphoma by using gMDA, gQNE, and other estimators. We only focus on the top 1,000 genes of the largest p-values in performing the marginal two-sample t-test.

First, we consider the goodness of fit test of the pair-elliptical. Our goal is to select genes such that their corresponding data from both tissue types are approximately pair-elliptical separately. Then we apply the linear discriminative classifier in (VII.9) to the selected genes for classification. Specifically, let $\{X_i\}_{i=1}^n$ be the dataset corresponding to the prostate tumor. We employ the goodness of fit test proposed in Section III-C by calculating the statistic $Z(\{(X_{ij}, X_{ik})\})$ for each pair of $1 \leq j \neq k \leq d$. The Holm's adjusted p-values $\{\pi_{jk}^H\}_{j,k=1}^d$ are evaluated according to (III.5). We delete any gene j if there exists another gene k such that the adjusted p-value $\pi_{jk}^H < \alpha = 0.05$. All the left ones are the selected genes for the prostate tumor. Same procedures are also applied to the dataset of B cell lymphoma. The final screened genes are the intersection of the selected genes for both categories. In the end, the number of selected genes is 225.

Figure 5 illustrates the results for the above procedure corresponding to the prostate tumor. In detail, Figure 5(A) reports the histogram of the test statistics $Z(\{(X_{ij}, X_{ik})\})$ for all pairs

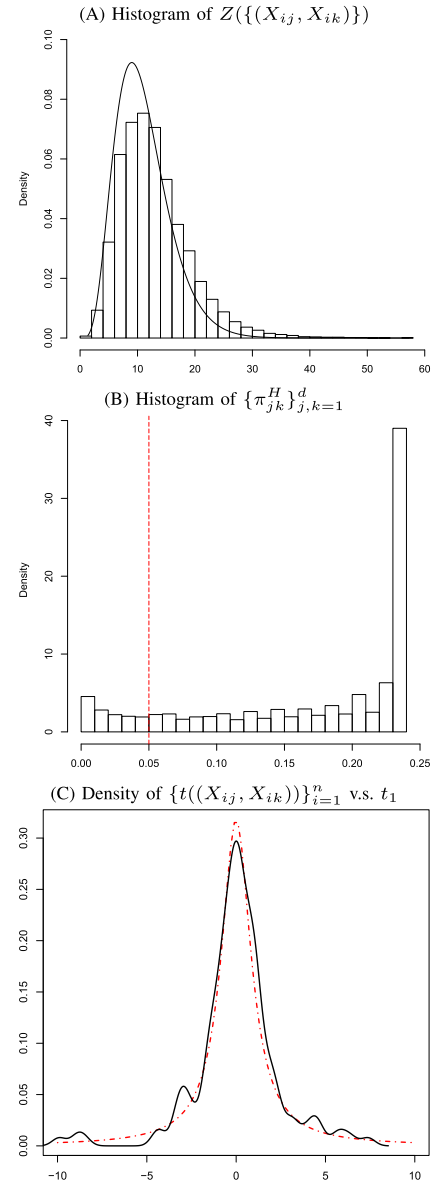


Fig. 5. (A) The histogram of $Z(\{(X_{ij}, X_{ik})\})$ for the pairs of genes for the prostate tumor and the curve is the density function of χ^2_{11} ; (B) The histogram of the Holm's adjusted p-values for the prostate tumor and the red dashed line is the significance level $\alpha = 5\%$; (C) The estimated density of $\{t((X_{ij^*}, X_{ik^*}))\}_{i=1}^n$ (black solid line) versus the density of t_1 (red dashed line) for the pair of genes with largest p-value of the statistic $Z(\{(X_{ij}, X_{ik})\})$.

of genes corresponding to the prostate tumor. The curve in the figure is the density of χ^2 distribution with degrees of freedom $[\sqrt{n}] - 1 = 11$. Figure 5(A) shows that the empirical distribution of the statistics is close to the χ^2 distribution. Figure 5(B) illustrates the histogram of the Holm's adjusted p-values for all pairs of genes. Figure 5(B) shows that we do not reject the goodness of fit test of the pair-elliptical for most pairs of genes in the dataset. We select the pair of genes with the largest p-value and denote the pair as (j^*, k^*) . Figure 5(C) reports the estimated density function of the t-statistics from the selected pair $\{t((X_{ij^*}, X_{ik^*}))\}_{i=1}^n$ versus the density of the t-distribution with degree of freedom 1. The density function is estimated by kernel density estimator with

TABLE II
MEANS (STANDARD DEVIATIONS IN THE PARENTHESES) OF
CLASSIFICATION ERRORS ON GPL96 FOR PROSTATE TUMOR
AND B CELL LYMPHOMA BY APPLYING THE gMAD, gQNE,
PEARSON'S, KENDALL'S TAU AND SPEARMAN'S RHO
COVARIANCE ESTIMATORS TO THE LINEAR
DISCRIMINATIVE CLASSIFIER WITH
100 REPLICATIONS

gMAD	gQNE	Pearson	Kendall	Spearman
0.055	0.059	0.134	0.099	0.098
(0.003)	(0.004)	(0.003)	(0.007)	(0.005)

bandwidth $h = 0.4$. We see that these two distributions are close to each other.

Secondly, we plug the samples of these selected genes to the linear discriminative classifier in (VII.9). To compare the classification errors for different covariance estimators, each time we randomly select 74 samples from the prostate tumor and 74 from the B cell lymphoma as the training dataset. For the training dataset, we divide each category into two parts randomly: One is to derive the classifiers from different covariance estimators by (VII.9) and the other is for tuning the parameters of the first part by minimizing its classification error. The rest samples are applied as the testing dataset to calculate the final classification error. The above steps are repeated for 100 times. We summarize the classification errors in Table II. The errors reported in Table II demonstrate that the gMAD and the gQNE estimators significantly outperform the other estimators in the GPL96 dataset. This indicates the power of quantile-based statistics in high dimensions.

IX. DISCUSSION

This paper studies estimating large scatter matrices for high dimensional data with possibly very heavy tails. We propose a new distribution family, the pair-elliptical, for modeling such data. The pair-elliptical is more flexible than the elliptical. We also characterize the relation between the pair-elliptical and several popular distribution families in high dimensions. Built on the pair-elliptical family, we advocate using the quantile-based approaches for estimating large scatter and covariance matrices. These procedures are both statistically efficient and robust compared to their Gaussian alternatives.

In the future, it is of interest to investigate the performance of the quantile-based methods when the data are not independent. Studies of moment-based and rank-based approaches under high dimensional dependent data include [74] and [75]. However, their proof techniques cannot be directly applied to analyze the quantile-based estimators.

All the results in this paper are confined in parameter estimation and focus on nonasymptotic analysis. In the future, we are also interested in studying the asymptotic properties of quantile-based estimators under the regime that both d and n go to infinity.

APPENDIX

Lemma A.1: For any continuous random variable $X \in \mathbb{R}$ with the distribution function F and any n independent

observations X_1, \dots, X_n of X , we have for any $t > 0$,

$$\begin{aligned} \mathbb{P}(|\hat{Q}(\{X_i\}; r) - Q(X; r)| \geq t) \\ \leq \exp\left(-2n[F\{F^{-1}(r) + t\} - r - n^{-1}]^2\right) \\ + \exp\left(-2n[r - F\{F^{-1}(r) - t\}]^2\right), \end{aligned}$$

whenever $F\{F^{-1}(r) + t\} > r + n^{-1}$.

Proof: Let F_n denote the empirical distribution of X_1, \dots, X_n and $F_n^{-1}(r) := \hat{Q}(\{X_i\}; r)$. By the definition of $\hat{Q}(\cdot; \cdot)$ in (IV.1), we have for any $\epsilon \in [0, 1]$,

$$\epsilon \leq F_n(F_n^{-1}(\epsilon)) \leq \epsilon + \frac{1}{n}.$$

This implies that

$$\begin{aligned} \mathbb{P}\{\hat{Q}(\{X_i\}; r) - Q(X; r) \geq t\} &= \mathbb{P}\{F_n^{-1}(r) - F^{-1}(r) \geq t\} \\ &\leq \mathbb{P}\left[r + \frac{1}{n} \geq F_n\{t + F^{-1}(r)\}\right] \\ &= \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n I\{X_i \leq F^{-1}(r) + t\} \leq r + \frac{1}{n}\right]. \end{aligned}$$

We have

$$\begin{aligned} \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n I\{X_i \leq F^{-1}(r) + t\} \leq r + \frac{1}{n}\right] \\ = \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n I\{X_i \leq F^{-1}(r) + t\} - F\{F^{-1}(r) + t\} \right. \\ \left. \leq r + \frac{1}{n} - F\{F^{-1}(r) + t\}\right]. \end{aligned}$$

Because $\mathbb{E}I\{X_i \leq F^{-1}(r) + t\} = \mathbb{P}(X \leq F^{-1}(r) + t) = F\{F^{-1}(r) + t\}$ and $I\{X_i \leq F^{-1}(r) + t\} \in [0, 1]$ is a bounded random variable, by Hoeffding's inequality, we have

$$\begin{aligned} \mathbb{P}\{\hat{Q}(\{X_i\}; r) - Q(X; r) > t\} \\ \leq \exp\left(-2n[F\{F^{-1}(r) + t\} - r - n^{-1}]^2\right), \end{aligned} \quad (\text{A.1})$$

as long as t is large enough such that $F\{F^{-1}(r) + t\} > r + n^{-1}$.

On the other hand, we have

$$\begin{aligned} \mathbb{P}\{\hat{Q}(\{X_i\}; r) - Q(X; r) \leq -t\} &= \mathbb{P}\{F_n^{-1}(r) - F^{-1}(r) \leq -t\} \\ &\leq \mathbb{P}\left[r \leq F_n\{F^{-1}(r) - t\}\right] \\ &= \mathbb{P}\left[\frac{1}{n} \sum_{i=1}^n I\{X_i \leq F^{-1}(r) - t\} \geq r\right]. \end{aligned}$$

Then again, exploiting the Hoeffding's inequality as above, we have

$$\begin{aligned} \mathbb{P}\{\hat{Q}(\{X_i\}; r) - Q(X; r) \leq -t\} \\ \leq \exp\left(-2n[r - F\{F^{-1}(r) - t\}]^2\right). \end{aligned} \quad (\text{A.2})$$

Combining (A.1) and (A.2), we have the desired result. \square

Lemma A.2 (gMAD Concentration Inequality): Letting $X \in \mathbb{R}$ be a one dimensional continuous random variable.

We denote F_1 and F_2 to be the distribution functions of X and $|X - Q(X; 1/2)|$. Then we have for any $t > 0$,

$$\begin{aligned} & \mathbb{P}\left(|\hat{\sigma}^M(X; r) - \sigma^M(X; r)| > t\right) \\ & \leq 2 \exp\left(-2n\left[F_1\{F_1^{-1}(1/2) + t/2\} - \frac{1}{2} - \frac{1}{n}\right]^2\right) + \\ & \quad 2 \exp\left(-2n\left[\frac{1}{2} - F_1\{F_1^{-1}(1/2) - t/2\}\right]^2\right) + \\ & \quad \exp\left(-2n\left[F_2\{F_2^{-1}(r) + t/2\} - r - \frac{1}{n}\right]^2\right) \\ & \quad + \exp\left(-2n\left[r - F_2\{F_2^{-1}(r) - t/2\}\right]^2\right), \end{aligned}$$

whenever $F_1\{F_1^{-1}(1/2) + t/2\} - 1/2 > 1/n$ and $F_2\{F_2^{-1}(r) + t/2\} - r > 1/n$.

Proof: By definition, we have for any $t > 0$,

$$\begin{aligned} & \mathbb{P}\left(\hat{\sigma}^M(X; r) - \sigma^M(X; r) > t\right) \\ & = \mathbb{P}\left(\hat{Q}\left[\left\{|X_i - \hat{Q}\left\{\{X_i\}; \frac{1}{2}\right\}\right\}; r\right] - Q\left\{X - Q\left(X; \frac{1}{2}\right)\right\}; r\right) > t. \end{aligned}$$

We denote $\nu := Q(X; 1/2)$ and $\hat{\nu} := \hat{Q}(\{X_i\}; 1/2)$. We then have

$$\begin{aligned} & \mathbb{P}\left(\hat{\sigma}^M(X; r) - \sigma^M(X; r) > t\right) \\ & = \mathbb{P}\left\{\hat{Q}\left(\{|X_i - \hat{\nu}|\}, r\right) - Q\left(|X - \nu|, r\right) > t\right\} \\ & \leq \mathbb{P}\left\{\hat{Q}\left(\{|X_i - \nu|\}, r\right) + |\hat{\nu} - \nu| - Q\left(|X - \nu|, r\right) > t\right\} \\ & \leq \underbrace{\mathbb{P}\left\{\hat{Q}\left(\{|X_i - \nu|\}, r\right) - Q\left(|X - \nu|, r\right) > t/2\right\}}_{A_1} \\ & \quad + \underbrace{\mathbb{P}\left(|\hat{\nu} - \nu| > t/2\right)}_B. \end{aligned}$$

On the other hand, we have

$$\begin{aligned} & \mathbb{P}\left(\hat{\sigma}^M(X; r) - \sigma^M(X; r) < -t\right) \\ & = \mathbb{P}\left\{\hat{Q}\left(\{|X_i - \hat{\nu}|\}, r\right) - Q\left(|X - \nu|, r\right) < -t\right\} \\ & \leq \mathbb{P}\left\{\hat{Q}\left(\{|X_i - \nu|\}, r\right) - |\hat{\nu} - \nu| - Q\left(|X - \nu|, r\right) < -t\right\} \\ & \leq \underbrace{\mathbb{P}\left\{\hat{Q}\left(\{|X_i - \nu|\}, r\right) - Q\left(|X - \nu|, r\right) < -t/2\right\}}_{A_2} \\ & \quad + \underbrace{\mathbb{P}\left(|\hat{\nu} - \nu| > t/2\right)}_B. \end{aligned}$$

Using Lemma A.1, we have

$$\begin{aligned} B & \leq \exp\left(-2n\left[F_1\{F_1^{-1}(1/2) + t/2\} - \frac{1}{2} - \frac{1}{n}\right]^2\right) \\ & \quad + \exp\left(-2n\left[\frac{1}{2} - F_1\{F_1^{-1}(1/2) - t/2\}\right]^2\right). \end{aligned}$$

Moreover, we have

$$\begin{aligned} A_1 + A_2 & = \mathbb{P}\left\{\left|\hat{Q}\left(\{|X_i - \nu|\}, r\right) - Q\left(|X - \nu|, r\right)\right| > \frac{t}{2}\right\} \\ & \leq \exp\left(-2n\left[F_2\{F_2^{-1}(r) + t/2\} - r - \frac{1}{n}\right]^2\right) \\ & \quad + \exp\left(-2n\left[r - F_2\{F_2^{-1}(r) - t/2\}\right]^2\right), \end{aligned}$$

where we remind that F_1 and F_2 represent the distribution functions of X and $|X - Q(X; 1/2)|$. Finally, using the fact that

$$\mathbb{P}\left(|\hat{\sigma}^M(X; r) - \sigma^M(X; r)| > t\right) \leq A_1 + A_2 + 2B,$$

we have the desired result. \square

Lemma A.3 (gQNE Concentration): Let $X \in \mathbb{R}$ be a one dimensional continuous random variable and \tilde{X} be an independent copy of X . Let $G(\cdot)$ be the distribution function of $|X - \tilde{X}|$. We then have

$$\mathbb{P}\left(|\hat{\sigma}^Q(X; r) - \sigma^Q(X; r)| > t\right) \leq \exp\left(-n[G\{G^{-1}(r) + t\} - r - n^{-1}]^2\right) + \exp\left(-n[r - G\{G^{-1}(r) - t\}]^2\right),$$

whenever $G\{G^{-1}(r) + t\} > r + n^{-1}$.

Proof: We denote $Z := |X - \tilde{X}|$. By definition, we have

$$\begin{aligned} & \mathbb{P}\left(|\hat{\sigma}^Q(X; r) - \sigma^Q(X; r)| > t\right) \\ & = \mathbb{P}\left(|\hat{Q}(\{|X_i - X_j|_{i < j}\}; r) - Q(Z; r)| > t\right). \end{aligned}$$

Similar as in the proof of Lemma A.1, we have

$$\begin{aligned} & \mathbb{P}\left\{\hat{Q}(\{|X_i - X_j|_{i < j}\}; r) - Q(Z; r) \geq t\right\} \\ & \leq \mathbb{P}\left[\frac{2}{n(n-2)} \sum_{i < j} I\{|X_i - X_j| \leq G^{-1}(r) + t\} \leq r + \frac{1}{n}\right]. \end{aligned}$$

By Hoeffding's inequality in U-statistics (c.f., Equation (5.7) in [76]), we have

$$\begin{aligned} & \mathbb{P}\left\{\hat{Q}(\{|X_i - X_j|_{i < j}\}; r) - Q(Z; r) > t\right\} \\ & \leq \exp\left(-n[G\{G^{-1}(r) + t\} - r - n^{-1}]^2\right), \end{aligned} \quad (\text{A.3})$$

where we remind that G represents the distribution function of Z . Similarly, we have

$$\begin{aligned} & \mathbb{P}\left\{(\hat{Q}(\{|X_i - X_j|_{i < j}\}; r) - Q(Z; r) \leq -t)\right\} \\ & \leq \exp\left(-n[r - G\{G^{-1}(r) - t\}]^2\right). \end{aligned} \quad (\text{A.4})$$

Combining (A.3) and (A.4), we have the desired result. \square

A. Proof of Propositions III.7, III.8, and III.9

Proof of Proposition III.7: Consider a pair-elliptically distributed random vector $\mathbf{X} = (X_1, \dots, X_d)^T \sim PE_d(\boldsymbol{\mu}, \mathbf{S}, \xi_1)$. If \mathbf{X} is also transelliptically distributed, by definition, there exists a set of univariate strictly increasing functions $f = \{f_j\}_{j=1}^d$ such that $f(\mathbf{X}) := (f_1(X_1), \dots, f_d(X_d))^T \sim EC_d(\mathbf{0}, \boldsymbol{\Sigma}^0, \xi_2)$ for some generating variable ξ_2 . Because \mathbf{X} has the same Kendall's tau correlation matrix as $f(\mathbf{X})$, without loss of generality, we can

assume $\mathbf{S} = \Sigma^0$. Accordingly, the margins of \mathbf{X} follow the same distribution. Combined with the fact that the margins of $f(\mathbf{X})$ are identically distributed, we have the transformation functions $f_1 = f_2 = \dots = f_d = f_0$. The desired result follows by considering the following two cases. If $|\Sigma_{jk}^0| = 1$ for all $j, k \in \{1, \dots, d\}$, then we have $X_1 = (\pm)X_2 = \dots = (\pm)X_d$ almost surely. This, combined with the fact that \mathbf{X} are pair-elliptically distributed, implies that \mathbf{X} is elliptically distributed. Otherwise, if there exists $j \neq k$ such that $|\Sigma_{jk}^0| \in (0, 1)$, then without loss of generality, we assume $j = 1, k = 2$ and let $\rho := \Sigma_{12}^0 \neq 0$. Because $(X_1, X_2)^T$ is elliptically distributed, we then have

$$X_2|X_1 = X_1 \sim EC_1(\rho X_1, 1 - \rho^2, \xi'_1),$$

for some generating variable ξ'_1 , which implies that

$$\begin{aligned} \text{median}(X_2|X_1 = X_1) &= \rho X_1 \\ \text{median}(f_0(X_2)|f_0(X_1) = f_0(X_1)) &= f_0(\rho X_1). \end{aligned} \quad (\text{A.5})$$

On the other hand, because $(f_0(X_1), f_0(X_2))^T$ is elliptically distributed, we also have

$$f_0(X_2)|f_0(X_1) = f_0(X_1) \sim EC_1(\rho f_0(X_1), 10\rho^2, \xi'_2),$$

for some generating variable ξ'_2 , which implies that

$$\text{median}(f_0(X_2)|f_0(X_1) = f_0(X_1)) = \rho f_0(X_1). \quad (\text{A.6})$$

Combining Equations (A.5) and (A.6), we have for all $X_1 \in \mathbb{R}$, $f_0(\rho X_1) = \rho f_0(X_1)$, implying that $f_0(x) = ax$ for some $a \in \mathbb{R}$. This further implies that $\mathbf{X} \sim EC_d(\mathbf{0}, a^2 \Sigma^0, \xi_2)$ is elliptically distributed. \square

Proof of Proposition III.8: By definition, any pair in \mathbf{Y} satisfies

$$(Y_j, Y_k)^T \stackrel{D}{=} \xi_G \mathbf{A} \mathbf{U} \sim N_2(\boldsymbol{\mu}_{\{j,k\}}, \mathbf{S}_{\{j,k\}, \{j,k\}}),$$

where $q = \text{rank}(\mathbf{S}_{\{j,k\}, \{j,k\}})$, $\mathbf{A} \in \mathbb{R}^{2 \times q}$ with $\mathbf{A} \mathbf{A}^T = \mathbf{S}_{\{j,k\}, \{j,k\}}$, $\mathbf{U} \in \mathbb{R}^q$ uniformly distributed in \mathbb{S}^{q-1} , and ξ_G^2 follows a chi square distribution with degree of freedom q . This further implies that any pair in \mathbf{X} satisfies

$$(X_j, X_k)^T = \xi(Y_j, Y_k)^T \sim \xi \cdot \xi_G \cdot \mathbf{A} \mathbf{U},$$

and accordingly is elliptically distributed. This verifies that \mathbf{X} is pair-elliptically distributed. \square

Proof of Proposition III.9: Because the margin of a pair-normal is normally distributed and the only elliptical distribution that has normal margins is the Gaussian, we have the first assertion is true. We then prove the second assertion. Suppose that $\mathbf{Y} = (Y_1, \dots, Y_d)^T$ is pair-normally as well as nonparanormlly distributed. Then for any $j \in \{1, \dots, d\}$, Y_j is normally distributed. Moreover, because \mathbf{Y} is nonparanormlly distributed, we have $f_j(Y_j) \sim N(0, 1)$. This implies that f_j is a linear transformation. Therefore, \mathbf{Y} is Gaussian distributed because f_1, \dots, f_d are linear. \square

B. Proofs of Lemma IV.1

Proof: Since $\tilde{\mathbf{R}}$ is the minimizer of (IV.4), we have

$$\|\tilde{\mathbf{R}} - \hat{\mathbf{R}}\| \leq \|\tilde{\mathbf{R}} - \hat{\mathbf{R}}\|,$$

which implies $\mathbb{P}(\|\tilde{\mathbf{R}} - \mathbf{R}\| \geq t) \leq \mathbb{P}(\|\tilde{\mathbf{R}} - \hat{\mathbf{R}}\| + \|\hat{\mathbf{R}} - \mathbf{R}\| \geq t) \leq \mathbb{P}(\|\hat{\mathbf{R}} - \mathbf{R}\| \geq t)$. \square

C. Proofs of Theorems II.1

Proof of Theorem II.1: Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be n independent observations of $\mathbf{X} \in \mathbb{R}^d$ with $\mathbf{X}_i = (X_{i1}, \dots, X_{id})^T$. By Chebychev's inequality, we have

$$\mathbb{P}(|X_{ij}| \geq t) \leq (t/\|X_{ij}\|_{L_p})^{-p}.$$

Let's cut X_{ij} into two parts:

$$Y_{ij} = X_{ij} \mathbb{1}(|X_{ij}| \geq \epsilon) \quad \text{and} \quad Y_{ij}^* = X_{ij} \mathbb{1}(|X_{ij}| \leq \epsilon).$$

We have $|Y_{ij}^* - \mathbb{E}Y_{ij}^*|$ is upper bounded by 2ϵ and its variance has the property:

$$\text{Var}(Y_{ij}^*) \leq \mathbb{E}(Y_{ij}^*)^2 \leq \|X_{ij}\|_{L_2}^2 \leq \|X_{ij}\|_{L_p}^2 = K^2.$$

Accordingly, by Bernstein's inequality, we have

$$\mathbb{P}\left(\frac{1}{n} \left(\sum Y_{ij}^* - \mathbb{E}Y_{ij}^*\right) > t/2\right) \leq \exp\left(-\frac{1/8 \cdot nt^2}{K^2 + 1/3 \cdot \epsilon t}\right).$$

For Y_{ij} , we have

$$\mathbb{P}(Y_{ij} \neq 0) = \mathbb{P}(|X_{ij}| \geq \epsilon) \leq (\epsilon/K)^{-p},$$

and

$$|\mathbb{E}Y_{ij}| = |\mathbb{E}X_{ij} \mathbb{1}(|X_{ij}| \geq \epsilon)| \leq \frac{\mathbb{E}|X_{ij}|^p}{\epsilon^{p-1}} = \frac{K^p}{\epsilon^{p-1}}.$$

Accordingly, for any ϵ such that

$$K^p/\epsilon^{p-1} \leq t/2, \quad (\text{A.7})$$

we have for any $j \in \{1, \dots, d\}$,

$$\begin{aligned} &\mathbb{P}\left(\left|\frac{1}{n} \sum X_{ij} - \mu_j\right| > t\right) \\ &\leq \mathbb{P}\left(\left|\frac{1}{n} \sum (Y_{ij}^* - \mathbb{E}Y_{ij}^*)\right| > t/2\right) + \mathbb{P}\left(\left|\frac{1}{n} \sum (Y_{ij} - \mathbb{E}Y_{ij})\right| > t/2\right) \\ &\leq 2 \exp\left(-\frac{1/8 \cdot nt^2}{K^2 + 1/3 \cdot \epsilon t}\right) + n\mathbb{P}(Y_{ij} \neq 0) \\ &\leq 2 \exp\left(-\frac{1/8 \cdot nt^2}{K^2 + 1/3 \cdot \epsilon t}\right) + n(\epsilon/K)^{-p}. \end{aligned}$$

This yields that

$$\begin{aligned} \mathbb{P}(\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_\infty > t) &\leq d \cdot \mathbb{P}\left(\left|\frac{1}{n} \sum X_{ij} - \mu_j\right| > t\right) \\ &\leq d \cdot \left(2 \exp\left(-\frac{1/8 \cdot nt^2}{K^2 + 1/3 \cdot \epsilon t}\right) + n(\epsilon/K)^{-p}\right). \end{aligned}$$

Taking $t = 12K \cdot \sqrt{\log d/n}$ and $\epsilon = K \cdot \sqrt{n/(\log d)}$, as long as $\log d/n = o(1)$, we have

$$2d \exp\left(-\frac{1/8 \cdot nt^2}{K^2 + 1/3 \cdot \epsilon t}\right) \leq 2d^{-2.5}.$$

It is also straightforward to verify that for such chosen t and ϵ , we have $K^p/\epsilon^{p-1} \leq t/2$ as long as $p \geq 2$.

For the second term, we have

$$\begin{aligned} \log(n\epsilon^{-p} \cdot K^p) &= \log n - p \log \epsilon + p \log K \\ &= \log n - p(\log K + 1/2 \log n - 1/2 \log \log d) + p \log K \\ &= (-p/2 + 1) \log n + p/2 \log \log d. \end{aligned}$$

So for large enough p and K , we need to have

$$\log d + (-p/2 + 1) \log n + p/2 \log \log d \rightarrow -\infty.$$

Thus, letting $p = 2 + 2\gamma + \delta$, for $d = O(n^\gamma)$, with probability no smaller than $1 - 2d^{-2.5} - (\log d)^{p/2}n^{-\delta/2}$, we have

$$\|\bar{\mu} - \mu\|_\infty \leq 12K \sqrt{\frac{\log d}{n}}.$$

This completes the proof. \square

Proof of Theorem II.2: Consider the distribution $\mathbf{X} = (X_1, \dots, X_d)$ with X_1, \dots, X_d independent and identically distributed. For reasons that will be clear later, we can set $C = 1$ and consider

$$\begin{aligned} \mathbb{P}(X_1 = \sqrt{n \log d}) &= \frac{1}{2}(n \log d)^{-p/2}, \\ \mathbb{P}(X_1 = -\sqrt{n \log d}) &= \frac{1}{2}(n \log d)^{-p/2}, \\ \text{and } \mathbb{P}(X_1 = 0) &= 1 - (n \log d)^{-p/2}. \end{aligned} \quad (\text{A.8})$$

We then have $\mu = \mathbf{0}$ and $\|X_1\|_{L_q} \rightarrow I(q = p) + \infty \cdot I(q > p)$. Moreover, letting

$$\begin{aligned} \alpha &:= \mathbb{P}(|\bar{\mu}_1| \geq \sqrt{\log d/n}) \\ &\geq 1 - \mathbb{P}(X_1 = \dots = X_n = 0) \\ &\geq n(n \log d)^{-p/2}(1 + o(1)), \end{aligned} \quad (\text{A.9})$$

we have

$$\begin{aligned} \mathbb{P}(\|\bar{\mu} - \mu\|_\infty \geq \sqrt{\log d/n}) &= d\alpha - \binom{d}{2}\alpha^2 + \binom{d}{3}\alpha^3 - \dots \\ &= 1 - (1 - \alpha)^d. \end{aligned}$$

When $p = 2 + 2\gamma$ and $d = n^{\gamma+\delta_0}$ with $\delta_0 > 0$ and some $0 < \delta_1 < \delta_0$, we have

$$\begin{aligned} (1 - \alpha)^d &\leq \left(1 - n^{-\gamma}(\log d)^{-p/2}\right)^{n^{\gamma+\delta_0}} \\ &\leq \left(1 - \frac{1}{n^{\gamma+\delta_1}}\right)^{n^{\gamma+\delta_0}} \leq (1/e)^{n^{\delta_0-\delta_1}} \rightarrow 0. \end{aligned}$$

Accordingly, with probability tending to 1, we have

$$\|\bar{\mu} - \mu\|_\infty \geq \sqrt{\log d/n}.$$

When $C \neq 1$, we can replace all terms $\sqrt{n \log d}$ in Equation (A.8) with $\sqrt{Cn \log d}$ and all the proofs can follow. \square

D. Proofs of Theorems V.1 and V.2

Proof of Theorem V.1: For $j = 1, \dots, d$, we denote $\mu_j := Q(X_j; 1/2)$. By definition, we have for $j = 1, \dots, d$,

$$\sigma^M(X_j; r) = Q(|X_j - Q(X_j; 1/2)|; r) = Q(|X_j - \mu_j|; r).$$

In other words, we have

$$\begin{aligned} r &= \mathbb{P}\{|X_j - \mu_j| \leq \sigma^M(X_j; r)\} \\ &= \mathbb{P}\left\{\frac{|X_j - \mu_j|}{\sigma(X_j)} \leq \frac{\sigma^M(X_j; r)}{\sigma(X_j)}\right\}, \end{aligned}$$

because all the quantiles are unique. Accordingly, we have

$$\begin{aligned} \frac{\sigma^M(X_j; r)}{\sigma(X_j)} &= Q\left(\frac{|X_j - \mu_j|}{\sigma(X_j)}; r\right) \\ \Rightarrow \sigma^M(X_j; r) &= \sigma(X_j) \cdot Q\left(\frac{|X_j - \mu_j|}{\sigma(X_j)}; r\right). \end{aligned}$$

Using a similar technique as above, we can further derive that

$$\begin{aligned} \sigma^M(X_j + X_k; r) &= \sigma(X_j + X_k) \cdot Q\left(\frac{|X_j + X_k - Q(X_j + X_k; 1/2)|}{\sigma(X_j + X_k)}; r\right), \\ \sigma^M(X_j - X_k; r) &= \sigma(X_j - X_k) \cdot Q\left(\frac{|X_j - X_k - Q(X_j - X_k; 1/2)|}{\sigma(X_j - X_k)}; r\right). \end{aligned}$$

Therefore, $\mathbf{R}^{M;r} = c^{M;r} \Sigma$ if (V.1) holds and by definition

$$\begin{aligned} c^{M;r} &= \left\{ \frac{\sigma^M(X_j; r)}{\sigma(X_j)} \right\}^2 \\ &= \left\{ Q\left(\frac{|X_j - \mu_j|}{\sigma(X_j)}; r\right) \right\}^2. \end{aligned}$$

This completes the proof. \square

Proof of Theorem V.2: The proof of Theorem V.2 is similar as Theorem V.1, and is accordingly omitted here. \square

E. Proof of Theorem V.3

Proof: We first prove the case for the gMAD estimator $\mathbf{R}^{M;r}$. By the definition of the pair-elliptical and the discussions in Remark III.2, we have any pair in $PE_d(\mu, \mathbf{S}, \xi)$ can be written as $EC_2(\mu_{\{j,k\}}, \mathbf{S}_{\{j,k\}, \{j,k\}}, \phi)$ where ϕ is a properly defined characteristic function only depending on ξ . On one hand, by Theorem 2.16 in [41], we have for any $j \in \{1, \dots, d\}$,

$$\frac{X_j - \mu_j}{\sqrt{\mathbf{S}_{jj}}} \sim EC_1(0, 1, \phi). \quad (\text{A.10})$$

Because \mathbf{X} is continuous (by the definition of the pair-elliptical) and accordingly all quantiles of $EC_d(0, 1, \xi)$ are uniquely defined, using the same proof techniques as exploited in the proof of Theorem V.1, (A.10) shows that

$$\mathbf{R}_{jj}^{M;r} = \left(Q\left(\left| EC_1(0, 1, \phi) \right|; r \right) \right)^2 \cdot \mathbf{S}_{jj}$$

for any $j \in \{1, \dots, d\}$. On the other hand, for those $j \neq k$ such that $\mathbf{S}_{jj} + \mathbf{S}_{kk} + 2\mathbf{S}_{jk} \neq 0$ and $\mathbf{S}_{jj} + \mathbf{S}_{kk} - 2\mathbf{S}_{jk} \neq 0$, we have

$$\begin{aligned} X_j + X_k &\sim EC_d(\mu_j + \mu_k, \mathbf{S}_{jj} + \mathbf{S}_{kk} + 2\mathbf{S}_{jk}, \phi), \\ \text{and } X_j - X_k &\sim EC_d(\mu_j - \mu_k, \mathbf{S}_{jj} + \mathbf{S}_{kk} - 2\mathbf{S}_{jk}, \phi), \end{aligned}$$

and accordingly $X_j + X_k$ and $X_j - X_k$ are both elliptically distributed with the same characterization function ϕ . Therefore, we have

$$\begin{aligned} &\frac{X_j + X_k - Q(X_j + X_k; 1/2)}{\sqrt{\mathbf{S}_{jj} + \mathbf{S}_{kk} + 2\mathbf{S}_{jk}}} \\ &\stackrel{D}{=} \frac{X_j - X_k - Q(X_j - X_k; 1/2)}{\sqrt{\mathbf{S}_{jj} + \mathbf{S}_{kk} - 2\mathbf{S}_{jk}}} \sim EC_1(0, 1, \phi), \end{aligned}$$

and accordingly we have

$$\begin{aligned} \left| \frac{X_j - \mu_j}{\sqrt{\mathbf{S}_{jj}}} \right| &\stackrel{D}{=} \left| \frac{X_j + X_k - Q(X_j + X_k; 1/2)}{\sqrt{\mathbf{S}_{jj} + \mathbf{S}_{kk} + 2\mathbf{S}_{jk}}} \right| \\ &\stackrel{D}{=} \left| \frac{X_j - X_k - Q(X_j - X_k; 1/2)}{\sqrt{\mathbf{S}_{jj} + \mathbf{S}_{kk} - 2\mathbf{S}_{jk}}} \right|. \end{aligned}$$

Thus, in parallel to the proof in Theorem V.1, we have

$$\mathbf{R}_{jk}^{M;r} = \left(Q\left(\left| EC_1(0, 1, \phi) \right|; r\right) \right)^2 \cdot \mathbf{S}_{jk}.$$

For those $j \neq k \in \{1, \dots, d\}$ such that $\mathbf{S}_{jj} + \mathbf{S}_{kk} - 2\mathbf{S}_{jk} = 0$, we have $\text{Cov}(X_j - X_k) = 0$, implying that $X_j = X_k$, a.s.. Accordingly, we have

$$\sigma^M(X_j, X_k; r) = \frac{1}{4} \{ \sigma^M(X_j + X_k; r) \}^2 = (\sigma^M(X_j + X_k; r))^2.$$

Accordingly, $\mathbf{R}_{jk}^{M;r} = \mathbf{R}_{jj}^{M;r} = \left(Q\left(\left| EC_1(0, 1, \phi) \right|; r\right) \right)^2 \cdot \mathbf{S}_{jj} = \left(Q\left(\left| EC_1(0, 1, \phi) \right|; r\right) \right)^2 \cdot \mathbf{S}_{jk}$. Similarly, for those $j \neq k \in \{1, \dots, d\}$ such that $\mathbf{S}_{jj} + \mathbf{S}_{kk} + 2\mathbf{S}_{jk} = 0$, we have $\mathbf{R}_{jk}^{M;r} = -\mathbf{R}_{jj}^{M;r} = -\left(Q\left(\left| EC_1(0, 1, \phi) \right|; r\right) \right)^2 \cdot \mathbf{S}_{jj} = \left(Q\left(\left| EC_1(0, 1, \phi) \right|; r\right) \right)^2 \cdot \mathbf{S}_{jk}$. This completes the proof of the first part.

Secondly, we switch to the gQNE estimator. We note that, by Remark III.2,

$$\mathbb{E} \exp(it^T (\mathbf{X}_{\{j,k\}} - \widetilde{\mathbf{X}}_{\{j,k\}})) = \phi^2(t^T \mathbf{S}_{\{j,k\}, \{j,k\}} t),$$

and accordingly $\mathbf{X} - \widetilde{\mathbf{X}} \sim PE_d(\mathbf{0}, 2\mathbf{S}, \phi^2)$ and is continuously distributed. Then by following the same proof as in the first part we have the desired result.

In the end, let's show the scale constant. Using the same argument as above, we have

$$\mathbf{R}^{M;r} = \left(Q\left(\left| X_0 \right|; r\right) \right)^2 \cdot \Sigma.$$

Because X_0 is continuous and symmetric, letting $q_r := Q\left(\left| X_0 \right|; r\right)$, we have

$$\begin{aligned} \mathbb{P}(|X_0| \leq q_r) &= 2\mathbb{P}(X_0 \leq q_r) - 1 = r \\ \Rightarrow \mathbb{P}(X_0 \leq q_r) &= \frac{1+r}{2} \Rightarrow q_r = Q\left(\left| X_0 \right|; \frac{1+r}{2}\right). \end{aligned}$$

Similarly, we have

$$\mathbf{R}^{Q;r} = \left(Q\left(\left| Z_0 \right|; \frac{1+r}{2}\right) \right)^2 \cdot \Sigma.$$

This finalizes the proof. \square

F. Proof of Theorem V.5

Proof: Using Lemma A.2 and Assumption (A1), we have for any $j \in \{1, \dots, d\}$,

$$\begin{aligned} \mathbb{P}(|\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r)| > t) \\ \leq 3 \exp(-2n(\eta_1 t/2 - 1/n)^2) + 3 \exp(-2n(\eta_1 t/2)^2), \end{aligned}$$

whenever $t/2 \leq \kappa_1$ and $\eta_1 t/2 > 1/n$. Accordingly, we have

$$\begin{aligned} \mathbb{P}(|\hat{\mathbf{R}}_{jj}^{M;r} - \mathbf{R}_{jj}^{M;r}| > t) \\ &= \mathbb{P}(|(\hat{\sigma}^M(X_j; r) + \sigma^M(X_j; r))(\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r))| > t) \\ &\leq \mathbb{P}(|(\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r))^2 + 2\sigma^M(X_j; r)(\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r))| > t) \\ &\leq \mathbb{P}\left(|\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r)| > \sqrt{\frac{t}{2}}\right) \end{aligned}$$

$$\begin{aligned} &+ \mathbb{P}\left(|\hat{\sigma}^M(X_j; r) - \sigma^M(X_j; r)| > \frac{t}{2\sigma^M(X_j; r)}\right) \\ &\leq 6 \exp(-2n(\eta_1 \sqrt{t/2}/2 - 1/n)^2) \\ &+ 6 \exp(-2n(\eta_1 t/(4\sigma^M(X_j; r)))^2). \end{aligned} \quad (\text{A.11})$$

Using Lemma A.2 and Assumption (A1), we have for any $j \neq k$,

$$\begin{aligned} \mathbb{P}(|\hat{\sigma}^M(X_j + X_k; r) - \sigma^M(X_j + X_k; r)| > t) \\ \leq 3 \exp(-2n(\eta_1 t/2 - 1/n)^2) + 3 \exp(-2n(\eta_1 t/2)^2), \\ \mathbb{P}(|\hat{\sigma}^M(X_j - X_k; r) - \sigma^M(X_j - X_k; r)| > t) \\ \leq 3 \exp(-2n(\eta_1 t/2 - 1/n)^2) + 3 \exp(-2n(\eta_1 t/2)^2). \end{aligned}$$

And accordingly, letting $\theta_{\max} := 2\sqrt{\|\mathbf{R}^{M;r}\|_{\max}}$, we have

$$\begin{aligned} \mathbb{P}(|\hat{\mathbf{R}}_{jk}^{M;r} - \mathbf{R}_{jk}^{M;r}| > t) \\ \leq \mathbb{P}(|(\hat{\sigma}^M(X_j + X_k; r))^2 - (\sigma^M(X_j + X_k; r))^2| > 2t) \\ + \mathbb{P}(|(\hat{\sigma}^M(X_j - X_k; r))^2 - (\sigma^M(X_j - X_k; r))^2| > 2t) \\ \leq \mathbb{P}\left(|\hat{\sigma}^M(X_j + X_k; r) - \sigma^M(X_j + X_k; r)| > \sqrt{t}\right) \\ + \mathbb{P}\left(|\hat{\sigma}^M(X_j + X_k; r) - \sigma^M(X_j + X_k; r)| > \frac{t}{\sigma^M(X_j + X_k; r)}\right) \\ + \mathbb{P}\left(|\hat{\sigma}^M(X_j - X_k; r) - \sigma^M(X_j - X_k; r)| > \sqrt{t}\right) \\ + \mathbb{P}\left(|\hat{\sigma}^M(X_j - X_k; r) - \sigma^M(X_j - X_k; r)| > \frac{t}{\sigma^M(X_j - X_k; r)}\right) \\ \leq 6 \exp(-2n(\eta_1 \sqrt{t}/2 - 1/n)^2) + 6 \exp(-n\eta_1^2 t/2) \\ + 6 \exp(-2n(\eta_1 t/(2\theta_{\max}) - 1/n)^2) + 6 \exp(-n\eta_1^2 t/(2\theta_{\max}^2)) \\ \leq 24 \max\{\exp(-2n(\eta_1 \sqrt{t}/2 - 1/n)^2), \\ \exp(-2n(\eta_1 t/(2\theta_{\max}) - 1/n)^2)\}. \end{aligned} \quad (\text{A.12})$$

Combining Equations (A.11) and (A.12), we have, with probability $1 - 24\alpha^2$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{M;r} - \mathbf{R}^{M;r}\|_{\max} &\leq \\ \max \left\{ \underbrace{\frac{6}{\eta_1^2} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2}_{T_1}, \right. \\ \left. \underbrace{\frac{4\sqrt{\|\mathbf{R}^{M;r}\|_{\max}}}{\eta_1} \left(\sqrt{\frac{\log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)}_{T_2} \right\}, \end{aligned}$$

whenever n is large enough such that $T_1 \leq 8\kappa_1^2$ and $T_2 \leq 2\kappa_1 \cdot \min_{j \neq k} \{2\sigma^M(X_j; r), \sigma^M(X_j + X_k; r), \sigma^M(X_j - X_k; r)\}$. Combining the above inequality with Theorem V.3, we complete the whole proof. \square

G. Proof of Theorem V.6

Proof: Using Lemma A.3 and Assumption (A2), we have for any $j \in \{1, \dots, d\}$,

$$\begin{aligned} \mathbb{P}(|\hat{\sigma}^Q(X; r) - \sigma^Q(X; r)| > t) &\leq \exp\left(-n[\eta_2 t - 1/n]^2\right) \\ &+ \exp\left(-n(\eta_2 t)^2\right), \end{aligned}$$

whenever $t \leq \kappa_2$ and $\eta_2 t > 1/n$. Using a similar proof technique as in the proof of Theorem V.5, we have

$$\begin{aligned} \mathbb{P}(|\hat{\mathbf{R}}_{jj}^{Q;r} - \mathbf{R}_{jj}^{Q;r}| > t) &\leq \mathbb{P}\left(|\hat{\sigma}^Q(X_j; r) - \sigma^Q(X_j; r)| > \sqrt{t/2}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}^Q(X_j; r) - \sigma^Q(X_j; r)| > \frac{t}{2\sigma^Q(X_j; r)}\right) \\ &\leq 2 \exp\left(-n[\eta_2 \sqrt{t/2} - 1/n]^2\right) \end{aligned} \quad (\text{A.13})$$

$$+ 2 \exp\left(-n(\eta_2 t / (2\sigma^Q(X_j; r)) - 1/n)^2\right). \quad (\text{A.14})$$

Similarly, we have

$$\begin{aligned} \mathbb{P}(|\hat{\mathbf{R}}_{jk}^{Q;r} - \mathbf{R}_{jk}^{Q;r}| > t) &\leq \mathbb{P}\left(|\hat{\sigma}^Q(X_j + X_k; r) - \sigma^Q(X_j + X_k; r)| > \sqrt{t}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}^Q(X_j - X_k; r) - \sigma^Q(X_j - X_k; r)| > \sqrt{t}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}^Q(X_j + X_k; r) - \sigma^Q(X_j + X_k; r)| > \frac{t}{\sigma^Q(X_j + X_k; r)}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}^Q(X_j - X_k; r) - \sigma^Q(X_j - X_k; r)| > \frac{t}{\sigma^Q(X_j - X_k; r)}\right) \\ &\leq 4 \exp(-n[\eta_2 \sqrt{t} - 1/n]^2) + 4 \exp(-n(\eta_2 t / \zeta_{\max} - 1/n)^2), \end{aligned} \quad (\text{A.15})$$

where $\zeta_{\max} := 2\sqrt{\|\mathbf{R}^{Q;r}\|_{\max}}$. Combining (A.13) and (A.15) leads to that, with probability larger than or equal to $1 - 8\alpha$,

$$\begin{aligned} \|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} &\leq \max \left\{ \underbrace{\frac{2}{\eta_2^2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)^2}_{T_3}, \right. \\ &\quad \left. \underbrace{\frac{2\sqrt{\|\mathbf{R}^{Q;r}\|_{\max}}}{\eta_2} \left(\sqrt{\frac{2 \log d + \log(1/\alpha)}{n}} + \frac{1}{n} \right)}_{T_4} \right\}, \end{aligned}$$

whenever $T_3 \leq 2\kappa_2^2$ and $T_4 \leq 2\kappa_1 \cdot \min_{j \neq k} \{2\sigma^Q(X_j; r), \sigma^Q(X_j + X_k), \sigma^Q(X_j - X_k)\}$. Finally, combining the above inequality with Theorem V.3, we complete the whole proof. \square

In this section we provide the proofs of the results presented in Section VII.

H. Proof of Corollary VII.1

Proof: We first prove that (VII.4) holds. Because $\mathbf{R}^{Q;r}$ is feasible to Equation (VII.2), we have

$$\|\hat{\mathbf{R}}^{Q;r} - \tilde{\mathbf{R}}^{Q;r}\|_{\max} \leq \|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max},$$

implying that

$$\begin{aligned} \mathbb{P}\left(\left|\tilde{\mathbf{R}}_{jk}^{Q;r} - \mathbf{R}_{jk}^{Q;r}\right| \geq t\right) &\leq \mathbb{P}\left(\left|\tilde{\mathbf{R}}_{jk}^{Q;r} - \hat{\mathbf{R}}_{jk}^{Q;r}\right| + \left|\hat{\mathbf{R}}_{jk}^{Q;r} - \mathbf{R}_{jk}^{Q;r}\right| \geq t\right) \\ &\leq \mathbb{P}\left(\|\tilde{\mathbf{R}}^{Q;r} - \hat{\mathbf{R}}^{Q;r}\|_{\max} + \|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \geq t\right) \\ &\leq \mathbb{P}(\|\hat{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \geq t/2). \end{aligned}$$

Combined with (A.13) and (A.15), the above inequality implies that

$$\begin{aligned} \mathbb{P}\left(\|\tilde{\mathbf{R}}_{jk}^{Q;r} - \mathbf{R}_{jk}^{Q;r}\|_{\max} \geq t\right) &\leq d^4 (4 \exp(-n[\eta_2 \sqrt{t/4} - 1/n]^2) \\ &\quad + 4 \exp(-n(\eta_2 t / (2\zeta_{\max}) - 1/n)^2)). \end{aligned}$$

Plugging $t = \zeta$ into the above equation, we have the desired result.

Secondly, we prove that (VII.5) holds. Because $\mathbb{E}X_j^4 \leq K$, by Chebyshev's inequality, with probability no smaller than $1 - n^{-2\xi}$, we have

$$|\hat{\sigma}_j^2 - \sigma^2(X_j)| \leq c_1 n^{-1/2+\xi}.$$

Moreover, using (A.13), we have for any given $j \in \{1, \dots, d\}$, by Markov inequality, with probability larger than or equal to $1 - 4\alpha$,

$$|\hat{\mathbf{R}}_{jj}^{Q;r} - \mathbf{R}_{jj}^{Q;r}| \leq \zeta_j.$$

For notation simplicity, we denote $\sigma_j := \sigma_j(X_j)$, $r_j := \mathbf{R}_{jj}^{Q;r}$, and $\hat{r}_j := \hat{\mathbf{R}}_{jj}^{Q;r}$. Accordingly, we have

$$\begin{aligned} \|\hat{\Sigma}^{Q;r} - \Sigma\|_{\max} &= \left\| \frac{\hat{\sigma}_j^2}{\hat{r}_j} \tilde{\mathbf{R}}^{Q;r} - \frac{\sigma_j^2}{r_j} \mathbf{R}^{Q;r} \right\|_{\max} \\ &\leq \left\| \frac{\hat{\sigma}_j^2}{\hat{r}_j} \tilde{\mathbf{R}}^{Q;r} - \frac{\sigma_j^2}{r_j} \tilde{\mathbf{R}}^{Q;r} \right\|_{\max} + \frac{\sigma_j^2}{r_j} \|\tilde{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \\ &\leq \left| \frac{\hat{\sigma}_j^2}{\hat{r}_j} - \frac{\sigma_j^2}{r_j} \right| \cdot (\|\tilde{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \\ &\quad + \|\mathbf{R}^{Q;r}\|_{\max}) + \frac{1}{c^{Q;r}} \|\tilde{\mathbf{R}}^{Q;r} - \mathbf{R}^{Q;r}\|_{\max} \end{aligned}$$

while noting that $c^{Q;r} = r_j / \sigma_j^2$. Finally, we have

$$\begin{aligned} \left| \frac{\hat{\sigma}_j^2}{\hat{r}_j} - \frac{\sigma_j^2}{r_j} \right| &= \left| \frac{\hat{\sigma}_j^2 - \sigma_j^2}{\hat{r}_j} + \sigma_j^2 \left(\frac{1}{\hat{r}_j} - \frac{1}{r_j} \right) \right| \\ &\leq \frac{|\hat{\sigma}_j^2 - \sigma_j^2|}{r_j - |\hat{r}_j - r_j|} + \frac{\sigma_j^2}{r_j} \cdot \frac{|\hat{r}_j - r_j|}{r_j - |\hat{r}_j - r_j|}. \end{aligned}$$

This implies that with probability no smaller than $1 - n^{-2\xi} - 12\alpha$, we have

$$\left| \frac{\hat{\sigma}_j^2}{\hat{r}_j} - \frac{\sigma_j^2}{r_j} \right| \leq \frac{c_1 n^{-1/2+\xi}}{r_j - \zeta_j} + \frac{1}{c^{Q;r}} \cdot \frac{\zeta_j}{r_j - \zeta_j},$$

which completes the proof of the second part. The third part can then be proved by combining Equation (VII.5) and the proof of Theorem 2 in [56]. \square

I. Proof of Corollary VII.4

Proof: According to Theorem 6 in [40], if the tuning parameter $\lambda \geq \|\Theta\|_{1,\infty} \|\hat{\Sigma}^{Q;r} - \Sigma\|_{\max}$, we have there exists two constants C_1, C_2 such that

$$\|\hat{\Theta}^{Q;r} - \Theta\|_2 \leq C_1 s \lambda^{1-q}, \|\hat{\Theta}^{Q;r} - \Theta\|_{\max} \leq \|\Theta\|_{1,\infty} \lambda.$$

and $\|\hat{\Theta}^{Q;r} - \Theta\|_F^2 \leq C_2 s d \lambda^{2-q}$. Combining the above results with Corollary VII.1, we prove the desire rate. \square

J. Proof of Corollary VII.5

Proof: According to Theorem 3.1 in [67], if $\lambda \geq \|\hat{\mathbf{R}} - \mathbf{R}\|_{\max}$ and $\max_{1 \leq k \leq m} \|\mathbf{u}_k(\boldsymbol{\Sigma})\|_0 \leq s$, we have

$$\|\hat{\boldsymbol{\Pi}}_m^{Q;r} - \boldsymbol{\Pi}_m\|_F \leq \frac{4s\lambda}{\lambda_1(\mathbf{R}) - \lambda_2(\mathbf{R})} \cdot \|\hat{\mathbf{R}} - \mathbf{R}\|_{\max}. \quad (\text{A.16})$$

where $\boldsymbol{\Pi}_m^{Q;r}$, \mathbf{R} , and $\hat{\mathbf{R}}$ can be the Q-PCA estimator, population and sample versions of the quantile-based scatter matrix estimators based on the gQNE. Combining the rate of $\|\hat{\mathbf{R}} - \mathbf{R}\|_{\max}$ in Theorem V.5 and Theorem V.6, the theorem is proved. \square

K. Proof of Corollary VII.6

Proof: We mainly follow the procedure in the proof of Theorem 2 in [71]. Suppose $\delta^T \beta^* \leq M$. Let $\delta = \mu_1 - \mu_2$, $\hat{\delta} = \hat{\mu}_1 - \hat{\mu}_2$. Since $(\hat{\beta}^{Q;r})^{Q;r}$ is the solution of (VII.9), we have

$$|(\beta^*)^T \hat{\Sigma}^{Q;r} \hat{\beta}^{Q;r} - (\beta^*)^T \delta| \leq (\lambda + \|\hat{\delta} - \delta\|_{\infty}) \|\beta^*\|_1 \leq 2\lambda \|\beta^*\|_1,$$

and

$$|(\beta^*)^T \hat{\Sigma}^{Q;r} \hat{\beta}^{Q;r} - (\hat{\beta}^{Q;r})^T \delta| \leq (\lambda + \|\hat{\delta} - \delta\|_{\infty}) \|\hat{\beta}^{Q;r}\|_1.$$

Combining the above two equations together, we have $|(\hat{\beta}^{Q;r} - \beta^*)^T \delta| \leq 4\lambda \|\beta^*\|_1$ which implies

$$\frac{((\hat{\beta}^{Q;r})^T \delta)^2}{(\delta^T \beta^*)^2} \geq 1 - 8M^{-1} \lambda \|\beta^*\|_1. \quad (\text{A.17})$$

On the other hand, we have

$$\begin{aligned} \|\Sigma \hat{\beta}^{Q;r} - \delta\|_{\infty} &\leq \|(\Sigma - \hat{\Sigma}^{Q;r}) \hat{\beta}^{Q;r}\|_{\infty} + 2\lambda \\ &\leq \|\beta^*\|_1 \|\Sigma - \hat{\Sigma}^{Q;r}\|_{\max} + 2\lambda, \end{aligned}$$

which implies that

$$\begin{aligned} &|(\hat{\beta}^{Q;r})^T \Sigma \hat{\beta}^{Q;r} - \delta^T \beta^*| \\ &\leq |(\hat{\beta}^{Q;r})^T \Sigma \hat{\beta}^{Q;r} - \delta^T \hat{\beta}^{Q;r}| + |(\hat{\beta}^{Q;r} - \beta^*)^T \delta| \\ &\leq \|\beta^*\|_1^2 \|\Sigma - \hat{\Sigma}^{Q;r}\|_{\max} + 6\lambda \|\beta^*\|_1, \end{aligned} \quad (\text{A.18})$$

where the last inequality is due to the definition of $\hat{\beta}^{Q;r}$. Therefore, for sufficiently large n , we have $(\hat{\beta}^{Q;r})^T \Sigma \hat{\beta}^{Q;r} \geq M/2$. We can rewrite (A.18) as

$$\frac{\delta^T \beta^*}{(\hat{\beta}^{Q;r})^T \Sigma \hat{\beta}^{Q;r}} \geq 1 - 2M^{-1} \|\beta^*\|_1^2 \|\Sigma - \hat{\Sigma}^{Q;r}\|_{\max}. \quad (\text{A.19})$$

Combining (A.17) and (A.19), we have

$$\frac{\text{Rq}((\hat{\beta}^{Q;r})^{Q;r})}{\text{Rq}(\beta^*)} \geq 1 - 2M^{-1} \|\beta^*\|_1^2 \|\Sigma - \hat{\Sigma}^{Q;r}\|_{\max} - 20M^{-1} \lambda \|\beta^*\|_1.$$

This completes the proof. \square

This section describes the algorithm to solve (IV.4) for the matrix element-wise supremum norm $\|\cdot\|_{\max}$. In particular, we aim to solve

$$\tilde{\mathbf{R}} = \arg \min_{\mathbf{R} \succeq \mathbf{0}} \|\mathbf{R} - \hat{\mathbf{R}}\|_{\max} \quad (\text{A.20})$$

by the algorithm proposed in [77]. Since for any matrix $\mathbf{A} \in \mathbb{R}^d$ we can reformulate $\|\mathbf{A}\|_{\max} = \max_{\mathbf{Z} \in B_1} \text{Tr}(\mathbf{Z}^T \mathbf{A})$,

where $B_1 = \{\mathbf{Z} \in \mathbb{R}^{d \times d} \mid \mathbf{Z} = \mathbf{Z}^T, \sum_{i,j=1,\dots,d} |\mathbf{Z}_{ij}| \leq 1\}$. We define $B_2 = \{\mathbf{Z} \in \mathbb{R}^{d \times d} \mid \mathbf{Z} \succeq \mathbf{0}\}$ and (A.20) can be rewritten as a minimax problem

$$\min_{\mathbf{R} \in B_2} \max_{\mathbf{Z} \in B_1} \text{Tr}(\mathbf{Z}^T (\mathbf{R} - \hat{\mathbf{R}})). \quad (\text{A.21})$$

In order to solve the above minimax problem, we need to first study two subproblems on matrix projection. The first is

$$P_{B_2}(\mathbf{A}) = \arg \min_{\mathbf{B} \in B_2} \|\mathbf{A} - \mathbf{B}\|_F^2, \quad (\text{A.22})$$

where $\|\cdot\|_F$ is the Frobenius norm. We can have the closed form solution to A.22 such that $P_{B_2}(\mathbf{A}) = \mathbf{U} \mathbf{\Lambda}_+ \mathbf{U}^T$, where $\mathbf{A} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^T$ is the spectral decomposition of \mathbf{A} and $\mathbf{\Lambda}_+ = \text{diag}(\mathbf{\Lambda}_{11} \vee 0, \dots, \mathbf{\Lambda}_{dd} \vee 0)$. The second subproblem we are interested in is $P_{B_1}(\mathbf{A}) = \arg \min_{\mathbf{B} \in B_1} \|\mathbf{A} - \mathbf{B}\|_F^2$. This is equivalent to the vectorized problem $P_{B_1}(\mathbf{A}) = \arg \min_{\|\mathbf{v}\|_1 \leq 1} \|\text{vec}(\mathbf{A}) - \mathbf{v}\|_2^2$. We denote $\mathbf{a} = \text{vec}(\mathbf{A})$ and $|\mathbf{a}| = \text{sign}(\mathbf{a}) \circ \mathbf{a}$, where \circ is the Hadamard product. Let $T_{|\mathbf{a}|}$ be the permutation transformation matrix of $|\mathbf{a}|$ such that $T_{|\mathbf{a}|}(|\mathbf{a}|)$ is in descending order. We now define $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$ under two kinds of cases: $\|\mathbf{a}\|_1 \leq 1$ and $\|\mathbf{a}\|_1 > 1$.

If $\|\mathbf{a}\|_1 \leq 1$, we let $(\tilde{\mathbf{x}}, \tilde{\mathbf{y}}) := (\tilde{\mathbf{a}}, 0)$. If $\|\mathbf{a}\|_1 > 1$, we define $\Delta \mathbf{a} = (\tilde{a}_1 - \tilde{a}_2, \dots, \tilde{a}_{d-1} - \tilde{a}_d, \tilde{a}_d)^T$. Since $\Delta \mathbf{a}_i > 0$ for all $i = 1, \dots, n$ and $\sum_{i=1}^d i \Delta \mathbf{a}_i = \|\mathbf{a}\|_1 > 1$. We choose the smallest integer K such that $\sum_{i=1}^K i \Delta \mathbf{a}_i \geq 1$. Let

$$\tilde{\mathbf{y}} := \frac{1}{K} \left(\sum_{i=1}^K \tilde{a}_i - 1 \right), \tilde{\mathbf{x}} := (\tilde{a}_1 - \tilde{y}, \dots, \tilde{a}_K - \tilde{y}, 0, \dots, 0)^T \in \mathbb{R}^d.$$

According to [77], we can write $P_{B_1}(\mathbf{A}) = \text{sign}(\mathbf{a}) \circ T_{|\mathbf{a}|}^{-1}(\tilde{\mathbf{x}})$. We set arbitrary $\mathbf{Z}^0 \in B_1$, $\mathbf{R}^0 \in B_2$ as the initializations of the algorithm, $\gamma \in (0, 2)$ as the step length of each iteration, $\epsilon > 0$ as the tolerance, N as the maximum number of iterations. Algorithm 1 provides the following convergence to the exact solution of (A.20).

Algorithm 1 Matrix Nearness Problem in the Maximum Norm in (A.20)

```

 $\tilde{\mathbf{R}} \leftarrow \text{MatrixMaxProj}(\hat{\mathbf{R}}, \mathbf{Z}^0, \mathbf{R}^0, \gamma, \epsilon, N)$ 
for  $t = 0, \dots, N$  do
   $\mathbf{R}_0^t \leftarrow \mathbf{R}^t - P_{B_2}(\mathbf{R}^t - \mathbf{Z}^t)$ 
   $\mathbf{Z}_0^t \leftarrow \mathbf{Z}^t - P_{B_1}(\mathbf{R}^t + \mathbf{Z}^t - \hat{\mathbf{R}})$ 
  if  $\max\{\|\mathbf{R}_0^t\|_{\max}, \|\mathbf{Z}_0^t\|_{\max}\} < \epsilon$ , then
    break
  else
     $\mathbf{R}^{t+1} \leftarrow \mathbf{R}^t - \gamma(\mathbf{R}_0^t - \mathbf{Z}_0^t)/2$ 
     $\mathbf{Z}^{t+1} \leftarrow \mathbf{Z}^t - \gamma(\mathbf{R}_0^t + \mathbf{Z}_0^t)/2$ 
  end if
end for
return  $\tilde{\mathbf{R}} = \mathbf{R}^t$ 

```

Theorem A.4 ([77]): If $(\mathbf{R}^{\text{opt}}, \mathbf{Z}^{\text{opt}})$ is the solution to (A.21). Let $\mathbf{R}^t, \mathbf{Z}^t, \mathbf{R}_0^t, \mathbf{Z}_0^t$ be the sequence obtained from Algorithm 1. We have

$$\begin{aligned} &\|\mathbf{R}^{t+1} - \mathbf{R}^{\text{opt}}\|_F^2 + \|\mathbf{Z}^{t+1} - \mathbf{Z}^{\text{opt}}\|_F^2 \\ &\leq \|\mathbf{R}^t - \mathbf{R}^{\text{opt}}\|_F^2 + \|\mathbf{Z}^t - \mathbf{Z}^{\text{opt}}\|_F^2 \\ &\quad + \frac{\gamma(2-\gamma)}{2} (\|\mathbf{R}_0^t\|_F^2 + \|\mathbf{Z}_0^t\|_F^2). \end{aligned}$$

REFERENCES

- [1] P. J. Bickel and E. Levina, "Covariance regularization by thresholding," *Ann. Statist.*, vol. 36, no. 6, pp. 2577–2604, Dec. 2008.
- [2] P. J. Bickel and E. Levina, "Regularized estimation of large covariance matrices," *Ann. Statist.*, vol. 36, no. 1, pp. 199–227, Feb. 2008.
- [3] 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, Aug. 2010.
- [4] L. Xue and H. Zou, "Discussion of minimax estimation of large covariance matrices under ℓ_1 -norm," *Stat. Sinica*, vol. 22, pp. 1319–1378, 2013.
- [5] K. Lounici, "High-dimensional covariance matrix estimation with missing observations," *Bernoulli*, vol. 20, no. 3, pp. 1029–1058, Aug. 2014.
- [6] W.-K. Chen, *The Circuits and Filters Handbook*. Boca Raton, FL, USA: CRC Press, 2002.
- [7] B. O. Bradley and M. S. Taqqu, "Financial risk and heavy tails," in *Handbook of Heavy Tailed Distributions in Finance*. 2003, pp. 35–103.
- [8] F. Han and H. Liu, "Scale-invariant sparse PCA on high-dimensional meta-elliptical data," *J. Amer. Stat. Assoc.*, vol. 109, no. 505, pp. 275–287, Jan. 2014.
- [9] H. Liu, F. Han, M. Yuan, J. Lafferty, and L. Wasserman, "High-dimensional semiparametric Gaussian copula graphical models," *Ann. Statist.*, vol. 40, no. 4, pp. 2293–2326, Aug. 2012.
- [10] L. Xue and H. Zou, "Regularized rank-based estimation of high-dimensional nonparanormal graphical models," *Ann. Statist.*, vol. 40, no. 5, pp. 2541–2571, Oct. 2012.
- [11] F. Han and H. Liu, "Semiparametric principal component analysis," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 171–179.
- [12] F. Han, T. Zhao, and H. Liu, "CODA: High dimensional copula discriminant analysis," *J. Mach. Learn. Res.*, vol. 14, no. 1, pp. 629–671, 2013.
- [13] F. Han and H. Liu, "Transelliptical component analysis," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 25, 2012, pp. 368–376.
- [14] H. Liu, F. Han, and C.-H. Zhang, "Transelliptical graphical models," in *Proc. Adv. Neural Inf. Process. Syst.*, 2012, pp. 809–817.
- [15] F. Han and H. Liu, "Statistical analysis of latent generalized correlation matrix estimation in transelliptical distribution," 2013, *arXiv:1305.6916*. [Online]. Available: <http://arxiv.org/abs/1305.6916>
- [16] Z. Wang, H. Liu, and T. Zhang, "Optimal computational and statistical rates of convergence for sparse nonconvex learning problems," *Ann. Statist.*, vol. 42, no. 6, pp. 2164–2201, Dec. 2014.
- [17] J. Fan, Z. T. Ke, H. Liu, and L. Xia, "QUADRO: A supervised dimension reduction method via Rayleigh quotient optimization," 2013, *arXiv:1311.5542*. [Online]. Available: <http://arxiv.org/abs/1311.5542>
- [18] J. Fan, F. Han, and H. Liu, "PAGE: Robust pattern guided estimation of large covariance matrix," Tech. Rep., 2014.
- [19] O. Catoni, "Challenging the empirical mean and empirical variance: A deviation study," *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, vol. 48, no. 4, pp. 1148–1185, Nov. 2012.
- [20] R.-Z. Li, K.-T. Fang, and L.-X. Zhu, "Some Q - Q probability plots to test spherical and elliptical symmetry," *J. Comput. Graph. Statist.*, vol. 6, no. 4, pp. 435–450, Dec. 1997.
- [21] V. Koltchinskii and L. Sakhanenko, "Testing for ellipsoidal symmetry of a multivariate distribution," in *High Dimensional Probability II*. Springer, 2000, pp. 493–510.
- [22] L. Sakhanenko, "Testing for ellipsoidal symmetry: A comparison study," *Comput. Statist. Data Anal.*, vol. 53, no. 2, pp. 565–581, Dec. 2008.
- [23] F. W. Huffer and C. Park, "A test for elliptical symmetry," *J. Multivariate Anal.*, vol. 98, no. 2, pp. 256–281, Feb. 2007.
- [24] A. Batsidis, N. Martin, L. Pardo, and K. Zografos, "A necessary power divergence-type family of tests for testing elliptical symmetry," *J. Stat. Comput. Simul.*, vol. 84, no. 1, pp. 57–83, Jan. 2014.
- [25] S. Holm, "A simple sequentially rejective multiple test procedure," *Scand. J. Statist.*, vol. 6, no. 2, pp. 65–70, 1979.
- [26] D. Donoho and J. Jin, "Higher criticism for detecting sparse heterogeneous mixtures," *Ann. Statist.*, vol. 32, no. 3, pp. 962–994, Jun. 2004.
- [27] P. Hall and J. Jin, "Innovated higher criticism for detecting sparse signals in correlated noise," *Ann. Statist.*, vol. 38, no. 3, pp. 1686–1732, Jun. 2010.
- [28] R. A. Maronna, R. D. Martin, and V. J. Yohai, *Robust Statistics: Theory and Methods*. Hoboken, NJ, USA: Wiley, 2006.
- [29] F. R. Hampel, "The influence curve and its role in robust estimation," *J. Amer. Stat. Assoc.*, vol. 69, no. 346, pp. 383–393, Jun. 1974.
- [30] P. J. Rousseeuw and C. Croux, "Alternatives to the median absolute deviation," *J. Amer. Stat. Assoc.*, vol. 88, no. 424, pp. 1273–1283, Dec. 1993.
- [31] C. Croux and A. Ruiz-Gazen, "High breakdown estimators for principal components: The projection-pursuit approach revisited," *J. Multivariate Anal.*, vol. 95, no. 1, pp. 206–226, Jul. 2005.
- [32] P. J. Huber and E. Ronchetti, *Robust Statistics*, 2nd ed. Hoboken, NJ, USA: Wiley, 2009.
- [33] R. Gnanadesikan and J. R. Kettenring, "Robust estimates, residuals, and outlier detection with multiresponse data," *Biometrics*, vol. 28, no. 1, pp. 81–124, 1972.
- [34] M. G. Genton and Y. Ma, "Robustness properties of dispersion estimators," *Statist. Probab. Lett.*, vol. 44, no. 4, pp. 343–350, Oct. 1999.
- [35] Y. Ma and M. G. Genton, "Highly robust estimation of dispersion matrices," *J. Multivariate Anal.*, vol. 78, no. 1, pp. 11–36, Jul. 2001.
- [36] R. A. Maronna and R. H. Zamar, "Robust estimates of location and dispersion for high-dimensional datasets," *Technometrics*, vol. 44, no. 4, pp. 307–317, Nov. 2002.
- [37] L. Wang, Y. Wu, and R. Li, "Quantile regression for analyzing heterogeneity in ultra-high dimension," *J. Amer. Stat. Assoc.*, vol. 107, no. 497, pp. 214–222, Mar. 2012.
- [38] A. Belloni and V. Chernozhukov, " ℓ_1 -penalized quantile regression in high-dimensional sparse models," *Ann. Statist.*, vol. 39, no. 1, pp. 82–130, 2011.
- [39] L. Wang, "The L_1 penalized LAD estimator for high dimensional linear regression," *J. Multivariate Anal.*, vol. 120, pp. 135–151, Sep. 2013.
- [40] T. Cai, W. Liu, and X. Luo, "A constrained ℓ_1 minimization approach to sparse precision matrix estimation," *J. Amer. Stat. Assoc.*, vol. 106, no. 494, pp. 594–607, Jun. 2011.
- [41] K.-T. Fang, S. Kotz, and K. W. Ng, *Symmetric Multivariate and Related Distributions* (Monographs on Statistics and Applied Probability). London, U.K.: Chapman & Hall, 1990.
- [42] J. Owen and R. Rabinovitch, "On the class of elliptical distributions and their applications to the theory of portfolio choice," *J. Finance*, vol. 38, no. 3, pp. 745–752, Jun. 1983.
- [43] J. B. Berk, "Necessary conditions for the CAPM," *J. Econ. Theory*, vol. 73, no. 1, pp. 245–257, Mar. 1997.
- [44] A. J. McNeil, R. Frey, and P. Embrechts, *Quantitative Risk Management: Concepts, Techniques and Tools*. Princeton, NJ, USA: Princeton Univ. Press, 2010.
- [45] P. Embrechts, A. McNeil, and D. Straumann, "Correlation and dependence in risk management: Properties and pitfalls," in *Risk Management: Value at Risk and Beyond*. 2002, pp. 176–223.
- [46] D. B. Marden and D. G. Manolakis, "Using elliptically contoured distributions to model hyperspectral imaging data and generate statistically similar synthetic data," *Proc. SPIE*, vol. 5425, pp. 558–572, Aug. 2004.
- [47] J. Frontera-Pons, M. Mahot, J. P. Ovarlez, F. Pascal, S. K. Pang, and J. Chanussot, "A class of robust estimates for detection in hyperspectral images using elliptical distributions background," in *Proc. IEEE Int. Geosci. Remote Sens. Symp. (IGARSS)*, Jul. 2012, pp. 4166–4169.
- [48] G. Frahm, "Generalized elliptical distributions: Theory and applications," Ph.D. dissertation, Universität zu Köln, Cologne, Germany, 2004.
- [49] H.-B. Fang, K.-T. Fang, and S. Kotz, "The meta-elliptical distributions with given marginals," *J. Multivariate Anal.*, vol. 82, no. 1, pp. 1–16, Jul. 2002.
- [50] H. Liu, J. Lafferty, and L. Wasserman, "The nonparanormal: Semiparametric estimation of high dimensional undirected graphs," *J. Mach. Learn. Res.*, vol. 10, no. 10, pp. 2295–2328, Oct. 2009.
- [51] A. B. Tsybakov, *Introduction to Nonparametric Estimation*. New York, NY, USA: Springer, 2009.
- [52] Z. D. Bai and Y. Q. Yin, "Limit of the smallest eigenvalue of a large dimensional sample covariance matrix," *Ann. Probab.*, vol. 21, no. 3, pp. 1275–1294, Jul. 1993.
- [53] I. M. Johnstone, "On the distribution of the largest eigenvalue in principal components analysis," *Ann. Statist.*, vol. 29, no. 2, pp. 295–327, Apr. 2001.
- [54] N. El Karoui, "Operator norm consistent estimation of large-dimensional sparse covariance matrices," *Ann. Statist.*, vol. 36, no. 6, pp. 2717–2756, Dec. 2008.
- [55] A. J. Rothman, E. Levina, and J. Zhu, "Generalized thresholding of large covariance matrices," *J. Amer. Stat. Assoc.*, vol. 104, no. 485, pp. 177–186, Mar. 2009.
- [56] L. Xue, S. Ma, and H. Zou, "Positive-definite ℓ_1 -penalized estimation of large covariance matrices," *J. Amer. Stat. Assoc.*, vol. 107, no. 500, pp. 1480–1491, Dec. 2012.
- [57] H. Liu, L. Wang, and T. Zhao, "Sparse covariance matrix estimation with eigenvalue constraints," *J. Comput. Graph. Statist.*, vol. 23, no. 2, pp. 439–459, Apr. 2014.

- [58] T. T. Cai and H. H. Zhou, "Optimal rates of convergence for sparse covariance matrix estimation," *Ann. Statist.*, vol. 40, no. 5, pp. 2389–2420, Oct. 2012.
 - [59] N. Meinshausen and P. Bühlmann, "High-dimensional graphs and variable selection with the Lasso," *Ann. Statist.*, vol. 34, no. 3, pp. 1436–1462, Jun. 2006.
 - [60] M. Yuan and Y. Lin, "Model selection and estimation in the Gaussian graphical model," *Biometrika*, vol. 94, no. 1, pp. 19–35, Feb. 2007.
 - [61] J. Friedman, T. Hastie, and R. Tibshirani, "Sparse inverse covariance estimation with the graphical Lasso," *Biostatistics*, vol. 9, no. 3, pp. 432–441, Jul. 2008.
 - [62] O. Banerjee, L. E. Ghaoui, and A. d'Aspremont, "Model selection through sparse maximum likelihood estimation for multivariate Gaussian or binary data," *J. Mach. Learn. Res.*, vol. 9, pp. 485–516, Jun. 2008.
 - [63] C. Lam and J. Fan, "Sparsistency and rates of convergence in large covariance matrix estimation," *Ann. Statist.*, vol. 37, no. 6B, pp. 4254–4278, Dec. 2009.
 - [64] M. Yuan, "High dimensional inverse covariance matrix estimation via linear programming," *J. Mach. Learn. Res.*, vol. 11, pp. 2261–2286, Aug. 2010.
 - [65] T. T. Cai, W. Liu, and H. H. Zhou, "Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation," *Ann. Statist.*, vol. 44, no. 2, pp. 455–488, Apr. 2016.
 - [66] I. M. Johnstone and A. Y. Lu, "On consistency and sparsity for principal components analysis in high dimensions," *J. Amer. Stat. Assoc.*, vol. 104, no. 486, pp. 682–693, Jun. 2009.
 - [67] V. Q. Vu, J. Cho, J. Lei, and K. Rohe, "Fantope projection and selection: A near-optimal convex relaxation of sparse PCA," in *Proc. Adv. Neural Inf. Process. Syst.*, 2013, pp. 2670–2678.
 - [68] Y. Guo, T. Hastie, and R. Tibshirani, "Regularized discriminant analysis and its application in microarrays," *Biostatistics*, vol. 1, no. 1, pp. 1–18, 2005.
 - [69] J. Fan and Y. Fan, "High-dimensional classification using features annealed independence rules," *Ann. Statist.*, vol. 36, no. 6, pp. 2605–2637, Dec. 2008.
 - [70] J. Shao, Y. Wang, X. Deng, and S. Wang, "Sparse linear discriminant analysis by thresholding for high dimensional data," *Ann. Statist.*, vol. 39, no. 2, pp. 1241–1265, Apr. 2011.
 - [71] T. Cai and W. Liu, "A direct estimation approach to sparse linear discriminant analysis," *J. Amer. Stat. Assoc.*, vol. 106, no. 496, pp. 1566–1577, Dec. 2011.
 - [72] J. Fan, Y. Feng, and X. Tong, "A ROAD to classification in high dimensional space: The regularized optimal affine discriminant," *J. Roy. Stat. Soc. B, Stat. Methodol.*, vol. 74, no. 4, pp. 745–771, 2012.
 - [73] M. N. McCall, B. M. Bolstad, and R. A. Irizarry, "Frozen robust multiarray analysis (fRMA)," *Biostatistics*, vol. 11, no. 2, pp. 242–253, Apr. 2010.
 - [74] J. Fan, Y. Liao, and M. Mincheva, "Large covariance estimation by thresholding principal orthogonal complements," *J. Roy. Stat. Soc. B, Stat. Methodol.*, vol. 75, no. 4, pp. 603–680, Sep. 2013.
 - [75] F. Han and H. Liu, "Principal component analysis on non-Gaussian dependent data," in *Proc. Int. Conf. Mach. Learn.*, vol. 30, 2013, pp. 240–248.
 - [76] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *J. Amer. Stat. Assoc.*, vol. 58, no. 301, pp. 13–30, Mar. 1963.
 - [77] M. H. Xu and H. Shao, "Solving the matrix nearness problem in the maximum norm by applying a projection and contraction method," *Adv. Oper. Res.*, vol. 2012, pp. 1–15, Jan. 2012.
- Junwei Lu** received the Ph.D. degree from the Department of Operations Research and Financial Engineering, Princeton University. He is currently an Assistant Professor of biostatistics with the Harvard T.H. Chan School of Public Health. His research interests include statistics, machine learning, and their applications to genomics and neuroscience.
- Fang Han** received the Ph.D. degree from the Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health. He is currently an Assistant Professor of statistics with the University of Washington. His main methodological/theoretical interests are in rank-based statistics, non-parametric and semi-parametric regressions, time series analysis, and random matrix theory.
- Han Liu** received the joint Ph.D. degree in machine learning and statistics from the Machine Learning Department, Carnegie Mellon University. He is currently an Associate Professor with the Department of Electrical Engineering and Computer Science and the Department of Statistics, Northwestern University. Before joining Northwestern University, he has directed the Center of Deep Reinforcement Learning, Tencent AI Lab. He had been a Faculty Member with Princeton University and Johns Hopkins University. His research direction is ubiquitous analytics, which deploys computational intelligence (AI) and computational trust (Blockchain) on edges and clouds to achieve analytical advantages.