

# High Dimensional Semiparametric Scale-Invariant Principal Component Analysis

Fang Han and Han Liu

**Abstract**—We propose a new high dimensional semiparametric principal component analysis (PCA) method, named Copula Component Analysis (COCA). The semiparametric model assumes that, after unspecified marginally monotone transformations, the distributions are multivariate Gaussian. COCA improves upon PCA and sparse PCA in three aspects: (i) It is robust to modeling assumptions; (ii) It is robust to outliers and data contamination; (iii) It is scale-invariant and yields more interpretable results. We prove that the COCA estimators obtain fast estimation rates and are feature selection consistent when the dimension is nearly exponentially large relative to the sample size. Careful experiments confirm that COCA outperforms sparse PCA on both synthetic and real-world data sets.

**Index Terms**—High dimensional statistics, principal component analysis, nonparanormal distribution, robust statistics

## 1 INTRODUCTION

In this paper we propose a new principal component analysis (PCA), named Copula Component Analysis (COCA), based on a semiparametric model for analyzing high dimensional non-Gaussian data. The semiparametric model assumes that, after marginal-wise unspecified strictly increasing transformations, the data are Gaussian distributed. This model is proposed by [1] and a rank-based estimator for inferring graphical models is proposed by [2]. In this paper, we generalize their results to estimate the leading eigenvectors of the correlation and covariance matrices. New estimation methods and their theoretical and empirical performances are provided.

Let  $X \in \mathbb{R}^d$  be the random vector with interest to us. PCA aims at recovering the top  $m$  leading eigenvectors  $u_1, \dots, u_m$  of  $\Sigma := \text{Cov}(X)$ . In practice,  $\Sigma$  is unknown and is replaced by the sample covariance matrix  $S$  using  $n$  independent realizations of  $X$ . For fixed  $d$ , PCA always achieves a consistent estimator and its asymptotic efficiency property is well addressed [3]. However, under a double asymptotic framework in which both the sample size  $n$  and dimensionality  $d$  can increase (with possibly  $d > n$ ), [4] showed that the leading eigenvector of  $S$  cannot converge to  $u_1 = (u_{11}, \dots, u_{1d})^T$ . A common remedy is to assume that  $s := \text{card}(\{j : u_{1j} \neq 0\}) < n$ . Different sparse PCA algorithms are developed to exploit this sparsity structure and we refer to, [5], [6], and [7], among others.

There are several drawbacks of PCA and sparse PCA: (i) Data are assumed to be Gaussian or sub-Gaussian distributed such that a fast convergence rate can be obtained;

(ii) They are not scale-invariant, i.e., changing the measurement scale of variables makes the estimates different [8]; (iii) They are not robust to data contaminations (outliers, for example). To address these concerns, we propose a high dimensional semiparametric scale-invariant principal component analysis method, named COpa Component Analysis, based on the nonparanormal family. Here we say that  $X = (X_1, \dots, X_d)^T$  is nonparanormally distributed if there exists a set of univariate strictly increasing functions  $f = \{f_j\}_{j=1}^d$  such that  $(f_1(X_1), \dots, f_d(X_d))^T \sim N_d(0, \Sigma^0)$ . By treating the monotone transformation functions  $\{f_j\}_{j=1}^d$  as a type of data contamination, COCA aims at recovering the leading eigenvectors of the latent correlation matrix  $\Sigma^0$ .

Compared with PCA and sparse PCA, COCA is scale-invariant and its estimating procedure is adaptive over the whole nonparanormal family. The nonparanormal family contains and is much larger than the Gaussian. By exploiting a rank-based regularized procedure for parameter estimation, the COCA estimator is not only robust to modeling and data contaminations, but can be consistent even when the dimensionality is nearly exponentially large relative to the sample size.

In this paper, to complete the story, a scale variant PCA method, named Copula PCA, is also proposed. Copula PCA estimates the leading eigenvector of the latent covariance matrix  $\Sigma$  (detailed definition provided in Section 2.2). To estimate  $\Sigma$ , instead of  $\Sigma^0$ , in a fast rate, we prove that extra conditions are required on the transformation functions.

[2] proposed a procedure called the nonparanormal SKEPTIC to estimate the graphical model via exploiting the nonparanormal distribution to model the data and rank based methods for estimation. COCA is different from the nonparanormal SKEPTIC in three aspects: (i) Their focus is on graph estimation, in contrast, this paper focuses on PCA and propose new estimation methods with thorough theoretical analysis provided; (ii) We provide a second step projection to make the estimated rank-based correlation and covariance matrices positive semidefinite, and prove that the same parametric rate can be preserved; (iii) Unlike the previous analysis, this paper provides extra conditions on

• F. Han is with the Department of Biostatistics, Johns Hopkins University, Baltimore, 21205 MD. E-mail: fhan@jhsph.edu.

• H. Liu is with the Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08540. E-mail: hanliu@princeton.edu.

Manuscript received 12 July 2012; revised 23 July 2013; accepted 17 Jan. 2014. Date of publication 23 Feb. 2014; date of current version 10 Sept. 2014. Recommended for acceptance by E.P. Xing.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below. Digital Object Identifier no. 10.1109/TPAMI.2014.2307886

the transformation functions to guarantee the fast rates of convergence for Copula PCA, and we discuss the advantages of COCA over Copula PCA.

The rest of the paper is organized as follows. In the next section, we briefly discuss the statistical model of the scale-invariant PCA and review the nonparanormal model and rank-based estimators shown in [1], [2]. In Section 3, we present the model of COCA and introduce the corresponding estimators and algorithms. We provide a theoretical analysis of COCA estimators in Section 4. In Section 5, we employ COCA on both synthetic and real-world data to show its empirical usefulness. Some of the results in this paper were first stated without proofs in a conference version [9].

## 2 BACKGROUND

We start with notations: Let  $M = [M_{jk}] \in \mathbb{R}^{d \times d}$  and  $v = (v_1, \dots, v_d)^T \in \mathbb{R}^d$ . Let  $v$ 's subvector with entries indexed by  $I$  be denoted by  $v_I$ . Let  $M$ 's submatrix with rows indexed by  $I$  and columns indexed by  $J$  be denoted by  $M_{IJ}$ . Let  $M_{I*}$  and  $M_{*J}$  be the submatrix of  $M$  with rows in  $I$ , and the submatrix of  $M$  with columns in  $J$ . For  $0 < q < \infty$ , we define the  $\ell_q$  and  $\ell_\infty$  vector norms as  $\|v\|_q := (\sum_{i=1}^d |v_i|^q)^{1/q}$  and  $\|v\|_\infty := \max_{1 \leq i \leq d} |v_i|$ , and we define  $\|v\|_0 := \text{card}(\text{supp}(v))$ . Here  $\text{card}(\cdot)$  represents the cardinality and  $\text{supp}(v) := \{j : v_j \neq 0\}$ . We define the matrix  $\ell_{\max}$  norm as the elementwise maximum value:  $\|M\|_{\max} := \max\{|M_{ij}|\}$ . We define  $\text{Tr}(M)$  to be the trace of  $M$ . Let  $\Lambda_j(M)$  be the  $j$ th largest eigenvalue of  $M$ . In particular,  $\Lambda_{\min}(M) := \Lambda_d(M)$  and  $\Lambda_{\max}(M) := \Lambda_1(M)$  are the smallest and largest eigenvalues of  $M$ . The vectorized matrix of  $M$ , denoted by  $\text{vec}(M)$ , is defined as  $\text{vec}(M) := (M_{*1}^T, \dots, M_{*d}^T)^T$ . Let  $S^{d-1} := \{v \in \mathbb{R}^d : \|v\|_2 = 1\}$  be the  $d$ -dimensional  $\ell_2$  sphere. For any two vectors  $a, b \in \mathbb{R}^d$  and any two square matrices  $A, B \in \mathbb{R}^{d \times d}$ , denote the inner product of  $a$  and  $b$ ,  $A$  and  $B$  by  $\langle a, b \rangle := a^T b$  and  $\langle A, B \rangle := \text{Tr}(A^T B)$ . Let  $\text{diag}(M) := (M_{11}, M_{22}, \dots, M_{dd})^T$ . we denote  $\text{sign}(a) := (\text{sign}(a_1), \dots, \text{sign}(a_d))^T$ , where  $\text{sign}(x) := x/|x|$  with the convention  $0/0 = 0$ .

### 2.1 The Models of PCA and Scale-Invariant PCA

PCA is not scale-invariant, meaning that variables measured in different scales will result in different estimators [10]. To attack this problem, PCA conducted on the sample correlation matrix  $S^0$  instead of the sample covariance matrix  $S$  is commonly used. We call the procedure of conducting PCA on  $S^0$  the scale-invariant PCA. It is realized that a large portion of works claiming doing PCA are actually doing the scale-invariant PCA [8], and the theoretical performance of the scale-invariant PCA in low dimensions has been studied [11], [12]. It is under debate whether PCA or the scale-invariant PCA are preferred in different circumstances and we refer to [13], [10], and [14] for more discussions on it.

Let  $\Sigma^0$  and  $\Sigma$  be the correlation and covariance matrices of a random vector  $X \in \mathbb{R}^d$ . Let  $\omega_1 \geq \omega_2 \geq \dots \geq \omega_d > 0$  and  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d > 0$  be the eigenvalues of  $\Sigma$  and  $\Sigma^0$ . Let  $u_1, \dots, u_d$  and  $\theta_1, \dots, \theta_d$  be the corresponding eigenvectors. The next proposition claims that the estimators  $\{\hat{u}_1, \dots, \hat{u}_d\}$  and  $\{\hat{\theta}_1, \dots, \hat{\theta}_d\}$ , which are the eigenvectors of the sample covariance and correlation matrices  $S$  and  $S^0$ , are the

maximum likelihood estimators (MLEs) of  $\{u_1, \dots, u_d\}$  and  $\{\theta_1, \dots, \theta_d\}$ :

**Proposition 2.1 ([3]).** *Let  $X \sim N_d(\mu, \Sigma)$  and  $\Sigma^0$  be the correlation matrix of  $X$ . Let  $x_1, \dots, x_n$  be  $n$  independent realizations of  $X$ . Then the estimators of PCA,  $\{\hat{u}_1, \dots, \hat{u}_d\}$ , and the estimators of the scale-invariant PCA,  $\{\hat{\theta}_1, \dots, \hat{\theta}_d\}$ , are the MLEs of  $\{u_1, \dots, u_d\}$  and  $\{\theta_1, \dots, \theta_d\}$ .*

The scale-invariant PCA is a safe procedure for dimension reduction when variables are measured in different scales. In this paper we further show that under a more general nonparanormal (or Gaussian copula) model, the scale-invariant PCA will pose less conditions than PCA to make the estimators achieve good theoretical performance.

### 2.2 The Nonparanormal Distribution

We first introduce the two definitions of the nonparanormal distribution separately shown in [1] and [2]. These two definitions will be used to define the models of COCA and Copula PCA in the next section.

**Definition 2.2 ([1]).** *A random vector  $X = (X_1, \dots, X_d)^T$  with means  $\mu = (\mu_1, \dots, \mu_d)^T$  and standard deviations  $\{\sigma_1, \dots, \sigma_d\}$  is said to follow a margin-preserved nonparanormal distribution  $MNPN_d(\mu, \Sigma, f)$  if and only if there exists a set of strictly increasing univariate functions  $f = \{f_j\}_{j=1}^d$  such that:*

$$f(X) = (f_1(X_1), \dots, f_d(X_d))^T \sim N_d(\mu, \Sigma),$$

where  $\text{diag}(\Sigma) = (\sigma_1^2, \dots, \sigma_d^2)^T$ . We call  $\Sigma$  the latent covariance matrix.

**Definition 2.3 ([2]).** *Let  $f^0 = \{f_j^0\}_{j=1}^d$  be a set of strictly increasing univariate functions. We say that a  $d$  dimensional random vector  $X = (X_1, \dots, X_d)^T$  follows a nonparanormal distribution  $NPN_d(\Sigma^0, f^0)$ , if*

$$f^0(X) := (f_1^0(X_1), \dots, f_d^0(X_d))^T \sim N_d(\mathbf{0}, \Sigma^0),$$

where  $\text{diag}(\Sigma^0) = 1$ . We call  $\Sigma^0$  the latent correlation matrix.

We have the following lemma, which proves that the two definitions of the nonparanormal are equivalent.

**Lemma 2.4.** *A random vector  $X \sim NPN_d(\Sigma^0, f^0)$  if and only if there exists  $\mu = (\mu_1, \dots, \mu_d)^T$ ,  $\Sigma = [\Sigma_{jk}] \in \mathbb{R}^{d \times d}$  with*

$$\mathbb{E}(X_j) = \mu_j, \text{Var}(X_j) = \Sigma_{jj} \quad \text{and} \quad \Sigma_{jk}^0 = \frac{\Sigma_{jk}}{\sqrt{\Sigma_{jj} \cdot \Sigma_{kk}}},$$

and a set of strictly increasing univariate functions  $f = \{f_j\}_{j=1}^d$  such that  $X \sim MNPN_d(\mu, \Sigma, f)$ .

**Proof.** Using the connection that  $f_j(\cdot) = \mu_j + \sigma_j f_j^0(\cdot)$ , for  $j \in \{1, 2, \dots, d\}$ .  $\square$

[1] proved that the nonparanormal family is equivalent to the continuous Gaussian copula family [15]. Definition 2.3 is more appealing because it emphasizes the correlation and hence matches the spirit of the copula. However, Definition 2.2 enjoys notational simplicity in analyzing the nonparanormal based linear discriminant analysis and scale-variant PCA methods.

Here we note that in Definition 2.2, the model is identifiable. Moreover, the parameters  $\mu$  and  $\Sigma$  in the latent

Gaussian random vector  $f(X) \sim N_d(\mu, \Sigma)$  are unique. The identifiability issue has been discussed in [1]. The uniqueness of  $\mu$  and  $\Sigma$  in  $f(X)$  are imposed by modeling assumption: We assume that the transformation function  $f$  preserves the first two marginal moments, i.e.,  $\mathbb{E}X_j = \mathbb{E}f_j(X_j)$  and  $\text{Var}(X_j) = \text{Var}(f_j(X_j))$  for  $j = 1, \dots, d$ . In this way, we can exploit the nonparanormal model in conducting the procedures that require more information besides the correlations.

### 2.3 Spearman's Rho Correlation and Covariance Matrices

Given  $n$  data points  $x_1, \dots, x_n \in \mathbb{R}^d$ , where  $x_i = (x_{i1}, \dots, x_{id})^T$ , we denote by

$$\hat{\mu}_j := \frac{1}{n} \sum_{i=1}^n x_{ij} \quad \text{and} \quad \hat{\sigma}_j = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \hat{\mu}_j)^2},$$

the marginal sample means and standard deviations. Let  $r_{ij}$  be the rank of  $x_{ij}$  among  $x_{1j}, \dots, x_{nj}$  and  $\bar{r}_j := \frac{1}{n} \sum_{i=1}^n r_{ij} = \frac{n+1}{2}$ , we consider the following statistics:

$$\hat{\rho}_{jk} = \frac{\sum_{i=1}^n (r_{ij} - \bar{r}_j)(r_{ik} - \bar{r}_k)}{\sqrt{\sum_{i=1}^n (r_{ij} - \bar{r}_j)^2 \cdot \sum_{i=1}^n (r_{ik} - \bar{r}_k)^2}},$$

and the correlation matrix estimators:

$$\hat{\mathbf{R}}_{jk} = \begin{cases} 2 \sin\left(\frac{\pi}{6} \hat{\rho}_{jk}\right), & j \neq k, \\ 1, & j = k. \end{cases} \quad (2.1)$$

The Equation (2.1) is inspired from Equation (6.4) in [16]. We denote by  $\hat{\mathbf{R}} := [\hat{\mathbf{R}}_{jk}]$  and  $\hat{\mathbf{S}} := [\hat{\mathbf{S}}_{jk}] = [\hat{\sigma}_j \hat{\sigma}_k \hat{\mathbf{R}}_{jk}]$  the Spearman's rho correlation and covariance matrices. Lemma 2.5, coming from [2], claims that  $\hat{\mathbf{R}}$  can approach  $\Sigma^0$  in the parametric rate.

**Lemma 2.5 ([2]).** When  $x_1, \dots, x_n \sim^{i.i.d} NPN_d(\Sigma^0, f^0)$ , for any  $n \geq \frac{21}{\log d} + 2$ , with probability at least  $1 - 1/d^2$ ,

$$\|\hat{\mathbf{R}} - \Sigma^0\|_{\max} \leq 8\pi \sqrt{\frac{\log d}{n}}. \quad (2.2)$$

## 3 METHODS

In this section, we first provide the statistical models of Copula Component Analysis and Copula PCA method. And then we introduce several algorithms to solve this problem.

### 3.1 Models

One of the intuition of PCA is coming from the Gaussian distribution. The principal components define the major axes of the contours of constant probability for the multivariate Gaussian [3]. However, such an interpretation does not exist when the distributions are away from the Gaussian. [17] constructed examples where PCA cannot preserve the structure of the data. Here we propose a toy example to show this phenomenon.

In Fig. 1, we randomly generate 10,000 samples from three different types of nonparanormal distributions. We suppose that  $X \sim NPN_2(\Sigma^0, f^0)$ . Here we set  $\Sigma^0 = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$  and transformation functions as follows: (A)  $f_1^0(x) = x^3$  and  $f_2^0(x) = x^{1/3}$ ; (B)  $f_1^0(x) = \text{sign}(x)x^2$  and  $f_2^0(x) = x^3$ ; (C)  $f_1^0(x) = f_2^0(x) = \Phi^{-1}(x)$ .

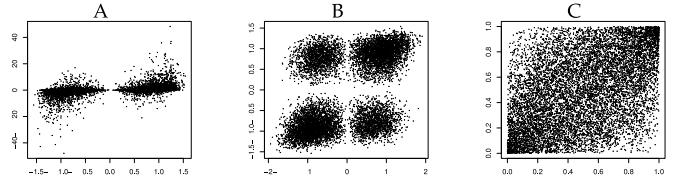


Fig. 1. Scatter plots of three nonparanormals,  $X \sim NPN_2(\Sigma^0, f^0)$ . Here  $\Sigma_{12}^0 = 0.5$  and the transformation functions have the form as follows: (A)  $f_1^0(x) = x^3$  and  $f_2^0(x) = x^{1/3}$ ; (B)  $f_1^0(x) = \text{sign}(x)x^2$  and  $f_2^0(x) = x^3$ ; (C)  $f_1^0(x) = f_2^0(x) = \Phi^{-1}(x)$ .

$f_1^0(x) = f_2^0(x) = \Phi^{-1}(x)$ , where  $\Phi$  is defined as the distribution function of the standard Gaussian distribution. Here, researchers might wish to conduct PCA separately on different clusters in (A) and (B). For (C), the data look very noisy and a nice major axis might be considered not existing.

However, considering the monotone transformation  $f^0$  as a type of data contamination, the geometric intuition of PCA comes back by estimating the principal components of the latent Gaussian distribution. In the next section, we will present the model of COCA and Copula PCA motivated from this observation.

#### 3.1.1 COCA Model

We first show the model of Copula Component Analysis method, where the idea of the scale-invariant PCA is exploited. We wish to estimate the leading eigenvector of the latent correlation matrix. In particular, let  $\theta_1$  be the leading eigenvectors of  $\Sigma^0$ . For  $0 \leq q \leq 1$ , the  $\ell_q$  ball  $\mathbb{B}_q(R_q)$  is defined as

$$\begin{aligned} \text{when } q = 0, \quad \mathbb{B}_0(R_0) &:= \{v \in \mathbb{R}^d : \text{card}(\text{supp}(v)) \leq R_0\}; \\ \text{when } 0 < q \leq 1, \quad \mathbb{B}_q(R_q) &:= \{v \in \mathbb{R}^d : \|v\|_q^q \leq R_q\}. \end{aligned}$$

Accordingly, the COCA model  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$  is considered:

$$\mathcal{M}^0(q, R_q, \Sigma^0, f^0) = \{X : X \sim NPN_d(\Sigma^0, f^0), \theta_1 \in \mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q)\}. \quad (3.1)$$

The  $\ell_q$  ball induces a (weak) sparsity pattern when  $0 \leq q \leq 1$  and has been analyzed in linear regression [18] and sparse PCA [7], [19]. Moreover, the data are assumed to come from a nonparanormal (or Gaussian copula) distribution, which contains and is a much larger distribution family than the Gaussian.

Inspired by the model  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$ , we consider the following estimator  $\tilde{\theta}_1$ , which is the global optimum to the following equation with the constraint that  $\tilde{\theta}_1 \in \mathbb{B}_q(R_q)$  for some  $0 \leq q \leq 1$ :

$$\begin{aligned} \tilde{\theta}_1 &= \arg \max_{v \in \mathbb{R}^d} v^T \hat{\mathbf{R}} v, \\ &\text{subject to } v \in \mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q). \end{aligned} \quad (3.2)$$

Here  $\hat{\mathbf{R}}$  is the estimated Spearman's rho correlation matrix. The corresponding estimator  $\tilde{\theta}_1$  can be considered as a nonlinear dimension reduction procedure and has the potential to gain more flexibility compared with PCA.

### 3.1.2 Copula PCA Model

In contrast, we provide another method called Copula PCA, where we wish to estimate the leading eigenvector of the latent covariance matrix. In particular, let  $u_1$  be the leading eigenvector of  $\Sigma$ . The following Copula PCA model  $\mathcal{M}(q, R_q, \Sigma, f)$  is considered:

$$\begin{aligned} \mathcal{M}(q, R_q, \Sigma, f) = & \{X : X \sim MNPN_d(\mu, \Sigma, f), \\ & u_1 \in \mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q)\}. \end{aligned} \quad (3.3)$$

An estimator corresponding to the above model is:

$$\begin{aligned} \tilde{u}_1 = \arg \max_{v \in \mathbb{R}^d} v^T \hat{S} v \\ \text{subject to } v \in \mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q), \end{aligned} \quad (3.4)$$

where  $\hat{S}$  is the Spearman's rho covariance matrix.

### 3.1.3 Attainability of the Proposed Estimators

The direct computation of estimators  $\tilde{\theta}_1$  and  $\tilde{u}_1$  as defined in Equation (3.2) and (3.4) might be time consuming. However, in the following section we show several algorithms which could approach these two global optima and have good empirical performance. In particular, in Section 4 we will provide the theoretical performance in terms of guarantees of convergence and convergence rate of parameter estimation for the proposed algorithms. We will show that the global optima proposed in Equations (3.2) and (3.4) can be well approached by using the truncated power algorithm. This algorithm has a (weaker) guarantee of convergence and under certain sufficient conditions the corresponding estimator can achieve the same convergence rate as the global optimum. Detailed theoretical analysis is provided in Section 4 as two new theorems (Theorems 4.12 and 4.14).

## 3.2 Algorithms

In this section we provide three sparse PCA algorithms, which the Spearman's rho correlation and covariance matrices  $\hat{R}$  and  $\hat{S}$  can be directly plugged in.

### 3.2.1 COCA and Copula PCA with PMD

Penalized matrix decomposition (PMD) is proposed by [20]. The main idea of PMD is a bi-convex optimization algorithm to the following problem:

$$\arg \max_{v,w} v^T \hat{\Gamma} w, \quad \text{s.t. } \|v\|_2^2 \leq 1, \|w\|_2^2 \leq 1, \|v\|_1 \leq \delta, \|w\|_1 \leq \delta.$$

COCA with PMD and Copula PCA with PMD are listed in the following:

- 1. Input: A symmetric matrix  $\hat{\Gamma}$ . Initialize  $w \in \mathbb{S}^{d-1}$ .
- 2. Iterate until convergence:
  - (a)  $v \leftarrow \arg \max_{v \in \mathbb{R}^d} v^T \hat{\Gamma} w$  subject to  $\|v\|_1 \leq \delta, \|v\|_2^2 \leq 1$ .
  - (b)  $w \leftarrow \arg \max_{w \in \mathbb{R}^d} w^T \hat{\Gamma} v$  subject to  $\|w\|_1 \leq \delta, \|w\|_2^2 \leq 1$ .
- 3. Output:  $w$ .

Here  $\hat{\Gamma}$  is either  $\hat{R}$  or  $\hat{S}$ , corresponding to COCA with PMD and Copula PCA with PMD.  $\delta$  is the tuning parameter. In practice, [20] suggested using the first leading eigenvector of

$\hat{\Gamma}$  to be the initial value. PMD can be considered as a solver to Equation (3.2) and Equation (3.4) with  $q = 1$ .

### 3.2.2 COCA and Copula PCA with SPCA

The SPCA algorithm is proposed by [21]. The main idea of SPCA is to exploit a regression approach to PCA and then utilize the lasso and elastic net [22] to calculate a sparse estimator. COCA with SPCA and Copula PCA with SPCA are listed as follows:

- 1. Input: A symmetric matrix  $\hat{\Gamma}$ . Initialize  $v \in \mathbb{S}^{d-1}$ .
- 2. Iterate until convergence:
  - (a)  $w \leftarrow \arg \min_{w \in \mathbb{R}^d} (v - w)^T \hat{\Gamma} (v - w) + \delta_1 \|w\|_2^2 + \delta_2 \|w\|_1$ ;
  - (b)  $v \leftarrow \hat{\Gamma} w / \|\hat{\Gamma} w\|_2$ .
- 3. Output:  $w / \|w\|_2$ .

Here  $\hat{\Gamma}$  is either  $\hat{R}$  or  $\hat{S}$ , corresponding to COCA with SPCA and Copula PCA with SPCA.  $\delta_1 \in \mathbb{R}$  and  $\delta_2 \in \mathbb{R}$  are two tuning parameters. In practice, [21] suggested using the first leading eigenvector of  $\hat{\Gamma}$  to be the initial value. SPCA can also be considered as a solver to Equation (3.2) and Equation (3.4) with  $q = 1$ .

### 3.2.3 COCA and Copula PCA with TPower

Truncated power method (TPower) is proposed by [5]. The main idea of TPower is to utilize the power method, but truncate the vector to a  $\ell_0$  ball in each iteration. Actually, TPower can be generalized to a family of algorithms to solve Equation (3.2) when  $0 \leq q \leq 1$ , as presented in Algorithm 3.1. We name it  $\ell_q$  Constraint Truncated Power Method (qTPM). In particular, when  $q = 0$ , the algorithm qTPM coincides with [5]'s method.

---

#### Algorithm 1 $\ell_q$ Constraint Truncated Power Method

**Input:** : symmetry matrix  $\hat{\Gamma}$ , initial vector  $\tilde{\theta}_{q,0} \in \mathbb{R}^d$

**Output:** :  $\theta_{q,\infty}$

Let  $t = 1$  and  $R_q$  be the tuning parameter

**repeat**

    compute  $x_t = \hat{\Gamma} \cdot \tilde{\theta}_{q,t-1} / \|\hat{\Gamma} \cdot \tilde{\theta}_{q,t-1}\|_2$

**if**  $\|x_t\|_q \leq R_q^{1/q}$  **then**

$\theta_{q,t} = x_t$

**else**

        Let  $A_{tk}$  be the indices of  $v_t$  with the largest  $k$  absolute values

        Compute  $1 \leq k \leq d - 1$  such that

$\|\text{TRC}(x_t, A_{tk})\|_q / \|\text{TRC}(x_t, A_{tk})\|_2 \leq R_q^{1/q}$  and

$\|\text{TRC}(x_t, A_{t(k+1)})\|_q / \|\text{TRC}(x_t, A_{t(k+1)})\|_2 > R_q^{1/q}$

$\theta_{q,t} = \text{TRC}(x_t, A_{tk}) / \|\text{TRC}(x_t, A_{tk})\|_2$

**end if**

$t \leftarrow t + 1$

**until** Convergence

---

More specifically, we use the classical power method, but in each iteration  $t$  we project the intermediate vector  $x_t$  to the intersection of the  $d$ -dimension sphere  $\mathbb{S}^{d-1}$  and the  $\ell_q$  ball with the radius  $R_q^{1/q}$ . The idea is to sort  $x_t$  from the highest to the lowest and find the highest  $k$  absolute values and truncate all the others to zero, such

that the resulting vector lies in  $\mathbb{S}^{d-1} \cap \mathbb{B}_q(R_q)$  and is closest to the boundary of  $\mathbb{B}_q(R_q)$ .

For any vector  $v = (v_1, \dots, v_d)^T$  and a index set  $J \subset \{1, \dots, d\}$ , we define the truncation function TRC to be

$$\text{TRC}(v, J) := (v_1 \cdot I(1 \in J), \dots, v_d \cdot I(d \in J))^T, \quad (3.5)$$

where  $I(\cdot)$  is the indicator function. Realizing that for any  $p > q > 0$  and  $v \in \mathbb{R}^d$ ,  $\|v\|_p \leq \|v\|_q \leq n^{1/q-1/p} \|v\|_p$ , we have that the  $\ell_q$  ball constraint is only active when  $R_q \leq d^{1-\frac{q}{2}}$ . In practice,  $R_q$  can be regarded as a tuning parameter. Lemma 3.1 states that, when  $R_q > 1$ , in each step of the iteration there exists a unique solution. In the following  $a^{1/0} := a$  for any  $a \in \mathbb{R}$ .

**Lemma 3.1.** Given  $v := (v_1, \dots, v_d)^T$  with

$$v_1 \geq v_2 \geq \dots \geq v_d \geq 0 \quad \text{and} \quad A_k = \{1, \dots, k\}$$

the for any  $0 < q \leq 1$  and  $k \in \{1, \dots, d-1\}$ ,

$$\frac{\|v_{A_{k+1}}\|_q}{\|v_{A_{k+1}}\|_2} \geq \frac{\|v_{A_k}\|_q}{\|v_{A_k}\|_2} \geq 1. \quad (3.6)$$

When  $q = 0$ , qTPM reduces to TPower algorithm proposed by [5]. Therefore, we can combine COCA estimation consistency result in the next section with [5, Theorem 1] to obtain a geometric convergence rate. Detailed theoretical analysis will be provided in Section 4. Because our main focus is on COCA instead of the sparse PCA algorithm, the general convergence rate for qTPM will be discussed in another paper. In practice, we use the estimator obtained from SPCA [21] as the initial starting point, as suggested by [5].

### 3.2.4 Generalization to the First $m$ Sparse Eigenvectors

We use the iterative deflation method to learn the first  $m$  instead of the first one leading eigenvectors, following the discussions of [5], [23], [24], [25]. In detail, a matrix  $\hat{\Gamma} \in \mathbb{R}^{d \times s}$  deflates a vector  $v \in \mathbb{R}^d$  and results to a new matrix  $\hat{\Gamma}'$ :

$$\hat{\Gamma}' := (\mathbf{I} - vv^T)\hat{\Gamma}(\mathbf{I} - vv^T). \quad (3.7)$$

In this way,  $\hat{\Gamma}'$  is orthogonal to  $v$ .

### 3.2.5 Projection to the Positive Semi-Definite Matrices Cone

To fit in the convex formulation in sparse PCA like semidefinite relaxation DSPCA [26], we project  $\hat{\mathbf{R}}$  into the cone of the positive semidefinite matrices and find solution  $\bar{R}$  to the following convex optimization problem:

$$\tilde{\mathbf{R}} = \arg \min_{\mathbf{M} \succeq 0} \|\hat{\mathbf{R}} - \mathbf{M}\|_{\max}. \quad (3.8)$$

Here  $\ell_{\max}$  norm is chosen such that the theoretical properties in Lemma 2.5 can be preserved. In particular, we have the following lemma:

**Lemma 3.2.** For all  $t \geq 16\pi\sqrt{\frac{\log d}{n \log 2}}$ , for any  $n \geq \frac{37\pi}{t} + 2$ , the minimizer  $\tilde{\mathbf{R}}$  to Equation (3.8) satisfies the following exponential inequality for all  $1 \leq j, k \leq d$ :

$$\mathbb{P}(|\tilde{\mathbf{R}}_{jk} - \Sigma_{jk}^0| \geq t) \leq 2 \exp\left(-\frac{nt^2}{128\pi^2}\right). \quad (3.9)$$

In practice, the optimization problem in Equation (3.8) can be formulated as the dual of a graphical lasso problem with the smallest possible tuning parameter that still guarantees a feasible solution [2]. And then we define  $\tilde{\mathbf{R}}$  and  $\tilde{\mathbf{S}} := [\tilde{S}_{jk}] = [\hat{\sigma}_j \hat{\sigma}_k \hat{R}_{jk}]$  to be the projected Spearman's rho correlation and covariance matrices. In practice we can always do such a projection and use  $\tilde{\mathbf{R}}$  and  $\tilde{\mathbf{S}}$  instead of  $\hat{\mathbf{R}}$  and  $\hat{\mathbf{S}}$ .

## 4 THEORETICAL PROPERTIES

In this section we provide the theoretical properties of COCA and Copula PCA methods. In particular, we are interested in the high dimensional case when  $d > n$  with both  $d$  and  $n$  increasing.

### 4.1 Rank-Based Correlation and Covariance Matrices Estimation

In this section we state the main result on quantifying the convergence rate of  $\hat{\mathbf{R}}$  to  $\Sigma^0$  and  $\hat{\mathbf{S}}$  to  $\Sigma$ . In particular, we establish the results on the  $\ell_{\max}$  convergence rates of the Spearman's rho correlation and covariance matrices to  $\Sigma$  and  $\Sigma^0$ .

For COCA, Lemma 2.5 is enough. For Copula PCA, however, we still need to quantify the convergence rate of  $\hat{\mathbf{S}}$  to  $\Sigma$ . The key to prove the leading eigenvector can be recovered in a fast rate is to show that the estimated covariance matrix  $\hat{\mathbf{S}}$  converges to  $\Sigma$  in the  $\ell_{\max}$  norm in a fast rate. To this end, we need extra conditions on the unknown transformation functions  $\{f_j\}_{j=1}^d$ . We define the *subgaussian transformation function class*. Let  $(\sigma_1^2, \dots, \sigma_d^2)^T := \text{diag}(\Sigma)$  and

$$\Sigma = \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d \end{pmatrix} \cdot \Sigma^0 \cdot \begin{pmatrix} \sigma_1 & 0 & \dots & 0 \\ 0 & \sigma_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_d \end{pmatrix}.$$

**Definition 4.1.** Let  $Z \in \mathbb{R}$  be a random variable following the standard Gaussian distribution. The subgaussian transformation function class  $\text{TF}(K)$  is defined as the set of functions  $\{g_0 : \mathbb{R} \rightarrow \mathbb{R}\}$  which satisfies that:

$$\mathbb{E}|g_0(Z)|^m \leq \frac{m!}{2} K^m, \quad \forall m \in \mathbb{Z}^+.$$

**Remark 4.2.** Here we note that for any function  $g_0 : \mathbb{R} \rightarrow \mathbb{R}$ , if there exists a constant  $L < \infty$  such that

$$g_0(z) \leq L \quad \text{or} \quad g'_0(z) \leq L \quad \text{or} \quad g''_0(z) \leq L, \quad \forall z \in \mathbb{R}, \quad (4.1)$$

then  $g_0 \in \text{TF}(K)$  for some constant  $K$ . To show that, we have the absolute moments of the standard Gaussian distribution satisfying,  $\forall m \in \mathbb{Z}^+$ :

$$\begin{aligned} \mathbb{E}|Z|^m &\leq (m-1)!! < m!! \\ \mathbb{E}|Z^2|^m &= (2m-1)!! < m! \cdot 2^m. \end{aligned} \quad (4.2)$$

Because  $g_0$  satisfies the condition in Equation (4.1), using Taylor expansion, we have for any  $z \in \mathbb{R}$ ,

$$\begin{aligned} g_0(z) &\leq |g_0(0)| + L \quad \text{or} \quad |g_0(z)| \leq |g_0(0)| + L|z|, \quad \text{or} \\ |g_0(z)| &\leq |g_0(0)| + |g'_0(0)z| + Lz^2. \end{aligned} \quad (4.3)$$

Combining Equations (4.2) and (4.3), we have  $\mathbb{E}|g_0(Z)|^m \leq \frac{m!}{2} K^m$  for some constant  $K$ . This proves the assertion.

Then we have the following result, which states that  $\Sigma$  can also be recovered in the parametric rate. The key of the proof is to show that the marginal sample means and standard deviations of the nonparanormal can converge to the population means and standard deviations in an exponential rate.

**Lemma 4.3.** *Let  $x_1, \dots, x_n$  be  $n$  independent realizations of a random vector  $X$ , where  $X \sim MNPN_d(\mu, \Sigma, f)$ . If  $g := \{g_j = f_j^{-1}\}_{j=1}^d$  satisfies for all  $j = 1, \dots, K$ ,  $g_j^2 \in TF(K)$  where  $K < \infty$  is some constant, we have for any  $1 \leq j, k \leq d$ , for any  $n \geq \frac{37\pi}{t} + 2$ ,*

$$\mathbb{P}(|\hat{\mathbf{S}}_{jk} - \Sigma_{jk}| > t) \leq 2 \exp(-c_1 nt^2), \quad (4.4)$$

$$\mathbb{P}(|\hat{\mu}_j - \mu_j| > t) \leq 2 \exp(-c_2 nt^2), \quad (4.5)$$

where  $c_1$  and  $c_2$  are two constants only depending on the choice of  $K$ .

**Remark 4.4.** Lemma 4.3 claims that, under certain constraint on the transformation functions, the latent covariance matrix  $\Sigma$  can be recovered using the Spearman's rho covariance matrix. However, in this case, the marginal distributions of the nonparanormal are required to be sub-gaussian and cannot be arbitrarily continuous. This makes Copula PCA a less favored method compared with COCA.

## 4.2 COCA and Copula PCA

In this section we provide the main result on the upper bound of the estimation error of COCA estimators and Copula PCA estimators. We say that the model  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$  holds if the data are drawn from an element in the model  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$ ; We say that the model  $\mathcal{M}(q, R_q, \Sigma, f)$  holds if the data are drawn from an element in the model  $\mathcal{M}(q, R_q, \Sigma, f)$ .

The next theorem provides an upper bound on the angle between the global optimum  $\tilde{\theta}_1$  to Equation (3.2) and the true parameter  $\theta_1$ .

**Theorem 4.5.** *Let  $\tilde{\theta}_1$  be the global optimum in Equation (3.2) and the model  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$  holds. For any two vectors  $v_1 \in \mathbb{S}^{d-1}$  and  $v_2 \in \mathbb{S}^{d-1}$ , let  $|\sin \angle(v_1, v_2)| := \sqrt{1 - (v_1^\top v_2)^2}$ . Then we have, for any  $n \geq \frac{21}{\log d} + 2$ , with probability at least  $1 - 1/d^2$ ,*

$$\sin^2 \angle(\tilde{\theta}_1, \theta_1) \leq \gamma_q R_q^2 \left( \frac{64\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}},$$

where  $\gamma_q = 2 \cdot I(q = 1) + 4 \cdot I(q = 0) + (1 + \sqrt{3})^2 \cdot I(0 < q < 1)$  and  $\lambda_j = \Lambda_j(\Sigma^0)$  for  $j = 1, 2$ .

**Proof.** The key idea of the proof is to utilize the  $\ell_{\max}$  norm convergence result of  $\hat{\mathbf{R}}$  to  $\Sigma^0$  as shown in Lemma 2.5,

then apply the proof of [7, Theorem 2.2]. For self-containedness, a proof is provided in Section B.4.  $\square$

It can be observed that the convergence rate of  $\tilde{\theta}_1$  to  $\theta_1$  will be faster when  $\theta_1$  lies in a more sparse ball. It makes sense because the effect of "the curse of dimensionality" will be decreasing when the parameters are more and more sparse. Generally, when  $R_q$  and  $\lambda_1, \lambda_2$  do not scale with  $(n, d)$ , the rate is  $O_P((\frac{\log d}{n})^{1-q/2})$ , which is the parametric rate [6], [7], [19] obtains.

Given Theorem 4.5, we can immediately obtain the following corollary, which quantifies the expected angle between  $\tilde{\theta}_1$  and  $\theta_1$ .

**Corollary 4.6.** *In the conditions of Theorem 4.5, we have*

$$\mathbb{E} \sin^2 \angle(\tilde{\theta}_1, \theta_1) \leq \gamma_q R_q^2 \left( \frac{64\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}} + \frac{1}{d^2}.$$

**Proof.** Define  $\epsilon = \sin \angle(\tilde{\theta}_1, \theta_1)$ . Because  $\sin^2(\cdot) \in [0, 1]$ , using Theorem 4.5, we have

$$\begin{aligned} \mathbb{E} \epsilon^2 &= \mathbb{E} \left[ \epsilon^2 \cdot I \left( \epsilon^2 \leq \gamma_q R_q^2 \left( \frac{64\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}} \right) \right] \\ &\quad + \mathbb{E} \left[ \epsilon^2 \cdot I \left( \epsilon^2 > \gamma_q R_q^2 \left( \frac{64\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}} \right) \right] \\ &\leq \gamma_q R_q^2 \left( \frac{64\pi^2}{c_1(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}} + \frac{1}{d^2}. \end{aligned}$$

This completes the proof.  $\square$

In the next corollary, we provide a sparsity recovery consistency result for  $\tilde{\theta}_1$ . It can be observed that the true sparsity pattern can be recovered in a fast rate given a constraint on the minimum absolute value of the signal part of  $\theta_1$ .

**Corollary 4.7.** *Let  $\tilde{\theta}_1$  be the global solution to Equation (3.2) and the model  $\mathcal{M}^0(0, R_0, \Sigma^0, f^0)$  holds. Let  $\Theta^0 := \text{supp}(\theta_1)$  and  $\hat{\Theta}^0 := \text{supp}(\tilde{\theta}_1)$ . If we further have  $\min_{j \in \Theta^0} |\theta_{1j}| \geq \frac{16\sqrt{2}R_0\pi}{\lambda_1 - \lambda_2} \sqrt{\frac{\log d}{n}}$ , then for any  $n \geq \frac{21}{\log d} + 2$ ,  $\mathbb{P}(\hat{\Theta}^0 = \Theta^0) \geq 1 - d^{-2}$ .*

**Proof.** The key of the proof is to construct a contradiction given Theorem 4.5 and the condition on the minimum absolute value of nonzero entries of  $\theta_1$ . Detailed proof can be found in Section B.5.  $\square$

Similarly, we can give an upper bound for the estimation rate of the Copula PCA estimator  $\tilde{u}_1$  to the true leading eigenvalue  $u_1$  of the latent covariance matrix  $\Sigma$ . The next theorem provides the detail result.

**Theorem 4.8.** *Let  $\tilde{u}_1$  be the global solution to Equation (3.4) and the model  $\mathcal{M}(q, R_q, \Sigma, f)$  holds. If  $g := \{g_j = f_j^{-1}\}_{j=1}^d$  satisfies  $g_j^2 \in TF(K)$  for all  $1 \leq j \leq d$ , then we have, for any  $n \geq \frac{21}{\log d} + 2$ , with probability at least  $1 - 1/d^2$ ,*

$$\sin^2 \angle(\tilde{u}_1, u_1) \leq \gamma_q R_q^2 \left( \frac{4}{c_1(\omega_1 - \omega_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}},$$

where  $\gamma_q = 2 \cdot I(q = 1) + 4 \cdot I(q = 0) + (1 + \sqrt{3})^2 \cdot I(0 < q < 1)$ ,  $\omega_j = \Lambda_j(\Sigma)$  for  $j = 1, 2$  and  $c_1$  is a constant defined in Equation (4.4), only depending on  $K$ .

**Proof.** Under the conditions that  $g := \{g_j = f_j^{-1}\}_{j=1}^d$  satisfies  $g_j^2 \in TF(K)$  for all  $1 \leq j \leq d$ , we can utilize Lemma 4.3 and have that

$$\mathbb{P}(|\hat{\mathbf{S}}_{jk} - \Sigma_{jk}| > t) \leq 2 \exp(-c_1 n t^2), \quad \forall j, k \in \{1, \dots, d\}.$$

Using this key observation, all the proofs in Theorem 4.5 can still proceed until Equation (B.19). In particular, let  $\epsilon_u := \sin \angle(\mathbf{u}_1, \tilde{\mathbf{u}}_1)$ , we have

$$\begin{aligned} \mathbb{P}(\epsilon_u^2 \geq t) &\leq \mathbb{P}\left(\frac{\gamma_q R_q^2}{(\omega_1 - \omega_2)^{2-q}} \|\text{vec}(\hat{\mathbf{S}} - \Sigma)\|_\infty^{2-q} \geq t\right) \\ &= \mathbb{P}\left(\|\hat{\mathbf{S}} - \Sigma\|_{\max} \geq \left(\frac{t(\omega_1 - \omega_2)^{2-q}}{\gamma_q R_q^2}\right)^{1/(2-q)}\right) \\ &\leq d^2 \exp\left(-c_1 n \left(\frac{t(\lambda_1 - \lambda_2)^{2-q}}{\gamma_q R_q^2}\right)^{2/(2-q)}\right). \end{aligned}$$

Choosing  $t = \gamma_q R_q^2 \left(\frac{4}{c_1(\omega_1 - \omega_2)^2} \frac{\log d}{n}\right)^{\frac{2-q}{2}}$ , we have the result.  $\square$

Given Theorem 4.8, we can immediately obtain the following corollary, which bounds the expected angle between  $\tilde{\mathbf{u}}_1$  and  $\mathbf{u}_1$ .

**Corollary 4.9.** *In the conditions of Theorem 4.8, we have*

$$\mathbb{E} \sin^2 \angle(\tilde{\mathbf{u}}_1, \mathbf{u}_1) \leq \gamma_q R_q^2 \left(\frac{4}{c_1(\omega_1 - \omega_2)^2} \cdot \frac{\log d}{n}\right)^{\frac{2-q}{2}} + \frac{1}{d^2}.$$

**Proof.** Using the same techniques in proving Corollary 4.6.  $\square$

Similarly, we can prove that under mild conditions  $\tilde{\mathbf{u}}_1$  can recover the support set of  $\mathbf{u}_1$ .

**Corollary 4.10.** *Let  $\tilde{\mathbf{u}}_1$  be the global solution to Equation (3.4) and the model  $\mathcal{M}(0, R_0, \Sigma, f)$  holds. Let  $\Theta := \text{supp}(\mathbf{u}_1)$  and  $\tilde{\Theta} := \text{supp}(\tilde{\mathbf{u}}_1)$ . If  $g := \{g_j = f_j^{-1}\}_{j=1}^d$  satisfies  $g_j^2 \in TF(K)$  for all  $1 \leq j \leq d$  and we further have  $\min_{j \in \Theta} |u_{1j}| \geq \frac{4\sqrt{2}R_0}{\sqrt{c_1}(\omega_1 - \omega_2)}$ , then for any  $n \geq \frac{21}{\log d} + 2$ ,  $\mathbb{P}(\tilde{\Theta} = \Theta) \geq 1 - \frac{1}{d^2}$ .*

**Proof.** Using the same techniques in proving Corollary 4.7.  $\square$

**Remark 4.11.** Assuming that the transformation function  $g$  satisfies that  $g_j^2 \in TF(K)$  for  $j = 1, \dots, d$  restricts the distribution families of the nonparanormal. We note that this constraint is close to claiming that the marginal distributions of the random vector  $X$  have sub-gaussian tails. However, Copula PCA is still an interesting procedure in estimating the leading eigenvectors in the sense that it provides a sparse PCA approach on a model strictly larger than the Gaussian, while consistently and robustly estimating the true latent leading eigenvector in a fast rate.

Let  $\check{\theta}_1$  denote the estimator derived using the Truncated Power method, as shown in Algorithm 3.1 by setting  $q = 0$  and the input matrix  $\Gamma$  to be  $\hat{\mathbf{S}}$ . In the next theorem we show that, under mild conditions,  $\check{\theta}_1$  can approach  $\theta_1$  in a fast near-optimal rate.

**Theorem 4.12.** *Let the tuning parameter in TPower be denoted by  $k := \text{card}(\text{supp}(\check{\theta}_1))$  such that  $k \geq 4R_0$  and the initial*

*starting point be denoted by  $v_0$  with  $\text{card}(\text{supp}(v_0)) \leq k$  and  $\|v_0\|_2 = 1$ . Let*

$$\begin{aligned} v_1 &:= \frac{n\lambda_2 + 8\pi(R_0 + 2k)\sqrt{\log d}}{n\lambda_1 - 8\pi(R_0 + 2k)\sqrt{\log d}} \quad \text{and} \\ v_2 &:= 8\sqrt{2}\pi(R_0 + 2k)\sqrt{\log d} \cdot (n(\lambda_1 - \lambda_2))^2 \\ &\quad - 32\pi(\lambda_1 - \lambda_2)(R_0 + 2k)\sqrt{n \log d} \\ &\quad + 320\pi^2(R_0 + 2k)^2 \log d)^{1/2}. \end{aligned}$$

*If the model  $\mathcal{M}^0(0, R_0, \Sigma^0, f^0)$  holds and the following three assumptions hold:*

- (A1)  $\lambda_1$  and  $\lambda_2$  scale with  $(n, d)$  such that  $\lambda_1 - \lambda_2 \geq 16\pi(R_0 + 2k)\sqrt{\frac{\log d}{n}}$ ;
- (A2)  $(1 + 3\sqrt{R_0/k})(1 - 0.45(1 - v_1^2)) < 1$ ;
- (A3) Letting  $\zeta_1 := |\theta_1^T v_0| - v_2$  be a fixed constant in  $[0, 1]$ , we have  $0 < (1 - v_1^2)\zeta_1(1 - \zeta_1^2)/2 - 2\nu_2 - \sqrt{R_0/k} < 1$ ,

We have, with probability larger than  $1 - d^{-2}$ ,

$$|\sin \angle(\check{\theta}_1, \theta_1)| \leq \frac{C}{\lambda_1 - \lambda_2} (R_0 + 2k) \cdot \sqrt{\frac{\log d}{n}},$$

for some generic constant  $C$  not scaled with  $(n, d)$ .

**Proof.** The key of the proof is to show that

$$\max_{\mathbf{v} \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(R_0+2k)} |\mathbf{v}^T (\hat{\mathbf{R}} - \Sigma^0) \mathbf{v}| \leq 8\pi(R_0 + 2k) \sqrt{\frac{\log d}{n}}$$

with large probability and for any  $\mathbf{v}_1, \mathbf{v}_2 \in \mathbb{S}^{d-1}$ ,

$$\sqrt{1 - |\mathbf{v}_1^T \mathbf{v}_2|} \leq |\sin \angle(\mathbf{v}_1, \mathbf{v}_2)| \leq 2\sqrt{1 - |\mathbf{v}_1^T \mathbf{v}_2|}.$$

Detailed proof can be found in Section B.6.  $\square$

**Remark 4.13.** Here Assumption (A1) is to control the difference between the top two leading eigenvalues  $\lambda_1$  and  $\lambda_2$ , such that  $\theta_1$  can be differentiated from  $\theta_2$ . Assumption (A3) is to control the closedness of the initial value  $v_0$  to  $\theta_1$ . This assumption makes sense because Truncated Power method is a nonconvex formulation in estimating  $\theta_1$ . The theory verifies that, when assumptions hold, the obtained estimator  $\check{\theta}_1$  can obtain the same convergence rate as the global optimum  $\theta_1$ .

Let  $\check{\mathbf{u}}_1$  denote the estimator derived using the Truncated Power method(TPower), as shown in Algorithm 3.1 by setting  $q = 0$  and the input matrix  $\Gamma$  to be  $\hat{\mathbf{S}}$ . In the next theorem we show that, under mild conditions,  $\check{\mathbf{u}}_1$  can approach  $\mathbf{u}_1$  in a fast near-optimal rate.

**Theorem 4.14.** *Let the tuning parameter in TPower be denoted by  $k := \text{card}(\text{supp}(\check{\mathbf{u}}_1))$  such that  $k \geq 4R_0$  and the initial starting point be denoted by  $v_0$  with  $\text{card}(\text{supp}(v_0)) \leq k$  and  $\|v_0\|_2 = 1$ . Let*

$$\begin{aligned} v_3 &:= \frac{n\lambda_2\sqrt{c_1} + 2(R_0 + 2k)\sqrt{\log d}}{n\lambda_1\sqrt{c_1} - 2(R_0 + 2k)\sqrt{\log d}} \quad \text{and} \\ v_4 &:= 2\sqrt{2}(R_0 + 2k)\sqrt{\log d} \cdot (nc_1(\lambda_1 - \lambda_2))^2 \\ &\quad - 8(\lambda_1 - \lambda_2)(R_0 + 2k)\sqrt{nc_1 \log d} + 20(R_0 + 2k)^2 \log d)^{1/2}. \end{aligned}$$

If the model  $\mathcal{M}(0, R_0, \Sigma, f)$  holds and the following three assumptions hold:

- (B1)  $\omega_1$  and  $\omega_2$  scale with  $(n, d)$  such that  $\omega_1 - \omega_2 \geq 4(R_0 + 2k)\sqrt{\frac{\log d}{nc_1}}$ ;
- (B2)  $(1 + 3\sqrt{R_0/k})(1 - 0.45(1 - v_3^2)) < 1$ ;
- (B3) Letting  $\zeta_2 := |\mathbf{u}_1^T \mathbf{v}_0| - v_4$  be a fixed constant in  $[0, 1]$ , we have  $0 < (1 - v_3^2)\zeta_2(1 - \zeta_2^2)/2 - 2v_4 - \sqrt{R_0/k} < 1$ ,

We have, with probability larger than  $1 - d^{-2}$ ,

$$|\sin \angle(\check{\mathbf{u}}_1, \mathbf{u}_1)| \leq \frac{C}{\omega_1 - \omega_2} (R_0 + 2k) \cdot \sqrt{\frac{\log d}{n}},$$

for some generic constant  $C$  not scaled with  $(n, d)$ .

### 4.3 Discussion on the Attainability of the Optimum

In Section 3.1.3 we show that the optimum to Equations (3.2) and (3.4) are hard to compute. To approximate the global optimum  $\tilde{\theta}_1$  and  $\tilde{\mathbf{u}}_1$ , we advocate using the Truncated Power method [5] and provide the theoretical analysis for the corresponding algorithm, shown in Theorems 4.12 and 4.14. To guarantee convergence of the proposed algorithm, we need to make sure that the initial vector  $\mathbf{v}_0$  is not too far away from the true vector  $\theta_1$  or  $\mathbf{u}_1$ . In this section we discuss two approaches in finding such a vector  $\mathbf{v}_0$  in light of the arguments in [5]:

(i) As suggested by [5] (Paragraph 2, Page 905), to find a proper initial vector  $\mathbf{v}_0$ , we can take a relatively large pilot tuning parameter  $\bar{k}$  so that the requirement on  $\theta_1^T \mathbf{v}_0 \gtrsim \sqrt{R_0/\bar{k}}$  is easier to be satisfied. Using  $\bar{k}$  we get a pilot estimator  $\bar{\mathbf{v}}$  and then plug it into the qTPM algorithm with a smaller tuning parameter  $k$ . [5] provided some theoretical justification for this procedure. They also provided thorough numerical experiments to show that this approach is practically effective in application.

(ii) An alternative way to choose the initial vector  $\mathbf{v}_0$  is to exploit the estimator obtained from other sparse PCA algorithms to initialize qTPM. For example, we can plug the Spearman's rho correlation and covariance matrices  $\hat{\mathbf{R}}$  and  $\hat{\mathbf{S}}$  into the Sparse PCA algorithm with the semidefinite programming formulation [26] (We call it the SDP algorithm). From the theory of [5], we know that if the SDP procedure provides a consistent estimator of  $\theta_1$ , we could use the SDP estimator to initialize the qTPM algorithm and achieve the desired rate.

### 4.4 Discussion on the Optimal Rate of Convergence of COCA

Many results have been established in understanding the sparse PCA problem. For example, under the Gaussian assumptions, [27] discussed the problem of support recovery of leading eigenvectors, [28] discussed the problem of sparse principal component detection and [7] proposed methods that obtain a  $\sqrt{R_0 \log d/n}$  rate of convergence for parameter estimation when  $\mathbf{u}_1$  is sparse with support set size  $R_0$  and showed that this rate is minimax optimal confined in the Gaussian family.

COCA is significantly different from the procedures in the above mentioned papers in the sense that: (i) In

methodology, we suggest using the Spearman's rho correlation matrix  $\hat{\mathbf{R}}$  to estimate  $\Sigma^0$ , instead of using the sample correlation matrix  $\mathbf{S}^0$ . Empirical results in the next section show that rank-based methods are more robust to modeling and data contaminations than the methods based on the Pearson sample correlation matrix. (ii) In theory, in terms of modeling flexibility, COCA gains more compared with the results in [6], [7], [19]: The nonparanormal family contains many heavy-tailed distributions with arbitrary margins, which cannot be handled by the Gaussian-based procedures. COCA is the optimal method when  $R_0$  is fixed. When not, it is unclear whether COCA is the optimal method confined in the nonparanormal family.

Addressing the optimal rate of convergence of COCA is challenging due to the reason that the data can be very heavy-tailed and the transformed rank-based correlation matrix has a much more complex structure than the Pearson's covariance/correlation matrix.

However, here we lay out a venue in attempt to prove a sharper rate of convergence of COCA. More specifically, we prove that COCA can attain the parametric  $\sqrt{R_0 \log d/n}$  rate of convergence if a condition called "third-order sign subgaussian condition" holds for the nonparanormally distributed random vector  $X$ .

**Definition 4.15 (Third-Order Sign Subgaussian Condition).** Let  $X_1$  be a random vector and  $X_2, X_3$  be two independent copies of  $X_1$ . For any random vector  $v \in \mathbb{S}^{d-1}$ , we let  $\mathbf{O} \in \mathcal{R}^{d \times d}$  be the population-wise Spearman's rho matrix with

$$\mathbf{O} := 3 \cdot \mathbb{E}\{\text{sign}(X_1 - X_2)(\text{sign}(X_1 - X_3))^T\},$$

and let

$$Y_v := v^T \text{sign}(X_1 - X_2)(\text{sign}(X_1 - X_3))^T v.$$

Then  $X_1$  is said to satisfy the third-order sign subgaussian condition if and only if there exists an absolute constant  $c$  such that for any  $v \in \mathbb{S}^{d-1}$

$$\mathbb{E} \exp\{t(Y_v - \mathbb{E} Y_v)\} \leq \exp(c(\|\Sigma_0\|_2 + \|\mathbf{O}\|_2)^2 t^2), \quad \text{for } |t| < t_0,$$

where  $t_0$  is a positive number such that  $t_0(\|\Sigma_0\|_2 + \|\mathbf{O}\|_2)^2$  is lower bounded by a fixed constant.

We then have the following theorem, which states that we can recover  $\theta_1$  in the parametric rate of convergence when  $X$  satisfies the third-order sign subgaussian condition.

**Theorem 4.16.** When the model  $\mathcal{M}^0(0, R_0, \Sigma^0, f^0)$  holds and the nonparanormally distributed random vector  $X$  satisfies Equation (4.6), we have

$$|\sin \angle(\hat{\theta}_1, \theta_1)| = O_P\left(\frac{\lambda_1 + \|\mathbf{O}\|_2}{\lambda_1 - \lambda_2} \sqrt{\frac{R_0 \log d}{n}}\right).$$

Theorem 4.16 can be shown to be correct in three steps and we sketch the proof as follows.

(i) By using the argument in [2] (Page 2319), we have

$$\hat{\rho}_{jk} = \frac{n-2}{n-1} U_{jk} + \frac{3}{n+1} \hat{\tau}_{jk},$$

where  $\hat{\tau}_{jk} \in [-1, 1]$  is the Kendall's tau correlation coefficient and

$$U_{jk} = \frac{3}{n(n-1)(n-2)} \sum_{i \neq s \neq t} \text{sign}(x_{ij} - x_{sj})(x_{ik} - x_{tk}).$$

(ii) We only focus on  $U_{jk}$  and then following the proof of Lemma 5.4 in [29] until Equation (5.21), where we substitute Equation (5.22) by (4.6), we can prove that

$$\|\hat{\mathbf{O}} - \mathbf{O}\|_2 = O_P\left((\lambda_1 + \|\mathbf{O}\|_2)\sqrt{\frac{R_0 \log d}{n}}\right),$$

where  $\hat{\mathbf{O}}$  is the empirical realization of  $\mathbf{O}$  with  $\hat{\mathbf{O}}_{jk} = \hat{\rho}_{jk}$  for  $j, k \in \{1, \dots, d\}$ .

(iii) Combining with the proof of Lemma C.2 in [30], we can show that the  $\sin(\cdot)$  transformation in  $\hat{\mathbf{R}}$  does not hurt the rate and hence we have

$$\|\hat{\mathbf{R}} - \Sigma_0\|_2 = O_P\left((\lambda_1 + \|\mathbf{O}\|_2)\sqrt{\frac{R_0 \log d}{n}}\right).$$

This completes the proof.

## 5 EXPERIMENTS

In this section we investigate the empirical performance of the COCA method. Three sparse PCA algorithms are considered: penalized matrix decomposition proposed by [20], SPCA proposed by [21] and Truncated Power method (TPower) proposed by [5]. The following three methods are considered:

- Pearson: the sparse PCA algorithm using the Pearson sample correlation matrix;
- Spearman: the sparse PCA algorithm using the Spearman's rho correlation matrix;
- Oracle: the sparse PCA algorithm using the Pearson sample correlation matrix of the latent Gaussian data (perfect without data contamination).

### 5.1 Numerical Simulations

In the simulation study we study the empirical performance for support recovery and parameter estimation for different estimators where samples are drawn from an element of the model  $\mathcal{M}^0(0, R_0, \Sigma^0, f^0)$ .

In detail, we sample  $n$  data points  $x_1, \dots, x_n$  from the nonparanormal distribution  $X \sim NPN_d(\Sigma^0, f^0)$ . Here we set  $d = 100$ . We follow the same generating scheme as in [31] and [5]. A covariance matrix  $\Sigma$  is first synthesized through the eigenvalue decomposition, where the first two eigenvalues are given and the corresponding eigenvectors are pre-specified to be sparse. In detail, we suppose that the first two leading eigenvectors of  $\Sigma$ ,  $u_1$  and  $u_2$ , are sparse in the sense that only the first  $s = 10$  entries of  $u_1$  and the second  $s = 10$  entries of  $u_2$  nonzero, i.e.,

$$u_{1j} = \begin{cases} \frac{1}{\sqrt{10}}, & 1 \leq j \leq 10, \\ 0, & \text{otherwise,} \end{cases} \quad \text{and } u_{2j} = \begin{cases} \frac{1}{\sqrt{10}}, & 11 \leq j \leq 20, \\ 0, & \text{otherwise,} \end{cases}$$

and  $\omega_1 = 5$ ,  $\omega_2 = 2$ ,  $\omega_3 = \dots = \omega_d = 1$ . The remaining eigenvectors are chosen arbitrarily. The correlation matrix  $\Sigma^0$  is accordingly generated from  $\Sigma$ , with  $\lambda_1 = 4$ ,  $\lambda_2 = 2.5$ ,

$\lambda_3, \dots, \lambda_d \leq 1$  and the two leading eigenvectors sparse:

$$\theta_{1j} = \begin{cases} \frac{-1}{\sqrt{10}}, & 1 \leq j \leq 10, \\ 0, & \text{otherwise,} \end{cases} \quad \text{and } \theta_{2j} = \begin{cases} \frac{-1}{\sqrt{10}}, & 11 \leq j \leq 20, \\ 0, & \text{otherwise.} \end{cases}$$

To sample data from the nonparanormal distribution, we also need the transformation functions:  $f^0 = \{f_j^0\}_{j=1}^d$ . Here two types of transformation functions are considered:

- Linear transformation (or no transformation):

$$f_{\text{linear}}^0 = \{h_0, h_0, \dots, h_0\}, \quad \text{where } h_0(x) := x.$$

- Nonlinear transformation: there exist five univariate monotone functions  $h_1, h_2, \dots, h_5 : \mathbb{R} \rightarrow \mathbb{R}$  and

$$f_{\text{nonlinear}}^0 = \{h_1, h_2, h_3, h_4, h_5, h_1, h_2, h_3, h_4, h_5, \dots\},$$

where

$$h_1^{-1}(x) := x, \quad h_2^{-1}(x) := \frac{\text{sign}(x)|x|^{1/2}}{\sqrt{\int |t|\phi(t) dt}},$$

$$h_3^{-1}(x) := \frac{x^3}{\sqrt{\int t^6 \phi(t) dt}},$$

$$h_4^{-1}(x) := \frac{\Phi(x) - \int \Phi(t)\phi(t) dt}{\sqrt{\int (\Phi(y) - \int \Phi(t)\phi(t) dt)^2 \phi(y) dy}},$$

and

$$h_5^{-1}(x) := \frac{\exp(x) - \int \exp(t)\phi(t) dt}{\sqrt{\int (\exp(y) - \int \exp(t)\phi(t) dt)^2 \phi(y) dy}}.$$

Here  $\phi$  and  $\Phi$  are defined to be the probability density and cumulative distribution functions of the standard Gaussian. We then generate  $n = 100, 200$  or  $500$  data points from:

- [Scheme 1]  $X \sim NPN_d(\Sigma^0, f_{\text{linear}}^0)$  where  $f_{\text{linear}}^0 = \{h_0, h_0, \dots, h_0\}$  and  $\Sigma_0$  is defined as above.
- [Scheme 2]  $X \sim NPN_d(\Sigma^0, f_{\text{nonlinear}}^0)$  where  $f_{\text{nonlinear}}^0 = \{h_1, h_2, h_3, h_4, h_5, \dots\}$  and  $\Sigma_0$  is defined as above.

To evaluate the robustness of different methods, we adopt a similar data contamination procedure as in [2]. Let  $r \in [0, 1]$  represent the proportion of samples being contaminated. For each dimension, we randomly select  $\lfloor nr \rfloor$  entries and replace them with either  $5$  or  $-5$  with equal probability. The final data matrix we obtained is  $\mathbf{X} \in \mathbb{R}^{n \times d}$ . PMD, SPCA and TPower are then employed on  $\mathbf{X}$  to computer the estimated leading eigenvector  $\tilde{\theta}_1$ .

To evaluate the empirical variable selection property of different methods, we define

$$\mathcal{S} := \{1 \leq j \leq d : \theta_{1j} \neq 0\}, \\ \hat{\mathcal{S}}_\delta := \{1 \leq j \leq d : \tilde{\theta}_{1j} \neq 0\},$$

to be the support sets of the true leading eigenvector  $\theta_1$  and the estimated leading eigenvector  $\tilde{\theta}_1$  using the tuning

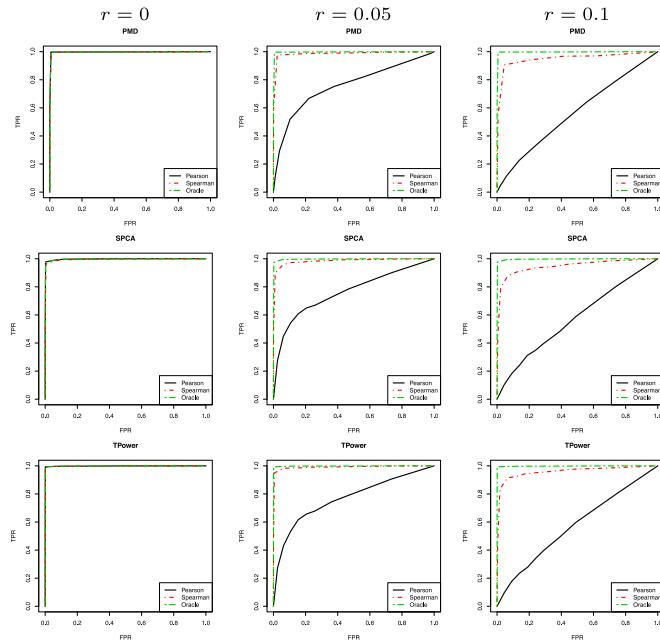


Fig. 2. ROC curves for PMD, SPCA and Truncated Power method (top, middle, bottom) with linear (no) transformation and data contamination at different levels ( $r = 0, 0.05, 0.1$ ). Here  $n = 100$  and  $d = 100$ .

parameter  $\delta$ . In this way, the false positive number (FPN) and false negative number (FNN) of  $\delta$  are defined as:

$$\begin{aligned} \text{FPN}(\delta) &:= \text{the number of features in } \hat{\mathcal{S}}_\delta \text{ not in } \mathcal{S}, \\ \text{FNN}(\delta) &:= \text{the number of features in } \mathcal{S} \text{ not in } \hat{\mathcal{S}}_\delta. \end{aligned}$$

Then we can further define the false positive rate (FPR) and false negative rate (FNR) corresponding to the tuning parameter  $\delta$  to be

$$\begin{aligned} \text{FPR}(\delta) &:= \text{FPN}(\delta)/(d - s), \\ \text{FNR}(\delta) &:= \text{FNN}(\delta)/s. \end{aligned}$$

Under the Scheme 1 and Scheme 2 with different levels of contamination ( $r = 0, 0.05$  or  $0.1$ ), we repeatedly generate the data matrix  $\mathbf{X}$  for 1,000 times and compute the averaged FPR and False Negative Rates using a path of tuning parameters  $\delta$ . The feature selection performances of different methods are then evaluated by plotting  $(\text{FPR}(\delta), 1 - \text{FNR}(\delta))$ . The corresponding ROC curves are presented in Figs. 2 and 3.

In Fig. 2, Scheme 1 is explored and it can be observed that under the most ideal case where there is no contamination ( $r = 0$ ) and  $\mathbf{X}$  is exactly Gaussian, Pearson, Spearman and Oracle can all recover the sparsity pattern perfectly.

However, when the data are contaminated where outliers exist, the performances of Pearson utilizing PMD, SPCA and TPower significantly decrease, while the rank-based method Spearman is still very close to Oracle.

In Fig. 3, Scheme 2 is explored and  $\mathbf{X}$  follows a nonparametric distribution and is non-Gaussian. It can be observed that, in Scheme 2, even without data contamination ( $r = 0$ ), Pearson cannot recover the support set of  $\theta_1$ , while Spearman can still recover the sparsity pattern almost perfectly. When the data are contaminated where outliers exist,

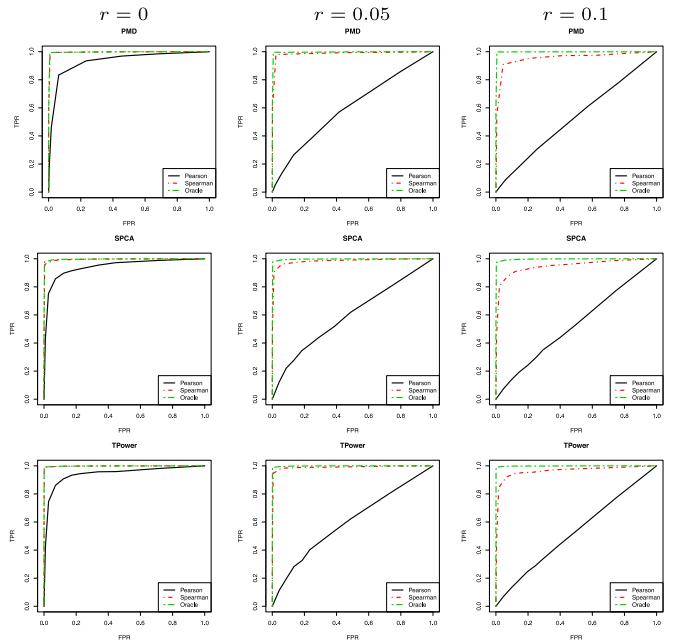


Fig. 3. ROC curves for PMD, SPCA and Truncated Power method (top, middle, bottom) with nonlinear transformation and data contamination at different levels ( $r = 0, 0.05, 0.1$ ). Here  $n = 100$  and  $d = 100$ .

the performance of the rank-based method Spearman utilizing PMD, SPCA and TPower is still very close to Oracle.

To explore the empirical performances of difference methods using different algorithms more, we define an oracle tuning parameter  $\delta^*$  to be the  $\delta$  with the lowest  $\text{FPR}(\delta) + \text{FNR}(\delta)$ :  $\delta^* := \arg \min_{\delta} (\text{FPR}(\delta) + \text{FNR}(\delta))$ . In this way, an estimator  $\hat{\theta}_1$  using the oracle tuning parameter  $\delta^*$  can be calculated and we computer the angle between  $\theta_1$  and  $\hat{\theta}_1$ :  $\sin \angle(\theta_1, \hat{\theta}_1)$  to quantify the estimation consistency.

In Tables 1 and 2, the averaged  $\sin \angle(\theta_1, \hat{\theta}_1)$  values for Scheme 1 and Scheme 2,  $n = 100, 200, 500$ , contamination levels  $r = 0, 0.05, 0.1$  and utilizing three algorithms (PMD, SPCA and TPower) are presented. There are mainly three observations drawn from the results:

- In the perfectly Gaussian data (Scheme 1 with  $r = 0$ ), Pearson performs slightly better than Spearman. However, the difference is not significantly. When  $r \neq 0$ , Spearman outperforms Pearson significantly.
- In Scheme 2 where the data are non-Gaussian, even when  $r = 0$  and  $n$  is large, Pearson's estimation error is significantly away from zero. Spearman can still achieve good performance here and perform much more robustly when  $r \neq 0$  compared with Pearson.
- In both Scheme 1 and Scheme 2, when  $r = 0$ , Spearman is close to Oracle and is tending to zero when  $n$  is large. When  $r \neq 0$ , the performance of Spearman drops, but significantly less than Pearson.

With regard to the comparison among the three algorithms (PMD, SPCA, TPower), we have two more comments restricted to what we observe:

- PMD's estimator  $\hat{\theta}_1$  seems not converging to  $\theta_1$  in our simulation studies. This might be due to the fact that PMD is more sensitive to the choice of initial values.

**TABLE 1**  
Quantitative Comparison on the Data Set under  
the Generating Scheme 1 with Linear Transformation

method	n	r	Pearson	Spearman	Oracle
PMD	100	0.00	0.2739(0.0337)	0.2806(0.0372)	0.2739(0.0337)
		0.05	0.6909(0.1307)	0.3720(0.0941)	0.2942(0.0620)
		0.10	0.9007(0.0852)	0.4414(0.0952)	0.2742(0.0293)
	200	0.00	0.2576(0.0000)	0.2577(0.0012)	0.2576(0.0000)
		0.05	0.4630(0.0969)	0.2610(0.0110)	0.2576(0.0000)
		0.10	0.7768(0.1089)	0.2871(0.0409)	0.2576(0.0000)
	500	0.00	0.2576(0.0000)	0.2576(0.0000)	0.2576(0.0000)
		0.05	0.2730(0.0234)	0.2576(0.0000)	0.2576(0.0000)
		0.10	0.4651(0.1033)	0.2576(0.0000)	0.2576(0.0000)
SPCA	100	0.00	0.3686(0.0879)	0.3952(0.0885)	0.3686(0.0879)
		0.05	0.6765(0.0910)	0.4412(0.0867)	0.3605(0.0814)
		0.10	0.8660(0.0800)	0.5173(0.0857)	0.3614(0.0977)
	200	0.00	0.1869(0.0489)	0.2060(0.0534)	0.1869(0.0489)
		0.05	0.4335(0.0892)	0.2451(0.0753)	0.1836(0.0523)
		0.10	0.7016(0.0998)	0.3236(0.0863)	0.1874(0.0558)
	500	0.00	0.0762(0.0178)	0.0833(0.0190)	0.0762(0.0178)
		0.05	0.2319(0.0676)	0.1045(0.0274)	0.0807(0.0225)
		0.10	0.3854(0.0925)	0.1362(0.0305)	0.0799(0.0199)
TPower	100	0.00	0.1126(0.0726)	0.1312(0.0913)	0.1126(0.0726)
		0.05	0.6513(0.1175)	0.2423(0.1452)	0.1132(0.0668)
		0.10	0.8726(0.0776)	0.3900(0.1551)	0.1096(0.0637)
	200	0.00	0.0709(0.0151)	0.0761(0.0169)	0.0709(0.0151)
		0.05	0.3730(0.1433)	0.0933(0.0281)	0.0683(0.0176)
		0.10	0.7310(0.0912)	0.1306(0.0714)	0.0667(0.0172)
	500	0.00	0.0424(0.0114)	0.0459(0.0112)	0.0424(0.0114)
		0.05	0.1210(0.0423)	0.0581(0.0120)	0.0420(0.0096)
		0.10	0.3858(0.1349)	0.0694(0.0167)	0.0422(0.0116)

The means of the  $\sin \angle(\theta_1, \theta_1)$  with their standard deviations in parentheses are presented. Here  $n$  is changing from 100 to 500 and  $d = 100$ .

- TPower performs generally better than SPCA. We also find that the computing time of TPower is less than SPCA.

## 5.2 Large-Scale Genomic Data Analysis

In this section we investigate the performance of Spearman compared with Pearson using one of the largest microarray data sets [32]. In summary, we collect in all 13,182 publicly available microarray samples from Affymetrix HGU133a platform. The raw data contain 20,248 probes and 13,182 samples belonging to 2,711 tissue types (e.g., lung cancers, prostate cancer, brain tumor etc.). There are at most 1,599 samples and at least 1 sample belonging to each tissue type. We merge the probes corresponding to the same gene. There are remaining 12,713 genes and 13,182 samples. The main purpose of this experiment is to compare the performance of Spearman with Pearson. We use the Truncated Power method proposed by [5] in this section.

We first show that the data are non-Gaussian. To this end, we randomly pick 16 genes and all samples from a certain tissue type, then the corresponding Quantile-to-Quantile plots (QQ plots) compared with the Gaussian are presented in Fig. 4 to illustrate their normality. It can be observed that all the sixteen marginal distributions are severely away from the Gaussian.

We adopt the same idea of data-preprocessing as in [2]. In particular, we first remove the batch effect by applying the surrogate variable analysis proposed by [33]. We then extract the top 2,000 genes with the highest marginal

**TABLE 2**  
Quantitative Comparison on the Data Set under  
the Generating Scheme 2 with Nonlinear Transformation

method	n	r	Normal	Spearman	Oracle
PMD	100	0.00	0.5076(0.1504)	0.2878(0.0451)	0.2778(0.0361)
		0.05	0.8729(0.1025)	0.3497(0.0820)	0.2814(0.0421)
		0.10	0.9514(0.0584)	0.4338(0.0952)	0.2775(0.0371)
	200	0.00	0.3272(0.0743)	0.2576(0.0000)	0.2576(0.0000)
		0.05	0.6867(0.1359)	0.2610(0.0139)	0.2576(0.0000)
		0.10	0.8910(0.0919)	0.2807(0.0370)	0.2576(0.0000)
	500	0.00	0.2582(0.0036)	0.2576(0.0000)	0.2576(0.0000)
		0.05	0.4439(0.1096)	0.2576(0.0000)	0.2576(0.0000)
		0.10	0.7055(0.1421)	0.2576(0.0000)	0.2576(0.0000)
SPCA	100	0.00	0.5210(0.0961)	0.4005(0.0946)	0.3768(0.1000)
		0.05	0.8453(0.0973)	0.4470(0.0851)	0.3673(0.0819)
		0.10	0.9245(0.0742)	0.5141(0.0949)	0.3556(0.0977)
	200	0.00	0.3583(0.0889)	0.1949(0.0532)	0.1788(0.0489)
		0.05	0.6448(0.1050)	0.2729(0.0782)	0.1847(0.0534)
		0.10	0.8502(0.1097)	0.3212(0.0852)	0.1927(0.0599)
	500	0.00	0.1744(0.0483)	0.0843(0.0252)	0.0780(0.0218)
		0.05	0.3699(0.0979)	0.1053(0.0257)	0.0788(0.0206)
		0.10	0.5546(0.1229)	0.1318(0.0318)	0.0779(0.0174)
TPower	100	0.00	0.4516(0.1216)	0.1346(0.0832)	0.1202(0.0746)
		0.05	0.8315(0.1094)	0.2372(0.1517)	0.1053(0.0513)
		0.10	0.9323(0.0730)	0.3608(0.1583)	0.1088(0.0629)
	200	0.00	0.1942(0.1056)	0.0740(0.0191)	0.0702(0.0190)
		0.05	0.6193(0.1263)	0.0900(0.0313)	0.0661(0.0172)
		0.10	0.8608(0.0926)	0.1266(0.0596)	0.0663(0.0185)
	500	0.00	0.1025(0.0310)	0.0465(0.0101)	0.0437(0.0086)
		0.05	0.3296(0.1293)	0.0586(0.0154)	0.0422(0.0101)
		0.10	0.6296(0.1157)	0.0708(0.0171)	0.0403(0.0099)

The means of the  $\sin \angle(\theta_1, \theta_1)$  with their standard deviations in parentheses are presented. Here  $n$  is changing from 100 to 500 and  $d = 100$ .

standard deviations. There are, accordingly, 2,000 genes left and the data matrix we are focusing is  $2,000 \times 13,182$ .

We then explore several tissue types with the largest sample size:

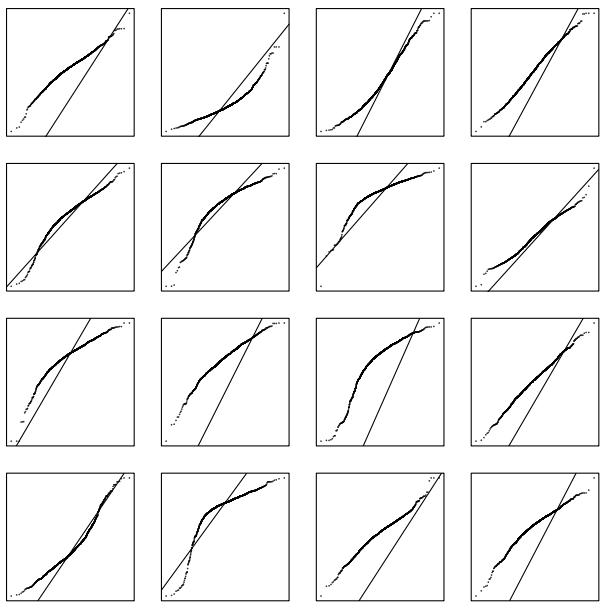


Fig. 4. Sixteen randomly picked genes' Quantile-to-Quantile plots. The x-axis represents the theoretical quantiles and the y-axis represents the sample quantiles.

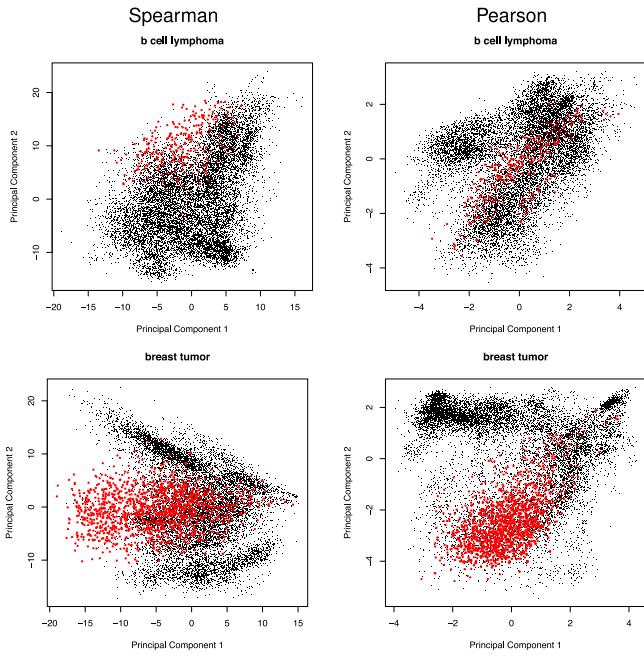


Fig. 5. The scatter plots of the first two principal components of the data set. Spearman and Pearson are compared (left to right) and b cell lymphoma and breast tumor are explored (top to bottom). Each black point represents a sample and each red point represents a sample belonging to the corresponding tissue type.

- Breast tumor, which has 1,599 samples;
- B cell lymphoma, which has 213 samples;
- Prostate tumor, which has 148 samples;
- Wilms tumor, which has 143 samples.

For each tissue type listed above, we apply Spearman and Pearson on the data belonging to this specific tissue type and obtain the first two leading sparse eigenvectors. Here we set  $R_0 = 100$  for both eigenvectors. For Spearman, we do a normal score transformation [15] on the original data set. We subsequently project the whole data set to the first two principal components using the obtained eigenvectors. The according two-dimension visualization is illustrated in Figs. 5 and 6.

In Figs. 5 and 6 each black point represents a sample and each red point represents a sample belonging to the corresponding tissue type. It can be observed that, in 2D plots learnt by Spearman, the red points are averagely more dense and more close to the border of the sample cluster. The first phenomenon indicates that Spearman has the potential to preserve more common information shared by samples from the same tissue type. The second phenomenon indicates that Spearman has the potential to differentiate samples from different tissue types more efficiently.

### 5.3 Brain Imaging Data

In this section we apply Spearman and Pearson to a brain imaging data: The ADHD 200 data set [34]. Here 776 subjects' functional scans were collected, where 491 of which are normal persons and 285 of which are diagnosed attention deficit hyperactive disorder (ADHD). The data are normalized and 264 voxels with biological interests are extracted. These voxels broadly cover the major functional regions of the cerebral cortex and cerebellum. We refer to [34] and [35] for details in

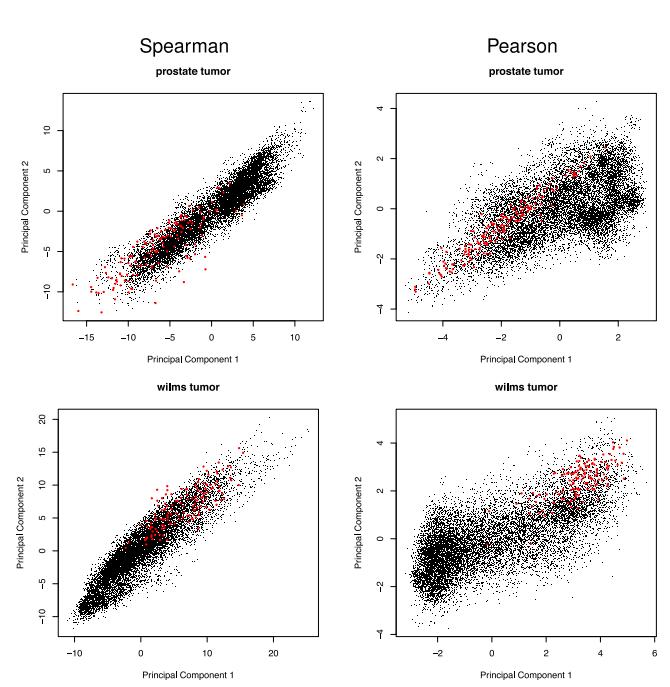


Fig. 6. The scatter plots of the first two principal components of the data set. Spearman and Pearson are compared (left to right) and prostate tumor and Wilms tumor are explored (top to bottom). Each black point represents a sample and each red point represents a sample belonging to the corresponding tissue type.

data preprocessing and voxel definitions. In this manuscript we are only interested in the normal persons, leading to a data matrix with 491 rows and 264 columns.

We apply Spearman and Pearson, with  $R_0$  set to be 20 in each sparse estimated eigenvector, to the ADHD data, and plot the first principal component against the second, third, and fourth principal components. Fig. 7 visualizes the results. Here similar as in Section 5.2, for Spearman, we conduct a normal score transformation on the original data. It can be observed that there are outliers in the principal components calculated by Pearson, which can make the inference based on the principal components very unstable. In contrast, the principal components calculated by Spearman are very concrete and present almost like a bivariate Gaussian distribution.

## 6 CONCLUSION

In this paper we propose a semiparametric scale-invariant principal component analysis named Copula Component Analysis. Several contributions we make include: (i) We generalize the Gaussian assumption used in justifying the high dimensional sparse PCA to the nonparanormal; (ii) We utilize the rank-based nonparametric correlation coefficient estimator, Spearman's rho, in estimating the latent correlation matrix; (iii) We provide sufficient conditions under which the estimation consistency and feature selection consistency for COCA can be achieved; (iv) We also explore sufficient conditions under which Copula PCA can achieve the same theoretical properties as COCA, and discuss the advantages of COCA over Copula PCA; (v) Careful experimental studies are conducted to confirm that COCA outperforms Copula PCA on both synthetic and real-world data sets.

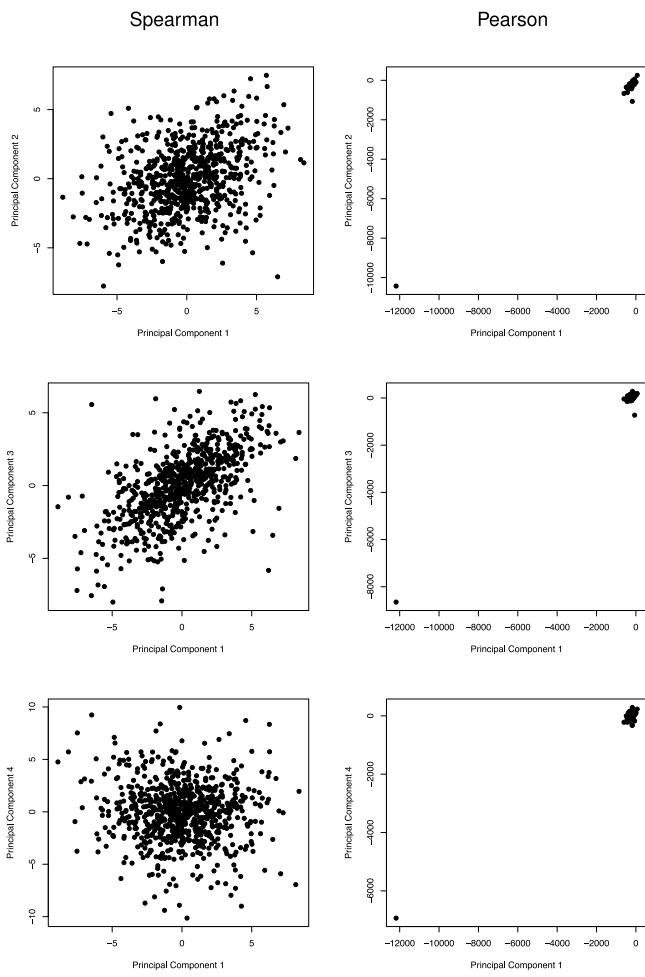


Fig. 7. The scatter plots of the first principal component against the second, third, and fourth principal components (from top to bottom) of the brain imaging data set. Spearman and Pearson are compared (from left to right)

## APPENDIX A SUPPORTING INEQUALITIES

**Lemma A.1.** Let  $\hat{\mathbf{R}}$  be the Spearman's rho correlation matrix. Then for any  $n \geq \frac{37\pi}{t} + 2$ , we have

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq 2 \exp\left(-\frac{nt^2}{16\pi^2}\right).$$

**Proof.** Using the notation and proof of [2, Theorem 4.1], we have whenever  $n \geq \frac{9\pi^2}{(1-\alpha)^2 c^2 \log d}$ ,

$$\mathbb{P}\left(|\hat{\mathbf{R}}_{jk} - \mathbb{E}\hat{\mathbf{R}}_{jk}| > \frac{2c}{\pi} \sqrt{\frac{\log d}{n}}\right) \leq 2 \exp\left(-\frac{2\alpha^2 c^2 \log d}{27\pi^2}\right),$$

and whenever  $t \geq \frac{6\pi}{t} + 2$  (A.1)

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq \mathbb{P}\left(|\hat{\mathbf{R}}_{jk} - \mathbb{E}\hat{\mathbf{R}}_{jk}| > \frac{2t}{\pi}\right). \quad (\text{A.2})$$

In Equation (A.1), letting  $t = c\sqrt{\frac{\log d}{n}}$  and  $\alpha = 3\sqrt{6}/8$  and applying (A.2), we have whenever  $n \geq \frac{37\pi}{t}$ ,

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq 2 \exp\left(-\frac{nt^2}{16\pi^2}\right).$$

This completes the proof.  $\square$

**Theorem A.2 (Bernstein Inequality).** Let  $x_1, \dots, x_n$  be  $n$  independent realizations of a random variable  $X$  with  $\mathbb{E}X = 0$ . Suppose that for some positive constant  $K$ , we have

$$\mathbb{E}|X^m| \leq \frac{m!}{2} K^{m-2}, \quad m = 2, 3, \dots$$

Then for all  $0 < t \leq \frac{1-2C}{2KC}$ ,

$$\mathbb{P}\left(\frac{1}{n} \sum_{i=1}^n x_i \geq t\right) \leq \exp(-nCt^2),$$

where  $C$  is a generic constant not scaled with  $(n, d)$ .

**Proof.** Using Lemma 5.7 of [36], we have for all  $a > 0$

$$\mathbb{P}\left(\sum x_i \geq \sqrt{na}\right) \leq \exp\left(-\frac{a^2}{2(aKn^{-1/2} + 1)}\right).$$

Letting  $t = a/\sqrt{n}$ , we have

$$\mathbb{P}\left(\frac{1}{n} \sum x_i \geq t\right) \leq \exp\left(-\frac{nt^2}{2(Kt + 1)}\right).$$

If  $t \leq \frac{1/C-2}{2K}$ , we further have  $2(Kt + 1) \leq 1/C$ , implying that

$$\mathbb{P}\left(\frac{1}{n} \sum x_i \geq t\right) \leq \exp(-nCt^2).$$

## APPENDIX B

### MAIN PROOFS

#### B.1 Proof of Lemma 3.1

**Proof.** Equation (3.6) holds if and only if

$$\frac{\|\mathbf{v}_{A_{k+1}}\|_2}{\|\mathbf{v}_{A_k}\|_2} \leq \frac{\|\mathbf{v}_{A_{k+1}}\|_q}{\|\mathbf{v}_{A_k}\|_q}.$$

This is equivalent to proving that

$$\left(1 + \frac{v_{k+1}^2}{\sum_{j=1}^k v_j^2}\right)^{1/2} \leq \left(1 + \frac{v_{k+1}^q}{\sum_{j=1}^k v_j^q}\right)^{1/q}. \quad (\text{B.1})$$

If  $v_{k+1} = 0$ , it is easy to see that Equation (B.1) holds. If not, denoting by  $m_j = \frac{v_j}{v_{k+1}} \geq 1$ , to prove that Equation (B.1) holds is equivalent to proving that

$$\left(1 + \frac{1}{\sum_{j=1}^k m_j^2}\right)^q \leq \left(1 + \frac{1}{\sum_{j=1}^k m_j^q}\right)^2.$$

Realizing that for any  $x \in \mathbb{R}^+ \cap \{0\}$  and  $0 < \alpha \leq 1$

$$x^\alpha \leq 1 + \alpha(x - 1),$$

we have

$$\begin{aligned} \left(1 + \frac{1}{\sum_{j=1}^k m_j^2}\right)^q &\leq 1 + q \cdot \frac{1}{\sum_{j=1}^k m_j^2} \leq 1 + 2 \cdot \frac{1}{\sum_{j=1}^k m_j^q} \\ &\leq \left(1 + \frac{1}{\sum_{j=1}^k m_j^q}\right)^2. \end{aligned}$$

This completes the proof.  $\square$

## B.2 Proof of Lemma 3.2

**Proof.** By using Lemma A.1, we have

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq 2 \exp\left(-\frac{nt^2}{16\pi^2}\right).$$

Because  $\Sigma^0$  is feasible to Equation (3.8),  $\tilde{\mathbf{R}}$  must satisfy that:

$$\|\hat{\mathbf{R}} - \tilde{\mathbf{R}}\|_{\max} \leq \|\hat{\mathbf{R}} - \Sigma^0\|_{\max}.$$

Using the triangular inequality, we then have

$$\begin{aligned} \mathbb{P}(|\tilde{\mathbf{R}}_{jk} - \Sigma_{jk}^0| \geq t) &\leq \mathbb{P}(|\tilde{\mathbf{R}}_{jk} - \hat{\mathbf{R}}_{jk}| + |\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| \geq t) \\ &\leq \mathbb{P}(\|\tilde{\mathbf{R}} - \hat{\mathbf{R}}\|_{\max} + \|\hat{\mathbf{R}} - \Sigma^0\|_{\max} \geq t) \\ &\leq \mathbb{P}(\|\hat{\mathbf{R}} - \Sigma^0\|_{\max} \geq t/2) \leq d^2 \exp\left(-\frac{nt^2}{64\pi^2}\right) \\ &\leq 2 \exp\left(\frac{2 \log d}{\log 2} - \frac{nt^2}{64\pi^2}\right). \end{aligned}$$

Using the fact that  $t \geq 16\pi\sqrt{\frac{\log d}{n \log 2}}$ , we have the result.  $\square$

## B.3 Proof of Lemma 4.3

**Proof.** Because  $g_j^2 \in TF(K)$ , where  $K$  is a constant not scaled with  $(n, d)$ , we have that  $X_j$ 's moments are controlled by  $K$  for  $j = 1, \dots, d$ . Therefore,  $\mu_j$  and  $\sigma_j$  are not scaled with  $(n, d)$ . Accordingly, we can assume that  $\mu = 0$  and  $\text{diag}(\Sigma) = 1$  without loss of generality. Let  $X = (X_1, \dots, X_d)^T \sim MNPN_d(\mu, \Sigma, f)$ . To prove that Equation (4.4) and Equation (4.5) hold, the key is to prove that the high order moments of each  $X_j$  and  $X_j^2$  will not grow very fast.

Generally, define  $Z := f_j(X_j) \sim N(0, 1)$ . We have for any  $m \in \mathbb{Z}^+$ , because  $g_j^2 \in TF(K)$  for some constant  $K$ , by definition

$$\mathbb{E}|X_j^2|^m = \mathbb{E}|g_j(Z)|^m \leq \frac{m!}{2} K^m.$$

Moreover, we have  $\mathbb{E}(X_j)^m$  can be bounded by a similar term, in detail,

$$\begin{aligned} \mathbb{E}|X_j|^m &= \mathbb{E}|X_j^{m/2}|^{m/2} \leq \frac{(m/2)!}{2} K^{m/2} < \frac{m!}{2} K^m, \text{ if } m \text{ is even,} \\ \mathbb{E}|X_j|^m &\leq 1 + \mathbb{E}|X_j|^m I(|X_j| \geq 1) \leq 1 + \mathbb{E}(|X_j|^{m+1} I(|X_j| \geq 1)) \\ &\leq 1 + \mathbb{E}|X_j|^{m+1} \leq 1 + \frac{(\frac{m+1}{2})!}{2} K^{\frac{m+1}{2}} < \frac{m!}{2} (2K + 2)^m, \\ &\quad \text{if } m \text{ is odd.} \end{aligned}$$

Therefore, realizing that  $\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}$ ,  $\hat{\sigma}_j^2 = \frac{1}{n} \sum_{i=1}^n x_{ij}^2$ , and  $\mathbb{E}X_j^m \leq \frac{m!}{2} (2K + 2)^m$  and  $\mathbb{E}(X_j^2)^m \leq \frac{m!}{2} K^m$ , we can

apply the Bernstein inequality (shown in Theorem A.2) to obtain a concentration inequality for  $\hat{\mu}_j$  and  $\hat{\sigma}_j^2$ . In particular, we have

$$\mathbb{P}(|\hat{\mu}_j - \mu_j| > t) \leq 2 \exp(-c_2 nt^2), \quad (\text{B.2})$$

$$\mathbb{P}(|\hat{\sigma}_j^2 - \sigma_j^2| > t) \leq 2 \exp(-c_3 nt^2), \quad (\text{B.3})$$

where  $c_2$  and  $c_3$  only depend on  $K$ . Using Equation (B.3), we further have

$$\mathbb{P}(|\hat{\sigma}_j - \sigma_j| > t) \leq \mathbb{P}(|\hat{\sigma}_j^2 - \sigma_j^2| > c_0 t) \leq 2 \exp(-c_4 nt^2),$$

where  $c_0$  is a generic constant and  $c_4 = c_3 \cdot c_0^2$  only depends on the choice the  $K$ . To finalize the proof, we need to show that combining  $\hat{R}$  with  $\{\hat{\sigma}_1, \dots, \hat{\sigma}_d\}$  will not hurt the rate. To show that, suppose that

$$\mathbb{P}(|\hat{\sigma}_j - \sigma_j| > t) \leq \eta_1(n, t), \quad (\text{B.4})$$

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq \eta_2(n, t). \quad (\text{B.5})$$

Letting  $\sigma_{\max}^2 := \max_j(\sigma_j^2)$  be controlled by  $K/2$ , we have

$$\begin{aligned} &\mathbb{P}(|\hat{\mathbf{S}}_{jk} - \Sigma_{jk}| > \epsilon) \\ &\leq \mathbb{P}\left(|(\hat{\sigma}_j \hat{\sigma}_k - \sigma_j \sigma_k) \hat{\mathbf{R}}_{jk}| > \frac{\epsilon}{2}\right) + \left(|\sigma_j \sigma_k (\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0)| > \frac{\epsilon}{2}\right) \\ &\leq \mathbb{P}\left(|\hat{\sigma}_j \hat{\sigma}_k - \sigma_j \sigma_k| > \frac{\epsilon}{2}\right) + \left(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > \frac{\epsilon}{2\sigma_{\max}^2}\right) \\ &\leq \mathbb{P}\left(|\hat{\sigma}_j - \sigma_j| > \sqrt{\frac{\epsilon}{6}}\right) + \mathbb{P}\left(|\hat{\sigma}_k - \sigma_k| > \sqrt{\frac{\epsilon}{6}}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}_j - \sigma_j| > \frac{\epsilon/6}{\sigma_{\max}}\right) \\ &\quad + \mathbb{P}\left(|\hat{\sigma}_k - \sigma_k| > \frac{\epsilon}{6\sigma_{\max}}\right) + \eta_2\left(n, \frac{\epsilon}{2\sigma_{\max}^2}\right) \\ &\leq 2\eta_1\left(n, \sqrt{\frac{\epsilon}{6}}\right) + 2\eta_1\left(n, \frac{\epsilon}{6\sigma_{\max}}\right) + \eta_2\left(n, \frac{\epsilon}{2\sigma_{\max}^2}\right). \end{aligned}$$

By using Lemma A.1, we have

$$\mathbb{P}(|\hat{\mathbf{R}}_{jk} - \Sigma_{jk}^0| > t) \leq 2 \exp\left(-\frac{nt^2}{16\pi^2}\right).$$

It means that  $\eta_1$  and  $\eta_2$  are both of parametric exponential decay rate. This completes the proof.  $\square$

## B.4 Proof of Theorem 4.5

**Proof.** For  $\mathcal{M}^0(q, R_q, \Sigma^0, f^0)$  with  $0 \leq q \leq 1$ , we define

$$\epsilon := \sin \angle(\theta_1, \tilde{\theta}_1) \quad \text{and} \quad \Sigma^0 = \lambda_1 \theta_1 \theta_1^T + \Psi_0, \quad (\text{B.6})$$

where  $\Psi_0 = \sum_{j=2}^d \lambda_j \theta_j \theta_j^T$  is perpendicular to  $\theta_1$ . For all  $\theta \in \mathbb{S}^{d-1}$ , we have

$$\begin{aligned} \langle \Sigma^0, \theta_1 \theta_1^T - \theta \theta^T \rangle &= \langle \Sigma^0, \theta_1 \theta_1^T \rangle - \langle \lambda_1 \theta_1 \theta_1^T + \Psi_0, \theta \theta^T \rangle \\ &= \lambda_1 - \lambda_1 \langle \theta_1, \theta \rangle^2 - \langle \Psi_0, \theta \theta^T \rangle, \end{aligned} \quad (\text{B.7})$$

$$\begin{aligned} \text{and } \langle \Psi_0, \theta \theta^T \rangle &= \theta^T \Psi_0 \theta = \theta^T (\mathbf{I}_d - \theta_1 \theta_1^T) \Sigma^0 (\mathbf{I}_d - \theta_1 \theta_1^T) \theta \\ &\leq \lambda_2 \|(\mathbf{I}_d - \theta_1 \theta_1^T) \theta\|_2^2 = \lambda_2 - \lambda_2 \langle \theta_1, \theta \rangle^2. \end{aligned} \quad (\text{B.8})$$

Moreover, by definition,

$$\sin^2 \angle(\theta_1, \theta) = 1 - (\theta_1^T \theta)^2 = 1 - \langle \theta_1, \theta \rangle^2. \quad (\text{B.9})$$

Combining Equation (B.7) with Equation (B.9), we have

$$\langle \Sigma^0, \theta_1 \theta_1^T - \theta \theta^T \rangle \geq (\lambda_1 - \lambda_2) \sin^2 \angle(\theta_1, \theta).$$

Therefore, letting  $\tilde{\theta}_1$  be the minimizer to Equation (3.2), we have

$$\begin{aligned} \epsilon^2 &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \Sigma^0, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle \\ &\leq \frac{1}{\lambda_1 - \lambda_2} (\langle \Sigma^0 - \hat{\mathbf{R}}, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle + \langle \hat{\mathbf{R}}, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle) \\ &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \Sigma^0 - \hat{\mathbf{R}}, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle. \end{aligned} \quad (\text{B.10})$$

The last inequality holds because

$$\langle \hat{\mathbf{R}}, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle = \theta_1^T \hat{\mathbf{R}} \theta_1 - \tilde{\theta}_1^T \hat{\mathbf{R}} \tilde{\theta}_1 \leq 0.$$

Therefore, using Equation (B.10),

$$\begin{aligned} \epsilon^2 &\leq \frac{1}{\lambda_1 - \lambda_2} \langle \Sigma^0 - \hat{\mathbf{R}}, \theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T \rangle \\ &= \frac{1}{\lambda_1 - \lambda_2} \langle \text{vec}(\Sigma^0 - \hat{\mathbf{R}}), \text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T) \rangle \\ &\leq \frac{1}{\lambda_1 - \lambda_2} \|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty \cdot \|\text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T)\|_1, \end{aligned} \quad (\text{B.11})$$

where the last inequality is by using Hölder Inequality.

When  $q = 1$ , we have

$$\begin{aligned} \|\text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T)\|_1 &\leq \|(\theta_1 \theta_1^T)\|_1 + \|(\tilde{\theta}_1 \tilde{\theta}_1^T)\|_1 \\ &= \sum_j \sum_k |\theta_{1j} \theta_{1k}| + \sum_j \sum_k |\tilde{\theta}_{1j} \tilde{\theta}_{1k}| \\ &= \|\theta_1\|_1^2 + \|\tilde{\theta}_1\|_1^2 \leq 2R_1^2. \end{aligned} \quad (\text{B.12})$$

The last inequality holds because both  $\tilde{\theta}_1$  and  $\theta_1$  belong to  $\mathbb{B}_1(R_1)$ . Therefore, we have

$$\epsilon^2 \leq 2R_1^2 \cdot \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2}. \quad (\text{B.13})$$

When  $0 < q < 1$ , denoting by  $\vartheta = \text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T)$ , we have

$$\begin{aligned} \|\vartheta\|_q^q &\leq \|\text{vec}(\theta_1 \theta_1^T)\|_q^q + \|\text{vec}(\tilde{\theta}_1 \tilde{\theta}_1^T)\|_q^q \\ &= \sum_j \sum_k (\theta_{1j} \theta_{1k})^q + \sum_j \sum_k (\tilde{\theta}_{1j} \tilde{\theta}_{1k})^q \\ &= \|\theta_1\|_q^q + \|\tilde{\theta}_1\|_q^q \leq 2R_q^2. \end{aligned} \quad (\text{B.14})$$

Therefore, denoting by  $S_\vartheta := \{j, |\vartheta_j| > \tau\}$  for some  $\tau$ , we have

$$\begin{aligned} \|\vartheta\|_1 &= \|\vartheta_{S_\vartheta}\|_1 + \sum_{j \in S_\vartheta} |\vartheta_j| \leq \sqrt{\text{card}(S_\vartheta)} \|\vartheta\|_2 + \tau \sum_{j \in S_\vartheta} \frac{|\vartheta_j|}{\tau} \\ &\leq \sqrt{\text{card}(S_\vartheta)} \|\vartheta\|_2 + \tau \sum_{j \in S_\vartheta} \left( \frac{|\vartheta_j|}{\tau} \right)^q \\ &\leq \sqrt{\text{card}(S_\vartheta)} \|\vartheta\|_2 + 2R_q^2 \tau^{1-q}. \end{aligned} \quad (\text{B.15})$$

The last inequality holds because

$$\begin{aligned} \text{card}(S_\vartheta) \cdot \tau^q &\leq \sum_{j \in S_\vartheta} |\vartheta_j|^q \leq \|\vartheta\|_q^q \leq 2R_q^2 \\ \text{and } \sum_{j \in S_\vartheta} |\vartheta_j|^q &\leq \|\vartheta\|_q^q \leq 2R_q^2. \end{aligned}$$

Therefore, letting  $\tau = \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2}$  and realizing that

$$\|\vartheta\|_2^2 = \|\text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T)\|_2^2 = 2(1 - (\theta_1^T \tilde{\theta}_1)^2) = 2\epsilon^2,$$

combining Equation (B.11) with (B.15), we have

$$\epsilon^2 \leq 2\tau^{1-q/2} R_q \epsilon + 2\tau^{2-q} R_q^2.$$

Therefore,  $\epsilon \leq (1 + \sqrt{3})\tau^{1-q/2} R_q$ , in other words,

$$\epsilon^2 \leq (1 + \sqrt{3})^2 R_q^2 \left( \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2} \right)^{2-q}. \quad (\text{B.16})$$

When  $q = 0$ , denoting by  $\vartheta = \text{vec}(\theta_1 \theta_1^T - \tilde{\theta}_1 \tilde{\theta}_1^T)$ ,

$$\|\vartheta\|_1 \leq \sqrt{\text{card}(\text{supp}(\vartheta))} \|\vartheta\|_2 \leq \sqrt{2R_0^2} \cdot \sqrt{2\epsilon^2} = 2R_0\epsilon. \quad (\text{B.17})$$

Therefore, combining Equation (B.11) with Equation (B.17), we have

$$\epsilon^2 \leq \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2} \cdot 2R_0\epsilon,$$

which is equivalent to stating that

$$\epsilon^2 \leq 4R_0^2 \left( \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2} \right)^2. \quad (\text{B.18})$$

Combining Equation (B.13), (B.16) and (B.18), we have that for all  $0 \leq q \leq 1$ :

$$\epsilon^2 \leq \gamma_q R_q^2 \left( \frac{\|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty}{\lambda_1 - \lambda_2} \right)^{2-q},$$

where  $\gamma_q = 2 \cdot I(q=1) + 4 \cdot I(q=0) + (1 + \sqrt{3})^2 \cdot I(0 < q < 1)$ . Then, using Lemma 4.3, we have

$$\begin{aligned} \mathbb{P}(\epsilon^2 \geq t) &\leq \mathbb{P} \left( \frac{\gamma_q R_q^2}{(\lambda_1 - \lambda_2)^{2-q}} \|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty^{2-q} \geq t \right) \\ &= \mathbb{P} \left( \|\hat{\mathbf{R}} - \Sigma^0\|_{\max} \geq \left( \frac{t(\lambda_1 - \lambda_2)^{2-q}}{\gamma_q R_q^2} \right)^{1/(2-q)} \right) \\ &\leq d^2 \exp \left( -\frac{n}{16\pi^2} \cdot \left( \frac{t(\lambda_1 - \lambda_2)^{2-q}}{\gamma_q R_q^2} \right)^{2/(2-q)} \right), \end{aligned} \quad (\text{B.19})$$

where in the last inequality the constant  $\frac{1}{16\pi^2}$  is derived by using Lemma A.1. Finally, choosing  $t = \gamma_q R_q^2 \left( \frac{64\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n} \right)^{\frac{2-q}{2}}$ , we have the result.  $\square$

## Proof of Corollary 4.7

**Proof.** Without loss of generality and for simplicity, we may assume that  $\tilde{\theta}_1^T \theta_1 \geq 0$ , because otherwise we can simply do appropriate sign changes in the proof. We first note that

$$\text{card}(\hat{\Theta}^0) = \text{card}(\Theta^0) = R_0. \quad (\text{B.20})$$

If  $\hat{\Theta}^0 \neq \Theta^0$ , then let  $\hat{\Theta}_d := (\hat{\Theta}^0 / \Theta^0) \cup (\Theta^0 / \hat{\Theta}^0)$ . We have

$$\|\tilde{\theta}_1 - \theta_1\|_2^2 \geq \|(\tilde{\theta}_1 - \theta_1)_{\hat{\Theta}_d}\|_2^2 = \sum_{j \in \hat{\Theta}_d} (\tilde{\theta}_{1j}^2 + \theta_{1j}^2).$$

The last equality holds because for any  $j \in \hat{\Theta}_d$ , either  $\tilde{\theta}_{1j}$  or  $\theta_{1j}$  are non-zero. Because of Equation (B.20), if  $\hat{\Theta}_d \neq \emptyset$ , there must exist  $j \in \hat{\Theta}_d$  such that  $\theta_{1j} \neq 0$ . Therefore,

$$\|\tilde{\theta}_1 - \theta_1\|_2 \geq \min_{j \in \Theta} |\theta_{1j}| \geq \frac{16\sqrt{2}R_0\pi}{\lambda_1 - \lambda_2} \sqrt{\frac{\log d}{n}}.$$

Then we have

$$\sin^2 \angle(\tilde{\theta}_1, \theta_1) = 1 - (\tilde{\theta}_1^T \theta_1)^2 \geq 1 - \tilde{\theta}_1^T \theta_1 = \frac{\|\theta_1 - \tilde{\theta}_1\|_2^2}{2},$$

implying that

$$\sin^2 \angle(\tilde{\theta}_1, \theta_1) \geq \frac{\|\tilde{\theta}_1 - \theta_1\|_2^2}{2} \geq \frac{256R_0^2\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n}.$$

Therefore, applying Theorem 4.5, we have

$$\mathbb{P}(\hat{\Theta}^0 \neq \Theta^0) \leq \mathbb{P}\left(\sin^2 \angle(\tilde{\theta}_1, \theta_1) \geq \frac{256R_0^2\pi^2}{(\lambda_1 - \lambda_2)^2} \cdot \frac{\log d}{n}\right) \leq \frac{1}{d^2}.$$

This completes the proof.  $\square$

## B.6 Proof of Theorem 4.12

**Proof.** We first prove that  $\max_{v \in \mathbb{S}^{d-1} \cap \mathbb{B}_0(R_0+2k)} |v^T(\hat{\mathbf{R}} - \Sigma^0)|$  can be bounded by  $8\pi(R_0 + 2k)\sqrt{\frac{\log d}{n}}$  with large probability. To show that, we have

$$\begin{aligned} |v^T(\hat{\mathbf{R}} - \Sigma^0)v| &= |\langle \hat{\mathbf{R}} - \Sigma^0, vv^T \rangle| \\ &= |\langle \text{vec}(\hat{\mathbf{R}} - \Sigma^0), \text{vec}(vv^T) \rangle| \\ &\leq \|\text{vec}(\hat{\mathbf{R}} - \Sigma^0)\|_\infty \cdot \|\text{vec}(vv^T)\|_1 \\ &\leq \|\hat{\mathbf{R}} - \Sigma^0\|_{\max} \cdot (R_0 + 2k). \end{aligned}$$

Using Lemma 2.5, we have the result. Then replace  $\rho(E, s)$  with  $8\pi(R_0 + 2k)\sqrt{\frac{\log d}{n}}$  in [5, Theorem 1]. Finally, realizing that for any  $v_1, v_2 \in \mathbb{S}^{d-1}$

$$\begin{aligned} \sqrt{1 - |v_1^T v_2|} &\leq \sqrt{1 - |v_1 v_2|^2} = \sin \angle(v_1, v_2) \\ &= \sqrt{1 + |v_1^T v_2|} \cdot \sqrt{1 - |v_1^T v_2|} \leq 2\sqrt{1 - |v_1^T v_2|}, \end{aligned}$$

we have the result.  $\square$

## ACKNOWLEDGMENTS

The authors thank the associate editor and two anonymous reviewers, who made numerous helpful suggestions for

improvements to the paper. The authors are supported by NSF Grants III-1116730 and NSF III-1332109, NIH R01MH102339, NIH R01GM083084, and NIH R01HG06841, and FDA HHSF223201000072C. Fang is also supported by a fellowship from Google.

## REFERENCES

- [1] H. Liu, J. Lafferty, and L. Wasserman, "The nonparanormal: Semiparametric estimation of high dimensional undirected graphs," *J. Mach. Learn. Res.*, vol. 10, pp. 2295–2328, 2009.
- [2] 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, 2012.
- [3] T. Anderson, *An Introduction to Multivariate Statistical Analysis*. New York, NY, USA: Wiley, 1958, vol. 2.
- [4] I. Johnstone and A. Lu, "On consistency and sparsity for principal components analysis in high dimensions," *J. Amer. Statist. Assoc.*, vol. 104, no. 486, pp. 682–693, 2009.
- [5] X.-T. Yuan and T. Zhang, "Truncated power method for sparse eigenvalue problems," *J. Mach. Learn. Res.*, vol. 14, pp. 899–925, 2013.
- [6] Z. Ma, "Sparse principal component analysis and iterative thresholding," *Ann. Statist.*, vol. 41, no. 2, pp. 772–801, 2013.
- [7] V. Vu and J. Lei, "Minimax rates of estimation for sparse PCA in high dimensions," in *Proc. 15th Int. Conf. Artif. Intell. Stat.*, 2012, vol. 22, pp. 1278–1286.
- [8] M. Borgognone, J. Bussi, and G. Hough, "Principal component analysis in sensory analysis: Covariance or correlation matrix?" *Food Quality Preference*, vol. 12, no. 5–7, pp. 323–326, 2001.
- [9] F. Han and H. Liu, "Semiparametric principal component analysis," in *Proc. Advances Neural Inf. Process. Syst.*, 2012, vol. 24, pp. 171–179.
- [10] B. Flury, *A First Course Multivariate Statistics*. New York, NY, USA: Springer, 1997.
- [11] S. Konishi, "Asymptotic expansions for the distributions of statistics based on the sample correlation matrix in principal component analysis," *Hiroshima Math. J.*, vol. 9, no. 3, pp. 647–700, 1979.
- [12] H. Nagao, "On the jackknife statistics for eigenvalues and eigenvectors of a correlation matrix," *Ann. Inst. Stat. Math.*, vol. 40, no. 3, pp. 477–489, 1988.
- [13] C. Chatfield and A. Collins, *Introduction to Multivariate Analysis*. Boca Raton, FL, USA: CRC, 1980, vol. 166.
- [14] R. Johnson and D. Wichern, *Applied Multivariate Statistical Analysis*, 6th ed. Englewood Cliffs, NJ, USA: Prentice hall, 2007.
- [15] C. Klaassen and J. Wellner, "Efficient estimation in the bivariate normal copula model: Normal margins are least favourable," *Bernoulli*, vol. 3, no. 1, pp. 55–77, 1997.
- [16] W. Kruskal, "Ordinal measures of association," *J. Amer. Stat. Assoc.*, vol. 53, no. 284, pp. 814–861, 1958.
- [17] M. Balasubramanian and E. Schwartz, "The isomap algorithm and topological stability," *Science*, vol. 295, no. 5552, p. 7, 2002.
- [18] G. Raskutti, M. Wainwright, and B. Yu, "Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -Balls," *IEEE Trans. Inf. Theory*, vol. 57, no. 10, pp. 6976–6994, Oct. 2011.
- [19] D. Paul and I. Johnstone, "Augmented sparse principal component analysis for high dimensional data," Arxiv preprint arXiv:1202.1242, 2012.
- [20] D. Witten, R. Tibshirani, and T. Hastie, "A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis," *Biostatistics*, vol. 10, no. 3, pp. 515–534, 2009.
- [21] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *J. Comput. Graph. Statist.*, vol. 15, no. 2, pp. 265–286, 2006.
- [22] H. Zou and T. Hastie, "Regularization and variable selection via the elastic net," *J. Royal Statist. Soc.: Series B (Stat. Methodol.)*, vol. 67, no. 2, pp. 301–320, 2005.
- [23] L. Mackey, "Deflation methods for sparse PCA," in *Proc. Adv. Neural Inf. Process. Syst.*, 2009, vol. 21, pp. 1017–1024.
- [24] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre, "Generalized power method for sparse principal component analysis," *J. Mach. Learn. Res.*, vol. 11, pp. 517–553, 2010.
- [25] Y. Zhang, A. d'Aspremont, and L. El Ghaoui, "Sparse PCA: Convex relaxations, algorithms and applications," *Handbook on Semidefinite, Conic and Polynomial Optimisation* New York, NY, USA: Springer, 2012 pp. 915–940.

- [26] A. d'Aspremont, L. El Ghaoui, M. Jordan, and G. Lanckriet, "A direct formulation for sparse PCA using semidefinite programming," *SIAM Rev.*, vol. 49, pp. 434–448, 2004.
- [27] A. A. Amini and M. J. Wainwright, "High-dimensional analysis of semidefinite relaxations for sparse principal components," *Ann. Statist.*, vol. 37, no. 5B, pp. 2877–2921, 2009.
- [28] Q. Berthet and P. Rigollet, "Optimal detection of sparse principal components in high dimension," *Ann. Statist.*, vol. 41, no. 4, pp. 1780–1815, 2013.
- [29] F. Han and H. Liu, "Optimal rates of convergence for latent generalized correlation matrix estimation in transelliptical distribution," arXiv preprint arXiv:1305.6916, 2013.
- [30] M. Wegkamp and Y. Zhao, "Adaptive estimation of the copula correlation matrix for semiparametric elliptical copulas," arXiv preprint arXiv:1305.6526, 2013.
- [31] H. Shen and J. Huang, "Sparse principal component analysis via regularized low rank matrix approximation," *J. Multivariate Anal.*, vol. 99, no. 6, pp. 1015–1034, 2008.
- [32] M. McCall, B. Bolstad, and R. Irizarry, "Frozen robust multiarray analysis (FRMA)," *Biostatistics*, vol. 11, no. 2, pp. 242–253, 2010.
- [33] J. Leek and J. Storey, "Capturing heterogeneity in gene expression studies by surrogate variable analysis," *PLoS Genetics*, vol. 3, no. 9, p. e161, 2007.
- [34] A. Eloyan, J. Muschelli, M. B. Nebel, H. Liu, F. Han, T. Zhao, A. D. Barber, S. Joel, J. J. Pekar, S. H. Mostofsky, and B. Caffo, "Automated diagnoses of attention deficit hyperactive disorder using magnetic resonance imaging," *Frontiers Syst. Neurosci.*, vol. 6, p. 61, 2012.
- [35] J. D. Power, A. L. Cohen, S. M. Nelson, G. S. Wig, K. A. Barnes, J. A. Church, A. C. Vogel, T. O. Laumann, F. M. Miezin, B. L. Schlaggar, and S. E. Petersen, "Functional network organization of the human brain," *Neuron*, vol. 72, no. 4, pp. 665–678, 2011.
- [36] S. Van De Geer, *Empirical Processes in M-Estimation*. Cambridge, U.K.: Cambridge Univ. Press, 2000, vol. 105.



**Fang Han** received the BSc degree in statistics at the Peking University, Beijing, China, in 2008, the MSc degree in biostatistics from the University of Minnesota in 2010. He is currently working toward the PhD degree in the Department of Biostatistics, Johns Hopkins University. He studies high dimensional statistics and learning theories, including graphical model estimation, principal component analysis, and discriminant analysis.



**Han Liu** received Joint PhD degree in machine learning and statistics from the Carnegie Mellon University in 2011. He is currently an assistant professor of Operations Research and Financial Engineering at Princeton University. His research interests include statistics, machine learning, optimization, and probability.

▷ For more information on this or any other computing topic, please visit our Digital Library at [www.computer.org/publications/dlib](http://www.computer.org/publications/dlib).