

# Principal component analysis for big data

Jianqing Fan<sup>\*</sup>, Qiang Sun<sup>†</sup>, Wen-Xin Zhou<sup>‡</sup> and Ziwei Zhu<sup>§</sup>

## Abstract

Big data is transforming our world, revolutionizing operations and analytics everywhere, from financial engineering to biomedical sciences. The complexity of big data has often made dimension reduction techniques necessary before conducting statistical inference. Principal component analysis, commonly referred to as PCA, has become an essential tool for multivariate data analysis and unsupervised dimension reduction, the goal of which is to find a lower dimensional subspace that captures most of the variation in the dataset. This article provides an overview of methodological and theoretical developments of PCA over the last decade, with focus on its applications to big data analytics. We first review the mathematical formulation of PCA and its theoretical development from the view point of perturbation analysis. We then briefly discuss the relationship between PCA and factor analysis as well as its applications to large covariance estimation and multiple testing. PCA also finds important applications in many modern machine learning problems, and we focus on community detection, ranking, mixture model and manifold learning in this paper.

**Keywords:** Big data; Covariance matrix; Dimension reduction; Distributed learning; Factor analysis; High dimensionality; Principal Components; Robust PCA.

---

<sup>\*</sup>Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544, USA. E-mail: [jqfan@princeton.edu](mailto:jqfan@princeton.edu).

<sup>†</sup>Department of Statistical Sciences, University of Toronto, Toronto, ON M5S 3G3. E-mail: [qsun.ustc@gmail.com](mailto:qsun.ustc@gmail.com).

<sup>‡</sup>Department of Mathematics, University of California, San Diego, La Jolla, CA 92093. E-mail: [wez243@ucsd.edu](mailto:wez243@ucsd.edu).

<sup>§</sup>Department of Operations Research and Financial Engineering, Princeton University, Princeton, NJ 08544, USA. E-mail: [ziweiz@princeton.edu](mailto:ziweiz@princeton.edu).

# 1 Introduction

Principal component analysis (PCA), first introduced by Karl Pearson (Pearson, 1901), is one of the most commonly used techniques for dimension reduction in many disciplines, such as neurosciences, genomics and finance (Izenman, 2008). We refer the readers to Jolliffe (2014) for a recent review. It extracts latent principal factors that preserve most of the variation in the dataset. Let  $\mathbf{x}$  be a random vector taking values in  $\mathbb{R}^d$  with mean zero and covariance matrix  $\Sigma$ . With this formalism, PCA seeks projection direction vectors,  $\mathbf{v}_1, \dots, \mathbf{v}_K \in \mathbb{R}^d$ , such that

$$\mathbf{v}_1 \in \operatorname{argmax}_{\|\mathbf{v}\|_2=1} \mathbf{v}^\top \Sigma \mathbf{v}, \quad \mathbf{v}_2 \in \operatorname{argmax}_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1} \mathbf{v}^\top \Sigma \mathbf{v}, \quad \mathbf{v}_3 \in \operatorname{argmax}_{\|\mathbf{v}\|_2=1, \mathbf{v} \perp \mathbf{v}_1, \mathbf{v}_2} \mathbf{v}^\top \Sigma \mathbf{v}, \dots$$

In other words,  $\{\mathbf{v}_j\}_{j=1}^K$  are the top  $K$  eigenvectors of  $\Sigma$ . Given  $\mathbf{V}_K \equiv (\mathbf{v}_1, \dots, \mathbf{v}_K)$ , we can then project the original high dimensional data onto the low dimensional space spanned by columns of  $\mathbf{V}_K$  to achieve the goal of dimensionality reduction. Since  $\mathbf{V}_K$  captures the most variation in the dataset, these projected data points approximately preserve the geometric properties of the original data, which are amenable to downstream statistical analysis. In real applications, the true covariance matrix  $\Sigma$  is typically unknown; we need to substitute  $\Sigma$  with its empirical version  $\hat{\Sigma}$ .

The high complexity of big data, such as massiveness, contamination, nonlinearity and decentralization, has posed fundamental challenges to statistical inference. PCA, an 100-year-old idea, has been and is still shining as a powerful tool for modern data analytics. Modern developments for PCA focus on attempts that address various challenges created by big data. For example, the massiveness of features in big data has been shown to create many notorious problems, making many conventional inferential procedures ill-posed. Recent study (Fan, Liao and Mincheva, 2013) shows that PCA is closely connected to factor analysis. This motivates new methodological developments in multiple testing problems when tens of thousands of possibly dependent statistical tests are evaluated simultaneously (Fan, Han and Gu, 2012; Fan et al., 2017), which partially solve the high dimensional inference problem. Moreover, big data are often contaminated by outliers or heavy-tailed errors (Fan, Li and Wang, 2017; Sun, Zhou and Fan, 2017), motivating the use of robust covariance inputs in the PCA formulation (Fan, Wang and Zhu, 2016; Minsker, 2016). This results in a form of robust PCA. Moreover, machine learning algorithms, such as those in clustering, community detection, ranking, matrix completion and mixture models, often involve solving a highly nonconvex system. This makes developing practical and efficient computational algorithms a grand challenge. PCA, or spectral method more generally, can often be used to solve a reduction of the highly nonconvex system, without losing much statistical efficiency. See Anandkumar et al. (2014) and Abbe (2017). Manifold sometimes can be used to approximate the nonlinearity structure of a dataset. Surprisingly, PCA finds applications in this setting and achieves a form of nonlinear dimension reduction (Torgerson, 1952; Roweis and Saul, 2000).

## 2 Covariance matrix estimation and PCA

We now begin the journey of PCA for big data. Given a small number  $K$ , the main goal of PCA is to estimate the  $K$ -dimensional principal eigenspace of  $\Sigma \in \mathbb{R}^{d \times d}$  such that it captures most of the variation in the dataset. In statistical applications, the true covariance matrix  $\Sigma$  is typically unknown and in practice, PCA is conducted on some estimator  $\widehat{\Sigma} = \widehat{\Sigma}(\mathbf{X}_1, \dots, \mathbf{X}_n)$  of  $\Sigma$ , where  $\mathbf{X}_1, \dots, \mathbf{X}_n$  are observations from  $\mathbf{X}$ . A significant error in recovering the eigenvalues and eigenvectors of  $\Sigma$  using those of  $\widehat{\Sigma}$  would lead to a substantial loss of information contained in the data by PCA projection. As direct applications of matrix perturbation theory (see Section 2.1), bounds on the estimation error  $\|\widehat{\Sigma} - \Sigma\|$  under some matrix norm  $\|\cdot\|$  can be effectively applied to relate the eigenvalues/eigenvectors of  $\Sigma$  and of  $\widehat{\Sigma}$  under the spectral gap condition. Therefore, it is important to build a good covariance estimator, say  $\widehat{\Sigma}$ , with small statistical error in the sense that for any given  $\delta > 0$ , we are interested in minimizing the value  $r$  that satisfies  $\mathbb{P}(\|\widehat{\Sigma} - \Sigma\| > r) \leq \delta$  for some matrix norm  $\|\cdot\|$ .

### 2.1 Perturbation theory

Considering  $\Delta = \widehat{\Sigma} - \Sigma$  as a perturbation, it is crucial to understand how the eigenspace of  $\widehat{\Sigma} = \Sigma + \Delta$  perturbs around that of  $\Sigma$ . This problem has been widely studied in the literature (Davis and Kahan, 1970; Wedin, 1972; Stewart, 1991, 2006; Yu, Wang and Samworth, 2015). Among these results, the  $\sin \Theta$  theorems, established by Davis and Kahan (1970) and Wedin (1972), have become fundamental tools and are commonly used in numerical analysis, statistics and machine learning. While Davis and Kahan (1970) focused on eigenvectors of symmetric matrices, Wedin's  $\sin \Theta$  theorem studies the singular vectors for asymmetric matrices and provides a uniform perturbation bound for both the left and right singular spaces in terms of the singular value gap and perturbation level. Over the years, various extensions have been made in different settings. For example, Vu (2011), Shabalin and Nobel (2013), O'Rourke, Vu and Wang (2013) and Wang (2015) considered the rotations of singular vectors after random perturbations; Cai and Zhang (2016) established separate perturbation bounds for the left and right singular subspaces. See also Dopico (2000) and Stewart (2006). Recently, Fan, Wang and Zhong (2016) derived new perturbation bounds, measured in the  $\ell_\infty$ -norm, for singular vectors (or eigenvectors in the symmetric case), which can be applied to robust covariance estimation in high-dimensional factor models and robust estimation of the false discovery proportion (Fan et al., 2017) when the sampling distributions are heavy-tailed. See Sections 3.2 and 3.3 for details on these statistical applications.

## 2.2 Robust covariance inputs and robust PCA

Due to matrix perturbation theory, upper bounds on the spectral norm  $\|\widehat{\Sigma} - \Sigma\|_2$  or the elementwise  $\ell_\infty$ -norm (also known as the max norm)  $\|\widehat{\Sigma} - \Sigma\|_{\max}$  can be used to establish corresponding bounds on the  $\ell_2$  distance and  $\ell_\infty$  distance between the population eigenvectors and their empirical counterparts obtained from  $\widehat{\Sigma}$ , respectively. Given independent observations  $\mathbf{X}_1, \dots, \mathbf{X}_n$  from  $\mathbf{X}$  with  $\mathbb{E}\mathbf{X} = \mathbf{0}$ , the sample covariance matrix, namely  $\widehat{\Sigma}_{\text{sam}} := n^{-1} \sum_{i=1}^n \mathbf{X}_i \mathbf{X}_i^\top$ , is arguably the most natural choice to estimate  $\Sigma \in \mathbb{R}^{d \times d}$  when the dimension  $d$  is smaller than the sample size  $n$ . The finite sample bound on  $\|\widehat{\Sigma}_{\text{sam}} - \Sigma\|_2$  or  $\|\widehat{\Sigma}_{\text{sam}} - \Sigma\|_{\max}$  has been well studied in the literature (Vershynin, 2012; Tropp, 2012). If  $\mathbf{X}$  has a sub-Gaussian distribution in the sense that for all unit vectors  $\mathbf{v} \in \mathbb{S}^{d-1}$ ,  $\mathbb{E} \exp(\lambda \mathbf{v}^\top \mathbf{X}) \leq \exp(c\lambda^2 \mathbf{v}^\top \Sigma \mathbf{v})$  for some constant  $c$ , then Remark 5.40 in Vershynin (2012) implies that for every  $t \geq 0$ ,  $\|\widehat{\Sigma}_{\text{sam}} - \Sigma\|_2 \leq \max(\delta, \delta^2) \|\Sigma\|_2$  with probability at least  $1 - 2e^{-t}$ , where  $\delta = C_1 \sqrt{(d+t)/n}$ . Similarly, it can be shown that with probability greater than  $1 - 2e^{-t}$ ,  $\|\widehat{\Sigma}_{\text{sam}} - \Sigma\|_{\max} \leq \max(\eta, \eta^2)$ , where  $\eta = C_2 \sqrt{(\log d + t)/n}$ . Here,  $C_1, C_2 > 0$  are constants depending only on  $c$ .

However, when the distribution is heavy-tailed, one cannot expect sub-Gaussian behaviors of the sample covariance in either the spectral or max norm (Catoni, 2012). Therefore, to perform PCA for heavy-tailed data, the sample covariance is a risky choice to begin with. Indeed, alternative robust estimators have been constructed to achieve better finite sample performance. In Fan, Li and Wang (2017), the authors constructed an elementwise robustified version of the sample covariance matrix using the adaptive Huber loss minimization (Sun, Zhou and Fan, 2017), and derived a sub-Gaussian-type deviation inequality in the max norm under finite fourth moment condition instead of sub-Gaussianity. Based on a novel shrinkage principle, Fan, Wang and Zhu (2016) and Minsker (2016) independently constructed global robustified variants of the sample covariance with sub-Gaussian behavior under the spectral norm as long as the fourth moment of  $\mathbf{X}$  is finite. A different robust method using the idea of median-of-means was proposed and studied by Minsker (2015). More recently, Giulini (2017) studied robust PCA in a more general setting where the data sample is made of independent copies of some random variable ranging in a separable real Hilbert space. Together, these results provide a new perspective on robustness from a nonasymptotic standpoint, and also represent a useful complement to the previous results on robust PCA. For instance, Candés et al. (2011) focused on a different notion of robustness and showed that it is possible to recover the principal components of a data matrix when the observations are contained in a low-dimensional space but arbitrarily corrupted by additive noise. See also Chandrasekaran et al. (2011), Zhang and Lerman (2014) and the references therein.

### 3 PCA and factor analysis

Principal component analysis and factor analysis are two important problems in their respective fields and are seemingly unrelated at first sight. Lately, it is shown in [Fan, Liao and Mincheva \(2013\)](#) that the high-dimensional factor model is innately related to PCA, which makes it different from the classical factor model with fixed dimensionality ([Lawley and Maxwell, 1971](#)) and helps one understand why PCA can be used for the factor analysis when the top eigenvalues are spikes. In addition, this observation has triggered a series of interesting studies on matrix perturbation theory and robust covariance estimation.

#### 3.1 Factor model and PCA

Let  $\mathbf{X} = (X_1, \dots, X_d)^\top$  be a random vector of outcomes of interest, which may represent financial returns, housing prices or gene expressions. The impact of dependence among outcomes is currently among the most discussed topics in various statistical problems, such as variable selection, covariance and precision matrix estimation and multiple testing. For example, financial returns depend on the equity market risks, housing prices depend on the economic health, gene expressions can be stimulated by cytokines, among others. Because of the presence of common factors, it is unrealistic to assume that many outcomes are uncorrelated. It is thus natural to assume a factor model structure, which relies on the identification of a linear space of random vectors capturing the dependence among the data. To do so, we consider an approximate factor model, which has been frequently used in economic and financial studies ([Fama and French, 1992](#); [Chamberlain and Rothschild, 1983](#); [Bai and Ng, 2002](#)) as well as genomics ([Leek and Storey, 2008](#)):

$$X_{ij} = \mu_j + \mathbf{b}_j^\top \mathbf{f}_i + \varepsilon_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, d, \quad (3.1)$$

or in a matrix form  $\mathbf{X}_i = \boldsymbol{\mu} + \mathbf{B}\mathbf{f}_i + \boldsymbol{\varepsilon}_i$  with  $\mathbf{B} = (\mathbf{b}_1, \dots, \mathbf{b}_d)^\top$ .

Here  $X_{ij}$  is the response for the  $j$ th feature of the  $i$ th observation  $\mathbf{X}_i = (X_{i1}, \dots, X_{id})^\top$ ,  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_d)^\top$  is the intercept,  $\mathbf{b}_j$  is a vector of factor loadings,  $\mathbf{f}_i$  is a  $K$ -dimensional vector of common factors and  $\boldsymbol{\varepsilon}_i = (\varepsilon_{i1}, \dots, \varepsilon_{id})^\top$  is the error term, which is typically called the idiosyncratic component, uncorrelated with or independent of  $\mathbf{f}_i$ . We emphasize that, in model (3.1), only  $X_{ij}$ 's are observable. Intuitively, the unobserved common factors can only be inferred reliably when  $d \rightarrow \infty$ . Under model (3.1),  $\boldsymbol{\Sigma} = \text{cov}(\mathbf{X})$  is given by

$$\boldsymbol{\Sigma} = \mathbf{B}\text{cov}(\mathbf{f}_i)\mathbf{B}^\top + \boldsymbol{\Sigma}_\varepsilon, \quad \boldsymbol{\Sigma}_\varepsilon = (\sigma_{\varepsilon,jk})_{1 \leq j,k \leq d} = \text{cov}(\boldsymbol{\varepsilon}_i). \quad (3.2)$$

The literature on approximate factor models typically assumes that the first  $K$  eigenvalues of  $\mathbf{B}\text{cov}(\mathbf{f}_i)\mathbf{B}^\top$  diverge at rate  $O(d)$ , whereas all the eigenvalues of  $\boldsymbol{\Sigma}_\varepsilon$  are bounded as  $d \rightarrow \infty$ . This assumption holds naturally when the factors are pervasive

in the sense that a non-negligible fraction of factor loadings should be non-vanishing. The decomposition (3.2) is then asymptotically identifiable as  $d \rightarrow \infty$ .

Now we elucidate why PCA can be used for the factor analysis in the presence of spiked eigenvalues. Note that the linear space spanned by the first  $K$  principal components of  $\mathbf{B}\text{cov}(\mathbf{f}_i)\mathbf{B}^\top$  coincides with that spanned by the columns of  $\mathbf{B}$  when  $\text{cov}(\mathbf{f}_i)$  is non-degenerate. Therefore, we can assume without loss of generality that the columns of  $\mathbf{B}$  are orthogonal and  $\text{cov}(\mathbf{f}_i)$  is the identity matrix. See also Bai and Li (2012) and Bai and Ng (2013). Let  $\bar{\mathbf{b}}_1, \dots, \bar{\mathbf{b}}_K$  be the columns of  $\mathbf{B}$ , ordered such that  $\|\bar{\mathbf{b}}_1\|_2 \geq \dots \geq \|\bar{\mathbf{b}}_K\|_2$ . Then,  $\{\bar{\boldsymbol{\xi}}_j := \bar{\mathbf{b}}_j / \|\bar{\mathbf{b}}_j\|_2\}_{j=1}^K$  are the top  $K$  eigenvectors of  $\mathbf{B}\mathbf{B}^\top$  with eigenvalues  $\{\bar{\lambda}_j := \|\bar{\mathbf{b}}_j\|_2^2\}_{j=1}^K$ . Let  $\{\lambda_j\}_{j=1}^d$  be the eigenvalues of  $\boldsymbol{\Sigma}$  in a non-increasing order with the corresponding eigenvectors being  $\{\boldsymbol{\xi}_j\}_{j=1}^d$ . From the *pervasiveness* assumption that the eigenvalues of  $d^{-1}\mathbf{B}^\top\mathbf{B}$  are distinct and bounded away from zero and infinity, it follows that  $\{\bar{\lambda}_j\}_{j=1}^K$  grow at rate  $O(d)$ . In addition, by Weyl's theorem and the  $\sin \Theta$  theorem of Davis and Kahan (Davis and Kahan, 1970),

$$\max_{1 \leq j \leq K} |\lambda_j - \bar{\lambda}_j| \bigvee \max_{j > K} |\lambda_j| \leq \|\boldsymbol{\Sigma}_\varepsilon\|_2 \quad \text{and} \quad \max_{1 \leq j \leq K} \|\boldsymbol{\xi}_j - \bar{\boldsymbol{\xi}}_j\|_2 = O(d^{-1}\|\boldsymbol{\Sigma}_\varepsilon\|_2).$$

These results state that PCA and factor analysis are approximately the same if  $\|\boldsymbol{\Sigma}_\varepsilon\|_2 = o(d)$ . This is assured through a sparsity condition on  $\boldsymbol{\Sigma}_\varepsilon$ , which is usually measured through  $m_d = \max_{1 \leq j \leq d} \sum_{k=1}^d |\sigma_{\varepsilon,jk}|^q$  for some  $q \in [0, 1]$ . The intuition is that, after subtracting out the common factors, many pairs of the cross-sectional units become weakly correlated. This generalized notion of sparsity was used by Bickel and Levina (2008), under which it holds that  $\|\boldsymbol{\Sigma}_\varepsilon\|_2 \leq \max_{1 \leq j \leq d} \sum_{k=1}^d |\sigma_{\varepsilon,jk}|^q (\sigma_{\varepsilon,jj}\sigma_{\varepsilon,kk})^{(1-q)/2} = O(m_d)$  if the variances  $\sigma_{\varepsilon,jj}$ 's are uniformly bounded. Therefore, in the approximate sparse setting where  $m_d = o(d)$ , the pervasiveness assumption implies that the principal components  $\{\lambda_j\}_{j=1}^K$  and the remaining components  $\{\lambda_j\}_{j=K+1}^d$  are well separated and the first  $K$  principal components are approximately the same as the standardized columns of the factor loading matrix. In this setting, PCA serves as a valid approximation to factor analysis only if  $d \rightarrow \infty$ .

### 3.2 Application to large covariance estimation

Covariance structure plays a particularly important role in high-dimensional data analysis. Apart from PCA, a large collection of fundamental statistical methods, such as linear and quadratic discriminant analysis, clustering analysis and regression analysis, require the knowledge of the covariance structure or certain aspects thereof. Realizing the importance of estimating large covariance matrices and the challenges that are brought by the high dimensionality, since the seminal work of Bickel and Levina (2008) and Fan, Fan and Lv (2008), various regularization techniques have been proposed to estimate  $\boldsymbol{\Sigma}$  consistently. See Cai, Ren and Zhou (2016) and Fan, Liao and Liu (2016) for two comprehensive surveys on this topic. One commonly as-

sumed low-dimensional structure is that the covariance matrix is sparse, namely many entries are zero or nearly so. In many applications, however, the sparsity assumption directly on  $\Sigma$  is not appropriate. A natural extension is conditional sparsity (Fan, Fan and Lv, 2008). Given the common factors, the outcomes are weakly correlated. In other words, in decomposition (3.2) we assume that  $\Sigma_\epsilon$  is approximately sparse as in Bickel and Levina (2008). Motivated by the innate connection between PCA and factor analysis in high dimensions, Fan, Liao and Mincheva (2013) proposed a principal orthogonal complement thresholding method (POET) to estimate  $\Sigma$  in model (3.1): (i) Run the singular value decomposition on the sample covariance matrix of  $\mathbf{X}$ ; (ii) keep the covariance matrix that is formed by the first  $K$  principal components; and (iii) apply the adaptive thresholding procedure (Rothman, Levina and Zhu, 2009; Cai and Liu, 2011) to the remaining components of the sample covariance matrix.

As discussed in Section 2.2, in the presence of heavy-tailed data, the use of the sample covariance matrix is controversial and therefore is in question to begin with. Tailored to elliptical factor models, Fan, Liu and Wang (2017+) proposed to use the marginal Kendall’s tau to estimate  $\Sigma$  and its top  $K$  eigenvalues, and separately, use the spatial Kendall’s tau to construct estimators for the corresponding leading eigenvectors. In more general settings where no shape constraints or parametric assumptions are imposed (normality, symmetry, elliptically symmetry, etc.), robust alternatives, such as the  $U$ -type covariance estimator and adaptive Huber covariance estimator considered in Fan et al. (2017), are preferred to be taken as the initial estimators in the POET procedure.

### 3.3 Application to factor-adjusted multiple testing

An iconic example of model (3.1) is the factor pricing model in financial economics, where  $X_{ij}$  is the excess return of fund/asset  $j$  at time  $i$ ,  $\mathbf{f}_i$ ’s are the systematic risk factors related to some specific linear pricing model, such as the capital asset pricing model (Sharpe, 1964), and the Fama-French three-factor model (Fama and French, 1993). Although the key implication from the multi-factor pricing theory is that the intercept  $\mu_j$  should be zero, known as the “mean-variance efficiency” pricing, for any asset  $j$ , an important question is whether such a pricing theory can be validated by empirical data. In fact, a very small proportion of  $\mu_j$ ’s might be nonzero according to the Berk and Green equilibrium (Berk and Green, 2004; Fan, Liao and Yao, 2015). It is practically important to identify those positive  $\mu_j$ ’s, which amounts to conducting the following  $d$  hypothesis testing problems simultaneously:

$$H_{0j} : \mu_j = 0 \quad \text{versus} \quad H_{1j} : \mu_j \neq 0, \quad j = 1, \dots, d. \quad (3.3)$$

In the presence of common factors,  $X_1, \dots, X_p$  are highly correlated and therefore directly applying classical false discovery rate (FDR) controlling procedures (Benjamini and Hochberg, 1995; Storey, Taylor and Siegmund, 2004) can lead to inaccurate false



discovery control and spurious outcomes. To improve the efficiency and to incorporate the strong dependence information, various factor-adjusted procedures have been proposed. See, for example, [Barras, Scaillet and Wermer \(2010\)](#), [Fan, Han and Gu \(2012\)](#) and [Lan and Du \(2017+\)](#). These works assume that both the factor and idiosyncratic noise follow multivariate normal distributions. However, the Gaussian assumption on the sampling distribution is often unrealistic in many applications, especially in genomics and finance. Also, as noted in [Perry and Owen \(2010\)](#), only non-Gaussian latent variables are detectable. To address the two challenges that are brought by strong dependence and heavy-tailedness, [Fan et al. \(2017\)](#) proposed a factor-adjusted robust multiple testing (FARM-Test) procedure with control of the false discovery proportion. Starting with a robust initial estimate of the covariance matrix, the FARM-Test method uses its first several principal components to estimate the factor loading matrix; next, using these estimated loading vectors, runs adaptive Huber regression to estimate the realized factors; robust test statistics are then constructed by subtracting out the realized common factors from the robust mean estimators; finally, the data-driven critical value is computed so that the estimated false discovery proportion is controlled at the prespecified level. The first step is motivated by the approximate equivalence between PCA and factor analysis as discussed in [Section 3.2](#).

## 4 Applications to statistical machine learning

PCA finds applications in many statistical machine learning problems. In this section, we focus on its applications to four major problems in the machine learning literature: clustering and community detection, ranking, mixture model and manifold learning.

### 4.1 Clustering and Community Detection

Clustering and community detection is an important task in data analysis ([Lowrimore and Manton, 2016](#)). We focus on the stochastic block model (SBM), which is widely regarded as a canonical model for this purpose. We refer the readers to [Abbe \(2017\)](#) for a review of recent developments on SBM. We start with the definition of SBM. Let  $n$  be the number of vertices,  $k$  the number of communities,  $\mathbf{p} = (p_1, \dots, p_k)^\top$  the probability vector and  $\mathbf{W}$  the  $k \times k$  symmetric transition probability matrix with entries in  $[0, 1]$ . We say a pair  $(\mathbf{X}, G)$  is sampled from  $\text{SBM}(n, \mathbf{p}, \mathbf{W})$  if  $\mathbf{X} = (X_1, \dots, X_n)^\top$  is an  $n$ -dimensional random vector with independent and identically distributed components distributed under  $\mathbf{p}$ , and  $G$  is an  $n$ -vertex simple graph where vertices  $i$  and  $j$  are connected with probability  $W_{X_i, X_j}$ , independently of other pairs of vertices. The  $i$ th community set is defined by  $\Omega_i := \{v : X_v = i, v = 1, \dots, n\}$ . Roughly speaking, a SBM is a random graph with planted clusters: the cluster sizes follow a multinomial distribution with probability vector  $\mathbf{p}$  and the probability that



a member in the  $i$  cluster and a member in the  $j^{\text{th}}$  group get connected is  $W_{i,j}$ .

To fix idea, we consider the SBM with two communities, where the inner-cluster probability is  $a$  and the across-cluster probability is  $b$  with  $b < a$ . Further assume for simplicity that the two clusters have exactly size  $n/2$ , and index the first cluster with the first  $n/2$  vertices. Given an observed graph with the adjacency matrix  $\mathbf{A}$  (indicating whether or not two vertices are connected), the expected value  $\mathbb{E}\mathbf{A}$  of this graph has four blocks given by

$$\mathbb{E}\mathbf{A} = \begin{bmatrix} a \cdot \mathbf{I}_{n/2} & b \cdot \mathbf{I}_{n/2} \\ b \cdot \mathbf{I}_{n/2} & a \cdot \mathbf{I}_{n/2} \end{bmatrix}.$$

It can be verified that the above matrix has three eigenvalues  $a + b$ ,  $a - b$  and  $0$ , where  $0$  has multiplicity  $n - 2$ , and the eigenvectors associated with the two largest eigenvalues are

$$\begin{bmatrix} \mathbf{1}_{n/2} \\ \mathbf{1}_{n/2} \end{bmatrix} \quad \text{and} \quad \begin{bmatrix} \mathbf{1}_{n/2} \\ -\mathbf{1}_{n/2} \end{bmatrix},$$

respectively. In other words, if one were to work with the expected adjacency matrix, communities could be simply recovered by taking the eigenvector associated with the second largest eigenvalue. In practice, we only observe a single-shot of the SBM graph  $\mathbf{A}$ , which can be viewed as a perturbation of the expected adjacency matrix:

$$\mathbf{A} = \mathbb{E}\mathbf{A} + \mathbf{Z},$$

where  $\mathbf{Z} = \mathbf{A} - \mathbb{E}\mathbf{A}$  is the perturbation. Therefore, the spectral method to achieve community detection can be formulated as solving the following optimization problem

$$\max_{\|\mathbf{x}\|_2^2=n, \mathbf{x}^\top \mathbf{1}_n=0} \mathbf{x}^\top \mathbf{A} \mathbf{x}, \quad \text{where } \mathbf{1}_n = (1, \dots, 1)^\top \in \mathbb{R}^n,$$

which is equivalent to finding the second eigenvector, as the first eigenvector is  $\mathbf{1}_n/\sqrt{n}$ .

However, as pointed by [Abbe, Bandeira and Hall \(2016\)](#), it has been an open problem whether the above spectral method is optimal for achieving the exact recovery property until very recently. We say that an algorithm achieves the exact recovery property if the community membership are recovered exactly with probability going to 1. [Abbe et al. \(2017\)](#) bridged this gap by providing a sharp entrywise eigenvector analysis for random matrices with low expected rank. Other than applications in SBM, the result can also be applied in other statistical machine learning problems, such as factor analysis and matrix completion, which we will not discuss in detail.

## 4.2 Ranking

Suppose we have a large collection of  $n$  items and we are given partially revealed comparisons between pairs of items. For example, player A defeats player B; video A

is preferred to video B when both are recommended. The goal is to identify the  $K$  items that receive the highest ranks based on these pair comparisons. This problem, which is called the top- $K$  rank aggregation, has wide applications including web search (Dwork et al., 2001), recommendation systems (Baltrunas et al., 2010), sports competition (Massey, 1997), etc. One of the most widely used parametric models is the Bradley-Terry-Luce (BTL) model (Bradley and Terry, 1952). In the BTL model, we have

$$\mathbb{P}(\text{item } j \text{ is preferred over item } i) = \frac{\omega_j^*}{\omega_i^* + \omega_j^*},$$

where  $\omega_i^*$  is the preference score of item  $i$  or the ability of  $i$ -th person. The task then boils down to finding the  $K$  items with the highest preference scores.

To see how PCA can be applied to this problem, we introduce one spectral method called “rank centrality” proposed by Negahban, Oh and Shah (2017). Consider each item as a node in a graph and construct a random walk on this graph where at each time, the random walk is possible to go from vertex  $i$  to vertex  $j$  if items  $i$  and  $j$  were ever compared; and if so, the likelihood of going from  $i$  to  $j$  depends on how often  $i$  lost to  $j$ . That is, the random walk is more likely to move to a neighbor who has more “wins”. The frequency this walk visits a particular node in the long run, or equivalently the stationary distribution, is the score of the corresponding item. Specifically, define the edge set  $\mathcal{E} := \{(i, j) : \text{item } i \text{ and } j \text{ were compared}\}$  and consider the transition matrix  $\mathbf{P}^* \in \mathbb{R}^{n \times n}$  such that

$$P_{i,j}^* = \begin{cases} \frac{1}{d} \frac{\omega_j^*}{\omega_i^* + \omega_j^*} & \text{if } (i, j) \in \mathcal{E}, \\ 1 - \frac{1}{d} \sum_{k:(i,k) \in \mathcal{E}} \frac{\omega_k^*}{\omega_i^* + \omega_k^*} & \text{if } i = j, \\ 0 & \text{otherwise,} \end{cases}$$

where  $d$  is a sufficiently large constant that makes every row nonnegative (a proper probability distribution). We can verify that the normalized score vector

$$\boldsymbol{\pi}^* := \frac{1}{\sum_{i=1}^n \omega_i^*} [\omega_1^*, \dots, \omega_n^*]^\top$$

is the stationary distribution of the Markov chain induced by  $\mathbf{P}^*$ , since  $\mathbf{P}^*$  and  $\boldsymbol{\pi}^*$  are in detailed balance, namely,  $\pi_i^* P_{i,j}^* = \pi_j^* P_{j,i}^*$ . Hence, by definition,  $\boldsymbol{\pi}^*$  is the top left singular vector of  $\mathbf{P}^*$ :  $\boldsymbol{\pi}^* \mathbf{P}^* = \boldsymbol{\pi}^*$ . This motivates the algorithm of “rank centrality” that uses the top left singular vector of the empirical transition probability matrix  $\hat{\mathbf{P}}$  as an estimate of the preference score.

As for the statistical guarantee, Negahban, Oh and Shah (2017) showed that  $\Omega(n \log n)$  pairs are needed for consistency of estimating  $\boldsymbol{\omega}^*$  in terms of the Euclidean norm. This is also the sample size needed for the Erdős-Rényi comparison graph to

get connected, which is the minimum condition that makes the identification of top  $K$  items possible. However, this does not lead to accurate identification of the top  $K$  items. A recent work [Chen et al. \(2017\)](#) showed that the same sample complexity ensures exact top- $K$  identification, thus matching the minimax lower bound established by [Chen and Suh \(2015\)](#) before. This was accomplished via optimal control of the entrywise error of the score estimates.

### 4.3 Mixture Model

PCA can also be applied to learning latent variable models or mixture regression models, which are important models for investigating heterogeneous data. To illustrate the idea, we consider a mixture of  $k$  Gaussian distributions with spherical covariances. Let  $w_i \in (0, 1)$  be the probability of choosing component  $i \in \{1, \dots, k\}$ ,  $\{\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k\} \subseteq \mathbb{R}^d$  be the component mean vectors, and  $\sigma^2 \mathbf{I}$  be the common covariance matrix. An observation in this model is given by

$$\mathbf{x} = \boldsymbol{\mu}_h + \mathbf{z}, \quad (4.1)$$

where  $h$  is a discrete random variable such that  $\mathbb{P}(h = j) = w_j$  for  $j = 1, \dots, k$ , and  $\mathbf{z} \sim \mathcal{N}(0, \sigma^2 \mathbf{I})$ . In other words,  $\mathbf{x} \sim w_1 \mathcal{N}(\boldsymbol{\mu}_1, \sigma^2 \mathbf{I}) + \dots + w_k \mathcal{N}(\boldsymbol{\mu}_k, \sigma^2 \mathbf{I})$  follows the mixture of the Gaussian distribution. We remark here that this is different from the topic model since every realization of  $\mathbf{x}$  corresponds to a different realization of  $h$ . The parameters of interest are the component mean vectors,  $\boldsymbol{\mu}_j$ 's.

Let  $\mathbf{M} = \mathbb{E}(\mathbf{x}\mathbf{x}^\top) - \sigma^2 \mathbf{I}$ . Then, it is easy to see ([Hsu and Kakade, 2013](#))

$$\mathbf{M} = \sum_{j=1}^k w_j \boldsymbol{\mu}_j \boldsymbol{\mu}_j^\top.$$

This indicates that running PCA on the population level would recover the  $k$ -dimensional linear subspace spanned by the mean component vectors. It can not fully identify all  $\boldsymbol{\mu}_j$ 's due to identifiability issues in general, but does help reduce the dimension of the parameter space, enabling the possibility of random initializations for a more delicate method in a second stage, such as higher-order tensor decomposition or likelihood based approaches. See [Anandkumar et al. \(2014\)](#) for a comprehensive review on tensor methods for latent variable models which includes the above discussed setting as a special case.

It is possible to extend the tensor decomposition method to study the finite mixture linear regression problems. One can replace the component mean vector in model (4.1) by  $\langle \mathbf{x}, \boldsymbol{\beta}_h \rangle$ , where  $h$  is again a random variable indicating the label of submodels. We refer to [Yi, Caramanis and Sanghavi \(2016\)](#) and the references therein for more discussions on such generalizations. We point out that there are still many open problems in this direction. For example, the global landscape of mixture models are not well understood in general.

## 4.4 Manifold Learning

Some of the most popular methods in manifold learning are also based on spectral methods. For example, the classical multidimensional scaling (MDS) is used as a numerical tool for manifold learning, which is frequently used in psychometrics as a means of visualizing the level of similarity (or dissimilarity) of individual cases in a dataset (Torgerson, 1952). It is also known as principal coordinate analysis (PCoA), emphasizing the fact that the classical MDS takes the dissimilarity matrix as an input and outputs a coordinate matrix by assigning each object a coordinate. The classical MDS uses the fact that the coordinate matrix  $\mathbf{U}$  can be derived by eigenvalue decomposition from  $\mathbf{B} = \mathbf{U}\mathbf{U}^\top$ , and the matrix  $\mathbf{B}$  can be computed from proximity matrix  $\mathbf{D}$  by using double centering:

- 1 Set up the matrix of squared dissimilarities  $\mathbf{D}^2 = [d_{ij}^2]$ .
- 2 Apply the double centering  $\mathbf{B} = -\frac{1}{2}\mathbf{J}\mathbf{D}^2\mathbf{J}$  using the centering matrix  $\mathbf{J} = \mathbf{I} - n^{-1}\mathbf{1}\mathbf{1}^\top$ , where  $n$  is the number of objects.
- 3 Extract the  $m$  largest positive eigenvalues  $\lambda_1, \dots, \lambda_m$  of  $\mathbf{B}$  and the corresponding eigenvectors  $\mathbf{v}_1, \dots, \mathbf{v}_m$ .
- 4 Let  $\mathbf{U} = \mathbf{V}_m\mathbf{\Lambda}_m^{1/2}$ , where  $\mathbf{V}_m = (\mathbf{v}_1, \dots, \mathbf{v}_m)$  and  $\mathbf{\Lambda}_m = \text{diag}(\lambda_1, \dots, \lambda_m)$ .

In the above algorithm, the dimension  $m$  can be chosen using criteria based on the eigenvalue ratios such as those in the factor analysis (Ahn and Horenstein, 2013; Lam and Yao, 2012). The classical MDS is a linear dimension reduction method which uses the euclidean distances between objects as the dissimilarity measures. It has many extensions by designing different input matrices based on dataset and objectives. These extensions include the metric MDS, the nonmetric MDS and the generalized MDS, largely extending the scope of MDS, especially to the nonlinear setting. We refer readers to Borg and Groenen (2005) for various forms of MDS and corresponding algorithms. Theoretical analysis in the literature of MDS concentrates on the case whether the objects from a higher dimensional space can be embedded into a lower dimensional space, Euclidean or non-Euclidean. See, for example, de Leeuw and Heiser (1982). However, we are not aware of any statistical results measuring the performance of MDS under randomness, such as perturbation analysis when the objects are sampled from a probabilistic model.

Other nonlinear methods for manifold learning includes ISOMAP (Tenenbaum, de Silva and Langford, 2000), local linear embedding (LLE) (Roweis and Saul, 2000) and diffusion maps (DM) (Coifman and Lafon, 2006), to name a few. Most of these procedures relies on PCA or local PCA. We illustrate this using the LLE, which is designed to be simple and intuitive, and can be computed efficiently. It mainly contains two steps: the first step is to determine the nearest neighbors for each data point, and catch the local geometric structure of the dataset through finding the

barycenter coordinate for those neighboring points; the second step is, by viewing the barycenter coordinates as the “weights” for the neighboring points, to evaluate the eigenvectors corresponding to the first several largest eigenvalues of the associated “affinity matrix” to locally embed the data to a lower dimensional Euclidean space. Surprisingly, despite its popularity in the manifold learning literature, [Wu and Wu \(2017\)](#) provided the asymptotic analysis of LLE until very recently.

## 5 Discussion

PCA is a powerful tool for data analytics ([Jolliffe, 2014](#)). Entering the era of big data, it is also finding many applications in modern machine learning problems. In this review article, we focus on clustering and community detection ([Lowrimore and Manton, 2016](#)), ranking, mixture model and manifold learning. Other applications, such as matrix completion ([Keshavan, Montanari and Oh, 2010](#)), phase synchronization ([Zhong and Boumal, 2017](#)), image segmentation ([Dambreville, Rathi and Tannen, 2006](#)) and functional data analysis ([Ramsay, 2016](#)), are not discussed here due to space limitations. Motivated by the fact that data are often collected and stored in distant places, many distributed algorithms for PCA have been proposed ([Qu et al., 2002](#); [Feldman, Schmidt and Sohler, 2013](#)) and shown to provide strong statistical guarantees ([Fan et al., 2017](#)).

Another important application of PCA is augmented principal component regression, which is an extension of the classical principal regression ([Kendall, 1957](#); [Hotelling, 1957](#); [Stock and Watson, 2002](#)). The basic idea is to assume that latent factors that impact on a large fraction of covariates also impact on the response variable. Therefore, we use PCA to extract latent factors from the covariates and then regress the response variable on these latent factors along with any augmented variables, resulting in an augmented factor models. An example of this is the multi-index prediction based on estimated latent factors in [Fan, Xue and Yao \(2017\)](#). A related topic to this is the supervised principal component regression ([Bair et al., 2006](#)).

Recently, a number of interesting theoretical results on the empirical principal components under weak signals have been developed ([Johnstone, 2001](#); [Hoyle and Rattray, 2003](#); [Baik, Ben Arous and P ech e, 2005](#); [Paul, 2007](#); [Koltchinskii and Lounici, 2017](#); [Wang and Fan, 2017](#)), which are closely related to the rapid advances in random matrix theory. Interested readers are referred to the literature above for details.

## References

ABBE, E. (2017). Community detection and stochastic block models: Recent developments. *arXiv preprint arXiv:1703.10146*.

- ABBE, E., BANDEIRA, A. S. and HALL, G. (2016). Exact recovery in the stochastic block model. *IEEE Trans. Inf. Theory* **62**, 471–487.
- ABBE, E., FAN, J., WANG, K. and ZHONG, Y. (2017). Entrywise eigenvector analysis of random matrices with low expected rank. *arXiv preprint arXiv:1709.09565*.
- AHN, S. C. and HORENSTEIN, A.Ř. (2013). Eigenvalue ratio test for the number of factors. *Econometrica* **81**, 1203–1227.
- ANANDKUMAR, A., GE, R., HSU, D., KAKADE, S. M. and TELGARSKY, M. (2014). Tensor decompositions for learning latent variable models. *J. Mach. Learn. Res.* **15**, 2773–2832.
- BAI, J. and LI, K. (2012). Statistical analysis of factor models of high dimension. *Ann. Statist.* **40**, 436–465.
- BAI, J. and NG, S. (2002). Determining the number of factors in approximate factor models. *Econometrica* **70**, 191–221.
- BAI, J. and NG, S. (2013). Principal components estimation and identification of static factors. *J. Econom.* **176**, 18–29.
- BAIK, J., BEN AROUS, G. and PÉCHÉ, S. (2005). Phase transition of the largest eigenvalue for nonnull complex sample covariance matrices. *Ann. Probab.* **33**, 1643–1697.
- BAIR, E., HASTIE, T., PAUL D. and TIBSHIRANI R. (2006). Prediction by supervised principal components. *J. Amer. Statist. Assoc.* **101**, 119–137.
- BALCAN, M.-F., KANCHANAPALLY, V., LIANG, Y. and WOODRUFF, D. (2014). Improved distributed principal component analysis. *Advances in Neural Information Processing Systems* **27**, 3113–3121.
- BALTRUNAS, L., MAKCINSKAS, T. and RICCI, F. (2010). Group recommendations with rank aggregation and collaborative filtering. In *Proc. ACM Conf. Recommender Syst.* 119–126.
- BARRAS, L., SCAILLET, O. and WERMERS, R. (2010). False discoveries in mutual fund performance: Measuring luck in estimated alphas. *J. Finance* **65**, 179–216.
- BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **57**, 289–300.
- BERK, J. B. and GREEN, R. C. (2004). Mutual fund flows and performance in rational markets. *J. Polit. Econ.* **112**, 1269–1295.



- BICKEL, P. J. and LEVINA, E. (2008). Covariance regularization by thresholding. *Ann. Statist.* **36**, 2577–2604.
- BORG, I. and GROENEN, P. J. (2005). *Modern Multidimensional Scaling: Theory and Applications*. Springer Science & Business Media.
- BOUTSIDIS, C., ZOUZIAS, A., MAHONEY, M. W. and DRINEAS, P. (2015). Randomized dimensionality reduction for  $k$ -means clustering. *IEEE Trans. Inf. Theory* **61**, 1045–1062.
- BRADLEY, R. A. and TERRY, M. E. (1952) Rank analysis of incomplete block designs: I. The method of paired comparisons. *Biometrika* **39**, 324–345
- CAI, T. T. and LIU, W. (2011). Adaptive thresholding for sparse covariance matrix estimation. *J. Amer. Statist. Assoc.* **106**, 672–684.
- CAI, T. T., REN, Z. and ZHOU, H. H. (2016). Estimating structured high-dimensional covariance and precision matrices: Optimal rates and adaptive estimation. *Electron. J. Stat.* **10**, 1–59.
- CAI, T. T. and ZHANG, A. (2016). Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *arXiv preprint arXiv:1605.00353*.
- CANDÉS, E. J., LI, X., MA, Y. and WRIGHT, J. (2011). Robust principal component analysis? *J. ACM* **58**, 11.
- CATONI, O. (2012). Challenging the empirical mean and empirical variance: A deviation study. *Ann. Inst. Henri Poincaré Probab. Stat.* **48**, 1148–1185.
- CHAMBERLAIN, G. and ROTHSCHILD, M. (1983). Arbitrage, factor structure, and mean variance analysis on large asset markets. *Econometrica* **51**, 1281–1304.
- CHANDRASEKARAN, V., SANGHAVI, S., PARRILO, P. A. and WILLSKY, A. S. (2011). Rank-sparsity incoherence for matrix decomposition. *SIAM J. Optim.* **21**, 572–596.
- CHEN, Y., FAN, J., MA, C. and WANG, K. (2017). Spectral method and regularized MLE are both optimal for top- $K$  ranking. *arXiv preprint arXiv:1608.08285*.
- CHEN, Y. and SUH, C. (2015). Spectral MLE: Top- $K$  rank aggregation from pairwise comparisons. In *Proceedings of the 32nd International Conference on Machine Learning* **37**, 371–380.
- COIFMAN, R. R. and LAFON, S. (2006). Diffusion maps. *Appl. Comput. Harmon. Anal.* **21**, 5–30.

- DAMBREVILLE, S. and RATHI, Y. and TANNEN, A. (2006). Shape-based approach to robust image segmentation using kernel PCA. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* **1**, 977–984.
- DAVIS, C. and KAHAN, W. M. (1970). The rotation of eigenvectors by a perturbation. III. *SIAM J. Numer. Anal.* **7**, 1–46.
- DE LEEUW, J. and HEISER, W. (1982). Theory of multidimensional scaling. In *Handbook of Statistics, Vol 2*, 285–316.
- DOPICO, F. M. (2000). A note on  $\sin \Theta$  theorems for singular subspace variations. *BIT Numer. Math.* **40**, 395–403.
- DWORK, C., KUMAR, R., NAOR, M. and SIVARKUMAR, D. (2001). Rank aggregation methods for the web. In *Proceedings of the 10th international conference on World Wide Web* 613–622.
- FAMA, E. and FRENCH, K. (1992). The cross-section of expected stock returns. *J. Finance* **47**, 427–465.
- FAMA, E. F. and FRENCH, K. R. (1993). Common risk factors in the returns on stocks and bonds. *J. Financ. Econ.* **33**, 3–56.
- FAN, J., FAN, Y. and LV, J. (2008). Large dimensional covariance matrix estimation using a factor model. *J. Econom.* **147**, 186–197.
- FAN, J., HAN, X. and GU, W. (2012). Estimating false discovery proportion under arbitrary covariance dependence. *J. Amer. Statist. Assoc.* **107**, 1019–1048.
- FAN, J., KE, Y., SUN, Q. and ZHOU, W.-X. (2017). FARM-Test: Factor-adjusted robust multiple testing with false discovery control. *arXiv preprint arXiv:1711.05386*.
- FAN, J., LI, Q. and WANG, Y. (2017). Estimation of high dimensional mean regression in the absence of symmetry and light tail assumptions. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **79**, 247–265.
- FAN, J., LIAO, Y. and LIU, H. (2016). An overview of the estimation of large covariance and precision matrices. *Econom. J.* **19**, C1–C32.
- FAN, J., LIAO, Y. and MINCHEVA, M. (2013). Large covariance estimation by thresholding principal orthogonal complements. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **75**, 1–44.
- FAN, J., LIAO, Y. and YAO, J. (2015). Power enhancement in high-dimensional cross-sectional tests. *Econometrica* **83**, 1497–1541.

- FAN, J., LIU, H. and WANG, W. (2017+). Large covariance estimation through elliptical factor models. *Ann. Statist.* To appear.
- FAN, J., WANG, D., WANG, K. and ZHU, Z. (2017). Distributed estimation of principal eigenspaces. *arXiv preprint arXiv:1702.06488* .
- FAN, J., WANG, W. and ZHONG, Y. (2016). An  $\ell_\infty$  eigenvector perturbation bound and its application to robust covariance estimation. *arXiv preprint arXiv:1603.03516*.
- FAN, J., WANG, W. and ZHU, Z. (2016). A shrinkage principle for heavy-tailed data: High-dimensional robust low-rank matrix recovery. *arXiv preprint arXiv:1603.08315*.
- FAN, J., XUE, L. and YAO, J. (2017). Sufficient forecasting using factor models. *J. Econom.* **201**, 292–306.
- FELDMAN, D., SCHMIDT, M. and SOHLER, C. (2013). Turning big data into tiny data: Constant-size coresets for  $k$ -means, PCA and projective clustering. In *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1434–1453. SIAM.
- GIULINI, I. (2017). Robust PCA and pairs of projections in a Hilbert space. *Electron. J. Stat.* **11**, 3903–3926.
- HOTELLING, H. (1957). The relations of the newer multivariate statistical methods to factor analysis. *The British Journal of Statistical Psychology* **10**, 69–79.
- HOYLE, D. and RATTRAY, M. (2003). Limiting form of the sample covariance eigenspectrum in PCA and kernel PCA. In *Proceedings of the 16th International Conference on Neural Information Processing Systems* 1181–1188.
- HSU, D. and KAKADE, S. M. (2013). Learning mixtures of spherical Gaussians: Moment methods and spectral decompositions. In *Proceedings of the 2013 ACM Conference on Innovations in Theoretical Computer Science* 11–19. ACM, New York.
- HUBER, P. J. (1981). *Robust Statistics*. Wiley.
- IZENMAN, A. J. (2008). *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*. Springer, New York.
- JOHNSTONE, I. M. (2001). On the distribution of the largest eigenvalue in principal components analysis. *Ann. Statist.* **29**, 295–327.
- JOLLIFFE, I. (2014). Principal component analysis. *Wiley StatsRef: Statistics Reference Online*.

- KE, Q. and KANADE, T. (2005). Robust  $L_1$  norm factorization in the presence of outliers and missing data by alternative convex programming. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.* **1**, 739–746.
- KENDALL, M. (1957). *A Course in Multivariate Analysis*. Griffin, London.
- KESHAVAN, R., MONTANARI, A. and OH, S. (2010) Matrix completion from a few entries. *IEEE Trans. Inf. Theory* **56**, 2980–2998
- KOLTCHINSKII, V. and LOUNICI, K. (2017). Normal approximation and concentration of spectral projectors of sample covariance. *Ann. Statist.* **45**, 121–157.
- LAM, C. and YAO, Q. (2012). Factor modeling for high-dimensional time series: Inference for the number of factors. *Ann. Statist.* **40**, 694–726
- LAN, W. and DU, L. (2017+). A factor-adjusted multiple testing procedure with application to mutual fund selection. *J. Bus. Econom. Statist.* To appear.
- LAWLEY, D. and MAXWELL, A. (1971). *Factor Analysis as a Statistical Method (2nd edition)*. Butterworth, London.
- LEEK, J. and STOREY, J. (2008). A general framework for multiple testing dependence. *Proc. Natn. Acad. Sci. USA* **105**, 18718–18723.
- LOWRIMORE, G. and MANTON, K. G. (2016). Cluster analysis: Overview. *Wiley StatsRef: Statistics Reference Online*. 1–19.
- MASSEY, K. (1997). *Statistical Models Applied to the Rating of Sports Teams*. Master’s thesis, Bluefield College, Bluefield, VA.
- MINSKER, S. (2015) Geometric median and robust estimation in Banach spaces. *Bernoulli* **21**, 2308–2335.
- MINSKER, S. (2016). Sub-Gaussian estimators of the mean of a random matrix with heavy-tailed entries. *arXiv preprint arXiv:1605.07129*.
- MISES, R. and POLLACZEK-GEIRINGER, H. (1929). Praktische verfahren der gleichungsauflösung. *ZAMM-Journal of Applied Mathematics and Mechanics/Zeitschrift für Angewandte Mathematik und Mechanik* **9**, 58–77.
- NEGAHBAN, S., OH, S. and SHAH, D. (2017). Rank centrality: Ranking from pairwise comparisons. *Operations Research* **65**, 266–287.
- O’ROURKE, S., VU, V and WANG, K. (2013). Random perturbation of low rank matrices: Improving classical bounds. *arXiv preprint arXiv:1311.2657*.

- PAUL, D. (2007). Asymptotics of sample eigenstructure for a large dimensional spiked covariance model. *Statist. Sinica* **17**, 1617–1642.
- PEARSON, K. (1901). On lines and planes of closest fit to systems of points in space. *The London, Edinburgh and Dublin Philosophical Magazine and Journal of Science*. **2**, 559–572.
- PERRY, P. O. and OWEN, A. B. (2010). A rotation test to verify latent structure. *J. Mach. Learn. Res.* **11**, 603–624.
- QU, Y., OSTROUCHOV, G., SAMATOVA, N. and GEIST, A. (2002). Principal component analysis for dimension reduction in massive distributed data sets. In *Proc. IEEE Int. Conf. Data Mining (ICDM)*.
- RAMSAY, J. O. (2016). Functional data analysis – Theory. *Wiley StatsRef: Statistics Reference Online*. 1–13.
- ROTHMAN, A. J., LEVINA, E. and ZHU, J. (2009). Generalized thresholding of large covariance matrices. *J. Amer. Statist. Assoc.* **104**, 177–186.
- ROWEIS, S. T. and SAUL, L. K. (2000). Nonlinear dimensionality reduction by locally linear embedding. *Science* **290**, 2323–2326.
- SHABALIN, A. A. and NOBEL, A. B. (2013). Reconstruction of a low-rank matrix in the presence of Gaussian noise. *J. Multivar. Anal.* **118**, 67–76.
- SHARPE, W. F. (1964). Capital asset prices: A theory of market equilibrium under conditions of risk. *J. Finance* **19**, 425–442.
- SHEN, H. and HUANG, J. Z. (2008). Sparse principal component analysis via regularized low rank matrix approximation. *J. Multivariate Anal.* **99**, 1015–1034.
- STEWART, G. W. (1991). Perturbation theory for the singular value decomposition. In *SVD and Signal Processing II: Algorithms, Analysis and Applications* (R. J. Vaccaro, ed.) 99–109.
- STEWART, M. (2006). Perturbation of the SVD in the presence of small singular values. *Linear Algebra Appl.* **419**, 53–77.
- STOCK, J. H. and WATSON, M. W. (2002). Forecasting using principal components from a large number of predictors. *J. Amer. Statist. Assoc.* **97**, 1167–1179.
- STOREY, J. D., TAYLOR, J. E. and SIEGMUND, D. (2004). Strong control, conservative point estimation and simultaneous conservative consistency of false discovery rates: A unified approach. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **66**, 187–205.

- SUN, Q., ZHOU, W.-X. and FAN, J. (2017). Adaptive Huber regression: Nonasymptotic optimality and phase transition. *arXiv preprint arXiv:1706.06991*.
- TENENBAUM, J. B., DE SILVA, V. and LANGFORD, J. C. (2000). A global geometric framework for nonlinear dimensionality reduction. *Science* **290**, 2319–2323.
- TROPP, J. A. (2012). User-friendly tail bounds for sums of random matrices. *Found. Comput. Math.* **12**, 389–434.
- TORGERSON, W. S. (1952). Multidimensional scaling: I. theory and method. *Psychometrika* **17**, 401–419.
- VERSHYNIN, R. (2012). Introduction to the non-asymptotic analysis of random matrices. In *Compressed Sensing* 210–268. Cambridge Univ. Press, Cambridge.
- VU, V. (2011). Singular vectors under random perturbation. *Random Structures & Algorithms* **39**, 526–538.
- WANG, R. (2015). Singular vector perturbation under Gaussian noise. *SIAM J. Matrix Anal. Appl.* **36**, 158–177.
- WANG, W. and FAN, J. (2017). Asymptotics of empirical eigen-structure for ultra-high dimensional spiked covariance model. *Ann. Statist.* **45**, 1342–1374.
- WEDIN, P.-Å. (1972). Perturbation bounds in connection with singular value decomposition. *Nordisk Tidskr. Informationsbehandling (BIT)* **12**, 99–111.
- WU, H.-T. and WU, N. (2017). Think globally, fit locally under the manifold setup: Asymptotic analysis of locally linear embedding. *arXiv preprint arXiv:1703.04058*.
- YI, X., CARAMANIS, C. and SANGHAVI, S. (2016). Solving a mixture of many random linear equations by tensor decomposition and alternating minimization. *arXiv preprint arXiv:1608.05749*.
- YU, Y., WANG, T. and SAMWORTH, R. J. (2015). A useful variate of the Davis–Kahan theorem for statisticians. *Biometrika* **102**, 315–323.
- ZHANG, T. and LERMAN, G. (2014). A novel  $M$ -estimator for robust PCA. *J. Mach. Learn. Res.* **15**, 749–808.
- ZHONG, Y. and BOUMAL, N. (2017). Near-optimal bounds for phase synchronization. *arXiv preprint arXiv:1703.06605*.