

A Selective Overview of Sparse Principal Component Analysis

Hui Zou and Lingzhou Xue
University of Minnesota and Pennsylvania State University

Abstract—Principal component analysis (PCA) is a widely used technique for dimension reduction, data processing and feature extraction. The three tasks are particularly useful and important in high-dimensional data analysis and statistical learning. However, the regular PCA encounters great fundamental challenges under high-dimensionality and may produce ‘wrong’ results. As a remedy, sparse PCA has been proposed and studied. Sparse PCA is shown to offer a ‘right’ solution under high-dimensions. In this article, we review methodological and theoretical developments of sparse PCA, as well as its applications in scientific studies.

I. PCA

Principal component analysis (PCA) was invented by (Pearson 1901). As a dimension reduction and feature extraction method, PCA has numerous applications in statistical learning, such as handwritten zip code classification (Hastie et al. 2009), human face recognition (Hancock et al. 1996), eigengenes analysis (Alter et al. 2000), gene shaving (Hastie et al. 2000), and so on. It would not be exaggerating to say that PCA is one of the most widely used and most important multivariate statistical techniques.

This review article focuses on the high-dimensional extension of the regular PCA, which is often called sparse PCA. There are several popular sparse PCA methods in the literature, which will be reviewed in Section 2. Their formulations are different but related, because the regular PCA has several equivalent definitions from different viewing angles. These definitions are equivalent without sparsity constraints, and differ with sparsity constraints. To be self-contained, we briefly discuss the several views of PCA in the following.

From a dimension reduction perspective, PCA can be described as a set of orthogonal linear transformations of the original variables such that the transformed variables maintain the information contained in the original variables as much as possible. Specifically, let \mathbf{X} be a $n \times p$ data matrix, where n and p are the number of observations and the number of variables, respectively. For ease of presentation, assume the column means of \mathbf{X} are all 0. The first principal component is defined as $Z_1 = \sum_{j=1}^p \alpha_{1j} X_j$ where $\alpha_1 = (\alpha_{11}, \dots, \alpha_{1p})^T$ is chosen to maximize the variance of Z_1 , i.e.,

$$\alpha_1 = \arg \max_{\alpha} \alpha^T \hat{\Sigma} \alpha \quad \text{subject to } \|\alpha_1\| = 1 \quad (1)$$

with $\hat{\Sigma} = \frac{\mathbf{X}^T \mathbf{X}}{n}$. The rest principal components can be defined sequentially as follows:

$$\alpha_{k+1} = \arg \max_{\alpha} \alpha^T \hat{\Sigma} \alpha \quad (2)$$

subject to

$$\|\alpha\| = 1 \quad \text{and} \quad \alpha^T \alpha_l = 0, \quad \forall 1 \leq l \leq k. \quad (3)$$

This definition implies that the first K loading vectors are the first K eigenvectors of $\hat{\Sigma}$.

The eigen-decomposition formulation of PCA also relates PCA to the singular value decomposition (SVD) of \mathbf{X} . Let the SVD of \mathbf{X} be

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

where \mathbf{D} is a diagonal matrix with diagonal elements d_1, \dots, d_p in a descending order, and \mathbf{U} and \mathbf{V} are $n \times p$ and $p \times p$ orthonormal matrices, respectively. Because the columns of \mathbf{V} are the eigenvectors of $\hat{\Sigma}$, \mathbf{V} is the loading matrix of the principal components. By $\mathbf{X} \mathbf{V} = \mathbf{U} \mathbf{D}$, we see that $Z_k = U_k d_k$ where U_k is the k th column of \mathbf{U} . Note that SVD can be interpreted as the best low rank approximation to the data matrix.

PCA has another geometric interpretation, as the closest linear manifold approximation of the observed data. This definition actually matches the construction of PCA considered by Pearson (1901). Let \mathbf{x}_i be the i th row of \mathbf{X} . Consider the first k principal components jointly $\mathbf{V}_k = [V_1 | \dots | V_k]$. By definition, \mathbf{V}_k is a $p \times k$ orthonormal matrix. Project each observation to the linear space spanned by $\{V_1, \dots, V_k\}$. The projection operator is $\mathbf{P}_k = \mathbf{V}_k \mathbf{V}_k^T$ and the projected data is $\mathbf{P}_k \mathbf{x}_i$, $1 \leq i \leq n$. One way to define the best projection is by minimizing the total ℓ_2 approximation error

$$\min_{\mathbf{V}_k} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{V}_k \mathbf{V}_k^T \mathbf{x}_i\|^2. \quad (4)$$

It is easy to show that the solution is exactly the first k principal components.

In applications variables can have different scales and units. Practitioners often standardize each variable such that its marginal sample variance is one. When this practice is applied to PCA, the resulting covariance matrix of standardized variables is the sample correlation matrix of the raw variables. Note that the eigenvalues and eigenvectors of the correlation matrix can be different from those of the covariance matrix.

II. METHODS FOR SPARSE PRINCIPAL COMPONENTS

Each principal component is a linear combination of all p variables, which makes it difficult to interpret the derived principal components as new features. Rotation techniques are commonly used to help practitioners to interpret principal components (Jolliffe 1995). Vines (2000) considered simple principal components by restricting the loadings to take values from a small set of allowable integers such as 0, 1 and -1. This restriction may be useful for certain applications but not all. Simple thresholding is an ad hoc way to achieve sparse loadings by setting the loadings with absolute values smaller than a threshold to zero. Although the simple thresholding is

frequently used in practice, it can be potentially misleading in various respects (Cadima & Jolliffe 1995). Sparse variants of PCA aim to achieve a good balance between variance explained (dimension reduction) and sparse loadings (interpretability).

Sparse learning is ubiquitous in high-dimensional data analysis. Prior to the sparse principal component problem, an important question is how to select variables in high-dimensional regression. For the regression problem, the *lasso* proposed in Tibshirani (1996) is a very promising technique for simultaneous variable selection and prediction. The lasso regression is an ℓ_1 penalized least squares method. The use of ℓ_1 penalization yields a sparse solution and also permits efficient computations.

A. SCoTLASS

Inspired by the lasso regression, Jolliffe et al. (2003) proposed a procedure called *SCoTLASS* to obtain sparse loadings by directly imposing an ℓ_1 constraint on the loading vector.

SCoTLASS extends the standard PCA by taking the variance maximization perspective of the PCA. It successively maximizes the variance

$$a_k^T \hat{\Sigma} a_k \quad (5)$$

subject to

$$a_k^T a_k = 1 \quad \text{and (for } k \geq 2) \quad a_h^T a_k = 0, \quad h < k; \quad (6)$$

and the extra ℓ_1 constraints

$$\sum_{j=1}^p |a_{kj}| \leq t \quad (7)$$

for some tuning parameter t . However, SCoTLASS is high computational cost which makes it an impractical solution for high-dimensional data analysis. It motivated researchers to consider more efficient proposals for sparse principal components. A related but more efficient approach is the generalized power method presented in Section 2.5.

B. SPCA

After SCoTLASSO, the first computational efficient sparse PCA algorithm for high-dimensional data was introduced by Zou et al. (2006). Their method is named *Sparse Principal Component Analysis* (SPCA). Before reviewing the technical details, let us consider the application of SPCA to the pitprops data (Jeffers 1967), a classical example showing the difficulty of interpreting principal components. The pitprops data has 180 observations and 13 measured variables. In Zou et al. (2006) the first six ordinary principal components and the first six sparse principal components are computed. Here we only cite the results of the first three principal components in Table 1. Compared with the standard PCA, SPCA generated very sparse loading structures without losing much variance.

In the original lasso paper Tibshirani used a quadratic programming solver to compute the lasso regression estimator, which is not very efficient for high-dimensional data. Efron et al. (2004) derived the first efficient algorithm named *LARS*

	PCA			SPCA		
	PC1	PC2	PC3	PC1	PC2	PC3
topdiam	-.404	.218	-.207	-.477		
length	-.406	.186	-.235	-.476		
moist	-.124	.541	.141		.785	
testsg	-.173	.456	.352		.620	
ovensg	-.057	-.170	.481	.177		.640
ringtop	-.284	-.014	.475			.589
ringbut	-.400	-.190	.253	-.250		.492
bowmax	-.294	-.189	-.243	-.344	-.021	
bowdist	-.357	.017	-.208	-.416		
whorls	-.379	-.248	-.119	-.400		
clear	.011	.205	-.070			
knots	.115	.343	.092		.013	
diaknot	.113	.309	-.326			-.015
variance	32.4	18.3	14.4	28.0	14.0	13.3

TABLE 1

COMPARE PCA AND SPCA ON THE PITPROPS DATA. EMPTY CELLS MEAN ZERO LOADINGS. THE VARIANCE OF SPCA IS EXPECTED TO BE SMALLER THAN THAT OF PCA, BY THE DEFINITION OF PCA. THE DIFFERENCES IN VARIANCE ARE SMALL.

for computing the entire solution path of the lasso regression model with high-dimensional data. Motivated by LARS, Zou et al. (2006) proposed to tackle the sparse principal component problem from a regression formulation. The resulting algorithm is SPCA.

SPCA extends the linear manifold approximation view of the PCA to derive sparse loadings. Recall that the first principal component can be defined as

$$\alpha_1 = \arg \min_{\alpha, \beta} \sum_{i=1}^n \|\mathbf{x}_i - \alpha \alpha^T \mathbf{x}_i\|^2 \quad (8)$$

subject to $\|\alpha\|^2 = 1$.

We reformulate (8) as

$$\arg \min_{\alpha, \beta} \sum_{i=1}^n \|\mathbf{x}_i - \alpha \beta^T \mathbf{x}_i\|^2 \quad (9)$$

subject to $\|\alpha\|^2 = 1$ and $\alpha = \beta$.

The following theorem says that we can drop the equality constraint in (9) and still recover the first loading vector exactly.

Theorem 1 (Zou et al. (2006)). *For any $\lambda_0 > 0$, let*

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n \|\mathbf{x}_i - \alpha \beta^T \mathbf{x}_i\|^2 + \lambda_0 \|\beta\|^2 \quad (10)$$

subject to $\|\alpha\|^2 = 1$.

Then $\hat{\beta} \propto V_1$.

In Theorem 1 the extra ℓ_2 term $\lambda_0 \|\beta\|^2$ is not needed when $p < n$. When $p > n$, any $\lambda_0 > 0$ should be used and it does not affect the normalized β_1 . By dropping the equality constraint $\alpha = \beta$, we can use an alternating minimization algorithm to optimize the criterion in (10) because α and β are separated variables. With a fixed α , the optimization problem over β is a regression problem.

Based on Theorem 1, we can impose a sparse penalty on β to gain zero loading because the normalizing step does not

change the support of β . The SPCA for the first principal component is defined as

$$(\hat{\alpha}, \hat{\beta}) = \arg \min_{\alpha, \beta} \sum_{i=1}^n \|\mathbf{x}_i - \alpha \beta^T \mathbf{x}_i\|^2 + \lambda_0 \|\beta\|^2 + \lambda_1 \|\beta\|_1$$

subject to $\|\alpha\|^2 = 1$.

and the output loading vector is $\hat{V}_1 = \hat{\beta} / \|\hat{\beta}\|$. For $n > p$, we can let $\lambda_0 = 0$ and solving β with a fixed α is a lasso regression problem, which can be done efficiently. When $n < p$, we need to use a positive λ_0 (e.g., $\lambda_0 = 10^{-3}$), solving β with a fixed α is an elastic net regression problem (Zou & Hastie 2005), which can be done efficiently as well.

Theorem 1 can be generalized to handle the first k principal components simultaneously, as stated in the next theorem.

Theorem 2 (Zou et al. (2006)). *Suppose we are considering the first k principal components. Let $\mathbf{A}_{p \times k} = [\alpha_1, \dots, \alpha_k]$ and $\mathbf{B}_{p \times k} = [\beta_1, \dots, \beta_k]$. For any $\lambda_0 > 0$, let*

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{A} \mathbf{B}^T \mathbf{x}_i\|^2 + \lambda_0 \sum_{j=1}^k \|\beta_j\|^2$$

subject to $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}$.

Then $\hat{\beta}_j \propto V_j$ for $j = 1, 2, \dots, k$.

Then the SPCA criterion for the first k sparse principal components is defined as

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{A} \mathbf{B}^T \mathbf{x}_i\|^2 + \lambda_0 \sum_{j=1}^k \|\beta_j\|^2 + \sum_{j=1}^k \lambda_{1,j} \|\beta_j\|_1$$

subject to $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}$,

where different $\lambda_{1,j}$ s are allowed for penalizing the loadings of different principal components.

Zou et al. (2006) proposed an alternating algorithm to minimize the SPCA criterion (13).

B given A: For each j , let $Y_j^* = \mathbf{X} \alpha_j$. It can shown that $\hat{\mathbf{B}} = [\hat{\beta}_1, \dots, \hat{\beta}_k]$ and each $\hat{\beta}_j$ is obtained via

$$\hat{\beta}_j = \arg \min_{\beta_j} \|Y_j^* - \mathbf{X} \beta_j\|^2 + \lambda_0 \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1. \quad (14)$$

One can use either the LARS-EN algorithm (Zou & Hastie 2005) or the cyclic coordinate descent algorithm (Friedman et al. 2007) to solve (14). Both algorithms are efficient for high-dimensional data.

A given B: is fixed, the optimization problem of \mathbf{A} is

$$\arg \min_{\mathbf{A}} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{A} \mathbf{B}^T \mathbf{x}_i\|^2 = \|\mathbf{X} - \mathbf{X} \mathbf{B} \mathbf{A}^T\|^2,$$

subject to $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}$. This is called a reduced rank *Procrustes rotation* problem in Zou et al. (2006) because when $k = p$ it is the Procrustes rotation problem (Mardia et al. 1979). Zou et al. (2006) derived an explicit solution to the reduced rank Procrustes rotation problem. We compute the SVD

$$(\mathbf{X}^T \mathbf{X}) \mathbf{B} = \mathbf{U} \mathbf{D} \mathbf{V}^T, \quad (15)$$

and set $\hat{\mathbf{A}} = \mathbf{U} \mathbf{V}^T$.

The SPCA algorithm iterates between the elastic net regression step and the SVD step till convergence. The output is the normalized \mathbf{B} matrix: $\hat{V}_j = \hat{\beta}_j / \|\hat{\beta}_j\|$, $1 \leq j \leq k$.

Zou et al. (2006) derived another SPCA criterion to further speed up the computation efficiency. The derivation is based on the observation that Theorem 2 is valid for all $\lambda_0 > 0$. It turns out that a thrifty solution emerges if λ_0 is taken to be a large constant.

Theorem 3 (Zou et al. (2006)). *Let $\hat{V}_j(\lambda_0) = \frac{\hat{\beta}_j}{\|\hat{\beta}_j\|}$ ($j = 1, \dots, k$) be the sparse loadings defined in (13). Let $(\hat{\mathbf{A}}, \hat{\mathbf{B}})$ be the solution of the optimization problem*

$$(\hat{\mathbf{A}}, \hat{\mathbf{B}}) = \arg \min_{\mathbf{A}, \mathbf{B}} -2\text{Tr}(\mathbf{A}^T \mathbf{X}^T \mathbf{X} \mathbf{B}) + \sum_{j=1}^k \|\beta_j\|^2 + \sum_{j=1}^k \lambda_{1,j} \|\beta_j\|_1$$

subject to $\mathbf{A}^T \mathbf{A} = \mathbf{I}_{k \times k}$.

When $\lambda_0 \rightarrow \infty$, $\hat{V}_j(\lambda_0) \rightarrow \frac{\hat{\beta}_j}{\|\hat{\beta}_j\|}$.

Solving (16) can also be done via an alternating minimization algorithm. Given \mathbf{A} , we have that for each j ,

$$\hat{\beta}_j = \arg \min_{\beta_j} -2\alpha_j^T (\mathbf{X}^T \mathbf{X}) \beta_j + \|\beta_j\|^2 + \lambda_{1,j} \|\beta_j\|_1, \quad (17)$$

and the solution is given by

$$\hat{\beta}_j = S(\mathbf{X}^T \mathbf{X} \alpha_j, \frac{\lambda_{1,j}}{2})$$

where $S(Z, \gamma)$ is the soft-thresholding operator on a vector $Z = (z_1, \dots, z_p)$ with thresholding parameter γ and

$$S(Z, \gamma)_j = (|z_j| - \gamma)_+ \text{sgn}(z_j), \quad 1 \leq j \leq p.$$

Given \mathbf{B} , the solution of \mathbf{A} is again $\hat{\mathbf{A}} = \mathbf{U} \mathbf{V}^T$ where \mathbf{U}, \mathbf{V} are from the SVD of $(\mathbf{X}^T \mathbf{X}) \mathbf{B}$: $(\mathbf{X}^T \mathbf{X}) \mathbf{B} = \mathbf{U} \mathbf{D} \mathbf{V}^T$.

C. A semidefinite programming approach

We introduce some necessary notation first. Use $\text{Card}(M)$ to denote the number of nonzero element of M , where M can be a vector of a matrix. The notation $|M|$ means that we replace each element of M with its absolute value. Let $\mathbf{1}_p$ be the p -vector of 1.

Consider the first k -sparse principal component with at most k nonzero loadings. A natural definition of the optimal k -sparse loading vector is

$$\arg \max_{\alpha} \alpha^T \hat{\Sigma} \alpha \quad (18)$$

$$\text{subject to } \|\alpha\| = 1, \quad \text{Card}(\alpha) \leq k.$$

When $k = p$, then the above definition gives the loadings of the first principal component. However, (18) is nonconvex and computationally difficult, especially when p is large. Convex relaxation is a standard technique used in operational research to handle difficult nonconvex problems. d'Aspremont et al. (2007) developed a convex relation of (18), which is expressed as a semidefinite programming problem.

Let $\mathbf{P} = \alpha \alpha^T$. We write $\alpha^T \hat{\Sigma} \alpha = \text{Tr}(\hat{\Sigma} \mathbf{P})$. The norm-1 constraint on α leads to a linear equality constraint on \mathbf{P} :

$\text{Tr}\mathbf{P} = 1$. Moreover, the cardinality constraint $\|\alpha\|_0 \leq k$ implies $\text{Card}(\mathbf{P}) \leq k^2$. Hence, we consider the following optimization problem of \mathbf{P} :

$$\begin{aligned} & \arg \max_{\mathbf{P}} \text{Tr}(\hat{\Sigma}\mathbf{P}) \\ & \text{subject to } \text{Tr}\mathbf{P} = 1, \quad \text{Card}(\mathbf{P}) \leq k^2, \\ & \quad \mathbf{P} \succeq 0 \quad \text{and} \quad \text{rank}(\mathbf{P}) = 1. \end{aligned} \quad (19)$$

The above formulation in (19) is still nonconvex and difficult to handle due to the cardinality constraint and the rank one constraint. By definition, \mathbf{P} is symmetric and $\mathbf{P}^2 = \mathbf{P}$. Observe that

$$\|\mathbf{P}\|_F^2 = \text{Tr}(\mathbf{P}^T\mathbf{P}) = \text{Tr}(\mathbf{P}) = 1.$$

By Cauchy-Schwartz,

$$\mathbf{1}_p^T |\mathbf{P}| \mathbf{1}_p \leq \sqrt{\text{Card}(\mathbf{P})} \|\mathbf{P}\|_F \leq k.$$

Therefore, d'Aspremont et al. (2007) suggested to relax the cardinality constraint in (19) to a linear inequality constraint $\mathbf{1}_p^T |\mathbf{P}| \mathbf{1}_p \leq k$. Furthermore, they dropped the rank one constraint and ended up with the *DSPCA* formulation:

$$\begin{aligned} & \arg \max_{\mathbf{P}} \text{Tr}(\hat{\Sigma}\mathbf{P}) \\ & \text{subject to} \\ & \quad \text{Tr}\mathbf{P} = 1, \\ & \quad \mathbf{1}_p^T |\mathbf{P}| \mathbf{1}_p \leq k, \\ & \quad \mathbf{P} \succeq 0. \end{aligned} \quad (20)$$

The above is recognized as a semidefinite programming problem and can be solved by software such as *SDPT3*.

DSPCA only solves for \mathbf{P} but not α . To compute the loading vector α , d'Aspremont et al. (2007) recommended truncating \mathbf{P} and retaining only the dominant (sparse) eigenvector of \mathbf{P} . For the second sparse principal component, it is suggested to replace $\hat{\Sigma}$ with $\hat{\Sigma} - (\alpha^T \hat{\Sigma} \alpha) \alpha \alpha^T$ in (20). The same procedure can be repeated to compute the rest sparse principal components.

For larger problems, d'Aspremont et al. (2007) discussed a Nesterov's smooth minimization technique to handle *DSPCA*. The computation complexity of the algorithm is $O(p^4 \sqrt{\log(p)}/\epsilon)$, where ϵ is the numerical accuracy of the solution. d'Aspremont et al. (2008) discussed a greedy algorithm to speed up the computation. An alternating direction method of multipliers was proposed in Ma (2013a).

DSPCA formulation generated many interests in the operational research and machine learning communities. Some follow-up works include Lu & Zhang (2012), Vu et al. (2013) and d'Aspremont (2011), among others.

D. Iterative thresholding methods

PCA can be done via the singular value decomposition (SVD) of the data matrix. Thus, it is natural to consider a sparse PCA algorithm based on the SVD of \mathbf{X} . This idea was explored in Shen & Huang (2008) and Witten et al. (2009).

Let the SVD of \mathbf{X} be $\mathbf{X} = \mathbf{U}\mathbf{D}\mathbf{V}^T$. Consider the first principal component. We know the loading vector is V_1 , the

first column of \mathbf{V} . It is a well known result that SVD of \mathbf{X} is related to the best lower rank approximation of \mathbf{X} (Eckart & Young 1936). Specifically, let \tilde{U} be a norm-1 n -vector and \tilde{V} be a p -vector. Consider $\tilde{U}\tilde{V}^T$ as a rank one approximation of \mathbf{X} . The best rank one approximation is defined as

$$\min_{\tilde{U}, \tilde{V}} \|\mathbf{X} - \tilde{U}\tilde{V}^T\|_F^2 \quad \text{subject to} \quad \|\tilde{U}\| = 1, \quad (21)$$

and the solution is $\tilde{U} = U_1$ and $\tilde{V} = d_1 V_1$ where d_1 is the first singular value.

Based on (21) Shen & Huang (2008) proposed the following optimization problem

$$(\hat{U}, \hat{V}) = \arg \min_{U, V} \|\mathbf{X} - UV^T\|_F^2 + \lambda \|V\|_1 \quad \text{subject to} \quad \|U\| = 1, \quad (22)$$

and the sparse loading vector is normalized \hat{V} , $\frac{\hat{V}}{\|\hat{V}\|}$. An alternating minimization algorithm is used to solve (22). Note that given V , the optimal U is $U = \mathbf{X}V/\|\mathbf{X}V\|$. Given U , the optimal V is

$$\arg \min_V -2\text{Tr}(\mathbf{X}^T UV^T) + \|V\|^2 + \lambda \|V\|_1$$

and the solution is given by the soft-thresholding operator:

$$V = S(\mathbf{X}^T U, \frac{\lambda}{2})$$

Thus, Shen and Huang's method is an iterative thresholding algorithm.

Note that the above procedure is similar in spirit to the SPCA algorithm in (16). The big difference is that SPCA solves k components simultaneously, but Shen and Huang's method only deals with one component at a time.

Shen & Huang (2008) proposed to sequentially compute the rest sparse principal components. Suppose that we have computed the first k (U, V) pairs, let $\mathbf{X}_{(k+1)} = \mathbf{X} - \sum_{i=1}^k U_i V_i^T$ and then the iterative thresholding algorithm is applied to $\mathbf{X}_{(k+1)}$ to get $(U_{(k+1)}, V_{(k+1)})$. The normalized $V_{(k+1)}$ is the loading vector of the $(k+1)$ th sparse principal component. The λ parameter is allowed to differ for different principal components.

In the same vein Witten et al. (2009) proposed a penalized matrix decomposition (PMD) criterion as follows

$$(\hat{U}, \hat{V}, \hat{d}) = \arg \min_{U, V, d} \|\mathbf{X} - dUV^T\|_F^2 \quad (23)$$

$$\text{subject to} \quad \|U\| = 1, \|U\|_1 \leq c_1; \quad \|V\| = 1, \|V\|_1 \leq c_2.$$

By straightforward calculation, it can be shown that (23) is equivalent to the following optimization problem

$$(\hat{U}, \hat{V}) = \arg \max_{U, V} U^T \mathbf{X} V \quad (24)$$

$$\text{subject to} \quad \|U\| = 1, \|U\|_1 \leq c_1; \quad \|V\| = 1, \|V\|_1 \leq c_2.$$

and $\hat{d} = \hat{U}^T \mathbf{X} \hat{V}$.

They also used an alternating minimization algorithm to compute (24). Given V , we update U by solving

$$\max_U U^T \mathbf{X} V \quad \text{subject to} \quad \|U\| = 1, \|U\|_1 \leq c_1. \quad (25)$$

Given U , we update V by solving

$$\max_V U^T \mathbf{X} V \quad \text{subject to} \quad \|V\| = 1, \|V\|_1 \leq c_2. \quad (26)$$

The equality constraint $\|U\| = 1, \|V\| = 1$ in (25) and (26) can be replaced with inequality constraint $\|U\| \leq 1, \|V\| \leq 1$ and the solutions remain the same. So, (25) and (26) are examples of the following convex optimization problem

$$\hat{Z} = \arg \max_Z Z^T R \quad \text{subject to} \quad \|Z\| \leq 1, \|Z\|_1 \leq c. \quad (27)$$

It is easy to see that the solution to (27) is

$$\hat{Z} = \frac{S(R, \Delta_c)}{\|S(R, \Delta_c)\|},$$

where S is the soft-thresholding operator and Δ_c is selected as follows: $\Delta_c = 0$ if $\|\frac{R}{\|R\|}\|_1 \leq c$, otherwise $\Delta_c > 0$ is chosen to satisfy $\|\hat{Z}\|_1 = c$.

E. A generalized power method

Consider the first principal component. By the variance maximization definition, a direct formulation of ℓ_1 constrained sparse principal component is

$$\begin{aligned} & \arg \max_{\|\alpha\|=1} \alpha^T \mathbf{X}^T \mathbf{X} \alpha \\ & \text{subject to } \|\alpha\|_1 \leq t. \end{aligned} \quad (28)$$

Equivalently, we can solve

$$\begin{aligned} & \arg \max_{\|\alpha\|=1} \sqrt{\alpha^T \mathbf{X}^T \mathbf{X} \alpha} \\ & \text{subject to } \|\alpha\|_1 \leq t. \end{aligned} \quad (29)$$

Journée et al. (2010) considered the Lagrangian form of (29)

$$\arg \max_{\|\alpha\|=1} \sqrt{\alpha^T \mathbf{X}^T \mathbf{X} \alpha} - \lambda \|\alpha\|_1. \quad (30)$$

They offered a *generalized power method* for solving (31). Their idea takes advantage of this simple observation: let $\tilde{U} = \arg \max_{\|U\|=1} U^T Z$, then $\tilde{U} = Z/\|Z\|$ and $\tilde{U}^T Z = \|Z\|$. Thus, an equivalent formulation of (31) is

$$\begin{aligned} (U^*, \alpha^*) &= \arg \max_{U, \alpha} U^T \mathbf{X} \alpha - \lambda \|\alpha\|_1 \\ & \text{subject to} \quad \|U\| = 1, \|\alpha\| = 1. \end{aligned} \quad (31)$$

Notice that the formulation (31) is the Lagrangian form of the PMD formulation (24) without imposing the ℓ_1 constraint on U .

For any U , the optimal α and $\mathbf{X}^T U$ must share the same sign for each component. Let $z_j = |\alpha_j|$, and $Z = |\alpha|$. Then the optimal Z^* must satisfy

$$\begin{aligned} Z^* &= \arg \max_Z \sum_{j=1}^p (|\mathbf{X}^T U|_j - \lambda) z_j \\ & \text{subject to} \quad z_j \geq 0, \quad \sum_{j=1}^p z_j^2 = 1. \end{aligned} \quad (32)$$

When $|\mathbf{X}^T U|_j - \lambda \leq 0$, $z_j^* = 0$. By Cauchy-Schwartz, it is easy to see that the solution to (32) is

$$z_j^* = \frac{(|\mathbf{X}^T U|_j - \lambda)_+}{\sqrt{\sum_{j=1}^p (|\mathbf{X}^T U|_j - \lambda)_+^2}}, \quad (33)$$

which yields

$$\alpha = \frac{S(\mathbf{X}^T U, \lambda)}{\|S(\mathbf{X}^T U, \lambda)\|}, \quad (34)$$

where S is the soft-thresholding operator. Plugging (33) back to the objective function in (31), we obtain a new optimization criterion of U :

$$U^* = \arg \max_{U: \|U\|=1} \sqrt{\sum_{j=1}^p (|\mathbf{X}^T U|_j - \lambda)_+^2},$$

or equivalently

$$U^* = \arg \max_{U: \|U\| \leq 1} \sum_{j=1}^p (|\mathbf{X}^T U|_j - \lambda)_+^2. \quad (35)$$

Once U^* is solved, we have

$$\alpha^* = \frac{S(\mathbf{X}^T U^*, \lambda)}{\|S(\mathbf{X}^T U^*, \lambda)\|}.$$

Solving U^* is a n -dimensional optimization problem, although the original formulation (31) is a p -dimensional optimization problem. When $p \gg n$, the generalized power method achieves great computational savings. Moreover, the objective function in (35) is differentiable and convex, and the constraint set is compact and convex. Journée et al. (2010) used an efficient gradient method to compute U^* and analyzed its convergence property. They also showed that the generalized power method can be extended to handle the first k principal components jointly.

There are other proposals for constructing sparse principal components such as the truncated power method in Yuan & Zhang (2013) and the exact and greedy algorithms in Moghaddam et al. (2006).

III. THEORETICAL RESULTS

Theoretical analysis of sparse PCA received considerable attention in the past decade. In what follows, we first discuss the inconsistency of the classical PCA in the high-dimensional setting, and then present recent theoretical developments of sparse PCA.

A. Inconsistency of PCA under high-dimensions

Statistical analysis of PCA views $\hat{\Sigma}$ as the empirical covariance matrix and there is the population PCA on the true covariance matrix Σ . In the conventional setting where the dimension is fixed and the sample size increases, the principal eigenvectors of the sample covariance matrix are the consistent estimates of the principal eigenvectors of the corresponding population covariance matrix (Anderson 2003).

However, the sample principal eigenvectors are inconsistent estimates of the corresponding population principal eigenvectors in the high-dimensional setting where the dimension is no longer fixed and may be much larger than the sample size. The inconsistency phenomenon was first observed in the unsupervised learning theory literature in physics (for example, Biehl & Mietzner (1994) and Watkin & Nadal (1994)). About a decade ago, a series of papers in the statistics

literature (for example, Baik & Silverstein (2006), Paul (2007), Nadler (2008), Johnstone & Lu (2009) and Jung & Marron (2009)) investigated the inconsistency results of the classical PCA when estimating the leading principal eigenvectors in the high-dimensional setting. Baik & Silverstein (2006), Paul (2007) and Nadler (2008) showed that when $\lim_{n \rightarrow \infty} p/n = \gamma \in (0, 1)$, the largest eigenvalue λ_1 is of unit multiplicity and $\lambda_1 \leq \sqrt{\gamma}$, the leading sample principal eigenvector \hat{v}_1 is asymptotically orthogonal to the leading population principal eigenvector v_1 almost surely, that is,

$$P\left(\lim_{n \rightarrow \infty} |\mathbf{v}_1^T \hat{v}_1| = 0\right) = 1.$$

Johnstone & Lu (2009) considered the rank-one case and gave the sufficient and necessary condition for the consistence estimation of the leading population principal eigenvector. Let $R(\hat{v}_1, v_1) = \cos \alpha(\hat{v}_1, v_1)$ be the cosine of the angle between \hat{v}_1 and v_1 , and let $\omega = \lim_{n \rightarrow \infty} \|\mathbf{v}_1\|^2/\sigma^2$ be the limiting signal-to-noise ratio. Johnstone & Lu (2009) proved that

$$P\left(\lim_{n \rightarrow \infty} R^2(\hat{v}_1, v_1) = R_\infty^2(\omega, c)\right) = 1$$

where $c = \lim_{n \rightarrow \infty} p/n$ and

$$R_\infty^2(\omega, c) = (\omega^2 - c)_+ / (\omega^2 + c\omega).$$

Note that $R_\infty^2(\omega, c) < 1$ if and only if $c > 0$. Thus, \hat{v}_1 is a consistent estimate of v_1 if and only if $c = 0$, which implies the inconsistency of the classical PCA in the high-dimensional setting. Jung & Marron (2009) further studied the strong inconsistency of the leading sample principal eigenvector in the high dimension and low sample size context where the sample size is fixed and the dimension increases.

These inconsistency results call for new formulation of principal components that are consistent estimators of the population principal components under high-dimensions.

B. Consistency of sparse PCA

In recent years, there is a series of papers to develop the theoretical properties of the sparse PCA in the statistics literature. The consistency results are established for various regularized estimators of the leading eigenvectors. Under the rank-one scenario with $n^{-1} \log(n \vee p) \rightarrow 0$ as $n \rightarrow \infty$, Johnstone & Lu (2009) established a consistency result for the classical PCA performed on a selected subset of variables satisfying $\hat{\sigma}^2 \geq \sigma^2(1 + \alpha_n)$, where $\alpha_n = \alpha(n^{-1} \log(n \vee p))^{1/2}$. Specifically, Johnstone & Lu (2009) proved that the estimated principal eigenvector \hat{v}_1^I obtained via the subset selection rule is consistent:

$$P\left(\lim_{n \rightarrow \infty} \alpha(\hat{v}_1^I, v_1) = 0\right) = 1$$

when the magnitudes of ordered coefficients of v_1 have rapid decay, i.e., the r -th largest magnitude of v_1 is no greater than $C r^{-1/q}$, $r = 1, 2, \dots$, for some $0 < q < 2$ and $0 < C < \infty$. This marginal variance selection method fails when the variables have equal or almost equal variance. Nevertheless, Johnstone & Lu (2009) proved the first theoretical justification for sparse PCA. Shen et al. (2013) established the consistency of the sparse PCA in the high dimension and low sample

size context. Amini & Wainwright (2009) studied the support recovery property of the semidefinite programming approach of d'Aspremont et al. (2007) under the k -sparse assumption for the leading eigenvector in the rank-1 spiked covariance model. Ma (2013b) proved the consistency of the iterative thresholding approach under a spiked covariance model. Lei & Vu (2015) provided the general sufficient conditions for sparsistency for the Fantope projection and selection method. In a very recent paper, Jankova & van de Geer (2018) proposed a de-biased sparse PCA estimator and studied the asymptotic inference of the sparse eigenvectors.

C. Minimax rates of convergence

The minimax rate of estimation is another important theoretical development for the sparse PCA. The seminal paper by Birnbaum et al. (2013) studied the minimax rates of convergence and adaptive estimation when the rank is a fixed number and the ordered coefficients of each principal eigenvector have rapid decay. Specifically, Birnbaum et al. (2013) established a lower bound on the minimax risk of estimators under various models of sparsity for the population eigenvectors. Ma (2013b) showed that the iterative thresholding estimator attains the minimax rate of convergence over a certain Gaussian class of distributions when the rank is treated as a fixed constant. By allowing the rank increase with the sample size, Cai et al. (2013) and Vu & Lei (2013) studied the minimax optimality and adaptive estimation of the principal subspace for the sparse PCA in the high-dimensional setting. Following Cai et al. (2013), we assume that the $n \times p$ data matrix \mathbf{X} is generated as follows:

$$\mathbf{X} = \mathbf{U} \mathbf{D} \mathbf{V}^T + \mathbf{Z}$$

where \mathbf{U} is the $n \times k$ random effects matrix with i.i.d. $N(0, 1)$ entries, $\mathbf{D} = \text{diag}(\lambda_1^{1/2}, \dots, \lambda_k^{1/2})$ is a diagonal matrix with ordered eigenvalues $\lambda_1 \geq \dots \geq \lambda_k > 0$, \mathbf{V} is an orthonormal matrix, \mathbf{Z} is a random matrix with i.i.d. $N(0, \sigma^2)$ entries, and \mathbf{U} and \mathbf{Z} are independent. Denote by Σ the covariance matrix of \mathbf{X} . Note that $\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T + \sigma^2 \mathbf{I}_p$ and also that the estimation of $\text{span}(\mathbf{V})$ is equivalent to the estimation of $\mathbf{V} \mathbf{V}^T$. Now, we consider the optimal estimation of the principal subspace $\text{span}(\mathbf{V})$ under the commonly used loss function $L(\mathbf{V}, \hat{\mathbf{V}}) = \|\mathbf{V} \mathbf{V}^T - \hat{\mathbf{V}} \hat{\mathbf{V}}^T\|_F^2$ and the following parameter space for Σ :

$$\begin{aligned} & \Theta(s, p, k, \lambda) \\ & = \{\Sigma = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T + \sigma^2 \mathbf{I}_p : \kappa \lambda \geq \lambda_1 \geq \dots \geq \lambda_k \geq \lambda > 0, \mathbf{V}^T \mathbf{V} = \mathbf{I}_k, \|\end{aligned}$$

where $\kappa > 1$, $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_k)$, and $\|\mathbf{V}\|_w = \max_{j=1, \dots, p} \|\mathbf{V}_{(j)*}\|_0$ is the weak ℓ_0 radius of \mathbf{V} . Note that the union of the column supports of \mathbf{V} is of size at most s . Cai et al. (2013) used the local metric entropy (LeCam 1973, Yang & Barron 1999) to construct the lower bound, and then obtain the minimax risk bound in the high-dimensional setting as follows:

$$\inf_{\hat{\mathbf{V}}} \sup_{\Sigma \in \Theta(s, p, k, \lambda)} E[L(\mathbf{V}, \hat{\mathbf{V}})] \asymp \left[\frac{\lambda/\sigma^2 + 1}{n(\lambda/\sigma^2)^2} \left(k(s - k) + s \log \frac{ep}{s} \right) \right] \wedge 1.$$

Cai et al. (2015) studied the minimax rates under the spectral norm, which is directly related to estimating the rank of the factor model.

D. Statistical and Computational Trade-off

It is important to point out that there is a fundamental trade-off between statistical and computational performance. In general, there are no known computationally efficient methods to obtain the minimax rate optimal estimators for the sparse PCA. Several seminal papers highlight the trade-off between computational and statistical efficiency for the sparse PCA, including Amini & Wainwright (2009), Berthet & Rigollet (2013), Krauthgamer et al. (2015), Wang et al. (2016), Gao et al. (2017) and others. Amini & Wainwright (2009) proved that no algorithm can reliably recover the sparse eigenvector under the single-spike covariance model when $k \geq Cn/\log p$ for some positive constant C and all sufficiently large n . Krauthgamer et al. (2015) further proved that the semidefinite programming approach (d'Aspremont et al. 2007) do not close the gap between computational and statistical efficiency as long as $k \geq C\sqrt{n}$ for some positive constant C and all the sufficiently large n . Berthet & Rigollet (2013) considered the optimal detection of sparse principal components in high dimension:

$$H_0 : \mathbf{x} \sim N(0, \mathbf{I}_p) \quad \text{versus} \quad H_1 : \mathbf{x} \sim N(0, \mathbf{I}_p + \theta \mathbf{v}_1 \mathbf{v}_1')$$

where \mathbf{v}_1 has a fixed number of nonzero components. To this end, Berthet & Rigollet (2013) studied a minimax optimal test based on the k -sparse largest eigenvalue of the empirical covariance matrix. The computation of this sparse eigenvalue statistic depends on a well-known decision problem associated to finding whether a graph contains a clique of size k , whose computational complexity is proved to be NP-complete in general (Karp 1972). In the follow-up paper, under the hardness assumption of the planted clique problem (Feldman et al. 2017), Wang et al. (2016) showed that there is an effective sample size regime in which no randomized polynomial time algorithm can achieve the minimax optimal rate for new and larger classes satisfying a restricted covariance concentration condition. Recently, Gao et al. (2017) obtained the first computational lower bounds for sparse PCA under the Gaussian single spiked covariance model and closed the gap in sparse PCA computational lower bounds left by Berthet & Rigollet (2013) and Wang et al. (2016).

IV. APPLICATIONS

Sparse PCA can be used in applications where PCA is normally used. For example, the use of sparse PCA in clustering can lead to sparse clustering algorithms (Chen et al. 2013). PCA is a part of the integrated omic-data analysis, where sparse PCA can be used to replace the regular PCA (Ritchie et al. 2015, Zang et al. 2016). We discuss a few recent applications of SPCA in medical imaging, ecology and neuroscience respectively.

Shape/image analysis Sjöstrand et al. (2007) applied SPCA to landmark-based shape analysis of the CC brain structure.

The authors extracted 5, 20, and 50 nonzero principal components out of the total 156 components corresponding to landmarks, and they also applied the standard PCA as a benchmark. In the subsequent analysis, they used the univariate regression to study the relationship between the resulting deformations based on extracted variables and four clinical outcome variables (gender, age, walking speed, and verbal fluency). Their findings confirmed the male/female mean shape differences and identified the deformation of the CC corresponding to the measure of walking speed. The results for verbal fluency were also meaningful anatomically. Sjöstrand et al. (2007) found that SPCA is useful to derive localized and interpretable patterns of variability while PCA did not provide much interpretational value.

Ecological study Motivated by generating meaningful combinations of the explanatory variables, Gravuer et al. (2008) applied SPCA to perform the dimension reduction before fitting the ABT model. The sparsity helps the interpretability of their model. Specifically, Gravuer et al. (2008) used SPCA to study a range of human, biogeographic, and biological influences on the invasion of *Trifolium* species into New Zealand. The sparse principal components were obtained from 29 categorical and continuous variables for three invasion stages (i.e., introduction, naturalization, and spread), and studied the relationship of sparse principal components to invasion success by using aggregated boosted trees. Specifically, the authors identified 8 sparse principal components on 22 variables for intentional introduction and unintentional introduction–naturalization stages, 7 sparse principal components on 25 variables for naturalization of intentionally introduced species and 7 sparse principal components on 28 variables for relative spread rate. Gravuer et al. (2008) found that SPCA simultaneously improve interpretability and maintain high explained variance.

Neuroscience study SPCA was used in Baden et al. (2016) to study the light-driven Ca^{2+} signals of the GCL cells given a set of standardized visual stimuli in a probabilistic clustering framework. Baden et al. (2016) first used SPCA to extract features that are localized in time and readily interpretable from the responses to the chirp, color, and moving bar stimulus, and then used a Gaussian mixture model on the extracted feature set for clustering. The authors extracted 20 features from the mean response to the chirp, 6 features from the mean response to the color stimulus, 8 features from the response time course and 4 features from its temporal derivative. Many classically used temporal response features were identified, including ON- and OFF responses with different kinetics or selectivity to different temporal frequencies. They also tried the standard PCA and found the results lead to inferior cluster quality.

SPCA is implemented in the R package `elasticnet` available from CRAN:

<http://cran.r-project.org/>.

The Matlab implementation of SPCA is available in the toolbox `SpaSM` from

<http://www2.imm.dtu.dk/projects/spasm/>

V. CONCLUDING REMARKS

In our discussion, we have only presented the use of ℓ_1 norm for sparsity, but there are other equally suitable penalty functions to be used in the sparse PCA methods, including SCAD (Fan & Li 2001, Fan et al. 2014) or ℓ_0 , among others.

We now have a good understanding of the role of sparsity in PCA and ways to effectively exploit the sparsity. There are still remaining issues. A very important question to be investigated further is *Automated sparse PCA?* By “automated” we mean that there is a principled but not overly complicated procedure to set these sparse parameters in sparse PCA. This question is particularly challenging when we solve several sparse principal components jointly. We would also like to have more empirical results to help us understand the pros and cons of each proposed sparse PCA technique, which may also inspire new and better approaches.

ACKNOWLEDGMENTS

We thank Professors Thierry Bouwmans, Yuejie Chi and Namrata Vaswani for inviting us to contribute this review paper to the special issue. We are grateful for the very helpful comments from three anonymous reviewers. Zou’s research is supported in part by the NSF grant DMS-1505111. Xue is supported by the National Science Foundation grant DMS-1505256.

REFERENCES

- Alter, O., Brown, P. & Botstein, D. (2000), ‘Singular value decomposition for genome-wide expression data processing and modeling’, *Proceedings of the National Academy of Sciences* **97**, 10101–10106.
- Amini, A. A. & Wainwright, M. J. (2009), ‘High-dimensional analysis of semidefinite relaxations for sparse principal components’, *The Annals of Statistics* **37**(5B), 2877–2921.
- Anderson, T. W. (2003), *An Introduction to Multivariate Statistical Analysis*, Wiley New York.
- Baden, T., Berens, P., Franke, K., Roseón, M., Bethge, M. & Euler, T. (2016), ‘The functional diversity of retinal ganglion cells in the mouse’, *Nature* **529**(7586), 345–350.
- Baik, J. & Silverstein, J. W. (2006), ‘Eigenvalues of large sample covariance matrices of spiked population models’, *Journal of Multivariate Analysis* **97**(6), 1382–1408.
- Berthet, Q. & Rigollet, P. (2013), ‘Optimal detection of sparse principal components in high dimension’, *The Annals of Statistics* **41**(4), 1780–1815.
- Biehl, M. & Mietzner, A. (1994), ‘Statistical mechanics of unsupervised structure recognition’, *Journal of Physics A* **27**(6), 1885–1897.
- Birnbaum, A., Johnstone, I. M., Nadler, B. & Paul, D. (2013), ‘Minimax bounds for sparse PCA with noisy high-dimensional data’, *The Annals of statistics* **41**(3), 1055–1084.
- Cadima, J. & Jolliffe, I. (1995), ‘Loadings and correlations in the interpretation of principal components’, *Journal of Applied Statistics* **22**, 203–214.
- Cai, T. T., Ma, Z. & Wu, Y. (2013), ‘Sparse PCA: Optimal rates and adaptive estimation’, *The Annals of Statistics* **41**(6), 3074–3110.
- Cai, T. T., Ma, Z. & Wu, Y. (2015), ‘Optimal estimation and rank detection for sparse spiked covariance matrices’, *Probability Theory and Related Fields* **161**, 781–815.
- Chen, G., Sullivan, P. F. & Kosorok, M. R. (2013), ‘Bi-clustering with heterogeneous variance’, *Proceedings of the National Academy of Sciences* **110**(30), 12253–12258.
- d’Aspremont, A. (2011), ‘Identifying small mean-reverting portfolios’, *Quantitative Finance* **11**(3), 351–364.
- d’Aspremont, A., Bach, F. & Ghaoui, L. E. (2008), ‘Optimal solutions for sparse principal component analysis’, *Journal of Machine Learning Research* **9**(Jul), 1269–1294.
- d’Aspremont, A., Ghaoui, L. E., Jordan, M. I. & Lanckriet, G. R. (2007), ‘A direct formulation for sparse PCA using semidefinite programming’, *SIAM Review* pp. 434–448.
- Eckart, C. & Young, G. (1936), ‘The approximation of one matrix by another of lower rank’, *Psychometrika* **1**, 211–218.
- Efron, B., Hastie, T., Johnstone, I. & Tibshirani, R. (2004), ‘Least angle regression’, *The Annals of Statistics* **32**, 407–499.
- Fan, J. & Li, R. (2001), ‘Variable selection via nonconcave penalized likelihood and its oracle properties’, *Journal of the American Statistical Association* **96**, 1348–1360.
- Fan, J., Xue, L. & Zou, H. (2014), ‘Strong oracle optimality of folded concave penalized estimation’, *The Annals of Statistics* **42**(3), 819–849.
- Feldman, V., Grigorescu, E., Reyzin, L., Vempala, S. & Xiao, Y. (2014), ‘Statistical algorithms and a lower bound for detecting planted cliques’, *Journal of the ACM* **64**(2), 8.
- Friedman, J., Hastie, T., Hoefling, H. & Tibshirani, R. (2007), ‘Pathwise coordinate optimization’, *The Annals of Applied Statistics* **1**(2), 302–332.
- Gao, C., Ma, Z. & Zhou, H. H. (2017), ‘Sparse CCA: Adaptive estimation and computational barriers’, *The Annals of Statistics* **45**(5), 2074–2101.
- Gravuer, K., Sullivan, J. J., Williams, P. A. & Duncan, R. P. (2008), ‘Strong human association with plant invasion success for trifolium introductions to new zealand’, *Proceedings of the National Academy of Sciences* **105**(17), 6344–6349.
- Hancock, P., Burton, A. & Bruce, V. (1996), ‘Face processing: human perception and principal components analysis’, *Memory and Cognition* **24**, 26–40.
- Hastie, T., Tibshirani, R., Eisen, M., Brown, P., Ross, D., Scherf, U., Weinstein, J., Alizadeh, A., Staudt, L. & Botstein, D. (2000), ‘gene shaving’ as a method for identifying distinct sets of genes with similar expression patterns’, *Genome Biology* **1**, 1–21.
- Hastie, T., Tibshirani, R. & Friedman, J. (2009), *The Elements of Statistical Learning: Data mining, Inference and Prediction, 2nd Edition*, Springer Verlag, New York.
- Jankova, J. & van de Geer, S. (2018), ‘De-biased sparse PCA: Inference and testing for eigenstructure of large covariance matrices’, *arXiv preprint arXiv:1801.10567*
- Jeffers, J. (1967), ‘Two case studies in the application of

- principal component', *Applied Statistics* **16**, 225–236.
- Johnstone, I. M. & Lu, A. Y. (2009), 'On consistency and sparsity for principal components analysis in high dimensions', *Journal of the American Statistical Association* **104**(486), 682–693.
- Jolliffe, I. (1995), 'Rotation of principal components: choice of normalization constraints', *Journal of Applied Statistics* **22**, 29–35.
- Jolliffe, I. T., Trendafilov, N. T. & Uddin, M. (2003), 'A modified principal component technique based on the lasso', *Journal of Computational and Graphical Statistics* **12**(3), 531–547.
- Journée, M., Nesterov, Y., Richtárik, P. & Sepulchre, R. (2010), 'Generalized power method for sparse principal component analysis', *Journal of Machine Learning Research* **11**(Feb), 517–553.
- Jung, S. & Marron, J. S. (2009), 'PCA consistency in high dimension, low sample size context', *The Annals of Statistics* **37**(6B), 4104–4130.
- Karp, R. M. (1972), 'Reducibility among combinatorial problems'. In *Complexity of Computer Computations (Proc. Sympos., IBM Thomas J. Watson Res. Center, Yorktown Heights, N.Y., 1972)*, 85–103.
- Krauthgamer, R., Nadler, B. & Vilenchik, D. (2015), 'Do semidefinite relaxations solve sparse PCA up to the information limit?' *The Annals of Statistics* **43**(3), 1300–1322.
- LeCam, L. (1973), 'Convergence of estimates under dimensionality restrictions', *The Annals of Statistics* **1**(1), 38–53.
- Lei, J. & Vu, V. (2015), 'Sparsistency and agnostic inference in sparse PCA', *The Annals of Statistics* **43**(1), 299–322.
- Lu, Z. & Zhang, Y. (2012), 'An augmented lagrangian approach for sparse principal component analysis', *Mathematical Programming* **135**(1-2), 149.
- Ma, S. (2013a), 'Alternating direction method of multipliers for sparse principal component analysis', *Journal of the Operations Research Society of China* **1**(2), 253–274.
- Ma, Z. (2013b), 'Sparse principal component analysis and iterative thresholding', *The Annals of Statistics* **41**(2), 772–801.
- Mardia, K., Kent, J. & Bibby, J. (1979), *Multivariate Analysis*, Academic Press.
- Moghaddam, B., Weiss, Y. & Avidan, S. (2006), Spectral bounds for sparse pca: Exact and greedy algorithms, in 'Advances in Neural Information Processing Systems', pp. 915–922.
- Nadler, B. (2008), 'Finite sample approximation results for principal component analysis: A matrix perturbation approach', *The Annals of Statistics* **36**(6), 2791–2817.
- Paul, D. (2007), 'Asymptotics of sample eigenstructure for a large dimensional spiked covariance model', *Statistica Sinica* **17**, 1617–1642.
- Pearson, K. (1901), 'Liii. 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**(11), 559–572.
- Ritchie, M. D., Holzinger, E. R., Li, R., Pendergrass, S. A. & Kim, D. (2015), 'Methods of integrating data to uncover genotype-phenotype interactions', *Nature Reviews. Genetics* **16**(2), 85–97.
- Shen, D., Shen, H. & Marron, J. S. (2013), 'Consistency of sparse PCA in high dimension, low sample size contexts', *Journal of Multivariate Analysis* **115**, 317–333.
- Shen, H. & Huang, J. Z. (2008), 'Sparse principal component analysis via regularized low rank matrix approximation', *Journal of Multivariate Analysis* **6**(99), 1015–1034.
- Sjöstrand, K., Rostrup, E., Ryberg, C., Larsen, R., Studholme, C., Baezner, H., Ferro, J., Fazekas, F., Pantoni, L., Inzitari, D. & Waldemar, G. (2007), 'Sparse decomposition and modeling of anatomical shape variation', *IEEE Transactions on Medical Imaging* **26**(12), 1625–1635.
- Tibshirani, R. (1996), 'Regression shrinkage and selection via the lasso', *Journal of the Royal Statistical Society, Series B* **58**, 267–288.
- Vines, S. (2000), 'Simple principal components', *Applied Statistics* **49**, 441–451.
- Vu, V. & Lei, J. (2013), 'Minimax sparse principal subspace estimation in high dimensions', *The Annals of Statistics* **41**(6), 2905–2947.
- Vu, V. Q., Cho, J., Lei, J. & Rohe, K. (2013), Fantope projection and selection: a near-optimal convex relaxation of sparse PCA, in 'Proceedings of the 26th International Conference on Neural Information Processing Systems', Curran Associates Inc., pp. 2670–2678.
- Wang, T., Berthet, Q. & Samworth, R. J. (2016), 'Statistical and computational trade-offs in estimation of sparse principal components', *The Annals of Statistics* **44**(5), 1896–1930.
- Watkin, T. & Nadal, J.-P. (1994), 'Optimal unsupervised learning', *Journal of Physics A* **27**(6), 1899–1915.
- Witten, D. M., Tibshirani, R. & Hastie, T. (2009), 'A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis', *Biostatistics* **10**(3), 515–534.
- Yang, Y. & Barron, A. (1999), 'Information-theoretic determination of minimax rates of convergence', *The Annals of Statistics* **27**(5), 1564–1599.
- Yuan, X. & Zhang, T. (2013), 'Truncated power method for sparse eigenvalue problems', *Journal of Machine Learning Research* **14**, 899–925.
- Zang, C., Wang, T., Deng, K., Li, B., Qin, Q., Xiao, T., Zhang, S., Meyer, C. A., He, H. H., Brown, M., Liu, J. S., Xie, Y. & Liu, X. S. (2016), 'High-dimensional genomic data bias correction and data integration using mancic', *Nature Communications* **7**, 11305.
- Zou, H. & Hastie, T. (2005), 'Regularization and variable selection via the elastic net', *Journal of the Royal Statistical Society: Series B* **67**(2), 301–320.
- Zou, H., Hastie, T. & Tibshirani, R. (2006), 'Sparse principal component analysis', *Journal of Computational and Graphical Statistics* **15**(2), 265–286.