

ORTHOGONAL SPARSE EIGENVECTORS: A PROCRUSTES PROBLEM

Konstantinos Benidis, Ying Sun, Prabhu Babu, and Daniel P. Palomar

Department of Electronic and Computer Engineering, Hong Kong University of Science and Technology, Hong Kong

ABSTRACT

The problem of estimating sparse eigenvectors of a symmetric matrix attracts a lot of attention in many applications, especially those with high dimensional data set. While classical eigenvectors can be obtained as the solution of a maximization problem, existing approaches formulated this problem by adding a penalty term into the objective function that encourages a sparse solution. Nevertheless, the resulting methods achieve sparsity at a sacrifice of the orthogonality property. In this paper, we develop a new method to estimate dominant sparse eigenvectors without trading off their orthogonality. The problem is highly non-convex and too hard to handle. We apply the minorization-maximization (MM) framework where we iteratively maximize a tight lower bound (surrogate function) of the objective function over the Stiefel manifold. The inner maximization problem turns out to be the rectangular Procrustes problem, which has a closed-form solution. Numerical experiments show that the proposed method matches or outperforms existing algorithms in terms of recovery probability and explained variance.

Index Terms— Sparse PCA, Rectangular Procrustes, Minorization - Maximization (MM).

1. INTRODUCTION

The extraction of the principal components of a matrix is a well studied problem. The main tool, Principal Component Analysis (PCA) [1], essentially finds the directions of maximum variance, so that we can achieve dimensionality reduction with minimum information loss. We can find these directions either by performing singular value decomposition (SVD) in a data matrix A , or via eigenvalue decomposition (EVD) in its corresponding covariance matrix Σ .

Nevertheless, the success of PCA is not only due to the capture of the directions of maximum variability of the principal components. Further, these directions are orthogonal to each other, i.e., they form an orthonormal basis. Finally, the PCs are uncorrelated which aids further statistical analysis.

In general, the resulting eigenvectors are dense vectors. Even if the underlying covariance matrix from which the samples are generated indeed has sparse eigenvectors, we do not

expect to get a sparse result due to estimation error. Further, in many applications, the principal components have an actual physical meaning (e.g. gene expression). Thus, a sparse eigenvector could help significantly the interpretability of the result.

Many different techniques have been proposed in this direction during the last two decades. One of the first approaches was to simply set to zero all the elements that their absolute value is smaller than a threshold [2]. In [3], the authors propose the SCoTLASS algorithm which maximizes the Rayleigh quotient of the covariance matrix, while sparsity is enforced with the Lasso penalty. Many recent approaches are based on reformulations or convex relaxations. For example in [4] Zou et al. formulate the sparse PCA problem as a ridge regression problem and they impose sparsity again using the Lasso penalty. In [5], the authors form a semidefinite program (SDP) after a convex relaxation of the sparse PCA problem, leading to the DSPCA algorithm. Low rank approximation of the data matrix is considered in [6], under sparsity penalties, while in [7], Journeé et al. reformulated the problem as an alternating optimization problem, resulting in the GPower algorithm. In [8], the support of the principal components is determined based on the desired cardinality. The entries that do not belong to the support are set to zero. Finally, in [9], the sparse generalized eigenvalue problem is considered only for the first principal component, where the MM framework is used.

In all the aforementioned algorithms the orthogonality property of the eigenvectors is sacrificed for sparse solutions. The advantages of an orthogonal basis are well known. For instance, an orthonormal basis can be extremely useful since it can reduce the potential computational cost of any post-processing phase; this may not seem much for vector spaces of small dimension but it is invaluable for high dimensional vector spaces or function spaces. Consider for example the solution of a linear system via Gaussian elimination. It requires $O(m^3)$ operations for a non-orthogonal basis, compared to $O(m)$ operations if the basis is orthogonal, where m is the dimension. This, among other optimal properties, motivates us to find sparse loading vectors that maintain their orthogonality.

Notation: \mathbf{R} denotes the real field and \mathbf{R}^m (\mathbf{R}_+^m) the set of (non-negative) real vectors of size m . Vectors are denoted by bold lower case letters and matrices by bold capital letters i.e.,

This work was supported by the Hong Kong RGC 16206315 research grant.

\mathbf{x} and \mathbf{X} , respectively. The i -th entry of a vector is denoted by x_i , the i -th column of matrix \mathbf{X} by \mathbf{x}_i , and the $(i\text{-th}, j\text{-th})$ element of a matrix by x_{ij} . A size m vector of ones is denoted by $\mathbf{1}_m$. $\text{vec}(\cdot)$ denotes the vectorized form of a matrix. The superscript $(\cdot)^T$ denotes the transpose of a matrix. $\text{diag}(\mathbf{x})$ is a diagonal matrix formed with \mathbf{x} at its principal diagonal. $[\mathbf{x}]_{m \times n}$ is an $m \times n$ matrix such that $\text{vec}([\mathbf{x}]_{m \times n}) = \mathbf{x}$. $\|\mathbf{x}\|_0$ denotes the number of nonzero elements of \mathbf{x} . $\mathbf{X} \succcurlyeq 0$ means that the matrix \mathbf{X} is positive semidefinite.

2. PROBLEM STATEMENT

Given a data matrix $\mathbf{A} \in \mathbf{R}^{n \times m}$, encoding n samples of m variables, we can extract the leading eigenvector of the scaled sample covariance matrix $\mathbf{S} = \mathbf{A}^T \mathbf{A}$ by solving the following optimization problem:

$$\begin{aligned} & \underset{\mathbf{u}}{\text{maximize}} && \mathbf{u}^T \mathbf{S} \mathbf{u} \\ & \text{subject to} && \mathbf{u}^T \mathbf{u} = 1. \end{aligned} \quad (1)$$

In order to get a sparse result, we include a regularization term in the objective that imposes sparsity, i.e.,

$$\begin{aligned} & \underset{\mathbf{u}}{\text{maximize}} && \mathbf{u}^T \mathbf{S} \mathbf{u} - \rho \|\mathbf{u}\|_0 \\ & \text{subject to} && \mathbf{u}^T \mathbf{u} = 1, \end{aligned} \quad (2)$$

where ρ is a regularization parameter.

Problem (2) can be generalized to extract multiple eigenvectors as follows:

$$\begin{aligned} & \underset{\mathbf{U}}{\text{maximize}} && \text{Tr}(\mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{D}) - \sum_{i=1}^q \rho_i \|\mathbf{u}_i\|_0 \\ & \text{subject to} && \mathbf{U}^T \mathbf{U} = \mathbf{I}_q. \end{aligned} \quad (3)$$

Here, q is the number of eigenvectors we wish to estimate, $\mathbf{U} \in \mathbf{R}^{m \times q}$ and $\mathbf{D} \succcurlyeq 0$ is a diagonal matrix.

The case presented in [9], is a special case of the above optimization problem, with $q = 1$. Nevertheless, it is not possible to follow the same procedure as in [9] to solve the problem due to the orthogonality constraint. Instead, we tackle this problem using the MM algorithm, which results in solving a sequence of Procrustes problem that has a closed-form solution given by SVD.

3. MAXIMIZATION USING THE MM FRAMEWORK

The MM algorithm is a way to handle optimization problems that are too difficult to face directly [10]. Consider a general optimization problem

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} && f(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \mathcal{X}, \end{aligned}$$

where \mathcal{X} is a closed set. At a given point $\mathbf{x}^{(k)}$, the MM algorithm finds a surrogate function $g(\mathbf{x}|\mathbf{x}^{(k)})$ of $f(\mathbf{x})$ satisfying the following properties: $f(\mathbf{x}^{(k)}) = g(\mathbf{x}^{(k)}|\mathbf{x}^{(k)})$ and $f(\mathbf{x}) \geq g(\mathbf{x}|\mathbf{x}^{(k)})$, $\forall \mathbf{x} \in \mathcal{X}$. Then \mathbf{x} is iteratively updated (with k denoting iterations) as:

$$\mathbf{x}^{(k+1)} = \arg \max_{\mathbf{x} \in \mathcal{X}} g(\mathbf{x}|\mathbf{x}^{(k)}). \quad (4)$$

It can be seen easily that $f(\mathbf{x}^{(k)}) \leq f(\mathbf{x}^{(k+1)})$ holds.

Return to the sparse eigenvalue problem (3). First, we approximate the ℓ_0 -norm by a differentiable function $g_p^\epsilon(\cdot)$ as [9], which leads to the following approximate problem:

$$\begin{aligned} & \underset{\mathbf{U}}{\text{maximize}} && \text{Tr}(\mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{D}) - \sum_{i=1}^q \rho_i \sum_{j=1}^m g_p^\epsilon(u_{ij}) \\ & \text{subject to} && \mathbf{U}^T \mathbf{U} = \mathbf{I}_q, \end{aligned} \quad (5)$$

where

$$g_p^\epsilon(u_{ij}) = \begin{cases} \frac{u_{ij}^2}{2\epsilon(p+\epsilon)\log(1+1/p)}, & |u_{ij}| \leq \epsilon \\ \frac{\log(\frac{p+|u_{ij}|}{p+\epsilon}) + \frac{\epsilon}{2(p+\epsilon)}}{\log(1+1/p)}, & |u_{ij}| > \epsilon, \end{cases} \quad (6)$$

and $0 < p \leq 1$ and $0 < \epsilon \ll 1$. Problem (5) has no closed-form solution. In the following, we apply the MM algorithm and derive a tight lower bound (surrogate function), $g(\mathbf{U}|\mathbf{U}^{(k)})$, for the objective function of (5), denoted by $f(\mathbf{U})$, at the $(k+1)$ -th iteration.

Proposition 1. *The function $f(\mathbf{U})$ is lowerbounded by the surrogate function*

$$g(\mathbf{U}|\mathbf{U}^{(k)}) = 2\text{Tr}\left(\left(\mathbf{G}^{(k)} - \mathbf{H}^{(k)}\right)^T \mathbf{U}\right) + c_1 - c_2, \quad (7)$$

where

$$\mathbf{G}^{(k)} = \mathbf{S} \mathbf{U}^{(k)} \mathbf{D}, \quad (8)$$

$$\mathbf{H}^{(k)} = \left[\text{diag}\left(\mathbf{w}^{(k)} - \mathbf{w}_{\max}^{(k)} \otimes \mathbf{1}_m\right) \text{vec}\left(\mathbf{U}^{(k)}\right) \right]_{m \times q}, \quad (9)$$

and c_1, c_2 are optimization irrelevant constants. The weights $\mathbf{w}^{(k)} \in \mathbf{R}_+^{mq}$ are given by

$$w_i^{(k)} = \begin{cases} \frac{\rho_i}{2\epsilon(p+\epsilon)\log(1+1/p)}, & |u_i^{(k)}| \leq \epsilon, \\ \frac{\rho_i}{2\log(1+1/p)|u_i^{(k)}|(|u_i^{(k)}|+p)}, & |u_i^{(k)}| > \epsilon, \end{cases} \quad (10)$$

where $\mathbf{u}^{(k)} = \text{vec}(\mathbf{U}^{(k)})$, and $\mathbf{w}_{\max}^{(k)} \in \mathbf{R}_+^q$, with $w_{\max,i}^{(k)}$ being the maximum weight that corresponds to the i -th eigenvector. Equality is achieved when $\mathbf{U} = \mathbf{U}^{(k)}$.

Proof. The first term of the objective is convex so a lower bound can be constructed by its first order Taylor expansion:

$$\text{Tr}(\mathbf{U}^T \mathbf{S} \mathbf{U} \mathbf{D}) \geq 2\text{Tr}\left(\left(\mathbf{S} \mathbf{U}^{(k)} \mathbf{D}\right)^T \mathbf{U}\right) + c_1, \quad (11)$$

where $c_1 = -\text{Tr} \left(\mathbf{U}^{(k)T} \mathbf{S} \mathbf{U}^{(k)} \mathbf{D} \right)$ is a constant.

Following the same approach as [9], we can bound the function $\sum_{i=1}^q \rho_i \sum_{j=1}^m g_p^\epsilon(u_{ij})$ with a weighted quadratic one. Based on the results of [9] and by incorporating the sparsity parameters ρ_i to the corresponding weights, it holds that

$$\sum_{i=1}^q \rho_i \sum_{j=1}^m g_p^\epsilon(u_{ij}) \leq \mathbf{u}^T \text{diag} \left(\mathbf{w}^{(k)} \right) \mathbf{u},$$

with $\mathbf{u} = \text{vec}(\mathbf{U})$, and the weights $\mathbf{w}^{(k)} \in \mathbf{R}_+^{mq}$ given by (10). This bound cannot lead to a closed form solution thus we apply the MM framework one more time. The idea is to create a concave term and linearize it since the linear approximation of a concave function is an upper bound of the function. It is easy to show that the following holds:

$$\mathbf{u}^T \text{diag} \left(\mathbf{w}^{(k)} \right) \mathbf{u} \leq 2 \text{Tr} \left(\mathbf{H}^{(k)T} \mathbf{U} \right) + c_2,$$

where $\mathbf{H}^{(k)}$ is given by (9), and $c_2 = 2 \left(\mathbf{1}_m^T \mathbf{w}_{\max}^{(k)} \right) - \mathbf{u}^{(k)T} \text{diag} \left(\mathbf{w}^{(k)} \right) \mathbf{u}^{(k)}$ is a constant. \square

Now, the optimization problem of every MM iteration takes the following form:

$$\begin{aligned} & \underset{\mathbf{U}}{\text{maximize}} \quad \text{Tr} \left(\left(\mathbf{G}^{(k)} - \mathbf{H}^{(k)} \right)^T \mathbf{U} \right) \\ & \text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}_q. \end{aligned} \quad (12)$$

Proposition 2. *The optimal solution of (12) is $\mathbf{U}^* = \mathbf{V}_L \mathbf{V}_R^T$, where $\mathbf{V}_L \in \mathbf{R}^{m \times q}$ and $\mathbf{V}_R \in \mathbf{R}^{q \times q}$ are the left and right singular vectors of the matrix $\left(\mathbf{G}^{(k)} - \mathbf{H}^{(k)} \right)$, respectively.*

Proof. Notice that problem (12) is equivalent to

$$\begin{aligned} & \underset{\mathbf{U}}{\text{minimize}} \quad \|\mathbf{U} - \mathbf{B}^{(k)}\|_F^2 \\ & \text{subject to} \quad \mathbf{U}^T \mathbf{U} = \mathbf{I}_q. \end{aligned} \quad (13)$$

with $\mathbf{B}^{(k)} = \mathbf{G}^{(k)} - \mathbf{H}^{(k)}$. Problem (13) is a rectangular Procrustes problem and its optimal solution is $\mathbf{U}^* = \mathbf{V}_L \mathbf{V}_R^T$ where $\mathbf{V}_L, \mathbf{V}_R$ are the left and right singular vectors of the matrix $\mathbf{B}^{(k)}$ [11]. \square

In Algorithm 1 we summarize the above iterative procedure. We will refer to it as IMRP.

4. NUMERICAL EXPERIMENTS

4.1. Random Data Drawn from a Sparse PCA Model

In the first experiment, we compare the performance of the proposed IMRP algorithm with a benchmark algorithm GPower $_{\ell_0}$ proposed in [7].

Algorithm 1 IMRP - Iterative Minimization of Rectangular Procrustes for the sparse eigenvector problem (5)

- 1: Set $k = 0$, choose $\mathbf{U}^{(0)} \in \{\mathbf{U} : \mathbf{U}^T \mathbf{U} = \mathbf{I}_q\}$
 - 2: **repeat**:
 - 3: Compute $\mathbf{G}^{(k)}, \mathbf{H}^{(k)}$ with (8),(9), respectively
 - 4: Compute $\mathbf{V}_L, \mathbf{V}_R$, the left and right singular vectors of $\left(\mathbf{G}^{(k)} - \mathbf{H}^{(k)} \right)$, respectively
 - 5: $\mathbf{U}^{(k+1)} = \mathbf{V}_L \mathbf{V}_R^T$
 - 6: $k \leftarrow k + 1$
 - 7: **until** convergence
 - 8: **return** $\mathbf{U}^{(k)}$
-

To illustrate the sparse recovering performance of our algorithm we generate synthetic data as in [7–9]. To this end, we construct a covariance matrix Σ through the eigenvalue decomposition $\Sigma = \mathbf{V} \text{diag}(\lambda) \mathbf{V}^T$, where the first q columns of $\mathbf{V} \in \mathbf{R}^{m \times m}$ have a pre-specified sparse structure. We consider a setup with $m = 500$, $n = 50$ and $q = 2$. We set the first two orthonormal eigenvectors to be

$$\begin{cases} v_{i1} = \frac{1}{\sqrt{10}} & \text{for } i = 1, \dots, 10, \\ v_{i1} = 0 & \text{otherwise,} \\ v_{i2} = \frac{1}{\sqrt{10}} & \text{for } i = 11, \dots, 20, \\ v_{i2} = 0 & \text{otherwise.} \end{cases}$$

The remaining eigenvectors are generated randomly, satisfying the orthogonality property. We set the eigenvalues to be $\lambda_1 = 400$, $\lambda_2 = 300$ and $\lambda_i = 1$ for $i = 3, \dots, 500$.

We randomly generate 500 data matrices $\mathbf{A} \in \mathbf{R}^{m \times n}$ by drawing n samples from a zero-mean normal distribution with covariance matrix Σ , i.e., $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$, for $i = 1, \dots, n$. Then we employ the two algorithms to compute the two leading sparse eigenvectors $\mathbf{u}_1, \mathbf{u}_2 \in \mathbf{R}^{500}$. We consider a successful recovery when both quantities $|\mathbf{u}_1^T \mathbf{v}_1|$ and $|\mathbf{u}_2^T \mathbf{v}_2|$ are greater than 0.99. For each sparsity parameter ρ_i , we use the normalizations presented in [9]. From Figure 1, we can see that our proposed algorithm IMRP achieves higher chance of exact recovery for a wide range of the parameters ρ_i .

4.2. Gene Expression Data

Now, we test the performance of the two algorithms on the gene expression dataset collected in the breast cancer study by Bild et al. [12]. The dataset contains 158 samples over 12,625 genes. We consider the 4,000 genes with the largest variances and we estimate the first 5 eigenvectors.

In order to avoid overestimation of the variance, we have used the notion of cumulative percentage of explained variance (CPEV), proposed in [6]. It is defined as:

$$\text{CPEV} = \text{Tr} \left(\mathbf{A}_q^T \mathbf{A}_q \right) / \text{Tr} \left(\mathbf{A}^T \mathbf{A} \right). \quad (14)$$

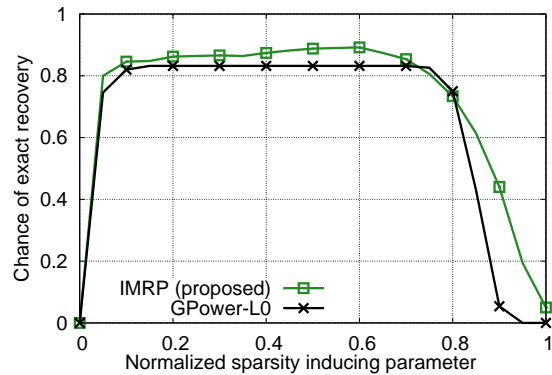


Fig. 1. Chance of exact recovery vs normalized regularization parameter.

where $A_q = AU(U^TU)^{-1}U^T$. Due to the orthogonality constraints, increasing the cardinality does not necessarily mean that the CPEV will increase. To this end, for a fixed cardinality, we depict the maximum variance being explained from the sparse eigenvectors up to this cardinality.

In Figure 2 we illustrate the cumulative percentage of explained variance versus the cardinality for the IMRP and GPower $_{\ell_0}$ algorithms. For maximum cardinality the percentage of explained variance becomes 1 for both algorithms. For fixed cardinality, the two algorithms can explain approximately the same amount of variance. For comparison we have also included the simple thresholding scheme which first computes the regular principal component and then keeps a required number of entries with largest absolute values.

5. CONCLUSION

We have proposed a new algorithm for sparse eigenvalue extraction. The algorithm is derived based on the minorization-majorization method that was applied after a smooth approximation of the ℓ_0 -norm. Unlike all the other existing methods, the resulting sparse eigenvectors from our proposed method maintain their orthogonality property. Numerical results have shown that IMRP matches or outperforms existing algorithms in terms of recovery and explained variance.

6. REFERENCES

- [1] I. T. Jolliffe, *Principal component analysis*. Wiley Online Library, 2002.
- [2] J. Cadima and I. T. Jolliffe, "Loading and correlations in the interpretation of principle compenents," *Journal of Applied Statistics*, vol. 22, no. 2, pp. 203–214, 1995.
- [3] I. T. Jolliffe, N. T. Trendafilov, and M. Uddin, "A modified principal component technique based on the LASSO," *Journal of computational and Graphical Statistics*, vol. 12, no. 3, pp. 531–547, 2003.

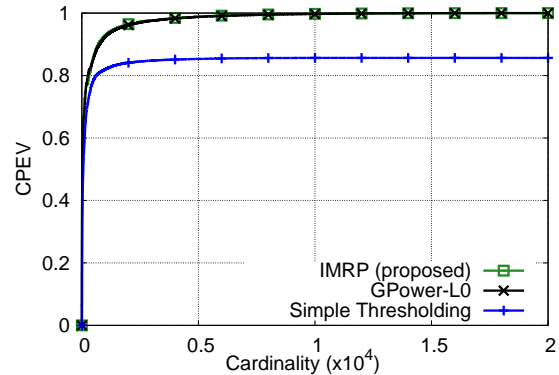


Fig. 2. CPEV vs cardinality.

- [4] H. Zou, T. Hastie, and R. Tibshirani, "Sparse principal component analysis," *Journal of computational and graphical statistics*, vol. 15, no. 2, pp. 265–286, 2006.
- [5] A. d'Aspremont, L. El Ghaoui, M. I. Jordan, and G. R. Lanckriet, "A direct formulation for sparse PCA using semidefinite programming," *SIAM review*, vol. 49, pp. 434–448, July 2007.
- [6] H. Shen and J. Z. Huang, "Sparse principal component analysis via regularized low rank matrix approximation," *Journal of multivariate analysis*, vol. 99, pp. 1015–1034, July 2008.
- [7] M. Journée, Y. Nesterov, P. Richtárik, and R. Sepulchre, "Generalized power method for sparse principal component analysis," *The Journal of Machine Learning Research*, vol. 11, pp. 517–553, Mar. 2010.
- [8] X.-T. Yuan and T. Zhang, "Truncated power method for sparse eigenvalue problems," *The Journal of Machine Learning Research*, vol. 14, no. 1, pp. 899–925, 2013.
- [9] J. Song, P. Babu, and D. P. Palomar, "Sparse generalized eigenvalue problem via smooth optimization," *IEEE Transactions on Signal Processing*, vol. 63, pp. 1627–1642, Apr. 2015.
- [10] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *The American Statistician*, vol. 58, pp. 30–37, Feb. 2004.
- [11] J. H. Manton, "Optimization algorithms exploiting unitary constraints," *IEEE Transactions on Signal Processing*, vol. 50, pp. 635–650, Mar. 2002.
- [12] A. H. Bild, G. Yao, J. T. Chang, Q. Wang, A. Potti, D. Chasse, M.-B. Joshi, D. Harpole, J. M. Lancaster, A. Berchuck, *et al.*, "Oncogenic pathway signatures in human cancers as a guide to targeted therapies," *Nature*, vol. 439, pp. 353–357, Jan. 2006.