

# Sparse Principal Component

## Analysis Formulations And Algorithms

Thomas Rusch    Norbert Walchhofer, Department of  
Finance, Accounting and Statistics

WU Vienna

June 20, 2011

- ▶ Background
  - ▶ Review of Principal Component Analysis (PCA)?
- ▶ Generalized Power Method for Sparse PCA
- ▶ Problem Formulations and Reformulations
  - ▶ Single-unit sparsePCA
    - ▶ Single-unit sparsePCA via  $\ell_1$ -penalty
    - ▶ Single-unit sparsePCA via  $\ell_0$ -penalty
  - ▶ Block sparsePCA
  - ▶ Power Method
- ▶ Proposed Algorithms and their Evaluation
  - ▶ Exemplary Algorithm

- ▶ Method for dimension reduction
- ▶ Orthogonal transformation of possibly correlated variables into uncorrelated principal components
- ▶ Project a centered data matrix  $A$  or a (sample) covariance matrix thereof  $\Sigma = A^T A$  from  $\mathbf{R}^p$  into  $\mathbf{R}^m$  where  $q \leq p$
- ▶ Aims at finding a few linear combinations the  $p$  variables, pointing in orthogonal directions explaining as much variance as possible.

## PCA - Formulation

$$z^* = \max_{z^T z \leq 1} z^T \Sigma z$$

Extracting the first principal component can be done in two ways:

- ▶ computing the first eigenvector of  $\Sigma$
- ▶ or the first right singular value of  $A$ .

Usually principal components are linear combinations of all input variables with loading vector  $z^*$  (score).

PCA aims to reduce complexity, however there are some drawbacks:

- ▶ principal components depend on many variables
- ▶ interpretation of components can be agonizing
- ▶ individual loadings can be negligible

Sparse PCA simplifies mass of loadings and therefore

- ▶ highlights the most essential structures,
- ▶ is easier to interpret,
- ▶ amount of input variables can be controlled for
- ▶ and it provides a reasonable *trade-off* between **explained variance** and **usability**.

Journée et al. (2010) provide following contributions:

- ▶ Formulations of for single-unit sparse PCA via  $\ell_1$  & cardinality ( $\ell_0$ )-penalty
- ▶ Formulations of for block sparse PCA via  $\ell_1$  & cardinality-penalty
- ▶ Reformulations to **convex** optimization problems
- ▶ Application of the Power Method for sparse PCA
- ▶ Development of algorithms to solve the reformulated optimization problems

Single-unit optimization tries to find sparse loadings for one principal component, before calculating the next one.

Consider following optimization problem

$$\Phi_{\ell_1}(\gamma) \stackrel{\text{def}}{=} \max_{z \in B^n} \sqrt{z^T \Sigma z} - \gamma \|z\|_1 \quad (1)$$

with sparsity-controlling parameter  $\gamma \geq 0$  and sample covariance matrix  $\Sigma = A^T A$ .

By setting  $\gamma = 0$  there can be shown that  $\Phi_{\ell_1}(0)$  leads to

$$\gamma < \|a_{i^*}\|_2, \quad (2)$$

defining the upper bound for  $\gamma$  where  $i^*$  is obtained by  $\max_i \|a_i\|_2$

## Reformulating the problem

$$\begin{aligned}\Phi_{\ell_1}(\gamma) &= \max_{z \in B^n} \|Az\|_2 - \gamma \|z\|_1 \\ &= \max_{z \in B^n} \max_{x \in B^n} x^T Az - \gamma \|z\|_1\end{aligned}\quad (3)$$

$$\begin{aligned}&= \max_{x \in B^n} \max_{z \in B^n} \sum_{i=1}^n z_i (a_i^T x) - \gamma \|z\|_1 \\ &= \max_{x \in B^n} \max_{z' \in B^n} \sum_{i=1}^n |z'_i| (|a_i^T x| - \gamma)\end{aligned}\quad (4)$$

where  $z_i = \text{sign}(a_i^T x) z'_i$ .

Equation 2 proves that there is a  $x \in B^n$  for which  $a_i^T x > \gamma$ .

## Further reformulating the problem

In view of 2, there is some  $x \in B^n$  for which  $a_i^T x > \gamma$ . By fixing  $x$ , solving the inner maximization problem for  $z'$  we obtain a closed solution for  $z^*$ :

$$z_i^* = z_i^*(\gamma) = \frac{\text{sign}(a_i^T x) [|a_i^T x| - \gamma]_+}{\sqrt{\sum_{k=1}^n [|a_k^T x| - \gamma]_+^2}}, \quad i = 1, \dots, n. \quad (5)$$

## Adjusting the objective function

Therefore Eq. 4 can be written as

$$\Phi_{\ell_1}^2(\gamma) = \max_{x \in S^p} \sum_{i=1}^n \left[ |a_i^T x| - \gamma \right]_+^2. \quad (6)$$

This results in a differentiable and **convex** objective function, where all local and global maximal must lie in den Euclidian sphere  $S^p$ , **reducing the search space** of our initial problem formulation (see Eq. 8) to dimension  $p$  with  $p \ll n$ !

What really happens...

- ▶ By introducing a vector  $x$  the optimization problem is split in two, solving  $x$  and  $z$ , respectively.
- ▶  $x$  is solved in Eq. 6 providing a sparsity pattern for  $z^*$ .
- ▶ This sparsity pattern indicates which  $z_i$  are active, i.e. are not 0.
- ▶ Therefore loadings only have to be calculated for  $p$  of the  $n$  variables of  $A$  (for one component).

In contrast to the  $\ell_1$ -penalty (soft constraint) the  $\ell_0$  or cardinality-penalty directly penalizes the number of non-zero components of vector  $z$  (hard constraint).

Optimization problem formulated in d'Aspremont et al. (2008)

$$\Phi_{\ell_0}(\gamma) \stackrel{\text{def}}{=} \max_{z \in B^n} \sqrt{z^T \Sigma z} - \gamma \|z\|_0 \quad (7)$$

Analogue to the  $\ell_1$  case with derive the boundary for  $\gamma$ , optimization for  $z_i^*$  and  $x$ :

$$\gamma < \|a_{i^*}\|_2^2,$$

$$z_i^* = z_i^*(\gamma) = \frac{[\text{sign}(a_i^T x)^2 - \gamma]_+ a_i^T x}{\sqrt{[\text{sign}(a_i^T x)^2 - \gamma]_+ (a_i^T x)^2}}, \quad i = 1, \dots, n,$$

$$\Phi_{\ell_1}^2(\gamma) = \max_{x \in S^p} \sum_{i=1}^n [(a_i^T x)^2 - \gamma]_+.$$

Block optimization tries to find sparse loadings for  $m$  principal components.

Consider following generalization of Eq. 3

$$\Phi_{\ell_1, m}(\gamma) \stackrel{\text{def}}{=} \max_{X \in \mathbb{S}_m^p} \max_{Z \in [S^n]^m} \text{Tr}(X^T A Z N) - \sum_{j=1}^m \gamma_j \sum_{i=1}^n |z_{ij}| \quad (8)$$

where  $\gamma = [\gamma_1, \dots, \gamma_m]^T \quad \forall \gamma_j \geq 0$  and  $N = \text{Diag}(\mu_1, \dots, \mu_m) \quad \forall \mu_j > 0$ .

Each  $\gamma_j$  controls the sparsity for the corresponding component. For positive  $\gamma_j$  columns of  $Z$  are not expected to be orthogonal anymore! Note that distinct values of  $\mu_j$  ensure the columns of  $X^*$  being the dominant  $m$  components, while also pursuing more sparse and orthogonal vectors.

Since the columns of  $Z$  are decoupled the reformulation can be done analogue to the single-unit case. Hence, for every column of  $X$  every row element is optimized, indicating the 'active-status' for each component of  $Z$  (i.e. variable of  $A$ ) of each row  $Z$ .

If  $\mu_j |a_i^T x_j^*| > \gamma_j$  is fulfilled  $z_{ij}^*$  is active.

The power method is an eigenvalue algorithm, given a matrix  $A$  trying to find the dominant eigenvalue  $\lambda$  and its corresponding eigenvector  $v$  such that  $Av = \lambda v$ . By avoiding a matrix decomposition it is very favorable for large sparse matrices since the computation is very low. The scalar  $q = x^T x$  converges linearly against the dominant eigenvalue.

$$x_{k+1} = \frac{Ax_k}{\|Ax_k\|} \quad (9)$$

$x_0$  can be an approximation or a random vector. The method works under the following assumptions:

- ▶  $A$  has an eigenvalue strictly greater than others
- ▶ Starting vector  $x_0$  has a non-zero component in the direction of the eigenvector of the dominant eigenvalue.

Based on gradient method for maximizing convex functions the authors show that a convex function

$$f^* = \max_{x \in Q} f(x) \quad (10)$$

can iteratively maximized by a subgradient, even if that  $f(x)$  is not assumed to be differentiable.

In our case we have to solve a quadratic objective function  $f(x) = \frac{1}{2}x^T Cx$  for  $C \in S_{++}^p$ , which can be solved by

$$x_{k+1} = \frac{C x_k}{\|C x_k\|}, \quad k \geq 0. \quad (11)$$

---

**Algorithm 4:** Block sparse PCA algorithm based on the  $\ell_1$ -penalty (16)

---

**input** : Data matrix  $A \in \mathbf{R}^{p \times n}$

Sparsity-controlling vector  $[\gamma_1, \dots, \gamma_m]^T \geq 0$

Parameters  $\mu_1, \dots, \mu_m > 0$

Initial iterate  $X \in S_m^p$

**output:** A locally optimal sparsity pattern  $P$

**begin**

**repeat**

**for**  $j = 1, \dots, m$  **do**

$x_j \leftarrow \sum_{i=1}^n \mu_j [\mu_j |a_i^T x_j| - \gamma_j]_+ \text{sign}(a_i^T x) a_i$

$X \leftarrow \text{Polar}(X)$

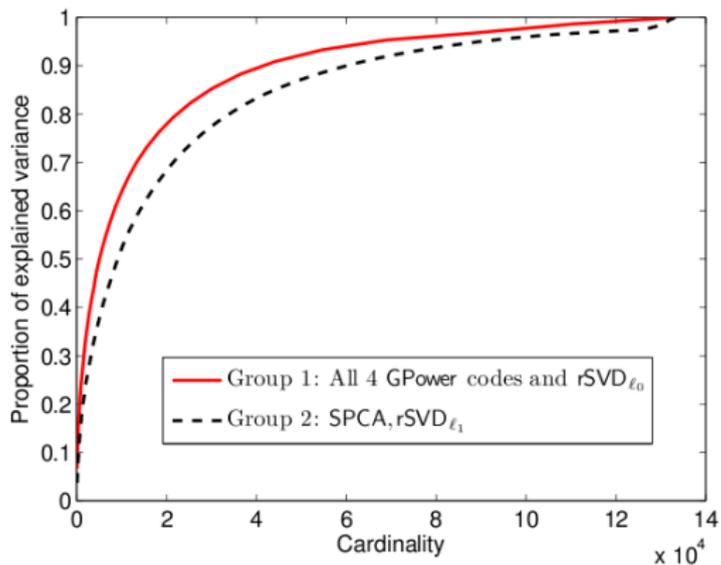
**until** a stopping criterion is satisfied

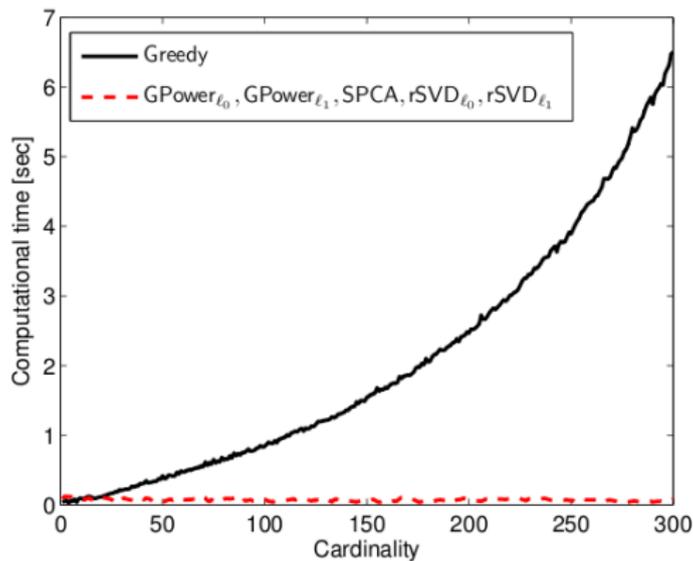
  Construct matrix  $P \in \{0, 1\}^{n \times m}$  such that  $\begin{cases} p_{ij} = 1 & \text{if } \mu_j |a_i^T x_j| > \gamma_j \\ p_{ij} = 0 & \text{otherwise.} \end{cases}$

**end**

---

- ▶ All four GPower algorithms, two single-unit and two block sparse PCA each with  $\ell_0$  and  $\ell_1$  penalty
- ▶ Greedy search algorithm of d'Aspremont et al. (2008) (non-convex)
- ▶ *SPCA* from Zhou et al. (2006) (lasso penalty)
- ▶  $rSVD_{\ell_0}$  and  $rSVD_{\ell_1}$  by Shen and Huang (2008)





$p \times n$	$50 \times 500$	$100 \times 1000$	$250 \times 2500$	$500 \times 5000$	$750 \times 7500$
GPower $_{\ell_1}$	0.22	0.56	4.62	12.6	20.4
GPower $_{\ell_0}$	0.06	0.17	2.15	6.16	10.3
GPower $_{\ell_1,m}$	0.09	0.28	3.50	12.4	23.0
GPower $_{\ell_0,m}$	0.05	0.14	2.39	7.7	12.4
SPCA	0.61	1.47	13.4	48.3	113.3
rSVD $_{\ell_1}$	0.29	1.12	7.72	22.6	46.1
rSVD $_{\ell_0}$	0.28	1.03	7.21	20.7	41.2

Table 8: Average computational time for the extraction of  $m = 5$  components (in seconds).

GPower algorithms show competitive behavior in terms of

- ▶ explained variance
- ▶ computation time
- ▶ control for sparsity pattern
- ▶ usability (data matrix and sample covariance matrix)

- A. d'Aspremont, F. Bach, and L. El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning*, 9:1269–1294, 2008.
- Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11: 517–553, 2010.
- Haipeng Shen and Jianhua Z. Huang. Sparse principal component analysis via regularized low rank matrix approximation. *J. Multivar. Anal.*, 99:1015–1034, July 2008. ISSN 0047-259X. doi: 10.1016/j.jmva.2007.06.007.
- H. Zhou, T. Hastie, and R. Tibshirani. Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, 15:265–286, 2006.

# Thank you for your attention

Thomas Rusch  
Norbert Walchhofer

WU Wirtschaftsuniversität Wien  
Augasse 2–6, A-1090 Wien