

# SPARSE LOADING NOISY PCA USING AN $l_0$ PENALTY

*M.O. Ulfarsson<sup>†</sup> and V. Solo<sup>‡</sup>*

<sup>†</sup>University of Iceland, Dept. Electrical Eng., Reykjavik, ICELAND

<sup>‡</sup>University of New South Wales, School of Electrical Eng., Sydney, AUSTRALIA

## ABSTRACT

In this paper we present a novel model based sparse principal component analysis method based on the  $l_0$  penalty. We develop an estimation method based on the generalized EM algorithm and iterative hard thresholding and an associated model selection method based on Bayesian information criterion (BIC). The method is compared to a previous sparse PCA method using both simulated data and DNA microarray data.

**Index Terms**— principal component analysis, sparsity, estimation.

## 1. INTRODUCTION

Sparse principal component analysis is a relatively recent extension of traditional principal component analysis (PCA). The idea is to seek sparse component loadings while explaining as much variance of the data set as possible.

There are two kinds of sparse PCA: sparse variable PCA (svPCA) [1, 2] and sparse loading PCA (slPCA) [3, 4, 5, 6, 7]. svPCA removes some of the original variables completely by simultaneously zeroing out all their loadings. slPCA keeps all of the original variables but zeroes out some of their loadings. Sparse PCA has, for example, been found useful for genomic data sets [4, 5, 6, 7] and medical imaging data [1, 2].

In this paper we develop a novel model based slPCA method using the  $l_0$  penalty which we call sparse loading noisy PCA (slnPCA). An estimation method based on the generalized EM algorithm is constructed. We modify BIC [8] to choose the associated tuning parameters. The method is compared to the sparse loading PCA method in [4] by using simulations, and a DNA expression microarrays data set.

In section 2 we review noisy PCA (nPCA). The slnPCA method is described in section 3, where we develop an estimation and model selection method. We also give formulas for computing the explained variance of the slnPCA. Section 4 gives simulation examples and a DNA microarray data example. Finally in section 5 conclusions are drawn.

### 1.1. Notation

Vectors and matrices are denoted by boldface letters.  $\mathbf{x} \sim N(\mathbf{m}, \mathbf{C})$  means that the random vector  $\mathbf{x}$  is Gaussian dis-

tributed with mean  $\mathbf{m}$  and covariance  $\mathbf{C}$ .  $I(x \neq 0)$  is the indicator function that is equal to 1 if  $x$  is not equal to zero and equal to 0 else. If  $\mathbf{G}$  is an  $M \times r$  matrix we denote its column vectors by  $\mathbf{g}_{(j)}, j = 1, \dots, r$  and its row vectors by  $\mathbf{g}_i^T, i = 1, \dots, M$ .  $\odot$  is the Hadamard product.

## 2. NOISY PCA

Noisy PCA (nPCA) [9] is a structured covariance model whose maximum likelihood estimator yields PCA. The model is given by

$$\mathbf{y}_t = \mathbf{G}\mathbf{u}_t + \boldsymbol{\epsilon}_t, \quad t = 1, \dots, T, \quad (1)$$

where  $\mathbf{y}_t$  is a zero-mean  $M \times 1$  vector,  $\mathbf{G}$  is an  $M \times r$  loading matrix,  $\boldsymbol{\epsilon}_t \sim N(\mathbf{0}, \sigma^2 \mathbf{I}_M)$ ,  $\mathbf{u}_t \sim N(\mathbf{0}, \mathbf{I}_r)$ , and  $\boldsymbol{\epsilon}_t$  and  $\mathbf{u}_t$  are uncorrelated. The normalized (divided by  $T$ ) log-likelihood function is given by

$$l_{\theta}(\mathbf{Y}) = -\frac{1}{2} \text{tr}(\mathbf{S}_y \boldsymbol{\Omega}^{-1}) - \frac{1}{2} \log |\boldsymbol{\Omega}|,$$

where  $\boldsymbol{\Omega} = \mathbf{G}\mathbf{G}^T + \sigma^2 \mathbf{I}_M$ ,  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_T]^T$ ,  $\mathbf{S}_y = \frac{1}{T} \mathbf{Y}\mathbf{Y}^T$  and  $\boldsymbol{\theta} = (\text{vec}(\mathbf{G})^T, \sigma^2)^T$ . It can be shown that the maximum likelihood estimates of  $\mathbf{G}$  and  $\sigma^2$  are given by

$$\begin{aligned} \hat{\mathbf{G}} &= \mathbf{P}_r(\mathbf{L}_r - \hat{\sigma}^2 \mathbf{I}_r)^{1/2} \mathbf{R}^T \\ \hat{\sigma}^2 &= \frac{1}{M-r} \sum_{j=r+1}^M l_j, \end{aligned}$$

where  $\mathbf{L}_r = \text{diag}(l_1, \dots, l_r)$  is a diagonal matrix containing the  $r$  largest eigenvalues of  $\mathbf{S}_y$ ;  $\mathbf{P}_r$  contains the  $r$  first eigenvectors of  $\mathbf{S}_y$  in its columns;  $\mathbf{R}$  is an arbitrary rotation matrix.

The estimated noisy principal components (nPCs)  $\hat{\mathbf{U}} = [\hat{\mathbf{u}}_1, \dots, \hat{\mathbf{u}}_T]^T$  are given by  $\hat{\mathbf{U}} = \mathbf{Y}\hat{\mathbf{G}}\mathbf{W}^{-1}$ , where  $\mathbf{W} = \hat{\mathbf{G}}^T \hat{\mathbf{G}} + \hat{\sigma}^2 \mathbf{I}_r$ .

## 3. SPARSE LOADING NPCA

We obtain sparse solutions by penalizing the number of entries in  $\mathbf{G}$ . The penalized likelihood is then

$$J_{\theta}(\mathbf{Y}) = -l_{\theta}(\mathbf{Y}) + \frac{h}{2} \|\mathbf{G}\|_0,$$

where  $\|\mathbf{G}\|_0 = \sum_{v=1}^M \sum_{j=1}^r I(g_{v,j} \neq 0)$  and  $g_{v,j}$  is the  $v, j$ -th element of  $\mathbf{G}$ . This should not be confused with our earlier work [2] where we zero out rows of  $\mathbf{G}$  all at once; here the penalty zeros out individual entries in  $\mathbf{G}$ .

### 3.1. Estimation

We maximize the sparse loading nPCA (slnPCA<sub>0</sub>) criterion with the generalized EM algorithm. The penalized normalized complete likelihood is

$$J_{\theta}(\mathbf{Y}, \mathbf{U}) = l_{\theta}(\mathbf{Y}, \mathbf{U}) - \frac{h}{2} \|\mathbf{G}\|_0,$$

where the complete normalized log-likelihood is given by

$$\begin{aligned} l_{\theta}(\mathbf{Y}, \mathbf{U}) &= -\frac{M}{2} \log \sigma^2 - \frac{1}{T} \sum_{t=1}^T \frac{\|\mathbf{y}_t - \mathbf{G} \mathbf{u}_t\|^2}{2\sigma^2} \\ &\quad - \frac{1}{2T} \sum_{t=1}^T \mathbf{u}_t^T \mathbf{u}_t \end{aligned}$$

In the E-step we compute the penalized EM functional

$$\text{EM}_p(\theta_0, \theta) = \text{EM}(\theta_0, \theta) - \frac{h}{2} \|\mathbf{G}\|_0,$$

where the EM functional is given by

$$\begin{aligned} \text{EM}(\theta_0, \theta) &= E_{\theta_0}(l_{\theta}(\mathbf{Y}, \mathbf{U}) | \mathbf{Y}) \\ &= -\frac{\text{tr}(\mathbf{S}_y)}{2\sigma^2} + \frac{\text{tr}(\mathbf{G} \mathbf{B}_0^T)}{\sigma^2} - \frac{\text{tr}(\mathbf{G} \mathbf{A}_0 \mathbf{G}^T)}{2\sigma^2} \\ &\quad - \frac{M}{2} \log \sigma^2, \end{aligned}$$

where  $\mathbf{A}_0 = \sigma_0^2 \mathbf{W}_0^{-1} + \frac{1}{T} \mathbf{U}_0^T \mathbf{U}_0$  and  $\mathbf{B}_0 = \frac{1}{T} \mathbf{Y}^T \mathbf{U}_0$ . In the M-step we maximize the penalized EM functional over  $\mathbf{G}$  which is equivalent to minimizing

$$J(\mathbf{G}) = \sum_{v=1}^M \left( \frac{1}{2} \mathbf{g}_v^T \mathbf{A}_0 \mathbf{g}_v - \mathbf{g}_v^T \mathbf{b}_{v,0} + \frac{h\sigma^2}{2} \|\mathbf{g}_v\|_0 \right),$$

where  $\mathbf{g}_v^T$  is the  $v$ -th row vector of  $\mathbf{G}$ . We use an iterative algorithm (see Algorithm 2 at the end of the paper) to solve this optimization problem, denoting the solution by  $\Gamma(\mathbf{A}_0, \mathbf{B}_0, \sigma^2, h)$ . Optimization of the penalized EM functional w.r.t.  $\sigma^2$  yields

$$\sigma_1^2 = \frac{1}{M} \left[ \text{tr}(\mathbf{A}_0 \mathbf{G}_1^T \mathbf{G}_1) - 2\text{tr}(\mathbf{B}_0^T \mathbf{G}_1) + \text{tr}(\mathbf{S}_y) \right].$$

We say the algorithm has converged if

$$1 - \min_j \left( \frac{|\mathbf{g}_{(j),k}^T \mathbf{g}_{(j),k-1}|}{\|\mathbf{g}_{(j),k}\| \|\mathbf{g}_{(j),k-1}\|} \right) < 1 \cdot 10^{-5}.$$

The slnPCA<sub>0</sub> algorithm is given in Algorithm 1.

---

#### Algorithm 1: The slnPCA<sub>0</sub> algorithm

---

**Input:** Data matrix  $\mathbf{Y}$ ,  $r$ , and  $h$

**Initialization:**  $\mathbf{G}_0$  and  $\sigma_0^2$

**while** (Not converged) **do**

$$\mathbf{W}_k = \mathbf{G}_k^T \mathbf{G}_k + \sigma_k^2 \mathbf{I}_r$$

$$\mathbf{U}_k = \mathbf{Y} \mathbf{G}_k \mathbf{W}_k^{-1}$$

$$\mathbf{A}_k = \sigma_k^2 \mathbf{W}_k^{-1} + \frac{1}{T} \mathbf{U}_k^T \mathbf{U}_k$$

$$\mathbf{B}_k = \frac{1}{T} \mathbf{Y}^T \mathbf{U}_k$$

$$\mathbf{G}_{k+1} = \Gamma(\mathbf{A}_k, \mathbf{B}_k, \sigma_k^2, h)$$

$$\sigma_{k+1}^2 =$$

$$\frac{1}{M} \left[ \text{tr}(\mathbf{A}_k \mathbf{G}_{k+1}^T \mathbf{G}_{k+1}) - 2\text{tr}(\mathbf{B}_k^T \mathbf{G}_{k+1}) + \text{tr}(\mathbf{S}_y) \right]$$

**Output:**  $\mathbf{G}$  and  $\sigma^2$

---

### 3.2. Model selection

We use the BIC statistic to select the number of principal components  $r$  and the sparseness tuning parameter  $h$ .

$$\text{BIC}_{r,h} = -2l_{\hat{\theta}}(\mathbf{Y}) + \frac{d_{r,h}}{T} \log T,$$

where  $d_{r,h}$  is the number of nonzero elements of  $\mathbf{G}$ ; note that this also depends on  $h$ . The model corresponding to the global minimum of the BIC surface is selected.

### 3.3. Variance explained and ordering

We define the variance explained by the sparse noisy principal components as

$$\text{tr}(\mathbf{Q}^T \mathbf{S}_y \mathbf{Q}) \quad (2)$$

where  $\mathbf{Q} = \mathbf{G}(\mathbf{G}^T \mathbf{G} + \sigma^2 \mathbf{I}_r)^{-1/2}$ . The columns of the sparse loading matrix  $\mathbf{G}$  are ordered according to their contribution in the sum (2).

## 4. EXAMPLES

In this section we provide two simulation examples comparing slnPCA<sub>0</sub> to the method in [4], which we will call ZHT, using an implementation by [10]. We also provide a DNA microarray data example. In simulation 2 and in DNA microarray data example the  $n \gg p$  (gene expression array) version of ZHT is used.

### 4.1. Simulation 1

We simulated  $T = 100$  data vectors from (1) with  $\mathbf{G}_{10 \times 2} = [\mathbf{g}_{(1)}, \mathbf{g}_{(2)}]$ , where  $\mathbf{g}_{(1)} = \sqrt{200} \cdot \tilde{\mathbf{g}}_{(1)} / \|\tilde{\mathbf{g}}_{(1)}\|$ ,  $\mathbf{g}_{(2)} = \sqrt{100} \cdot \tilde{\mathbf{g}}_{(2)} / \|\tilde{\mathbf{g}}_{(2)}\|$ ,

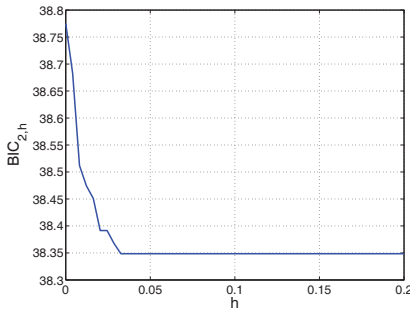
$$\tilde{\mathbf{g}}_{(1)} = [1, 1, 1, 1, 0, 0, 0, 0, 0, 0]^T$$

$$\tilde{\mathbf{g}}_{(2)} = [0, 0, 0, 0, 1, 1, 1, 1, 0, 0]^T$$

	angle 1 (deg)	angle 2 (deg)
slnPCA <sub>0</sub>	2.15	3.10
ZHT	2.49	3.49

**Table 1.** Simulation 1. The angles (in degrees) between the estimated and true loadings.

and  $\sigma^2 = 10$ . The simulation was run 10 times. For each run BIC was used to select the number of principal components  $r$  and the sparseness tuning parameter  $h$ . BIC picked the correct model each time. Fig. 1 shows a sample BIC profile versus the sparseness parameter  $h$ . The model corresponding to the local minimum at  $h = 0.039$  was selected. The ZHT

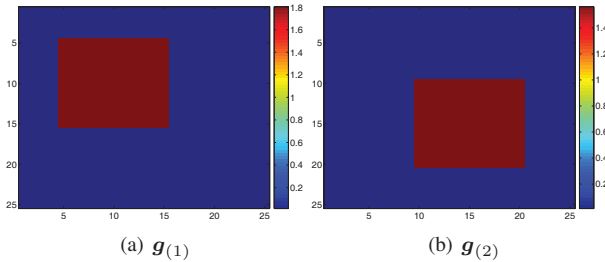


**Fig. 1.** Simulation 1. The  $BIC_{2,h}$  profile for slnPCA<sub>0</sub>.

method is not model based so the BIC cannot be used to select its tuning parameter. To compare we used the same  $r$  and sparsity level as BIC selected for the slnPCA<sub>0</sub> method. We computed the mean angle (in degrees) between the estimated loading vectors  $\hat{g}_{(1)}$ ,  $\hat{g}_{(2)}$  and the truth over 10 simulations. Table 1 shows that angle is lower for the slnPCA<sub>0</sub> method.

#### 4.2. Simulation 2

We simulated  $T = 100$  data vectors from the (1) with  $G_{625 \times 2} = [g_{(1)}, g_{(2)}]$ , where the loading vectors  $g_{(1)}$  and



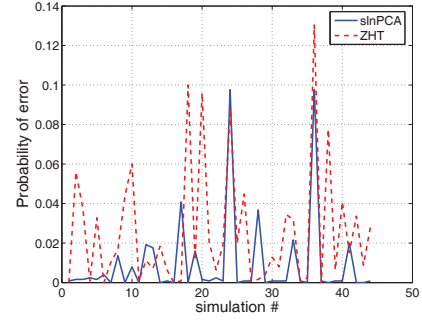
**Fig. 2.** Simulation 2. The loading vector  $g_{(1)}$  and  $g_{(2)}$  depicted as  $25 \times 25$  images.

$g_{(2)}$  are depicted (as  $25 \times 25=625$  images) in Fig. 2 and  $\sigma^2 = 5$ . We simulated the data 44 times and each time used

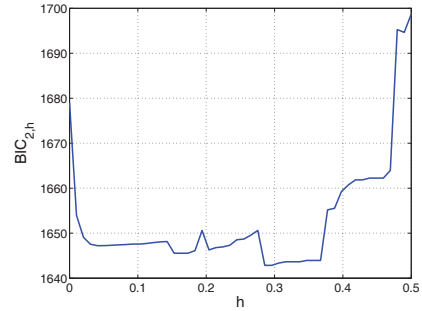
the BIC to select the tuning parameters. Fig. 3 shows a plot of the probability of error for each of those simulations where probability of error was defined as

$$P_E = \alpha P_F + (1 - \alpha)(1 - P_D)$$

, where  $P_D = \Pr(\text{decide signal}|\text{signal})$ ,  $P_F = \Pr(\text{decide signal}|\text{no signal})$ , and  $\alpha$  is the fraction of zero loadings. For each simulation BIC was used to select the slnPCA<sub>0</sub> model. The model corresponding to the lowest probability of error was selected for ZHT. The figure shows that probability of error is almost always much lower for slnPCA<sub>0</sub>. Fig. 4 shows a sample BIC profile vs  $h$  where the model corresponding to  $h = 0.28$  was selected.



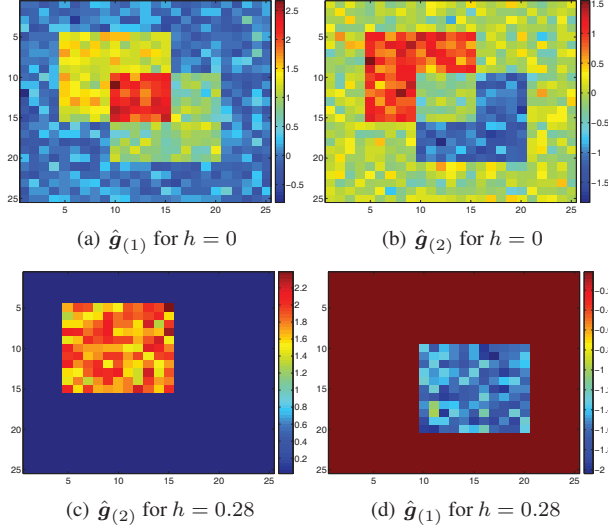
**Fig. 3.** Simulation 2. Probability of error vs simulation number.



**Fig. 4.** Simulation 2. The  $BIC_{2,h}$  profile for slnPCA<sub>0</sub>.

Fig. 5 shows the slnPCA<sub>0</sub> solution for  $h = 0$  (no penalty) which is equal to the nPCA solution described in section 2 and the slnPCA<sub>0</sub> solution selected by BIC. It can be seen that BIC selects the correct degree of sparseness. Interestingly, since the true loading vectors are correlated, the nPCA solution is not able to separate them. On the contrary the slnPCA<sub>0</sub> solution separates them.

The ZHT method, with the degree of sparseness, and  $r$  corresponding to minimum probability of error was used to estimate the loadings. The angle between the loading vectors and the truth can be seen in Table 2. Just as in simulation 1 the slnPCA<sub>0</sub> gives lower values.



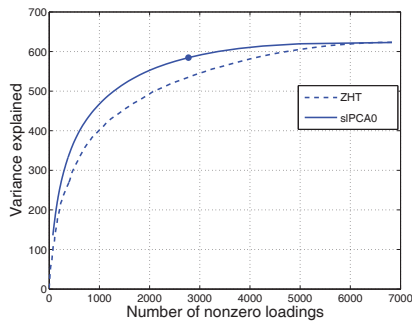
**Fig. 5.** Simulation 2. The estimated loadings

	angle 1 (deg)	angle 2 (deg)
slnPCA <sub>0</sub>	7.6	11.7
ZHT	14.8	23.9

**Table 2.** Simulation 2. The angles between the estimated and true loadings.

### 4.3. DNA microarray expression data

In this section we analyze DNA microarray expression data matrix of  $M=6830$  genes and  $T = 64$  samples [11]. We ex-



**Fig. 6.** The variance explained for the microarray expression data.

tract the first principal component using slnPCA<sub>0</sub> and also the ZHT method. Fig. 6 shows the variance explained vs the level of sparsity, the level of sparsity chosen by BIC is indicated with a dot on the figure. slnPCA<sub>0</sub> clearly explains more variance for a given level of sparsity than ZHT.

## 5. CONCLUSIONS

We have developed a new sparse loading PCA method based on an  $l_0$  penalized likelihood framework. Optimization is facilitated by a generalized EM algorithm which yields a simple thresholding procedure. A simulation study and real data analysis show that the method outperforms the standard procedure.

---

**Algorithm 2:** Iterations for  $\Gamma(A_0, B_0, \sigma^2, h)$ .

---

**Input:** Data matrix  $A_0 = [A_{ij}]$ ,  $B_0 = [b_{(1)}, \dots, b_{(r)}]$ ,  $\sigma^2, h$ .

**Initialization:**  $G_0 = [g_{(1),0}, \dots, g_{(r),0}]$

**while** (Not converged) **do**

**for**  $i = 1, \dots, r$  **do**

$r_{(i),k+1} =$   
         $b_{(i)} - \sum_{j < i} A_{ij} g_{(j),k+1} - \sum_{j > i} A_{ij} g_{(j),k}$   
         $g_{(i),k+1} = \frac{r_{(i),k+1}}{A_{ii}} I \left( \frac{r_{(i),k+1} \odot r_{(i),k+1}}{A_{ii}} > h\sigma^2 \right)$

**Output:**  $G$

---

## 6. REFERENCES

- [1] M.O. Ulfarsson and V. Solo, “Sparse variable PCA using geodesic steepest descent,” *IEEE Transactions on Signal Processing*, vol. 56, no. 12, pp. 5823–5832, 2008.
- [2] M.O. Ulfarsson and V. Solo, “Vector  $l_0$  sparse variable pca,” *IEEE Trans. Signal Proc.*, vol. 59, no. 5, pp. 1949–1958, 2011.
- [3] I.T. Jolliffe and M. Uddin, “A modified principal component technique based on the lasso,” *J. Comput. Graphical Stat.*, vol. 12, no. 3, pp. 531–547, 2003.
- [4] H. Zou, T. Hastie, and R. Tibshirani, “Sparse principal component analysis,” *J. Comput. Graphical Stat.*, vol. 15, no. 2, pp. 265–286, 2006.
- [5] A. D’Aspremont, F. Bach, and L.E. Ghaoui, “Optimal solutions for sparse principal component analysis,” *J. Mach. Learn. Res.*, vol. 9, pp. 1269–1294, 2008.
- [6] J. Huang H. Shen, “Sparse principal component analysis via regularized low rank matrix approximation,” *J. Multivariate Anal.*, vol. 99, pp. 1015–1034, 2008.
- [7] Y. Guan and J. Dy, “Sparse probabilistic principal component analysis,” in *AISTATS 2009*.
- [8] G. Schwartz, “Estimating the dimension of a model,” *Ann. Stat.*, vol. 6, no. 2, pp. 461–464, 1978.
- [9] D.N. Lawley, “Tests of significance of the latent roots of the covariance and correlation matrices,” *Biometrika*, vol. 43, pp. 128–136, 1956.
- [10] K. Sjostrand, “Matlab implementation of LASSO, LARS, the elastic net and SPCA,” 2005, Version 2.0.
- [11] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, Springer Verlag, New York, NY, 2001.