

Robust Sparse Principal Component Analysis

Christophe Croux

Faculty of Business and Economics

K.U. Leuven

3000 Leuven, Belgium

Christophe.Croux@econ.kuleuven.be

Peter Filzmoser and Heinrich Fritz

Department of Statistics and Probability Theory

Vienna University of Technology

1040 Vienna, Austria

P.Filzmoser@tuwien.ac.at; heinrich@fritz.cc

Abstract

A method for principal component analysis is proposed that is sparse and robust at the same time. The sparsity delivers principal components that have loadings on a small number of variables, making them easier to interpret. The robustness makes the analysis resistant to outlying observations. The principal components correspond to directions that maximize a robust measure of the variance, with an additional penalty term to take sparseness into account. We propose an algorithm to compute the sparse and robust principal components. The algorithm computes the components sequentially, and thus it can handle data sets with more variables than observations. The method is applied on several real data examples, and diagnostic plots for detecting outliers and for selecting the degree of sparsity are provided. A simulation experiment studies the effect on statistical efficiency by requiring both robustness and sparsity.

Keywords: dispersion measure, outliers, projection-pursuit, variable selection

1 Introduction

Principal component analysis (PCA) is a standard tool for dimension reduction of multivariate data. PCA searches for linear combinations of the variables, called principal components (PCs), that summarize well the data. The PCs correspond to directions maximizing the variance of the data projected on them (see, e.g., Jolliffe, 2002). The transformation matrix defining the principal components is called the loadings matrix, and it may be used to interpret the PCs. In general, PCA does not deliver easily interpretable components. Good interpretability of PCs is related to rather large or small (absolute) values in the loadings matrix yielding either quite strong or very weak contributions of the variables to the PCs. Loadings matrices with many values exactly

equal to zero, which we call *sparse loadings matrices*, are preferred. The interpretation of a particular principal component then only depends on a small subset of variables. This yields a *sparse PCA*, which is especially helpful for analyzing high dimensional data sets. In this paper we introduce a method for PCA that yields both sparse and *robust* results. Outliers frequently occur in multivariate data sets, and any multivariate procedure should take the possible presence of outliers into account.

Different approaches for computing sparse loadings matrices have been proposed in the literature. Vines (2000) and Anaya-Izquierdo *et al.* (2011) use a restriction on the loadings to integer values. Chipman and Gu (2005) introduce the constraint that as many zeros as possible should occur in the loadings matrix, with the aim of better interpretability. Jolliffe *et al.* (2003) introduce the SCoTLASS, related to the Lasso estimator (Tibshirani, 1996). Here the principal components maximize the variance but under an upper bound on the sum of the absolute values of the loadings. It is shown that such an approach yields better results than a two-step procedure, where after a standard PCA rotation techniques are performed (Jolliffe, 1995). Zou *et al.* (2006) use the elastic net to obtain a version of sparse PCA. Modifications and improvements of this method are made in Leng and Wang (2009). Witten *et al.* (2009) develop a general procedure for penalized matrix decomposition, and they show how to apply it for sparse PCA. Finally, Guo *et al.* (2010) introduce a fusion penalty to capture block structures within the variables. Other types of structured sparse PCA are considered in Jenatton *et al.* (2009) and Bien *et al.* (2010). All these methods, however, are not robust to outliers.

This paper proposes a PCA method that is robust and sparse at the same time. Several robust, but non sparse, PCA methods have been introduced in the literature (see, e.g., Filzmoser, 1999; Hubert *et al.*, 2005; Maronna, 2005), and robustness properties were investigated (Croux and Haesbroeck, 2000). Recently, fast procedures for robust PCA have been proposed in the machine learning literature, e.g. Candès *et al.* (2009); Xu *et al.* (2010, 2012). While these methods are suitable for high dimensional problems, and have good properties in practice, they do not deliver sparse loadings matrices. In this paper we focus on the *projection-pursuit* approach to PCA, where the PCs are extracted from the data by searching for directions that maximize a robust measure of variance of the data projected on it (Li and Chen, 1985; Croux and Ruiz-Gazen, 2005). Using a robust measure of variance prevents the PCs from being attracted by the outliers, since outliers

inflate the standard non-robust variance. An efficient algorithm for computing the projection-pursuit based PCs is the *Grid algorithm*, introduced in Croux *et al.* (2007). An implementation is available in the R package `pcaPP` (Filzmoser *et al.*, 2010). To the best of our knowledge, the PCA method we propose is the first one combining the properties of robustness and sparsity.

The paper is organized as follows: Section 2 defines the robust sparse principal components as the solution of a non-convex optimization problem. Section 3 shows how the Grid algorithm can be extended to find an approximate solution of this problem. The selection of tuning parameters is discussed in Section 4. Simulation results are presented in Section 5, and real data examples are shown in Section 6. The final Section 7 concludes.

2 Method

For n multivariate observations $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^p$, collected in the rows of the data matrix \mathbf{X} , the first PCA direction is given by

$$\mathbf{a}_1 = \underset{\|\mathbf{a}\|=1}{\operatorname{argmax}} V(\mathbf{a}^t \mathbf{x}_1, \dots, \mathbf{a}^t \mathbf{x}_n), \quad (1)$$

where V is a variance measure. In the standard non-robust case, V is the empirical variance (Var), and the resulting optimal direction \mathbf{a}_1 corresponds to the first eigenvector of the sample covariance matrix. Equation (1) is the projection-pursuit formulation for finding the first PC, with V the projection-pursuit index. Robust PCA directions can easily be obtained by taking a robust variance measure for V , like the squared Median Absolute Deviation or the more efficient squared Q_n estimator proposed in Rousseeuw and Croux (1993).

For a univariate data set y_1, \dots, y_n , the Q_n scale estimator is defined as the first quartile of all pairwise distances $|y_i - y_j|$, for $1 \leq i < j \leq n$. This first quartile is then multiplied by the constant 2.219 to get a consistent estimator for the scale of a normal distribution, and Q_n^2 is a consistent and robust estimator for the variance. Croux and Ruiz-Gazen (2005) show that using the Q_n^2 estimator as projection index yields robust and efficient estimates for the principal components. In the remainder of this paper, we use the Q_n^2 as robust variance estimator.

Suppose the first $j - 1$ PCA directions have already been found ($j > 1$), then the j th direction ($j \leq p$) is defined as

$$\mathbf{a}_j = \underset{\|\mathbf{a}\|=1, \mathbf{a} \perp \mathbf{a}_1, \dots, \mathbf{a} \perp \mathbf{a}_{j-1}}{\operatorname{argmax}} V(\mathbf{a}^t \mathbf{x}_1, \dots, \mathbf{a}^t \mathbf{x}_n), \quad (2)$$

imposing an orthogonality constraint to all previously found directions. The j th principal component is then the vector containing the PCA *scores*

$$z_{ij} = \mathbf{a}_j^t \mathbf{x}_i \quad \text{for } i = 1, \dots, n. \quad (3)$$

The loadings matrix for the first k PCs is denoted by \mathbf{A}_k , and it contains in its columns the optimal directions or loadings vectors \mathbf{a}_j , for $1 \leq j \leq k$. The loadings determine the contribution of each variable to the principal components. The matrix containing the principal component scores is then

$$\mathbf{Z}_k = \mathbf{X} \mathbf{A}_k. \quad (4)$$

Sparsity can be imposed on the PCA directions by adding an L_1 penalty in the objective function. As such, Jolliffe *et al.* (2003) introduced the SCoTLASS criterion,

$$\max_{\|\mathbf{a}\|=1, \mathbf{a} \perp \mathbf{a}_1, \dots, \mathbf{a} \perp \mathbf{a}_{j-1}} \mathbf{a}^t \hat{\Sigma} \mathbf{a}, \quad \text{subject to } \|\mathbf{a}\|_1 \leq t, \quad (5)$$

for obtaining the j th PCA direction, with $1 \leq j \leq p$. Here, $\hat{\Sigma}$ is the empirical covariance matrix, and the L_1 norm $\|\mathbf{a}\|_1 = \sum_{j=1}^p |\mathbf{a}_j|$ takes the sum of the absolute values of the components of the vector \mathbf{a} . It is more convenient to work with the dual formulation of the above problem, given by

$$\max_{\|\mathbf{a}\|=1, \mathbf{a} \perp \mathbf{a}_1, \dots, \mathbf{a} \perp \mathbf{a}_{j-1}} \mathbf{a}^t \hat{\Sigma} \mathbf{a} - \lambda_1 \|\mathbf{a}\|_1, \quad (6)$$

where λ_1 is a tuning parameter. The larger λ_1 , the more the components of \mathbf{a} are shrunken towards zero. Due to the use of the L_1 penalty, some of the loadings will even become exactly zero, similar as for the Lasso estimator in regression. The approach of Jolliffe *et al.* (2003) requires an estimated covariance matrix $\hat{\Sigma}$ as input of the maximization problem (5), which can be solved using the algorithm detailed in Trendafilov and Jolliffe (2006) or in Journée *et al.* (2010).

An obvious way for obtaining a sparse robust PCA method would be to simply replace the empirical covariance matrix by a robust covariance estimator. This is referred to as the plug-in approach, which is common in robust multivariate data analysis (Hubert *et al.*, 2008). But computing robust covariance matrices in high dimensions, and particularly if $p > n$, is cumbersome—the estimator may even not exist—and time consuming. There are, however, some non affine equivariant robust covariance matrix estimators, that are suitable for this plug-in approach. As such, the Spatial Sign Covariance matrix (Locantore *et al.*, 1999; Oja, 2010) is fast to compute in

high dimensions, and has a high breakdown point (Croux *et al.*, 2010a). In Section 5 we compare our proposed method with this plug-in approach.

In this paper, we employ the projection-pursuit approach where the PCs are directly obtained without estimating a covariance matrix. Adding the L_1 constraint in definition (1) yields

$$\tilde{\mathbf{a}}_1 = \underset{\|\mathbf{a}\|=1}{\operatorname{argmax}} V(\mathbf{a}^t \mathbf{x}_1, \dots, \mathbf{a}^t \mathbf{x}_n) - \lambda_1 \|\mathbf{a}\|_1. \quad (7)$$

The vector $\tilde{\mathbf{a}}_1$ is the first sparse PCA direction, and its sparsity is controlled by the tuning parameter λ_1 . Setting $\lambda_1 = 0$ results in the unconstrained first PCA direction \mathbf{a}_1 , but for increasing values of λ_1 , sparsity gains importance compared to robust variance maximization. Similarly, the j th sparse PCA direction ($1 < j \leq p$) is defined by

$$\tilde{\mathbf{a}}_j = \underset{\|\mathbf{a}\|=1, \mathbf{a} \perp \tilde{\mathbf{a}}_1, \dots, \mathbf{a} \perp \tilde{\mathbf{a}}_{j-1}}{\operatorname{argmax}} V(\mathbf{a}^t \mathbf{x}_1, \dots, \mathbf{a}^t \mathbf{x}_n) - \lambda_j \|\mathbf{a}\|_1, \quad (8)$$

with λ_j a tuning parameter, possibly different from λ_1 . Definition (7) and (8) are very elegant and simple, and maintain the basic interpretation of the principal components: we look for directions maximizing a robust variance, under the constraint that the loadings should not become too large. This constraint is considered by penalizing with the L_1 norm of \mathbf{a} . If $V = \operatorname{Var}$, then definitions (6) and (7) are the same. Note that most often one does not need all possible PCs, but only the first few. An advantage of the projection-pursuit approach is that the PCs are computed sequentially, reducing the computation time if the interest is only in the first few PCs.

3 Algorithm

To find the optimal directions in (1) and (2), maximization over a p -dimensional space is required. This is a non-convex optimization problem, even for $V = \operatorname{Var}$, and it is not possible to find analytical solutions for the optimal directions. Moreover, since V may be non-differentiable in its arguments, gradient-based methods are not always possible. Several proposals for finding good approximations of the projection-pursuit based PCs, applicable for any choice of the projection index V , have been made (Hubert *et al.*, 2002; Croux and Ruiz-Gazen, 2005; Croux *et al.*, 2007). In this paper we extend the Grid algorithm of Croux *et al.* (2007) for obtaining sparse solutions, i.e. to solve (7) and (8). The algorithm is fast to compute and was shown to be quite accurate,

even in larger dimensions. A complete implementation of the algorithm is available in the R package `pcaPP` (Filzmoser *et al.*, 2010). Below we give an outline of the algorithm.

Let k be the number of sparse PCs that need to be computed. Assume that the first $j - 1$ sparse PCA directions $\tilde{\mathbf{a}}_{j-1}$ are already obtained and are collected in the first $j - 1$ columns of the loadings matrix $\tilde{\mathbf{A}}_{j-1}$, with $1 \leq j \leq k - 1$. Now we want to compute $\tilde{\mathbf{a}}_j$. For notational consistency, set $\tilde{\mathbf{A}}_0^\perp$ equal to the identity matrix. For $j > 1$, let $\tilde{\mathbf{A}}_{j-1}^\perp$ be a matrix containing in its columns an orthonormal basis for the subspace orthogonal to the space spanned by the first $j - 1$ sparse PCA directions, and denote $\mathbf{x}_i^{(j-1)} = (\tilde{\mathbf{A}}_{j-1}^\perp)^t \mathbf{x}_i$, for $i = 1, \dots, n$. Solving problem (8) is then equivalent to maximizing the objective function

$$f(\mathbf{a}) = V(\mathbf{a}^t \mathbf{x}_1^{(j-1)}, \dots, \mathbf{a}^t \mathbf{x}_n^{(j-1)}) - \lambda_j \|\tilde{\mathbf{A}}_{j-1}^\perp \mathbf{a}\|_1, \quad (9)$$

under the restriction that $\|\mathbf{a}\| = 1$, with \mathbf{a} belonging to the lower-dimensional space \mathbb{R}^{p-j+1} . As sparseness relates to the components of a direction in the space of the original variables, and not to the lower dimensional space \mathbf{a} belongs to, we need to back-transform the vector \mathbf{a} to the original space before taking the L_1 norm in (9).

For optimizing (9) the Grid algorithm is used. The basic idea of this algorithm is to reduce the problem to a sequence of optimizations in a two-dimensional plane under the unit norm constraint. This boils down to a sequence of maximizations of a function over the unit circle, which is simply a univariate maximization problem that can be solved by a grid search over $[-\pi, \pi]$. Consider the optimization of (9) for a given value of $1 \leq j \leq k$. We take the following steps:

- Let $\mathbf{X}^{(j)}$ be the matrix having the vectors $\mathbf{x}_i^{(j-1)}$ in its rows. Denote $p_j = p - j + 1$ the number of columns of this matrix. Compute the value of the projection index $V_s = V(x_{1s}^{(j)}, \dots, x_{ns}^{(j)})$ for every column s of the matrix, and permute the column indices such that the projection indices are decreasing in s : $V_1 \geq V_2 \dots \geq V_{p_j}$. The starting value of the algorithm is then $\mathbf{a} = (1, 0, \dots, 0)$, a vector with p_j components.
- For $l = 1, \dots, \text{maxiter}$ (the outer loop), perform an improvement step in which the currently best solution $\mathbf{a} = (a^1, \dots, a^{p_j})^t$ is updated according to the following scheme:

For $1 \leq s \leq p_j$ (the inner loop)

(i) Maximize by univariate grid search the function

$$\gamma \rightarrow f(a^1 b(\gamma), \dots, a^{s-1} b(\gamma), \cos \gamma, a^{i+1} b(\gamma), \dots, a^{p_j} b(\gamma)),$$

where $b(\gamma) = \sin(\gamma)/\sqrt{1 - (a^s)^2}$ is such that the unit norm condition on the argument of the functions f continues to hold. The value of γ ranges over the interval $[\arccos(a^s) - \pi/(2^{l-1}), \arccos(a^s) + \pi/(2^{l-1})]$. Hence, for $l = 1$, the grid search is over the whole unit circle. As the iteration step l increases, we perform a more restricted search over a finer grid. This function is maximized by a grid search using $Ngrid$ evaluation points, with $Ngrid$ an odd integer. As such, $\gamma = \arccos(a^s)$, corresponding to the currently best solution, is always a grid point. Since $Ngrid$ remains constant, we are increasing the precision of the grid search in every iteration step of the outer loop.

(ii) Update the currently best solution as

$$\mathbf{a} \leftarrow f(a^1 b(\gamma^*), \dots, a^{s-1} b(\gamma^*), \cos \gamma^*, a^{s+1} b(\gamma^*), \dots, a^{p_j} b(\gamma^*))^t,$$

with γ^* corresponding to the maximum of the grid search. The new currently best solution can never have a lower value of the objective function f , since the previously best solution belongs to the finite grid over which we maximized.

- Finally, the optimal sparse direction \mathbf{a} found for the j th PC by the grid algorithm has to be back-transformed into the original space, yielding $\tilde{\mathbf{a}}_j = \tilde{\mathbf{A}}_{j-1}^\perp \mathbf{a}$.

The grid algorithm is a variant of a coordinate descent method. As follows from the description of the algorithm, the value of the objective function is guaranteed to be non decreasing over the iteration steps. Hence, convergence is guaranteed, and if $Ngrid$ and $maxiter$ are taken sufficiently large it can be expected that the algorithm converges to a local maximum of the objective function. Since the problem is nonconvex, one cannot expect to have guaranteed convergence to a global optimum; see Tseng (2001) for convergence properties of coordinate descent methods. In the simulations and numerical experiments, we take $Ngrid = 25$ and $maxiter = 10$ by default. We carried out many numerical experiments, and the grid algorithm always gave a robust and sparse solution, within reasonable computing time. Using these default parameters, the CPU time on a quad-core Intel Core i7 at 2GHz for various values of n and p is reported in Table 1 (in seconds). The computation time increases about linearly in n and exponentially in p .

Table 1: CPU time on a quad-core Intel Core i7 at 2GHz (in seconds), using the default parameters, and $k = 2$, $\lambda = 1$, scale measure MAD.

	$p = 10$	$p = 50$	$p = 100$	$p = 500$	$p = 1000$
$n = 10$	0.005	0.017	0.056	3.663	30.041
$n = 50$	0.012	0.059	0.185	4.120	29.092
$n = 100$	0.021	0.108	0.229	4.897	31.862
$n = 500$	0.083	0.417	0.896	9.029	44.142
$n = 1000$	0.178	0.934	1.823	14.089	56.193

4 Selection of sparsity parameter λ

The tuning parameter λ_j regulates the degree of sparseness. The larger λ_j , the less weight is given to the robust variance measure V in the objective function (8), for $j = 1, \dots, k$. To make the relative importance of the penalty term in (8) comparable across the different PCs, i.e. to have a similar degree of sparsity over the different principal components, we take

$$\lambda_j := \lambda \mathcal{V}(\mathbf{X}^{(j)}), \quad (10)$$

where the matrix $\mathbf{X}^{(j)}$ is defined in the previous section, and contains the data vectors projected on the orthogonal complement of the space spanned by the first $j-1$ optimal directions. Furthermore, \mathcal{V} denotes the total robust variance of a data matrix, and is for any n by p matrix \mathbf{Y} defined as

$$\mathcal{V}(\mathbf{Y}) = \sum_{s=1}^p V(\mathbf{y}_s), \quad (11)$$

where \mathbf{y}_s stands for the s th column of \mathbf{Y} and V is the robust variance measure used as projection index. Using this measure of total variation in (10) as a scaling factor, there is only one tuning parameter λ to be selected.

We propose to select the λ to minimize a BIC type criterion (see also Guo *et al.*, 2010; Leng and Wang, 2009)

$$\text{BIC}(\lambda) = \frac{\widetilde{\text{RV}}}{\text{RV}} + \text{df}(\lambda) \frac{\log(n)}{n}, \quad (12)$$

where $\widetilde{\text{RV}}$ and RV refer to the total robust variance of the residuals matrix obtained from a sparse PCA and an unconstrained PCA. Here, $\text{df}(\lambda)$ is the number of non-zero loadings when using λ

as the penalty parameter, as in Guo *et al.* (2010). The first term in the BIC is a measure for the quality of the fit, while the second term penalizes for model complexity. The calculation of \widetilde{RV} and RV is immediate, since they are given by

$$\widetilde{RV} = \mathcal{V}(\mathbf{X} - \mathbf{X} \tilde{\mathbf{A}}_k \tilde{\mathbf{A}}_k^t) \text{ and } RV = \mathcal{V}(\mathbf{X} - \mathbf{X} \mathbf{A}_k \mathbf{A}_k^t),$$

where \mathbf{X} stands for the data matrix, and \mathbf{A}_k and $\tilde{\mathbf{A}}_k$ denote the loadings matrices containing the first k PC directions (in the columns) for unconstrained and constrained PCA, respectively. Note that, for $V = \text{Var}$, the BIC criterion in (12) is equal to the one in Guo *et al.* (2010). In practice, the selection of λ is carried out by minimizing $\text{BIC}(\lambda)$ over a grid $[0, \lambda_{max}]$, where λ_{max} results in full sparseness of the sparse PCA solution with k components (i.e. every loading vector contains only one non-zero element).

In addition to λ , one also needs to choose the number of components k . The BIC criterion (12) is also depending on the choice of k . Appropriate selection of k is an old and common problem in principal components analysis, and many proposals have been made for it. In this paper we select the number k from the scree-plot of an unconstrained PCA (see Cattell, 1966). Such a scree-plot represents the percentage of explained (robust) variance (EV) by the PCs versus the number of principal components. Mathematically, the explained (robust) variance is given by

$$EV_k = \frac{\mathcal{V}(\mathbf{Z}_k)}{\mathcal{V}(\mathbf{X})}, \tag{13}$$

with \mathbf{Z}_k the matrix containing the principal component scores, see (4). For $V = \text{Var}$, EV_k equals the ratio of the sum of the k largest eigenvalues to the sum of all eigenvalues of the sample covariance matrix. When running this sparse PCA method, the same number of PCs is maintained. For this value of k , a selected λ should result in a sparser loadings matrix, with less (robust) variance. In Section 6 we present the so-called *tradeoff curve*, where the percentage of explained variance of the k sparse PCs is plotted as a function of λ . This graphical tool may be used, in addition to the BIC, for selecting an appropriate value of λ .

5 Simulation experiments

In this section we present two simulation experiments. The sparse method should (i) result in increased estimation precision when the true loadings matrix is sparse, and (ii) succeed in detecting

those variables that do not contribute to the principal components, i.e. true zero loadings are exactly estimated as zero. We contrast the standard sparse method, with $V = \text{Var}$, with the robust sparse method, with $V = Q_n^2$. If no outliers are present, then the two properties above hold for both methods. But it will be shown that, in presence of outliers, the standard sparse method does not meet its objectives anymore.

Experiment 1

We generate data sets of $n = 50$ observations in p dimensions. The true loadings matrix is

$$\mathbf{A} = \begin{pmatrix} \sqrt{0.5} & 0 & \sqrt{0.5} & 0 & 0 & \cdots & 0 \\ \sqrt{0.5} & 0 & -\sqrt{0.5} & 0 & 0 & \cdots & 0 \\ 0 & \sqrt{0.5} & 0 & \sqrt{0.5} & 0 & \cdots & 0 \\ 0 & \sqrt{0.5} & 0 & -\sqrt{0.5} & 0 & \cdots & 0 \\ 0 & 0 & 0 & 0 & 1 & & 0 \\ \vdots & \vdots & \vdots & \vdots & & \ddots & \\ 0 & 0 & 0 & 0 & 0 & & 1 \end{pmatrix}$$

and the eigenvalues are $\mathbf{l} = (1, 0.5, 0.1, 0.1, \dots, 0.1)$. The observations are generated from a multivariate normal distribution $N_p(\mathbf{0}, \mathbf{A}\mathbf{L}\mathbf{A}^t)$, with $\mathbf{L} = \text{diag}(\mathbf{l})$. Contamination is added by replacing a portion of δ observations by outliers, generated from the distribution $N_p(\boldsymbol{\mu}_{out}, \mathbf{I}_p)$ with $\boldsymbol{\mu}_{out} = (2, 4, 2, 4, 0, -1, 1, 0, 1, -1, \dots, 0, 1, -1)^t$. From the generated data set the loadings matrix is estimated, with $k = 2$. The resulting $\hat{\mathbf{A}}_2$ is compared to the true \mathbf{A}_2 , containing the first two columns of \mathbf{A} , by computing the angle φ between the subspaces spanned by the columns of the matrices.

Both the standard procedure and the robust sparse PCA procedure are applied to $m = 100$ simulated data sets. Figure 1 shows the average value of $2\varphi/\pi$ over the m simulations, called the average deviation in Hubert *et al.* (2005), as a function of the tuning parameter λ . Values close to one imply that the estimated subspace is almost orthogonal to the true one. Different outlier proportions, ranging from no contamination to 40% of outliers are considered. The number of variables taken is $p = 10$.

If no outliers are present ($\delta = 0$, solid line), we get the expected pattern. Starting with $\lambda = 0$ (i.e. non sparse PCA) the estimation error decreases until a minimum is reached at about $\lambda = 0.7$.

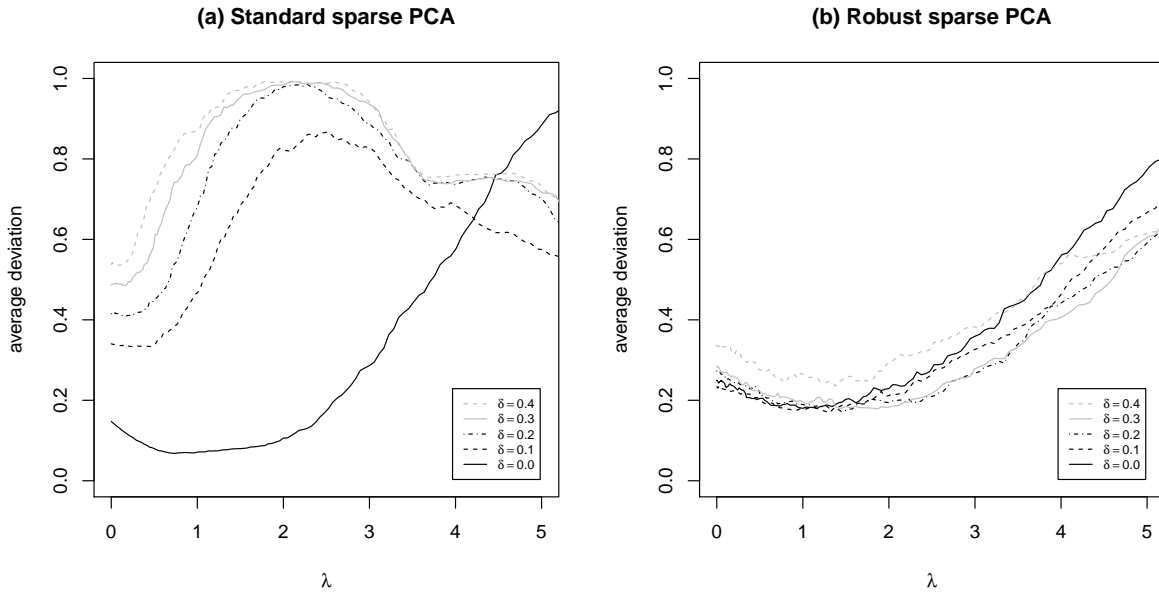


Figure 1: Average deviation between estimated and true loadings for (a) the standard and (b) robust sparse PCA methods for different levels of contamination δ and different values of λ . The sample size is $n = 50$ and the dimension $p = 10$.

Penalizing the loadings further yields again an increasing estimation error. If the true model is sparse (here about 80% of the true loadings are zero), sparse estimation methods may improve the precision of the maximum likelihood method. For the robust sparse method a similar pattern is observed. There is a slight loss in precision using the robust instead of the standard method. However, the robust method remains fairly accurate under contamination, as can be seen from the other curves in Figure 1 (b). This is in contrast with the standard method, where the estimation error increases substantially and exceeds the robust counterpart by a large amount. Finally, note that in presence of outliers the advantage of penalizing disappears for the standard method, since $\lambda = 0$ yields the smallest average deviation $2\varphi/\pi$. This does not happen for robust sparse PCA.

We repeat the same simulation experiment for $p = 200 > n = 50$. Here the number of variables largely exceeds the number of observations. The average deviations using standard and robust sparse PCA are given in Figure 2, showing that the robust sparse method is suitable for high dimensions as well. The results are comparable to Figure 1, but standard sparse PCA performs now even worse under contamination. We conclude that requiring robustness leads to a small

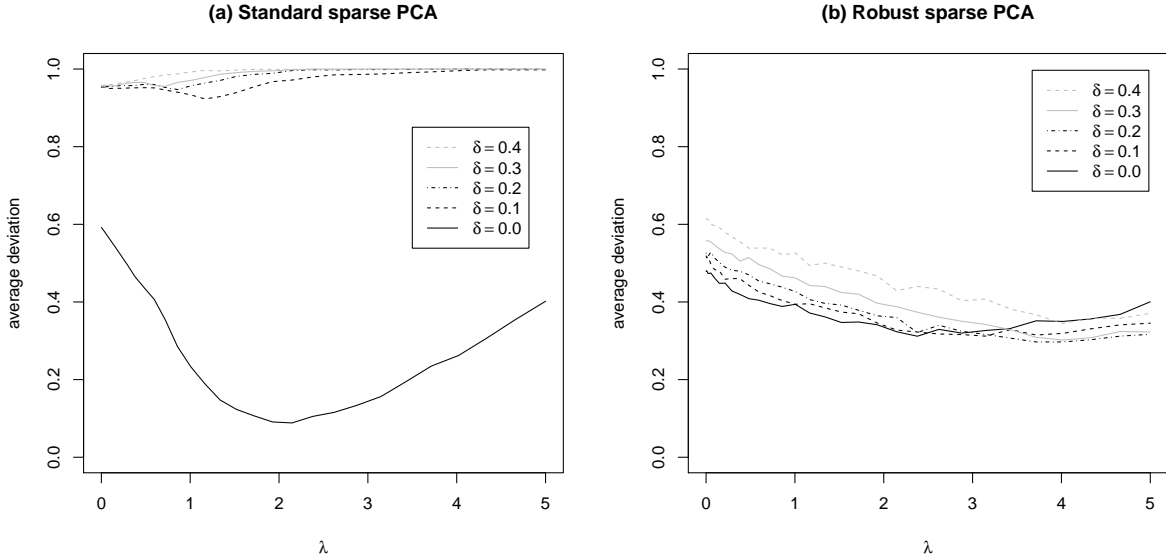


Figure 2: As Figure 1, but for sample size $n = 50$ and $p = 200$.

loss of statistical efficiency in absence of outliers, but to a major improvement when outliers are present. Requiring sparsity does not necessarily lead to a loss of precision. For a large range of values of λ , adding sparsity to robust PCA is even increasing the statistical efficiency.

We now compare the performance of the robust projection-pursuit approach with a plug-in approach. Here the sample covariance matrix in (6) is replaced by the robust Spatial Sign Covariance matrix (Locantore *et al.*, 1999; Oja, 2010). In Figure 3 the simulated average deviation is pictured as a function of λ , for different percentages of contamination, $n = 50$, and dimensions $p = 10$ (left panel) and $p = 200$ (right panel). The average deviations of the plug-in approach are much higher for all λ values, indicating that the plug-in approach does not deliver robust solutions. A reason might be that a large $p \times p$ robust covariance matrix needs to be estimated before the sparse PCs are computed. If the sample size is low or the dimension large, this may induce large estimation errors.

We repeated the above simulation exercise with outlier means $2\boldsymbol{\mu}_{out}$ and $\boldsymbol{\mu}_{out}/2$. We also made the outliers more concentrated by dividing the covariance matrix of the outlier generating distribution by 10. For all these configurations, the obtained results were quite similar to the ones reported above. In most but not all cases, changing the direction of $\boldsymbol{\mu}_{out}$ leads to bad performance of the standard method. However, when the outliers are in the subspace of the first k principal

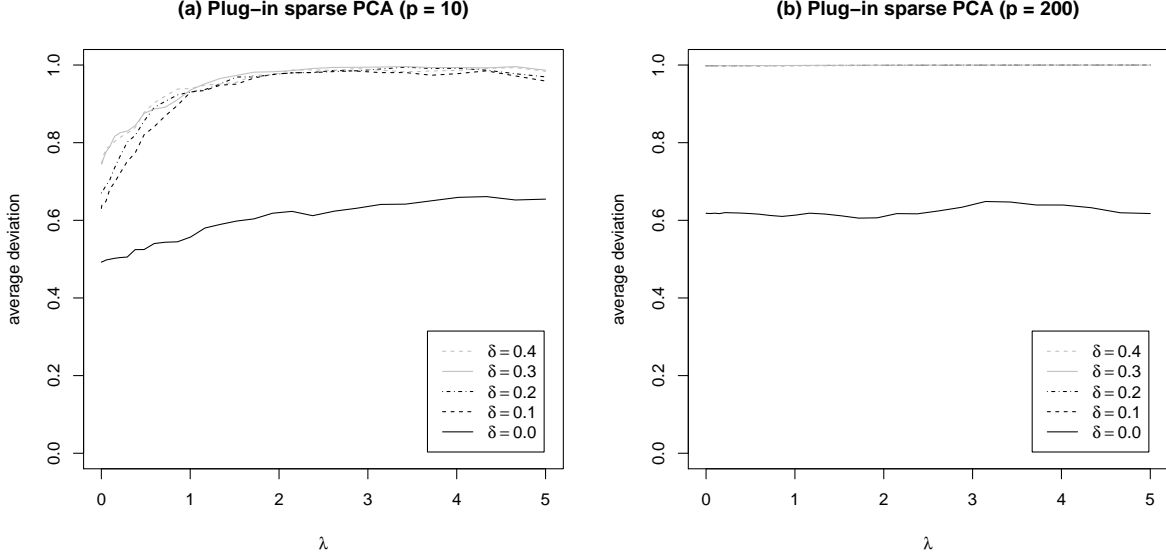


Figure 3: Average deviation between estimated and true loadings for the plug-in approach to robust sparse PCA for different levels of contamination δ and different values of λ . The sample size is $n = 50$ and the considered dimensions are $p = 10$ (left panel) and $p = 200$ (right panel)

components, then standard (sparse) PCA is able to cope with them. The reason is that this type of outliers will not attract the estimates of the columns of the loadings matrix in the wrong direction.

Experiment 2

We consider the same design as introduced by Zou *et al.* (2006), and subsequently used by Farcomeni (2009) and Guo *et al.* (2010) in the same context of sparse PCA. We have $n = 20$ observations and $p = 10$ variables driven by two latent variables

$$U_1 \sim N(0, 290), \quad U_2 \sim N(0, 300),$$

where $\varepsilon \sim N(0, 1)$, and U_1 , U_2 and ε are independent. The variances of the two latent variables are 290 and 300, respectively. The observed variables are constructed as

$$X_j = \begin{cases} U_1 + \varepsilon_j, & \text{if } 1 \leq j \leq 4 \\ U_2 + \varepsilon_j, & \text{if } 5 \leq j \leq 8 \\ -0.3U_1 + 0.925U_2 + \varepsilon + \varepsilon_j, & \text{if } j = 9, 10. \end{cases}$$

The error terms ε and ε_j , for $1 \leq j \leq 10$, are i.i.d. $N(0, 1)$. The first two principal components correspond to U_2 and U_1 , respectively, and in this order. The first block of variables, X_1 to X_4 is expected to have a high loading on the second PC, but zero loadings on the first one. The second block, X_5 to X_8 , should have high loadings on the first PC, but a zero loading on the second one. The remaining variables X_9 and X_{10} have a larger loading on the first PC than on the second one, and a sparse PCA could shrink this second loading to zero. We will add outliers generated from the distribution $N(\boldsymbol{\mu}_{out}, \sigma_{out}^2 \mathbf{I}_{10})$, with $\boldsymbol{\mu}_{out} = (0, -100, 100, 50, 0, 100, -100, 50, 75, -75)^t$, and $\sigma_{out}^2 = 20$. These added data are not univariate outliers, and hence are not detectable by making boxplots of the individual variables. The outliers do not follow the factor structure described above.

We generate $m = 100$ samples according to the simulation design, using outlier portions 0%, 10%, and 20%, and apply the standard and the robust version of the sparse PCA algorithm. For every sample, an optimal value of the tuning parameter was selected according to the BIC criterion. Then loadings of each of the 10 variables on the first two PCs are computed, as well as the percentage of explained (robust) variance EV. The reported values correspond to the median and median absolute deviation (MAD, between parenthesis) over the 100 replications, and are presented in Table 2, in a similar way as in Table 1 of Guo *et al.* (2010). The MAD is defined for a sample y_1, \dots, y_n as $\text{MAD}(y_1, \dots, y_n) = \text{median}_i |y_i - \text{median}_j y_j|$.

Without contamination (0%), the results are according to the expectations, and very much comparable to those of Guo *et al.* (2010). For both the standard and the robust sparse method, variables X_5 through X_{10} are solely represented in the first PC, variables X_1 to X_4 in the second PC, and the loadings of the last two variables for the second PC are shrunk to zero. When adding contamination it is seen from Table 2 that the standard PCA becomes distorted, and does not recover the sparsity in the data generating process. The robust method, however, still delivers sparse solutions. The price the robust method pays for the resistance with respect to outliers is an increased variability, as measured by the MAD values.

The standard sparse PC directions are attracted by the outliers and do no longer explain the actual structure of the majority of observations. The last row of the upper part of Table 2 indicates the explained variance by the first principal component increases substantially with an increasing level of contamination. This is a misleading outcome, since it is only caused by the use of the

Table 2: Second simulation experiment: estimated loadings of the 10 variables on the first two PCs using standard and robust sparse PCA. The reported values are the medians (upper part of the table) and MADs (lower part of the table) over 100 simulation runs. The last line in each part presents the percentage of explained variance EV (median in the upper part, MAD in the lower part). The different columns correspond to the percentage of outliers in the data.

		Standard Estimation						Robust Estimation					
		PC 1			PC 2			PC 1			PC 2		
		0%	10%	20%	0%	10%	20%	0%	10%	20%	0%	10%	20%
Block 1	X_1	0	0	0	0.5	0.14	0.12	0	0	0	0.46	0.34	0.30
	X_2	0	-0.41	-0.42	0.5	0.15	0.14	0	0	0	0.46	0.30	0.26
	X_3	0	0.42	0.42	0.5	0.15	0.15	0	0	0	0.44	0.28	0.25
	X_4	0	0.21	0.2	0.5	0.13	0.13	0	0	0	0.47	0.38	0.33
Block 2	X_5	0.41	0.01	0	0	-0.25	-0.27	0.39	0.33	0.32	0	0	0
	X_6	0.42	0.44	0.43	0	-0.24	-0.28	0.38	0.33	0.36	0	0	0
	X_7	0.41	-0.4	-0.41	0	-0.24	-0.27	0.39	0.23	0.26	0	0	0
	X_8	0.42	0.22	0.21	0	-0.26	-0.28	0.4	0.35	0.34	0	0	0
	X_9	0.39	0.33	0.32	0	-0.33	-0.34	0.31	0.31	0.22	0	0	0
	X_{10}	0.39	-0.3	-0.3	0	-0.33	-0.36	0.3	0.25	0.22	0	0	0
EV (%)		61.4	67.4	80.3	35.7	21.2	13.0	58.9	51.2	50.4	31.9	30.1	28.3
Block 1	X_1	0	0.01	0	0.01	0.15	0.15	0	0	0	0.12	0.24	0.28
	X_2	0	0.03	0.03	0.01	0.16	0.16	0	0.10	0.12	0.10	0.26	0.22
	X_3	0	0.03	0.02	0.01	0.15	0.13	0	0.07	0	0.12	0.40	0.36
	X_4	0	0.03	0.03	0.01	0.15	0.15	0	0.1	0	0.10	0.27	0.31
Block 2	X_5	0.02	0.02	0	0	0.25	0.22	0.01	0.15	0.11	0	0.13	0.06
	X_6	0.02	0.03	0.02	0	0.25	0.17	0.14	0.27	0.25	0	0.17	0.21
	X_7	0.02	0.03	0.03	0	0.36	0.3	0.12	0.27	0.35	0	0.25	0.14
	X_8	0.02	0.03	0.03	0	0.23	0.18	0.12	0.21	0.22	0	0.13	0.18
	X_9	0.04	0.03	0.02	0	0.13	0.1	0.18	0.24	0.32	0	0.1	0.25
	X_{10}	0.03	0.03	0.03	0	0.2	0.13	0.19	0.21	0.33	0	0.11	0.21
EV (%)		10.9	5.75	3.09	9.30	4.90	3.80	12.5	8.90	8.90	11.4	7.10	6.00

sample variance estimator, which is inflated by the outliers. Without outliers, the PCs are no longer representative of the bulk of the data. When using robust sparse PCA, we see that the percentage of explained variance remains about the same when the outliers are added.

6 Real data examples

The method is used for two differently structured data sets. The first example has $n > p$ and shows how the robust method is capable of spotting groups of outliers. The second example points out the method’s applicability on high-dimensional data sets, where $p > n$.

Example 1

The car data set (Kibler *et al.*, 1989) consists of 26 variables containing technical and insurance-related data for 205 different car models. Only continuous variables, and observations without missing values are considered here, resulting in a data set of size 195×26 . To make the scale of the variables comparable, we divide each column of the data matrix by its standard deviation (if $V = \text{Var}$) or by a robust scale measure (if $V = Q_n^2$). Figure 4 gives a scree-plot for non-sparse standard and robust PCA, plotting the explained variance, as defined in (13), versus the number of components. Based on this scree-plot we decide to retain the first four PCs, explaining about 80% of the total (robust) variance, for both approaches.

Figure 5 shows the *tradeoff curve*, discussed in Section 4, plotting the percentage of explained variance as a function of λ . The explained variances are computed over a grid of 100 different values of the tuning parameter λ , ranging from $\lambda = 0$ (no sparseness) up to full sparseness (exactly one non-zero loading per PC). This plot shows how increasing the sparsity parameter leads to a decrease in explained variance. The idea is that the selected λ should be such that the sharpest decline of tradeoff curve occurs afterwards. The selected λ should be close to the end of the first, relatively flat, part of the tradeoff curve. Using the BIC criterion from equation (12), minimized over the same grid of 100 values, we get $\lambda = 2.36$, corresponding to the vertical dashed line in the plot. From the tradeoff curve we conclude that this is an acceptable value. The sharper decline of the tradeoff curve occurs for a tuning parameter larger than 3.

Table 3 shows the resulting loadings for robust non-sparse PCA and robust sparse PCA,

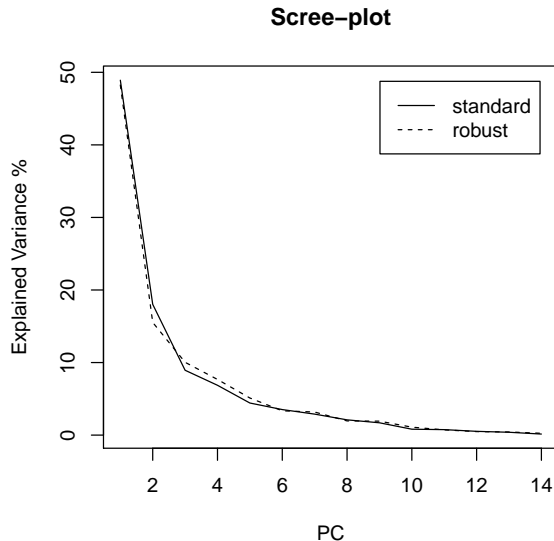


Figure 4: Scree-plots for standard (solid line) and robust (dashed line) PCA ($\lambda = 0$) for the car data set.

derived with $\lambda = 2.36$. While PC1 is somewhat similar in the robust and the sparse robust version, the remaining PCs are quite different. By adding the penalty term in the objective function, the number of non-zero loadings is reduced from 56 to 16, whereas the total amount of explained variance in the first four PCs drops from 81% to 64%. We do find this decrease in explained variance acceptable, given the gained sparsity in the loadings matrix. This could facilitate interpretation, in particular for the higher order principal components. For instance, the fourth principal component is uniquely determined by `peak-rpm`.

Further exploratory data analysis can be done by making distance-distance plots (see Hubert *et al.*, 2002). Such a plot presents two different distance measures: the score distance of each observation in the space of the first k PCs, and the orthogonal distance of each observation to this space. The score distance is a Mahalanobis-like measure of distance of an observation to the center within the PC space. The orthogonal distance describes the orthogonal distance of an observation to the space spanned by the first k PCs. Using cut-off values for both types of distances, outliers can be identified. For details on the construction of these plots, see Hubert *et al.* (2002). Figure 6 shows distance-distance plots for the car data, using standard and robust PCA, and their sparse versions, resulting in four different plots. As before, the first $k = 4$ PCs

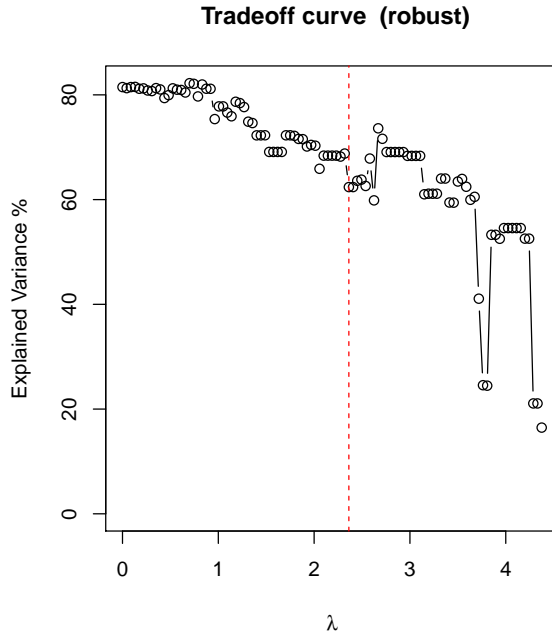


Figure 5: Tradeoff curve for robust sparse PCA computed on the car data set. The dashed line represents the λ selected by the BIC criterion.

are retained, and λ is selected according to the BIC. The robust distance-distance plot (Figure 6b) points out a very distinct outlier group (denoted by symbols \times) which in fact represents all car-models running on diesel. The robust sparse model (Figure 6d) is also able to clearly identify this particular group of outliers. Note that the engine type (diesel/normal) was not one of the variables used in the PCA. These points are identified as outliers because of other characteristics of their engines, which differentiate the groups, and were used in the analysis. In contrast, when considering the standard non-sparse (Figure 6a) and sparse (Figure 6c) distance-distance plots, these outliers cannot be identified, since their presence is masked by the use of a non-robust diagnostic measure. We conclude that in this example only the robust procedure is able to detect the group of outliers, and that adding the sparsity condition did not affected the diagnostic power of the robust distance-distance plot.

Table 3: Loadings of the variables on the first four robust non-sparse ($\lambda = 0$) and robust sparse ($\lambda = 2.36$) PCs of the car data set.

	Robust PCA				Robust sparse PCA			
	PC1	PC2	PC3	PC4	PC1	PC2	PC3	PC4
symboling	-0.03	-0.04	0.03	-0.17	0	0	0	0
wheel-base	0.24	0.25	0.08	0.16	0	0.50	0	0
length	0.29	0.18	-0.05	0.04	0.24	0	0.85	0
width	0.26	0.16	0.14	0.03	0.21	0	0	0
height	0.08	0.39	-0.26	0.32	0	0.87	0	0
curb-weight	0.24	0.13	0.12	0.00	0.32	0	0	0
bore	0.24	0.16	-0.25	0.04	0.21	0	0.03	0
stroke	0.00	-0.24	0.29	-0.58	0	0	0	0
compression-ratio	-0.47	0.61	0.49	-0.11	-0.45	0	0.53	0
horsepower	0.36	-0.01	0.16	-0.20	0.43	0	0	0
peak-rpm	0.08	-0.38	0.60	0.64	0	0	0	1.00
city-mpg	-0.31	0.04	-0.02	0.14	-0.30	0	0	0
highway-mpg	-0.33	0.07	-0.04	0.14	-0.35	0	0	0
price	0.33	0.31	0.34	-0.12	0.40	0	0.06	0
EV %	49.20	15.54	10.12	5.97	45.73	8.32	6.03	4.16
Cumulative EV %	49.20	64.74	74.85	80.82	45.73	54.05	60.08	64.24

Example 2

The yarn data set (see Swierenga *et al.*, 1999) contains near-infrared (NIR) spectra of 21 PET yarns of different density. The data are available in the R package “pls” as data frame “yarn”. This data set was used in Izenman (2008) and ter Braak (2009) for a prediction exercise, but we focus here on the use of PCA as an exploratory technique. A total of 268 different wavelengths were measured, yielding a data set of size 21×268 . As the algorithm discussed in Section 3 computes one (sparse) PC at a time and may stop after computing the k th component, it is especially useful in high-dimensional applications, where the actual information is restricted to a comparatively low-dimensional subspace. Due to this characteristic, computation time can be reduced tremendously, as in such settings usually only a few PCs are important. In the data set $k = 2$ PCs already explain more than 85% of the total (robust) variance, thus the iteration can be stopped after obtaining the first two principal components, rather than computing all $\min(n, p)$

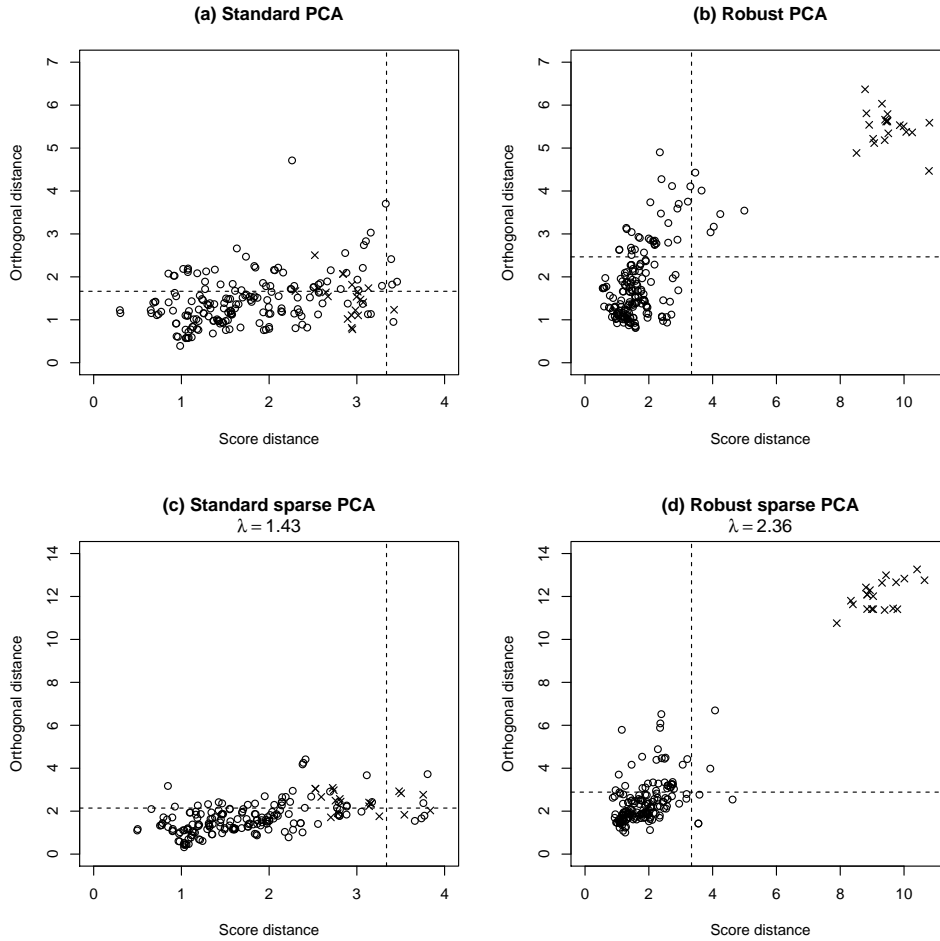


Figure 6: Distance-distance plots for standard and robust PCA and their sparse versions. In the robust plots vehicles running on diesel (×) are clearly distinguishable from vehicles using gasoline (○).

loadings vectors. In this particular example this reduces computation time by 90% (from 41 to 4 seconds for standard and from 135 to 13 seconds for robust PCA on an AMD Athlon x64 X2 4200+ running at 2.2GHz).

Figure 7 shows the spectral lines of the 21 observations (black). Three spectral intervals A, B and C are marked, as the variables in these areas show a higher variance than in other regions. In interval B the single yarns are grouped together to 5 “clusters”, whereas in region A and C this pattern cannot be observed and the yarns are more homogeneously structured. We add three outlying spectra (see Figure 7, in grey) in order to test the algorithm’s robustness properties in

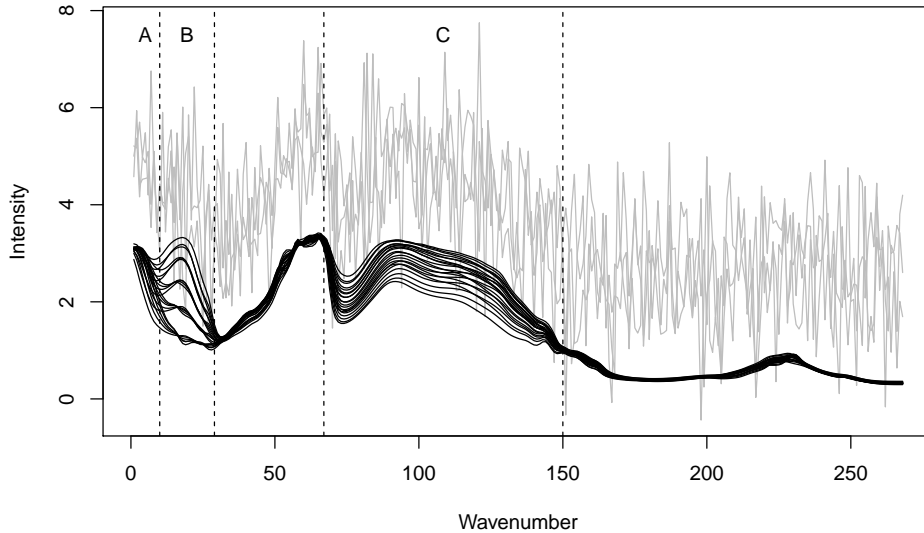


Figure 7: The yarn data set. The NIR spectrum of 21 different PET yarns (black), with intensity measured at 268 wavelengths. Outliers (in grey) were added for challenging the robust sparse PCA estimator.

high-dimensional scenarios. The outliers are different in location, but also have a larger variability, to make the robustness exercise more challenging.

We start by selecting an appropriate value for the number k of PCs to retain. The screeplot in Figure 8 conforms that $k = 2$ is a good choice, explaining most of the (robust) variance. Note that the large value for EV_1 for the standard method is mainly due to the fact that the sample variance is inflated by the outliers. The screeplot for standard PCA on the data set without the outlying spectra does resemble the behavior of the robust version (dashed line in Figure 8). Then, we use the tradeoff curve in Figure 9 for selecting a value of λ keeping a sufficiently large percentage of explained variance. For robust PCA we take $\lambda = 16$, a value at the end of the flat part of the curve and well before the sharp decrease in the tradeoff curve. For that value of λ we explain still 85% of the robust variance. The BIC criterion gives us a value of 19.55, which is not that different, but leads to an unacceptable value of 0% of explained robust variance. For standard PCA we take $\lambda = 12.77$ explaining 75% the total variance.

Figure 10 shows the loadings of the 268 variables, labeled with wavenumbers one to 268 for

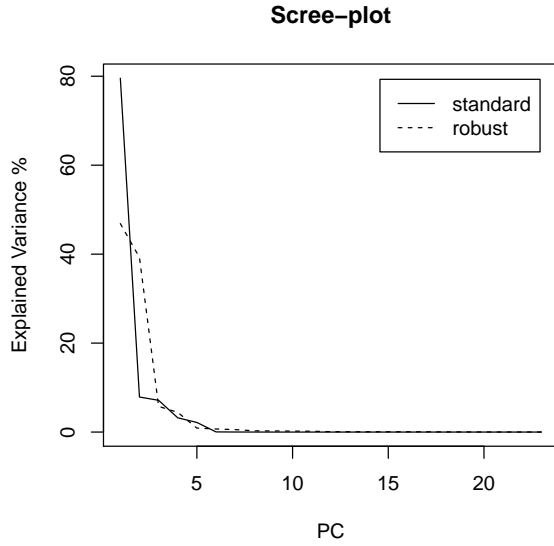


Figure 8: Scree-plots for a standard (solid line) and robust (dashed line) PCA ($\lambda = 0$) for the yarn data set.

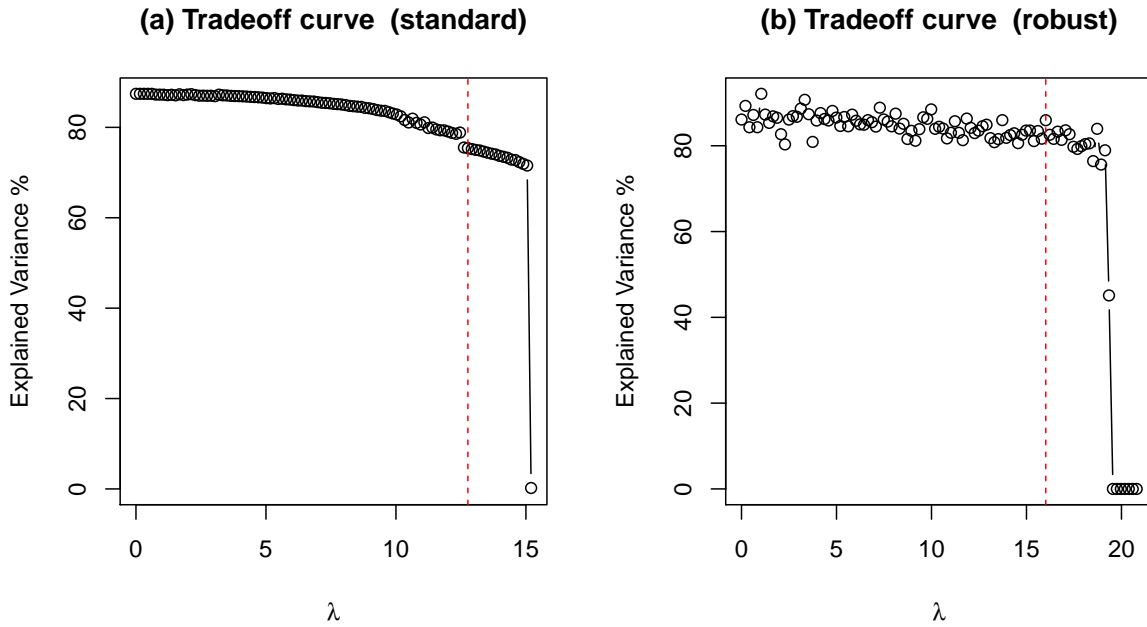


Figure 9: Tradeoff curves for standard and robust sparse PCA computed on the yarn data set. The dashed lines represent the selected value of the tuning parameter λ .

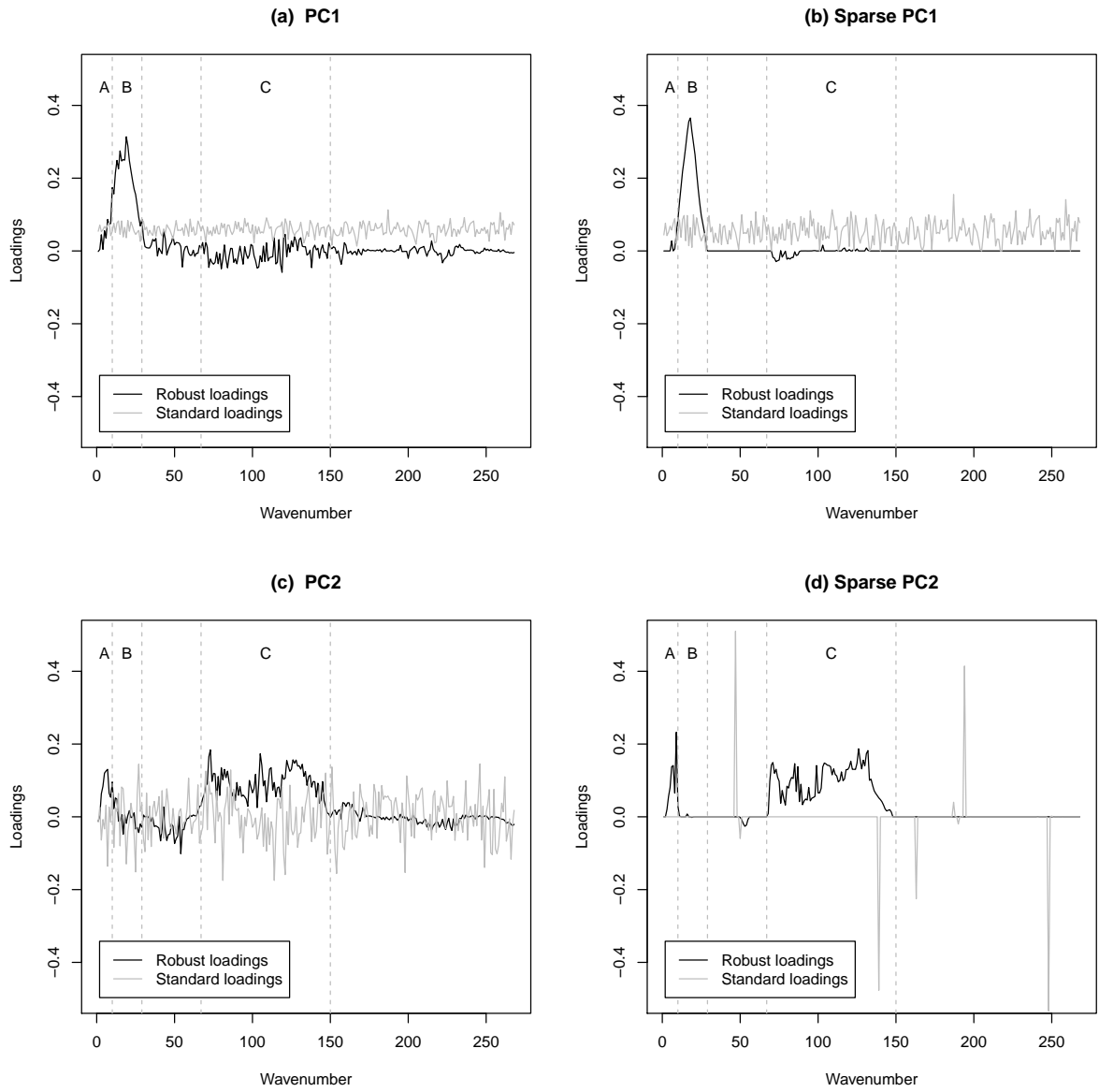


Figure 10: Loadings of the 268 variables on the first two principal components using standard (grey) and robust (black) PCA. Results are given for both sparse (right) and non-sparse (left) PCA for the yarn data set.

standard and robust PCA, and their sparse version. For standard PCA, the loadings in general do not seem to contain any interpretable structure and are heavily influenced by the outliers. The first standard sparse PC (panel b, grey line), hardly contains any zeros, whereas the second (panel d, grey line) only contains 11 non-zero loadings. However, this second sparse standard PC does not

point out specific spectral ranges, but is mainly made up of single spikes, describing the outlier’s random pattern. In contrast to this, robust PCA shows distinct features in all four plots. The first non-sparse robust PC (panel a, solid line) points out a peak at the spectrum’s lower end (region B). This peak is even more clearly detected by the robust sparse model (panel b, black line) and corresponds to the spectral range B in which the yarns reveal a rather “clustered” structure. Most of the loadings outside of the interval B are reduced to zero, illustrating that a sparse approach make interpretation easier. The second robust PC (panel c, black line) is mainly made up of the wavelengths in spectral ranges A and C, corresponding to the wavelengths with high variability but without “cluster structure” among the yarns. Wavelengths outside of these intervals A and C contribute less to the second PC, as their (absolute) loadings are quite low. The loadings of the second sparse robust PC (panel d, black line) do even better in separating the wavelengths in intervals A or C from the others; almost all loadings outside of these ranges are exactly equal to zero. As we can see from the tradeoff curve in Figure 9 (b), the sparse robust solution only explains 1% less variance than the non-sparse ($\lambda = 0$), whereas the number of non-zero loadings decreases from $2 \times 268 = 536$ to 159. Despite the noise added by the three outlying spectra, the robust sparse method is capable of finding distinct structures in the data.

Finally, we mention that in this example there is a natural ordering in the variables, suggesting that adding *fusion penalties* in the objective function, as in Guo *et al.* (2010), might further improve the performance of the sparse PCA method.

7 Concluding remarks

Sparse PCA delivers components that can be considered as a compromise between maximizing variance and simplifying interpretability. Robust sparse PCA keeps the goal of simple interpretability, but the determination of the PCA directions is not affected by outlying observations. The proposed approach is based on the idea of projection-pursuit, maximizing a robust variance for finding the directions. Projection-pursuit based PCA has the further advantage that the components are extracted sequentially, making it possible to stop the algorithm after a desired number of components. This is especially attractive for the analysis of high dimensional data with possibly fewer observations than variables.

The optimal level of the tuning parameter λ , optimal in terms of both interpretability and

explained variance, can be determined by an information criterion like the BIC. The simulations and the data examples have demonstrated that the robust sparse PCs are resistant with respect to data outliers, and that the resulting sparsity patterns are useful. The tradeoff curve, visualizing the tradeoff between explained variance and sparsity, can be used as an exploratory tool for obtaining more guidance on an optimal sparsity level. A complete implementation of the algorithm is available in the R package `pcaPP` (Filzmoser *et al.*, 2010).

There are several questions we did not address and which are left for future research. For instance, one could think of a joint selection criterion for the number of principal components and the tuning parameter λ , as opposed to the two-step approach followed in this paper. While we did study in this paper the robustness of the method with respect to outliers, the stability of the method with respect to small changes in the data had not been investigated. It would be interesting to compute influence functions of the considered estimators. Xu *et al.* (2012) state that instability to small changes in the data is the price to pay to get a sparse estimator.

Another limitation of the paper is that we only considered the L_1 norm in the constraint on the loadings. In regression analysis one frequently uses the L_2 norm, e.g. Maronna (2011) for regularized robust regression, but this does not lead to sparse solutions. On the other hand, using the L_0 norm, counting the number of non-zero components of a vector, does yield sparsity (see Farcomeni, 2009). The method proposed in this paper can also be applied to get a robust version of the latter approach. It remains to be seen whether the L_0 penalty is a better choice than the more standard L_1 penalty considered in this paper. Another generalization would be to add a supplementary penalty on the norm of the score vectors, given in (3), to get both sparse loadings coefficients and score vectors, as in Witten *et al.* (2009). This would yield a sparse variant of robust low-rank approximations of a data matrix, as in Maronna and Yohai (2008).

A naive approach to robust sparse PCA would be to estimate a sparse robust covariance matrix, and then compute the eigenvectors of it. While sparse robust covariance matrices have been proposed recently (Croux *et al.*, 2010b), this is not a useful approach since the eigenvectors will not inherit the sparsity of the matrix. A projection-pursuit approach, as undertaken in this paper, avoids this pitfall. Projection-pursuit approaches to sparse discriminant analysis and sparse canonical correlation analysis have recently proposed (see Witten and Tibshirani, 2011; Lykou and Whittaker, 2010), and can be made robust along similar lines as outlined in this paper.

Acknowledgment

The authors are grateful to the chief editor, an associate editor, and the reviewers for their helpful comments and suggestions.

References

- Anaya-Izquierdo, K., Critchley, F., and Vines, K. (2011). Orthogonal simple component analysis: a new, exploratory approach. *Annals of Applied Statistics*, **5**(1), 486–522.
- Bien, J., Xu, Y., and Mahoney, M. (2010). CUR from a sparse optimization viewpoint. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 23, pages 217–225.
- Candès, E., Li, X., Ma, Y., and Wright, J. (2009). Robust principal component analysis? *Journal of the ACM*, **58**(1), 1–37.
- Cattell, R. (1966). The scree test for the number of factors. *Multivariate Behaviour Research*, **1**, 245–276.
- Chipman, H. A. and Gu, H. (2005). Interpretable dimension reduction. *Journal of Applied Statistics*, **32**, 969–987.
- Croux, C. and Haesbroeck, G. (2000). Principal components analysis based on robust estimators of the covariance or correlation matrix: Influence functions and efficiencies. *Biometrika*, **87**, 603–618.
- Croux, C. and Ruiz-Gazen, A. (2005). High breakdown estimators for principal components: The projection-pursuit approach revisited. *Journal of Multivariate Analysis*, **95**, 206–226.
- Croux, C., Filzmoser, P., and Oliveira, M. (2007). Algorithms for projection-pursuit robust principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, **87**, 218–225.
- Croux, C., Dehon, C., and Yadine, A. (2010a). The k-step spatial sign covariance matrix. *Advances in Data Analysis and Classification*, **4**(2), 137–150.

- Croux, C., Gelper, S., and Haesbroeck, G. (2010b). Robust scatter regularization. In G. Saporta and Y. Lechevallier, editors, *Compstat 2010, Book of Abstracts*, page 138, Paris. Conservatoire National des Arts et Métiers (CNAM) and the French National Institute for Research in Computer Science and Control (INRIA).
- Farcomeni, A. (2009). An exact approach to sparse principal component analysis. *Computational Statistics*, **24**(4), 583–604.
- Filzmoser, P. (1999). Robust principal components and factor analysis in the geostatistical treatment of environmental data. *Environmetrics*, **10**, 363–375.
- Filzmoser, P., Fritz, H., and Kalcher, K. (2010). *pcaPP: Robust PCA by Projection Pursuit*. R package version 1.9-0.
- Guo, J., James, G., Levina, E., Michailidis, G., and Zhu, J. (2010). Principal component analysis with sparse fused loadings. *Journal of Computational and Graphical Statistics*, **19**, 930–946.
- Hubert, M., Rousseeuw, P. J., and Verboven, S. (2002). A fast method for principal components with application to chemometrics. *Chemometrics and Intelligent Laboratory Systems*, **60**, 101–111.
- Hubert, M., Rousseeuw, P. J., and Vanden Branden, K. (2005). Robpca: A new approach to robust principal component analysis. *Technometrics*, **47**, 64–79.
- Hubert, M., Rousseeuw, P., and Van Aelst, S. (2008). High-breakdown robust multivariate methods. *Statistical Science*, **23**(1), 92–119.
- Izenman, A. (2008). *Modern Multivariate Statistical Techniques: Regression, Classification, and Manifold Learning*. Springer-Verlag, New York.
- Jenatton, R., Obozinski, G., and Bach, F. (2009). Structured sparse principal component analysis. Technical report, arXiv:0909.1440.
- Jolliffe, I. T. (1995). Rotation of principal components: choice of normalization constraints. *Journal of Applied Statistics*, **22**, 29–35.
- Jolliffe, I. T. (2002). *Principal Component Analysis*. Springer-Verlag, New York, second edition.

- Jolliffe, I. T., Trendafilov, N. T., and Uddin, M. (2003). A modified principal component technique based on the lasso. *Journal of Computational and Graphical Statistics*, **12**, 531–547.
- Journée, M., Nesterov, Y., Richtárik, P., and Sepulchre, R. (2010). Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, **11**, 517–553.
- Kibler, D., Aha, D., and Albert, M. (1989). Instance-based prediction of real-valued attributes. *Computational Intelligence*, **5**, 51–57.
- Leng, C. and Wang, H. (2009). On general adaptive sparse principal component analysis. *Journal of Computational and Graphical Statistics*, **18**(1), 201–215.
- Li, G. and Chen, Z. (1985). Projection-pursuit approach to robust dispersion matrices and principal components: Primary theory and Monte Carlo. *Journal of the American Statistical Association*, **80**(391), 759–766.
- Locantore, N., Marron, J., Simpson, D., Tripoli, N., Zhang, J., and Cohen, K. (1999). Robust principal components for functional data. *Test*, **8**, 1–73.
- Lykou, A. and Whittaker, J. (2010). Sparse CCA using a Lasso with positivity constraints. *Computational Statistics & Data Analysis*, **54**(12), 3144–3157.
- Maronna, R. (2005). Principal components and orthogonal regression based on robust scales. *Technometrics*, **47**(3), 264–273.
- Maronna, R. (2011). Robust ridge regression for high-dimensional data. *Technometrics*, **53**(1), 44–53.
- Maronna, R. and Yohai, V. (2008). Robust low-rank approximation of data matrices with elementwise contamination. *Technometrics*, **50**(3), 295–304.
- Oja, H. (2010). *Multivariate Nonparametric Methods with R*. Springer-Verlag, New York.
- Rousseeuw, P. and Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of the American Statistical Association*, **88**(424), 1273–1283.

- Swierenga, H., de Weijer, A. P., van Wijk, R. J., and Buydens, L. M. C. (1999). Strategy for constructing robust multivariate calibration models. *Chemometrics and Intelligent Laboratory Systems*, **49**, 1–17.
- ter Braak, C. (2009). Regression by L1 regularization of smart contrasts and sums (roscas) beats pls and elastic net in latent variable model. *Journal of Chemometrics*, **23**(5), 217–228.
- Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, **58**, 267–288.
- Trendafilov, N. T. and Jolliffe, I. T. (2006). Projected gradient approach to the numerical solution of the scotlass. *Computational Statistics & Data Analysis*, **50**(1), 242–253.
- Tseng, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimization. *Journal of Optimization Theory and Applications*, **109**(3), 475–494.
- Vines, S. (2000). Simple principal components. *Applied Statistics*, **49**, 441–451.
- Witten, D. and Tibshirani, R. (2011). Penalized classification using Fisher’s linear discriminant. *Journal of the Royal Statistical Society, Series B*. In press.
- Witten, D., Tibshirani, R., and Hastie, T. (2009). A penalized matrix decomposition, with applications to sparse principal components and canonical correlation analysis. *Biostatistics*, **10**(3), 515–534.
- Xu, H., Caramanis, C., and Sanghavi, S. (2010). Robust PCA via outlier pursuit. In J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta, editors, *Advances in Neural Information Processing Systems*, volume 23, pages 2496–2504.
- Xu, H., Caramanis, C., and Mannor, S. (2012). Sparse algorithms are not stable: A no-free-lunch theorem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **34**(1), 187–193.
- Zou, H., Hastie, T., and Tibshirani, R. (2006). Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, **15**(2), 265–286.