

# **Robust sparse PCA based on projection-pursuit**

**Peter Filzmoser**

**Department of Statistics and Probability Theory  
Vienna University of Technology, Austria**

*Graz, Austria*

September 8, 2011



---

Vienna University of Technology

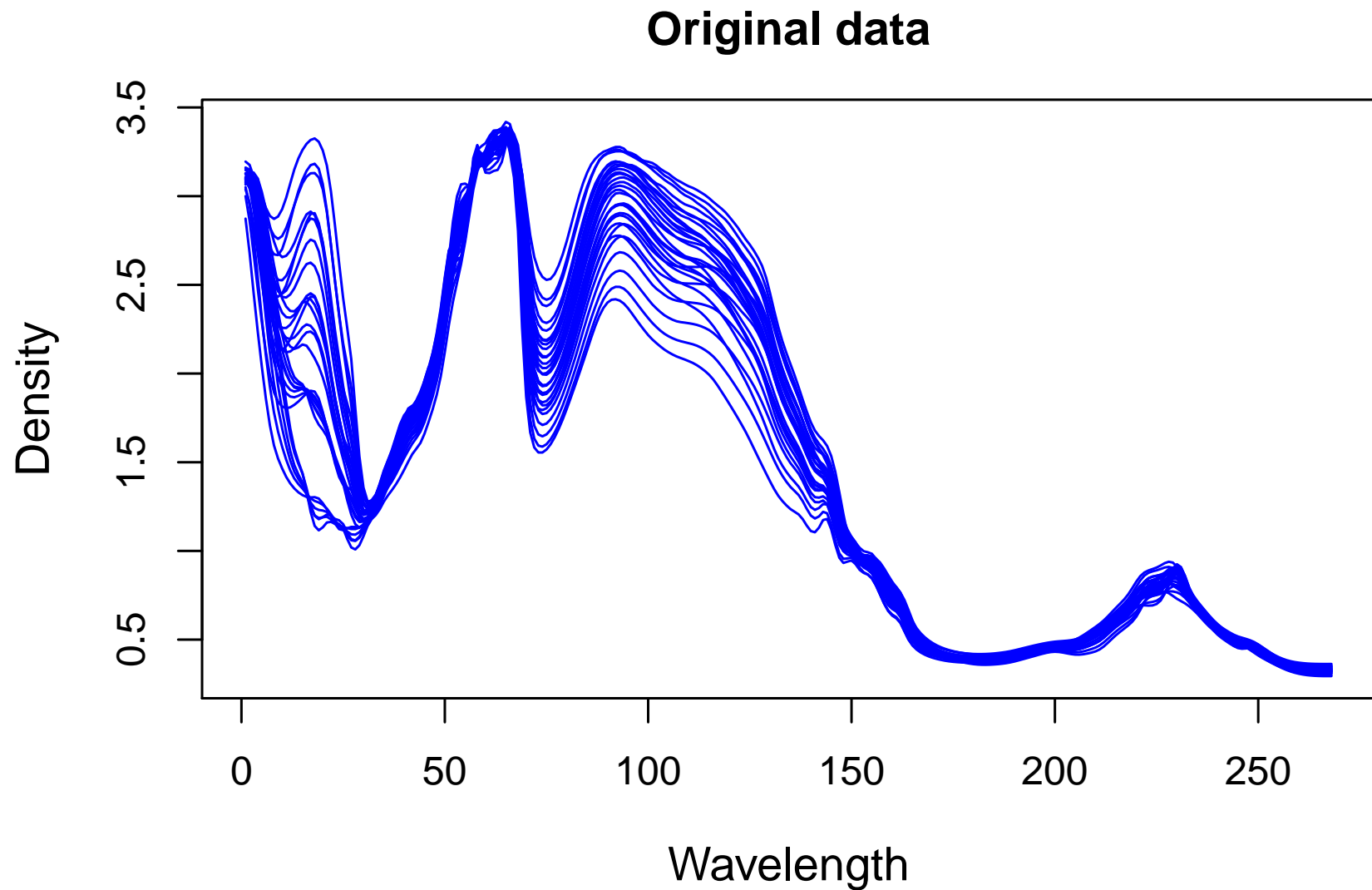
**Christophe Croux**, K.U. Leuven, Belgium

**Heinrich Fritz**, Vienna University of Technology

- **Projection-pursuit PCA**
- **PCA and sparsity**
- **Algorithm**
- **Examples**
- **Summary**

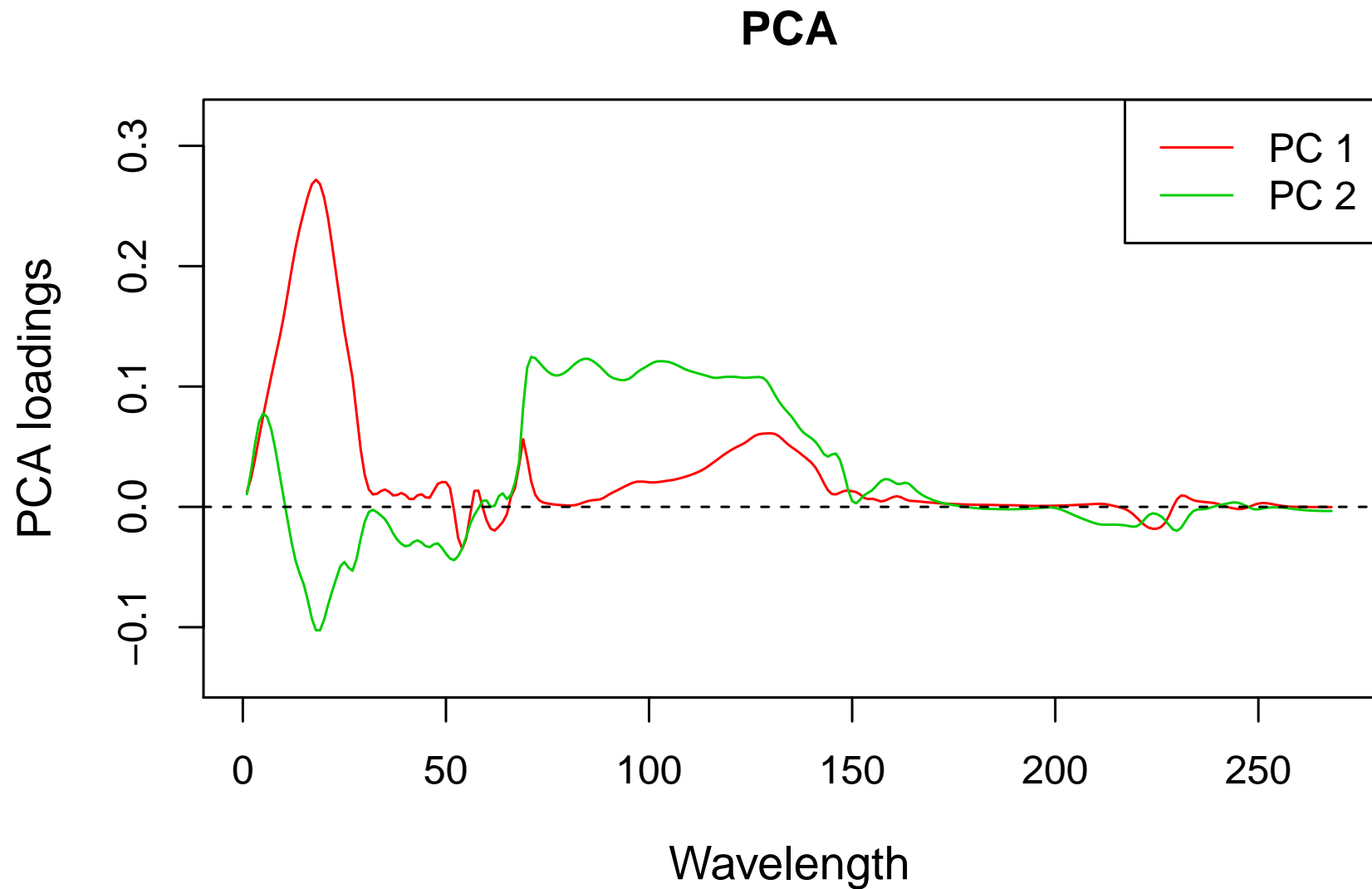
# Motivating example

Data set “yarn”: density of 28 pet yarns, measured at 268 wavelength



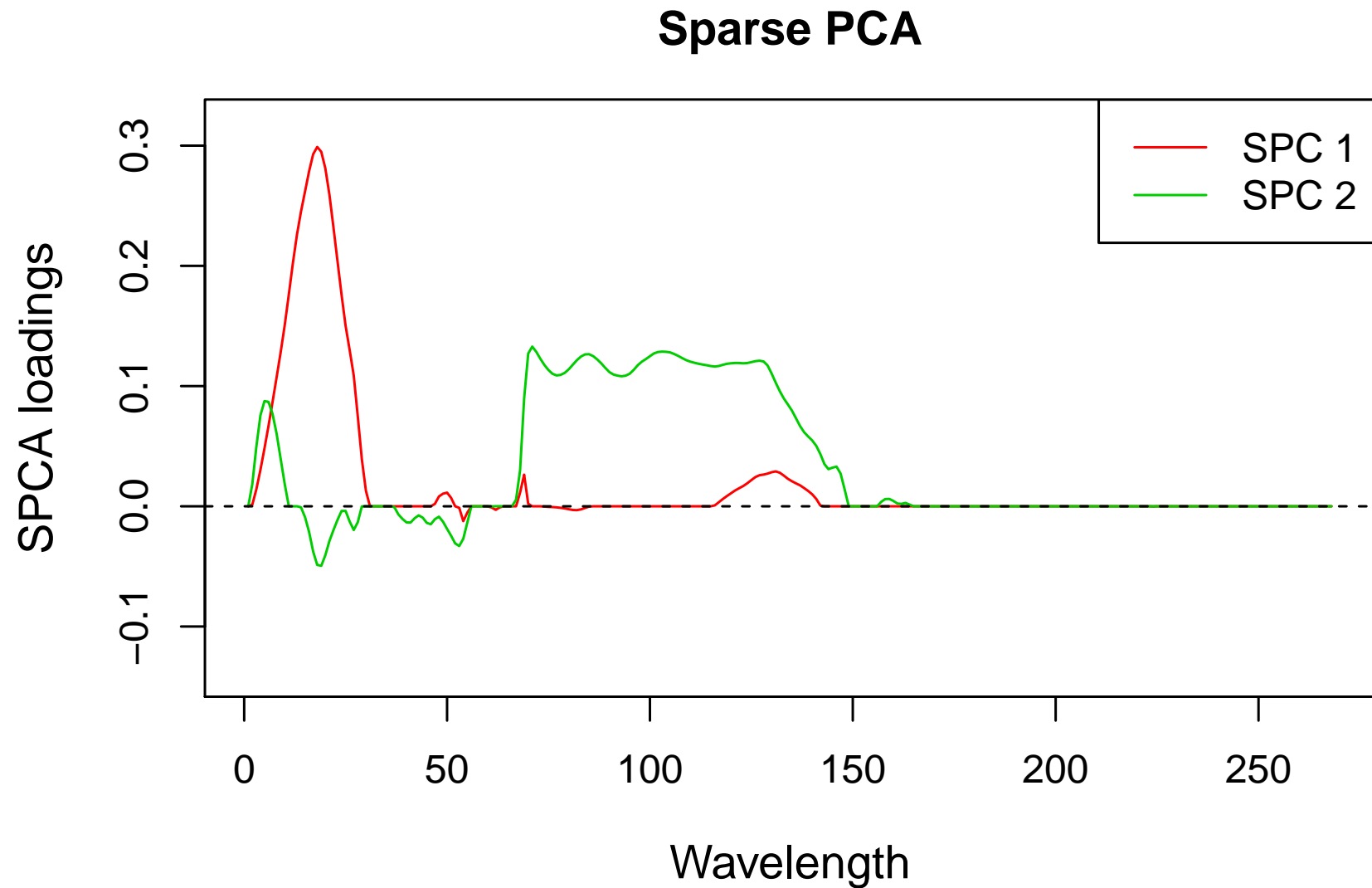
# Motivating example

First two PCA loadings:



# Motivating example

First two Sparse PCA loadings:



Given:  $n$  observations  $x_1, \dots, x_n \in \mathbb{R}^p$ , collected in the rows of  $X$ .  
The first PCA direction is given by

$$a_1 = \operatorname{argmax}_{\|a\|=1} V(a^t x_1, \dots, a^t x_n).$$

$V$  is a variance measure:

- **classical case:**  $V$  is the empirical variance (Var)  
 $\implies a_1$  corresponds to the first eigenvector of the sample covariance matrix.
- **robust case:**  $V$  is squared MAD or squared  $Q_n$  estimator (Rousseeuw and Croux, 1993)  
 $\implies a_1$  is the direction of the first **robust** PC.

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

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

**Loadings matrix** for the first  $k$  PCs:  $\mathbf{A}_k = (\mathbf{a}_1, \dots, \mathbf{a}_k)$ ,

**Scores matrix of the first  $k$  PCs:**  $\mathbf{Z}_k = \mathbf{X} \mathbf{A}_k$  with elements  $z_{ij} = \mathbf{a}_j^t \mathbf{x}_i$



...add an  $L_1$  penalty in the objective function (see Tibshirani, 1996):

**SCoTLASS criterion** (Jolliffe et al., 2003):

$j$ th sparse PCA direction ( $1 \leq j \leq p$ ):

$$\max_{\|a\|=1, a \perp a_1, \dots, a \perp a_{j-1}} a^t \hat{\Sigma} a, \quad \text{subject to } \|a\|_1 \leq t,$$

with  $L_1$  norm  $\|a\|_1 = \sum_{j=1}^p |a_j|$ , and  $\hat{\Sigma}$  the empirical covariance matrix.

**Equivalent formulation:**

$$\max_{\|a\|=1, a \perp a_1, \dots, a \perp a_{j-1}} a^t \hat{\Sigma} a - \lambda_1 \|a\|_1,$$

with the **tuning parameter**  $\lambda_1$ .

- $L_1$  penalty forces some of the loadings to become exactly zero;
- $\hat{\Sigma}$  necessary as input  $\longrightarrow$  robust estimation???  $p > n$ ???

First sparse PCA direction:

$$\tilde{a}_1 = \operatorname{argmax}_{\|a\|=1} V(a^t x_1, \dots, a^t x_n) - \lambda_1 \|a\|_1.$$

- Setting  $\lambda_1 = 0 \longrightarrow$  **unconstrained** first PCA direction  $a_1$ ;
- increasing  $\lambda_1 \longrightarrow$  **sparsity gains importance** compared to (robust) variance maximization.

*j*th sparse PCA direction ( $1 < j \leq p$ ):

$$\tilde{a}_j = \underset{\|a\|=1, a \perp \tilde{a}_1, \dots, a \perp \tilde{a}_{j-1}}{\operatorname{argmax}} V(a^t x_1, \dots, a^t x_n) - \lambda_j \|a\|_1$$

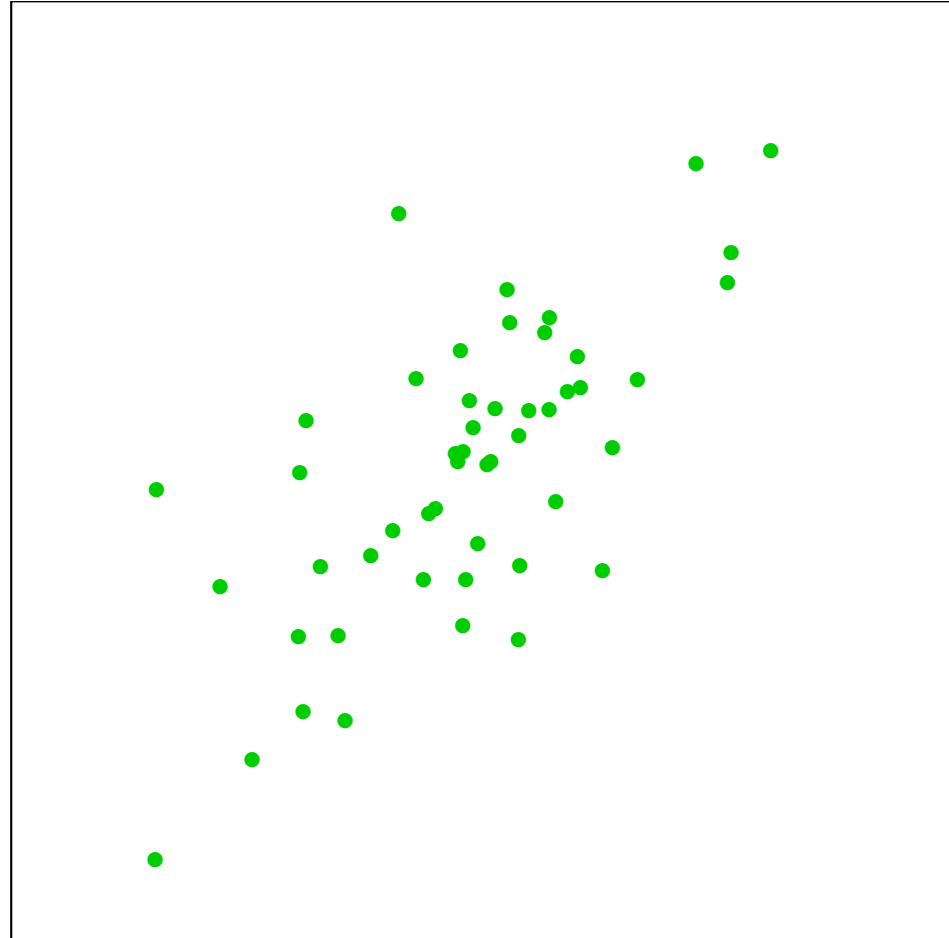
$\lambda_j$  is a **tuning parameter**, possibly different from  $\lambda_1$ .

If  $V = \text{Var}$   $\longrightarrow$  SCoTLASS and projection-pursuit definitions are the same.

**Projection-pursuit approach:** PCs are computed sequentially; stop at a desired number  $k < p$ .

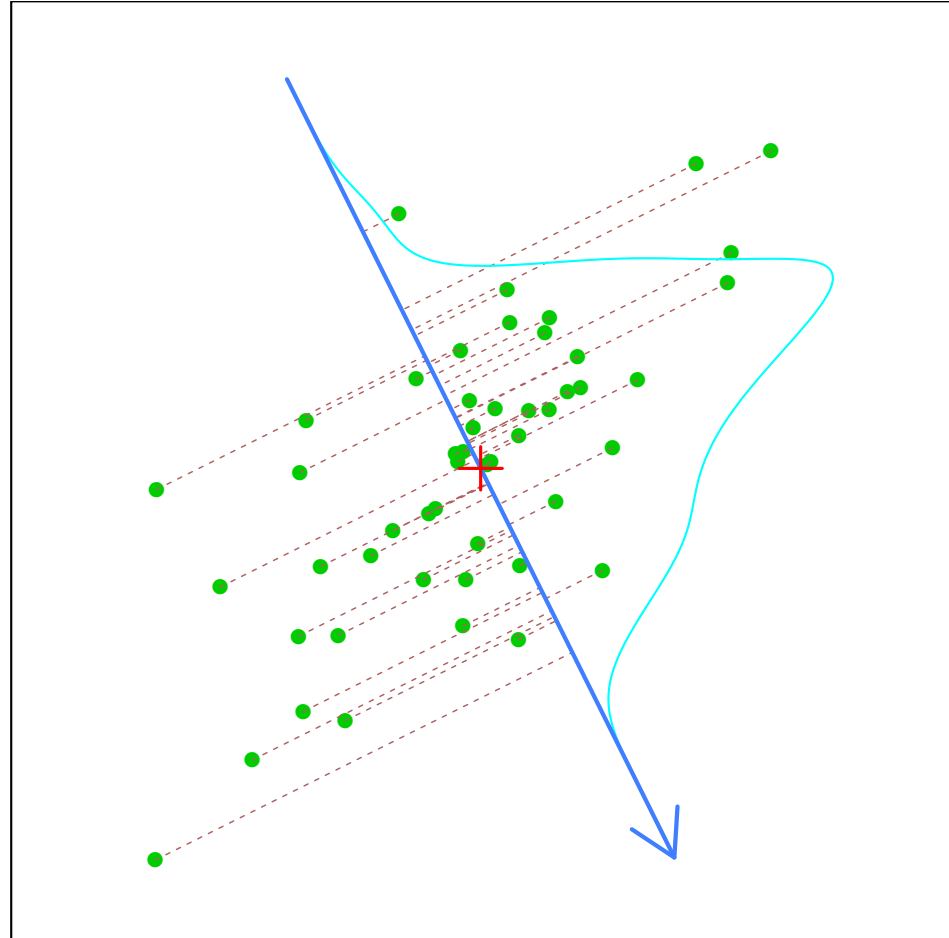
# Projection-pursuit algorithm

- project data on a direction



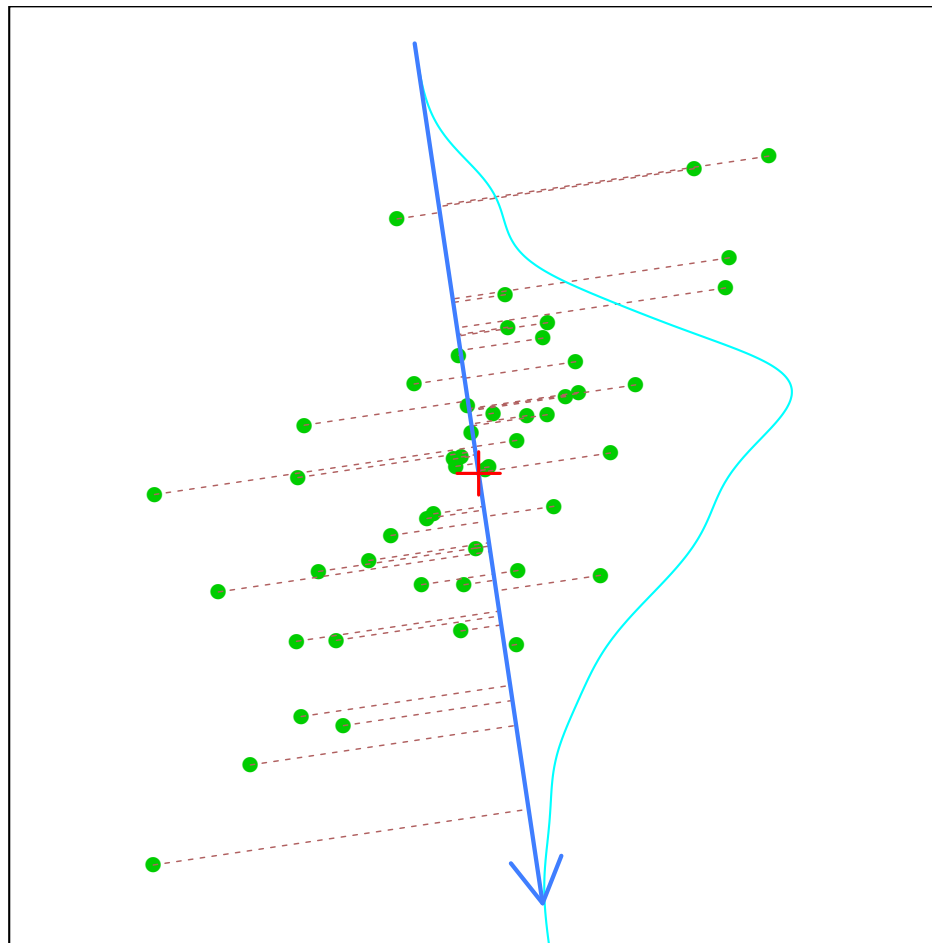
# Projection-pursuit algorithm

- project data on a direction
- compute (robust) variance



# Projection-pursuit algorithm

- project data on a direction
- compute (robust) variance
- repeat for “many” directions

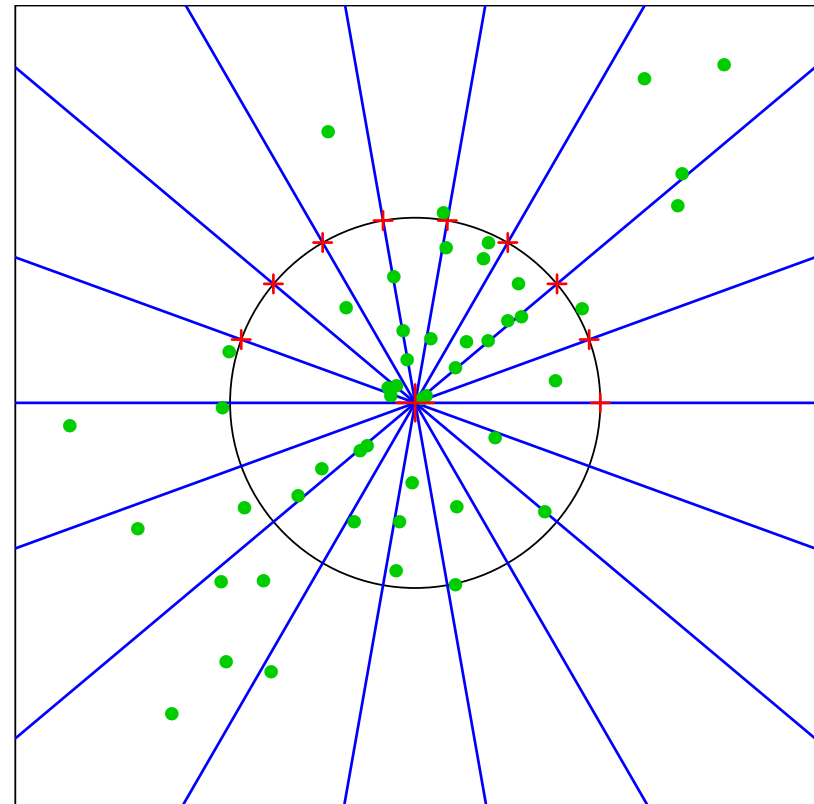


Optimization is done on a regular grid in the plane:

- select  $N_g$  regular grid points on the unit circle
- consider  $N_g$  candidate directions of each grid point through the center
- maximize the variance of the projected data, measured by a measure  $V$

$\Rightarrow$  direction which gives maximum (robust) variance

(Croux, Filzmoser, Oliveira, 2007)



Assume that the first  $j - 1$  sparse PCA directions  $\tilde{a}_{j-1}$  are already obtained and collected in the columns of  $\tilde{A}_{j-1}$ , with  $1 \leq j \leq k - 1$ .

Compute  $\tilde{a}_j$ :

- $\tilde{A}_{j-1}^\perp$  contains orthonormal basis of the subspace orthogonal to  $\tilde{A}_{j-1}$ .
- Project data into this space by  $x_i^{(j-1)} = (\tilde{A}_{j-1}^\perp)^t x_i$ , for  $i = 1, \dots, n$ .
- Use objective function

$$f(a) = V(a^t x_1^{(j-1)}, \dots, a^t x_n^{(j-1)}) - \lambda_j \|\tilde{A}_{j-1}^\perp a\|_1$$

with  $\|a\| = 1$ .

- Cycle through all **pairs of variables** using the **Grid algorithm**; update optimal direction iteratively.



# Selection of the tuning parameter

The tuning parameter  $\lambda_j$  regulates the degree of sparseness.

We want to have similar degree of sparseness in the different PCs

$$\implies \text{take } \lambda_j := \lambda \mathcal{V}(\mathbf{X}^{(j)}),$$

where  $\mathbf{X}^{(j)} = \mathbf{X} \tilde{\mathbf{A}}_{j-1}^\perp$ , and  $\mathcal{V}$  denotes the total robust variance, defined for any matrix  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_p)$  as

$$\mathcal{V}(\mathbf{Y}) = \sum_{i=1}^p V(\mathbf{y}_i).$$

## BIC type criterion:

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

with

- $\widetilde{\text{RV}} = \mathcal{V}(\mathbf{X} - \mathbf{X} \tilde{\mathbf{A}}_k \tilde{\mathbf{A}}_k^t)$  ... total robust variance of residuals from **sparse PCA**
- $\text{RV} = \mathcal{V}(\mathbf{X} - \mathbf{X} \mathbf{A}_k \mathbf{A}_k^t)$  ... total robust variance of residuals from **unconstrained PCA**
- $\text{df}(\lambda)$  ... number of **non-zero loadings** when using  $\lambda$ .

Select  $\lambda$  by minimizing  $\text{BIC}(\lambda)$  over a grid  $[0, \lambda_{max}]$ , where  $\lambda_{max}$  results in full sparseness.

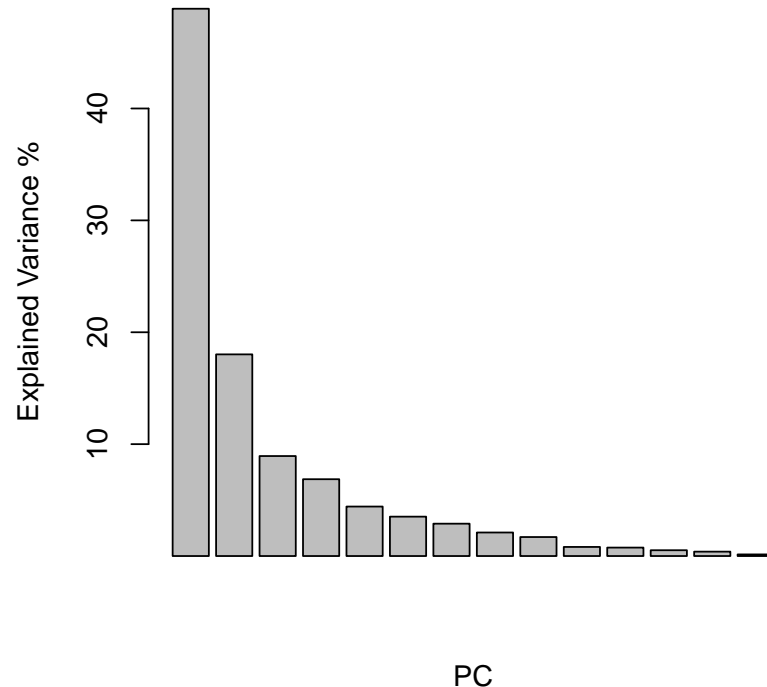
# Example car data set

**Car data** (Kibler et al., 1989): 205 different car models; 26 variables containing technical and insurance-related data

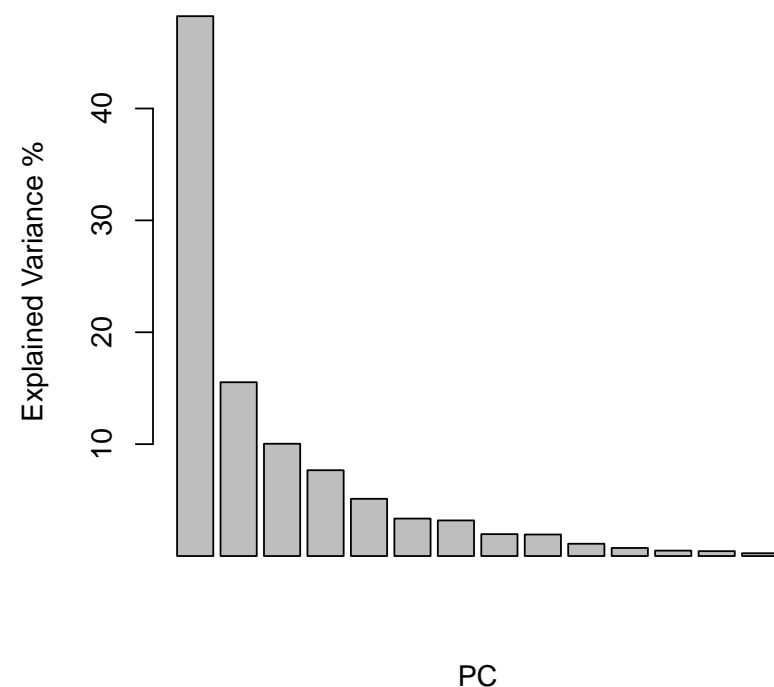
→ only continuous variables, no missings: data set of size  $195 \times 14$

Compare classical ( $V = \text{Var}$ ) and robust ( $V = Q_n^2$ ) approach.

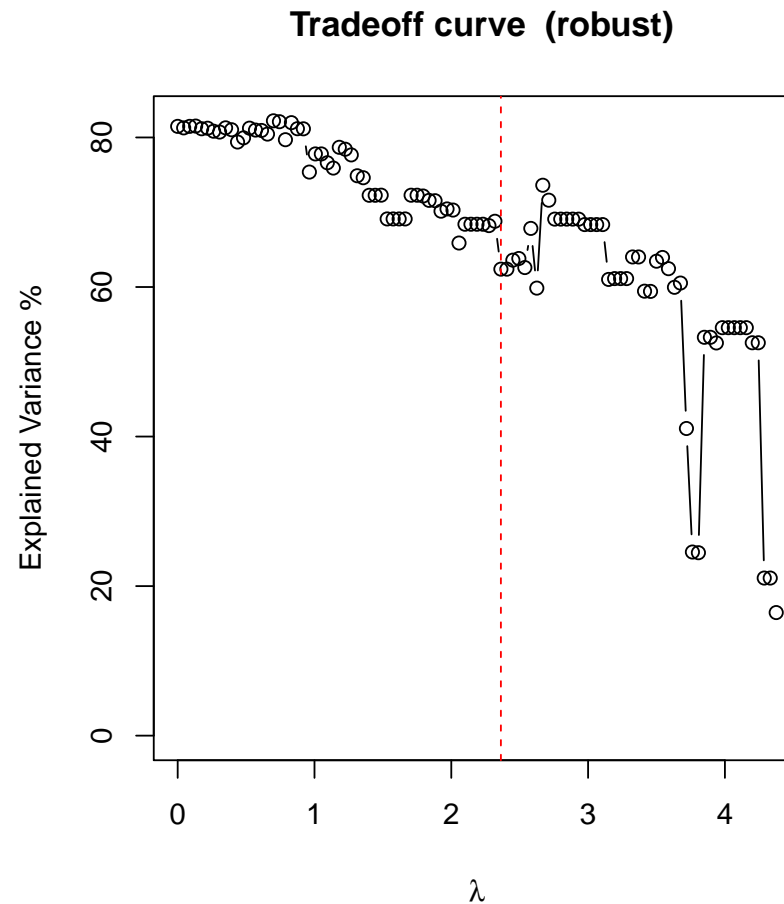
(a) Standard Scree-plot



(b) Robust Scree-plot



“Tradeoff curve” for sparse robust PCA:



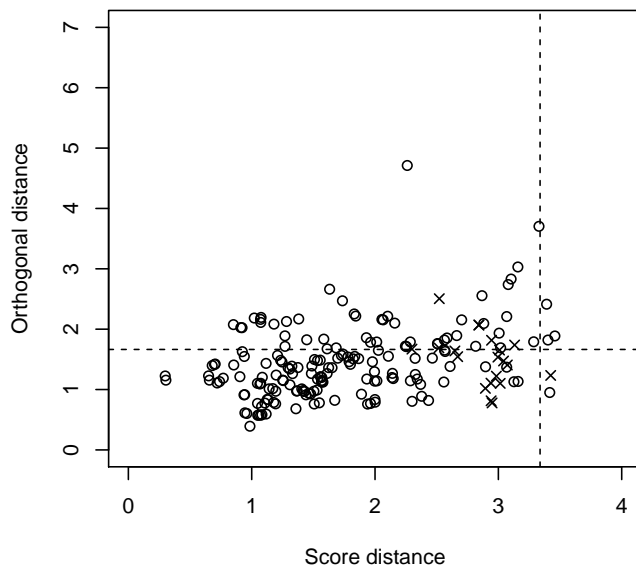
BIC criterion  $\Rightarrow \lambda = 2.36$

# Example car data set

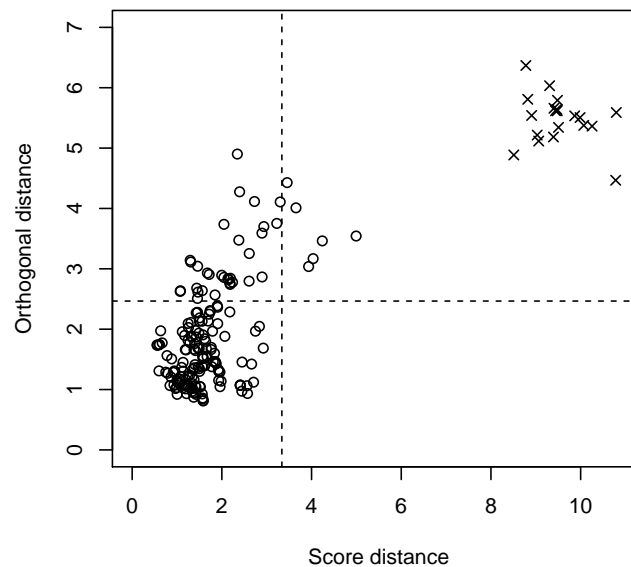
	Robust PCA ( $\lambda = 0$ )				Robust sparse PCA ( $\lambda = 2.36$ )			
	PC1	PC2	PC3	PC4	PC1	PC2	PC3	PC4
symboling	-0.03	-0.04	0.03	-0.17	0	0	0	0
wheel-base	0.24	0.25	0.08	0.16	0	0.50	0	0
length	0.29	0.18	-0.05	0.04	0.24	0	0.85	0
width	0.26	0.16	0.14	0.03	0.21	0	0	0
height	0.08	0.39	-0.26	0.32	0	0.87	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 car data set

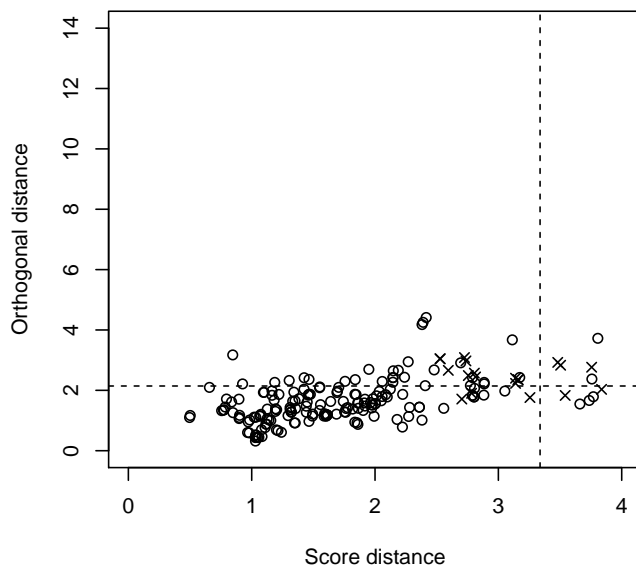
(a) Standard PCA



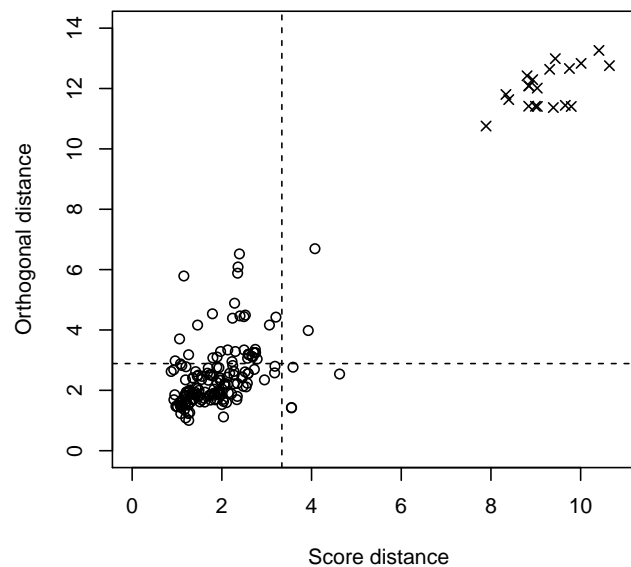
(b) Robust PCA



(c) Standard sparse PCA  
 $\lambda = 1.43$

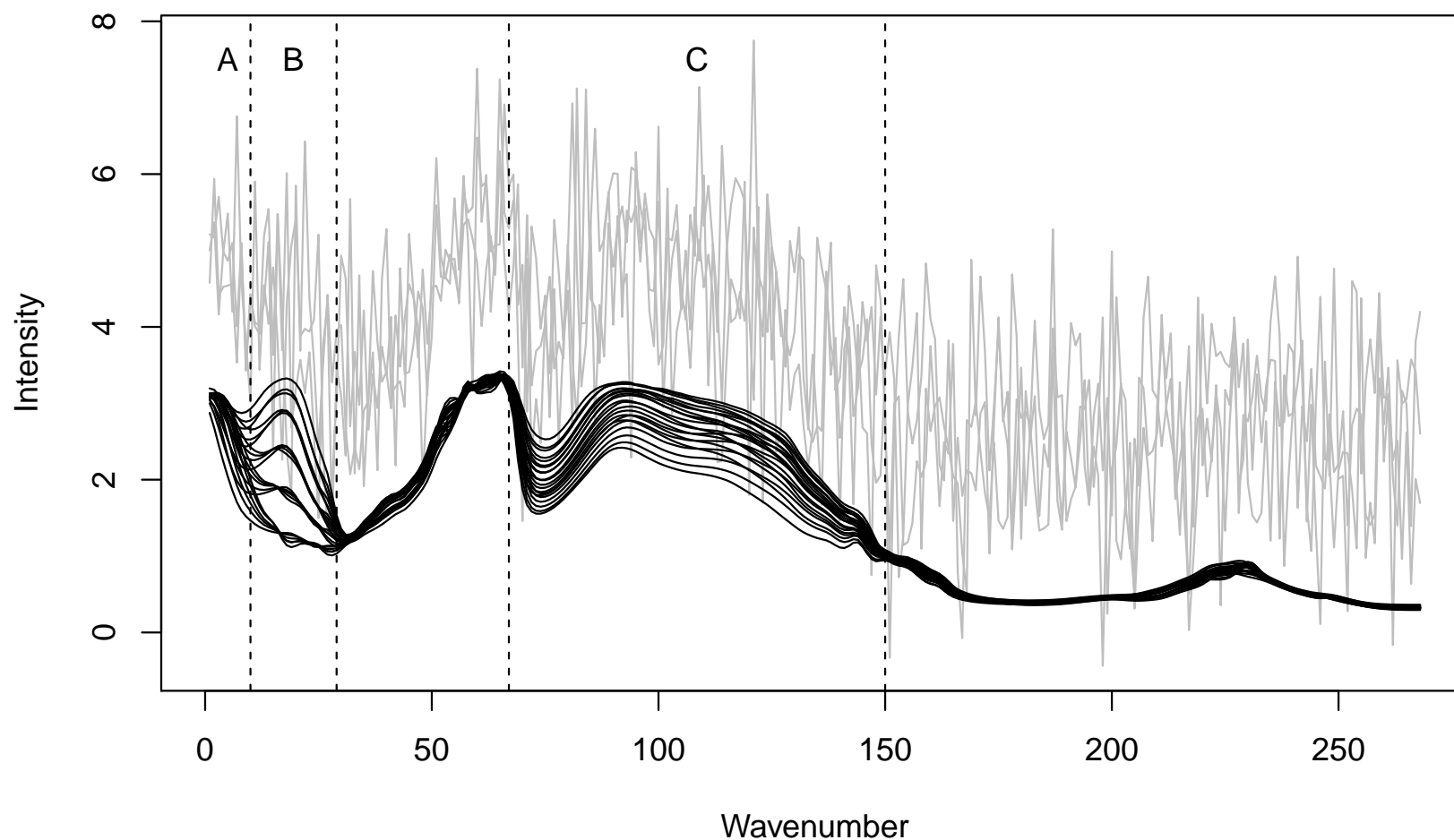


(d) Robust sparse PCA  
 $\lambda = 2.36$



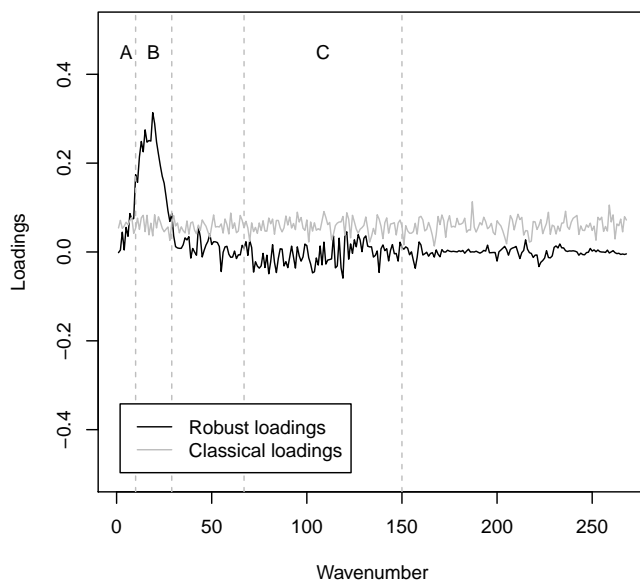
# Example yarn data set

**Yarn data**, with 3 outlying spectra added

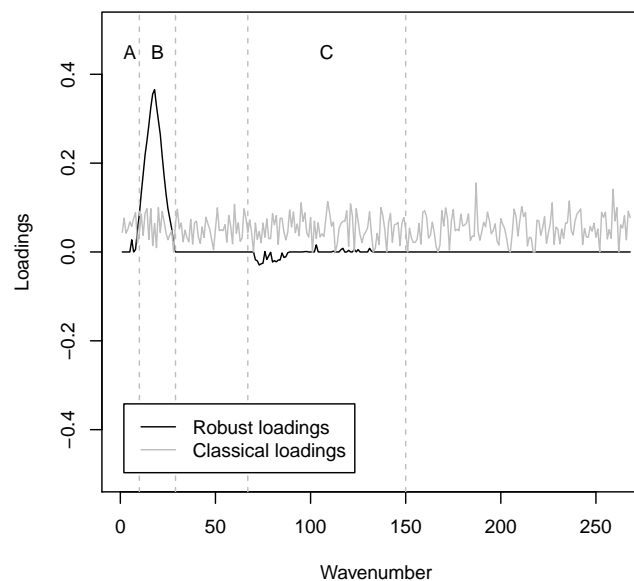


# Example yarn data set

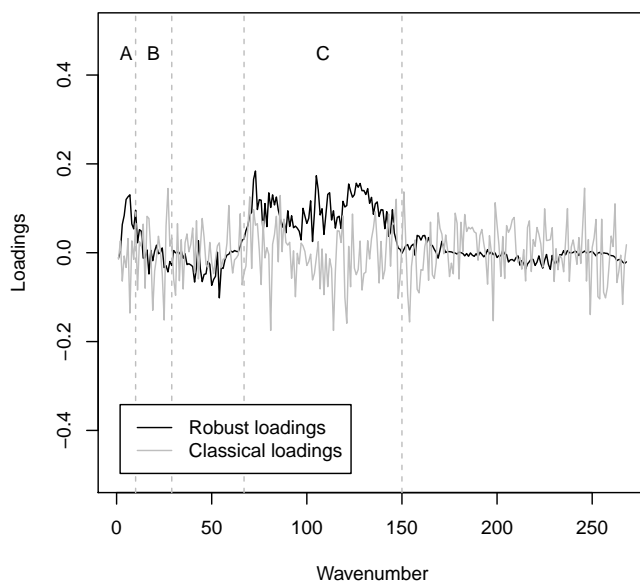
(a) PC1



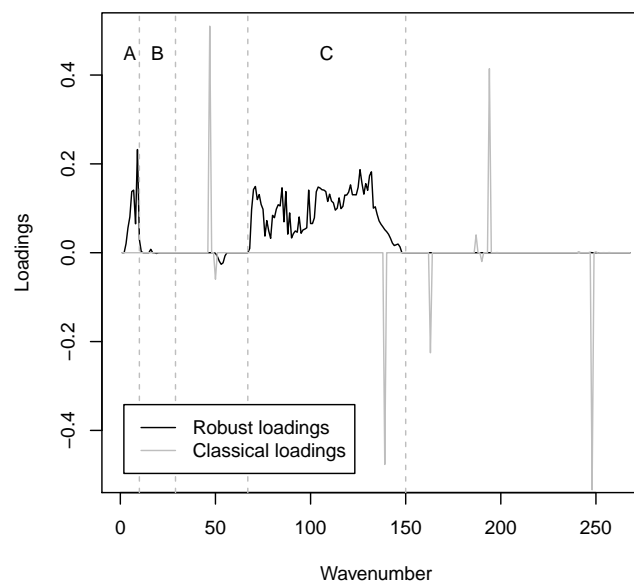
(b) Sparse PC1



(c) PC2



(d) Sparse PC2





Implementation in R package [pcaPP](#):

- Unconstrained PCA with Grid algorithm:

```
PCAgrid(x, k = 2, method = c ("mad", "sd", "qn"),  
maxiter = 10, splitcircle = 25, scores = TRUE, zero.tol = 1e-16,  
center = l1median, scale, trace = 0, store.call = TRUE, control, ...)
```

- Sparse PCA with Grid algorithm:

```
sPCAgrid(x, k = 2, method = c ("mad", "sd", "qn"), lambda = 1,  
maxiter = 10, splitcircle = 25, scores = TRUE, zero.tol = 1e-16,  
center = l1median, scale, trace = 0, store.call = TRUE, control, ...)
```

## Sparse robust PCA ...

- compromise between maximizing robust variance and simplifying interpretability
- determination of PCA directions is not affected by outliers

## Projection-pursuit approach ...

- robust estimation also possible for  $p > n$
- components are extracted sequentially (stop after  $k$ )
- robust variance estimator and sparsity criterion can easily be changed

Croux, C., Filzmoser, P., and Fritz, H. (2011). Robust sparse principal component analysis. Technical Report, SM-2011-2, Vienna University of Technology, Austria; (submitted for publication).

Croux, C., Filzmoser, P., and Oliveira, M. (2007). Algorithms for projection-pursuit robust principal component analysis. *Chemometrics and Intelligent Laboratory Systems*, **87**, 218–225.

Guo, J., James, G., Levina, E., Michailidis, G., and Zhu, J. (2011). Principal component analysis with sparse fused loadings. *Journal of Computational and Graphical Statistics*. To appear.

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.

Rousseeuw, P. and Croux, C. (1993). Alternatives to the median absolute deviation. *Journal of the American Statistical Association*, **88**(424), 1273–1283.

Tibshirani, R. (1996). Regression shrinkage and selection via the Lasso. *Journal of the Royal Statistical Society, Series B*, **58**, 267–288.

Zou, H., Hastie, T., and Tibshirani, R. (2006). Sparse principal component analysis. *Journal of Computational and Graphical Statistics*, **15**(2), 265–286.