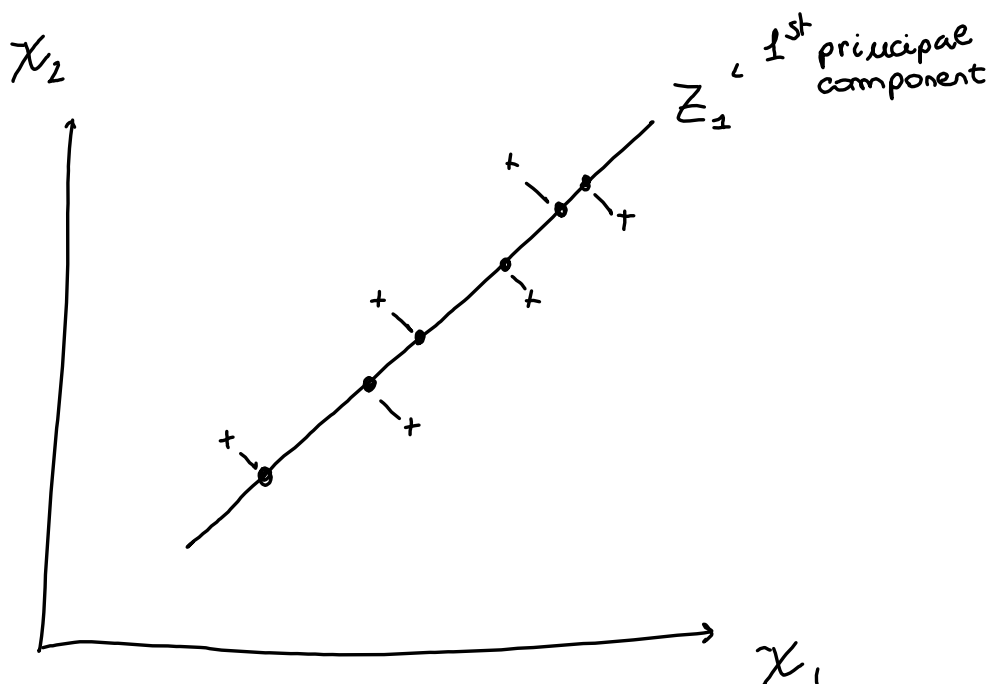


PCA:

projecting the data onto a low-dimensional space, without losing most of the information within the data



$$Z_1 := \varphi_1^T \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \Rightarrow \text{direction of the feature space along which the data are most variable}$$

In General:

$\mathbb{R}^n = (X_1, \dots, X_n)$ feature space

$\mathbb{R}^p = (Z_1, \dots, Z_p)$ principal component space ($p < n$)

• m training instances $\in \mathbb{R}^n$

$$\text{score} \rightarrow Z_{i,1} = \varphi_1^T X_i = \varphi_1^T \begin{pmatrix} X_{i1} \\ \vdots \\ X_{in} \end{pmatrix}$$

$$X \in \mathbb{R}^{m \times n}$$

data matrix

$$\Phi \in \mathbb{R}^{n \times p}$$

PC matrix

$$\begin{bmatrix} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & x_1 & \text{---} & \text{---} & \text{---} \\ \text{---} & x_2 & \text{---} & \text{---} & \text{---} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{---} & x_m & \text{---} & \text{---} & \text{---} \end{bmatrix} = \begin{bmatrix} \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{---} & \varphi_1^T & \text{---} & \text{---} & \text{---} \\ \text{---} & \varphi_2^T & \text{---} & \text{---} & \text{---} \\ \text{---} & \vdots & \text{---} & \text{---} & \text{---} \\ \text{---} & \varphi_n^T & \text{---} & \text{---} & \text{---} \end{bmatrix}$$

$\underbrace{\hspace{10em}}_{n \text{ features}} \qquad \underbrace{\hspace{10em}}_{p \text{ princ. comp}}$

PCA:

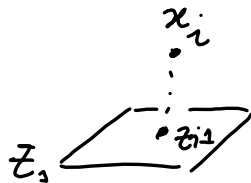
1st interpretation: Directions of highest variance

$X \in \mathbb{R}^{n \times p}$ data set ; p features.

$$E[X_i] = 0 \quad \forall i = 1, \dots, p \text{ (centered features)}$$

then:

$$\begin{aligned} \text{SCORE} \quad z_{i1} &= \phi_{11} x_{i1} + \phi_{21} x_{i2} + \dots + \phi_{p1} x_{ip} = \\ &= \phi_1^T \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} = \end{aligned}$$



projection of the i -th training instance onto the 1st princ. comp.

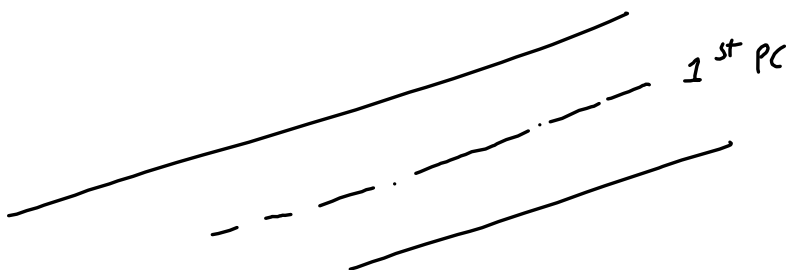
GOAL:

$$\max_{\phi_1^T} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\sum_{j=1}^p \phi_{j1} x_{ij} \right)^2 \right\} \text{ with } \|\phi_1\|^2 = 1$$

$$\max_{\phi_1^T} \left\{ \frac{1}{n} \sum_{i=1}^n \left(\phi_1^T x_i \right)^2 \right\} = \max_{\phi_1^T} \underbrace{V[\phi_1^T X]}$$

projection variance

searching the direction upon which the project points have highest variance.



PCA:

1st interpretation, geometric mess.

Repeat until you find the needed number of PC, adding the constraint the every new PC must be \perp others (\perp). After we have found k principal components the whole dataset can be projected upon the lower-dimensional principle component space.

$$\underbrace{\begin{pmatrix} x_1 & \dots & x_p \\ x_{11}, \dots, x_{1p} \\ x_{21}, \dots, x_{2p} \\ \vdots \\ x_{n1}, \dots, x_{np} \end{pmatrix}}_{\text{original space}} \underbrace{\begin{pmatrix} \phi_1 & \phi_2 & \dots & \phi_k \end{pmatrix}}_{\text{feature weight matrix}} = \underbrace{\begin{pmatrix} z_{11}, \dots, z_{1k} \\ z_{21}, \dots, z_{2k} \\ \vdots \\ z_{n1}, \dots, z_{nk} \end{pmatrix}}_{\text{principal component space}} \quad \boxed{k < p}$$

all projections on the 1st component.

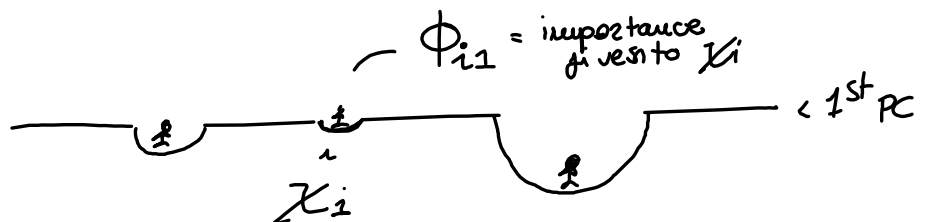
Score

$$\begin{pmatrix} \phi_{11} & \phi_{12} & \dots & \phi_{1k} \\ \phi_{21} & \phi_{22} & \dots & \phi_{2k} \\ \vdots & \vdots & & \vdots \\ \phi_{p1} & \phi_{p2} & \dots & \phi_{pk} \end{pmatrix}$$

loading vector of the 1st component

ϕ_{pk} = p-th feature weight (x_p) for the k-th PC.

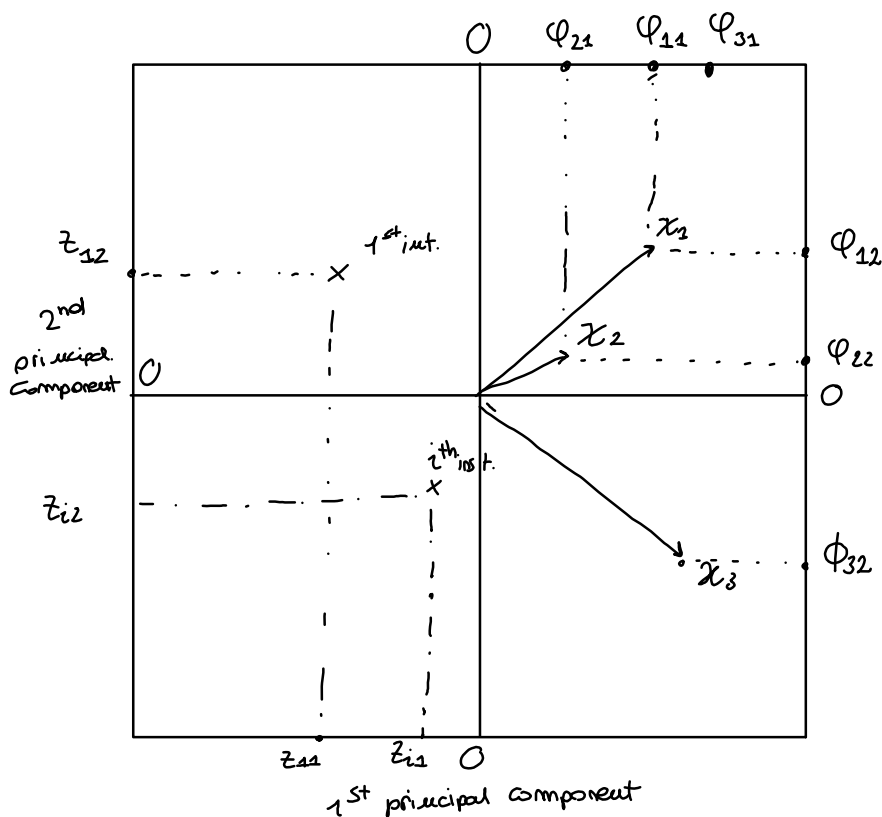
Tells the importance given from that PC to the original p features



PCA

1st interpretation: plotting.

The BIROT is a plot that shows the original data wrt the 1st and 2nd principal comp. It also shows the loadings for each feature.



$$\begin{array}{ccc}
 \text{original data} & \text{load. matrix} & \text{projected point} \\
 \begin{pmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ x_{n1} & x_{n2} & x_{n3} \end{pmatrix} & \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \\ \vdots & \vdots \\ \phi_{31} & \phi_{32} \end{pmatrix} & = \begin{pmatrix} z_{11} & z_{12} \\ z_{21} & z_{22} \\ \vdots & \vdots \\ z_{n1} & z_{n2} \end{pmatrix}
 \end{array}$$

ϕ_1 loading vector of the first component
 projection of x_1 onto the 2nd component
 new coordinates of the x_n train. instance

PCA:

2nd int: the principal components can be seen to be also the directions in the feature space that best approximate the data.

Specifically the same loading (ϕ_1, \dots, ϕ_k) vectors can be shown to be the solution of the following optimization problem:

$$\min_{A, B} \left\{ \sum_{j=1}^p \sum_{i=1}^n \left(x_{ij} - \sum_{m=1}^M z_{im} b_{im} \right)^2 \right\}$$

IDEA:

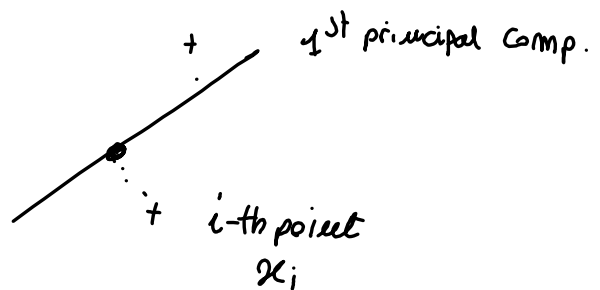
each point of the training set x_i can be approximated by M products $(z_{im} \cdot b_{im})$

SOLUTION

The first M principal components scores $(z_{im} = z_{im})$ and loadings $(\phi_{im} = b_{im})$ solve the problem.

$$x_{ij} \approx z_{i1}\phi_{i1} + z_{i2}\phi_{i2} + \dots$$

Meaning that the projected point (z_1) onto the principal components space is the best least square approx of x_i



PCA

% of variance explained

Each feature X_j can be seen as a random variable:

\tilde{X}_j with $E[X_j] = 0$ (if centered) and
a variance $V[\tilde{X}_j] = E[\tilde{X}_j^2] - (E[\tilde{X}_j])^2$

So the total variance in the data is:

$$\sum_{j=1}^p \text{Var}(\tilde{X}_j) = \frac{1}{n} \sum_{j=1}^p \sum_{i=1}^n x_{ij}^2$$

Also a princ. comp. can be seen as a rand. var. derived from $\tilde{X}_1, \dots, \tilde{X}_p$.

$$\tilde{Z}_m = \Phi_m^T \begin{pmatrix} \tilde{X}_1 \\ \vdots \\ \tilde{X}_p \end{pmatrix} \text{ with realization: } z_{im} = \Phi_m^T \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix}$$

$$\begin{aligned} \text{So we have: } V[\tilde{Z}_m] &= E[\tilde{Z}_m^2] - (E[\tilde{Z}_m])^2 \\ &= \frac{1}{n} \sum_{i=1}^n z_{im}^2 - \left(E\left[\Phi_m^T \begin{pmatrix} \tilde{X}_1 \\ \vdots \\ \tilde{X}_p \end{pmatrix} \right] \right)^2 \end{aligned}$$

So the % of variance explained by the m -th principal comp is:

$$\frac{V[\tilde{Z}_m]}{\sum_{j=1}^p V[\tilde{X}_j]} = \frac{\sum_{i=1}^n z_{im}^2}{\sum_{j=1}^p \sum_{i=1}^n x_{ij}^2}$$

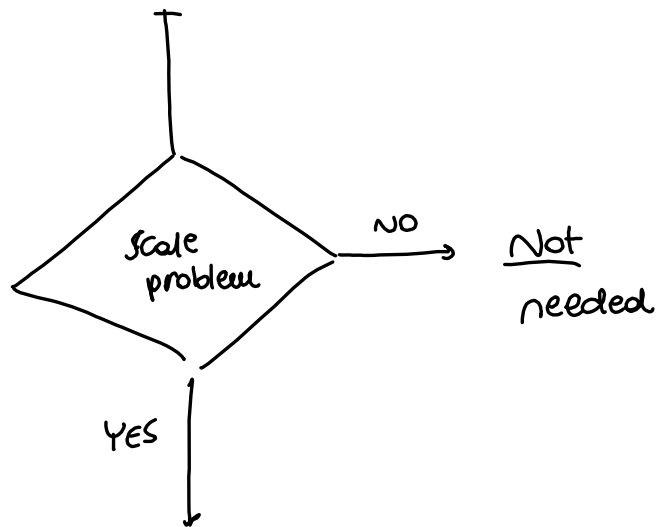
PCA

data preparation

1. Centering each feature at \emptyset .

Fundamental to make everything work

2. Feature scaling



Then it could be that $V[\tilde{X}_j] \geq V[\tilde{X}_i] \ (i \neq j)$ for ~~some~~ and not for any other reason. So an idea could be standardizing each feature: so that

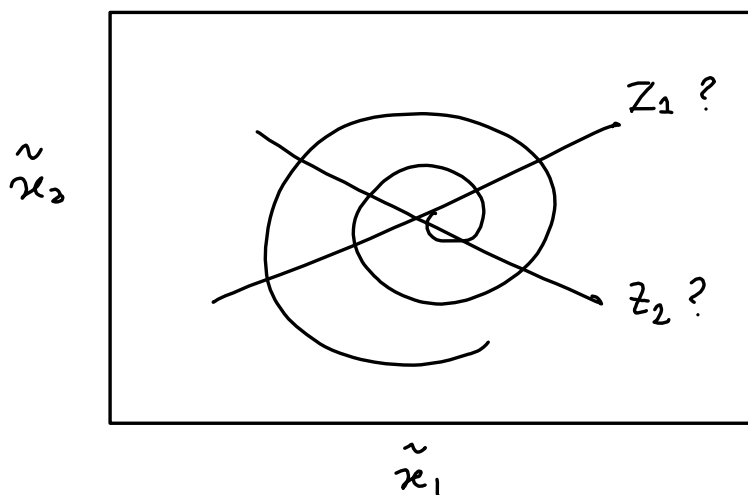
$$\tilde{X}_j \sim (E[\tilde{X}_j] = \emptyset; V[\tilde{X}_j] = 1)$$

(Unless the 1st p.c. could capture the most variance just of some features, cause of their scales and not for some interesting phenomena)

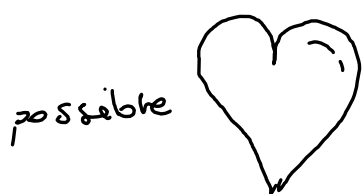
Kernel PCA

Normal PCA assumes that the variance within the data can be explained and viewed also as a lower-dimensional hyperplane.

What if the variance of data could be best understood as a non-linear lower-dimensional surface within the feature space?



We could apply a kernel trick, hoping that in a higher-dim. space than the original feature space, the data variance could be best understood and explained on a linear surface. (look SVM)



kernels

- Linear
- RBF Gaussian
- Sigmoid [...]