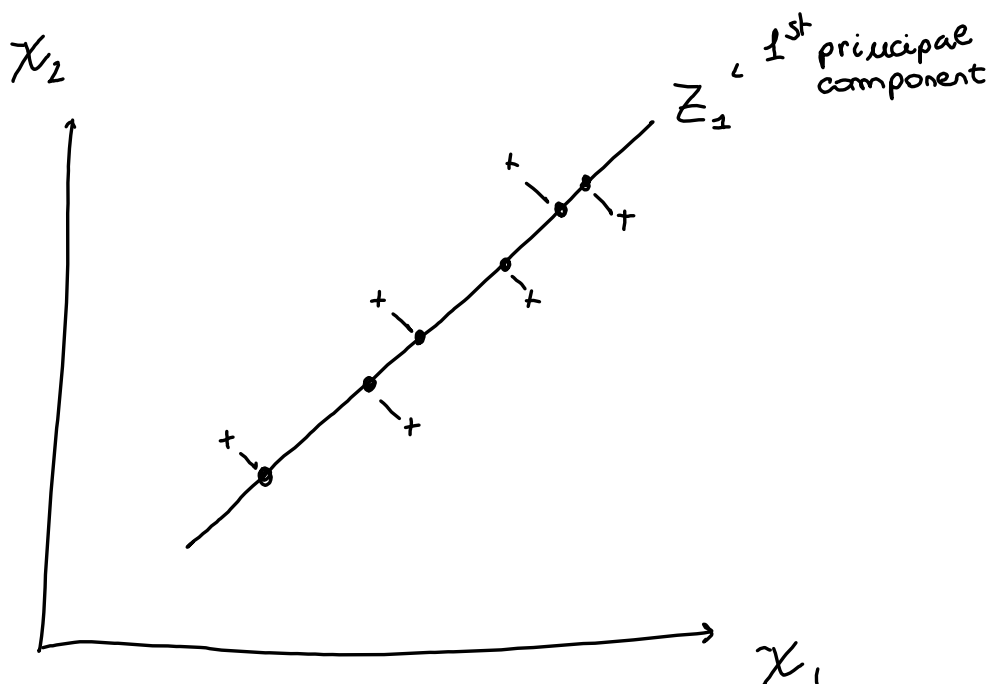# PCA:

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



$$Z_1 := \varphi_1^T \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \Rightarrow \text{direction of the feature space along which the data are \underline{most} variable}$$

## 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 $\in \mathbb{R}^n$

$$\text{score} \rightarrow Z_{i1} = \varphi_1^T x_1 = \varphi_1^T \begin{pmatrix} x_{11} \\ \vdots \\ x_{n1} \end{pmatrix}$$

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

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

$$\left[ \phantom{xxxxx} \right. = \begin{bmatrix} \text{---} & x_1 & \text{---} \\ \text{---} & x_2 & \text{---} \\ & \vdots & \\ \text{---} & x_m & \text{---} \end{bmatrix} \begin{bmatrix} \varphi_1^T \\ \varphi_2^T \\ \vdots \\ \varphi_n^p \end{bmatrix}$$

$n$ features

$p$ princ. comp

# PCA:

1st interpretation: Directions of highest variance

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

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

then:

$$Z_{i1} = \phi_{11} x_{i1} + \phi_{21} x_{i2} + \dots + \phi_{p1} x_{ip} =$$

SCORE

$$= \phi_1^T \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} =$$

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

GOAL:

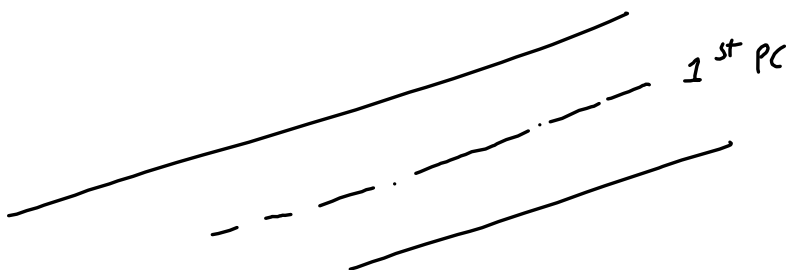$$\max_{\phi_1^T} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{j=1}^{p} \phi_{j1} x_{ij} \right)^2 \right\} \quad \text{with} \quad \|\phi_1\|^2 = 1$$

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

$\underbrace{\qquad\qquad\qquad\qquad\qquad\qquad}$

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

$\underbrace{\qquad\qquad}$ projection variance

1st PC

# 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 ($\perp\!\!\perp$). After we have found $k$ principal components the whole dataset can be projected upon the lower-dimensional prin. component space.

$$
\underbrace{\begin{pmatrix} x_{11}, \cdots, x_{1P} \\ x_{21}, \cdots, x_{2P} \\ \vdots \\ x_{n1}, \cdots, x_{nP} \end{pmatrix}}_{\substack{\text{original} \\ \text{space}}}
\underbrace{\begin{pmatrix} \varphi_1 \; \varphi_2 \sim \varphi_k \end{pmatrix}}_{\substack{\text{feature} \\ \text{weight} \\ \text{matrix}}}
=
\underbrace{\begin{pmatrix} z_{11}, \cdots, z_{1k} \\ z_{21}, \cdots, z_{2k} \\ \vdots \\ z_{n1}, \cdots, \boxed{z_{nk}} \end{pmatrix}}_{\substack{\text{principal component} \\ \text{space}}}
$$

$$
\begin{matrix} x_1 & \cdots & x_p \end{matrix}
$$

all projections on the 1st component.

$\boxed{k < P}$

$\longrightarrow$ score

$\lambda$

loading vector of the 1st component

$$
\begin{pmatrix} \ddot{\varphi}_{11} & \varphi_{12} & \cdots & \varphi_{1k} \\ \varphi_{21} & \varphi_{12} & \cdots & \varphi_{2k} \\ \vdots & \vdots & & \vdots \\ \varphi_{p1} & \varphi_{p2} & \cdots & \boxed{\varphi_{pk}} \end{pmatrix}
$$

$\longrightarrow$ p-th feature weight ($x_p$) for the k-th PC.

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

$\varphi_{11} = $ importance given to $x_i$



$< 1^{st}$ PC

$x_i$

# PCA

1st interpretation: plotting.

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



original data

$$\begin{pmatrix} X_{11} & X_{12} & X_{13} \\ X_{21} & & \\ \vdots & \vdots & \vdots \\ X_{n1} & X_{n2} & X_{n3} \end{pmatrix} \begin{pmatrix} \phi_{11} & \phi_{12} \\ \phi_{21} & \phi_{22} \\ \phi_{31} & \phi_{32} \end{pmatrix} = \begin{bmatrix} z_{11} & z_{12} \\ z_{21} & z_{12} \\ & \vdots \\ z_{n1} & z_{n2} \end{bmatrix}$$

load. matrix

projected point

~ projection of $x_1$ onto the 2nd component

< new coordinates of the $x_n$ train. instance

$\phi_1$ loading vector of the first component

# 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 $(\phi_1, ..., \phi_k)$ vectors can be shown to be the solution of the following optimization problem:

$$\min_{A,B} \left\{ \sum_{j=1}^{P} \sum_{i=1}^{n} \left( x_{ij} - \sum_{m=1}^{M} a_{im} b_{im} \right)^2 \right\}$$
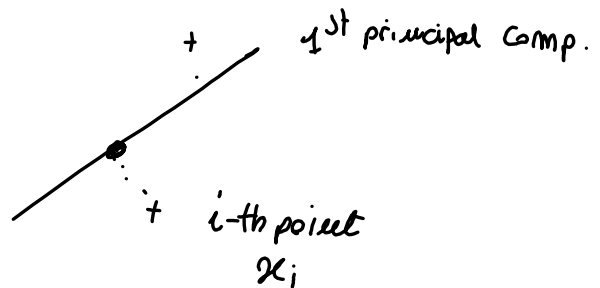
## IDEA :

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

## SOLUTION

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

$$x_{ij} \simeq z_{i1}\phi_{i1} + z_{i2}\phi_{i2} + \cdots$$

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



1st principal comp.

+ i-th point $x_i$

# PCA

% of variance explained

Each feature $X_j$ can be seen as a random variable:

$\tilde{X}_j$ with $E[X_j] = \emptyset$ (if centered) and

an variance $V[\tilde{X}_j] = E[\tilde{x}_j]^2 - (E[X_j])^{2^0}$

So the ~~total~~ variance in the data is:

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

Also a princ. comp. can be seen as a rand. var. derived from $\tilde{X}_2, ..., \tilde{X}_p$.

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

So we have: $V[\tilde{Z}_m] = E[\tilde{Z}_m^2] - (E[\tilde{Z}])^2$

$$= \frac{1}{n} \sum_{i=1}^{n} Z_{im}^2 - \left( E[\phi_m^T \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}]\right)^{2^0}$$

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

$$\frac{V[\hat{Z}_m]}{\sum_{j=1}^{p} V[\hat{X}_j]} = \frac{\sum_{i=1}^{n} Z_{im}^2}{\sum_{j=1}^{p} \sum_{i=1}^{n} x_{ij}^2}$$
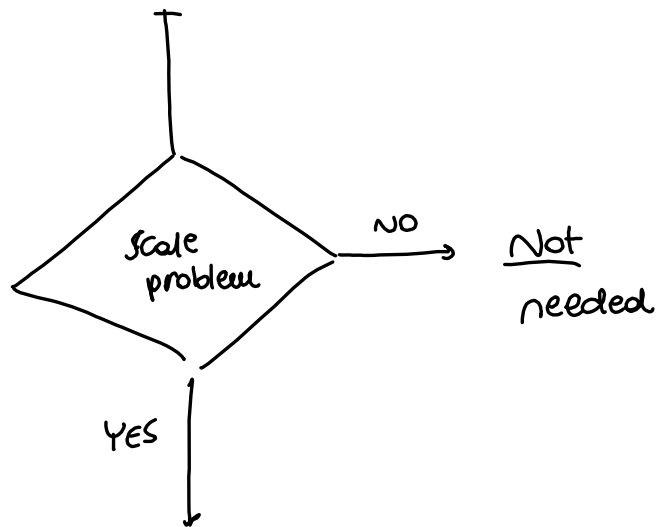
# PCA
data preparation

1. Centering each feature at $\emptyset$.

Fundamental to make everything work

2. Feature scaling



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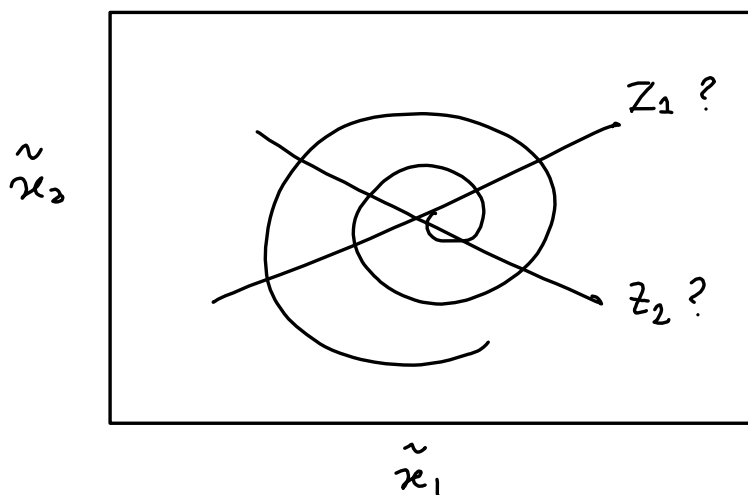be standardizing each feature: do that

$$\breve{X}_j \sim \left( E[\breve{X}_j] = \emptyset \; ; \; \mathbb{V}[\breve{X}_j] = 1 \right)$$

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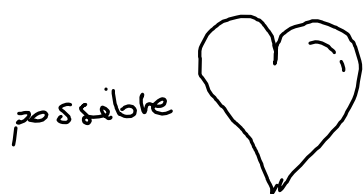some features, cause of their scals and not for
dome interesting phenomena

# Kernel PCA

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

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



We could apply a kernel trick, hoping that in an higher-dim. space than the original feature space, the data variance could be best understood and explained on a linear surface. (look SVM)
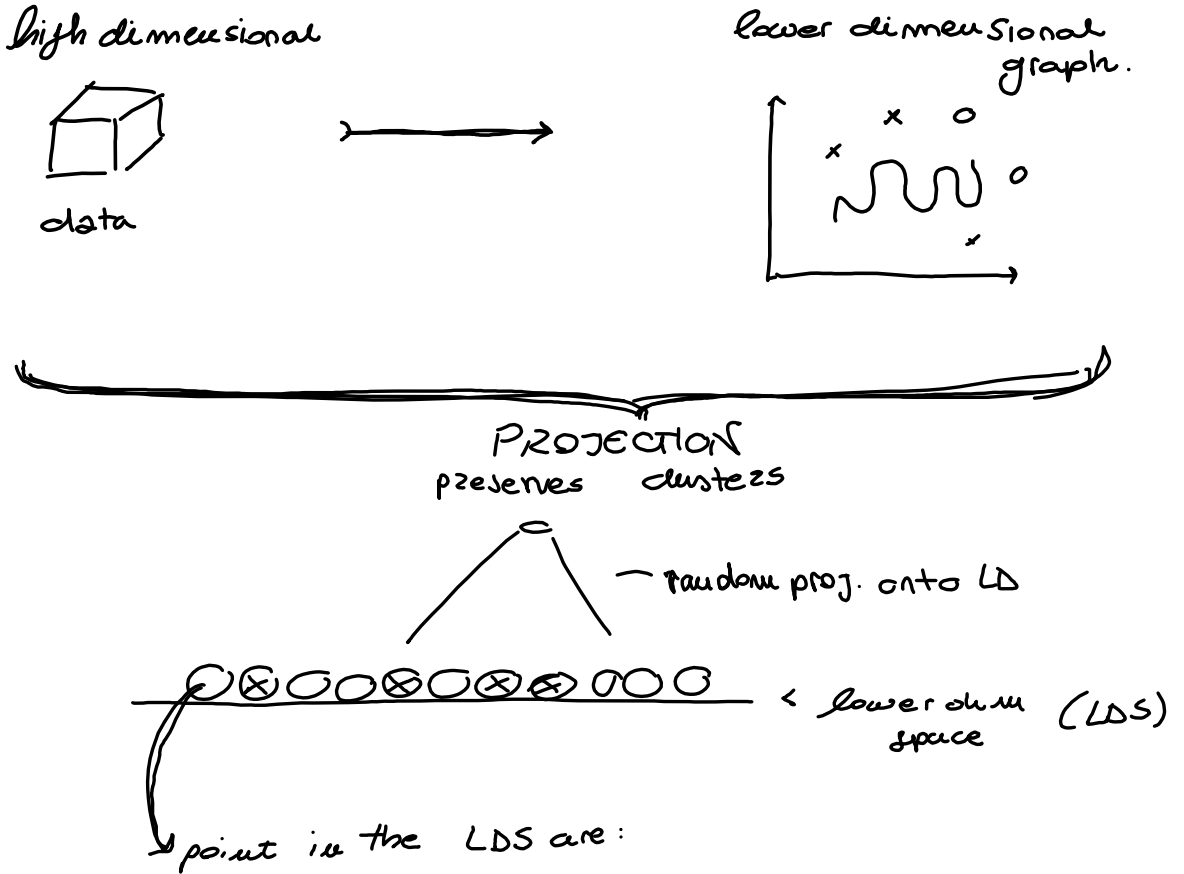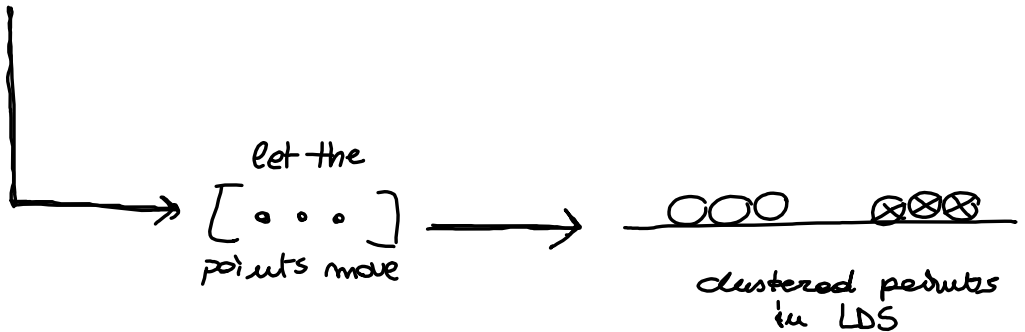
possible ♡   kernels → Linear
                      → RBF Gaussian
                      → Sigmoid    [...]

t - SNE : dim reduction algorithm widely used
for imaging

high dimensional

data

lower dimensional
graph.

PROJECTION
preserves clusters

— random proj. onto LD

< lower dim (LDS)
space

point in the LDS are:

repelled
by points far
from...

~~> ◯ ~~> attracted : points near in the original
space

let the
[ ● ● ● ]
points move

clustered points
in LDS

# t-SNE: moth behind

∀ space:

point near/far in the scatter plot ⟷ $Sim(x_j, x_i)$ using:



t is wider than $\mathcal{N}$ in general.

$$\mathcal{N}\left(x_j, \sigma^2\left(\frac{1}{f_x(x)}\right)\right) \text{ or } t$$

high-density

$\sigma^2 <<$

overall higher similarities

low-density

$\sigma^2 >>$

overall lower similarities.

similarity matrix

$$[S]$$

working:

LDS

points moving to make [S] or $\not A$ ≈ [S] LDS.

$$[S]_{orig. space} \triangleright [S]_{low-dim space}$$

NOTE  [S] LDS computed ~ t to avoid cluot clumping in LDS