

FoDA

L21

• Principal Component
• Analysis (PCA)

centering, PCA, MDS

Dimensionality Reduction

Input $A \subset \mathbb{R}^d$
 $a_1, a_2 \dots a_n \in \mathbb{R}^d$
 $d(a_i, a_j) = \|a_i - a_j\|_2$

Goal Low-dimensional Representation $a_i \mapsto b_i$

$B \subset \mathbb{R}^k$ $b_1, b_2 \dots b_n$

B represents k -dim subspace of \mathbb{R}^d

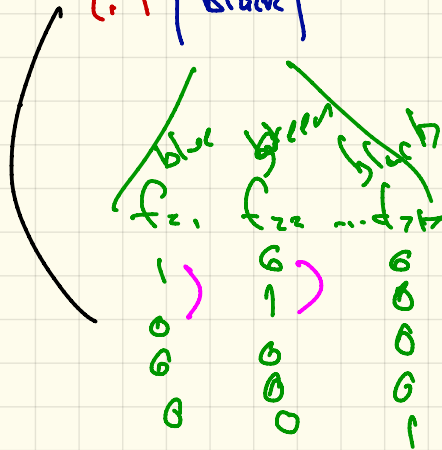
$V_B = \{v_1, v_2 \dots v_k\}$

Goal minimize $\| \pi_B(A) - A \|^2 = \sum_{i=1}^n \| \pi_B(a_i) - a_i \|^2$

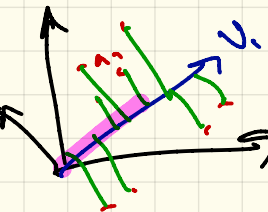
Alternativer Problem

data points	features	f_1	f_2	f_3	f_4	...	f_n
x_1		7.1	blue				
x_2		8.3	green				
...		.					
x_n		9.4	black				

can't mix
numerical
categorical
in unsupervised



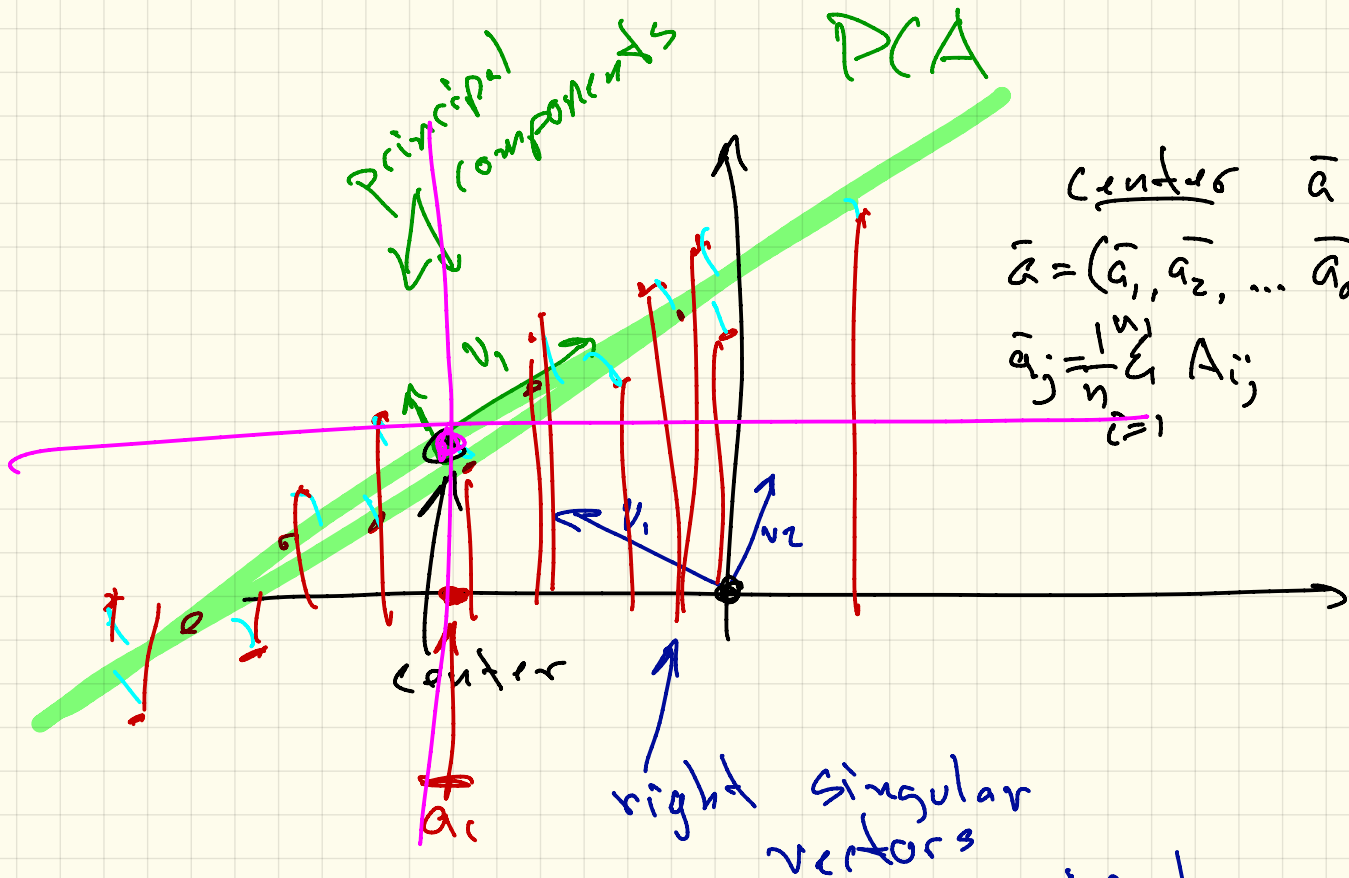
High Dimensional Data
 $a_1, a_2, \dots, a_n \in \mathbb{R}^d$



How do I sort?

↳ The first dimension v_1
provides the best 1-dim ordering

$$\{\langle v_1, a_1 \rangle, \langle v_1, a_2 \rangle, \dots, \langle v_1, a_n \rangle\}$$



center $\bar{a} \in \mathbb{R}^d$

$$\bar{a} = (\bar{a}_1, \bar{a}_2, \dots, \bar{a}_d)$$

$$\bar{a}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$$

Not Good!!

Centering $A \in \mathbb{R}^{n \times p}$

Two ways \rightarrow

1. compute center

$$\bar{a} = (\bar{a}_1, \dots, \bar{a}_p)$$

$$\bar{a}_j = \frac{1}{n} \sum_{i=1}^n A_{ij}$$

$$\bar{A}_{ij} = A_{ij} - \bar{a}_j$$

2. centering matrix

$$C_n = I_n - \frac{1}{n} \mathbb{1} \mathbb{1}^T$$

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix}_{n \times n}$$

$$= \frac{1}{n} \begin{bmatrix} 1 & 1 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

$$= \begin{bmatrix} 1-\frac{1}{n} & -\frac{1}{n} & & -\frac{1}{n} \\ -\frac{1}{n} & 1-\frac{1}{n} & & -\frac{1}{n} \\ \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{n} & -\frac{1}{n} & & 1-\frac{1}{n} \end{bmatrix}$$

$$A - \frac{1}{n} [\mathbb{1} \mathbb{1}^T] A = \bar{A}_{ij} = C_n A$$

Principal Component Analysis

Input $A \in \mathbb{R}^{n \times d}$, param $k < d < n$

Output $B: V_B = \{v_1, \dots, v_k\}$ minimizes $\sum_{i=1}^n \|\Pi_B(a_i) - a_i\|^2$

1. $\tilde{A} = C_n A$ ← centering

2. $U S V^T = \text{svd}(\tilde{A})$

3. $V_B = \{v_1, v_2, \dots, v_k\}$

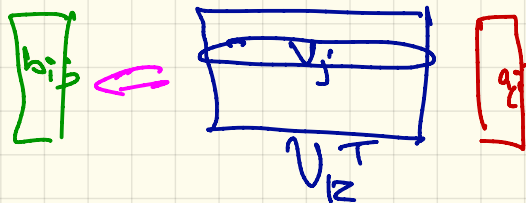
$B = \{b_1, \dots, b_n\}$

$B = V_k^T A \in \mathbb{R}^k$

If I want k -dimensional pts

$$B = \{b_1, \dots, b_n\} \in \mathbb{R}^k$$

$$B = V_k^T A = b_i = V_k^T a_i$$



$$b_{ij} = \langle v_j, a_i \rangle$$

If I want d -dim pts
subspace

1. Centering $\bar{A} \in \mathbb{R}^{n \times d}$

2. SVD $U S V^T$

$$\tilde{A}_k = U_k S_k V_k^T \in \mathbb{R}^{n \times d}$$

$$[A_k]_{ij} = [\tilde{A}_k]_{ij} + a_j$$

undo

SVD

best subspace
through origin

PCA

best subspace

(not necessarily)
w/ origin)

Variance
 RV. X

V_1 highest variance
 subspace

$$\text{Var}[X] = E[(X - E[X])^2]$$

$$P_r[X] = \frac{1}{|X|}$$

Sample expected value $E[X] = \frac{1}{|X|} \sum_{x \in X} x \cdot P_r[x]$

$$= \frac{1}{|X|} \sum_{x \in X} x$$

Sample variance

centering makes 0

$$\frac{1}{n} \sum_{i=1}^n \left(\langle v, a_i \rangle - \left(\frac{1}{n} \sum_{i=1}^n \langle v, a_i \rangle \right) \right)^2 = \frac{1}{n} \sum_{i=1}^n \langle v, a_i \rangle^2$$

$X = \{ \langle v, a_1 \rangle, \langle v, a_2 \rangle, \dots \}$

Multi-Dimensional Scaling (MDS)

Input (Ia) Set n objects and a distance d

(Ib) $D \in \mathbb{R}^{n \times n}$ $D_{ij} = d(x_i, x_j)$

Goal : Embedding of n objects in

$$\mathbb{R}^k \text{ as } G = \{g_1, \dots, g_n\}$$

$$\text{s.t. } \|g_i - g_j\| = D_{ij} = d(x_i, x_j)$$

Classical MDS

Input $D \in \mathbb{R}^{n \times n}$, k

1. Square $D^{(2)} : D_{ij}^{(2)} = D_{ij}^2$

2. $M = -\frac{1}{2} C_n D^{(2)} C_n$

3. $[L, U] = \text{eig}(M)$

4. Return $Q = U_k L_k^{1/2}$

$$D^{(2)} \approx -AA^T$$

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

← double centering

top k \uparrow
eigen
vectors \leftarrow

$\sqrt{}$ eigen
values \rightarrow
 \approx sing. values

Why does Classical MDS work?

Assume $\exists A = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix} \in \mathbb{R}^{n \times p} \parallel \{a_i - a_j\}$
 $d(a_i, a_j) = D_{ij}$

$$\|a_i - a_j\|^2 = \|a_i\|^2 + \|a_j\|^2 - 2 \langle a_i, a_j \rangle$$

$$(a - b)^2 = a^2 + b^2 - 2ab$$

$$\langle a_i, a_j \rangle = \frac{1}{2} (\|a_i\|^2 + \|a_j\|^2 - \|a_i - a_j\|^2)$$

$$[AA^T]_{ij} = \langle a_i, a_j \rangle = \frac{1}{2} (\|a_i\|^2 + \|a_j\|^2 - D_{ij}^2)$$

assume $a_1 = \text{origin } (0, \dots, 0)$

$$[AA^T]_{ii} = \frac{1}{2} (D_{i1}^2 + D_{j1}^2 - D_{ij}^2)$$