# FoDA   L21

---

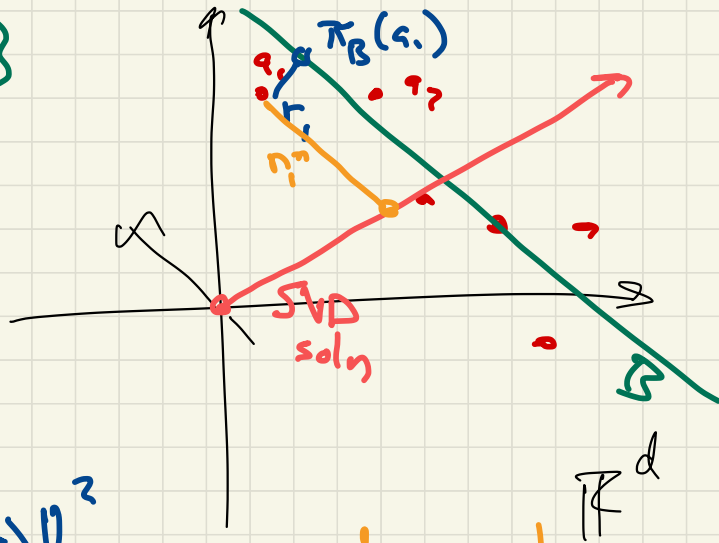## PCA (Principal Component Analysis)

and

## MDS (MultiDimensional Scaling)

# Principal Component Analysis

Input $A \in \mathbb{R}^{n \times d}$ $\qquad A = \{a_1, a_2, \ldots a_n\} \subset \mathbb{R}^d$

Goal $\quad B \qquad V_B = \{v_1, \ldots, v_k\}$

$\qquad\qquad\qquad\qquad$ orthogonal basis

$x \in \mathbb{R}^d$

$$\pi_B(x) = \sum_{j=1}^{k} v_j \langle v_j, x \rangle$$

Find minimize $\| A \cdot \pi_B(A) \|_F^2$
B

$$= \sum_{i=1}^{n} \underbrace{\| a_i - \pi_B(a_i) \|^2}_{r_i}$$

PCA fixes
by first
centering data



$\pi_B(a_1)$

$r_1$

$\eta_i^{vv}$

SVD
soln

$B$

$\mathbb{R}^d$

SVD does not
directly give soln
since includes 0

# A : Centering the data



$A_i$ ; $a_i$ ; $A_{ij}$ ; $\dot{b} \dot{2} \Rightarrow \bar{a}_j$

Shift of data so its average is the origin $\underline{O} = (0, 0, \dots 0)$

PCA = centering then SVD

$$\bar{a}_j = \frac{1}{n} \sum_{i=1}^{n} A_{ij} \qquad \text{average of } j\text{th coord.}$$
on all data.

$$\bar{a} = (\bar{a}_1, \bar{a}_2, \dots \bar{a}_d) \in \mathbb{R}^d \quad \leftarrow \text{average point of all data.}$$

$\tilde{A}$ defined so $\boxed{\tilde{A}_{ij} = A_{ij} - \bar{a}_j}$

# Centering Matrix $\quad C_n \in \mathbb{R}^{n \times n}$

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

$$\begin{bmatrix} 1 & & 0 \\ & 1 & \\ 0 & & \ddots \end{bmatrix} - \begin{bmatrix} \frac{1}{n} & \frac{1}{n} & \cdots \\ \frac{1}{n} & \ddots & \\ \cdots & & \frac{1}{n} \end{bmatrix}$$

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

$$\mathbb{1} \in \mathbb{R}^{n \times 1}$$

$$\mathbb{1} = \begin{bmatrix} \vdots \\ \vdots \end{bmatrix}$$

$$\mathbb{1} \mathbb{1}^T = \begin{bmatrix} \vdots & \vdots & \cdots & \vdots \end{bmatrix} \in \mathbb{R}^{n \times n}$$

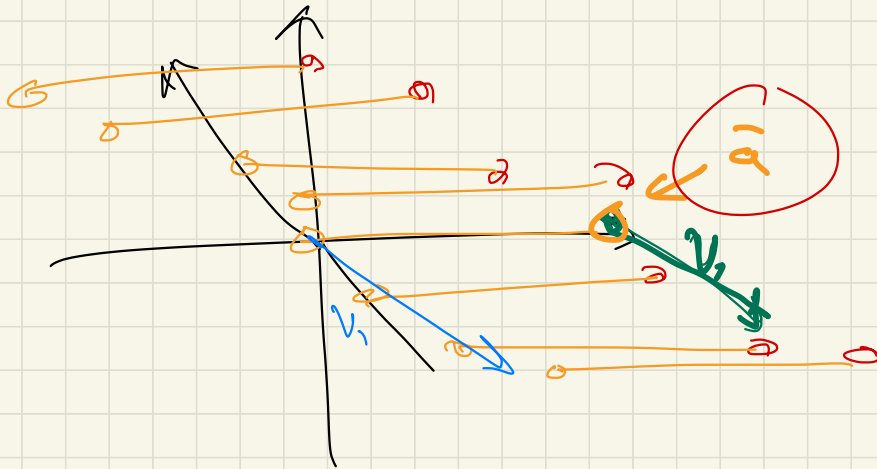$$\tilde{A} = C_n A = A - \frac{1}{n} \mathbb{1} \mathbb{1}^T A$$

# PCA

1. $\bar{A} = C_n A$

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

3. Take top $k$ sing. values $\sigma_1 = S_{11}, \sigma_2 = S_{22} ... \sigma_k$

   sing vectors $v_1, ... v_k$

<span style="color:blue">principal components</span>

<span style="color:red">Useful for</span>
<span style="color:red">• $k = 1, 2, 3$</span>

# Supervised vs. Unsupervised Learning

| Supervised | vs. | Unsupervised Learning |
|---|---|---|
| Input $(X, y)$ | | $X$ |

Goal   Predict
       $y$ using $X$

Find "shape" of
the data
(simple)

Endgoal

Subgoal

Exporadirg

# Difference PCA and Regression



PCA

$r_B(a_i)$

B

$(a_i, g^*)$

Y

$a_1$

$r_i$

$(a_i, g)$

$a_2$

$a_i$

$x_1$

$x_2$

$\mathbb{R}^d$

PCA

B rank $k$

$k < d$

Regression

$f: \mathbb{R}^d \to \mathbb{R}$

# MultiDimensional Scaling (MDS)

Input    distance matrix $D \in \mathbb{R}^{n \times n}$

$n$ objects    $P_1, P_2 \text{ (p) } P_n$   ← cities

$D_{ij} = \text{distance}(P_i, P_j)$

↑ cost of airplane ticket

$q_i := u(p_i)$
$\in \mathbb{R}^k$

**Goal**   Embedding $\phi$ $P_1 \dots P_n$ into $\mathbb{R}^k$

$u$    so   $\| u(p_i) - u(p_j') \| \approx \text{distance}(p_i, p_j)$

# Classical MDS

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

Centering Matrix $C_n$

$$\boxed{C_n^\top = C_n}$$

1. Square all distance

$$D^{(2)} \quad \left(D^{(2)}\right)_{ij} = \left(D_{ij}\right)^2$$

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

double centering

$$M \simeq \tilde{A} \tilde{A}^\top$$

eigenvalues

3. $[L, V] = eig(M)$

$V_k \in \mathbb{R}^{n \times k}$  $L_k \in \mathbb{R}^{k \times k}$ $\begin{bmatrix} \lambda_1 \lambda_2 & 0 \\ 0 & \ddots \lambda_k \end{bmatrix}$

4. Return $Q = V_k L_k^{1/2} = \{g_1, g_2 \dots g_n\} \in \mathbb{R}^k$

# Why Classical MDS work?

Similarity Matrix $S \in \mathbb{R}^{n \times n}$

$$S_{ij} = \text{similarity}(p_i, p_j)$$

$$= \langle a_i, a_j \rangle$$

$$[A A^T]_{ij} = \langle a_i, a_j \rangle$$

$$[D^{(2)}]_{ij} = \|a_i - a_j\|^2 = \|a_i\|^2 + \|a_j\|^2 - 2 \langle a_i, a_j \rangle$$

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

pretend
there
exist

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

$$A = \begin{bmatrix} a_1 \\ a_2 \\ a_n \end{bmatrix}$$

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

Solve for $\|a_i\|^2$ and $\|a_j\|^2$

if assume $a_1 = (0, 0, 0, \ldots 0)$

$$\|a_1\|^2 = 0$$
$$\|a_i\|^2 = \|a_i - a_1\|^2 = \left[ D^{(2)} \right]_{i,1}$$

$$\langle a_i, a_j \rangle = \frac{1}{2} \left( \left[ D^{(2)} \right]_{i,1} + \left[ D^{(2)} \right]_{j,1} - \left[ D^{(2)} \right]_{ij} \right)$$