# L18: Multidimensional Scaling (MDS),

# Linear Discrimenent Analysis (LDA),

## and

# Distance Metric Learning (DML)

# Principal Component Analysis

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

map $\quad u : \mathbb{R}^d \to \mathbb{R}^{k}$

goal $\quad \boxed{B \in \mathbb{R}^{n \times k}} \quad b_i = u(a_i)$

$\qquad\qquad\qquad\qquad\qquad a_i \in$ row in $A$

$\quad n = $ data points

in low dimensional $\mathbb{R}^{k}$

so $\quad d(i,j) = \| b_i - b_j \|$

# Multidimensional Scaling (MDS)

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

$$D_{ij} = d(i, j)$$

**Exemples**

- cities

$d(i,j) = $ cost of airline flight between $i, j$

- more abstractly just be given $D$

# Classical MDS

1. Convert $D$ into $D^{(2)}$ : $D^{(2)}_{ij} = (D_{ij})^2$

2. Double Centering

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

   (turns into $n \times n$ inner products)

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

3. Eigendecomposition $[L, V] = \text{eigs}(M)$

   $M = V L V^T = (V L^{1/2})(V L^{1/2})^T$

4. Project onto top $k$ eigenvectors

   return $B = V_k L_k^{1/2} \in \mathbb{R}^{n \times k}$

   embedded data points $b_i = V_k L_k^{1/2}$

   $\in \mathbb{R}^{k \times n} \quad \mathbb{R}^{n \times n}$

# Why does MDS work?

Like instead of dist matrix $D$

$\hookrightarrow$ similarity matrix $S : S_{ij} = \langle a_i, a_j \rangle$

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

well $S = A A^T \in \mathbb{R}^{n \times n}$

$$S_{ij} = \langle a_i, a_j \rangle$$

best embedding of $A$, from $S = A A^T$

top $k$ eigenvectors $S$, or top $k$ left singular vectors $A$

$$\underbrace{\|a_i - a_j\|^2}_{D_{ij}^2} = \|a_i\|^2 + \|a_j\|^2 - 2 \boxed{\langle a_i, a_j \rangle} = S_{ij}$$

| trick set $a_1 = (0, 0, \ldots 0) \Rightarrow \|a_i\|^2 = \|a_i - a_1\|^2 = D_{i1}^2$
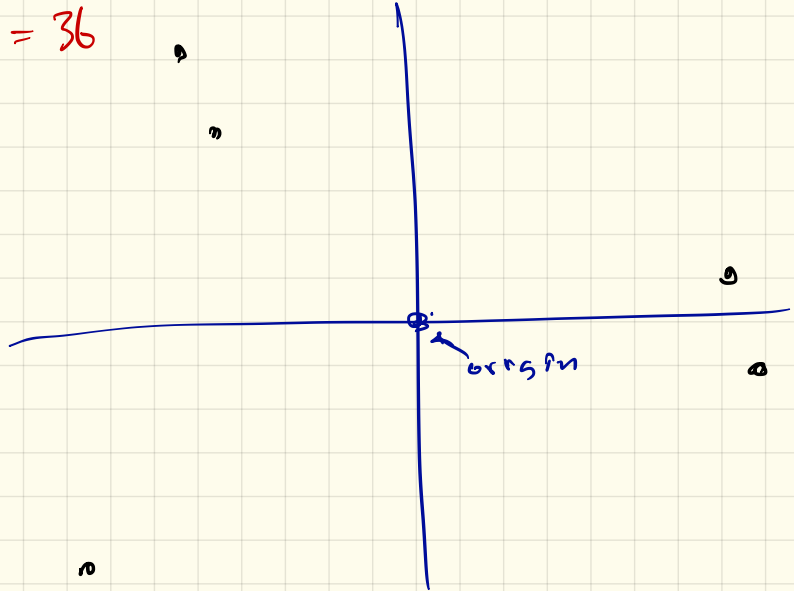
$$S_{ij} = \langle a_i, a_j \rangle = -\frac{1}{2} \boxed{\left( D_{ij}^2 \right.} - \left( D_{i1}^2 \right) - \left( D_{j1}^2 \right) \left. \right)$$

$\leftarrow$ average over all $a_i = 0$

$$D = \begin{bmatrix} 0 & 4 & 3 & 7 & 8 \\ 4 & 0 & 1 & 6 & 7 \\ 3 & 1 & 0 & 5 & 7 \\ 7 & 6 & 5 & 0 & 1 \\ 8 & 7 & 7 & 1 & 0 \end{bmatrix} \in \mathbb{R}^{5 \times 5}$$

$D_{4,2}$        $D_{4,2}^2 = 36$

# Linear Discriminant Analysis (LDA)

**Input** $A \in \mathbb{R}^{n \times d}$, also clusters

$$S_1, S_2, \dots S_k$$
$$\bigcup S_j = A \qquad S_i \cap S_j = \phi$$
$$i \neq j$$

**Goal:** Find the best

linear embedding to preserve

clusters

$\Bigg($ **Aside:**

t-SNE: find best embedding (non-linear)

that preserves cluster structure $\Bigg)$

$$\mu_i = \frac{1}{|S_i|} \sum_{x \in S_i} x \quad \text{mean} \quad \in \mathbb{R}^d$$

$$\Sigma_i = \frac{1}{|S_i|} \sum_{x \in S_i} (x - \mu_i)(x - \mu_i)^T \in \mathbb{R}^{d \times d}$$
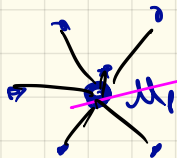
<u>covariance</u>

$$\mu = \frac{1}{|X|} \sum_{x \in X} x$$
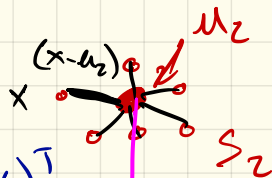
within class covariance

$$\Sigma_W = \frac{1}{|X|} \sum_{i=1}^{k} |S_i| \Sigma_i$$

$$= \frac{1}{|X|} \sum_{i=1}^{k} \sum_{x \in S_i} (x - \mu_i)(x - \mu_i)^T$$

$$\Sigma_B = \frac{1}{|X|} \sum_{i=1}^{k} |S_i| (\mu_i - \mu)(\mu_i - \mu)^T$$



$(x - \mu_2)$    $\mu_2$    $S_2$

$x$

$\mu_1$    $S_1$

$\mu$

$S_3$    $\mu_3$

$\underline{LDA}$ : 1. top $k'$ eigenvectors of

$$\boxed{\Sigma_W^{-1} \Sigma_B} \in \mathbb{R}^{d \times d}$$

$\longmapsto V_{k'}$

2. Project $\hat{X} \leftarrow V_{k'}^T X$

$$\hat{x} = V_{k'}^T x = \left( \langle x, v_1 \rangle, \langle x, v_2 \rangle, \ldots, \langle x, v_{k'} \rangle \right)$$

$\in \mathbb{R}^{k'}$      $\uparrow$ orig data point $\in \mathbb{R}^d$

top eigen vector

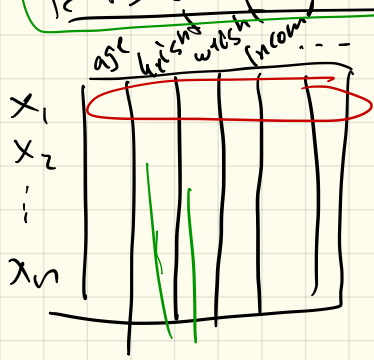$$v_1 = \underset{\|u\|=1}{\arg\max} \quad \frac{u^T \Sigma_B u}{u^T \Sigma_W u}$$

$\leftarrow$ for top eig $\Sigma_B$

$\leftarrow$ bottom eig $\Sigma_W$
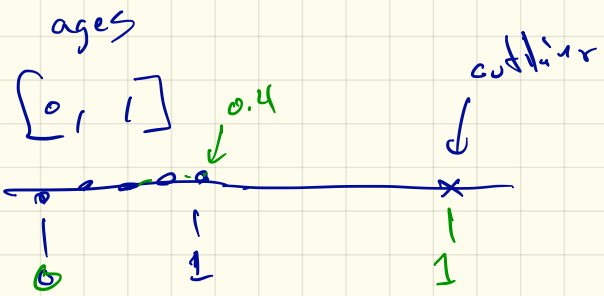
age  hrisht  wdsh  Income...

$x_1$
$x_2$
$\vdots$
$x_n$

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

$$\| a_i - a_j \| = \sqrt{(a_i - a_j)^T (a_i - a_j)}$$

$$= \sqrt{\sum_{l=1}^{d} (a_{il} - a_{jl})^2}$$

mixing units

fine

ages

$[0, 1]$   0.4

cutlier

0    1
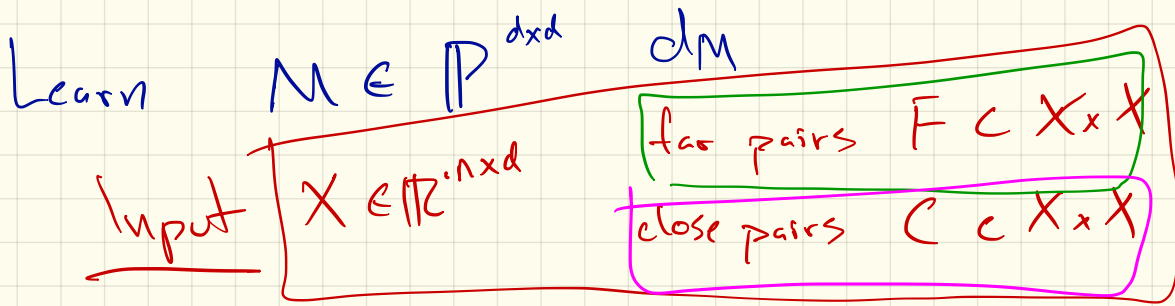
0    1    1

Learn Mahalanobis dist

$$d_M(a_i, a_j) = \| a_i - a_j \|_M = \sqrt{(a_i - a_j)^T M (a_i - a_j)}$$

$M \in \mathbb{R}^{d \times d}$
positive def.
$M \in \mathbb{P}$

Learn $M \in \mathbb{P}^{d \times d}$ $d_M$

Input $X \in \mathbb{R}^{n \times d}$

far pairs $F \subset X \times X$

close pairs $C \subset X \times X$

$$M^* = \frac{\max}{M \in \mathbb{P}} \; \frac{\min}{\{x_i, x_j\} \in F} \; d_M(x_i, x_j)^2$$

restrict
$Tr(M) = d$

$$\text{s.t.} \quad \sum_{\{x_i, x_j\} \in C} d_M(x_i, x_j)^2 \leq K$$

$$H = \sum_{\{x_i, x_j\} \in C} (x_i - x_j)(x_i - x_j)^T \in \mathbb{R}^{d \times d}$$

$$\Delta = \left\{ \alpha \in \mathbb{R}^{|F|} \;\middle|\; \sum_{i=1}^{|F|} \alpha_i = 1, \; \alpha_i \geq 0 \right\}$$

probability dist on $F$.

$H = H + \sigma I \in$ makes $H$ full rank

$T_{ij} \in F$

$X_{T_{ij}} = (x_i - x_j)(x_i - x_j)^T \in \mathbb{R}^{d \times d}$

$$\tilde{X}_T = H^{-1/2} X_T H^{-1/2}$$

argmax
$M$

$\min_{\alpha \in \Delta} \sum_{T \in F} \alpha_T \langle \hat{X}_T, M \rangle$

$$\underset{M \in \mathbb{P}}{\arg\max} \quad \underset{\alpha \in \Delta}{\min} \quad \sum_{\tau \in F}^{1} \alpha_\tau \langle \tilde{x}_\tau, M \rangle$$

$$\langle x, M \rangle = \sum_{s,t}^{1} x_{s,t} M_{s,t} \quad (\text{think of } d_M(x_\tau))$$

<u>optimize</u>   (Frank-Wolfe)

gradient

smoothing para

$\sigma = \text{eq. } 10^{-5}$

$$g_\sigma(M) = \frac{\sum_{\tau \in F}^{1} \exp\left(-\langle \tilde{x}_\tau, M \rangle / \sigma\right) \tilde{x}_\tau}{\sum_{\tau \in F}^{1} \exp\left(-\langle \tilde{x}_\tau, M \rangle / \sigma\right)}$$

$\tilde{x}_\tau \in \mathbb{R}^{d \times d}$

$\omega_\tau$

$$v_{\sigma,M} = \text{top eig}\left(g_\sigma(M)\right)$$

1. Init $M \in \mathbb{P}$ (arbitrarily $M = I$)

2. $v_t = v_{\sigma,M}$    $\leftarrow$ gradient

step 3

3. Update $M_t = \frac{t-1}{t} M_{t-1} + \frac{1}{t} v_t v_t^T$   $\in \mathbb{R}^{d \times d}$

Return $M$