

lecture 3 : data preprocessing

Principal component analysis

given a dataset $X \in \mathbb{R}^{3 \times 2}$, project the data using PCA

$$X = \begin{bmatrix} \downarrow & \downarrow \\ 2 & 4 \\ 1 & 3 \\ 0 & 0 \end{bmatrix}$$

feature A feature B

#1 center X

$$\mu_A = \frac{2+1+0}{3} = 1 \quad \mu_B = \frac{4+3+0}{3} = \frac{7}{3}$$

$$X_{\text{center}} = X - \mu = \begin{bmatrix} 1 & 5/3 \\ 0 & 2/3 \\ -1 & -7/3 \end{bmatrix}$$

#2 calculate covariance matrix $\Sigma = \frac{1}{n-1} X_{\text{center}}^T X_{\text{center}}$

$$\Sigma = \frac{1}{3-1} \cdot \begin{bmatrix} 1 & 5/3 \\ 0 & 2/3 \\ -1 & -7/3 \end{bmatrix} \cdot \begin{bmatrix} 1 & 0 & -1 \\ 5/3 & 2/3 & -7/3 \end{bmatrix}$$

$$\Sigma = \frac{1}{2} \begin{bmatrix} 2 & 4 \\ 4 & 26/3 \end{bmatrix} = \begin{bmatrix} 1 & 2 \\ 2 & 13/3 \end{bmatrix}$$

#3

compute eigenvalues & eigenvectors
(eigenvalue decomposition)

$$\Sigma = \lambda \cdot V \rightarrow \text{eigenvectors}$$

\downarrow
eigenvalues

$$\det(\Sigma - \lambda I) = 0$$

$$\Sigma - \lambda I = \begin{bmatrix} 1 - \lambda & 2 \\ 2 & \frac{13}{3} - \lambda \end{bmatrix}$$

$$\begin{aligned}\det(\Sigma - \lambda I) &= (1 - \lambda)(\frac{13}{3} - \lambda) - (2)(2) \\ &= \frac{13}{3} + \lambda^2 - \frac{16}{3}\lambda - 4 \\ &= \lambda^2 - \frac{16}{3}\lambda + \frac{1}{3}\end{aligned}$$

$$\Rightarrow \lambda_1 = -\frac{\sqrt{61} + 8}{3} \quad \text{and} \quad \lambda_2 = \frac{\sqrt{61} + 8}{3}$$

$$= 0.06325 \quad \quad \quad = 5.27008$$

(too long for him to ask us for v_1 and v_2 so im skipping to the end)

$$v_1 = \begin{bmatrix} -\frac{\sqrt{61} - 5}{6} \\ 1 \end{bmatrix}$$

$$v_2 = \begin{bmatrix} \frac{\sqrt{61} - 5}{6} \\ 1 \end{bmatrix}$$

$$\Rightarrow \Sigma v_1 = \lambda_1 v_1 \quad \text{and} \quad \Sigma v_2 = \lambda_2 v_2$$

#4

choose k largest eigenvalues (lets just do $k=1$)

$$\lambda_2 = \frac{\sqrt{61} + 8}{3} \quad \text{and} \quad v_2^T = \begin{bmatrix} \frac{\sqrt{61} - 5}{6} & 1 \end{bmatrix}$$

#5

project original data to k -dimensional space

$$Y = X_{\text{center}} \cdot w_k \quad \text{s.t. } Y \in \mathbb{R}^{n \times 1}$$

$\overset{3 \times 2}{X_{\text{center}}} \cdot \overset{2 \times 1}{w_k} \quad \overset{3 \times 1}{Y}$
 $\overset{d \times K}{w_k}$

$$w_k \in \mathbb{R}^d \quad w_k = [v_1 \quad v_2 \quad \dots]$$

$$w_k = \begin{bmatrix} \frac{\sqrt{61} - 5}{6} \\ \vdots \\ 1 \end{bmatrix}$$

$$Y = \begin{bmatrix} 1 & 5/3 \\ 0 & 2/3 \\ -1 & -7/3 \end{bmatrix} \cdot \begin{bmatrix} \frac{\sqrt{61} - 5}{6} \\ 1 \end{bmatrix}$$

$$Y = \begin{bmatrix} 2 \cdot 13504 \\ 2/3 \\ -2 \cdot 80171 \end{bmatrix} \rightarrow \text{new data w/ one feature}$$

back to #3 but with singular value decomposition

singular values

$$X_{\text{center}} = U \Sigma V^T$$

↑
left singular vectors
↓ right singular vectors (PC's)

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

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

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

$V^T \in \mathbb{R}^{d \times d}$

from the same eigenvalue
decomposition process

$$\Sigma = \begin{bmatrix} \sqrt{\lambda_1} & 0 \\ 0 & \sqrt{\lambda_2} \end{bmatrix} \quad U = \begin{bmatrix} v_1 & v_2 \end{bmatrix}$$

↗ PC's so essentially
its the same thing

$$U = X_{\text{center}}^{-1} \Sigma V^T$$

lecture 4: similarity measures

Pearson Correlation

$$A = [1, 2, 3, 4, 5] \quad B = [2, 4, 6, 8, 10]$$

$$\rho(A, B) = \frac{\text{cov}(A, B)}{\sigma_A \sigma_B}$$

$$\bar{A} = 15/5 = 3 \quad \bar{B} = 30/5 = 6$$

$$\sigma_A^2 = \frac{1}{5} \sum (a_i - \bar{A})^2 = \frac{10}{5} = 2$$

$$\sigma_B^2 = \frac{1}{5} \sum (b_i - \bar{B})^2 = \frac{40}{5} = 8$$

$$\text{cov}(A, B) = \frac{1}{5} \sum (a_i - \bar{A})(b_i - \bar{B}) = \frac{20}{5} = 4$$

$$\rho(A, B) = \frac{4}{\sqrt{2} \cdot \sqrt{8}} = 1 \rightarrow \text{perfect positive linear relationship}$$

there's also the non-linear case which we're skipping

euclidean distance

$$d(A, B) = \sqrt{\sum (a_i - b_i)^2}$$

$$A = [1 \ 2 \ 3] \quad B = [4 \ 5 \ 6]$$

$$\sum (a_i - b_i)^2 = 3^2 + 3^2 + 3^2 = 27$$

$$\sqrt{\sum (a_i - b_i)^2} = \sqrt{27} = 3\sqrt{3}$$

manhattan distance

$$d(A, B) = \sum |a_i - b_i|$$

$$A = [1 \ 2 \ 3] \quad B = [4 \ 5 \ 6]$$

$$d(A, B) = |-3| + |-3| + |-3| = 9$$

cosine similarity

$$d(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

$$A = [1 \ 2 \ 3] \quad B = [4 \ 5 \ 6]$$

$$A \cdot B = (1)(4) + (2)(5) + (3)(6) = 32$$

$$\|A\| = \sqrt{1^2 + 2^2 + 3^2} = \sqrt{14}$$

$$\|B\| = \sqrt{4^2 + 5^2 + 6^2} = \sqrt{77}$$

$$d(A, B) = \frac{32}{\sqrt{14} \cdot \sqrt{77}} \approx 0.975 \rightarrow \text{very similar}$$

if $B = [2 \ 4 \ 6] = 2A$, then $d(A, B) = 1$

jaccard similarity

$$d(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$A = \{1, 2, 3, 4, 5\} \quad B = \{4, 5, 6, 7\}$$

$$d(A, B) = \frac{2}{7} = 0.287$$