

Outline: PCA

Clustering

$$\min_{\mu \in \mathbb{R}^d} \sum_{i=1}^n \|x_i - \mu - A\theta_i\|_2^2 \quad \text{if } k=1 \Rightarrow \mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$\theta_i = u_i^\top (x_i - \mu)$

$$A = [u_1, \dots, u_n] \text{ orthonormal}$$

$$u_i = \arg \max \frac{1}{n} \sum (u_i^\top (x_i - \bar{x}))^2$$

$$u_i^\top u_i = 1$$

$$\max_{u_i} \frac{1}{n} \sum_i (u_i^\top (x_i - \bar{x}))^2 + \lambda \cdot (u_i^\top u_i - 1)$$

$$u_i^\top \underbrace{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^\top}_{S} \cdot u_i + \lambda (u_i^\top u_i - 1)$$

$$\frac{\partial}{\partial u_i} (\quad) = 2 \cdot S \cdot u_i + \lambda \cdot 2 u_i = 0 \Rightarrow S u_i = -\lambda u_i$$

General Case ($k \geq 1$)

Solution

$$\mu = \bar{x}$$

Eigenvectors
of the
sample covariance

$$\theta_i = A^\top (x_i - \bar{x})$$

$$A = [u_1, \dots, u_k]$$

$$\text{Define } \tilde{X} = \begin{bmatrix} \tilde{x}_1^\top \\ \vdots \\ \tilde{x}_n^\top \end{bmatrix}, \tilde{x}_i = x_i - \bar{x}$$

$$\min \sum_{i=1}^n \|x_i - \mu - A\theta_i\|_2^2$$

$$S = \frac{1}{n} \tilde{X}^\top \tilde{X}$$

$$A = U \Sigma V^\top \quad \text{SVD}$$

$$A^\dagger = V \Sigma^{-1} U^\top \quad \text{pseudo-inverse}$$

$$= (V^\top \Sigma^{-1} U)^\top$$

$$\Rightarrow A A^\dagger = U \Sigma V^\top V \Sigma^{-1} U^\top = U \Sigma^{-1} U^\top$$

$$I - A A^\dagger = I - U \Sigma^{-1} U^\top$$

$$(I - A A^\dagger)^2 = (I - U \Sigma^{-1} U^\top)(I - U \Sigma^{-1} U^\top) = I - U \Sigma^{-2} U^\top$$

1) Eliminate θ_i

$$\theta_i^+ = (A^\top A)^\dagger A^\top (x_i - \mu) = \boxed{A^\dagger (x_i - \mu)}$$

2) Eliminating μ

$$\sum_i \|x_i - \mu - A \cdot A^+(x_i - \mu)\|_2^2 = \sum_i \|(I - AA^+)(x_i - \mu)\|_2^2 \\ = \sum_i (x_i - \mu)^T (I - AA^+)^T (I - AA^+) (x_i - \mu) \\ \frac{\partial}{\partial \mu} (\cdot) = \sum_i 2 \cdot (I - AA^+) \cdot (x_i - \mu) = 2(I - AA^+) (\sum_i (x_i) - n\mu)$$

$$\mu^* = \frac{1}{n} \sum_i x_i.$$

3) Optimize A :

$$\min \sum_i \|(I - \tilde{A}\tilde{A}^+)(x_i - \mu)\|_2^2 \quad \text{where } \tilde{x}_i = x_i - \bar{x}.$$

$$= \min_{P: \text{projection}} \| \tilde{X} - P \tilde{X}^T \|_F^2 \quad \left(\|B\|_F = \sqrt{\sum_{i,j} B_{ij}^2} \right) \sqrt{\sum_i \|B_i\|_2^2} \\ \tilde{X} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} \quad \frac{1}{n} \tilde{X}^T \tilde{X} = \begin{bmatrix} \tilde{x}_1 & \tilde{x}_2 & \dots & \tilde{x}_n \end{bmatrix} \cdot \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_n^T \end{bmatrix} = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i \tilde{x}_i^T$$

$$P \tilde{X}^T = P \cdot \begin{bmatrix} \tilde{x}_1 & \tilde{x}_2 & \dots & \tilde{x}_n \end{bmatrix} = \begin{bmatrix} P\tilde{x}_1 & P\tilde{x}_2 & \dots & P\tilde{x}_n \end{bmatrix}$$

$$\tilde{X}^T - P \tilde{X}^T = \begin{bmatrix} \tilde{x}_1 - P\tilde{x}_1 & \dots & \tilde{x}_n - P\tilde{x}_n \end{bmatrix}$$

Theorem: Let Y have rank $r \geq k$, $Y = U \Sigma V^T$ is the SVD
Eckart-Young

$$\arg \min_{\substack{Z \in \mathbb{R}^{d \times n} \\ \text{rank}(Z)=k}} \|Y - Z\|_F = U \Sigma_k V^T \text{ where } \Sigma_k \text{ is } \Sigma \text{ with } \sigma_{k+1}, \dots, \sigma_r \text{ set to zero.}$$

(choose $A = [u_1, \dots, u_k]$ where u_i are eigenvectors of $\tilde{X}^T \tilde{X}$
 (right-singular vectors of \tilde{X})

Proof:

$$\|Y - Z\|_F = \|U \Sigma V^T - Z\|_F = \|U^T (U \Sigma V^T - Z)\|_F \\ = \|\Sigma V^T - U^T Z\|_F \\ = \|\Sigma - \frac{U^T Z V}{N}\|_F$$

$$\|Y - Z\|_F^2 = \|\Sigma - N\|_F^2 \quad \bullet \quad \bullet$$

$$= \sum_{i=1}^r (\pi_i - N_{ii})^2 + \underbrace{\sum_{i>r} N_{ii}}_{\text{error}} + \underbrace{\sum_{i \neq i} N_{ii}^2}_{\text{noise}}$$

Choose $N_{ii} = \sigma_i$ for $i=1, \dots, k$ and all others $N_{ij}=0$

$$Z = U \Sigma_k V^T \Rightarrow U^T Z V = \Sigma_k$$

(Eigen) Face Recognition

- Face similarity
 - in the reduced space
 - insensitive to lighting, expression, orientation
- Projecting new “faces”
 - everything is a face



new face

projected to eigenfaces

K-means

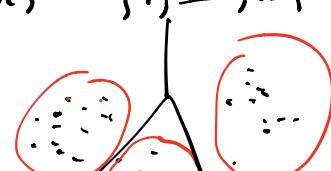
Clustering

Let $x_1, \dots, x_n \in \mathbb{R}^d$

Goal : Partition $\{x_1, \dots, x_n\}$ into disjoint subsets called “clusters”,
such that the points in the same cluster are more similar to each other

Cluster map $C: \{1, \dots, n\} \rightarrow \{1, \dots, k\}$

k is the number of clusters.



~~7.3.1~~

K-means Criterion

$$W(C) = \sum_{l=1}^k \sum_{i: C(i)=l} \|x_i - \bar{x}_l\|_2^2$$

where $\bar{x}_l = \frac{1}{n_l} \sum_{j: C(j)=l} x_j$, $n_l = \#\{i: C(i)=l\}$

K-means Algorithm minimizes $W(C)$ approximately

combinatorial optimization problem (hard)

iterative and sub-optimal algorithm

Initialize $m_1, \dots, m_k \in \mathbb{R}^d$

Repeat

For $i=1, \dots, n$

$C(i) = \arg \min_l \|x_i - m_l\|_2$

End

For $l=1, \dots, k$

$m_l = \frac{1}{|\{i: C(i)=l\}|} \sum_{i: C(i)=l} x_i$

End

Until clusters don't change

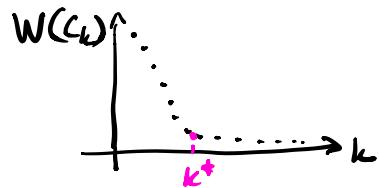
Also known as Lloyd-Max algorithm used in quantization.

Initialization: Common to randomly initialize. Also common to run the k-means algorithm several times and take the run with min $W(C)$

Model selection

How to set k ?

Simple



Idea: If $k < k^*$ $W(c_k) - W(c_{k+1})$ will be relatively large

If $k \geq k^*$ $W(c_k) - W(c_{k+1})$ will be relatively small

This suggests choosing the knee of the curve

- means ++

kernel k-means : $W(c) = \sum_{l=1}^k \sum_{i: c(i)=l} \|\Phi(x_i) - \mu_l\|_h^2$

where $\mu_l = \frac{1}{n_l} \cdot \sum_{j: c(j)=l} \Phi(x_j)$

$$\|\Phi(x_i) - \mu_l\|_h^2 = \|\Phi(x_i)\|_h^2 - 2 \langle \Phi(x_i), \mu_l \rangle + \|\mu_l\|_h^2$$

$$= \|\Phi(x_i)\|_h^2 - 2 \underbrace{\langle \Phi(x_i), \frac{1}{n_l} \cdot \sum_{j: c(j)=l} \Phi(x_j) \rangle}_{\sum_{j: c(j)=l} \langle \Phi(x_i), \Phi(x_j) \rangle} + \underbrace{\left\| \frac{1}{n_l} \cdot \sum_{j: c(j)=l} \Phi(x_j) \right\|_h^2}$$

$$\sum_{j: c(j)=l} \langle \Phi(x_i), \Phi(x_j) \rangle \\ K_{ij} = k(x_i, x_j)$$