

K-means Clustering

Supervised learning: y_i : labels (response)

Unsupervised learning: no y_i response

Note: y_i 's were used in loss, in testing, ROC/PR

$\{x_i\}_{i=1}^n$, $x_i \in \mathbb{R}^P$ is the data

Task: Learn $z_i \in \{1, \dots, K\}$ (cluster assignments)

eg clusters $C_1 = \{x_1, x_3, x_7\}$ $z_1 = 1, z_3 = 1, z_7 = 1$
 $C_2 = \{x_2, x_5, x_6\}$ $z_2 = 2, \dots$
 $C_3 = \{x_4, x_8, x_9\}$ $z_4 = 3, \dots$

also learn $m_k \in \mathbb{R}^P$, cluster centers, $k = 1, \dots, K$

Compression interpretation: summarize data

w/ only $\{z_i\}_{i=1}^n$, $\{m_k\}_{k=1}^K$

Exploration interpretation: cluster assignment

"mean something" about dist^2 of data

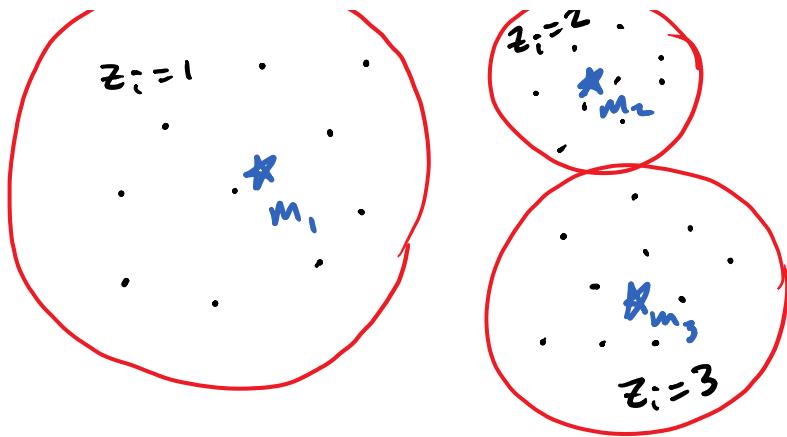
Semi supervised int.: clustering to form features for supervised learning.

K-means alg

Objective: $\min_{z, m} \sum_{i=1}^n \|x_i - m_{z_i}\|^2$

cluster assig. \uparrow centers \uparrow Distortion $J(z, m)$





$$\text{Notice, } J(z, m) = \sum_{k=1}^K \sum_{i: z_i=k} \|x_i - m_k\|_2^2$$

$$\text{So fixed } z_i, \min_{m_k} J(z, m) \Rightarrow \hat{m}_k = \frac{1}{|C_k|} \sum_{i \in C_k} x_i$$

$$\text{fixing } m_k, \min_{z_i} J(z, m) \Rightarrow z_i = \arg \min_k \|x_i - m_k\|^2$$

Lloyd's Algorithm

(1) Init m_k (randomly)

(2) Alternate

(a) Update $z_i \leftarrow \arg \min_k \|x_i - m_k\|_2^2$ for all i

(b) Update $m_k \leftarrow \frac{1}{|C_k|} \sum_{i \in C_k} x_i, C_k = \{i: z_i=k\}$

▷ J is non-increasing in iterations

▷ Finite # of configurations

⇒ Lloyd's will terminate

Hierarchical Clustering

Two main types

Agglomerative : bottom-up

Divisive : top-down

Agglomerative

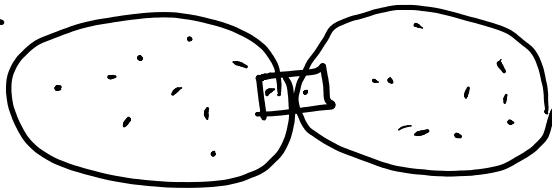
(1) Start w/ $K=n$ and $C_i = \{x_i\}$

(2) Find clusters C' & C most similar (\star)

(3) Merge clusters $C \cup C'$ repeat (2) ($K \leftarrow K-1$)

Cluster similarities

Single linkage :



$$d_{se}(C_1, C_2) = \min_{x \in C_1, y \in C_2} d(x, y)$$

Average linkage



$$d_{av}(C_1, C_2) = \frac{1}{|C_1||C_2|} \sum_{x \in C_1, y \in C_2} d(x, y)$$

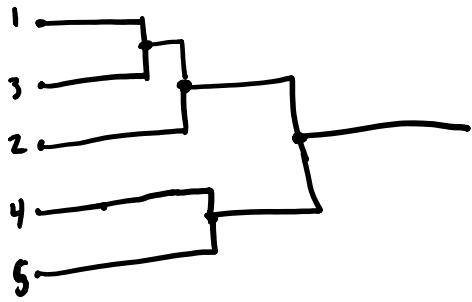
Complete linkage :

$$d_{cl}(C_1, C_2) = \max_{x \in C_1, y \in C_2} d(x, y)$$

▷ Typically, sl tends to produce unbalanced clusters

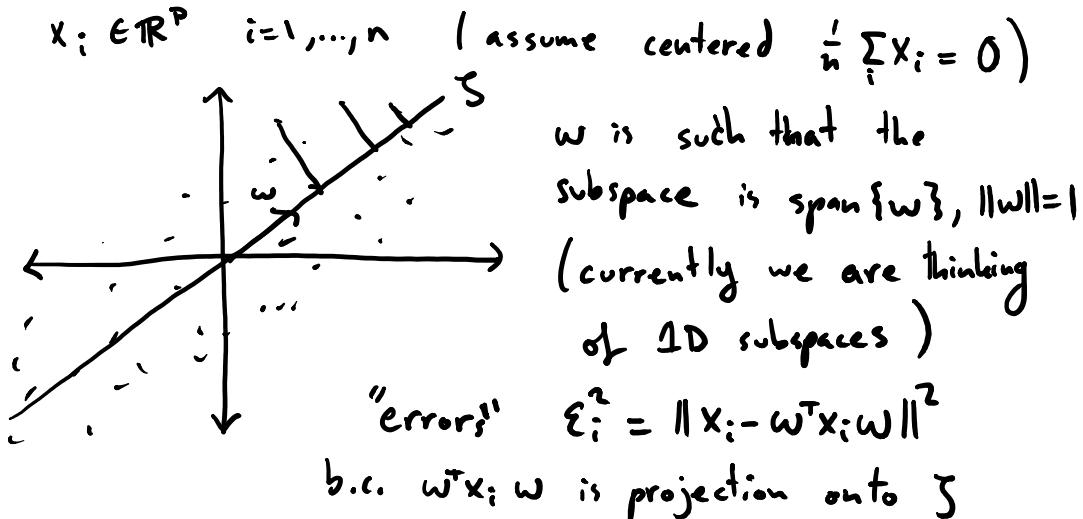
Dendrogram : visualization tool

Dendrogram : visualization tool



Principle Component Analysis

Wednesday, May 1, 2019 1:49 PM



$$\begin{aligned}\varepsilon_i^2 &= x_i^T x_i - 2 x_i^T (w^T x_i w) + (w^T x_i w)^2 \\ &= x_i^T x_i - 2 (w^T x_i) (x_i^T w) + (w^T x_i)^2 (w^T w) \\ &= x_i^T x_i - (w^T x_i)^2\end{aligned}$$

$$\min_{\|w\|=1} \sum_i \varepsilon_i^2 \Leftrightarrow \max_{\|w\|=1} \sum_i (w^T x_i)^2 \quad (+) \quad (\star)$$

$$\frac{1}{n} \sum_i (w^T x_i)^2 = \frac{1}{n} \sum_i w^T (x_i x_i^T) w = w^T \left(\frac{\sum x_i x_i^T}{n} \right) w$$

Because x_i 's are centered (\star) is covariance matrix

$$\hat{\Sigma} = \frac{1}{n} \sum_i x_i x_i^T = \frac{1}{n} X^T X$$

$$(+ \Leftrightarrow \max_{\|w\|=1} w^T \hat{\Sigma} w \quad \left(\begin{array}{l} w \text{ is eigenvector} \\ \text{corresp. to largest eigenvalue} \end{array} \right) \text{ of } \hat{\Sigma})$$

Recall Singular Value Decomp. (SVD) :

$$X = \underbrace{U}_{n \times n} \underbrace{D}_{n \times p} \underbrace{V^T}_{p \times p} \quad U^T U = I, \quad V^T V = I, \quad D \text{ is diagonal}$$

$$\hat{\Sigma} = \frac{1}{n} X^T X = \frac{1}{n} (U D V^T)^T (U D V^T) = \frac{1}{n} (V D^T U^T \underbrace{U D V^T}_I)$$

$$= \frac{1}{n} V D^T D V^T$$

d_{ii} - singular values of X
s.t. $d_{11}^2 \geq d_{22}^2 \geq \dots$

Eigenvalues of $\hat{\Sigma}$ are right eigenvectors of X
and eigenvalues of $\hat{\Sigma}$ are square of " of X

Let $\Lambda = D^T D$ then we have that λ_{ii} is the i^{th} eigenvalues of $\hat{\Sigma}$.

$$\hat{X}_k = U D_k V^T \quad \text{or} \quad D_k = \begin{pmatrix} d_{11} & & & & \\ & d_{22} & & & \\ & & \ddots & & \\ & & & d_{kk} & \\ 0 & & & & \ddots 0 \end{pmatrix}$$

result \hat{X}_k is the best rank k approximation of X
in the sense that $\sum_i \|x_i - \hat{x}_{ki}\|^2$ is smallest

$$\begin{aligned} \hat{x}_{ki} &= u_i^T D_k V^T = \sum_j u_{ij} (D_k V^T) \underset{\text{column}}{\downarrow} \\ &= \sum_j u_{ij} (\Lambda_k)_{jj} v_j = \sum_{j=1}^k u_{ij} \lambda_{jj} v_j \end{aligned}$$

K=1 $\omega = v_1$ and $x_i^T v_1$ is 1st PC loading

K=2 v_1, v_2 are basis $x_i^T v_1, x_i^T v_2$ are 1st & 2nd PC loadings

$$(XV)_{ij} = x_i^T v_j \quad X V = U D V^T V = U D$$

alg ($n > p$)

Compute $\hat{\Sigma} = \frac{1}{n} \tilde{X}^T \tilde{X}$ ($O(np^2)$)

Compute $\hat{\mathbf{X}} = \frac{1}{n} \mathbf{X}' \mathbf{X}$ ($O(n p^2)$)

Compute $\hat{\mathbf{Z}} = \mathbf{V} \Lambda \mathbf{V}^T$ (spectrum)

Compute loadings $\mathbf{x}_i^T \mathbf{v}_j$ $i = 1, \dots, n$
 $j = 1, \dots, k$

PLs have loading \downarrow and reconstruction

$$\hat{\mathbf{x}}_i = \sum_{j=1}^k (\mathbf{x}_i^T \mathbf{v}_j) \mathbf{v}_j$$