

5.2 Cluster Analysis

→ The aim is to group n objects into g classes.

- $g \ll n$

- Metric: success of clustering algorithms.

This is an example of unsupervised Learning algorithm.

5.2.1. K means clustering

* How to find the number of groups g ?

We assume g is given.

$\lambda_1, \dots, \lambda_n \in \mathbb{R}^P$: The purpose of K-means

clustering is to partition the data set into clusters

C_1, \dots, C_g with centers μ_1, \dots, μ_g as solution to

$$\min_{\substack{C_1, \dots, C_g \\ \mu_1, \dots, \mu_g \in \mathbb{R}^P}} \sum_{l=1}^g \sum_{i \in C_l} \|\lambda_i - \mu_l\|_2^2$$

When C_1, \dots, C_g are given; then

$$\mu_l = \bar{x}_l = \frac{1}{n_l} \sum_{i \in C_l} x_i$$

$$\{x_1, \dots, x_n\} \rightarrow \{1, \dots, n\}$$

$$\{1, \dots, n\} = \bigcup_{l=1}^g C_l \quad C_l \cap C_{l'} = \emptyset \quad l \neq l'$$

Algorithm. (K-means clustering) - Lloyd's algorithm

[- Given centers μ_1, \dots, μ_g , each point x_i]

is assigned to the cluster l given by

$$l = \arg \min_j \|x_i - \mu_j\|_2$$

- Update the centers as $\mu_l = \frac{1}{n_l} \sum_{i \in C_l} x_i$.

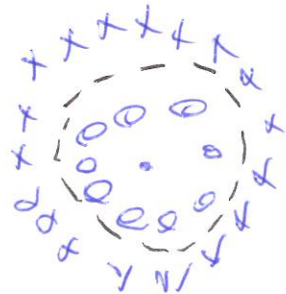
L

* Problem: we need to know the number of clusters g .

* We need a Euclidean space.

* The iterations might end up in a sub-optimal solution.

* You always end up with convex clusters.



5.2.2. Spectral Clustering

$$x_1, \dots, x_n \in \mathbb{R}^p, \quad G = (V, E, W)$$

$$\downarrow$$

$$v_i$$

$$\downarrow$$

$$v_n$$

$$w_{ij} = \frac{k}{\epsilon} (\|x_i - x_j\|) \begin{cases} \text{sym.} \\ \text{positive} \\ \text{locality} \end{cases}$$

e.g., $k_{\epsilon}(u) = \exp\left(-\frac{1}{2\epsilon} u^2\right)$

$$M = D^{-1}W \quad \text{transition matrix}$$

$$D = \text{diag}(\text{deg}(v_1), \dots, \text{deg}(v_n))$$

$$\text{deg}(v_i) = \sum_j w_{ij}$$

$$\mathbb{P}(X_{t+1}=j | X_t=i) = M_{ij} = \frac{w_{ij}}{\text{deg}(v_i)}$$

$$\mathbb{P}(X_t=j | X_0=i) = (M^t)_{ij}$$

$$M = D^{-1}W \quad \mathcal{S} = D^{1/2} M D^{-1/2} = V \Lambda V^T$$

$$\Rightarrow M = \Phi \Lambda \Psi^T$$

$$M = \Phi \Lambda \Psi^T = \sum \lambda_k \phi_k \psi_k^T$$

(Φ, Ψ) - biorthonormal.

$$v_i \rightarrow e_i^T M^t = \sum \lambda_k^t \phi_{k, si} \psi_k^T$$

ψ_k^T
= vectors

$$v_i \rightarrow \begin{pmatrix} \lambda_1^t \phi_{1, si} \\ \vdots \\ \lambda_n^t \phi_{ni} \end{pmatrix}$$

$$v_i = \begin{pmatrix} \lambda_2^t \phi_{2, si} \\ \vdots \\ \lambda_{d+1}^t \phi_{d+1, si} \end{pmatrix}$$

Spectral Clustering:

→ Use diffusion map to find the d -dimensional embedding of data.

→ Run k -means clustering on top of diffusion maps.

[Example next time for spectral clustering]

5.2.3. Hierarchical Clustering

Hierarchical clustering $\begin{cases} \rightarrow \text{agglomerative (bottom up)} \\ \rightarrow \text{divisive (top down)} \end{cases}$

Agglomerative clustering: assign a cluster

to each point - reduce the number of clusters

by iteratively merging smaller clusters into larger ones.

Divisive methods start with one large cluster containing all the datapoints and iteratively divide it into smaller clusters.

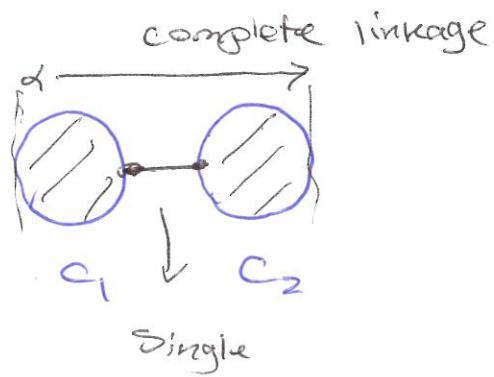
n objects v_1, \dots, v_n , pairwise dissimilarity

$$\delta_{ij} = \delta_{ji}$$

$$\Delta = (\delta_{ij}) \quad i, j = 1, \dots, n$$

A linkage function between two clusters C_1 and C_2 is defined as

$$d(C_1, C_2) = \begin{cases} \min_{i \in C_1, j \in C_2} \delta_{ij} & \text{Single linkage} \\ \max_{i \in C_1, j \in C_2} \delta_{ij} & \text{complete linkage} \\ \frac{1}{|C_1| \cdot |C_2|} \sum_{i \in C_1, j \in C_2} \delta_{ij} & \text{average linkage} \end{cases}$$



Alg. Agglomerative Clustering

1. Initialize clusters as ~~so~~ singletons
2. Initialize the set of available clusters.
 $S \leftarrow \{1, 2, \dots, n\}$

Repeat.

- Pick the two most similar clusters to merge

$$j, k = \underset{j, k \in S}{\text{argmin}} d(C_j, C_k)$$

- Merge C_k into C_j as $C_j \leftarrow C_j \cup C_k$

- Mark k as unavailable $S \leftarrow S - \{k\}$.

- Update dissimilarities $d(C_i, C_j)$

- ~~so~~ continue until no more cluster is available.

6 Support Vector Machine.

Given a training set $\{(x_1, y_1), \dots, (x_n, y_n)\}$

$$x_i \in \mathbb{R}^p \quad y_i \in \{-1, +1\}$$

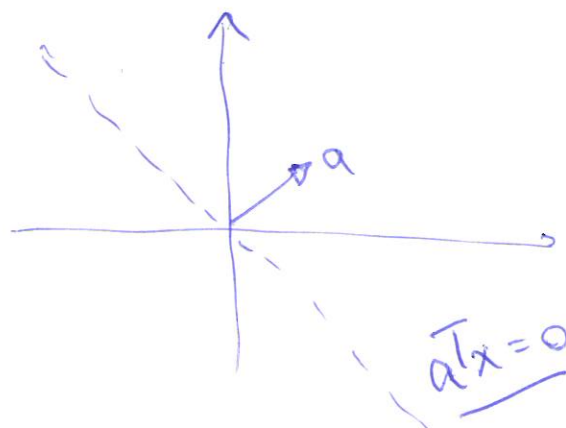
To select a hyperplane that separates the points into two classes and maximize the margin, i.e., the distance between the hyperplane and the closest points of the training set from each class.

6.1 Hyperplanes and Margins

$$a \in \mathbb{R}^p, b \in \mathbb{R}$$

a) Given $a \in \mathbb{R}^p$ $\{x \in \mathbb{R}^p \mid a^T x = 0\}$

is a $(p-1)$ -dimensional linear subspace orthogonal to a .



b) Given $a \in \mathbb{R}^p$, $b \in \mathbb{R}$

$$\{x \in \mathbb{R}^p \mid a^T x - b = 0\}$$

the linear space shifted by the vector $\frac{b}{\|a\|^2} a$.

$$a^T x - b = 0 \quad a^T(x - x_0) = a^T x - b$$

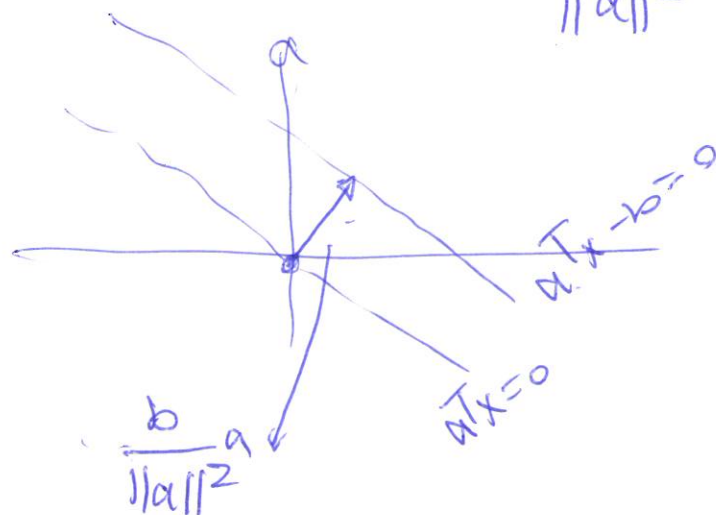
$$x_0 = \frac{ba}{\|a\|^2} \rightarrow \text{vector}$$

Hence, the linear space $\{a^T x = 0\}$ shifted by

$\frac{b}{\|a\|^2} a$ is given by $\{a^T x - b = 0\}$ and is

a hyperplane of distance $\frac{b}{\|a\|^2}$ from

$\{a^T x = 0\}$

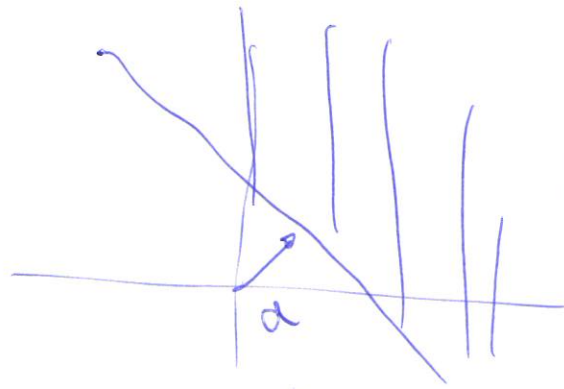


c) A half space of points lying on one side

of the hyperplane $\{a^T x - b = 0\}$ is given by

$$\{x \in \mathbb{R}^p \mid a^T x \geq b\}.$$

$$\{a^T x \geq b\}$$



$$\{a^T x \geq b\}$$

d) $H_1 = \{a^T x - b_1\}$

$$H_2 = \{a^T x - b_2\}$$

H_1, H_2 are parallel.

The distance between

two hyperplanes is

$$\text{given by } \left\| \frac{b_1}{\|a\|^2} a - \frac{b_2}{\|a\|^2} a \right\| = \frac{|b_1 - b_2|}{\|a\|}$$

$\frac{b_1}{\|a\|^2} a$

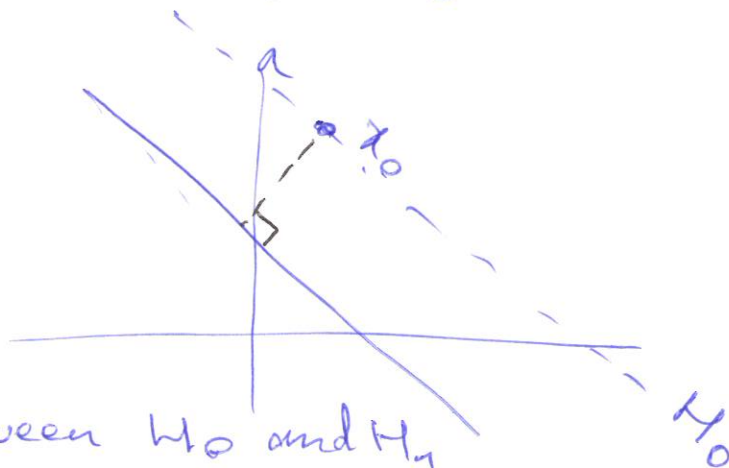
$\frac{b_2}{\|a\|^2} a$

$a^T x - b_1 = 0$ $a^T x - b_2 = 0$

e) $H_1 = \{x \mid a^T x - b = 0\} \quad x_0 \in \mathbb{R}^p$

$$H_0 = \{x \mid a^T x - b_0 = 0\}$$

$$a^T x_0 - b_0 = 0$$



The distance between H_0 and H_1

$$\frac{|b_1 - b_0|}{\|a\|} = \frac{|b_1 - a^T x_0|}{\|a\|}$$