

LECTURE 6 : INTRODUCTION TO CLUSTERING 04/11/2024

Clustering \neq Classification

↳ No-labels

↳ We have labels

Feature Selection (Supervised Learning)

① Variable Ranking

① Compare each feature with the Ground Truth classification & quantify the level of correlation

② Sort according with the correlation & choose those more correlated

• linear correlation

$$R(i) = \frac{\sum_k (x_i^k - \bar{x}_i) (y^k - \bar{y})}{\sqrt{\sigma_{x_i}^2 \sigma_y^2}}$$

R_1, \dots, R_m

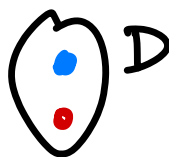
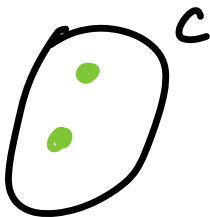
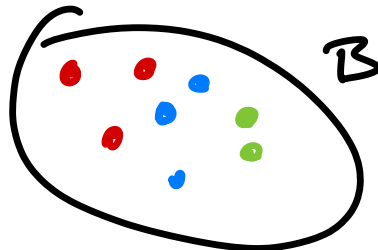
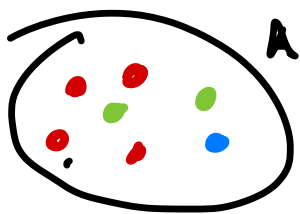
• Mutual Information

$$MI(i) = \iint p(x_i, y) \log \frac{p(x_i, y)}{p(x_i) p(y)}$$

$$= \sum_y \sum_{x_i} p(x_i, y) \log \frac{p(x_i, y)}{p(x_i) p(y)}$$

19 (Green, Blue, Red) y

(A, B, C, D) Feature x:



$$p(\text{Red}) = \frac{8}{19}$$

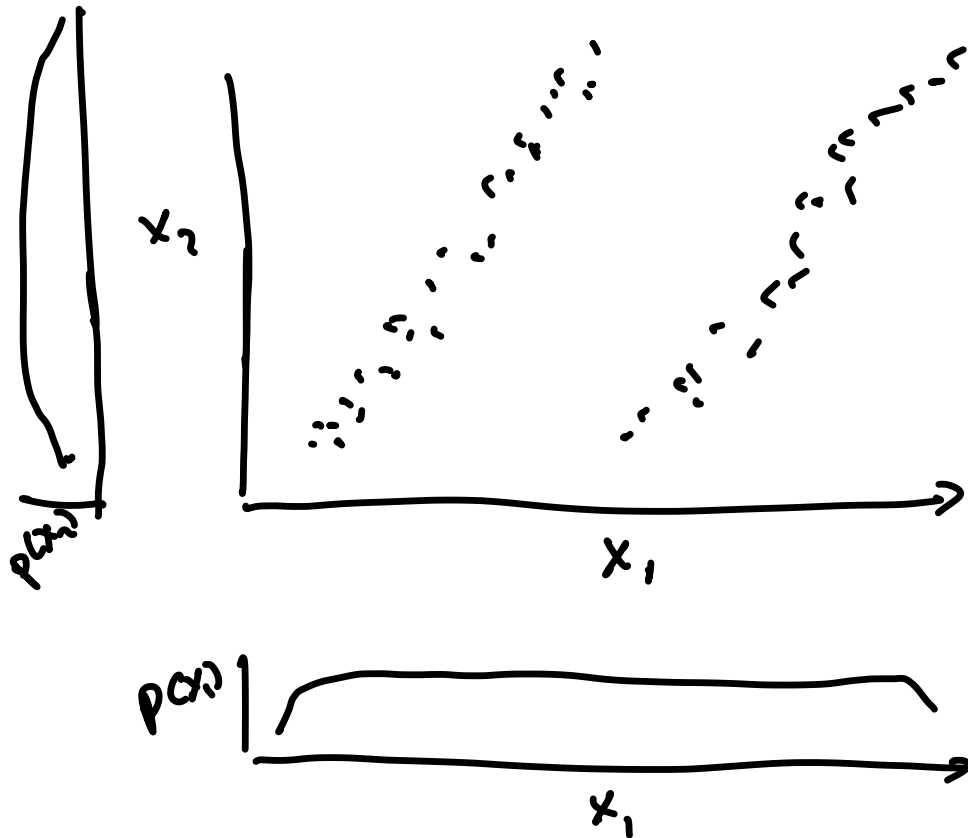
$$p(A) = \frac{7}{19}$$

$$p(\text{Red}, A) = \frac{4}{19}$$

$$\sum p \frac{4}{19} \cdot \log \frac{4/19}{7/19 \cdot 8/19} + \dots$$

Problems of variable ranking

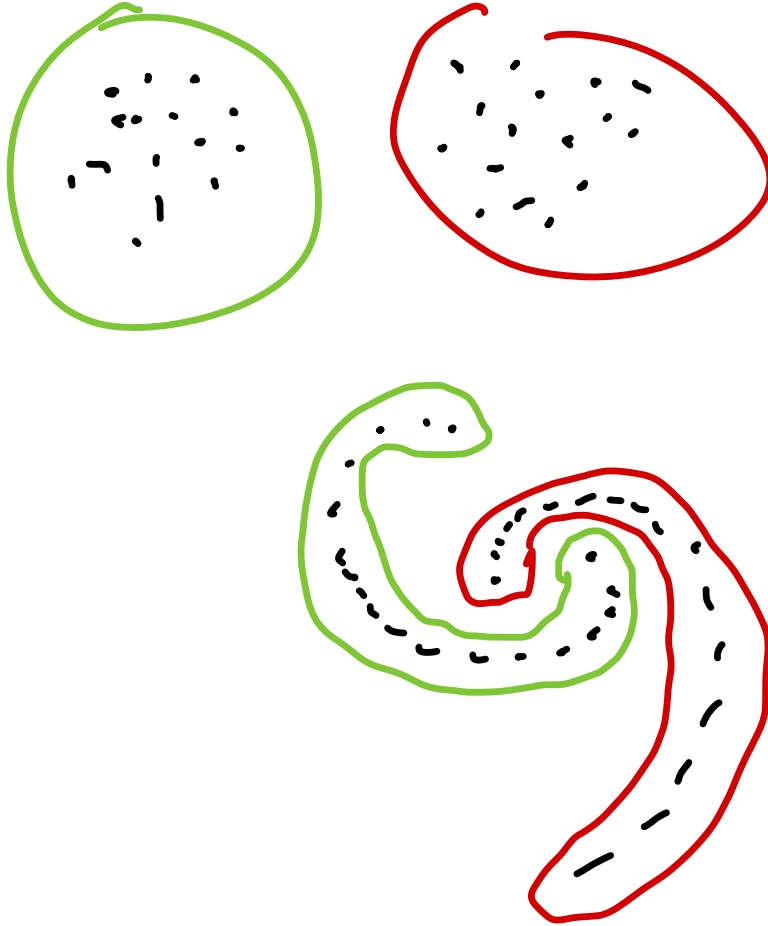
① Redundant variables



② Subset selection

- ① Explore subsets of features
- ② Rank

CLUSTERING



CLUSTERING METHODS CLASSIFIED ACCORDING TO THEIR OUTPUT

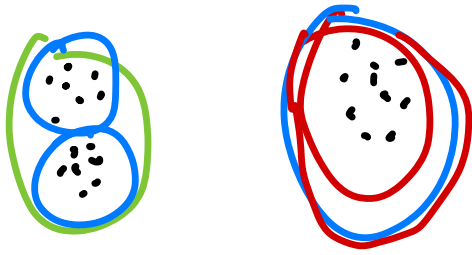
FLAT CLUSTERING : $Cl(i) = z^i$
 \uparrow integer

Fuzzy Clustering: $CL(i) = \vec{\mu}^i$
 \uparrow
 vector with k components

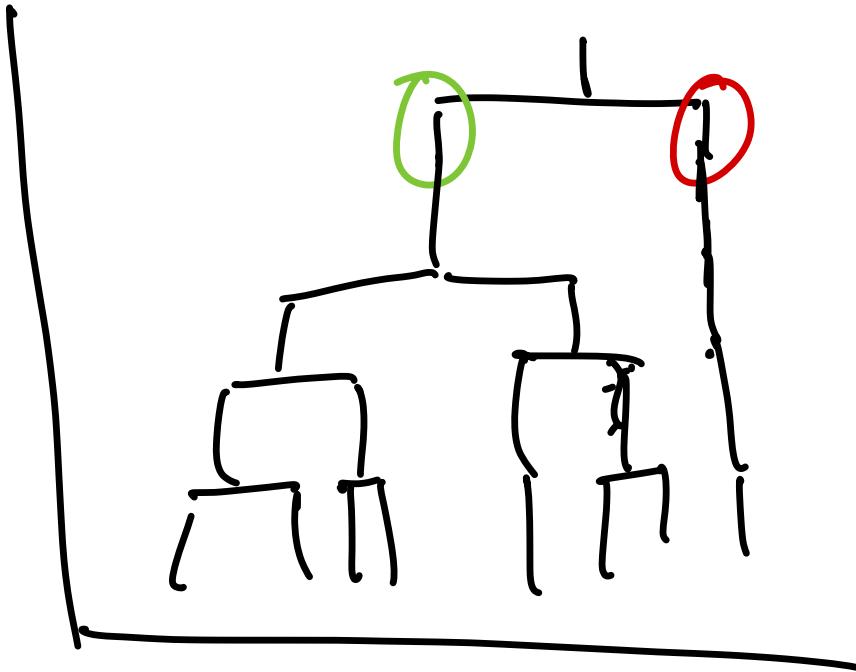
$\mu_j^i \rightarrow$ defines the degree of membership of point i to cluster j

$$\sum_{j=1}^K \mu_j^i = 1$$

Hierarchical clustering



• Dendrogram
Binary tree



K-means clustering

maximize inter cluster dissimilarity
minimizing intra cluster dissimilarity

$$c^l = \frac{\sum_{i=1}^n \delta(z^i, l) x^i}{\underbrace{\sum_{i=1}^n \delta(z^i, l)}}_n$$

Average coordinates
of points belonging
to cluster l

$$z^i = c_l(i) \quad ; \quad \delta(z^i, l) \begin{cases} 1 & \text{if } i \in l \\ 0 & \text{otherwise} \end{cases}$$

$$L(z) = \sum_{\ell}^k \sum_i^N \delta(z_i, \ell) \|x^i - c^\ell\|^2$$

Loss
of k-
means

Target: find z that minimizes $L(z)$
(NP-complete problem)

k-means algorithm

① Pick k -centers (among the data points)

② Iterate

Ⓐ Assign each data point to its nearest center
minimize intra cluster dis.

Ⓑ Recompute the centers
$$c^\ell = \frac{\sum_i \delta(z^i, \ell) x^i}{\sum \delta(z^i, \ell)}$$

maximize the distance between centers

Kernel K-means

① Find distances in the feature space

$$d_{i\ell}^2 = \|\vec{\phi}^i - \vec{\phi}^\ell\|^2 = (\vec{\phi}^i - \vec{\phi}^\ell) \cdot (\vec{\phi}^i - \vec{\phi}^\ell) =$$

$$\vec{\phi}^i \cdot \vec{\phi}^i + \vec{\phi}^\ell \cdot \vec{\phi}^\ell - 2 \vec{\phi}^i \cdot \vec{\phi}^\ell = K^{ii} + K^{\ell\ell} - 2 K^{i\ell}$$

$$d_{i\ell}^2 = \|\vec{\phi}^i - \vec{\nu}^\ell\|^2 = (\vec{\phi}^i - \vec{\nu}^\ell) \cdot (\vec{\phi}^i - \vec{\nu}^\ell)$$

↑
Centroid

$$= \vec{\phi}^i \cdot \vec{\phi}^i - 2 \vec{\phi}^i \cdot \vec{\nu}^\ell + \vec{\nu}^\ell \cdot \vec{\nu}^\ell =$$

$$K^{ii} - 2 \vec{\phi}^i \cdot \frac{\sum_j \delta(z^j, \ell) \vec{\phi}^j}{\sum_j \delta(z^j, \ell)} + \frac{\sum_m \delta(z^m, \ell) \vec{\nu}^m}{\sum_n \delta(z^n, \ell)} \cdot \frac{\sum_n \delta(z^n, \ell) \vec{\nu}^n}{\sum_n \delta(z^n, \ell)}$$

$$d_{i\ell}^2 =$$

$$K^{ii} - 2 \frac{\sum_j \delta(z^j, \ell) K^{ij}}{\sum_j \delta(z^j, \ell)} + \frac{\sum_{n,m} \delta(z^m, \ell) \delta(z^n, \ell) K_{nm}}{\left(\sum_n \delta(z^n, \ell)\right)^2}$$

① Compute the Kernel (Gram) matrix

② Choose randomly K centers

③ Iterate $z^i = \underset{\ell \in \text{cluster}}{\operatorname{argmin}} (d_{i\ell}^2)$