

CLUSTERING (Ch 10)

11/05/21

Aim of clustering: divide the dataset into subgroups / clusters

↓
homogeneous

- ↓
• similarity within groups
• dissimilarity between groups

Unsupervised learning (unknown groups)

Many applications:

- ① Genomics → Find unknown subtypes of breast cancer
→ Find groups of genes / proteins with similar functions
- ② Marketing (Segmentation) → Find clusters of shoppers in terms of their shopping history
- ③ Communities in social networks

MANY METHODS

- ① K-means clustering
- ② Hierarchical clustering
- ③ Mixture modelling : assume an underlying probability distribution



- ④ Kernel density clustering : clusters identified by a high density of points separated by areas of low density

K-MEANS CLUSTERING

↑
number of clusters fixed

Objective Split the dataset D into K clusters, C_1, \dots, C_K , that form a partition with meaningful clusters

$$\begin{array}{c} C_1 \cup \dots \cup C_K = D \\ C_i \cap C_j = \emptyset \quad \forall i \neq j \end{array}$$

So we need to define a within-cluster similarity (ideally associated to between-cluster dissimilarity)

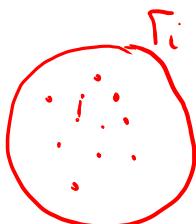
WITHIN-CLUSTER VARIATION

$$(c_1, \dots, c_k) = \underset{\substack{\Gamma_1, \dots, \Gamma_k \\ \text{partitions} \\ \text{of } D}}{\operatorname{argmin}} \sum_{i=1}^k w(\Gamma_i)$$

with

$$w(\Gamma_i) = \frac{1}{2|\Gamma_i|} \sum_{\substack{x \in \Gamma_i \\ y \in \Gamma_i}} \sum_{e=1}^p (x_e - y_e)^2$$

$\|x - y\|^2$ (Euclidean distance)



$$\text{with } \underline{x} = (x_1, \dots, x_p), \underline{y} = (y_1, \dots, y_p)$$

Why variation?

We can show that

$$W(\Gamma_i) = \sum_{x \in \Gamma_i} \sum_{e=1}^p (x_e - \bar{x}_{ie})^2 \quad \text{with } \bar{x}_{ie} = \underbrace{\frac{\sum_{x \in \Gamma_i} x_e}{|\Gamma_i|}}_{\text{sample mean of } x_e \text{ in cluster } \Gamma_i}$$

i.e. within-cluster variance.

Proof: ($p = 1$)

$$\begin{aligned} \frac{1}{|\Gamma_i|} \sum_{x \in \Gamma_i} \sum_{y \in \Gamma_i} (x-y)^2 &= \frac{1}{|\Gamma_i|} \sum_{x \in \Gamma_i} \sum_{y \in \Gamma_i} (x^2 - 2xy + y^2) \\ &= \frac{1}{|\Gamma_i|} \sum_{x \in \Gamma_i} \left[|\Gamma_i| x^2 - 2x \underbrace{\sum_{y \in \Gamma_i} y}_{\overline{|\Gamma_i| \bar{x}_i}} + \underbrace{\sum_{y \in \Gamma_i} y^2}_{|\Gamma_i| \bar{x}_i^2} \right] = \\ &= \frac{1}{|\Gamma_i|} \left[|\Gamma_i| \sum_{x \in \Gamma_i} x^2 - 2 |\Gamma_i| \bar{x}_i \underbrace{\sum_{x \in \Gamma_i} x}_{|\Gamma_i| \bar{x}_i} + |\Gamma_i| \sum_{y \in \Gamma_i} y^2 \right] \\ &= 2 \sum_{x \in \Gamma_i} x^2 - 2 |\Gamma_i| \bar{x}_i^2 \\ &= 2 \left(\sum_{x \in \Gamma_i} x^2 - |\Gamma_i| \bar{x}_i^2 \right) = 2 \left(\sum_{x \in \Gamma_i} (x - \bar{x}_i)^2 \right) \end{aligned}$$

Minimizing within-cluster variance \Leftrightarrow Maximising between-cluster variation

In order to see this:

- ① Since the overall sum of squared distances in the data is fixed, minimizing within-cluster distances means maximising between-cluster distances
- ② This is connected to the variance decomposition formula:

$$\text{Var}(X) = \underbrace{\mathbb{E}(\text{Var}(X|C))}_{\text{average of variances within clustering}} + \underbrace{\text{Var}(\mathbb{E}(X|C))}_{\text{variance of averages within groups} \rightarrow \text{distance between centroids}}$$

with
 $C = \{1, \dots, K\}$
cluster assignment

(unexplained variance) (explained variance)



IMPLEMENTATION

For too many possible clusterings $\rightarrow \approx k^n$

ALGORITHM (Fixed K)

① Random Initialization : randomly assign each observation to one of K clusters

Typically :

- ① Uniformly for each observation (most common)
- ② Randomly in space

② Iterate until convergence :

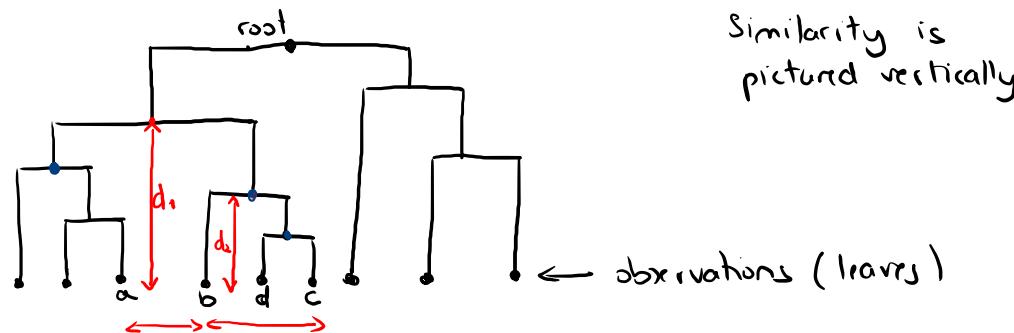
a) For each cluster, compute centroid \bar{x}_i

b) Assign each observation 1,..,n to the cluster whose centroid is closest (in Euclidean distance)

HIERARCHICAL CLUSTERING

A family of methods

- There is no prior assumption on the number of clusters
- The output consists of a hierarchy of nested clusters, represented by a dendrogram
- It can be built bottom-up (agglomerative approach) or top-down (divisive approach)
- Typical dendrogram:



Two main ingredients :

① A measure of similarity between points

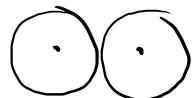
② A linkage criterion : a measure of distance between clusters

The results will depend on these choices, which can be also driven by the data-domain

CHOICE OF METRIC

Different choices:

- ① Euclidean distance (L_2 distance): $\sum_{e=1}^p (x_e - y_e)^2$
- ② Taxicab distance (L_1 distance): $\sum_{e=1}^p |x_e - y_e|$
- ③ Max (or L_∞ distance): $\max_e |x_e - y_e|$
- ④ Mahalanobis distance: $\sqrt{(\underline{x}-\underline{y})^T \Sigma^{-1} (\underline{x}-\underline{y})}$
- ⑤ Correlation: $\text{Cor}(\underline{x}, \underline{y})$



Text mining

- Levenshtein distance : # edits to mutate one string into another (between words)

CUP $\xrightarrow{\substack{\downarrow \\ \text{one edit}}}$ CAR

- Distance between documents :

$$d(x, y) = \cos(t(x), t(y)) = \frac{t(x)^T t(y)}{\|t(x)\| \|t(y)\|}$$

where

$$t(x) = (\underbrace{tf\text{-}idf}_{\text{across a list of words }}(\omega, x, C))$$

$$tf(\omega, x) = \frac{\# \text{ occurrences of word } \omega \text{ in a document } x}{\# \text{ words in a document}}$$

$$idf(\omega, C) = \log \left(\frac{\# \text{ documents in corpus } C}{\# \text{ documents in } C \text{ containing } \omega} \right)$$

$$tf\text{-}idf(\omega, x, C) = tf(\omega, x) \cdot idf(\omega, C)$$

CHOICE OF LINKAGE

Common choices:

① Single linkage

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

↓
choose metric between observations

→ tends to produce long thin clusters (chain)



② COMPLETE LINKAGE

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

→ compact clusters but very sensitive to outliers

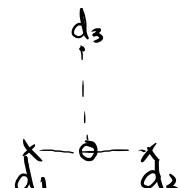
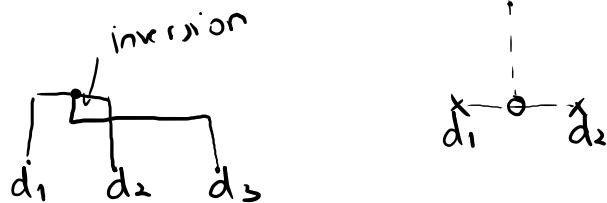
③ Average Linkage

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

④ Centroid Linkage

$$d_{CE}(C_1, C_2) = d(\bar{x}_1, \bar{x}_2) \quad \text{with } \bar{x}_i \text{ centroid of cluster } C_i$$

→ Non monotone



IMPLEMENTATION

AGGLOMERATIVE

- Start with each observation in one cluster
- At each step join observations / clusters that are closer to each other
- Stop after n steps

DIVISIVE

- Start with everything in one cluster C
- Pick observation most dissimilar from the rest and put it into a new cluster \tilde{C}
- Move any observation in C that is closest to \tilde{C} to \tilde{C} → Two clusters
- Repeat until you have one observation in each cluster.

HOW MANY CLUSTERS?

Several choices:

- ① "Lazy": Cut tree at pre-defined level of cluster similarity
- ② "Lazy-K": fix number K of clusters
- ③ LARGEST GAP: cut the tree before the joinings that decreases the quality of the clustering the most (a bit like a scree plot)

IN PRACTICE

- ① Careful with scaling
- ② Careful with outliers
- ③ Only choose clustering as a pre-processing /exploratory analysis step