

Unsupervised clustering

$$X = \begin{bmatrix} x_1 & e_1 \\ \vdots & \vdots \\ x_n & e_n \end{bmatrix} \quad x_i \in \mathbb{R}^p$$

$e_i \in \{1, \dots, g\}$

But: e_i are hidden & isn't known

- Goal:
- estimate g
 - estimate e_1, \dots, e_n

A general idea: units belonging to the same cluster (sharing the same label) are more similar than units belonging to different

To implement this idea \Rightarrow specify a dissimilarity

$d(x, y)$ dissimilarity between $x, y \in \mathbb{R}^p$

$d: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty]$

Properties:

1. $d(x, x) = 0 \quad \forall x \in \mathbb{R}^p$

2. $d(x, y) = 0 \iff x = y \quad \forall x, y \in \mathbb{R}^p$

3. $d(x, y) = d(y, x) \quad \forall x, y \in \mathbb{R}^p$

4. $d(x, y) \leq d(x, z) + d(z, y) \quad \forall x, y, z \in \mathbb{R}^p$

If d satisfies 1, 2, 3 metric

4. $d(x, y) \leq \max \{d(x, z), d(z, y)\}$

If d satisfies 1, 2, 4 ultra metric

$$x = \{\alpha, \beta, \gamma, \delta, \varepsilon, \zeta, \theta\}$$

$$y = \{\alpha, \beta, \gamma, \omega, \nu, \mu\}$$

$\hookrightarrow a = n$ first symbol where the 2 seq x and y are diff.

$$d(x, y) = \frac{1}{2n} \quad \text{ULTRAMETRIC}$$

Short list of trivial distances

$$x, y \in \mathbb{R}^p$$

$$\bullet \quad d(x, y) = \sqrt{(x-y)^T(x-y)} \quad \text{Euclidean distance}$$

Often x, y have been std.

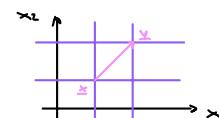
$$\bullet \quad d(x, y) = \sqrt{(x-y)^T \Sigma (x-y)}$$

Σ covariance $p \times p$ estimated by $\hat{\Sigma}$ pooled

$$\bullet \quad d^r(x, y) = \sum_{i=1}^p |x_i - y_i|^r \quad \text{Minkowski's dist } (r)$$

$r=2 \rightarrow$ Euclidean dist

$r=\infty \rightarrow$ Manhattan dist



$$x, y \in (\mathbb{R}^+)^p$$

$$\bullet \quad d(x, y) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i} \quad \text{Cauchy dist}$$

$$x, y \in \{0, 1\}^p \quad x = (0, 1, 1, 0, 0)$$

$$\bullet \quad d_{\text{cc}}(x, y) = \sum_{i=1}^p (x_i - y_i)^2 = \# \text{ discordances}$$

Example

Row data = {Blue, Brown, Green}

$$x = (x_1, x_2) \quad (0, 1) - \text{Blue}$$

$$(1, 0) - \text{Brown}$$

$$(0, 0) - \text{Green}$$

$$x, y \in \{0, 1\}^p - \text{Contingency Matrix}$$

	1	0	
1	a	b	
0	c	d	
			e

$$d^2_{\text{Eclud}}(x, y) = a+b$$

$$d(x, y) = e - \frac{a+b}{e}$$

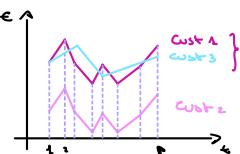
$$d_Q: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty)$$

$$d_C: \{0, 1\}^p \times \{0, 1\}^p \rightarrow [0, +\infty)$$

$$d = \lambda d_Q + (1-\lambda) d_C \quad \lambda \in [0, 1]$$



EXAMPLE

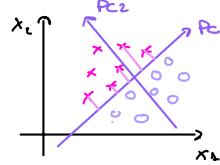
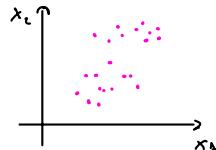


Studler exclusion distance, we might think that they are similar, but if we want to know how they respond to a sale, Cost 1 and 2 are similar instead. Depends on what we want to know.

$$\mathbf{x} = (x_1, \dots, x_p)$$

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{2(1 - \text{corr}(\mathbf{x}, \mathbf{y}))} \quad \leftarrow \text{Clust variables?}$$

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1q} \\ x_{21} & & & \\ \vdots & & & \\ x_{n1} & & & x_{nq} \end{bmatrix} \quad \left\{ \begin{array}{l} \text{clust units} \\ \text{clust var.} \end{array} \right.$$



$$d \text{ is specified} \implies D = [d_{ij}] \quad n \times n \quad \text{dist matrix}$$

$$d_{ij} = d(\mathbf{x}_i, \mathbf{x}_j) \quad i, j = 1, \dots, n$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_i \in \mathbb{R}^p$$

$$D = \begin{bmatrix} d_{ii} & & \\ & \ddots & d_{12} \\ d_{21} & & d_{ii} \end{bmatrix} \quad \text{if } d \text{ is reflexive, } \text{is symmetric}$$

Distance between finite subsets of \mathbb{R}^p

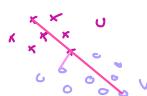
U, V finite subsets of \mathbb{R}^p , ? $d(U, V)$

- single linkage

$$d(U, V) = \min \{ d(\mathbf{x}_i, \mathbf{y}_j) : \mathbf{x}_i \in U, \mathbf{y}_j \in V \}$$

- Complete linkage

$$d(U, V) = \max \{ d(\mathbf{x}_i, \mathbf{y}_j) : \mathbf{x}_i \in U, \mathbf{y}_j \in V \}$$



- Average linkage

$$d(U, V) = \frac{1}{|U| \cdot |V|} \sum_{\mathbf{x}_i \in U} \sum_{\mathbf{y}_j \in V} d(\mathbf{x}_i, \mathbf{y}_j)$$

Hierarchical Agglomeration (Clust. Agglomeration)

- D distance

- linkage

Initialization:

- every unit is a cluster

Until convergence repeat

Step 1 Cluster together the two closest clusters

Step 2 Update D

Example

$n=5$

$$D = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 1 & 0 & & & \\ 2 & 2 & 0 & & \\ 3 & 6 & 5 & 0 & \\ 4 & 10 & 3 & 4 & 0 \\ 5 & 9 & 8 & 5 & 3 & 0 \end{bmatrix}$$

Round 1: Cluster $\{1, 2\}$ & $\{3\}$ in $\{1, 2, 3\}$, $d(\{1, 2\}, \{3\}) = 2$

$$D_{1,2} = \begin{bmatrix} \{1, 2\} & 3 & 4 & 5 \\ \{1, 2\} & 0 & & \\ 3 & 5 & 0 & \\ 4 & 3 & 4 & 0 \\ 5 & 8 & 5 & 3 & 0 \end{bmatrix}$$

Round 2: Cluster $\{4\}$ & $\{5\}$, $d(\{4\}, \{5\}) = 3$

$$D_2 = \begin{bmatrix} \{1, 2\} & 3 & \{4, 5\} \\ \{1, 2\} & 0 & \\ 3 & 5 & 0 & \\ \{4, 5\} & 8 & 4 & 0 \end{bmatrix}$$

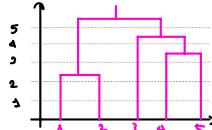


Round 3: Cluster $\{3,7\}$ & $\{4,5\}$, $d(\{3,7\}, \{4,5\}) = 4$

$\{3,7\}$	0	$\{3,4,5\}$
$\{3,4,5\}$	5	0

Round 4: Cluster $\{4,2\}$ & $\{3,4,5\}$ $d(\{4,2\}, \{3,4,5\}) = 5$

Ardenogram to summarize the entire process:



Obs: jitter the data first i.e. introduce some small noise

Cophenetic dist: $d_c(x, y) = \text{dist } d \text{ when } x \approx y \text{ merge in the cluster}$
 (ultra metric)



Ward's method for hierarchical clustering

$$\mathbb{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad \text{training set} \quad x_i \in \mathbb{R}^p$$

$$d: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty)$$

$$d^2(x, y) = (x - y)^T (x - y)$$

C cluster of points in \mathbb{R}^p (finite)

$$\bar{x} = \arg \min_{x \in \mathbb{R}^p} \sum_{z \in C} d^2(x, z) = \frac{1}{|C|} \sum_{z \in C} z \quad (\text{barycenter})$$

for $C_i \quad i=1, \dots, k$:

- C_i cluster of points in \mathbb{X}

- \bar{x}_i barycenter of C_i

- $ESS_i = \sum_{z \in C_i} d^2(z, \bar{x}_i) = \sum_{z \in C_i} (z - \bar{x}_i)^T (z - \bar{x}_i)$

- If \mathbb{X} has been split in clusters C_1, \dots, C_k : $ESS = ESS_1 + ESS_2 + \dots + ESS_k$

Algorithm

1. Initialization: Every point in the cluster $ESS = 0$

2. Until convergence, aggregate the two clusters that produce the minimum increase in ESS

⇒ dendrogram

Non-hierarch. clustering - K-means

Q: Cluster the training set in K (given) subsets C_1, \dots, C_K .

$$1. \bigcup_{i=1}^K C_i = \mathbb{X}$$

$$2. C_i \cap C_j = \emptyset \quad i, j = 1, \dots, K \quad i \neq j$$

d: $\mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty)$ metric or pseudo metric

C cluster of points in \mathbb{R}^p (finite)

$$\bar{x} = \arg \min_{x \in \mathbb{R}^p} \sum_{z \in C} d^2(z, x) \quad \text{centroid}$$

Goal: given K find C_1, \dots, C_K satisfying 1+2 and s.t. $\sum_{i=1}^K \sum_{z \in C_i} d^2(z, \bar{x}_i)$ is minimal OPTIMIZATION PROBLEM

Algorithm (sols proposed by K-means)

1. Initialization step: - Split at random \mathbb{X} in K subset (run different initializations because we need to know if it's a local min)

- or average at random K centroids

2. Until convergence repeat: S1. Compute $\bar{x}_1, \dots, \bar{x}_K$ s.t. $\bar{x}_i = \arg \min_{x \in \mathbb{R}^p} \sum_{z \in C_i} d^2(z, x) \quad i=1, \dots, K$

S2. If $\mathbb{X} \approx \bar{x}$, assign z to cluster $i=1, \dots, K$ if $d^2(z, \bar{x}_i) = \min \{d^2(z, \bar{x}_1), \dots, d^2(z, \bar{x}_K)\}$

stop when $\bar{x}_1, \dots, \bar{x}_K$ do not change

Obs: Step 2 is the difficult one

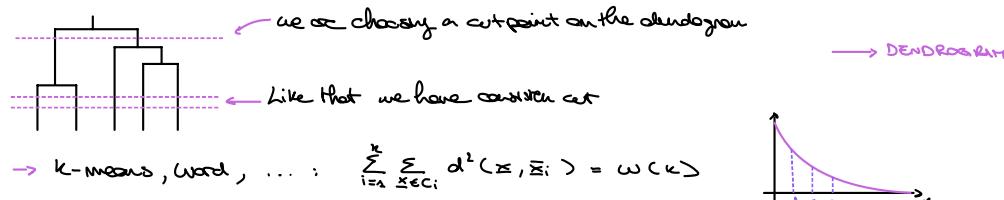
Simplification: for $i=1, \dots, K$ $\bar{x}_i = \arg \min_{x \in \mathbb{R}^p} \sum_{z \in C_i} d(z, x) \quad (\text{K-medoids})$



Pros and Cons

- Complete link, aver. link, k-means, Ward's, ...
 - ⇒ Create ellipsoidal clusters (wrt d)
- Single link
 - ⇒ chain effect, but explores structures that are not ellipsoids

How to choose k?



Non parametric density based Clustering - DBSCAN

Idea: Clusters are regions of high density

Low density regions are boundaries between clusters, or noise, or outliers.

$$\mathbb{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ training set } x_i \in \mathbb{R}^p$$

$$d: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty) \text{ metric}$$

For $\varepsilon > 0$, $x \in \mathbb{R}^p$ $N_\varepsilon(x) = \{y \in \mathbb{R}^p : d(x, y) < \varepsilon\}$ ε -neigh of x

For $i = 1, \dots, n$ $N_\varepsilon(x_i)$ $|N_\varepsilon(x_i)| = n$ of points in \mathbb{X} falling in $N_\varepsilon(x_i)$

Two parameters to be fixed in advance: 1. $\varepsilon > 0$ the radius

2. minPts $\geq \varepsilon$ an integer

Point classes in \mathbb{X} :

- x_i is a core point if $|N_\varepsilon(x_i)| \geq \text{minPts}$
- x_i is a border point if it's not a core point but it belongs to an ε -neigh of a core point
- x_i is noise if it's not core or border



Definition

To identify dense contiguous regions:

- x_i is directly density reachable from x_j if:
 - x_i is core
 - $x_i \in N_\varepsilon(x_j)$
- x_i is density reachable from $x_j \in \mathbb{X}$ if there is y_1, \dots, y_k $k \geq 1$, s.t.:
 - y_1, \dots, y_{k-1} are core
 - $y_k = x_i$, $y_{k-1} = x_j$
 - $y_k \in N_\varepsilon(y_{k-1})$
- $x_i \in \mathbb{X}$ and $x_j \in \mathbb{X}$ are density connected if there is an $x \in \mathbb{X}$ s.t. both x_i and x_j are density reachable from x .

DBSCAN identifies the $C \subseteq \mathbb{X}$ s.t.:

1. If x_j is density reach from x_i and $x_i \in C \Rightarrow x_j \in C$
2. For all couples x_i and x_j in C , x_i and x_j must be density connected

Multidimensional scaling (MDS)

$$\mathbb{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \text{ training set, } x_i \in \mathbb{R}^p \quad i = 1, \dots, n$$

$$d: \mathbb{R}^p \times \mathbb{R}^p \rightarrow [0, +\infty), \text{ dis} = d(x_i, x_j) \text{ for } i, j = 1, \dots, n$$

$D(n \times n)$ dissimilarity matrix

Question: $\mathbb{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad x_i \in \mathbb{R}^q \quad i = 1, \dots, n$

$d: \mathbb{R}^q \times \mathbb{R}^q \rightarrow [0, +\infty), \text{ dis}(x_i, x_j) = \delta_{ij} \quad \text{for } i, j = 1, \dots, n \quad \text{s.t. } \delta_{ii} = \text{dis}$

One-dimensional space with same distance as before

Solutions not unique, rotation, translation, etc., so we talk about family of configurations.

Classical MDS: find $\mathbb{Z}_1, \dots, \mathbb{Z}_n$ s.t. $\sum_{i,j} (d(x_i, x_j) - \delta_{ij})^2$ (if d is euclidean = PCA)

Kruskale: find $\mathbb{Z}_1, \dots, \mathbb{Z}_n$ s.t. STRESS = $\sum_{i,j} \frac{(\Theta(d_{ij}) - \delta_{ij})^2}{\sum \delta_{ij}}$ is minimized (also w.r.t. $\Theta: \mathbb{R} \rightarrow \mathbb{R}$ monotonically ↑)

