# Clustering

- major method in unsupervised learning:
  - find groups of similar instances $X_i \approx X_j$ $\Rightarrow$ put them into clusters
    $\hookrightarrow$ "similar"
  - two benefits: • we can analyze data one cluster at a time
    $\Rightarrow$ if alg. complexity $O(N^P)$ $P > 1$ $O(c \cdot N_1^{\frac{P}{1}})$
    $\uparrow$ entire data
    for clusters $\sum_k O(N_k^P) < O(N^P)$
    • we can analyze global behaviour by looking at cluster representatives instead of individual members
    $\Rightarrow$ simpler problem structure, faster, better interpretability for humans

- idea: $TS$ $\{X_i\}_{i=1}^{N}$, assume that labels $Y_i$ exist, but are unknown
  "latent" or "hidden" labels

  $\Rightarrow$ tasks: ① determine labels and clustering simultaneously
  ② find a good representative for each cluster, e.g.
  - mean: $\bar{X}_k = \arg\min_X \sum_{i \in C_k} (X_i - x)^2 = \frac{1}{N_k} \sum_{i \in C_k} X_i$

  - median: $\text{median}(C_k) = \arg\min_X \sum_{i \in C_k} \|X_i - X\|_1$    no simple solution in dimensions $D \geq 2$
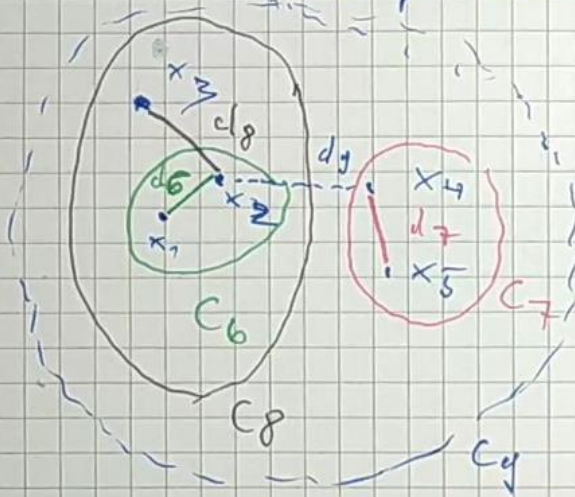
  - medoid: $\text{medoid}(C_k) = \arg\min_{X \in \{X_i : i \in C_k\}} \sum_{i \in C_k} \|X_i - X\|_1$

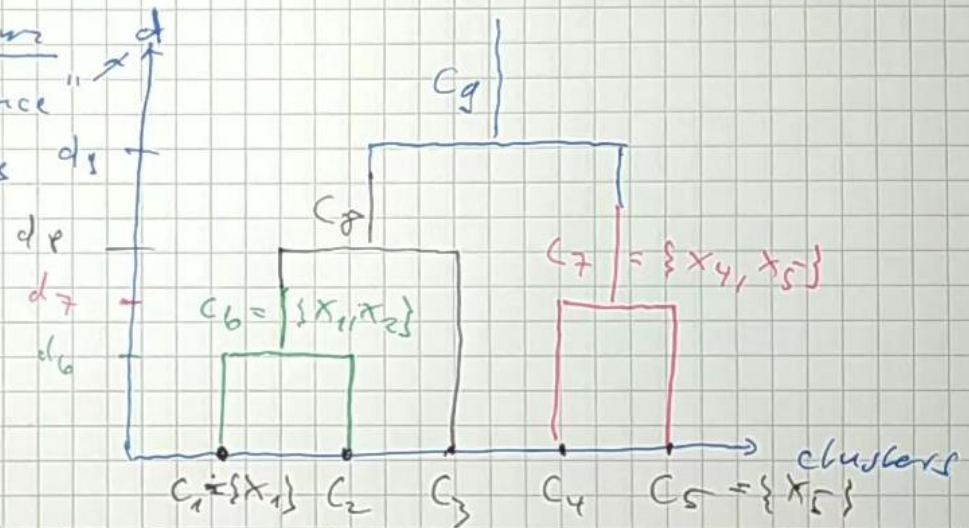  $\Rightarrow$ representative is chosen only among the cluster members

# Hierarchical Clustering

⓪ init: each instance is a cluster of its own : $k = 1, ..., N$    $C_k = \{X_k\}$

① repeat until all instance are in the same cluster: $(N-1)$ iterations
   - find the two existing clusters closest to each other
   - merge these two clusters into a new one    $(k = N+1, ..., 2N-1)$



### Dendrogram

"merge distance"
= dist. where a
pair of clusters
got merged

example:
single linkage
clustering

By stopping the algorithm early (when a minimum number of clusters is reached, or when the next merge edge exceeds a given threshold) ⟹ control the granularity of the clustering (hyperparameter)

Define distance between clusters: given (hyperparameter): distance between instances $d(X_i, X_{i'})$

- single linkage: $d(C_k, C_{k'}) = \min_{i \in C_k, i' \in C_{k'}} d(X_i, X_{i'})$

- complete linkage: $d(C_k, C_{k'}) = \max_{i \in C_k, i' \in C_{k'}} d(X_i, X_{i'})$

- average linkage: $d(C_u, C_{u'}) = \dfrac{1}{N_u N_{u'}} \sum\limits_{i \in C_u} \sum\limits_{i' \in C_{u'}} d(x_i, x_{i'})$

- representative linkage $d(C_u, C_{u'}) = d(\text{repr.}(C_u), \text{repr.}(C_{u'}))$

properties:
- single linkage is not robust against outliers $\Rightarrow$ "chaining"



$\Rightarrow$ complete linkage more robust? against this problem

chain incorrectly connects the clusters

- complete linkage faits may fail on the "closeness property" $=$ within cluster distance are less than between cluster distances $\Rightarrow$ single linkage better

- average and representative linkage are in between

- if data clusters well, all criteria give similar results

---

### K-means clustering  $K \stackrel{\wedge}{=}$ predefined number of clusters (here: C)

- consider total distance of TS:  $d_{total} = \dfrac{1}{2} \sum\limits_{i=1}^{N} \sum\limits_{i'=1}^{N} d(x_i, x_{i'})$

- any given clustering separates edges into "within" and "between" cluster edges

$$d_{total} = d_{in} + d_{betw.} \qquad d_{in} = \dfrac{1}{2} \sum\limits_{u=1}^{C} \sum\limits_{i \in C_u} \sum\limits_{i' \in C_u} d(x_i, x_{i'})$$

$$d_{betw.} = \dfrac{1}{2} \sum\limits_{u=1}^{C} \sum\limits_{i \in C_u} \sum\limits_{i' \notin C_u} d(x_i, x_{i'})$$

- a good clustering minimizes $d_{in}$ or maximize $d_{betw.} = d_{total} - d_{in}$ (both objectives are equivalent)

- exact optimization requires exhaustive search over all possible clustering
  ⟹ too expensive
- k-means alg. simple heuristic $d_h$ to minimize $d_{in}$ for special choice

$$d(x_i, x_{i'}) = \| x_i - x_{i'} \|_2^2$$

⟹ $d_{in}$ simplifies:

$$d_{in} = \sum_{k=1}^{c} N_k \sum_{i \in C_k} \| x_i - \bar{x}_k \|_2^2$$

$$\bar{x}_k = \frac{1}{N_k} \sum_{i \in C_k} x_i \qquad \text{mean of cluster } k$$

- introduce hidden labels: $Y_i$ : $Y_i = k$ means $x_i \in C_k$

  when means were known: $Y_i = \arg\min_k \| x_i - \bar{x}_k \|_2^2$

  ⟹ assign $x_i$ to the nearest cluster representative

⟹ optimize alternatingly:

⓪ define initial guess for cluster centers $\bar{x}_k^{(0)}$

① repeat until labels $Y_i$ do no longer change (always happens after finite steps)

    ⓐ update labels $\quad Y_i^{(t)} = \arg\min_k \| x_i - \bar{x}_k^{(t-1)} \|_2^2$

    ⓑ update means $\quad \bar{x}_k^{(t)} = \frac{1}{N_k^{(t)}} \sum_{i \in Y_i^{(t)} = k} x_i \qquad N_k^{(t)} = \# \{ Y_i^{(t)} = k \}$

- converges to a local optimum of $d_{in}$ (not the global one in general)
  ⟹ quality of result critically depends on initial guess

- k-means++ : improved initial guess

  ① choose $\bar{X}_1^{(0)} \sim \{X_i\}_{i=1}^N$ uniformly at random

  ② for $k = 2, \ldots, C$ : (other cluster centers)

  - define distances to closest existing center :

    $$\forall i : \quad d_i = \min_{k' \in 1, \ldots, k-1} d(X_i, \bar{X}_{k'}^{(0)})$$

  - define probabilities by normalization: $\bar{d} = \frac{\not{d_i}}{\not{\sum} d_i}$ $p_i = \frac{\not{d}}{\not{d}} \frac{d_i}{\sum_i d_i}$

  - choose $\bar{X}_k^{(0)} \sim$ discrete ( $\{p_i\}$ )

    ( already chosen points have $d_i = 0$, $p_i = 0 \Rightarrow$ will not be picked again

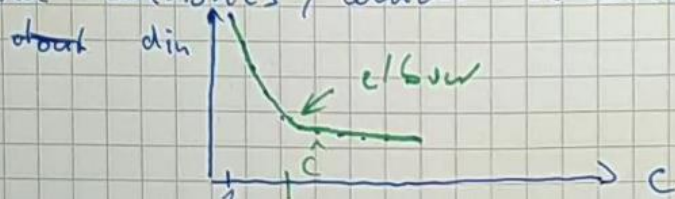    $p_i$ large $\hat{=}$ $d_i$ relatively large $\hat{=}$ point $i$ far away from existing centers

    $\Rightarrow \bar{X}_k^{(0)}$ are well spread out over the training distribution )

- varient: k-medoids : instead of mean, represent each cluster by its

  medoid $\Rightarrow$ works for any distance $d(X_i, X_{j'})$

  but more expensive than k-means

- unsolved: how to choose optimal number $C$ of clusters $\triangleq$ hyperparameter

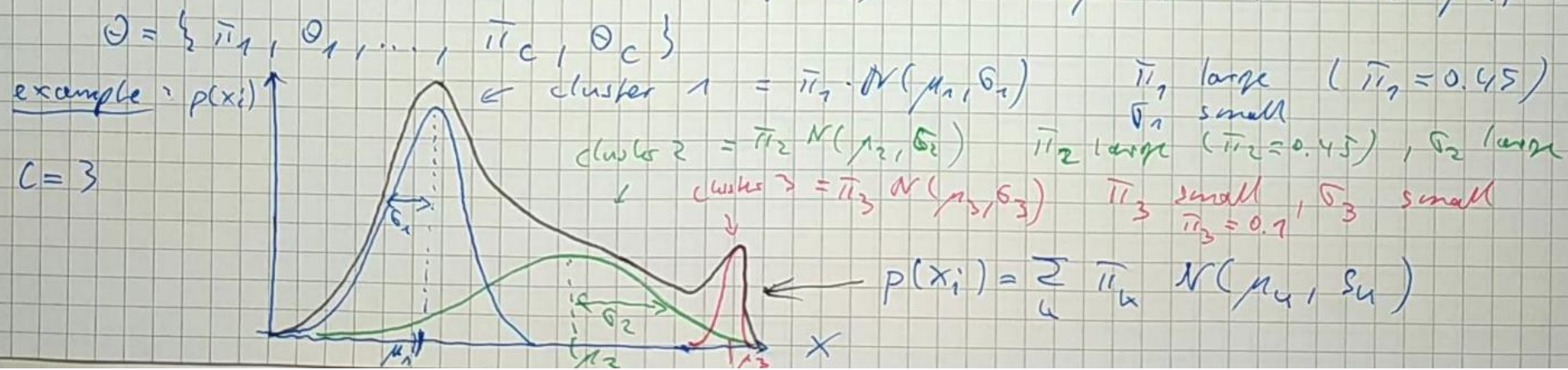  - several heuristics, which sometimes work, for example elbow method

## EM algorithm  (expectation - maximization)

- k-means recovers <u>hard cluster assignments</u> $Y_i$ = analogous to discriminative model in classification

- EM alg corresponds to a generative model:
  - soft assignments: for each $X_i$, define posterior $p(Y_i = 4 \mid X_i)$
  - model posterior by learning RHS of Bayes formula $p(Y_i \mid X_i) = \frac{p(X_i \mid Y_i) p(Y_i)}{p(X_i)}$

$\Rightarrow$ <u>mixture model</u>: $p(Y_i = 4)$ probab. that a point belongs to cluster k

$$p(X_i \mid Y_i = 4) \quad \text{shape of cluster } 4$$

as in QDA, but the labels $Y_i$ are now hidden in TS

most common: $p(X_i \mid Y_i = 4) = N(\mu_4, S_4) \Rightarrow$ Gaussian mixture model
GMM

$\Rightarrow$ solution has parameters $\underbrace{p(Y = 4) = \pi_4}_{\text{mixture weights}} \approx \frac{N_4}{N} \qquad \sum_4 \pi_4 = 1$

$$p(X \mid Y = 4; \Theta_4) = N(\mu_4, S_4) \qquad \Theta_4 = \{\mu_4, S_4\}$$

$$\Theta = \{\pi_1, \Theta_1, \dots, \pi_c, \Theta_c\}$$

example: $p(x_i)$

$C = 3$



cluster 1 $= \pi_1 \cdot N(\mu_1, S_1)$    $\pi_1$ large $(\pi_1 = 0.45)$
                                         $S_1$ small
cluster 2 $= \pi_2 N(\mu_2, S_2)$    $\pi_2$ large $(\pi_2 = 0.45)$, $S_2$ large
cluster 3 $= \pi_3 N(\mu_3, S_3)$    $\pi_3$ small, $S_3$ small
                                         $\pi_3 = 0.1$

$$p(x_i) = \sum_4 \pi_4 \, N(\mu_4, S_4)$$

optimize via maximum likelihood principle:

$$\hat{\Theta} = \arg\max_{\Theta} \prod_i p(x_i; \Theta) \iff \hat{\Theta} = \arg\max_{\Theta} \sum_i \log p(x_i; \Theta)$$

$$= \arg\max_{\Theta} \sum_i \log \underbrace{\sum_u \pi_u N(\mu_u, S_u)}_{\Theta_u}$$

$$\frac{\partial \text{Loss}}{\partial \Theta_u} = \sum_i \frac{1}{\partial \Theta_u} \log p(x_i; \Theta) = \sum_i \frac{1}{p(x_i; \Theta)} \frac{\partial}{\partial \Theta_u} p(x_i; \Theta)$$

$$= \sum_i \frac{1}{p(x_i; \Theta)} \frac{\partial}{\partial \Theta_u} \cancel{\log} \pi_u p(x_i; \Theta_u) = \sum_i \frac{\pi_u}{p(x_i; \Theta)} \frac{\partial}{\partial \Theta_u} p(x_i; \Theta_u)$$

$$= \sum_i \frac{\pi_u \, p(x_i; \Theta_u)}{p(x_i; \Theta)} \cdot \frac{1}{p(x_i; \Theta_u)} \frac{\partial}{\partial \Theta_u} p(x_i; \Theta_u)$$

$$= \sum_i \frac{\pi_u \, p(x_i; \Theta_u)}{p(x_i; \Theta)} \frac{\partial}{\partial \Theta_u} \log p(x_i; \Theta_u)$$

$$\left[ \text{Bayes:} \quad \pi_u = p(Y_i = u) \qquad \pi_u \, p(x_i; \Theta_u) = p(x_i, Y_i = u; \Theta_u) = p(x_i, Y_i = u; \Theta) \right.$$

$$\left. = p(Y_i = u \mid x_i; \Theta_u) \, p(x_i \mid \Theta) \right]$$

$$\boxed{\frac{\partial \text{Loss}}{\partial \Theta_u} = \sum_i p(Y_i = u \mid x_i; \Theta) \cdot \frac{\partial}{\partial \Theta_u} \log p(x_i \mid \Theta_u) \overset{!}{=} 0}$$

we cannot solve analytically for $\Theta_u$, because $\Theta_u \in \Theta$, the two dependencies cannot be rewritten in closed form $\Rightarrow$ <u>alternating optimization</u>

EM alg.: alternating optimization of the loss:

(0) initialization: choose number of clusters $C$ (hyperparameter)

initial cluster parameters $\Theta_k^{(0)}$ $(= \mu_k^{(0)}, S_k^{(0)})$

$\pi_k^{(0)}$

(1) for $t = 1, ..., T$ or until convergence

(a) define "responsibilities": auxiliary variables $\gamma_{ik}$ = how well is instance $i$ explained by cluster $k$ in the current guess

$$\gamma_{ik} = p(Y_i = k \mid x_i, \Theta_k^{(t-1)}) = \frac{\pi_k^{(t-1)} \cdot N(x_i \mid \mu_k^{(t-1)}, S_k^{(t-1)})}{\sum_{k'=1}^{C} \pi_{k'}^{(t-1)} \cdot N(x_i \mid \mu_{k'}^{(t-1)}, S_{k'}^{(t-1)})}$$

(b) E-step: optimize $\pi_k$ as the expectation of $\gamma_{ik}$

$$\pi_k^{(t)} = \frac{1}{N} \sum_i \gamma_{ik}$$

(c) M-step: maximize data likelihood with $\pi_k^{(t)}$ fixed

$\Rightarrow$ choose $\Theta_k = \{\mu_k, S_k\}$ to fit the data as well as possible

(equivalent to QDA, but with guessed soft class labels $\gamma_{ik}$)

$\Rightarrow$ weighted mean and covariance:

$$\mu_k^{(t)} = \frac{\frac{1}{N} \sum_i \gamma_{ik} x_i}{\underbrace{\frac{1}{N} \sum_i \gamma_{ik}}_{\pi_k^{(t)}}} \qquad S_k^{(t)} = \frac{\frac{1}{N} \sum_i \gamma_{ik} (x_i - \mu_k^{(t)})^T (x_i - \mu_k^{(t)})}{\underbrace{\frac{1}{N} \sum_i \gamma_{ik}}_{\pi_k^{(t)}}}$$