

## k-means optimisation

The  $\checkmark$  step makes sense:

(Step 1+2) Given clustering  $C = \{C_\ell\}_{\ell=1}^K$

$$W = \frac{1}{2} \sum_{\ell=1}^K \frac{1}{|C_\ell|} \sum_{i,j \in C_\ell} \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2$$

with centroids  $\vec{m}_\ell = \frac{1}{|C_\ell|} \sum_{i \in C_\ell} \vec{x}^{(i)} \quad \ell = 1, \dots, K$

$$\begin{aligned} W &= \frac{1}{2} \sum_{\ell=1}^K \frac{1}{|C_\ell|} \sum_{i,j \in C_\ell} \|(\vec{x}^{(i)} - \vec{m}_\ell) - (\vec{x}^{(j)} - \vec{m}_\ell)\|^2 \\ &= \sum_{\ell=1}^K \frac{1}{2} \frac{1}{|C_\ell|} \sum_{i,j \in C_\ell} \left[ \|\vec{x}^{(i)} - \vec{m}_\ell\|^2 + \|\vec{x}^{(j)} - \vec{m}_\ell\|^2 - 2(\vec{x}^{(i)} - \vec{m}_\ell) \cdot (\vec{x}^{(j)} - \vec{m}_\ell) \right] \\ &= \sum_{\ell=1}^K \left[ \sum_{i \in C_\ell} \|\vec{x}^{(i)} - \vec{m}_\ell\|^2 - (\vec{x}^{(i)} - \vec{m}_\ell) \cdot \underbrace{\frac{1}{|C_\ell|} \sum_{j \in C_\ell} (\vec{x}^{(j)} - \vec{m}_\ell)}_{=0} \right] \end{aligned}$$

In summary:

$$\begin{aligned} W &= \frac{1}{2} \sum_{\ell=1}^K \frac{1}{|C_\ell|} \sum_{i,j \in C_\ell} \|\vec{x}^{(i)} - \vec{x}^{(j)}\|^2 \\ &= \sum_{\ell=1}^K \sum_{i \in C_\ell} \|\vec{x}^{(i)} - \vec{m}_\ell\|^2 \end{aligned}$$

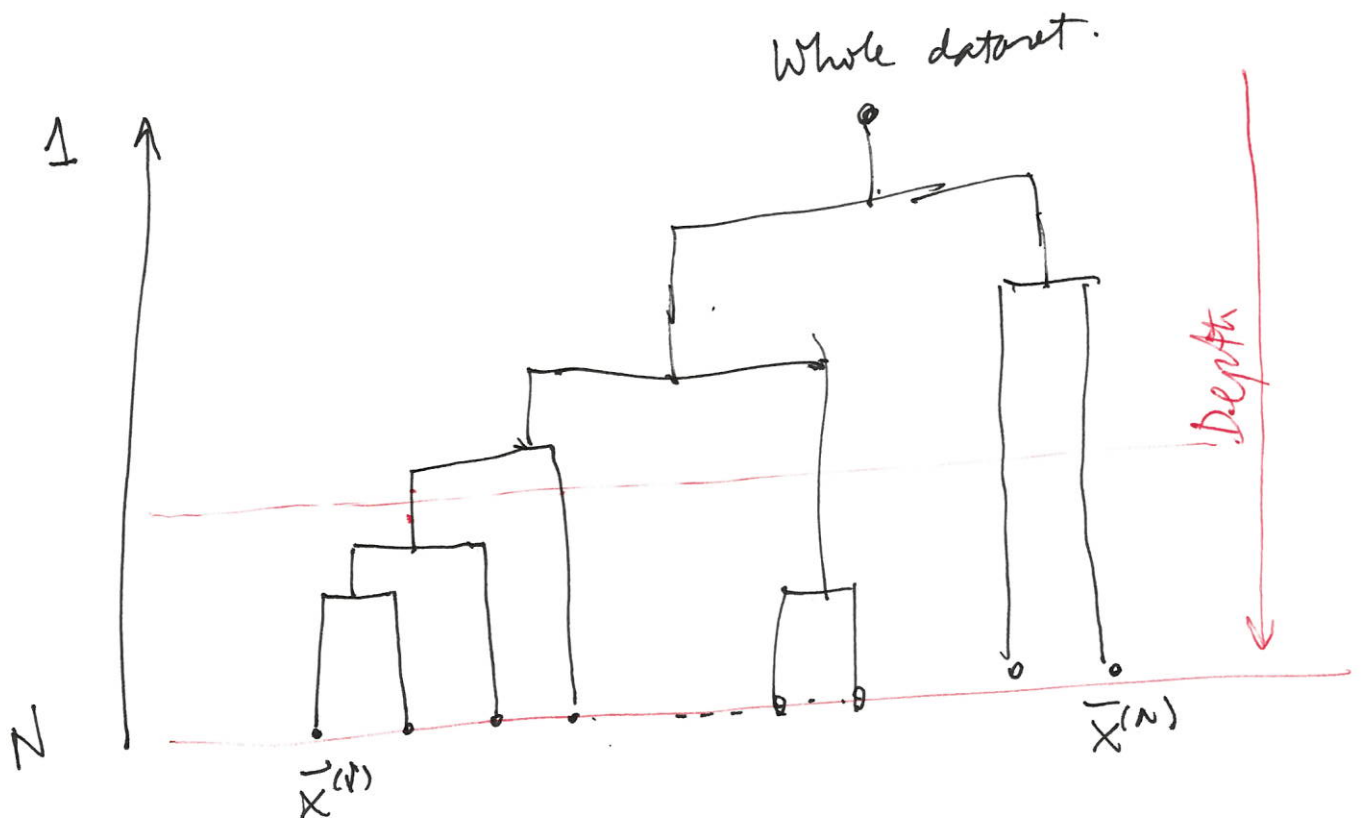
So trying to minimise the distance to the centroids is our objective.

## (2) Hierarchical clustering

$$\{\vec{X}^{(i)}\} \quad i=1, \dots, N$$

$$D_{N \times N} = (D_{ij}) \quad D_{ij} = d(\vec{X}^{(i)}, \vec{X}^{(j)})$$

Represents the whole data set as a binary tree ('dendrogram') where leaves are samples and the root is the whole dataset.



- Hierarchical structure imposed on the data because of the requirement of binary splits.
- Monotonicity relating <sup>height</sup> ~~depth~~ and dissimilarity within cluster.

Agglomerative schemes:

Reduce the number of clusters 1 by 1  
from  $N$  to 1.

(1) Simple linkage (SL):

Two clusters  $G, H$  :  $d_{SL}(G, H) = \min_{\substack{i \in G \\ j \in H}} D_{ij}$   
Nearest neighbour

(2) Complete linkage (CL):

Farthest neighbour :  $d_{CL}(G, H) = \max_{\substack{i \in G \\ j \in H}} D_{ij}$

(3) Group average:

$$d_{GA}(G, H) = \frac{1}{N_G N_H} \sum_{\substack{i \in G \\ j \in H}} D_{ij}$$

Disadvantage:

Recursive  $k$ -means with  $k=2$  gives a binary tree also but:

(1) monotonicity is not preserved

(2) it depends on the initialisations at every split.

So it can be quite variable.

Better schemes exist but, in general, agglomerative algorithms are used more extensively.