

Expectation - Maximization Clustering

$$\mathcal{X} = \{x_i\}_{i=1}^N$$

$$\text{likelihood} \Rightarrow L(\Phi | \mathcal{X}) = \left[\prod_{i=1}^N p(x_i | \Phi) \right]$$

$$\log L(\Phi | \mathcal{X}) = \sum_{i=1}^N \log \left[\underbrace{\sum_{k=1}^K p(x_i | C_k) \cdot \text{Pr}(C_k)}_{\text{mixture densities}} \right]$$

two sets of random variables

Z = cluster memberships

Φ = parameters $[\hat{\mu}_1, \hat{\mu}_2, \dots, \hat{\mu}_K, \hat{\Sigma}_1, \hat{\Sigma}_2, \dots, \hat{\Sigma}_K]$

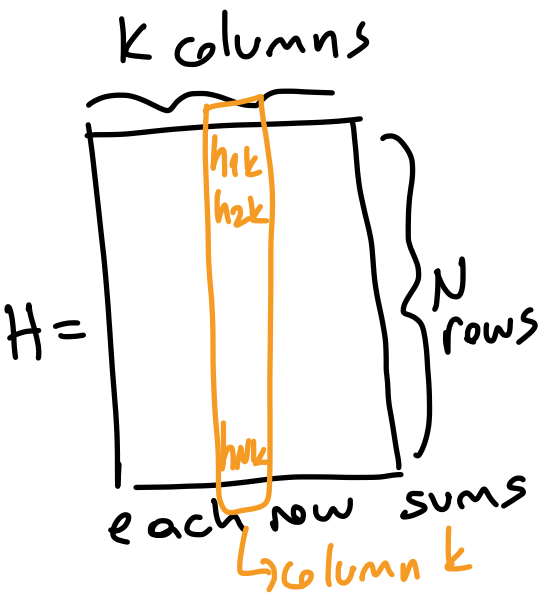
E-STEP: $E[L_c(\Phi | \mathcal{X}, Z) | \mathcal{X}, \Phi^{(+)}]$

M-STEP: $\Phi^{(t+1)} = \arg \max_{\Phi} E[L_c(\Phi | \mathcal{X}, Z) | \mathcal{X}, \Phi^{(+)}]$

E-STEP:

$$h_{ik} = E[z_{ik} | \mathcal{X}, \Phi^{(+)}] = \frac{p(x_i | C_k, \Phi^{(+)}) Pr(C_k)}{\sum_{c=1}^K p(x_i | C_c, \Phi^{(+)}) Pr(C_c)}$$

multivariate
Gaussian
density



$$h_{ik} \geq 0 \quad \forall (i, k)$$

$$\sum_{k=1}^K h_{ik} = 1 \quad \forall i$$

M-STEP:

$$\hat{Pr}^{(+1)}(C_k) = \frac{\sum_{i=1}^N h_{ik}}{N} \quad \forall k$$

$$\hat{\mu}_k^{(+1)} = \frac{\sum_{i=1}^N h_{ik} \cdot x_i}{\sum_{i=1}^N h_{ik}} \quad \forall k$$

$$\hat{\Sigma}_k^{(+1)} = \frac{\sum_{i=1}^N h_{ik} (x_i - \hat{\mu}_k^{(+1)}) (x_i - \hat{\mu}_k^{(+1)})^T}{\sum_{i=1}^N h_{ik}}$$

Hierarchical clustering

- finding groups such that instances (data points) in a group are more similar to each other than instances in different groups.
- Closer

COMPONENT #1: The distance function between pairs of data points

distance \Rightarrow dissimilarity

squared distance

distance $\uparrow \Rightarrow$ similarity \downarrow
distance $\downarrow \Rightarrow$ similarity \uparrow

$$k(x_i, x_j) = \exp \left[- \frac{\|x_i - x_j\|_2^2}{2\sigma^2} \right]$$

\downarrow similarity = 1
similarity = 0
distance = 0

distance = ∞

$$\left(\sum_{d=1}^D (x_{id} - x_{jd})^p \right)^{1/p}$$

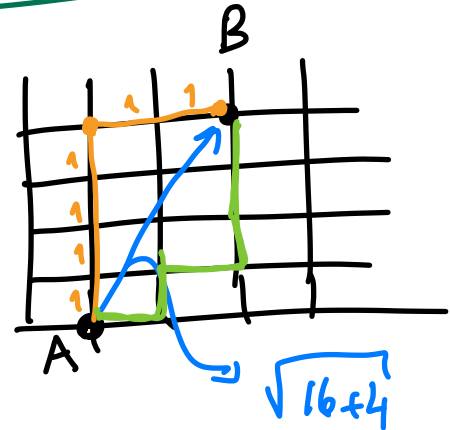
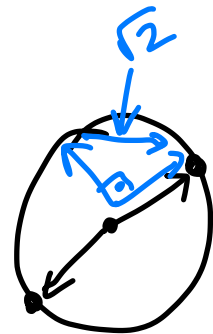
$$d(x_i, x_j) = \sum_{d=1}^D |x_{id} - x_{jd}|$$

Manhattan Distance

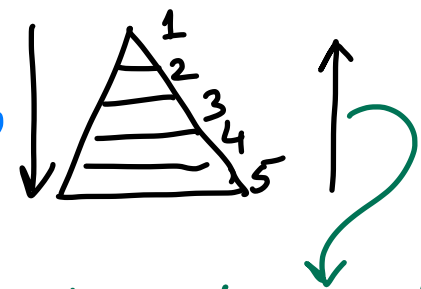
Euclidean Distance

$$d(x_i, x_j) = \|x_i - x_j\|_2 = \sqrt{(x_i - x_j)^T (x_i - x_j)}$$

$$= \sqrt{x_i^T x_i - 2x_i^T x_j + x_j^T x_j}$$



COMPONENT #2: The direction to proceed



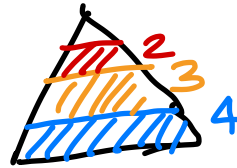
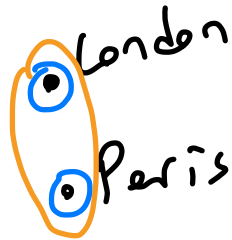
Divisive (top-to-bottom)

- divides big clusters into smaller ones
- starts with "1" cluster

Agglomerative (bottom-to-top)

- combines small clusters into bigger ones
- starts with "N" clusters

COMPONENT #3: The distance function between groups of data points



distance ({ London, Paris }, New York)

distance ({ London, Paris }, Rome)

distance (New York, Rome)

Centroid clustering $d(G_A, G_B) = \left\| \frac{\sum_{x_i \in G_A} x_i}{|G_A|} - \frac{\sum_{x_j \in G_B} x_j}{|G_B|} \right\|_2$

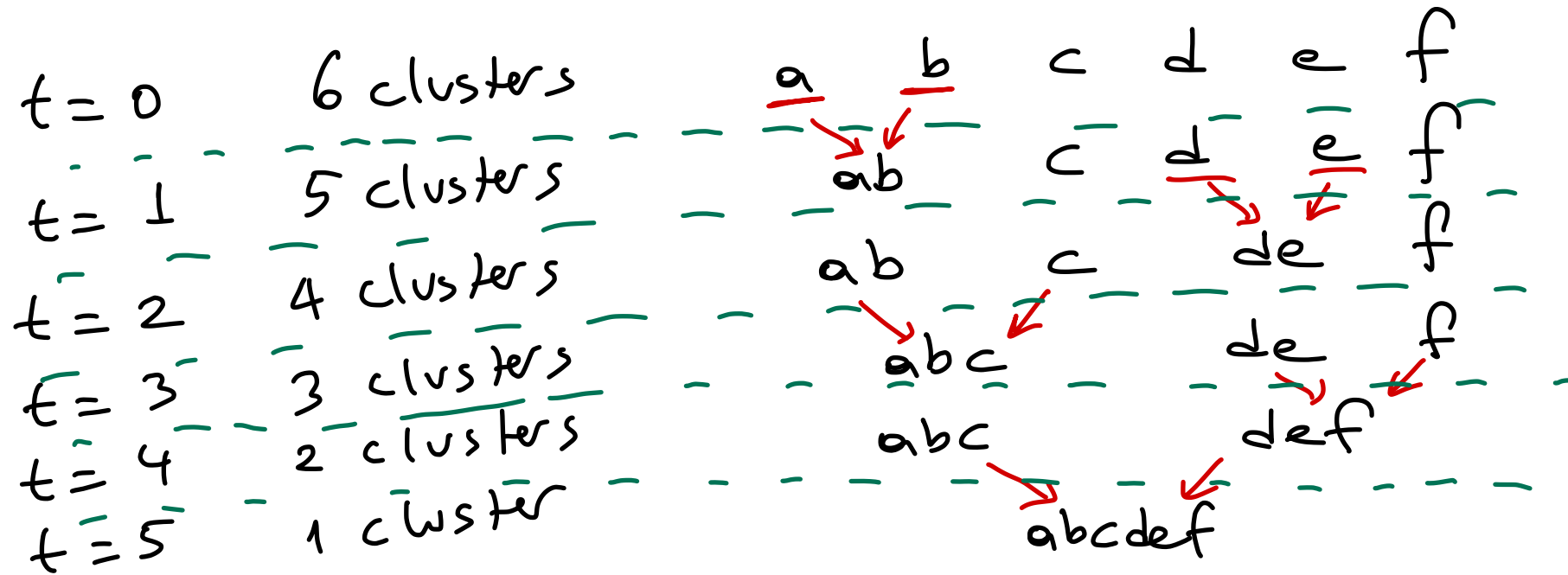
midpoint \rightarrow cardinality

Single-Link clustering: $d(G_A, G_B) = \min_{\substack{x_i \in G_A \\ x_j \in G_B}} d(x_i, x_j)$

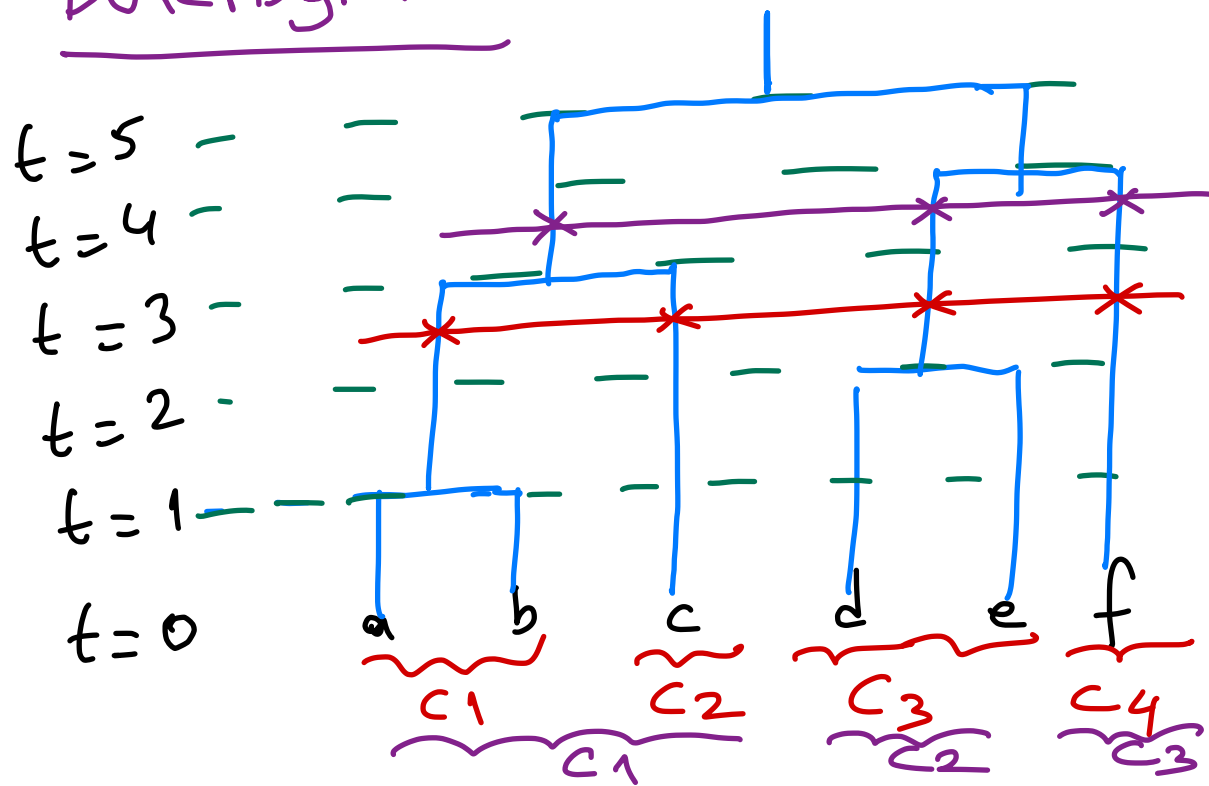
Complete-Link clustering: $d(G_A, G_B) = \max_{\substack{x_i \in G_A \\ x_j \in G_B}} d(x_i, x_j)$

Average-Link clustering: $d(G_A, G_B) = \frac{\sum_{x_i \in G_A} \sum_{x_j \in G_B} d(x_i, x_j)}{|G_A| |G_B|}$

•
•
•



Dendrogram



$K=4$ clusters

$K=3$ clusters