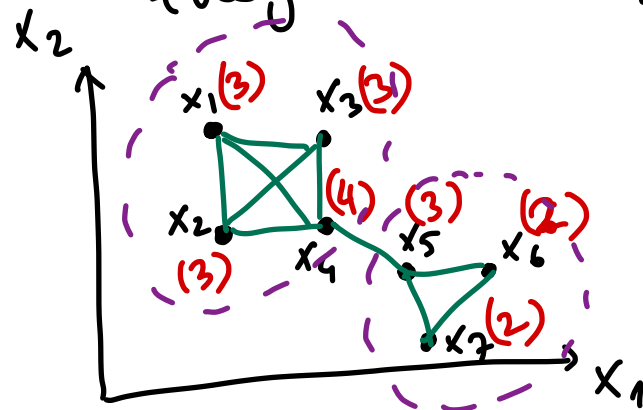# SPECTRAL CLUSTERING

- define local neighborhoods
- if the distance between $x_i$ & $x_j$ is smaller than a threshold they are neighbors.

$$b_{ij} = \begin{cases} 1 & \text{if } \|x_i - x_j\|_2 < \underline{\delta} \\ 0 & \text{otherwise} \end{cases}$$

$$b_{ij} = \begin{cases} \exp\left[-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right] & \text{if } \|x_i - x_j\|_2 < \underline{\delta} \\ 0 & \text{otherwise} \end{cases}$$

$$b_{ii} = 0 \quad \forall i$$

$$d_{ii} = \sum_{j \neq i} b_{ij} \quad \forall i$$

$\hookrightarrow$ # of neighbors of data point $i$

$$d_{ij} = 0 \quad \forall (i, j \neq i)$$

$$B = \begin{array}{c c} & \begin{array}{c c c c c c c} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 \end{array} \\ \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \end{array} & \left[\begin{array}{c c c c c c c} 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 \end{array}\right] \end{array}$$

$$D = \begin{array}{c c} & \begin{array}{c c c c c c c} x_1 & x_2 & x_3 & x_4 & x_5 & x_6 & x_7 \end{array} \\ \begin{array}{c} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \\ x_7 \end{array} & \left[\begin{array}{c c c c c c c} 3 & & & & & & \\ & 3 & & & & & \\ & & 3 & & & & \\ & & & 4 & & & \\ & & & & 3 & & \\ & & & & & 2 & \\ & & & & & & 2 \end{array}\right] \end{array}$$

mass $\leftarrow$ outgoing

$\hookrightarrow [3 \quad -1 \quad -1 \quad -1 \quad 0 \quad 0 \quad 0]$

connectivity or adjacency matrix

## Laplacian Matrix:

$$L_{N \times N} = D_{N \times N} - B_{N \times N}$$

$\longrightarrow$ each row (column) sums up to $0$.

$$L_{RANDOM-WALK} = D^{-1} \cdot L = D^{-1} \cdot (D-B) = \boxed{I - D^{-1} \cdot B}$$

$$L_{SYMMETRIC} = D^{-1/2} \cdot L \cdot D^{-1/2} = D^{-1/2} \cdot (D-B) D^{-1/2} = \boxed{I - D^{-1/2} \cdot B \cdot D^{-1/2}}$$

**SPECTRAL CLUSTERING**

**STEP #1:** Find the eigenvectors of normalized $L_{N \times N}$ matrix.

**STEP #2:** Pick R <u>smallest</u> eigenvectors.

**STEP #3:** Construct Z matrix as follows:

$$Z = \begin{bmatrix} v_1 & v_2 & \cdots & v_R \end{bmatrix}_{N \times R}$$

$\longrightarrow$ 1st smallest eigenvector

$\longrightarrow$ $R^{th}$ smallest eigenvector

**STEP #4:** Run k-means clustering algorithm on Z matrix to find K clusters.

**PARAMETERS:**

$\delta$: threshold.

R: # of eigenvectors to be included

K: # of clusters to be found.

# 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 Data Points

distance $\Rightarrow$ dissimilarity

distance $\uparrow$ similarity $\downarrow$

distance $\downarrow$ similarity $\uparrow$

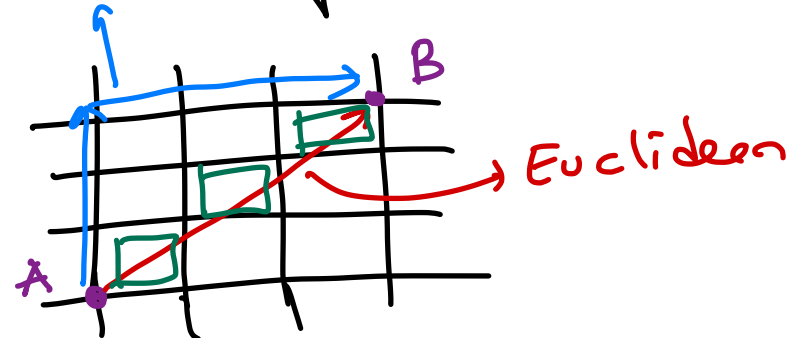$$k(x_i, x_j) = \exp\left[-\frac{\|x_i - x_j\|_2^2}{2\sigma^2}\right]$$

$0 \longleftrightarrow 1$

dissimilar $\longleftrightarrow$ similar

$0$ similar

$+\infty$ dissimilar

### Euclidean Distance

$$d(x_i, x_j) = \|x_i - x_j\|_2$$

$$= \sqrt{\sum_{d=1}^{D}(x_{id} - x_{jd})^2}$$

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

Manhattan

B

A

Euclidean

### Manhattan Distance (city-Block Distance)

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

# Component #2: The Direction to Proceed

**Agglomerative** (bottom-to-top)   **Divisive** (top-to-bottom)

⇒ Combines small clusters into bigger ones

⇒ starts with "N" clusters.

⇒ divides big clusters into smaller ones

⇒ starts with "1" cluster

# Component #3: The Distance Function Between Groups of Data Points.

$$\text{Distance}\left[\{Paris, London\} \quad , \quad \{New York\}\right]$$

$$\text{Distance}\left[\{Paris, London\} \quad , \quad \{Rome\}\right]$$



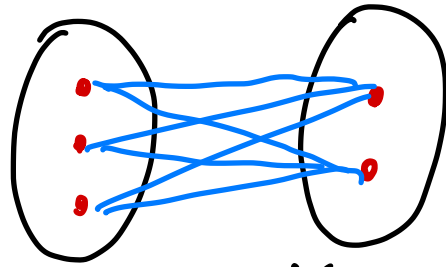$$\text{Distance}\left[\{London, Paris\}, \{Berlin, Rome\}\right]$$

**Centroid Clustering:**

$$d(C_A, C_B) = \left\| \frac{\sum\limits_{x_i \in C_A} x_i}{|C_A|} - \frac{\sum\limits_{j \in C_B} x_j}{|C_B|} \right\|_2$$

Cardinality of CA (# of members)

centroid of CA

centroid of CB

**Single-Link Clustering:**

$$d(C_A, C_B) = \min_{\substack{x_i \in C_A \\ x_j \in C_B}} d(x_i, x_j)$$
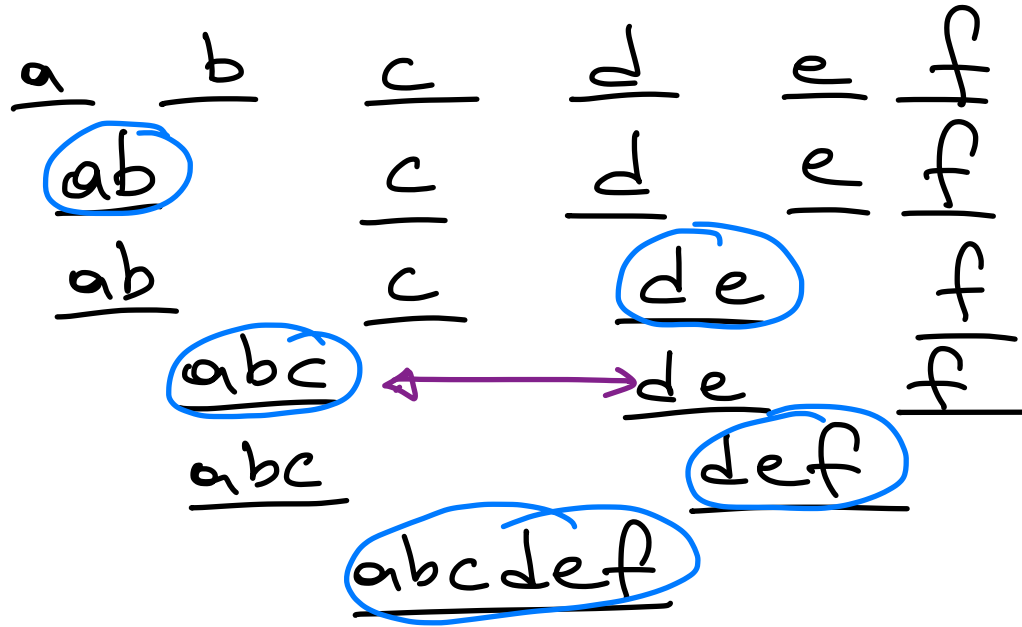
**Complete-Link Clustering:**

$$d(C_A, C_B) = \max_{\substack{x_i \in C_A \\ x_j \in C_B}} d(x_i, x_j)$$
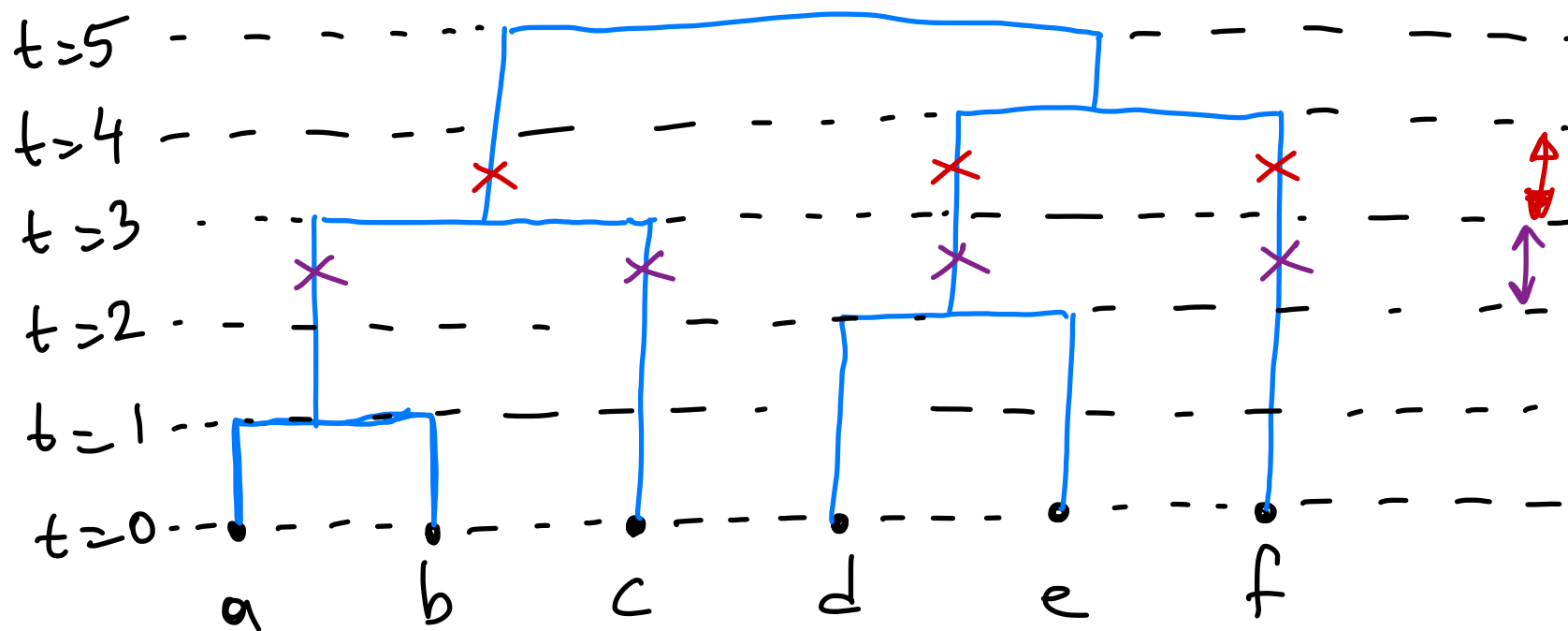
**Average-Link Clustering:**

$$d(C_A, C_B) = \frac{\sum\limits_{x_i \in C_A} \sum\limits_{x_j \in C_B} d(x_i, x_j)}{|C_A| \, |C_B|}$$

$t = 0$   6 clusters   a   b   c   d   e   f

$t = 1$   5 clusters   (ab)   c   d   e   f

$t = 2$   4 clusters   ab   c   (de)   f

$t = 3$   3 clusters   (abc) ⟷ de   f

$t = 4$   2 clusters   abc   (def)

$t = 5$   1 cluster   (abcdef)

Agglomerative

# Dendrogram



$K = 3$ clusters

$C_1 = \{a, b, c\}$

$C_2 = \{d, e\}$

$C_3 = \{f\}$

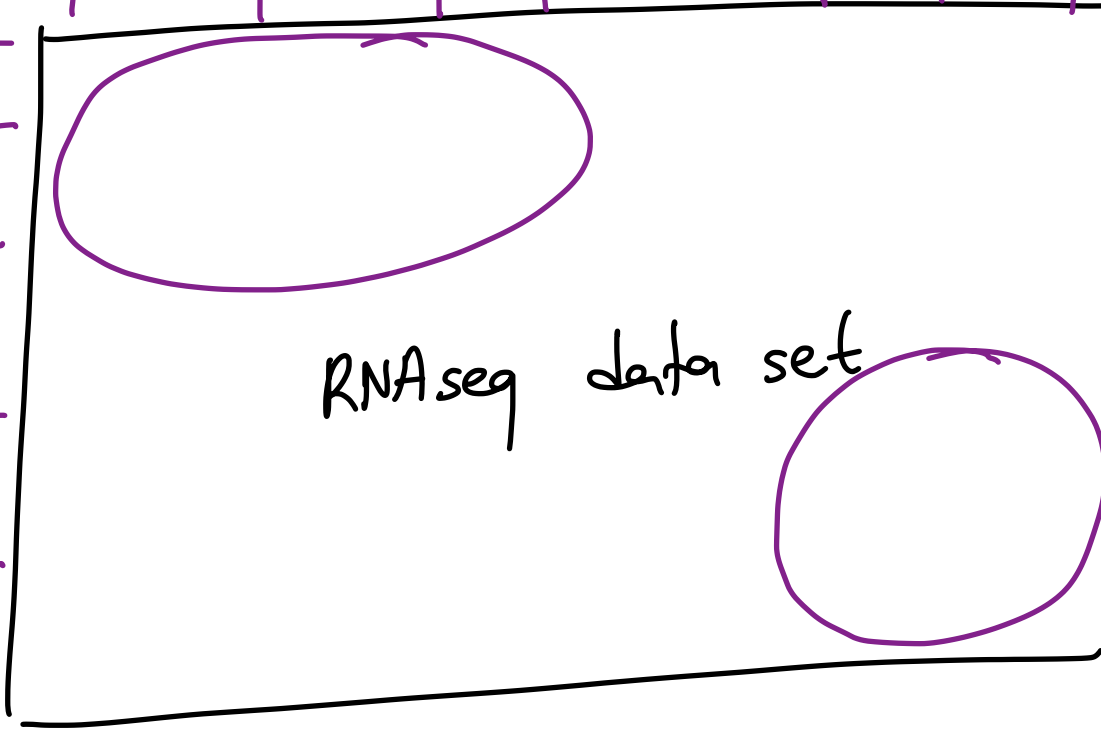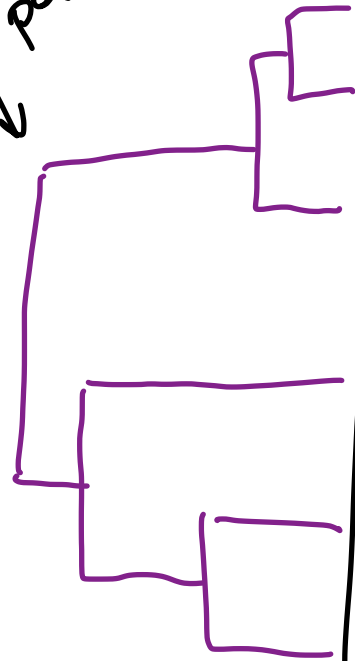$K = 4$ clusters

$C_1 = \{a, b\}$

$C_2 = \{c\}$

$C_3 = \{d, e\}$

$C_4 = \{f\}$

clustering of genes

clustering of patients

RNAseq data set

patients

genes

"bi-clustering"