# SPECTRAL CLUSTERING
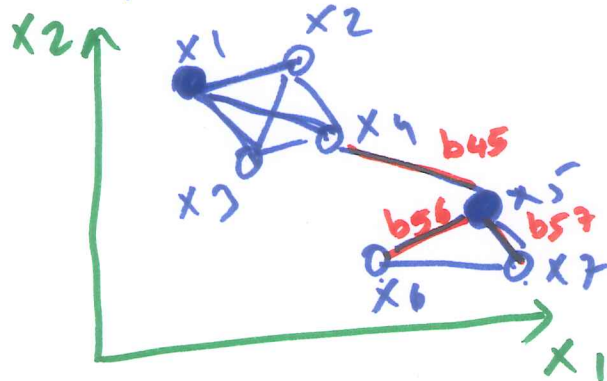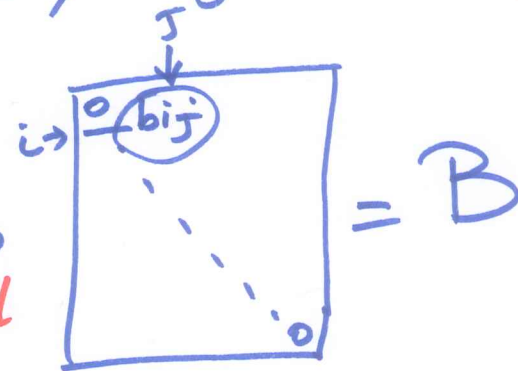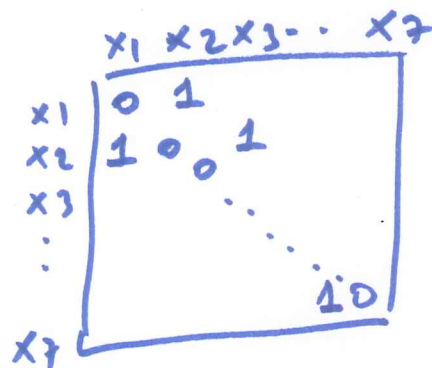
— define local neighborhoods

— distance between $x_i$ & $x_j$ is smaller then a threshold, they are neighbors.



$x_1$ $x_2$ $x_3$ $\cdots$ $x_7$

matrix with entries:
$x_1$: 0 1
$x_2$: 1 0 1
$x_3$: $\cdots$
$\vdots$
$x_7$: 1 0

$\implies$ = B (matrix with $b_{ij}$, $i \to$, $j \downarrow$)

threshold
$\downarrow$
$\|x_i - x_j\| > \delta$

$$b_{ij} = \begin{cases} 0 & \text{if the } \|x_i - x_j\| > \delta \\ \exp\left(-\dfrac{\|x_i - x_j\|^2}{2\sigma^2}\right) & \text{otherwise.} \end{cases}$$
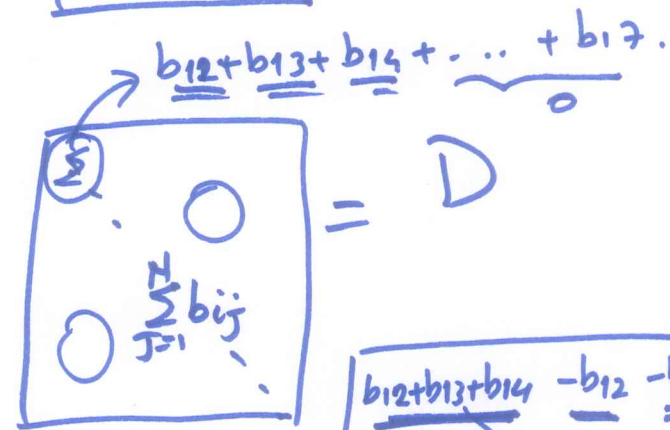
$\underbrace{\phantom{\exp\left(-\dfrac{\|x_i-x_j\|^2}{2\sigma^2}\right)}}_{k(x_i, x_j)}$

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

$\dfrac{b_{12} + b_{13} + b_{14} + \cdots + b_{17}}{0}$

= D

$\displaystyle\sum_{j=1}^{N} b_{ij}$

$b_{12} + b_{13} + b_{14} \;\; -b_{12} \;\; -b_{13} \cdots$

Laplacian matrix $\implies$ $\boxed{L = D - B}$

$\hookrightarrow$ each row sums up to 0

①

$$\mathcal{L}_{\text{RANDOM-WALK}} = \bar{D}^{-1} \cdot L = \bar{D}^{-1}(D-B) = \underline{I - \bar{D}^{-1} \cdot B}$$

$$\mathcal{L}_{\text{SYMMETRIC}} = D^{-1/2} \cdot L \, D^{-1/2} = I - \underline{D^{-1/2} \cdot B \cdot D^{-1/2}}$$

$\left.\right\}$ normalized $\mathcal{L}$ matrices.

$N \times N$ (under RANDOM-WALK)

$N \times N$ (under SYMMETRIC)

**SPECTRAL CLUSTERING**

**Step 1:** Find the eigenvectors of normalized $\mathcal{L}$ matrix

**Step 2:** Pick $R$ largest eigenvectors.

**Step 3:**
$$Z = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix} \begin{bmatrix} v_1 & v_2 & \cdots & v_R \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_N \end{bmatrix} \rightarrow N \times R$$

$v_1 \downarrow N \times 1 \qquad v_2 \downarrow N \times 1 \qquad v_R \downarrow N \times 1$

**Step 4:** Run k-means algorithm on $Z$ to find $k$ clusters.

**PARAMETERS**

$\delta$: distance threshold

$R$: # of eigenvects to be included

$K$: # of clusters.

$k(x_i, x_j) \quad d(x_i, x_j)$

②

# HIERARCHICAL CLUSTERING

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

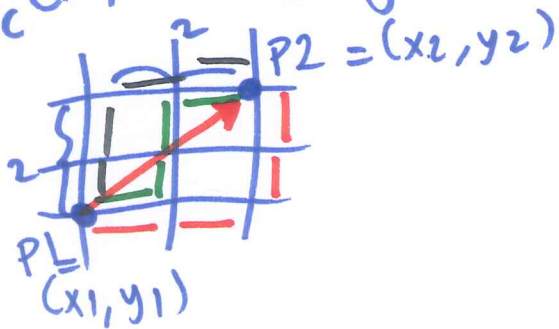First Component: Distance function between data points.

distance $\Rightarrow$ dissimilarity.

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

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

Euclidean distance $= \|x_i - x_j\|_2 = \sqrt{\sum_{d=1}^{D}(x_{id} - x_{jd})^2} = \sqrt{(x_i - x_j)^T (x_i - x_j)}$

$d_E(P_1, P_2) = \frac{2\sqrt{2}}{4}$

$d_C(P_1, P_2) = 4$

City-block distance $= \sum_{d=1}^{D}|x_{id} - x_{jd}|$

$\begin{bmatrix} 0 \\ 0 \end{bmatrix} \begin{bmatrix} \vdots \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} = 0$

$P2 = (x_2, y_2)$

PL
$(x_1, y_1)$

③

# Second Component: Direction to proceed.

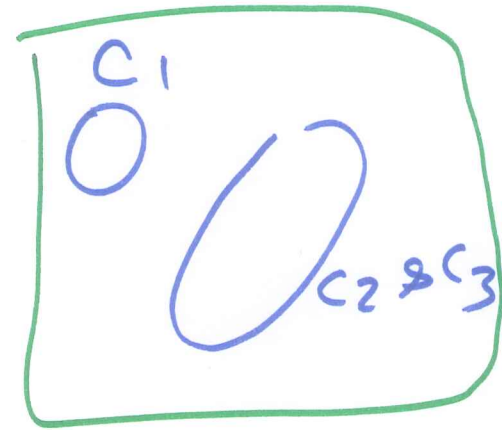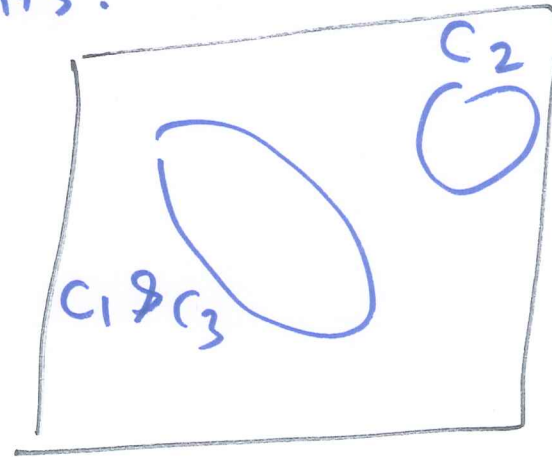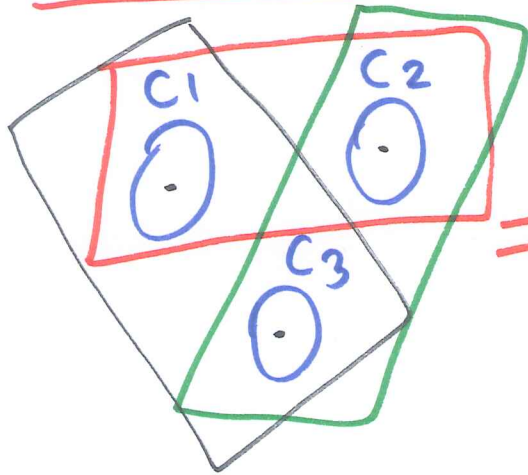Agglomerative.
→ Start with N clusters
→ Combine small clusters
into bigger ones.

Divisive
→ Start with one cluster
→ Divide big cluster into smaller ones.

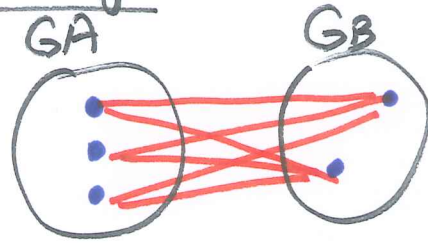# Third Component: Distance function between groups of data points.



$d(C_1, C_2)$
$d(C_1, C_3)$     } pick the smallest.
$d(C_2, C_3)$

④

Single-link clustering: $d(G_A, G_B) = \min\limits_{x_i \in G_A, x_j \in G_B} d(x_i, x_j)$

GA GB

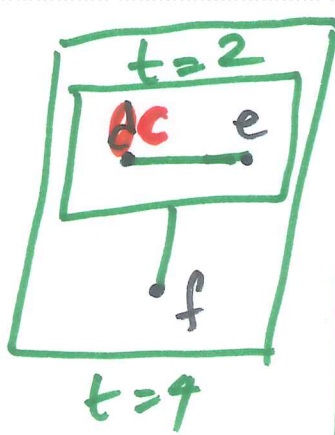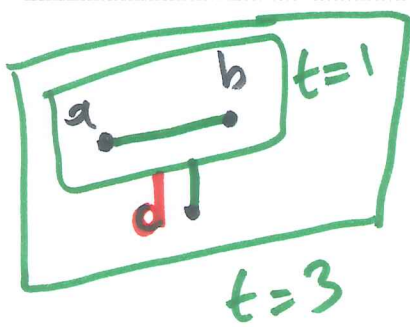Complete-link clustering: $d(G_A, G_B) = \max\limits_{x_i \in G_A, x_j \in G_B} d(x_i, x_j)$

Average-link clustering: $d(G_A, G_B) = \dfrac{\sum\limits_{x_i \in G_A} \sum\limits_{x_j \in G_B} d(x_i, x_j)}{|G_A| \cdot |G_B|}$

# of data points in GA
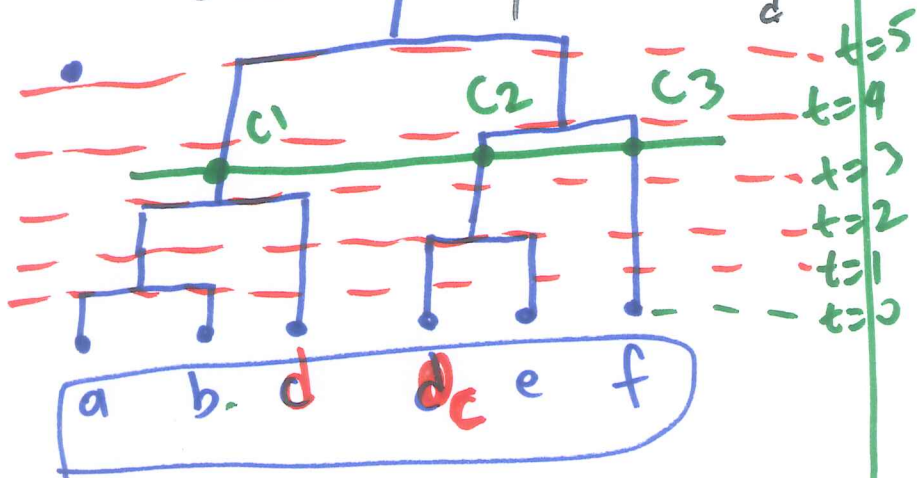
↳ # of data points in GB (cardinality)

Centroid clustering: $d(G_A, G_B) = \left\| \dfrac{\sum x_i}{|G_A|} - \dfrac{\sum x_j}{|G_B|} \right\|_2$

.
.
.

t=1, t=2, t=3, t=4 diagrams (a, b, c, d, e, f)

$t = 0$  6 clusters  $\{a\}, \{b\}, \{c\}, \{d\}, \{e\}, \{f\}$

$t = 1$  5 clusters  $\{a,b\}, \{c\}, \{d\}, \{e\}, \{f\}$

$t = 2$  4 clusters  $\{a,b\}, \{c\}, \{d,e\}, \{f\}$

$t = 3$  3 clusters  $\{a,b,c\}, \{d,e\}, \{f\}$

$t = 4$  2 clusters  $\{a,b,c\}, \{d,e,f\}$

$t = 5$  1 cluster  $\{a,b,c,d,e,f\}$

**Dendrogram**

abcd  ef | abcf  de
abce  df | abdf  ce
abde  cf | acdf  be
acde  bf | bcdf  ae
bcde  af |



a  b  c  d  e  f

$\dfrac{n!}{k!(n-k)!}$

$6 \rightarrow$ 5 vs 1
$15 \rightarrow$ 4 vs 2
$20 \rightarrow$ 3 vs 3
$\dfrac{\quad}{41}$  2 vs 4

$\binom{6}{?} = \dfrac{6!}{5!}$

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

$C_2 = \{d, e\}$

$C_3 = \{f\}$