

Clustering (Grouping  
together  
similar  
entities)  
 $(X, d(, ))$

$$X = \{x_1, x_2, \dots, x_N\}$$

$$d: X \times X \rightarrow \mathbb{R}^+$$

$$x_i \in \mathbb{R}^n$$

### • Clustering Hypothesis

Objects which are similar must be  
grouped together in same cluster

Objects which are not similar should be  
assigned to different clusters.

Take home quiz.

$$X = \{x_1, \dots, x_N\}$$

Q.  $N=4$  how  
— many distinct

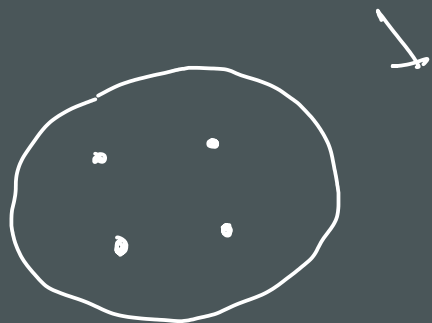
$$P = \{x_1, x_2, \dots, x_k\} \rightarrow \text{Partition of } X$$

Partitions of  
 $X$  ?

$$\forall i, j \quad x_i \subseteq X, \quad x_i \cap x_j = \emptyset, \quad \bigcup_{i=1}^k x_i = X$$

$$x_1, x_2, x_3, x_4 \rightarrow P_1 = \{\{ \}, \{x_1, x_2, x_3, x_4\}\}$$

$$P_2 = \{\{x_1\}, \{x_2, x_3, x_4\}\} \quad P_3 = \{\{x_1\}, \{x_2\}, \{x_3\}, \{x_4\}\}$$

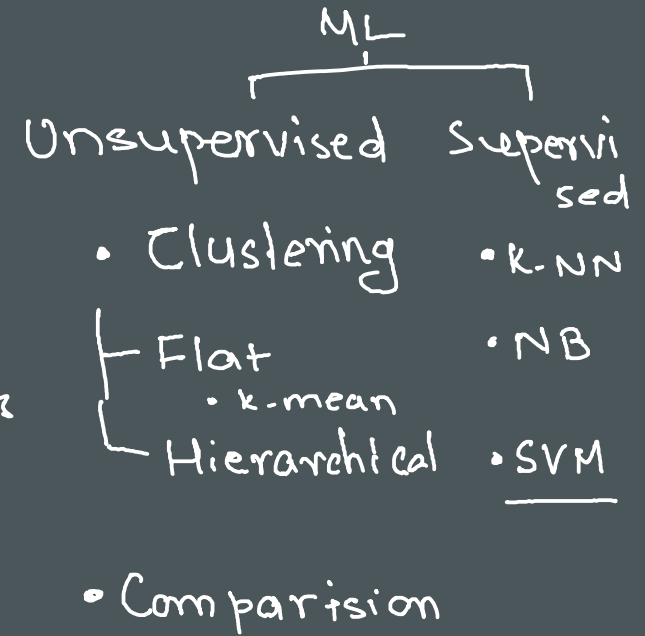
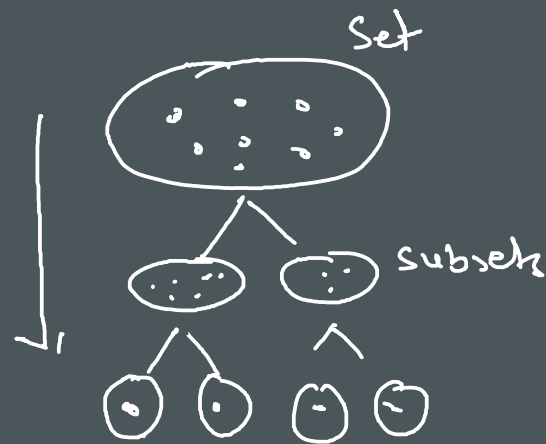
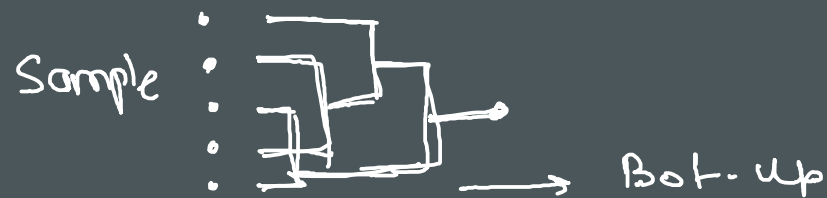


$$\cdot \{x_1, x_2, x_3, x_4\}$$

$$\{x_1\} \quad \{x_2\} \quad \{x_3\} \quad \{x_4\}$$

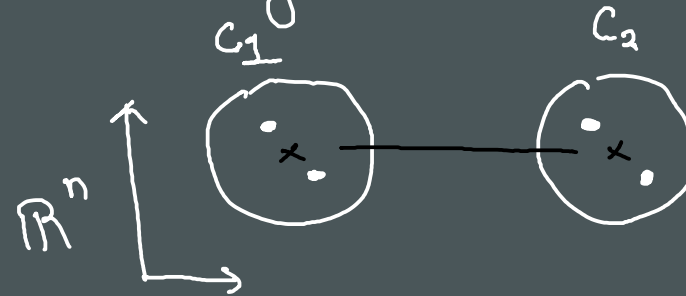
# Hierarchical Clustering

- Top-down
- Bottom-up



• Hierarchical Agglomerative Clustering (HAC) BU↑

→ Merge — Similarity



$(X, d(\cdot, \cdot)) \Rightarrow$  Metric Space

$\rightarrow$  distance metric in  $\mathbb{R}^n$

$X = \{x_1, \dots, x_n\}$   $d(x_i, x_j)$

1.  $d \geq 0$

2.  $d(x, y) = 0 \Rightarrow x = y$

3.  $d(x, z) \leq d(x, y) + d(y, z)$

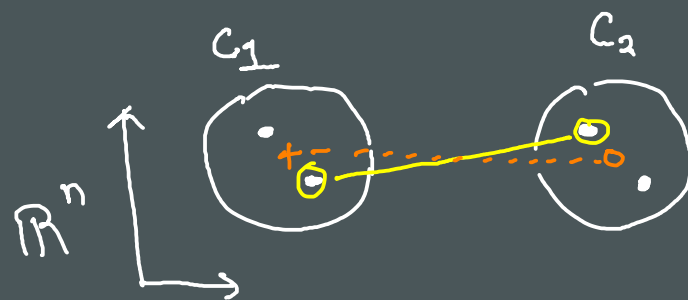
$\|t_1 - t_2\|$

triangle

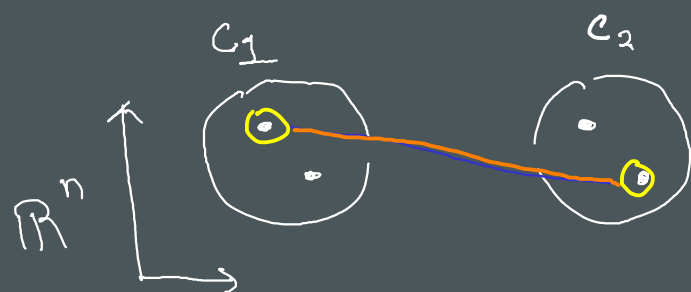
$\langle t_1 | t_2 \rangle$  inequality

(i) Centroids  $\text{sim}(\bar{c}_1, \bar{c}_2)$   $\bar{c}_1, \bar{c}_2 \in \mathbb{R}^n$

(ii) Single-link



(iii) Complete-link

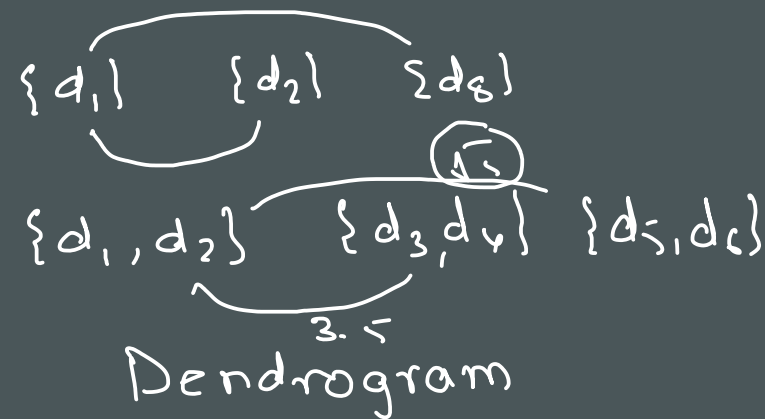
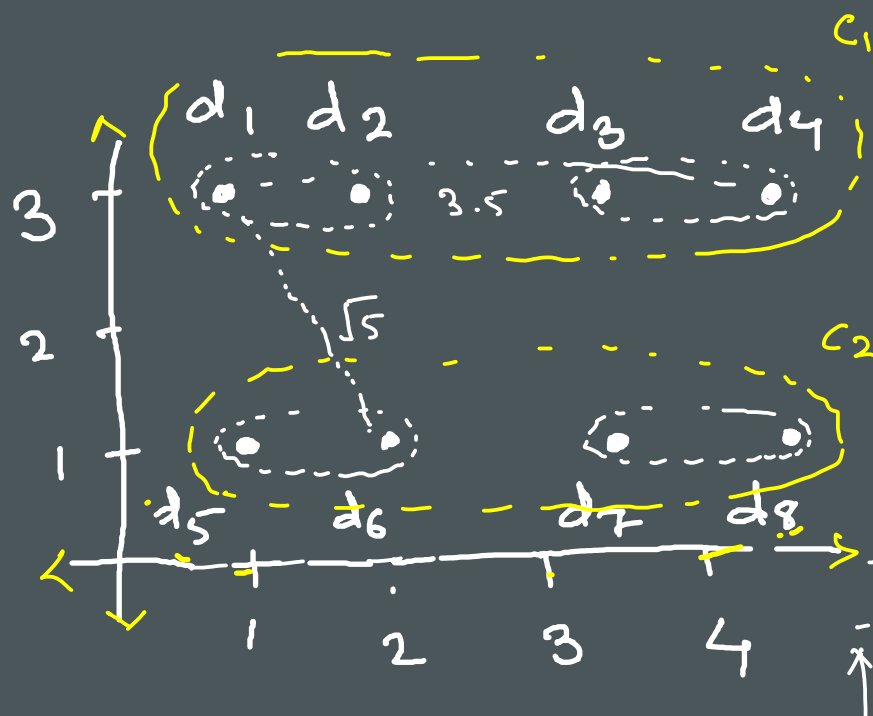


$$d(C_1, C_2) = \min_{\substack{x_1 \in C_1 \\ x_2 \in C_2}} d(x_1, x_2)$$

$$\max d(x_1, x_2)$$

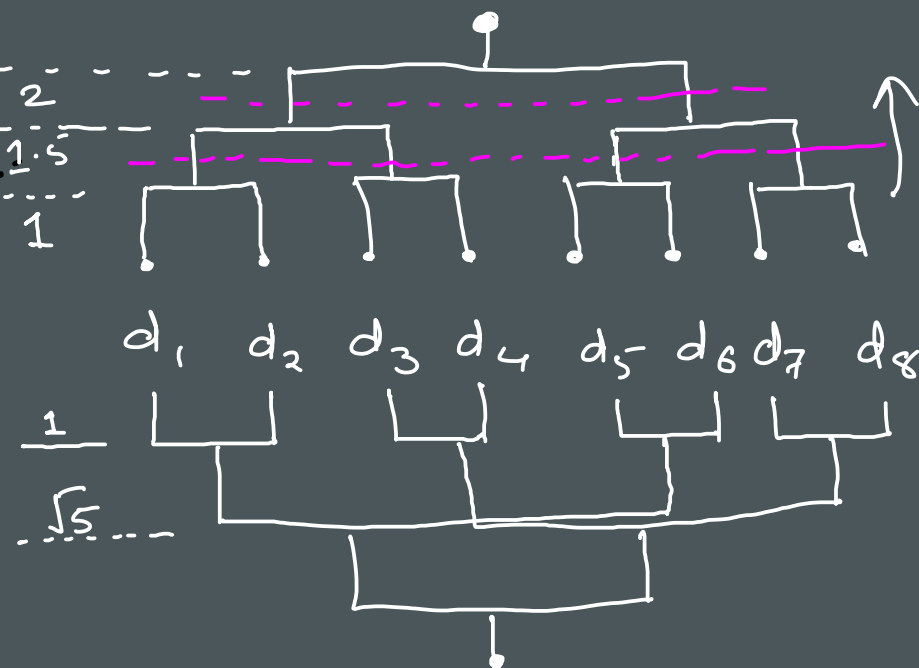
$$d(C_1, C_2) = \max_{\substack{x_1 \in C_1 \\ x_2 \in C_2}} d(x_1, x_2)$$

"dendrogram"

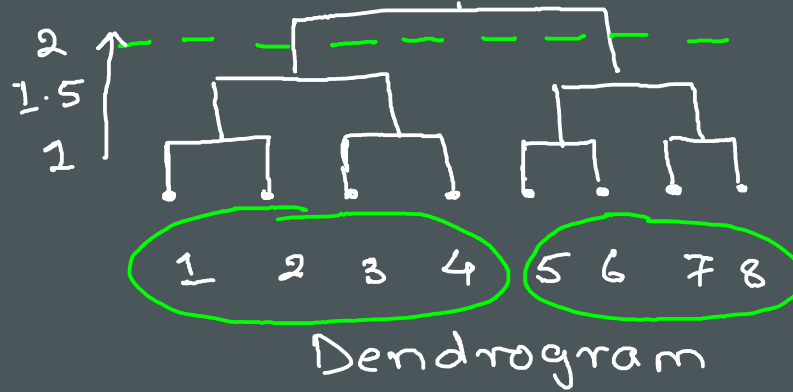
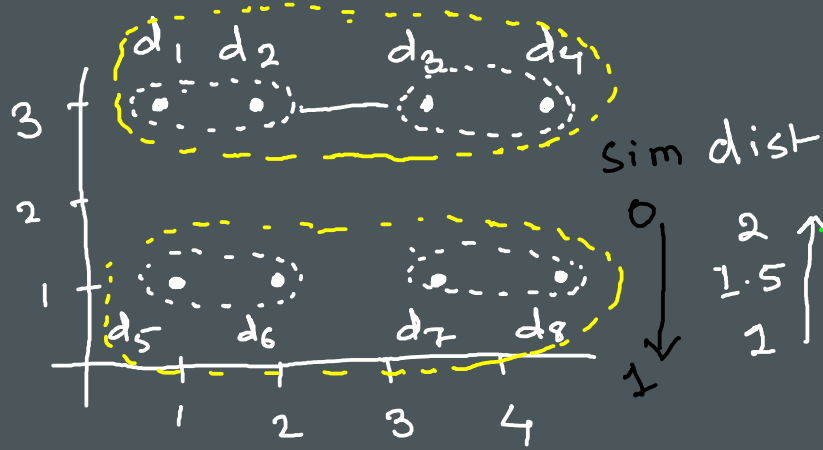


$$d(c_i, c_j) = \min_{d_i \in c_i, d_j \in c_j} d(d_i, d_j)$$

$$d(c_i, c_j) = \max_{d_i \in c_i, d_j \in c_j} d(d_i, d_j)$$

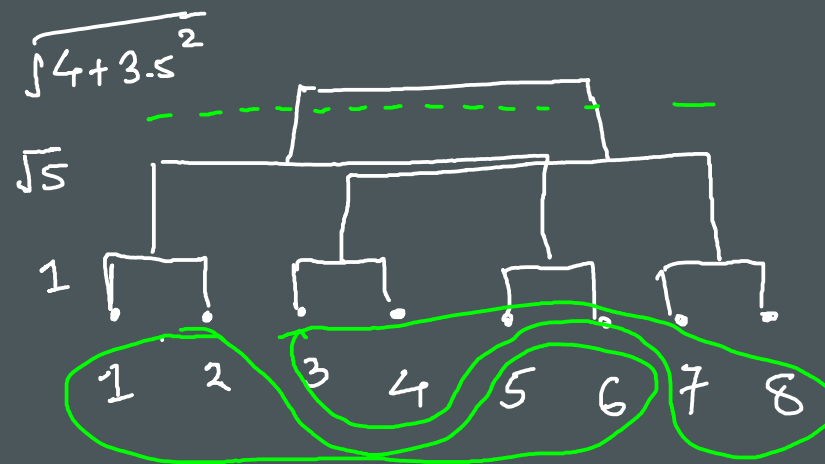
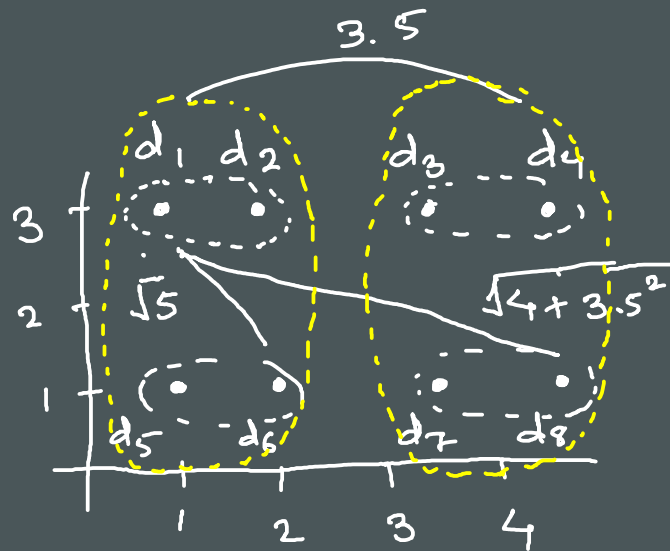


k-means



Single-Link

$$\min d(d_i, d_j) \\ d_i \in C_i \\ d_j \in C_j$$

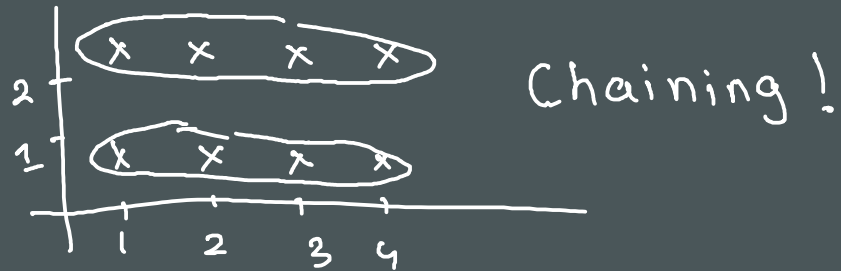


Complete-Link

$$\max d(d_i, d_j) \\ d_i \in C_i \\ d_j \in C_j$$

- Reduce the assessment of cluster quality to a single similarity bet<sup>n</sup> a pair of docs.
  - two most similar
  - two most dissimilar

Single link .



• Time Complexity ?

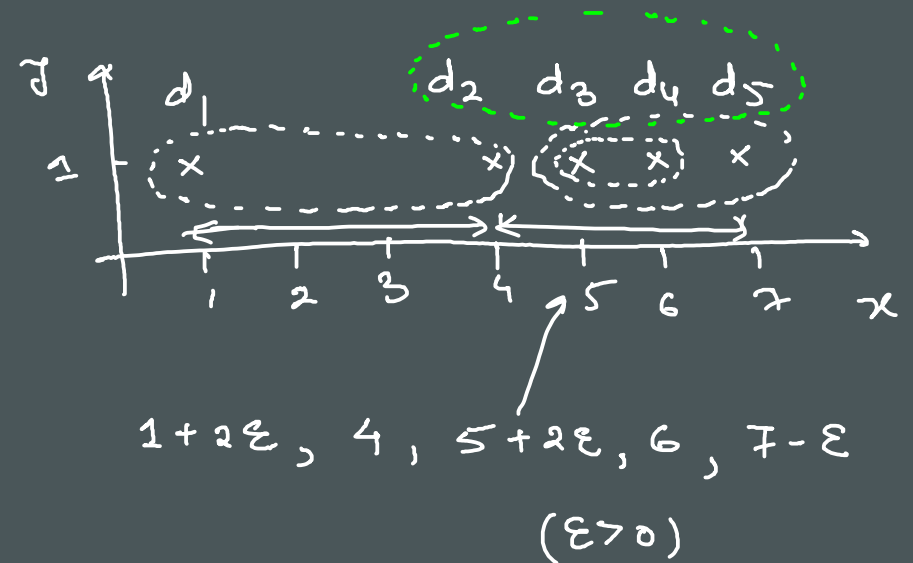
$$N - O(N^2) \sim O(N^3)$$

dist                      merge

$$\sim O(N^2 \log N)$$



• Complete link - susceptible to outliers



## Exercise

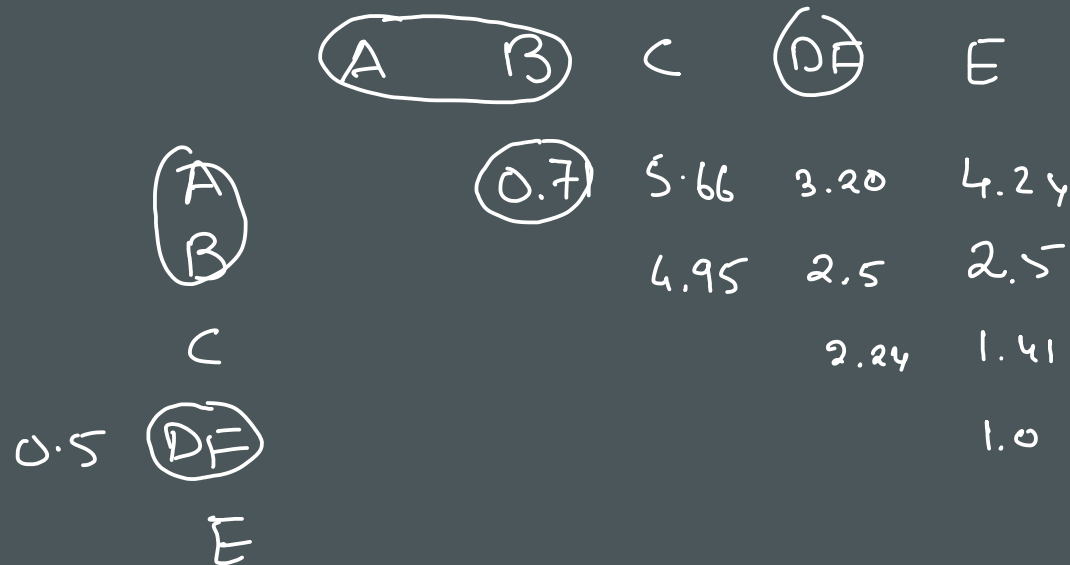
Explain Hierarchical Agglomerative Clustering in detail with three different 'merge' operations. We are given an input distance matrix of size 6 by 6 in the following Table. Entries in this matrix are calculated based on the geometry of feature vectors corresponding to the six observed document vectors. For the given distance matrix, sketch dendrograms and compare clustering results for at least two of the 'merge' strategies.

Dist	A	B	C	D	E	F
A	0.00	0.71	5.66	3.61	4.24	3.20
B		0.00	4.95	2.92	3.54	2.50
C			0.00	2.24	1.41	2.50
D				0.00	1.00	0.50
E					0.00	1.12
F						0.00

single link

— min over

all pairwise  
dist.



iterate similarly till  
you get one cluster!



Take home

Q: Why names  $\leftarrow$  single link ? (Graph Theory)  
complete link