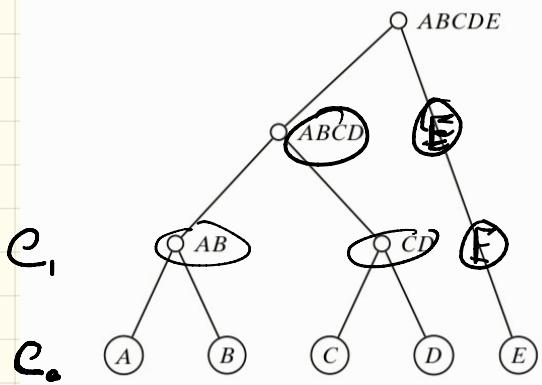


# Hierarchical Clustering

Hierarchical Clustering problem:

Given  $n$  points in  $\mathbb{R}^d$   
create a seq  
of nested partitions

Often! seq is viz as a dendrogram



course grained  
(one cluster w/ all pts) ↗  
trivial  
clustering

↓

(fine grained clustering) ↘  
one point per cluster

2 types of algos for Hieral clustering

- divisive algos (top down)
- agglomerative (bottom up)

Benefits:

- Viz tool
- find a "meaningful" level
- given  $k$  find  $k$ -clusters

Recall

= Given dataset  $D = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$   $\vec{x}_i \in \mathbb{R}^d$   
Clustering is a partition of  $D$

Given a clustering  $A = \{A_1, A_2, \dots, A_r\}$   
is nested in clustering  $B = \{B_1, B_2, \dots, B_s\}$   
iff  $\forall A_i \in A \exists B_j \in B \ni A_i \subseteq B_j$

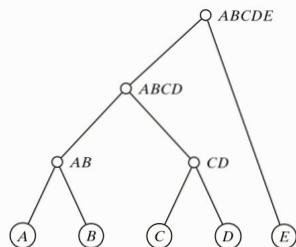
Hierarchical Clustering output

seq of clusterings  
 $(C_1 = \{\{\vec{x}_1\}, \{\vec{x}_2\}, \dots, \{\vec{x}_n\}\}, C_2, \dots, C_n = \{\{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}\})$ ,  
w/  $C_{i-1}$  is nested in  $C_i$

depends  
on algo & data

Q.2 w/ dendrogram

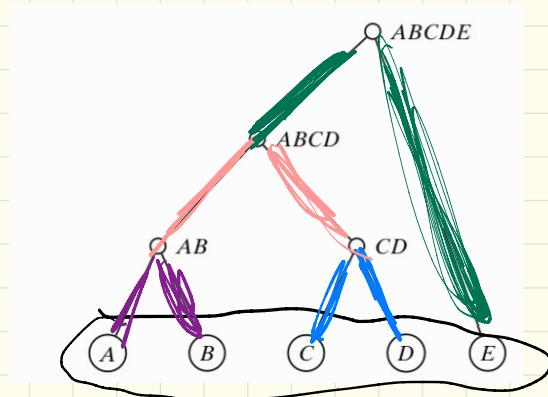
Clustering	Clusters
$C_1$	$\{A\}, \{B\}, \{C\}, \{D\}, \{E\}$
$C_2$	$\{AB\}, \{C\}, \{D\}, \{E\}$
$C_3$	$\{AB\}, \{CD\}, \{E\}$
$C_4$	$\{ABCD\}, \{E\}$
$C_5$	$\{ABCDE\}$



# Agglomerative Hierarchical Clustering

High level

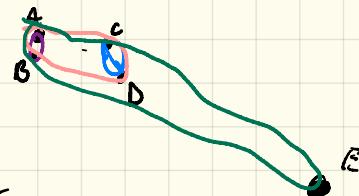
- start w/  $n$  clusters  
each w/ 1 point
- at each step  
merge 2 clusters  
(cluster centers, closest pair,  
closest-furthest pair)  
based on "closest"
- repeat w/ 1 less cluster



Sanity check question

How many iterations

- start w/  $n$  clusters
- each iter we remove 1 cluster  
⇒  $n-1$  iteration



If we want  $k$  clusters, stop at step  $n-k$

---

#### Algorithm 14.1: Agglomerative Hierarchical Clustering Algorithm

---

##### AGGLOMERATIVECLUSTERING( $\mathbf{D}, k$ ):

- 1  $\mathcal{C} \leftarrow \{C_i = \{\mathbf{x}_i\} \mid \mathbf{x}_i \in \mathbf{D}\}$  // Each point in separate cluster
  - 2  $\Delta \leftarrow \{\|\mathbf{x}_i - \mathbf{x}_j\| : \mathbf{x}_i, \mathbf{x}_j \in \mathbf{D}\}$  // Compute distance matrix
  - 3 **repeat**
  - 4     Find the closest pair of clusters  $C_i, C_j \in \mathcal{C}$
  - 5      $C_{ij} \leftarrow C_i \cup C_j$  // Merge the clusters
  - 6      $\mathcal{C} \leftarrow (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{ij}\}$  // Update the clustering
  - 7     Update distance matrix  $\Delta$  to reflect new clustering
  - 8 **until**  $|\mathcal{C}| = k$
-

Distances between clusters

$$L_2 \text{ dist of } \vec{x}, \vec{y} \in \mathbb{R}^d \quad \| \vec{x} - \vec{y} \| = \left( \sum_{i=1}^d (x_i - y_i)^2 \right)^{\frac{1}{2}}$$

Given clusters  $C_i$  &  $C_j$

Single Link - closest pair

$$\delta(C_i, C_j) = \min \left\{ \| \vec{x} - \vec{y} \| \mid \vec{x} \in C_i \text{ & } \vec{y} \in C_j \right\}$$

Mean distance - dist between centers

$\mu_i$  is center of  $C_i$

$\mu_j$  is center of  $C_j$

$$\delta(C_i, C_j) = \| \mu_i - \mu_j \|$$

Group avg - all pairs



$$n_i = |C_i|$$

$$n_j = |C_j|$$

$$\delta(C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} \| \vec{x} - \vec{y} \|}{n_i * n_j}$$

$$n_i * n_j$$

## Min Variance (Ward's method)

- increase in the sum of squared errors when two clusters merged together

Recall Cluster  $C_i$  w/ center  $\mu_i$

$$SSE(C_i) = \sum_{\vec{x} \in C_i} \|\vec{x} - \mu_i\|^2$$

given clustering  $C = \{C_1, C_2, \dots, C_m\}$   
 $\exists \mu_i$  is the center of  $C_i$

$$SSE(C) = \sum_{i=1}^m SSE(C_i) = \sum_{i=1}^m \sum_{\vec{x} \in C_i} \|\vec{x} - \mu_i\|^2$$

Ward's method

$$S(C_i, C_j) = SSE(\{C_i \cup C_j\}) - SSE(C_i) - SSE(C_j)$$

after some alg

$$\text{let } n_i = |C_i|$$

$$n_j = |C_j|$$

$\mu_i$ : center of  $C_i$

$\mu_j$ : center of  $C_j$

$$S(C_i, C_j) = \frac{n_i n_j}{n_i + n_j} \|\mu_i - \mu_j\|^2$$

weighted sq dist  
between individual cluster  
means and  
 $\frac{1}{2}$  of harmonic mean  
of cluster sizes

$$\mu_h(a, b) = \frac{2}{\frac{1}{a} + \frac{1}{b}}$$

$$= \frac{2ab}{a+b}$$

Eg w/ single link

S	B	C	D	E
A	1	3	2	4
B	3	2	3	
C		1	3	
D			5	

+ find smallest



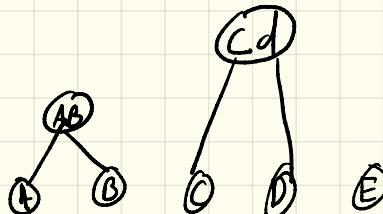
S	B	C	D	E
A	1	3	2	4
B	3	2	3	
C		1	3	
D			5	

↓ update dist matrix



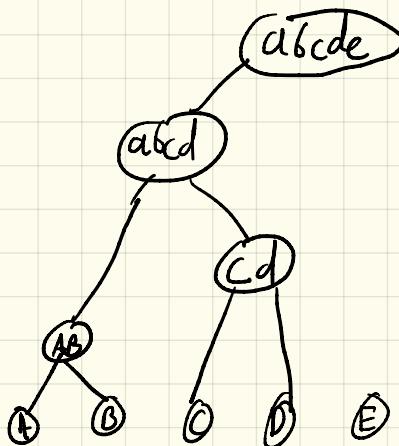
S	C	D	E
AB	3	2	3
C	1	3	
D		5	

↓ update dist matrix



S	C	D	E
AB	2	3	
CD	3		

↓ - - -



Update  $C_i \& C_j \rightarrow C_{ij}$  update dist to all other  $C_r$

Lance-Wittingham formulation

cluster  $C_r$

$$S(C_{ij}, C_r) = d_i S(C_i, C_r) + \alpha_j S(C_j, C_r)$$

$$+ \beta S(C_i, C_j)$$

$d_i, d_j, \beta, \gamma$  defined +  $\gamma |S(C_i, C_j) - S(C_j, C_r)|$   
in table

Measure	$\alpha_i$	$\alpha_j$	$\beta$	$\gamma$
Single link	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
Complete link	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$
Group average	$\frac{n_i}{n_i+n_j}$	$\frac{n_j}{n_i+n_j}$	0	0
Mean distance	$\frac{n_i}{n_i+n_j}$	$\frac{n_j}{n_i+n_j}$	$\frac{-n_i \cdot n_j}{(n_i+n_j)^2}$	0
Ward's measure	$\frac{n_i+n_r}{n_i+n_j+n_r}$	$\frac{n_j+n_r}{n_i+n_j+n_r}$	$\frac{-n_r}{n_i+n_j+n_r}$	0