

Hierarchical Clustering

- It produces a hierarchy (dendrogram) of nested clusters that can be visualized / analyzed
- It produces a set of possible solutions

Clustering	Partition
C ₁	{A}, {B}, {C}, {D}, {E}
C ₂	{A, B}, {C}, {D}, {E}
C ₃	{A, B}, {C}, {D, E}
C ₄	{A, B}, {C, D, E}
C ₅	{A, B, C, D, E}

- Strength: No need to assume any particular # of clusters

- Two approaches:
 - » Agglomerative
 - » Divisive

1) Divisive Hierarchical Clustering:

- Start with one cluster containing all data points → at each step, split a cluster contains a point (or k clusters)

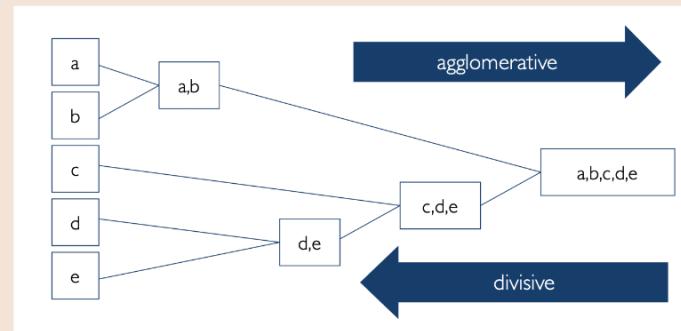
2) Agglomerative Hierarchical Clustering:

- Start with one cluster for each item → then each step merge similar clusters, until generate one (or k) clusters

Algorithm 14.1: Agglomerative Hierarchical Clustering Algorithm

AGGLOMERATIVECLUSTERING(D, k):

- 1 $\mathcal{C} \leftarrow \{C_i = \{x_i\} \mid x_i \in D\}$ // Each point in separate cluster
- 2 $\Delta \leftarrow \{\|x_i - x_j\| \mid x_i, x_j \in D\}$ // 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 Δ to reflect new clustering
- 8 until $|\mathcal{C}| = k$



- Distance between clusters;

① Start with cluster of individual points and distance matrix

② After some merging step, we have some clusters

③ We decide to merge two closest clusters C₂ and C₄ → now need to update distance matrix

④ How to update the distance matrix?



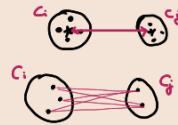
• Single Linkage: $\delta(C_i, C_j) = \min\{\|x - y\| \mid x \in C_i, y \in C_j\}$



• Complete Linkage: $\delta(C_i, C_j) = \max\{\|x - y\| \mid x \in C_i, y \in C_j\}$



• Mean / Centroid Linkage: $\delta(C_i, C_j) = \|\mu_i - \mu_j\|$



• Group Average: $\delta(C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} \|x - y\|}{n_i \cdot n_j}$

$$= \nabla SSE_{ij} = SSE_{ij} - SSE_i - SSE_j$$

• Minimum Variance / Ward's: $\delta(C_i, C_j) = \left(\frac{n_i n_j}{n_i + n_j} \right) \|\mu_i - \mu_j\|^2$

$$\delta(C_i, C_j) = \frac{\sum_{x \in C_i} \sum_{y \in C_j} (x - \mu_{ij})^2}{n_i + n_j}$$

• Lance - Williams Formula:

$$\delta(C_{ij}, C_r) = \alpha_i \cdot \delta(C_i, C_r) + \alpha_j \cdot \delta(C_j, C_r) + \beta \cdot \delta(C_i, C_j) + \gamma \cdot |\delta(C_i, C_r) - \delta(C_j, C_r)|$$

↳ whenever two clusters C_i and C_j merged into C_{ij}, update distance matrix by recomputing the distances C_{ij} to all other clusters C_r ($r \neq i, j$)

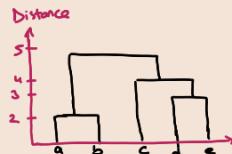
Measure	α_i	α_j	β	γ
Single link	$\frac{1}{2}$	$\frac{1}{2}$	0	$-\frac{1}{2}$
Complete link	$\frac{1}{2}$	$\frac{1}{2}$	0	$\frac{1}{2}$
Group average	$\frac{n_j}{n_i + n_j}$	$\frac{n_i}{n_i + n_j}$	0	0
Mean distance	$\frac{n_j}{n_i + n_j}$	$\frac{n_i}{n_i + n_j}$	$-\frac{n_i \cdot n_j}{(n_i + n_j)^2}$	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

-Ex: $\begin{bmatrix} a & b & c & d & e \\ 0 & 2 & 0 & 0 & 0 \\ b & 0 & 6 & 5 & 0 \\ c & 2 & 10 & 9 & 4 & 0 \\ d & 9 & 8 & 5 & 3 & 0 \\ e & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$

$\xrightarrow{\text{a and b closer}}$

a, b	c	d	e
0	5	0	0
9	4	0	0
8	5	3	0

to update distances we used single-linkage between (a,b) and others
 $d(a,2) = \min[d(1,2), d(2,3)] = d(2,3) = 5$
 $d(1,2) = \min[d(1,2), d(2,4)] = d(2,4) = 9$
 $d(1,2) = \min[d(1,2), d(2,5)] = d(2,5) = 8$



- Clustering Quality Measures

Internal Validation Measures

- Employ criteria that derived from data itself
- For instance; intracluster and intercluster distances to measure cohesion and separation

Cohesion: How similar are the points in same cluster?

$$WSS(C) = \sum_{i=1}^k \sum_{x_j \in C_i} d(x_j, \mu_i)^2$$

Centroid of C_i (n Euclidean dist)

Separation: How far apart are points in different clusters?

$$BSS(C) = \sum_{i=1}^k |C_i| \cdot d(\mu_i, \mu_j)^2$$

Centroid of whole dataset

External Validation Measures

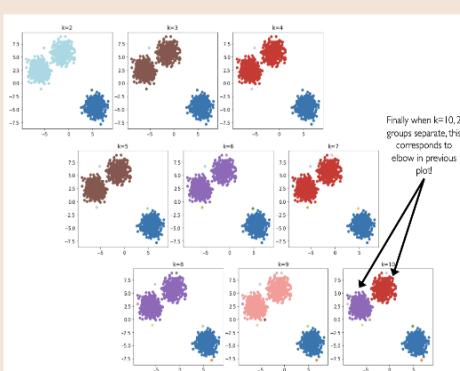
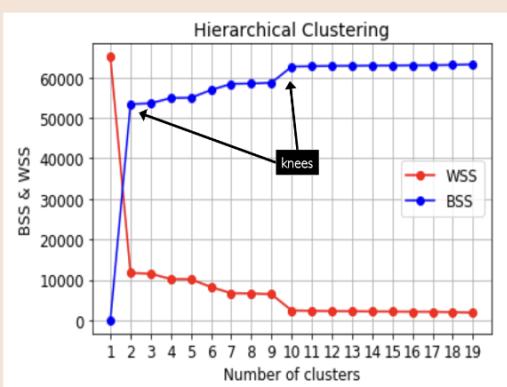
- Use prior or expert-specified info about the clusters.
- For ex.: iris dataset; evaluate the clustering using known class labels
- Employ criteria that are not inherent to the dataset (disjoint not alone)

Relative Validation Measures

- Aim to directly compare different solutions (obtained via different parameter setting for the same algorithm.)

» Evaluation of Hierarchical Clustering using Knee/Elbow Analysis

- Plot BSS and WSS for every clustering
- Look for a knee/elbow in the plot that show a significant variation of the evaluation metrics



!! We should use plots as guides

but also evaluate multiple solutions around knees/elbows because; it is limited, there is no single best criteria to select # of clusters

Apply more methods, look for an agreement

if methods disagree and do not provide clear indications none of result is good

» Clustering Evaluation Metrics:

- Purity
- Variation of information
- Pairwise Measures
- Dunn Index
- Silhouette Coefficient

generate t bootstrap samples

Apply clustering algorithm for different values of k

Computes distance between clusterings for each pair of datasets D_i and D_j for each k

Computes expected pairwise distance $\mu_d(k)$ to select k corresponding to minimum $\mu_d(k)$

```
Algorithm 17.1: Clustering Stability Algorithm for Choosing k
CLUSTERINGSTABILITY ( $A, t, k^{\max}, D$ ):
1  $n \leftarrow |D|$ 
2 // Generate t samples
3 for  $i = 1, 2, \dots, t$  do
4    $D_i \leftarrow$  sample  $n$  points from  $D$  with replacement
5 // Generate clusterings for different values of k
6 for  $i = 1, 2, \dots, t$  do
7   for  $k = 2, 3, \dots, k^{\max}$  do
8      $D_{ij} \leftarrow D_i \cap D_j$  // create common dataset using Eq. (17.52)
9     for  $k = 2, 3, \dots, k^{\max}$  do
10        $C_k(D_{ij}) \leftarrow$  cluster  $D_{ij}$  into  $k$  clusters using algorithm A
11       $d_{ij}(k) \leftarrow d(C_k(D_i), C_k(D_j))$  // distance between clusterings
12     $\mu_d(k) \leftarrow \frac{2}{t(t-1)} \sum_{i=1}^{t-1} \sum_{j>i} d_{ij}(k)$  // expected pairwise distance
13     $k^* \leftarrow \arg\min_k \{\mu_d(k)\}$ 
```

- Cluster Stability: Clustering obtained from datasets sampled from same distribution should be "stable"

It can be used to find good parameters for given clustering algorithm