

# Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki<sup>1</sup>    Wagner Meira Jr.<sup>2</sup>

<sup>1</sup>Department of Computer Science  
Rensselaer Polytechnic Institute, Troy, NY, USA

<sup>2</sup>Department of Computer Science  
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

## Chapter 14: Hierarchical Clustering

# Hierarchical Clustering

The goal of hierarchical clustering is to create a sequence of nested partitions, which can be conveniently visualized via a tree or hierarchy of clusters, also called the cluster *dendrogram*.

The clusters in the hierarchy range from the fine-grained to the coarse-grained – the lowest level of the tree (the leaves) consists of each point in its own cluster, whereas the highest level (the root) consists of all points in one cluster.

Agglomerative hierarchical clustering methods work in a bottom-up manner. Starting with each of the  $n$  points in a separate cluster, they repeatedly merge the most similar pair of clusters until all points are members of the same cluster.

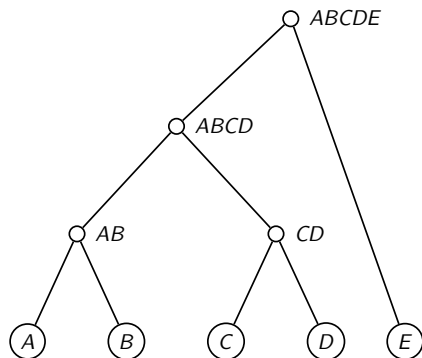
# Hierarchical Clustering: Nested Partitions

Given a dataset  $\mathbf{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ , where  $\mathbf{x}_i \in \mathbb{R}^d$ , a clustering  $\mathcal{C} = \{C_1, \dots, C_k\}$  is a partition of  $\mathbf{D}$ .

A clustering  $\mathcal{A} = \{A_1, \dots, A_r\}$  is said to be nested in another clustering  $\mathcal{B} = \{B_1, \dots, B_s\}$  if and only if  $r > s$ , and for each cluster  $A_i \in \mathcal{A}$ , there exists a cluster  $B_j \in \mathcal{B}$ , such that  $A_i \subseteq B_j$ .

Hierarchical clustering yields a sequence of  $n$  nested partitions  $\mathcal{C}_1, \dots, \mathcal{C}_n$ . The clustering  $\mathcal{C}_{t-1}$  is nested in the clustering  $\mathcal{C}_t$ . The cluster dendrogram is a rooted binary tree that captures this nesting structure, with edges between cluster  $C_i \in \mathcal{C}_{t-1}$  and cluster  $C_j \in \mathcal{C}_t$  if  $C_i$  is nested in  $C_j$ , that is, if  $C_i \subset C_j$ .

# Hierarchical Clustering Dendrogram



The dendrogram represents the following sequence of nested partitions:

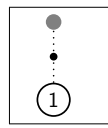
Clustering	Clusters
$C_1$	$\{A\}, \{B\}, \{C\}, \{D\}, \{E\}$
$C_2$	$\{AB\}, \{C\}, \{D\}, \{E\}$
$C_3$	$\{AB\}, \{CD\}, \{E\}$
$C_4$	$\{ABCD\}, \{E\}$
$C_5$	$\{ABCDE\}$

with  $C_{t-1} \subset C_t$  for  $t = 2, \dots, 5$ . We assume that  $A$  and  $B$  are merged before  $C$  and  $D$ .

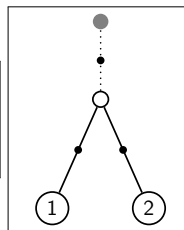
# Number of Hierarchical Clusterings

The total number of different dendrograms with  $n$  leaves is given as:

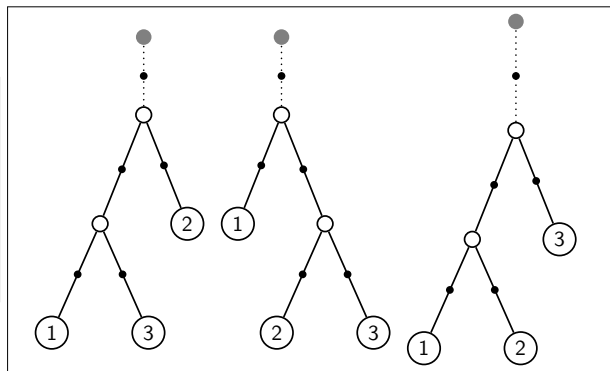
$$\prod_{m=1}^{n-1} (2m-1) = 1 \times 3 \times 5 \times 7 \times \dots \times (2n-3) = (2n-3)!!$$



(a)  $n=1$



(b)  $n=2$



(c)  $n=3$

# Agglomerative Hierarchical Clustering

In agglomerative hierarchical clustering, we begin with each of the  $n$  points in a separate cluster. We repeatedly merge the two closest clusters until all points are members of the same cluster.

Given a set of clusters  $\mathcal{C} = \{C_1, C_2, \dots, C_m\}$ , we find the *closest* pair of clusters  $C_i$  and  $C_j$  and merge them into a new cluster  $C_{ij} = C_i \cup C_j$ .

Next, we update the set of clusters by removing  $C_i$  and  $C_j$  and adding  $C_{ij}$ , as follows  $\mathcal{C} = (\mathcal{C} \setminus \{C_i, C_j\}) \cup \{C_{ij}\}$ .

We repeat the process until  $\mathcal{C}$  contains only one cluster. If specified, we can stop the merging process when there are exactly  $k$  clusters remaining.

# Agglomerative Hierarchical Clustering Algorithm

## AgglomerativeClustering( $D, k$ ):

- 1  $\mathcal{C} \leftarrow \{C_i = \{\mathbf{x}_i\} \mid \mathbf{x}_i \in \mathbf{D}\}$  // Each point in separate cluster
- 2  $\Delta \leftarrow \{\delta(\mathbf{x}_i, \mathbf{x}_j) : \mathbf{x}_i, \mathbf{x}_j \in \mathbf{D}\}$  // 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**  $|\mathcal{C}| = k$

# Distance between Clusters

## Single, Complete and Average

A typical distance between two points is the Euclidean distance or  $L_2$ -norm

$$\|\mathbf{x} - \mathbf{y}\|_2 = \left( \sum_{i=1}^d (x_i - y_i)^2 \right)^{1/2}$$

**Single Link:** The minimum distance between a point in  $C_i$  and a point in  $C_j$

$$\delta(C_i, C_j) = \min\{\|\mathbf{x} - \mathbf{y}\| \mid \mathbf{x} \in C_i, \mathbf{y} \in C_j\}$$

**Complete Link:** The maximum distance between points in the two clusters:

$$\delta(C_i, C_j) = \max\{\|\mathbf{x} - \mathbf{y}\| \mid \mathbf{x} \in C_i, \mathbf{y} \in C_j\}$$

**Group Average:** The average pairwise distance between points in  $C_i$  and  $C_j$ :

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



# Distance between Clusters: Mean and Ward's

**Mean Distance:** The distance between two clusters is defined as the distance between the means or centroids of the two clusters:

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

**Minimum Variance or Ward's Method:** The distance between two clusters is defined as the increase in the sum of squared errors (SSE) when the two clusters are merged, where the SSE for a given cluster  $C_i$  is given as

$$\delta(C_i, C_j) = \Delta SSE_{ij} = SSE_{ij} - SSE_i - SSE_j$$

where  $SSE_i = \sum_{\mathbf{x} \in C_i} \|\mathbf{x} - \mu_i\|^2$ . After simplification, we get:

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

Ward's measure is therefore a weighted version of the mean distance measure.

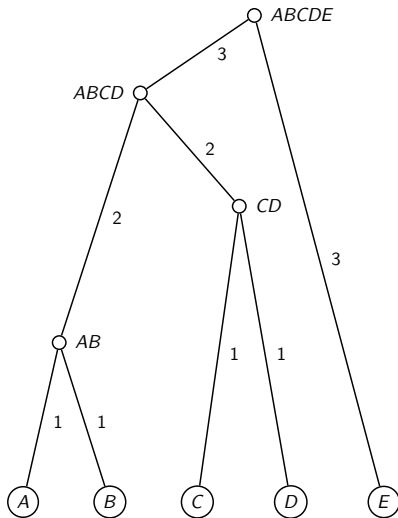
# Single Link Agglomerative Clustering

$\delta$	E
ABCD	3

$\delta$	CD	E
AB	2	3
CD		3

$\delta$	C	D	E
AB	3	2	3
C		1	3
D			5

$\delta$	B	C	D	E
A	1	3	2	4
B		3	2	3
C			1	3
D				5



# Lance–Williams Formula

Whenever two clusters  $C_i$  and  $C_j$  are merged into  $C_{ij}$ , we need to update the distance matrix by recomputing the distances from the newly created cluster  $C_{ij}$  to all other clusters  $C_r$  ( $r \neq i$  and  $r \neq j$ ).

The Lance–Williams formula provides a general equation to recompute the distances for all of the cluster proximity measures

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

The coefficients  $\alpha_i, \alpha_j, \beta$ , and  $\gamma$  differ from one measure to another.

# Lance–Williams Formulas for Cluster Proximity

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

**Single link:** Arithmetical trick to find the minimum.

$$\delta(C_{ij}, C_r) = \frac{\delta(C_i, C_r)}{2} + \frac{\delta(C_j, C_r)}{2} - \frac{|\delta(C_i, C_r) - \delta(C_j, C_r)|}{2}$$

**Complete link:** Arithmetical trick to find the maximum.

$$\delta(C_{ij}, C_r) = \frac{\delta(C_i, C_r)}{2} + \frac{\delta(C_j, C_r)}{2} + \frac{|\delta(C_i, C_r) - \delta(C_j, C_r)|}{2}$$

**Group average:** Weight the distance by the cluster size.

$$\delta(C_{ij}, C_r) = \frac{n_i}{n_i + n_j} \cdot \delta(C_i, C_r) + \frac{n_j}{n_i + n_j} \cdot \delta(C_j, C_r)$$

# Lance–Williams Formulas for Cluster Proximity

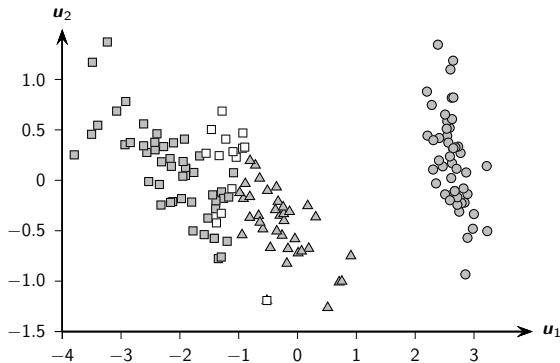
**Mean distance:** The new centroid is in the line defined by  $\mu_i$  and  $\mu_j$ , and its distance to  $\mu_r$  has to be adjusted by  $\frac{n_i \cdot n_j}{(n_i + n_j)^2}$ .

$$\delta(C_{ij}, C_r) = \frac{n_i}{n_i + n_j} \cdot \delta(C_i, C_r) + \frac{n_j}{n_i + n_j} \cdot \delta(C_j, C_r) + \frac{-n_i \cdot n_j}{(n_i + n_j)^2} \cdot \delta(C_i, C_j)$$

**Ward's measure:** The  $\Delta SSE$  of the new cluster is a weighted sum of the  $\Delta SSEs$  of the original clusters, adjusted by the fact that  $n_r$  was considered twice.

$$\delta(C_{ij}, C_r) = \frac{n_i + n_r}{n_i + n_j + n_r} \cdot \delta(C_i, C_r) + \frac{n_j + n_r}{n_i + n_j + n_r} \cdot \delta(C_j, C_r) + \frac{-n_r}{n_i + n_j + n_r} \cdot \delta(C_i, C_j)$$

# Iris Dataset: Complete Link Clustering



Contingency Table:

	iris-setosa	iris-virginica	iris-versicolor
$C_1$ (circle)	50	0	0
$C_2$ (triangle)	0	1	36
$C_3$ (square)	0	49	14

# Data Mining and Machine Learning: Fundamental Concepts and Algorithms

dataminingbook.info

Mohammed J. Zaki<sup>1</sup>    Wagner Meira Jr.<sup>2</sup>

<sup>1</sup>Department of Computer Science  
Rensselaer Polytechnic Institute, Troy, NY, USA

<sup>2</sup>Department of Computer Science  
Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

## Chapter 14: Hierarchical Clustering