

# Hierarchical Clustering

Kuldeep Dhole

Supervised By: Dr. Pollett

# What is clustering?

The process of grouping a set of objects into classes of similar objects

# What is Hierarchical clustering?

It is a method of cluster analysis which seeks to build a hierarchy of clusters.

## Strategies:

- Agglomerative (Bottom-up)
- Divisive (Top - down)

Agglomerative approach is preferred.

# Agglomerative (Bottom Up) Approach

-Given  $n$  points  $p_1, p_2, \dots, p_n$ ; We assume all as different clusters  $c_1, c_2, c_3, \dots, c_n$

-Algorithm is pretty intuitive and simple:

num = #n

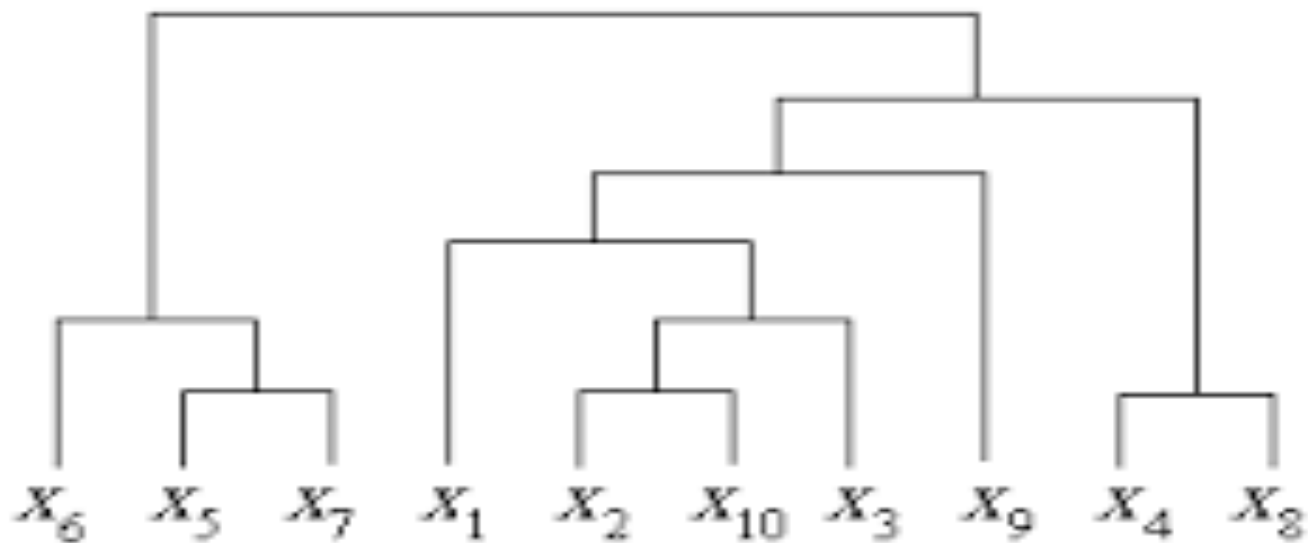
while (num > 1):

    find the pair of **closest distance** points  $p_i, p_j$  i.e.  $c_i, c_j$

    form a new cluster  $c_{i,j} = c_i + c_j$

    remove  $c_i$  and  $c_j$

# Dendrogram



# Closest Distance

-Let's assume we are in 2D space, and we have points  $p_1$ ,  $p_2$

-Approaches to get closest distances:

-Euclidean Distance:

$$[(x_1-x_2)^2 + (y_1-y_2)^2]^{0.5}$$

In multidimensional space,

$$[(a_1-a_2)^2 + (b_1-b_2)^2 + (c_1-c_2)^2 + \dots]^{0.5}$$

-Manhattan Distance:

$$|x_1-x_2| + |y_1-y_2|$$

In multidimensional space,

$$|a_1-a_2| + |b_1-b_2| + |c_1-c_2| + \dots$$

# Distance Calculations continued...

How to calculate distance between:

- a point  $p$  and cluster  $c_2$  OR
- a cluster  $c_1$  and cluster  $c_2$

Approaches:

- Single-link: We pick the point of closest distance from  $c_2$
- Complete-link: We pick the point of furthest distance from  $c_2$
- Centroid: We pick the centre of gravity from  $c_2$

# Computational Complexity

- Initially, we need to compute distances of all pairs of  $n$  individual points  $\Rightarrow O(m n^2)$
- In each loop of the algorithm, we compute the closest distance between most recently created cluster and other clusters.
- Using heap, we can achieve this by  $O(m n^2 \log(n^2))$



# Applications

- Clustering search results for efficient navigation e.g. “jaguar” -> c1 = car, c2 = animal, c3 = apple inc
- Citation ranking e.g. google scholar