

- Recap
- Hierarchical Clustering
- Proximity Matrix
- Implementation Scipy / sklearn
- Advantage Disadvantage

## Hierarchical Clustering:

- ① Agglomerative Clustering
- ② Divisive Clustering

### Agglomerative Clustering

➤ Bottom up approach



$n \rightarrow \infty$

①  $n \rightarrow n\_clusters_0$

②  $n\_clusters_1$

∧

$n\_clusters_0$

[ we merge the  
clusters which  
are most  
similar ]

③ Repeat till  
only 1 cluster  
is left

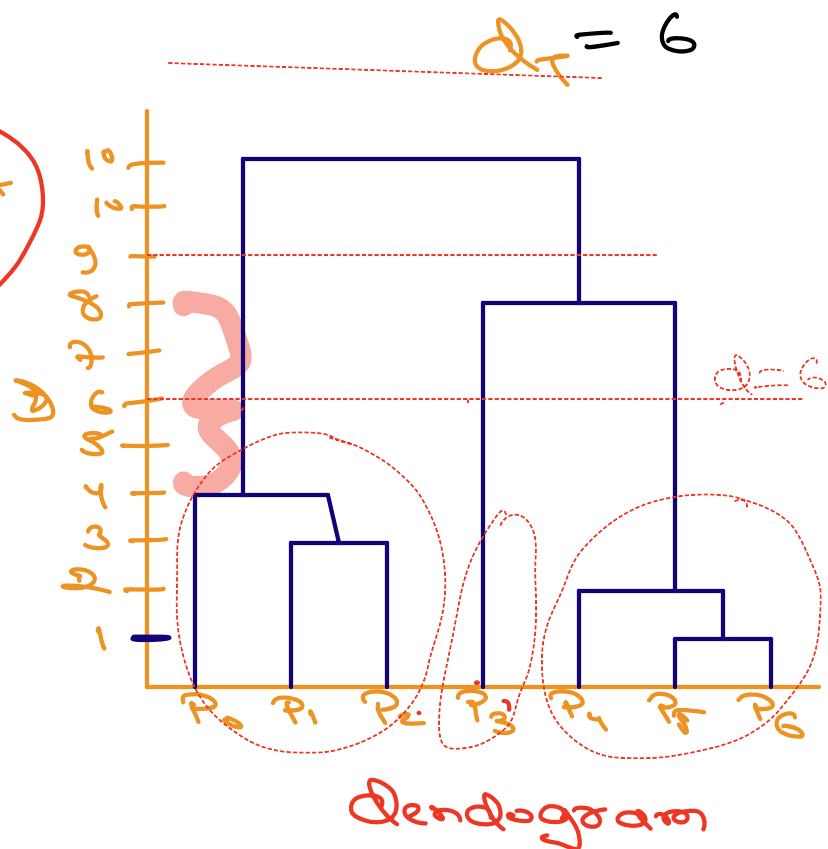
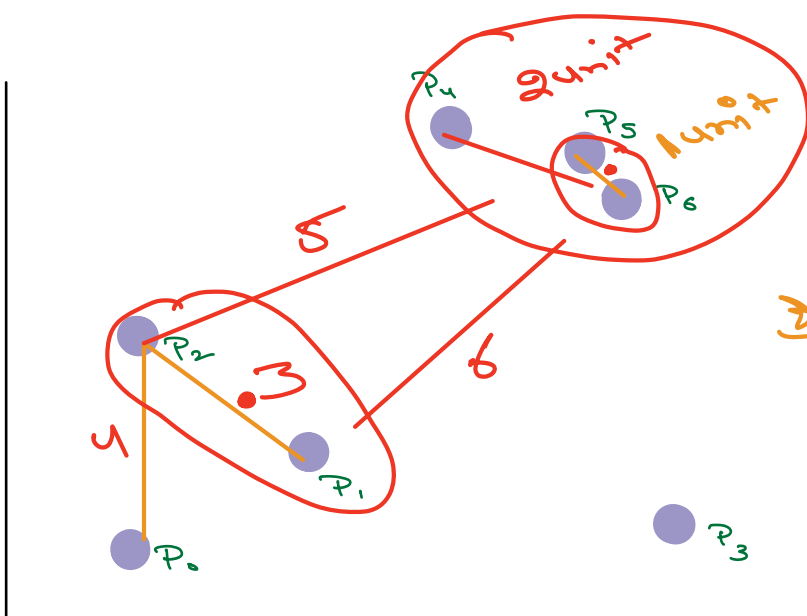
Step 1: Assume Every single point is a cluster

Step 2: Calculate proximity Matrix  $P(n \times n)$

Step 3: Using proximity Matrix

a) Merge the closest points

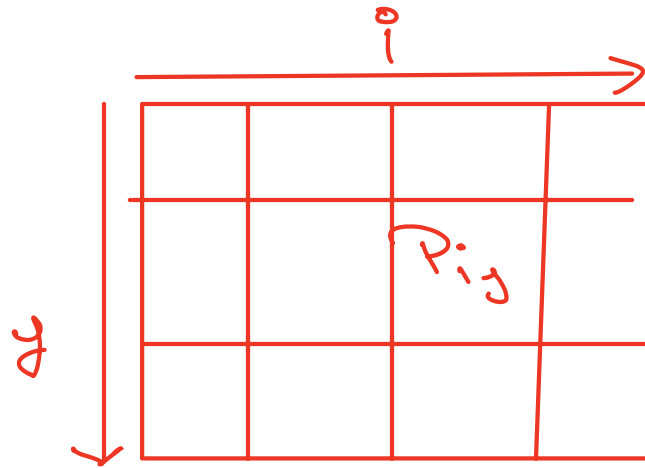
b) Update proximity matrix



# Proximity Matrix

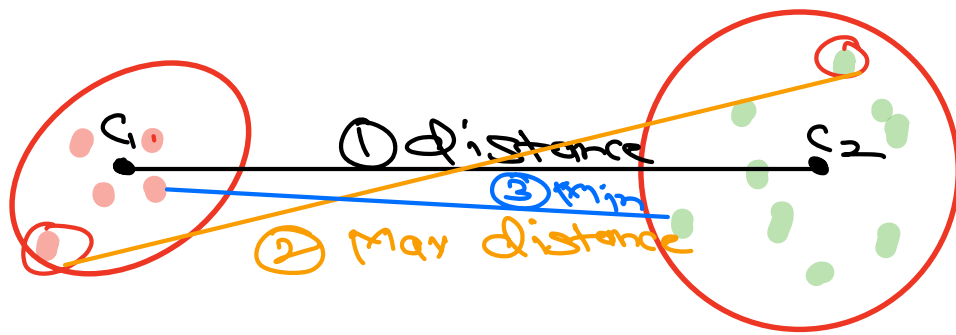
→ A matrix containing similarity values

→  $P_{n \times n}$  →  $P_{ij} \rightarrow \text{sim}(C_i, C_j)$



## Methods to Calculate Similarity for Proximity Matrix

- ① Distance between centroids of Cluster (Euclidean)
- ② Maximum Distance between 2 Cluster
- ③ Minimum Distance b/w 2 Cluster
- ④ Average Distance  $\rightarrow \frac{\sum_{x_i \in C_i} \sum_{x_j \in C_j} \text{dist}(x_i, x_j)}{|C_i| |C_j|}$
- ⑤ Ward's Distance  $\rightarrow \frac{\sum_{x_i \in C_i} \sum_{x_j \in C_j} \text{dist}(x_i, x_j)^2}{|C_i| |C_j|}$



Iteration ①

	1	2	3	4	5	6
1						
2						
3						0.11
4						
5						
6						



Iteration ②

	1	2	3	4	5
1			0		
2					
3					
4					
5					

Continue iteration till  $P_{n \times n} \Rightarrow P_{1 \times 1}$

# Limitation of Agglomerative

① Complexity: large Dataset

Space Complexity  $\Rightarrow O(n^2)$

Time Complexity

$\Rightarrow$  Proximity matrix  
Calculation

W.C  $\Rightarrow O(n^3)$   $(n \times n) \times n \times d$

k-means  $\Rightarrow$  S.C  $\Rightarrow O(k)$

$\downarrow$   
n-clusters  
T.C  $\Rightarrow O(n \times k \times d)$   
 $\downarrow$   $\downarrow$   $\downarrow$   
n-clusters  $\downarrow$   $\downarrow$   $\downarrow$   
dim

kmeans++  $\Rightarrow O(n \times k \times d \times i)$   
 $\downarrow$   
init-cent