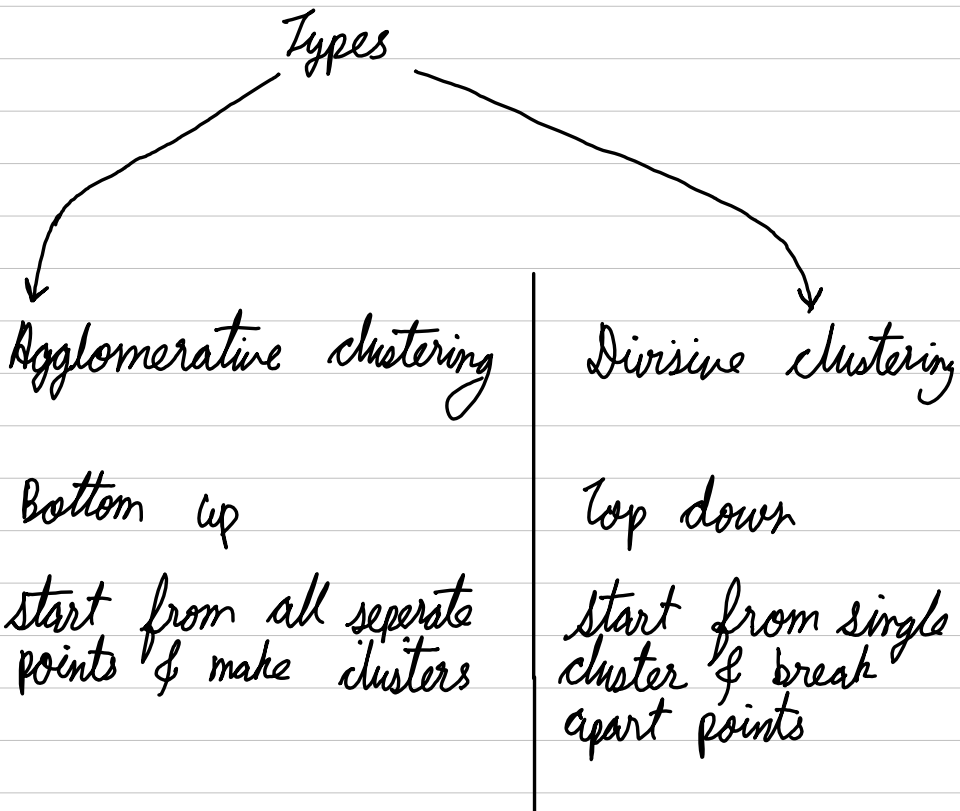


Heirchical Clustering

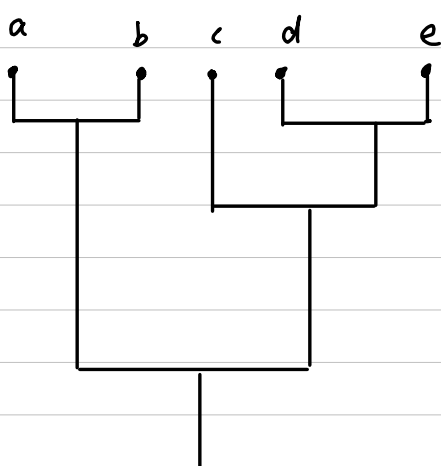
Unsupervised clustering algorithm

It works via grouping data into a tree of clusters

- steps →
1. Identify clusters close together
 2. Merge the clusters



Hierarchical Clustering is based on dendograms



Points that are close together form clusters together

The graph is called dendogram

At first level, the group formed is (a, b) c (d, e)

The similarity is based on distance criterion

For that, all points are compared

In next step, c is checked.
It is more close with cluster (d, e)
so it is merged.

At level 2

(a, b) , (c, d, e)

At level 3

(a, b, c, d, e)

At last level, all points are in single cluster

The process is hierarchical hence called hierarchical clustering

Q1) solve using hierarchical clustering

	X	Y
P1	0.4	0.53
P2	0.22	0.38
P3	0.35	0.32
P4	0.26	0.19
P5	0.08	0.41
P6	0.45	0.30

→

	P1	P2	P3	P4	P5	P6
P1	0					
P2	0.234	0				
P3	0.22	0.14	0			
P4	0.37	0.19	0.13	0		
P5	0.34	0.143	0.27	0.284	0	
P6	0.23	0.24	0.1	0.219	0.38	0

P3, P6 have smallest value

∴ (P3, P6) will be first cluster

so we combine P3 P6

	P1	P2	(P3 P6)	P4	P5
P1	0				
P2	0.234	0			
(P3 P6)	0.22	0.14	0		
P4	0.37	0.19	0.13	0	
P5	0.34	0.143	0.28	0.284	0

$$d((A, B), C) = \min \{d(A, C), d(B, C)\}$$

Less value is considered. Average value is not taken

$$\text{eg } \min((P_3 P_6), P_4) = \min(0.13, 0.22) = 0.13 \text{ etc...}$$

$$\min((P_3 P_6), P_5) = \min(0.28, 0.39) = 0.28$$

For old values, don't update

Minimum value of the table is 0.13

Combine P4 with P3 P6 (P4, (P3, P6))

	P1	P2	(P4 (P3, P6))	P5
P1	0			
P2	0.23	0		
(P4, (P3, P6))	0.22	0.14	0	
P5	0.34	0.14	0.23	0

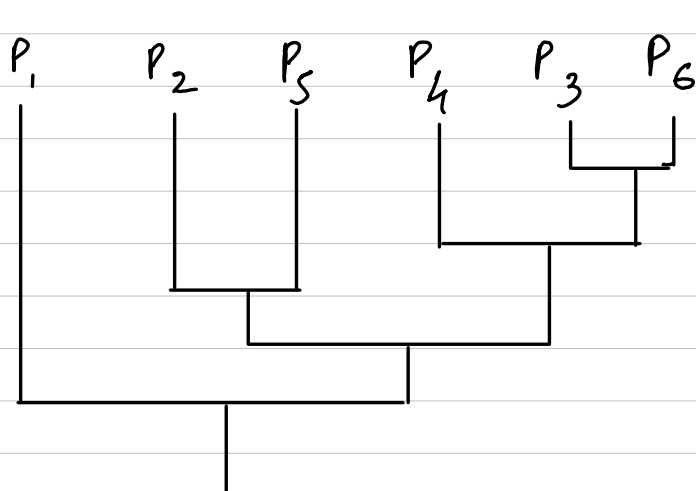
There are two Minimum values. let's consider P2 P5 combine

	P1	(P2 P5)	(P4 (P3, P6))
P1	0		
(P2 P5)	0.23	0	
(P4 (P3, P6))	0.22	0.14	0

Now combine clusters (P2 P5) and (P4 (P3, P6))

	P1	((P2 P5) (P4 (P3, P6)))
P1	0	
((P2 P5) (P4 (P3, P6)))	0.22	0

lastly combine P1 to make the last cluster



steps for Hierarchical clustering

step 1: Make adjacency table

step 2: Combine points with smallest value

step 3: Make new adjacency table where distance of point from cluster is minimum of all distances of point from the points in the cluster

step 4: Repeat steps 2 & 3 N times

step 5: Make Dendrogram

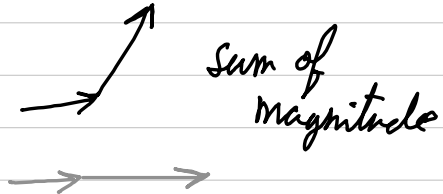
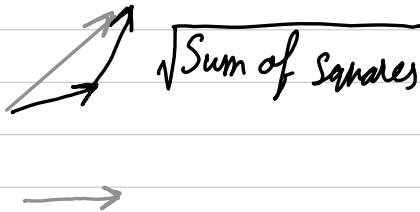
Types of distance

Euclidian

$$\sqrt{\sum x^2}$$

Manhattan

$$\sum ||x||$$



In last example we have used the euclidian distance

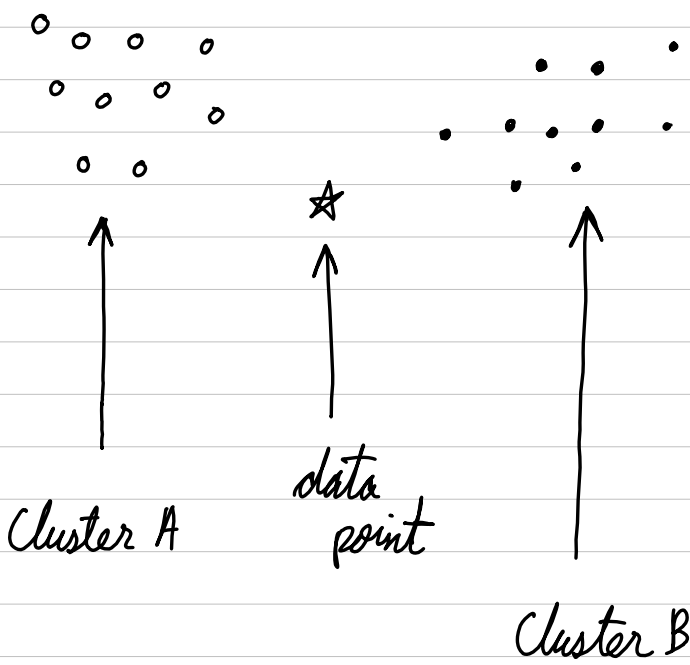
Both distances will give different clusters

Closeness Measures

We need to define what "closeness" means

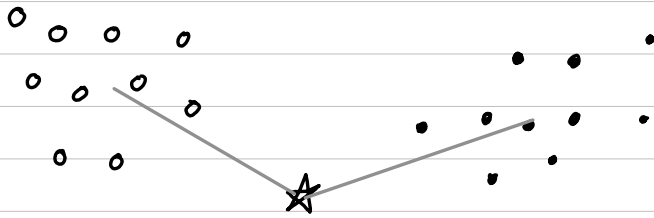
In last example, we used closeness as distance.

Other Methods also can be used



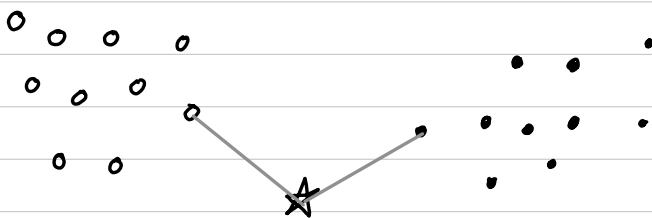
① Centroid Method

Distance of data from a cluster is distance of data from the average of all data points (ie centroid)



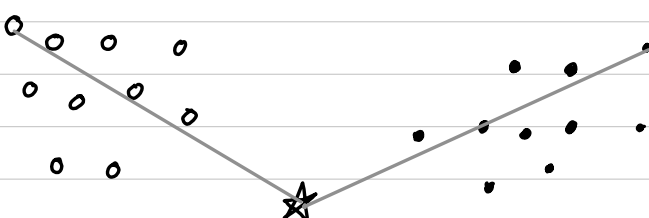
② Single linkage

Closest point to cluster

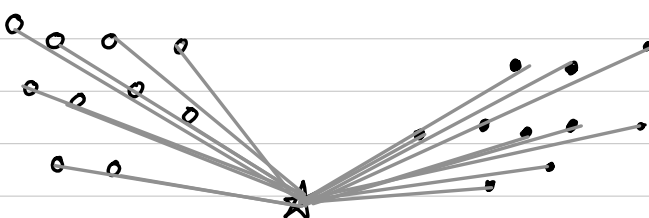


③ Complete linkage

Furthest point from cluster



④ Average linkage



Average of all the points

Advantages

Handles non convex clusters

Handles clusters of different sizes & densities

Handles Missing data & noise

Reveals the hierarchical structure that can be used to understand relationships between data

Deterministic results \rightarrow No need for initial seed

Drawbacks

Need to define stopping criterion

High computational cost & memory requirement

Time complexity $O(n^3)$

Too slow even for medium data sets