

we can divide Machine Learning into two parts

- ① Supervised ML ② Unsupervised ML

→ when ever we talk about Supervised ML algo we will have target/dependent/Supervisor.

Algorithms inside unsupervised

① k-means ~~→~~ ^{updated} → k-mean++ cluster

② hierarchical clustering

③ DBScan clustering

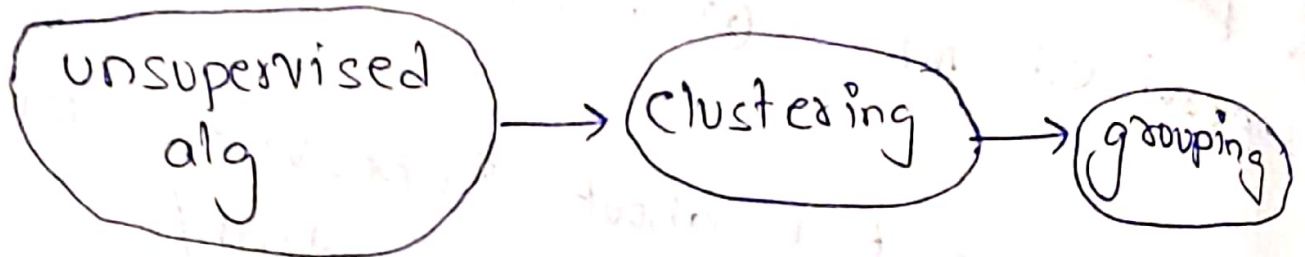
→ Most of time we use this three algo for clustering

Agenda behind taking Unsupervised data

→ If we see below example based on "Height" and "weight" we can predict the "BMI" because "BMI" is Supervisor

Height	weight	BMI
170	60	21
180	65	22
160	70	20

→ In unsupervised algo we will not perform prediction. we will do grouping the dataset.



→ In the below example we are clustering the data based on Country.

ex:-

Height	weight	BMI	Country
170	60	21	IND
180	65	22	UK
160	70	20	USA

→ To do this clustering we will use mathematical eqⁿ

- ① k-mean
- ② hierarchical
- ③ DB-Scan

cluster formation

→ Each cluster is formed based on the distance

① Euclidean distance measure

② Manhattan

③ Cosine

④ Tanimoto

⑤ Squared Euclidean

K-mean

→ we will have dataset. we need to findout similarity b/w dataset by using distance formula like Euclidean

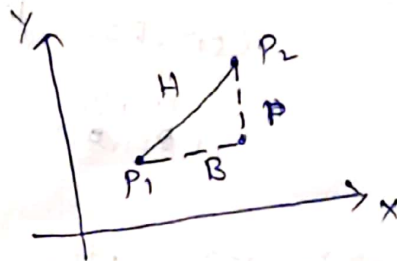
Equation for Euclidean distance

Pythagoras

$$H^2 = P^2 + B^2$$

$$H = \sqrt{P^2 + B^2}$$

$$D(P_1, P_2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$



The main Consider is k-mean are

① Centroid

② Distance

③ mean

→ In k-mean "k" is no. of Centroids

→ "k" is decides based on

① Elbow

② WCSS (within cluster sum of square)

→ ① Inter cluster

② Intra cluster

Evaluation matrix

① dunn index

② Silhouette Index

process

step 1:- Select two data points as Centroids. $k=2$

$C_1(x_1, y_1)$

$C_2(x_2, y_2)$

step 2:- Find the E.D from C_1, C_2 to remaining points. which is nearest C_1 we add it into that cluster.

ED = C_1, NP_1

ED = $C_2, NP(x_3, y_3)$

$$ED = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

step 3:- After adding new points into cluster we update the centroid of that cluster.

$$N_c = \left(\frac{x_1 + x_2}{2}, \frac{y_1 + y_2}{2} \right)$$

EL - blow method

→ This method is used to find the "k" value

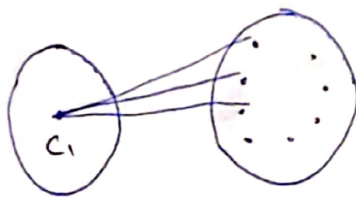
Intra cluster distance:-

→ The distance b/w a data item and the cluster centroid within a cluster



Inter cluster distance

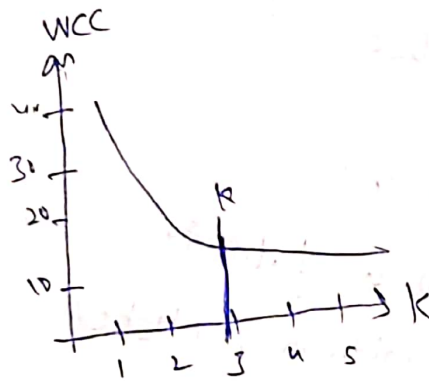
→ The distance b/w centroid of one cluster to data point in another cluster



WCSS:- within cluster sum of square

$$WCC = \sum_{i=1}^n d(c, x_i)^2$$

- To find "k" we will start with one cluster and cal WCC then one more cluster and so.. on
- Finally we see graph like below and choose "k" value



Diff blw k-mean & k-mean++

- In k-mean Centroid initialization is random
- In k-mean++, we take ~~one~~ ~~one~~ point based on largest distance
- In k-mean++ it may choose ~~no~~ outliers

How to validate cluster

These two ways to validate

① Dunn index (Explain in next class)

② Silhouette Score

Silhouette Score

$$S(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}, \text{ if } |C_i| > 1$$

→ The Score Come in blw "-1 to +1"

→ If Score near to +1 it's best model

→ If Score near to -1 it's worst model

$a(i)$ = mean distance from Centroid of the cluster to another point in same cluster

$b(i)$ = mean distance from Centroid of one cluster to another points in different cluster

For more detail check wiki

How to create best model

Custom model (or) Custom learning =
Supervised + Unsupervised

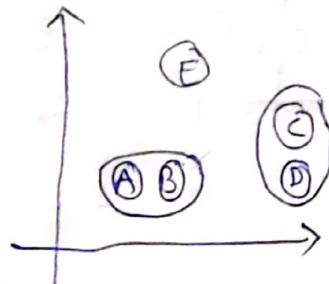
weight	height	Gender
--------	--------	--------

- Create cluster based on Gender
- Then build supervised model on weight, height.

② Hierarchical clustering

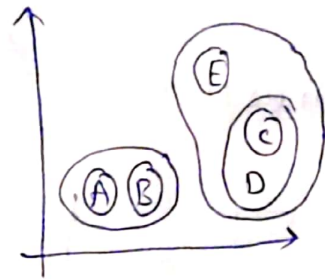
→ Let consider few data point
(A), (B), (C), (D), (E)

→ Group the nearest point

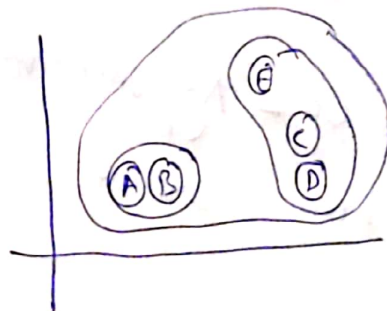


Point matrix
↓
?

→ 'E' is near to (C, D) cluster. So push 'E' to (C, D) cluster



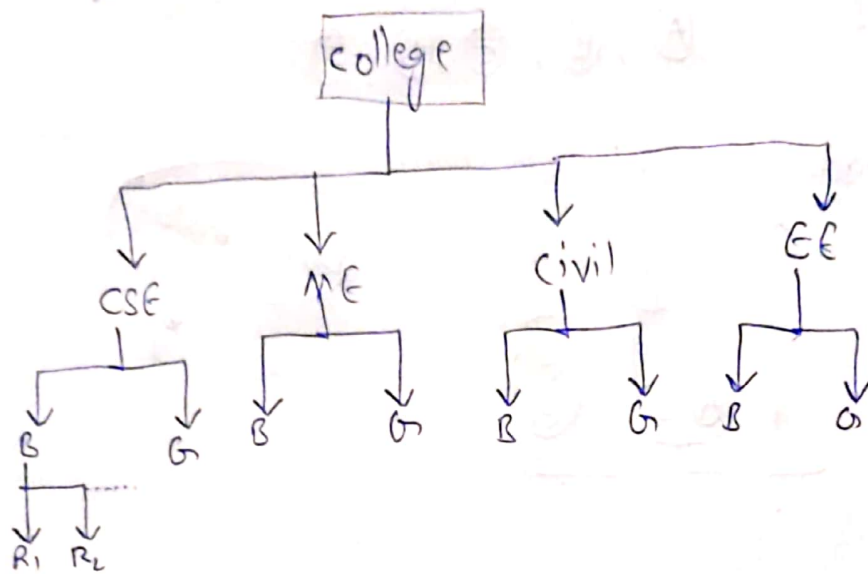
→ (E, C, D) cluster is near to (A, B) cluster. So group all the point into one cluster



→ We will create cluster at each and every point.

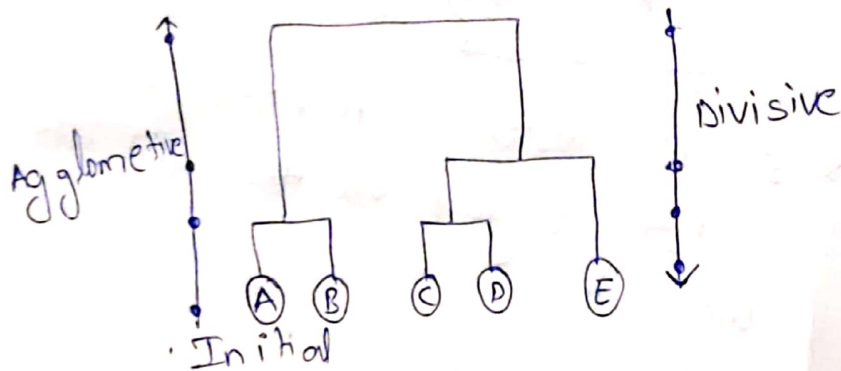
→ Here we don't need 'K' value.

Example of Hierarchical cluster



Den do gram

→ Representation of Hierarchical cluster.



- ★ going down to up - Agglomerative
- ★ going up to down - Divisive

- ① each point is cluster
- ② Bottom to up approach
- ③ Combining all the points as a single cluster

Note: Generally we will prefer ~~Agglomerative~~ Agglomerative method than Divisive

③ DB-Scan (Density Base spatial clustering with application with noise)

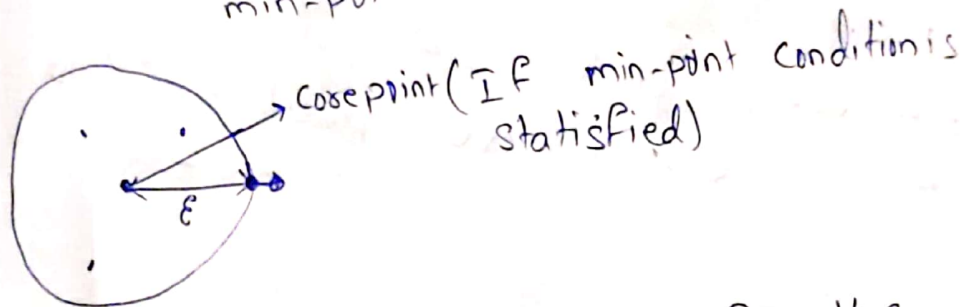
→ It is Density based approach
→ we can find outliers in it
we will create cluster based on below terminology

- ① eps (or) epsilon distance
- ② core point
- ③ Border point
- ④ Noise point-(outlier)
- ⑤ min-point

DB-Scan steps

Step-1 : Initialize ϵ , min point (Hyper)

Ex:- ϵ - Distance = ϵ {Hyperparameter}
min-point = 4



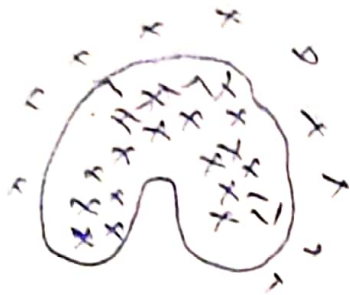
Step-2:- we will take one point from the circle and draw circle with ' ϵ ' distance.

→ IF does not contain min-point in circle.
either it will be noise point (or) Border point

noise point :- IF circle does not have core point

Border point :- IF circle have core point

→ we will use DBScan where k-mean
not able to find cluster



[Read wiki DBScan]

[check sklearn cluster to find best algo
to do cluster]

overview of algos

① k-mean

→ centroid based approach

② Hierarchical cluster

→ check each and every possibility of clusters

③ DB-Scan

→ Density base

point matrix

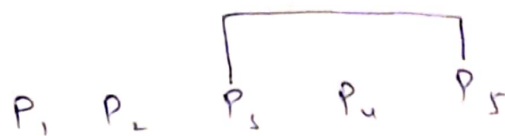
→ Consider 5 points

P_1, P_2, P_3, P_4, P_5

find distance b/w the point and
create cluster

step 1

	P_1	P_2	P_3	P_4	P_5
P_1	0				
P_2	9	0			
P_3	3	7	0		
P_4	6	5	9	0	
P_5	11	10	2		0



step 2

	P_1	P_2	$[P_3, P_5]$	P_4
P_1	0			
P_2	9	0		
P_3, P_5	3	7	0	
P_4	6	5	8	0

$$d(p_1, [p_3, p_5]) = \min[d(p_1, p_3), d(p_1, p_5)]$$