

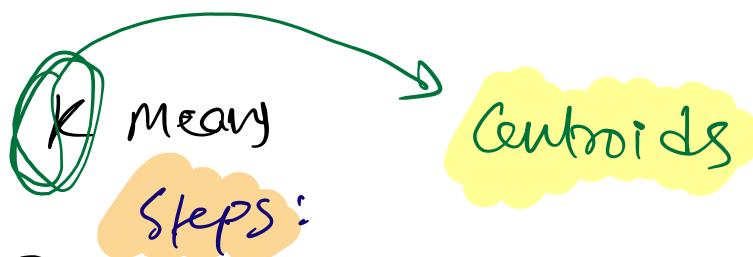
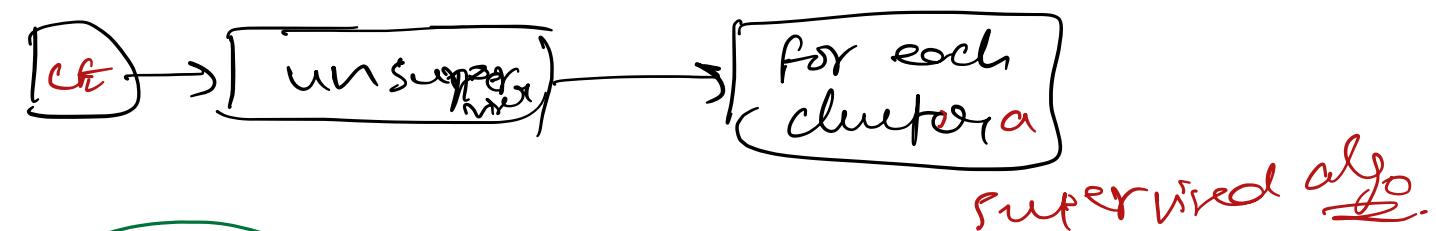
Unsupervised:

- ① K-Means clustering
- ② Hierarchical clustering
- ③ Silhouette score
- ④ DBScan clustering

⇒ we don't have any specific o/p in unsupervised ML

f_1 - f_n → based on these features we divide them into clusters
similar kind of data

→ Can be used in custom Ensemble modelling.

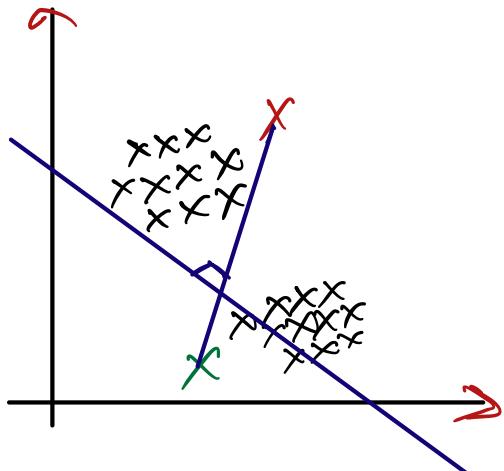


Steps:

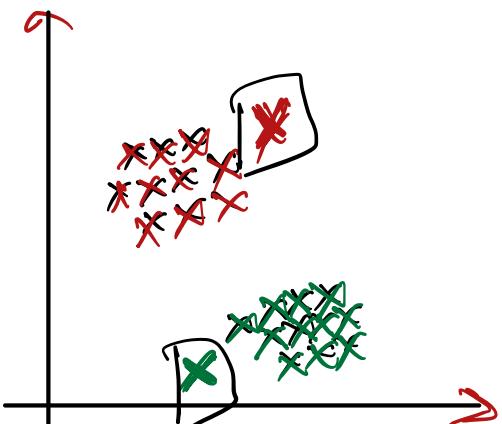
- ① we try diff 'k' values → select suitable
- ② we initialize 'k' number of centroids.
- ③ Compute the avg to update the centroid

Sit back, and enjoy
the show :D

GREAT
JOB!

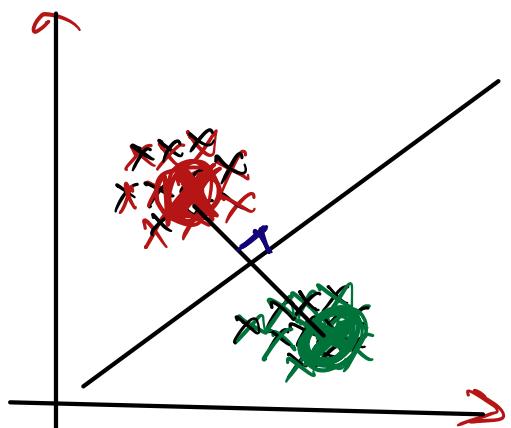


- ① first we consider plotting the whole data which we have.
- ② then we take '2' random points and try to find a perpendicular line.
- ③ After that we use Euclidean dist to find the centroid.



→ So, from the above results, we consider the points near to red as Red Category and green as green category.

→ after this, we find \bar{x}_{avg} of every category and then we will get a new centroid which is inside that particular category.



→ now, we can see that, all the points are categorised properly

→ we call them as clusters, and those centroids are represented by the 'K' values

Note:

→ If there are '2' clusters then value of $k=2$

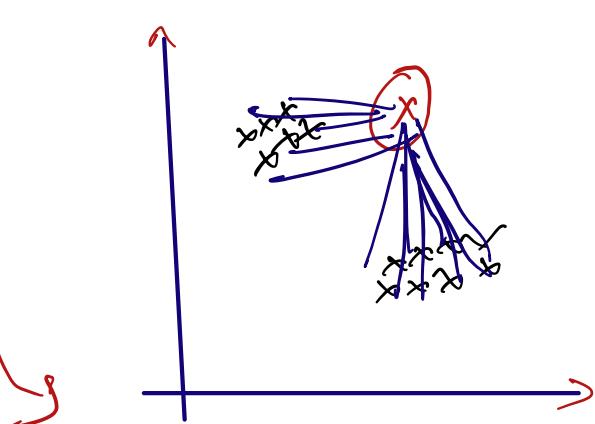
[how do we decide the k-value?
→ we use Elbow method] → optimize "K" value

→ we take iterations, and for every iteration we will construct a graph with respect to K-value and WCSS

WCSS

within cluster sum of square ---

→ before plotting the graph, let's understand diff graphs w.r.t to 'K' value.

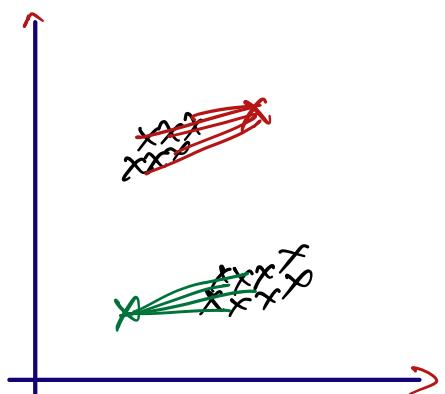


$\boxed{K=1}$

→ we took one centroid, and we try to find dist from all the points.

so, after considering distances we sum up and find squares.

→ so, this is going to give us a very high $\boxed{\text{WCSS}}$



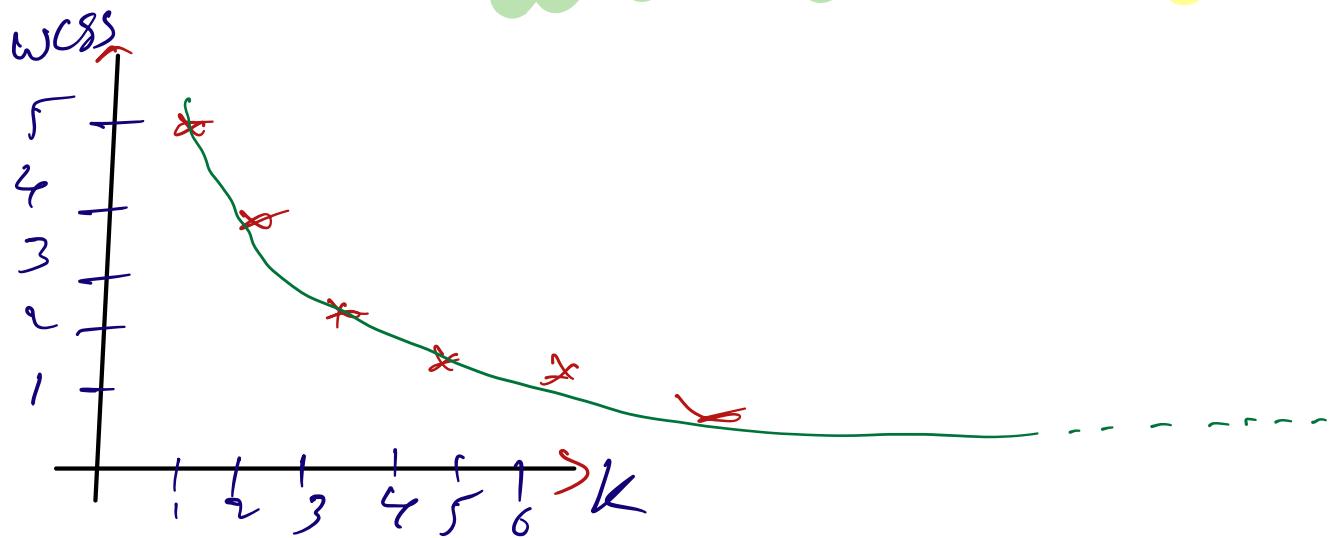
$\boxed{K=2}$

→ in this case, we have 2' centroids and their individual distances have been

been $\boxed{\text{WCSS}}$ value compared to $\boxed{K=1}$

→ So, the conclusion is, as we take more 'K' values WCSS value will decrease

→ and we plot the graph for 'k' and and we call it **Elbow method**



So, for what value of 'K' WCSS is giving an elbow style of graph (the drop in the graph)

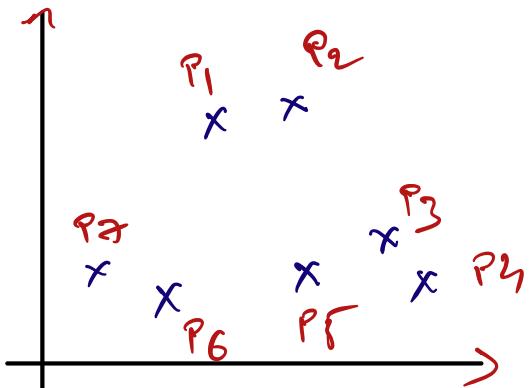
we prefer that particular value of K for dividing values into clusters.

Note: we can also say, where there is abrupt change, we consider that 'K' value.---

$$WCSS = \sum_{i=1}^n (c_i + x_i)^2$$

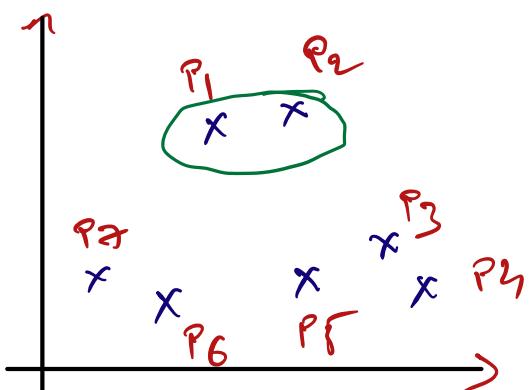
Hierarchical clustering:

→ let, we have various points plotted as shown below.

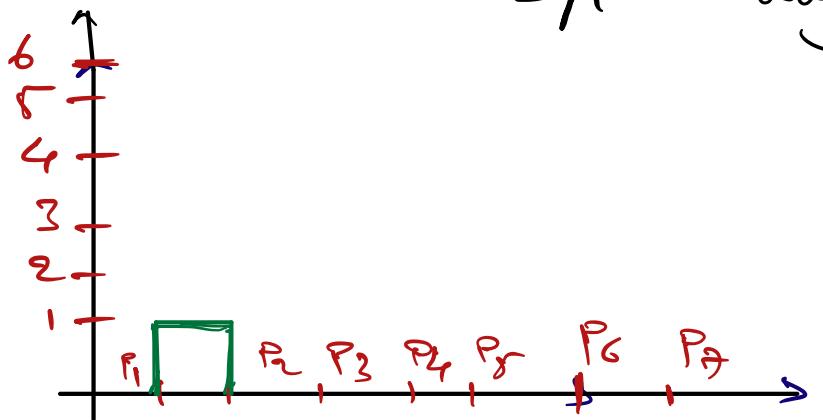
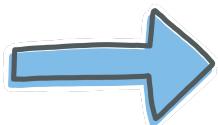


→ now, according to the distance b/w points, we group point P₁ & P₂.

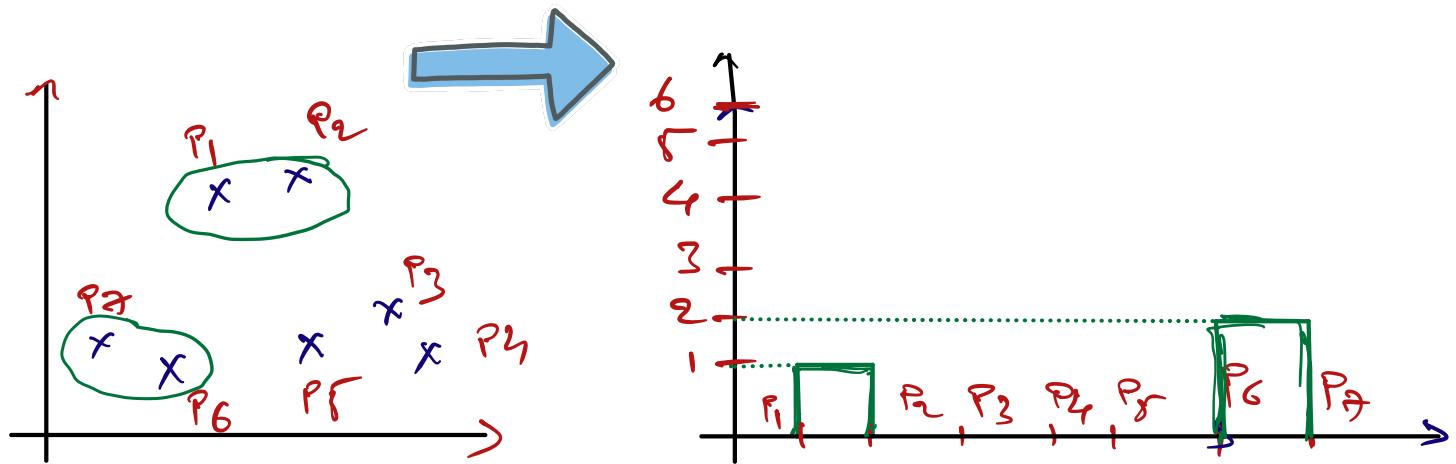
→ firstly, let's consider P₁ and P₂, as the distance b/w them is less, we group both of them as one entity.



→ now, depending on the grouping, in the backend we create another graph in the below shown way.



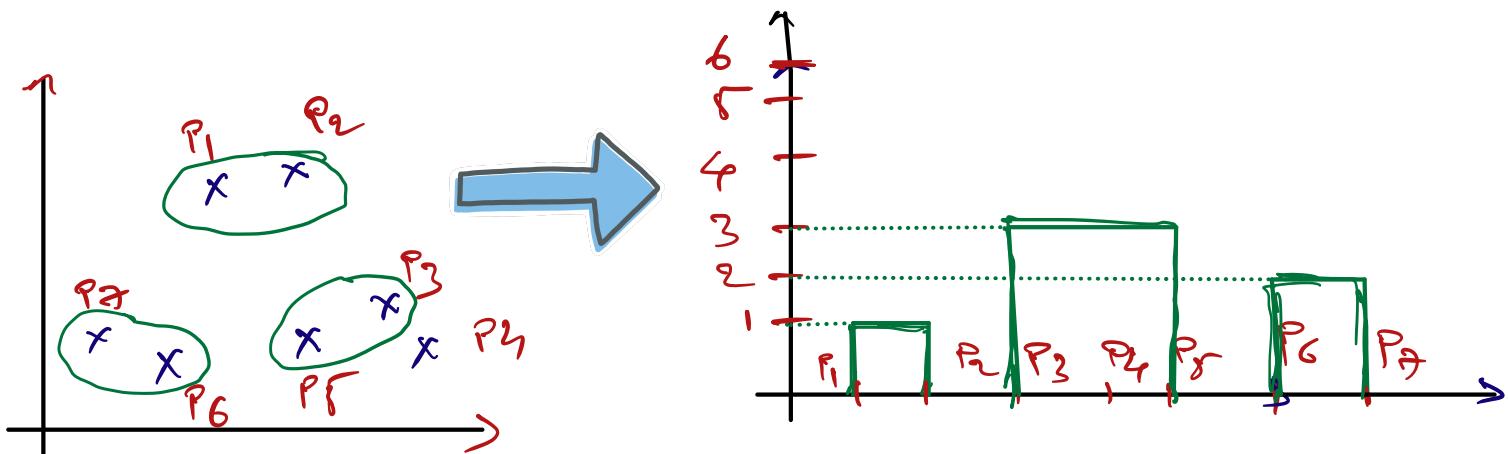
→ in the same way, we will combine P_6 and P_7 as well.



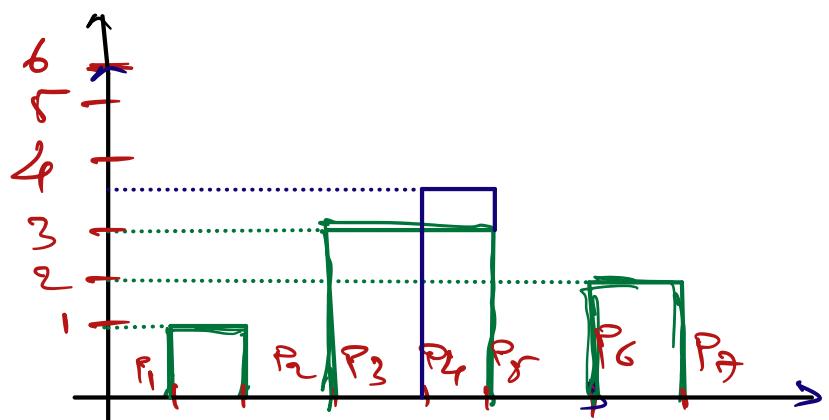
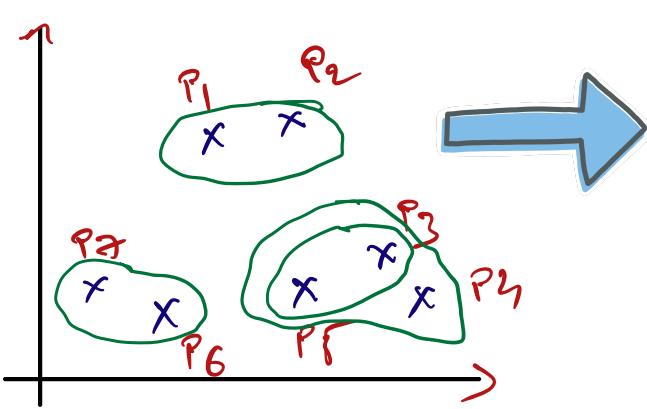
→ let's assume that, the distance b/w P_6 and P_7 is greater than P_1 and P_2 , for ex:- 9

→ next, we will group P_3 and P_4 and also

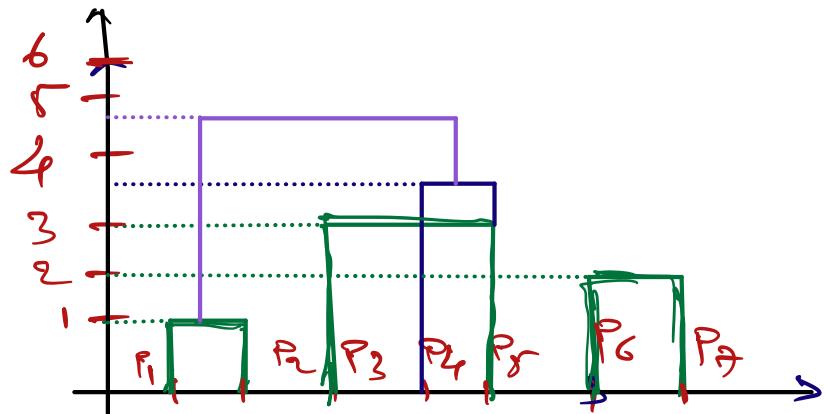
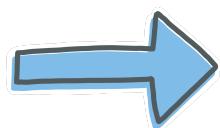
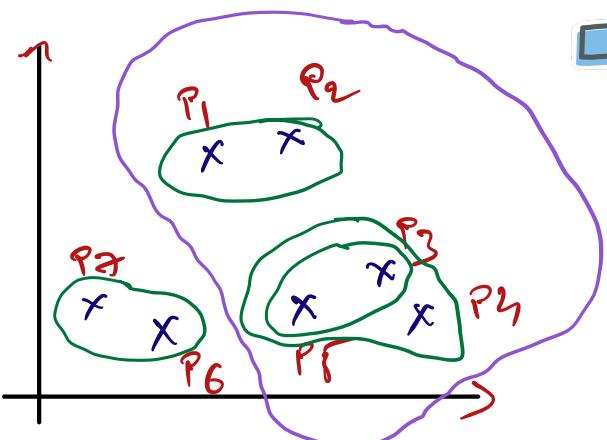
lets assume its dist is even more greater than the previous two



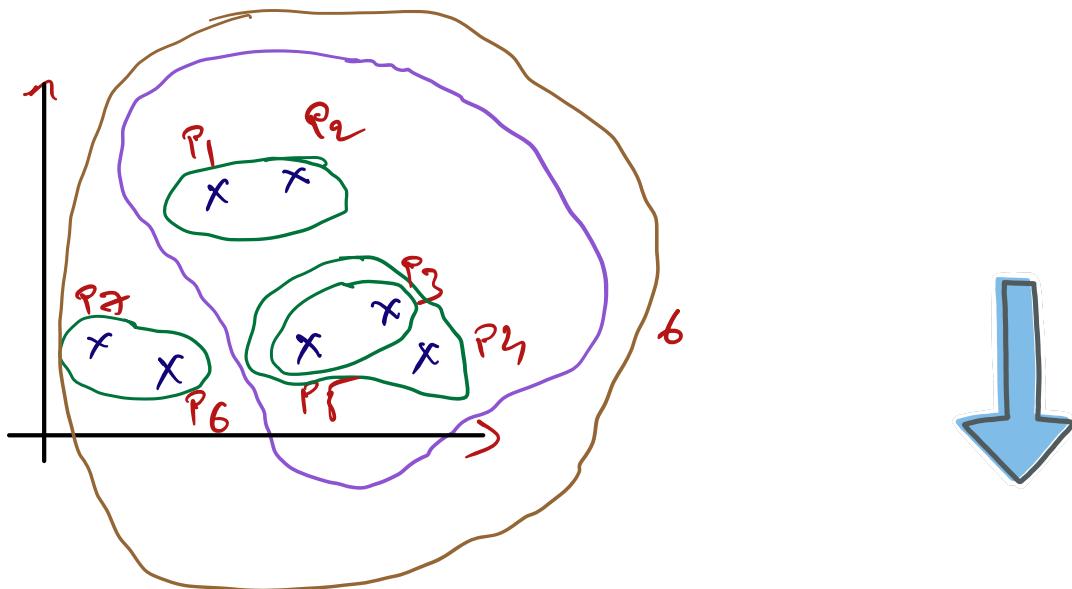
→ Now, we will combine P_4 with P_3 and P_7



→ as we can see, P_3, P_4, P_7 are nearest to P_1, P_2 group, we combine them.

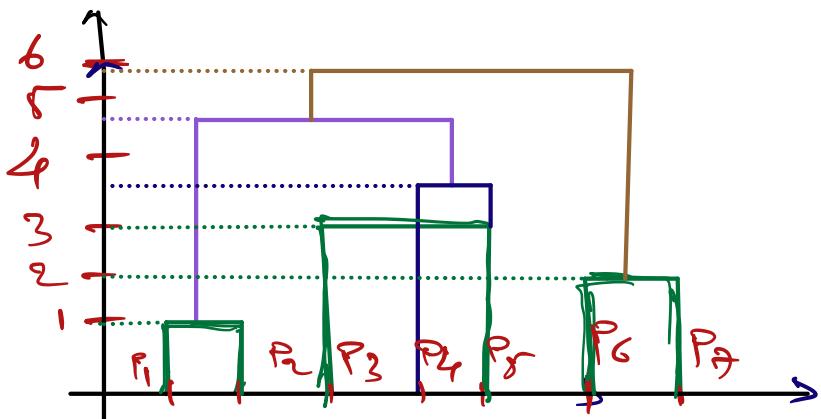


→ now, as we can see P_6, P_7 is near to the purple group, we join them.



This is !!
called as

Dendrogram



How do we find out, how many groups should we consider for the given number of points???

Ans: You need to find the longest vertical line that has no horizontal line passed through it in a Dendrogram

**DON'T
FORGET**

note: there should not be any horizontal line even if we extend any " "

TEST:

Maximum time is taken by which clustering?

Ans:

Hierarchical clustering.

Dataset is small \Rightarrow Dendrogram

Dataset is large \Rightarrow K-means

note:

\Rightarrow we should take centroids very far from the data points, else it might create confusion with no. of clusters.

\Rightarrow we solve this error with the help of

K means ++

↳ it will make sure, centroids are far.

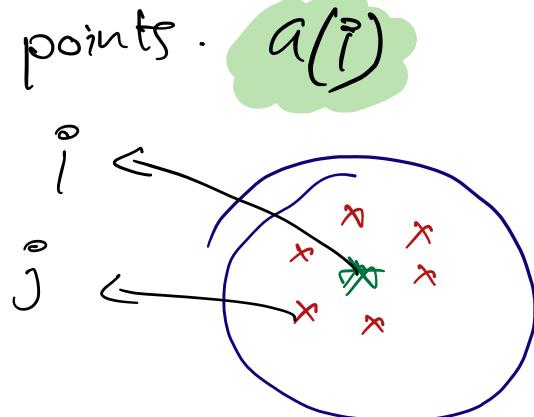
Validating clustering model:

→ we use Silhouette Score.

Steps:

(1) In a cluster, we compute the distance b/w centroid and the points. $a(i)$

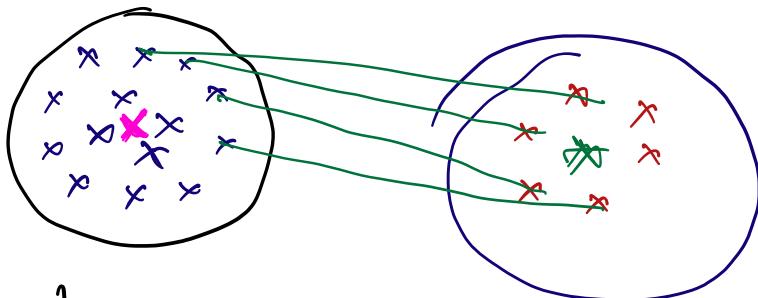
⇒ Summation
⇒ Avg



formula:-

$$a(i) = \frac{1}{|C_i| - 1} \sum_{\substack{j \in C_i, \\ i \neq j}} d(i, j)$$

(2) we will calculate distance from each and every point from one cluster to the another cluster
⇒ avg



TEST

If we have a good cluster, whether $a(i)$ will be greater or $b(i)$

Ans:

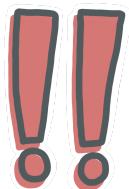
$$b(i) \gg a(i)$$

→ Value of silhouette clustering will be in range $[-1 \text{ to } +1]$

formula:

$$b(i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{j \in C_j} d(i, j)$$

formula:



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

and

$$s(i) = 0 \quad \text{if } |C_i| = 1$$

This page is for future notes :—

DB Scan Clustering :-

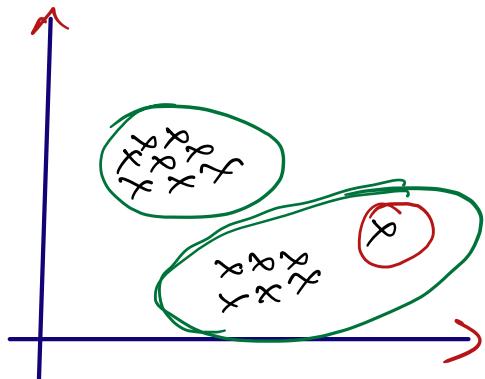
Density Based Spatial Clustering of Applications

* Epsilon

with Noise

- ① Min Points
- ② Core points
- ③ Border points
- ④ Noise points

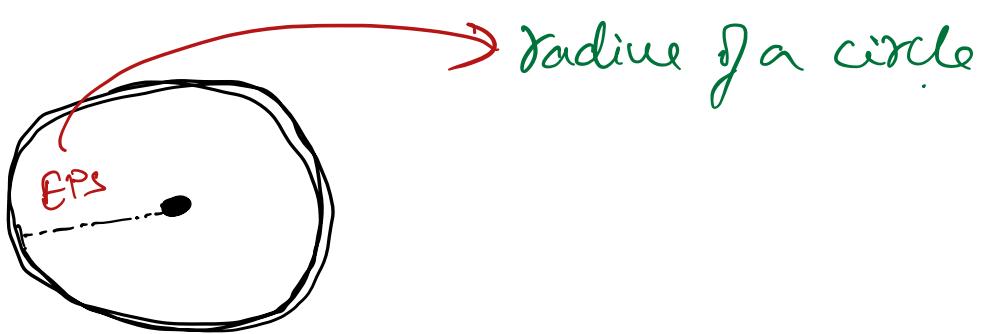
why we use DBScan?



→ In K-Means, we consider this core as 2 clusters
→ But this extra point is called as an outlier (or) noisy point, as it is separate from the similar data.

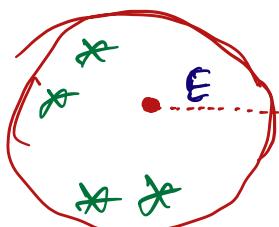
→ To overcome this error we use DBScan, which leaves out the outliers when clustering.

Epsilon:



Min point:

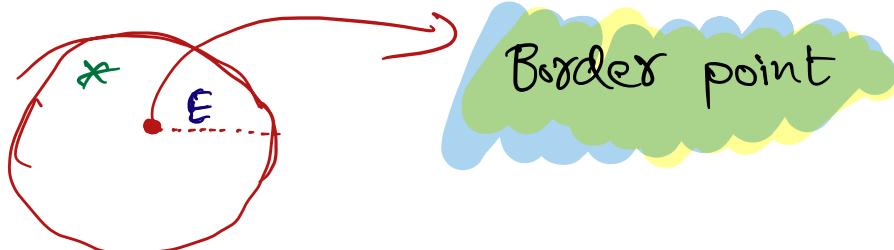
If it is a hyperparameter



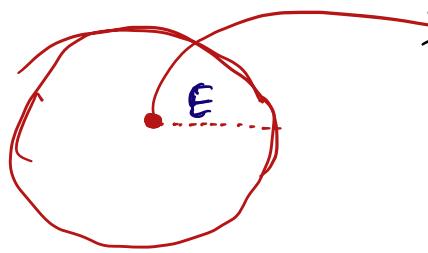
If $\text{min pts} = 4$

→ If in this Epsilon distance, we have the same min no. of min points, then the red point is called a Core point

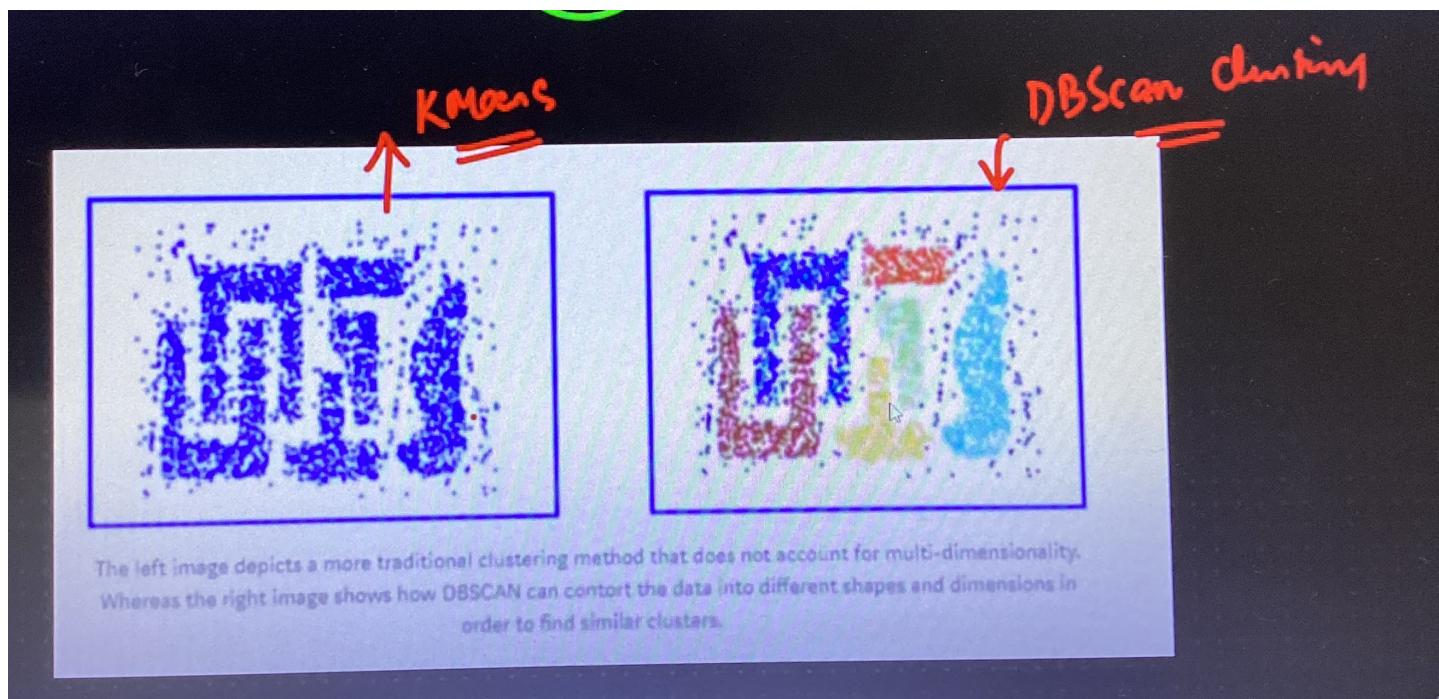
→ If we don't have points equal to our min point, we call it



→ If we don't have any point inside our epsilon distance, we call



⇒ So, when we have a noise point it will get neglected.



Agglomerative Clustering:-

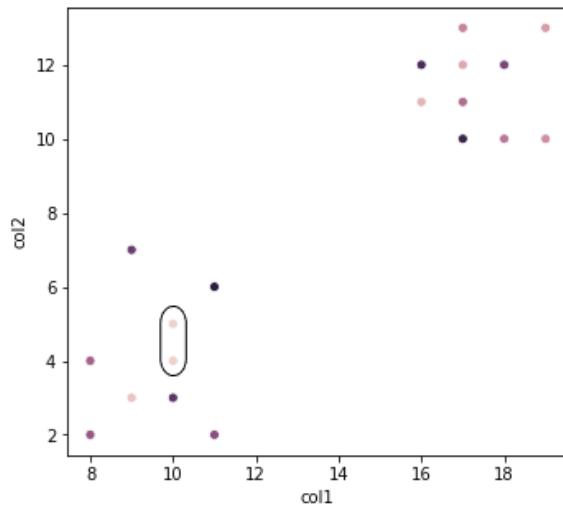
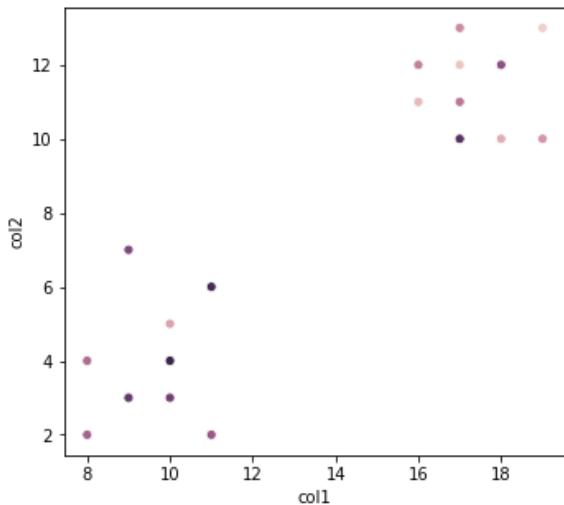
- it is a part of hierarchical clustering and we create dendograms as well for this
- More points are considered as one cluster and far points are divided into diff clusters

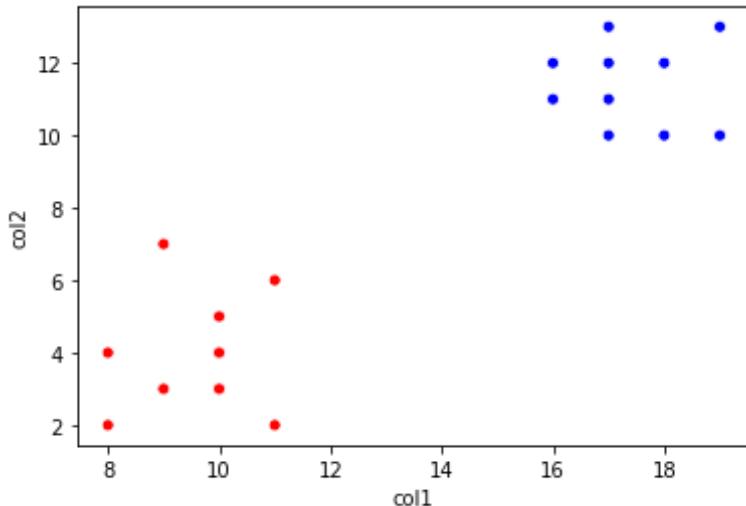
intuition:

→ It is a bottom-up approach, initially, each data point is a cluster of its own, further pairs of clusters are merged at one move up the hierarchy

Steps:

- ① Initially, all the data-points are a cluster of its own
- ② Take two nearest clusters and join them to form one single cluster.
- ③ proceed recursively step 2 until you obtain the desired number of clusters.





How?

- To obtain the desired number of clusters, the no. of clusters need to be reduced from initially being (n) cluster (n equals the total number of data points)
- Two clusters are combined by computing the similarity b/w them.

There are some methods which are used to calculate the similarity b/w two clusters,

- ① Distance b/w two closest points in two clusters
- ② Distance b/w two farthest points
- ③ The avg distance b/w all points
- ④ Distance b/w centroids of two clusters.

⇒ There are several pros and cons of choosing any of the above similarity metrics.

Inverse:

- The inverse of agglomerative clustering is **divisive clustering** which is also known as **DIANA** (Divise Analysis) and it works in "top-down" manner.
- It begins with the ~~root~~, in which all objects are included in a single cluster.
- At each step of iteration, the most heterogeneous cluster is divided in two.



Association Rule Mining:

- it is a method for identifying frequent patterns, correlations, associations (or) causal structures in data sets found in numerous databases such as relational databases, transactional databases, and other types of data repositories...
- it can also deal with categorical/object data...
- Given a set of transactions the goal of association rule mining is to find the rules that allow us to predict the occurrence of a specific item based on the occurrences of other items in the transaction.

It consists of two parts:

- ① an antecedent (if) and
- ② a consequent (then)

it is found in data

located in conjunction with the antecedent.

ex) "if a customer buys bread, he's 70% likely of buying milk"

antecedent = bread

consequent = milk

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

Image Source

ex of association rule mining

⑧ Association Rule Mining: Basic Definition:

Support Count (σ):

It accounts for the frequency of occurrence of an itemset.

In the above example:-

$$\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$$

Frequent Itemset:

→ It represents an itemset whose support is greater than (or) equal to the min threshold.

Association Rule:

→ It represents an implication expression of the form $\boxed{x \rightarrow y}$

→ here, X and Y represent any '2' itemsets....

Exr $\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$

Rule Evaluation Metric:

Support(S): It is the no. of transactions that include items from the $\{x\}$ and $\{y\}$ parts of the

Rule as a percentage of total transaction.

$$\text{Support} = \frac{\text{Count}(X+Y)}{\text{total}}$$

It is a fraction of transactions that include both \boxed{X} and \boxed{Y}

Confidence (C):

This ratio represents the total no. of transactions of all the items in $\{A\}$ and $\{B\}$ to the no. of transactions of the items in $\{A\}$

$$\text{Conf}(x \Rightarrow y) = \frac{\text{Supp}(x \text{ and } y)}{\text{Supp}(x)}$$

It counts the no. of items each time in \boxed{Y} appears in transactions that also include items in \boxed{X}

lift (L): The lift of the rule $X \Rightarrow Y$ is the confidence of the rule divided by the expected confidence.

$$\text{Lift}(X \Rightarrow Y) = \text{Conf}(X \Rightarrow Y) \div \text{Supp}(Y)$$

Lift values near 1 indicate that 'X' and 'Y' almost always appear together as expected.

$\boxed{\text{lift} > 1}$ they appear together more than expected....

$\boxed{\text{lift} < 1}$ less than expected....

Applications of Association Rule Mining:

1. Market-Basket Analysis
2. Medical Diagnosis
3. Census Data

Algorithms of Association Rule mining:

- ① Aprori Algo
- ② Eclat Algo
- ③ FP - Growth Algo

Apriori Algo:

- It is for finding frequent itemsets in a dataset for boolean Association rule.
- It uses prior knowledge of frequent itemset properties.
- we apply an iterative approach (or) level-wise search where k-frequent itemsets are used to find $k+1$ itemsets.
- we use Apriori property which helps by reducing the search space.

