



# Clustering HandWritten Notes

- Darshan R M

"Learn, Share,  
and  
Collaborate"

Week - 19

Unsupervised Learning

clustering

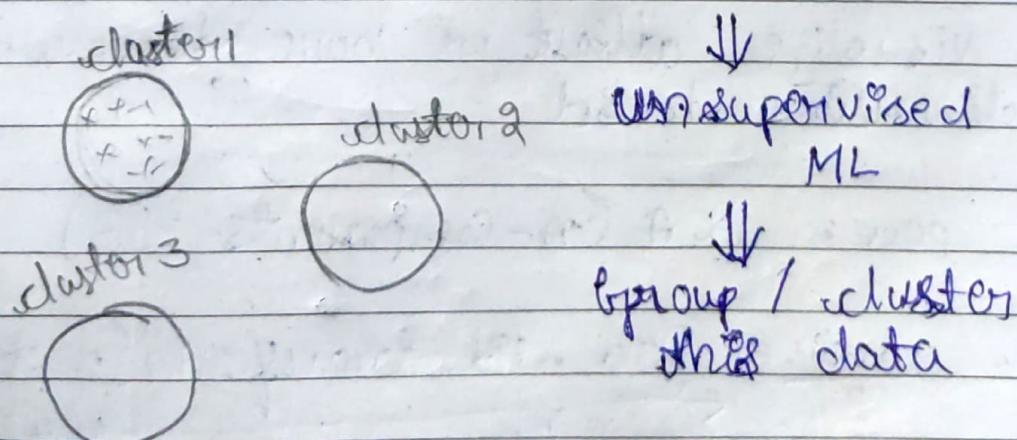
clustering & grouping your data  
similar clusters}

\* Supervised Learning → we have both  
I/p and O/p  
[Independent] [Dependent]

↓  
solve ↗ classification.  
↗ Regression.

\* Consider,

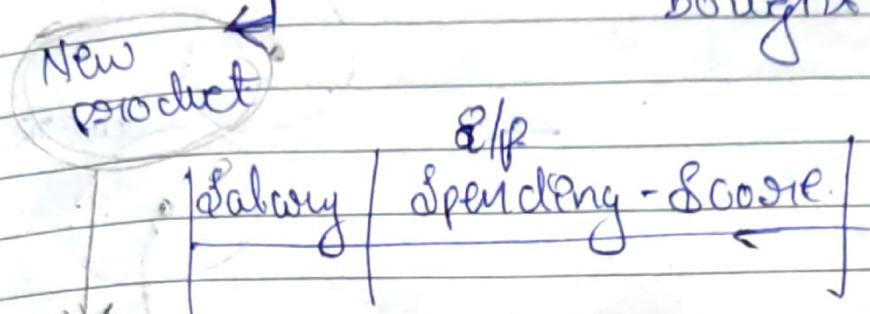
Age	Experience (salary)	c/p	No o/p feature
-	-	-	
-	-	-	
-	-	-	
-	-	-	



Q1:

## Customer segmentation

Product = data of customer bought product.



From existing data

I may get know who  
may buy the product

we can give  
discounts to  
those customers  
and "increase"  
the rate of buying.

15%      20%

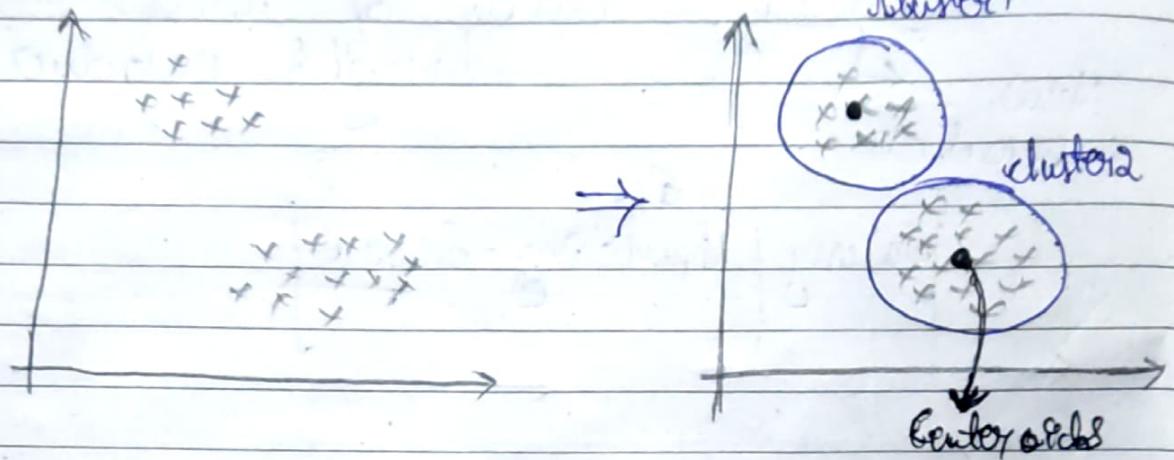
→ more if we  
give 20% then more people  
buy.

## \* Unsupervised ML ↗ Algorithm

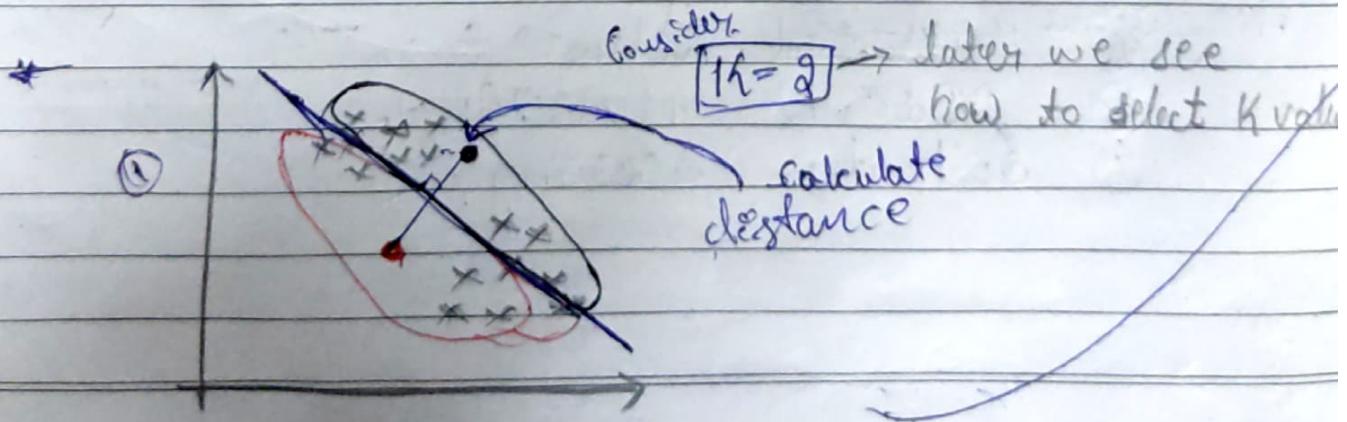
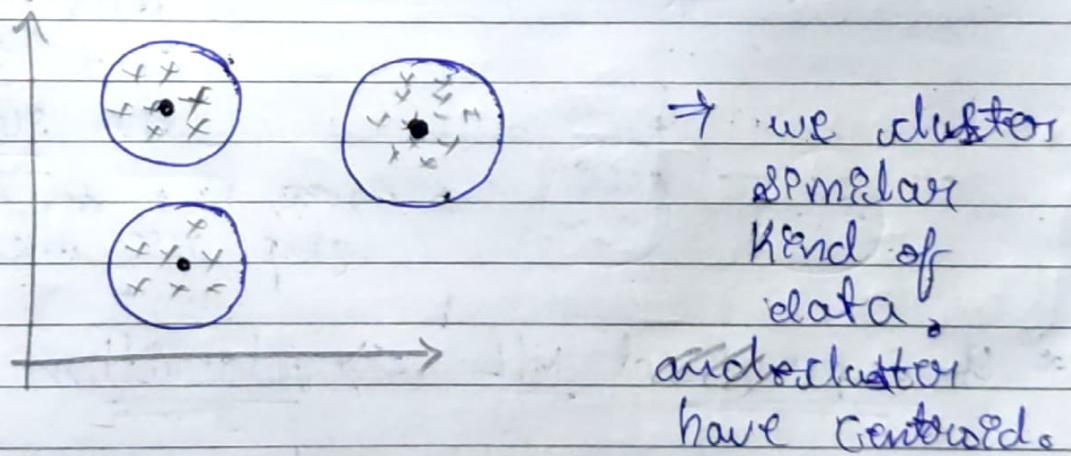
- ① K Means algorithm.
- Hierarchical ~~Hierarchical~~ Clustering ↗ Validate
- ③ DBSCAN clustering
- ④ Silhouette Scoring

## \* K-Means clustering algorithm

+ consider,



\* of 3. clusters



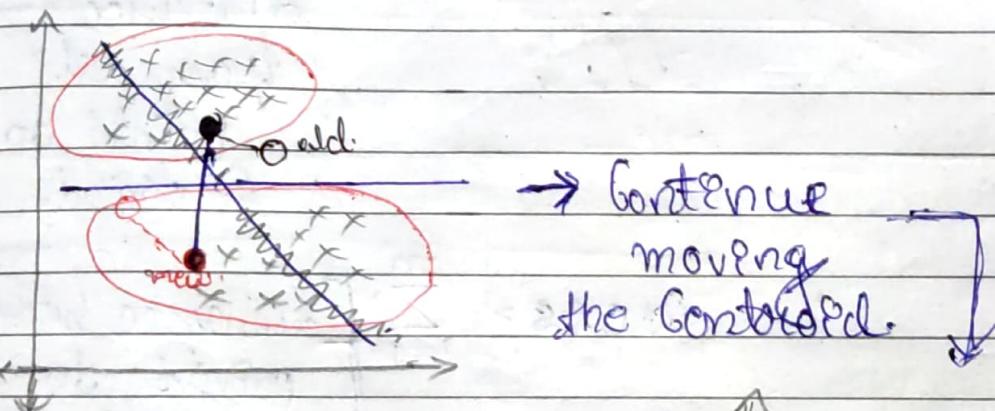
It means

steps ① Internalize some  $K \rightarrow$  centroids

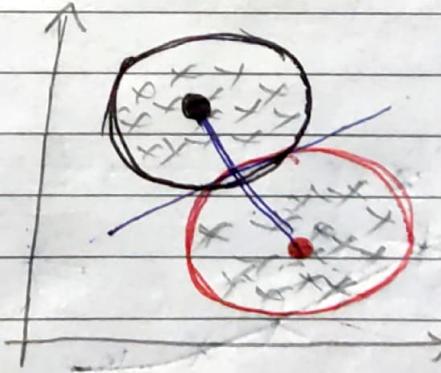
② Find out points that are nearest to the centroid, group them.

splat

③ Move the centroids  $\rightarrow$  average distance.



$\rightarrow$  Continue moving the Centroid.

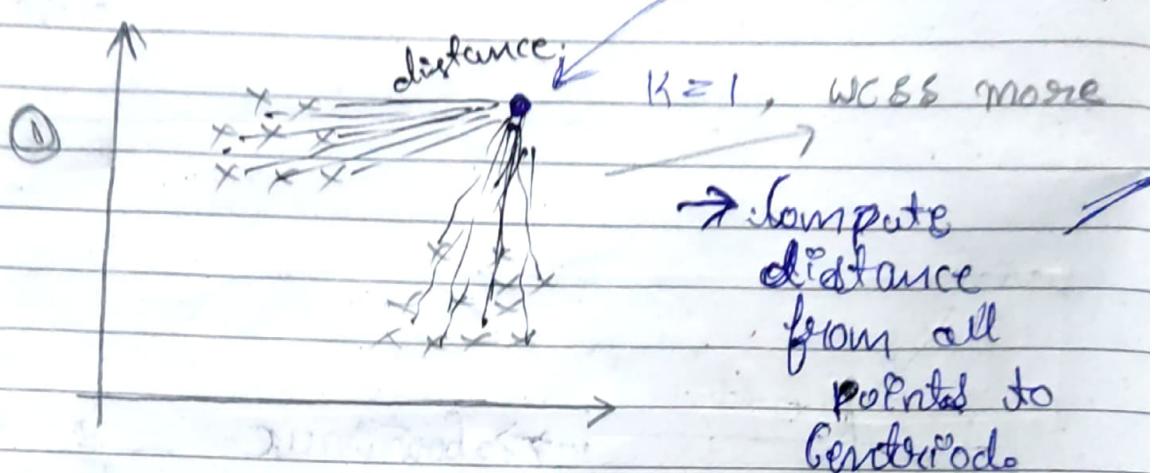


Q) How do we select K values?



WCSS = within cluster sum of squares

initial K = 1 to 20 random point initialization



WCSS graph

$$WCSS = \sum_{i=1}^n \text{distance between points to nearest centroid}$$

Elbow method

where we have to find a point where there is abrupt decrease in WCSS.

we have to select that point.

it will be high.

②



$$K=2,$$

WCSS is less

good one

~~so~~

WCSS

Elbow method.

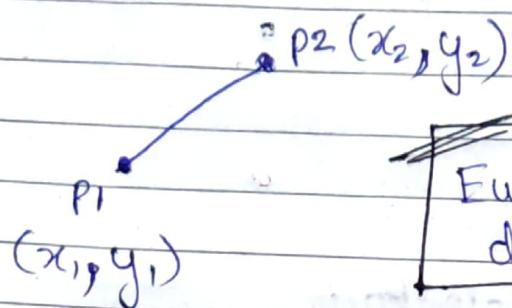
→ abrupt decrease in WCSS

not much [stabilizing]  
change in WCSS

1 2 3 4 5 6 7 K

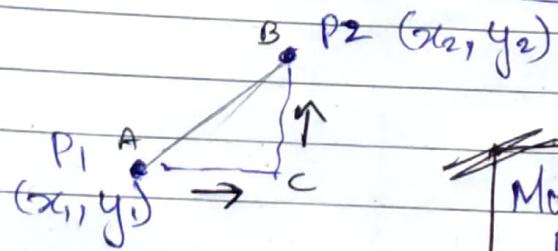
select this  
K value

## \* Eucleadian distance



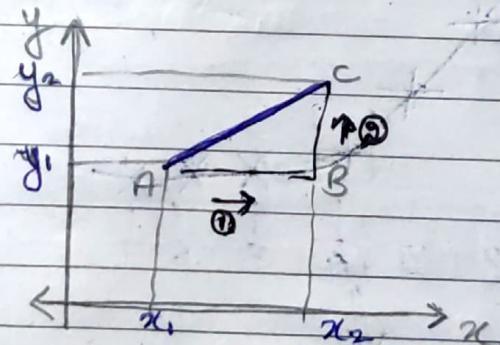
$$\text{Eucleadian distance} = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

## \* Manhattan distance



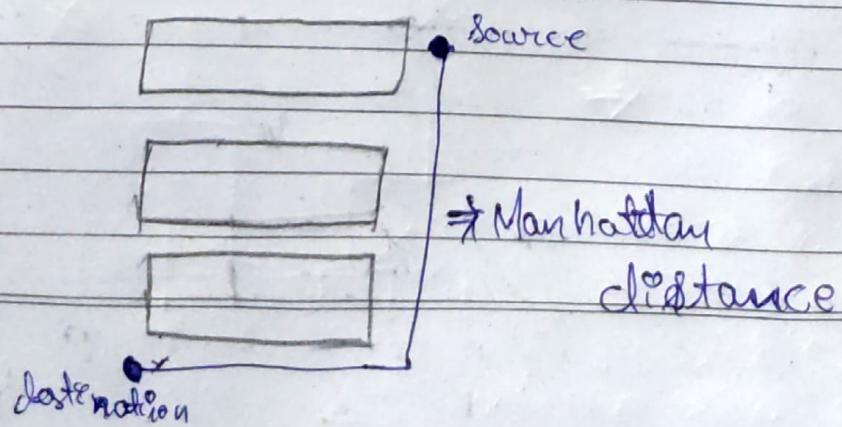
$$\text{Manhattan distance} = |x_2 - x_1| + |y_2 - y_1|$$

we have to take AC & CB



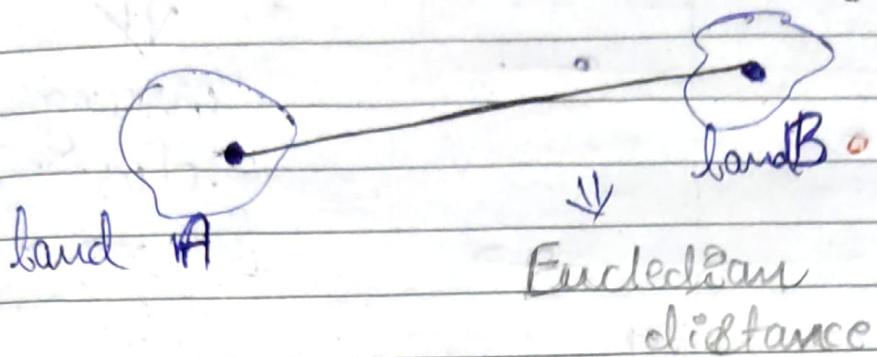
Q) when to use ~~and~~ what?

⇒ \* In US City have building



⇒ Manhattan distance

## ④ Aisi traffic: control.

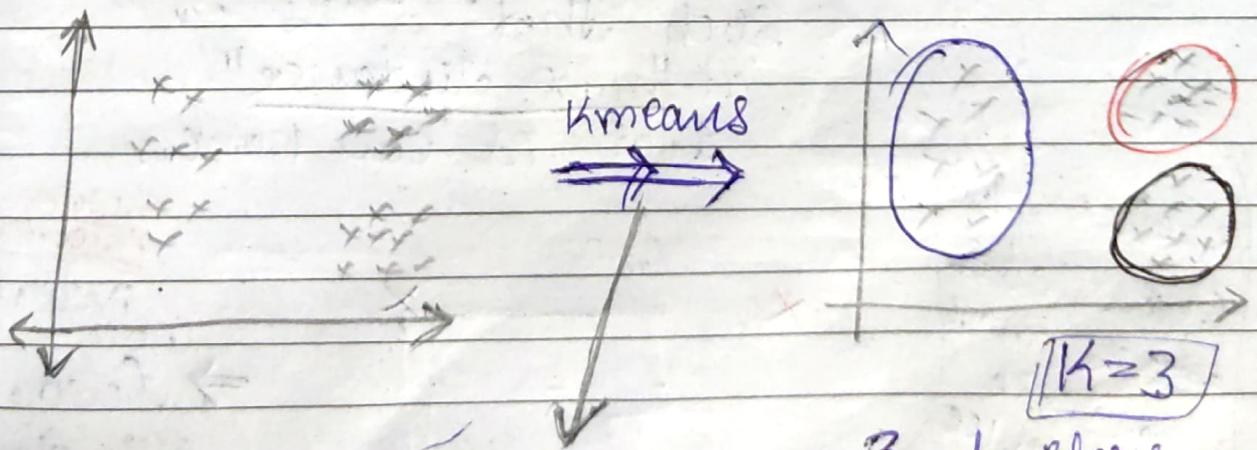


- \* what happens when 2 points are initialized near?



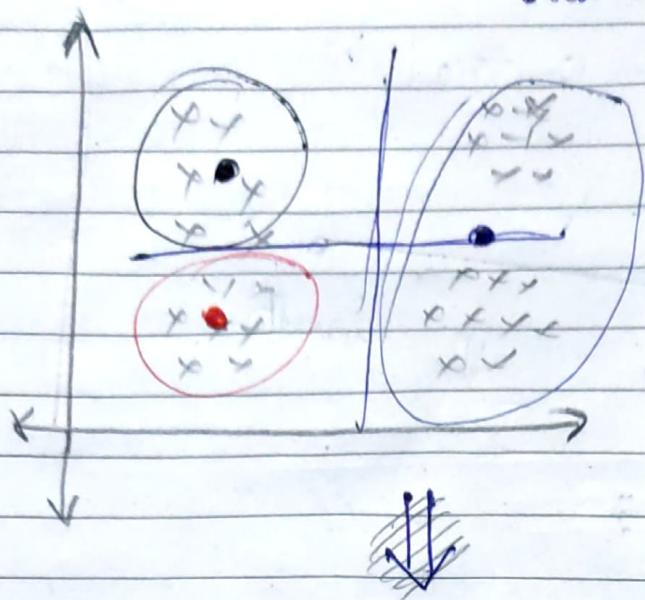
- \* Random Initialization Trap [ $K$  means++]

- \* Consider, data points



During this  
when we do  
random initialization  
what if initialization  
goes "wrong." ↪

randomly initialized centre



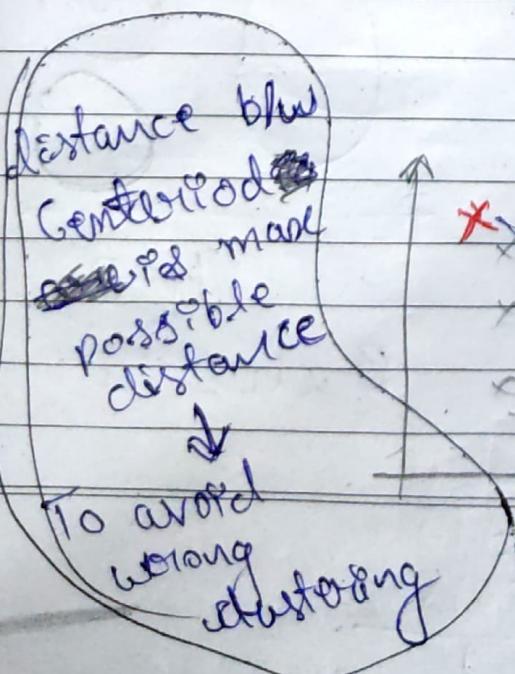
wrong clustering

[We need cluster which we have behind]

So, to avoid wrong initialization.

We use "K means++" initialization Technique.

⇒ Here, we initialize all the "Centroid" such that it is at "max. distance" that it can have.



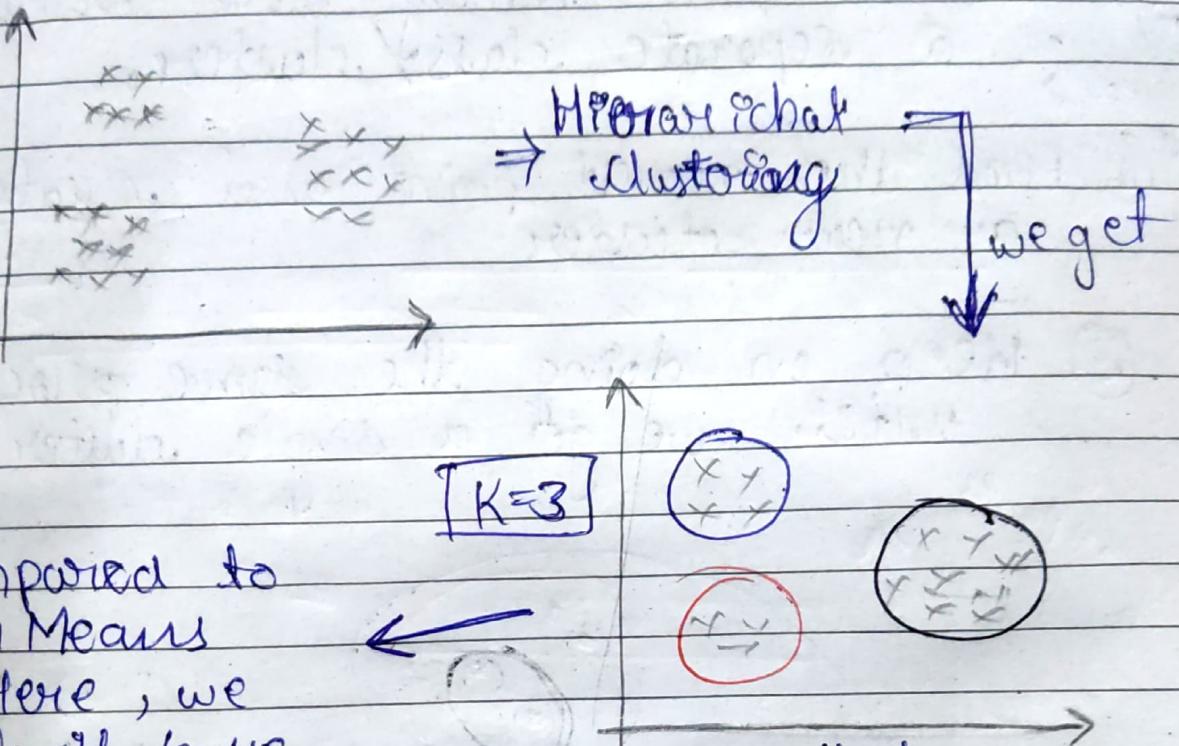
(mark)

answer interview

Initially here, ⇒ Centroid

are at maximum distance between each other.

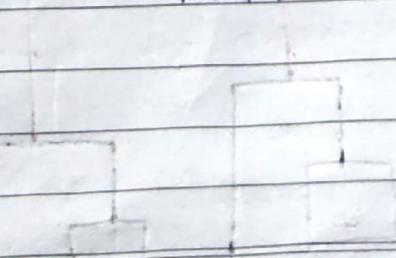
## \* Hierarchical clustering



## \* 2 Types of Hierarchical clustering.

- enough to understand ① Agglomerative.  $\Rightarrow$  Combining  $\rightarrow$  Reverse  
one type ② Divisive.  $\Rightarrow$  dividing

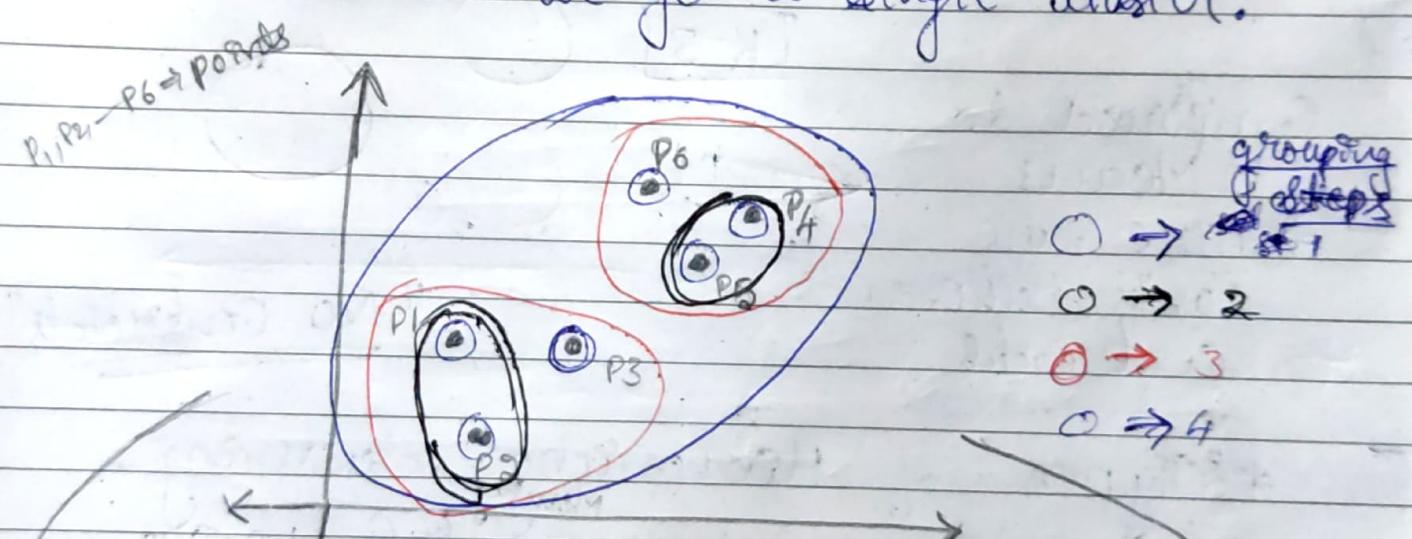
P.T.O.



## \* Agglomerative

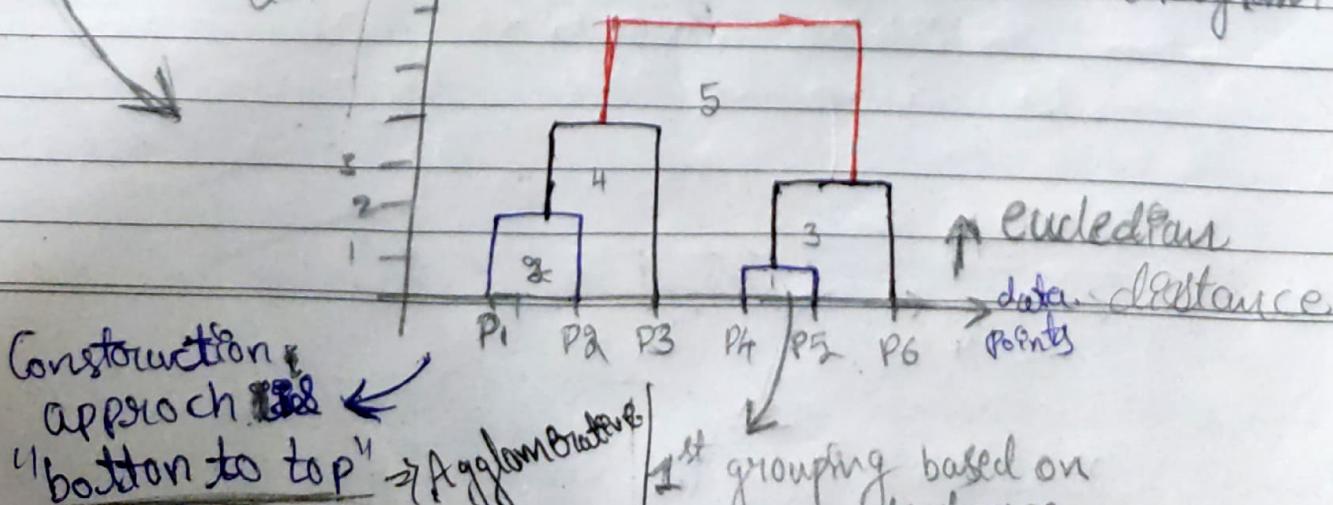
Step 8

- ① For each point initially we will consider as a separate class / cluster.
- ② Find the nearest point and create a new cluster.
- ③ Keep on doing the same process until we get a single cluster.



## \* Dendrogram

Euclidean distance



w.r.t  $K=2$   
but how can we find from Dendrogram?

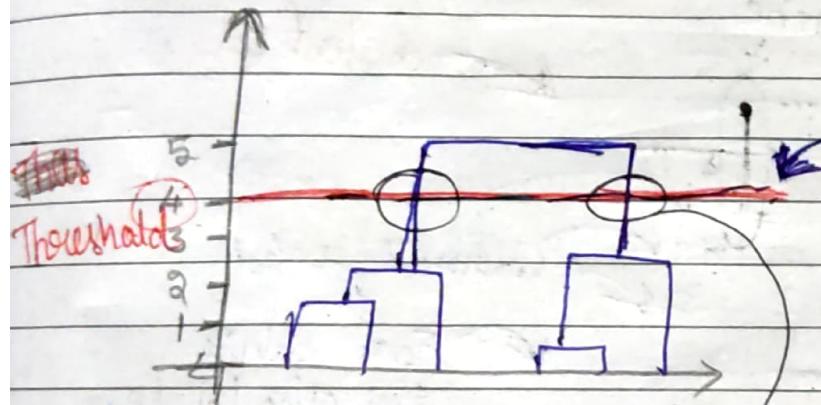
\* In dendrogram,

↑ bottom to top approach  $\Rightarrow$  Agglomerative

↓ top to bottom approach  $\Rightarrow$  Divisive.

\* How do we find K values?

→ found based on "threshold of euclidean distance"



Create a line  
tells that  
euclidean distance  
b/w all the  
points/cluster  
should not be  
more than 4.

Threshold line  
passes through  
2 points  $\Rightarrow$   $K=2$

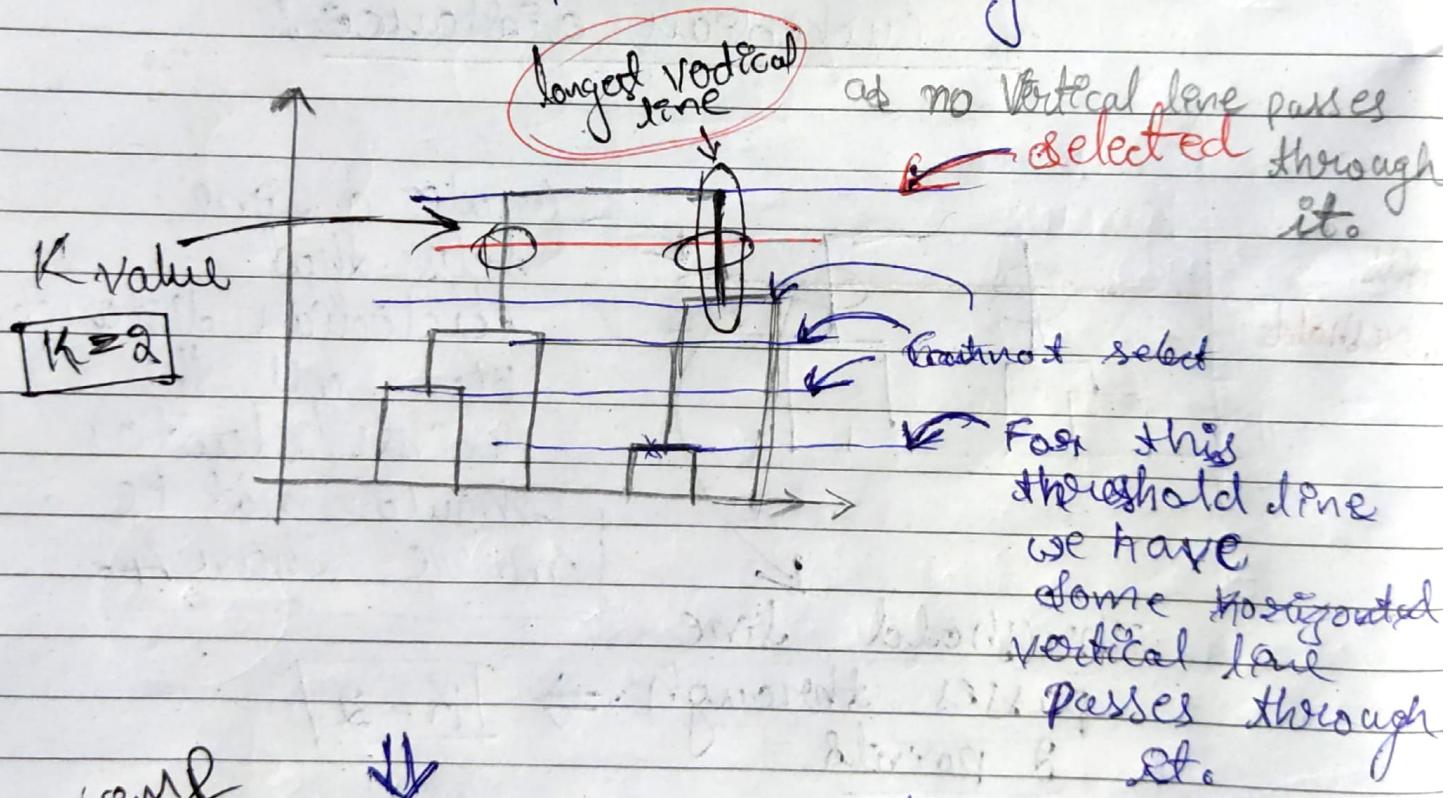
.. if we "decrease"  
threshold  
euclidean  
distance  
 $\xrightarrow{\text{then}}$  "No of clusters  
increases"

threshold  
No of clusters

\* But How do we find threshold value?  $\downarrow$

With the help of dendrogram

$\Rightarrow$  we have to find the longest horizontal vertical line so that none of the horizontal lines passes through it.



- ~~Ques~~  $\downarrow$
- ① Select the longest ~~horizontal~~ vertical line such that no ~~vertical~~ horizontal passes through it.

~~Advantages~~

## \* Hierarchical vs K-Means

Based on Scalability and Flexibility

- ① Dataset size  $\rightarrow$  Huge  $\rightarrow$  K-Means  
Small  $\rightarrow$  Hierarchical
- as we have to  
construct 'Dendrogram'.

- ② K-means  $\rightarrow$  Numerical data

Hierarchical  
clustering

Cuisine

Variety ~~& Similarity~~  
of data

Cuisine  
Similarity  
 $\rightarrow$   $\exists$

- ③ Centroids  $\rightarrow$  Elbow method.

difficult to  
find No of centroid

Hierarchical  $\rightarrow$  Can find easily.

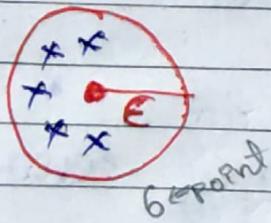
## \* DBSCAN

- $\rightarrow$  Core point
  - $\rightarrow$  border point
  - $\rightarrow$  outlier.
- } able to handle non linear data.

Example, ~~minpts = 3~~ minpts = 4,  $\epsilon$  = radius

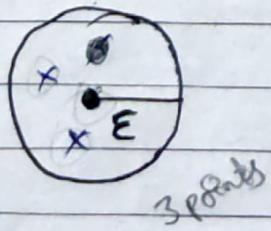
Select value using Hyper parameter tuning & also by Silhouette score

### \* Core point



④ No of points within the  $\epsilon$  ~~should~~ should be  $\geq \text{min points}$

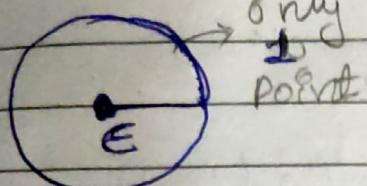
### \* Border point



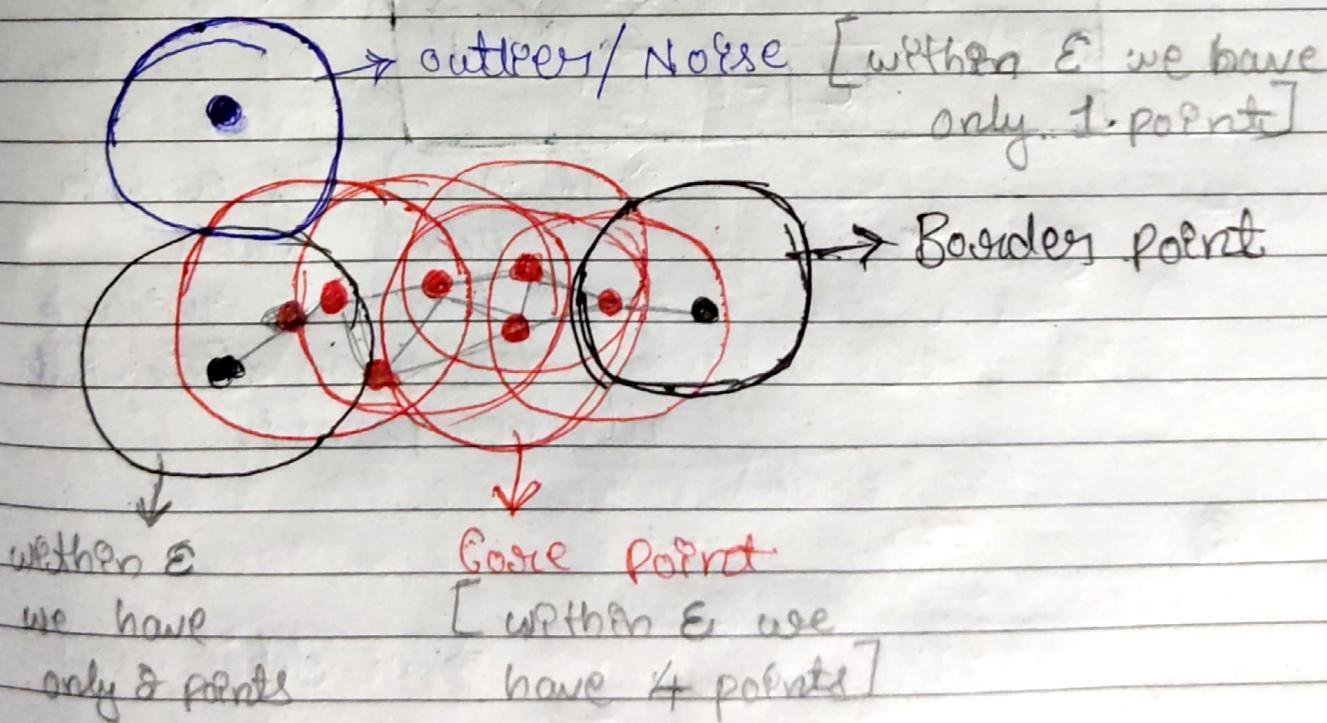
⑤ No of points within this  $\epsilon$  ~~should~~ will be less than min point.  
 $(\leq \text{min point})$

## \* outlier / Noise

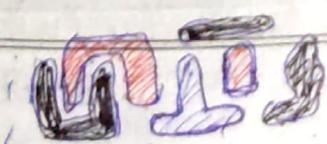
(\*) No. of points within  $E$  is 1.



\* In DBSCAN based on density in different regions we can group our Core point, Border point and outliers.



↓  
outlier / Noise will be removed/handled.



→ can handle (non linear data) + (noise)

## \* Silhouette clustering

How do we validate  $K=4$  flow method.

Steps

①

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



C1 cluster

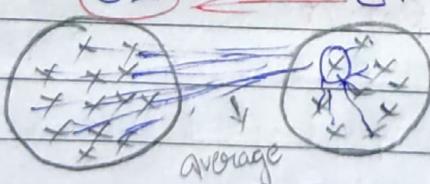
→ if  $a(i)$  is a data point then

except from a selected point

→ we compute average distance ~~between~~ with all other points.

(nearest cluster)

②



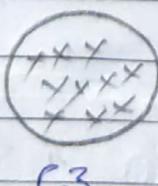
→ we define mean

clustering

to point  $i$  to

some cluster  $C_j$

nearest cluster



C3

→ we calculate the average distance to the points of cluster  $C_j$

nearest cluster

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

$\Rightarrow$  number of points in nearest cluster.

w.r.t

$a(i) \Leftarrow b(i)$   $\rightarrow$  if clustering  $\&$  done well

"average distance" with all other points "within" "cluster."

"average distance" with the points of the "nearest cluster."

### ③ Calculate clustering score

$$\delta(i) = \frac{b(i) - a(i)}{\max[a(i), b(i)]}$$

Value range  $\Rightarrow [-1 \text{ to } 1]$

$\nabla$  More near to 1 the better clustering model we have created.

(or)

$$\delta(i) = \begin{cases} 1 - \frac{a(i)}{b(i)}, & \text{if } a(i) < b(i) \\ 0, & \text{if } a(i) = b(i) \\ \frac{b(i)}{a(i)} - 1, & \text{if } a(i) > b(i) \end{cases}$$

$\approx +1$        $\approx -1$

Again revise  
step

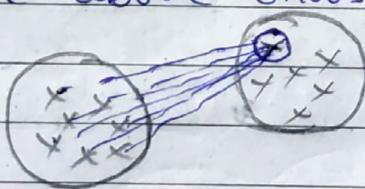
Revise

- ① Calculate average distance within cluster from one point.



$a(i)$

- ② Calculate average distance of all points ~~near~~ of nearest cluster from same above chosen points in  $a(i)$ .



$b(i)$

- ③ Calculate silhouette scoring and withdraw output.

Value near to +1  $\Rightarrow$  better the clustering

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

$$② b(i) = \min_{j \neq i} \frac{1}{|C_j|} \sum_{j \in C_j} d(i, j)$$

$$③ \delta(i) = \frac{b(i) - a(i)}{\max[b(i), a(i)]}$$

## \* K Means Implementation

only independent feature

① Create dataset

> from `sklearn.datasets import make_blobs`

↳ binary large objects

generates isotropic gaussian blobs for clustering

2 features

can give average also

>  $x, y = \text{make_blobs}(n\_samples=1000, n\_centers=3, n\_features=2)$

as unsupervised ML  
we won't use

[No of centroid]  
cluster

## \* Scatter plot

$x \rightarrow [[-1, -1], [-1, -2], \dots]$

`plt.scatter(x[:, 0], x[:, 1])`

② Train test split.

## ③ clustering

> from `sklearn.cluster import KMeans`

[  
max distance] To avoid wrong clustering

Parameters

~~n\_clusters, n\_init = "k-means++"~~

No of centroid

\* Elbow method → Find K value  
↓  
wcss ↓

> wcss = [ ]

for K in range(1, 11):

multilize the centroid Kmeans = KMeans(n\_clusters=K, random\_state=0)  
to our data ← Kmeans.fit(x\_train)  
wcss.append(Kmeans.inertia\_)

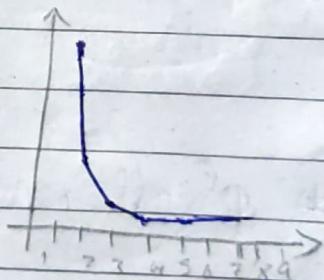
wcss

\* plot elbow curve  
↓

> plt.plot(range(1, 11), wcss)

\* axis label > plt.xlabel("range(1, 11)")

> plt.show()



#### ④ Model

> Kmeans = KMeans(n\_clusters=3,  
random\_state=0)

> y = Kmeans.fit\_predict(x\_train)

# To locate knee value for large data

Knee locator

"k-means++") > ! pip install kneed

> from ~~kneed~~ kneed import KneeLocator

>  $KL = \text{Kneelocator}(\text{range}(1, 10), w=8)$   
Curve = ~~convex~~  
direction = ~~decreasing~~) ~~Convex/Concave~~

>  $KL.$  elbow

> (5) Silhouette score

> from sklearn.metrics import silhouette\_score.

loop  
> for i in

> score = silhouette\_score(X\_train,  
kmeans.labels\_)

## Implementation

\* ~~Hierarchical~~ - Agglomerative clustering

① Import data

② Standardization

> from ~~sklearn.~~ ~~preprocessing~~ import StandardScaler

③ To visualize → PCA  
we do

> from ~~sklearn.~~ decomposition import PCA

> pca = PCA.fit(x\_scaled)

④ Clustering - Agglomerative.

⑤ \* Create dendrogram.

> import scipy.cluster.hierarchy as hc

> hc.dendrogram(hc.linkage(pca\_scaled,  
method='ward'))

fixed

## ⑤ ~~Ag~~-clustering

> from sklearn.cluster. Import AgglomerativeClustering

> cluster = AgglomerativeClustering(n\_clusters=2, affinity='euclidean', linkage='ward')

> cluster.fit(pca\_scaled)

cluster → .cluster\_labels\_

plt.scatter(pca\_scaled[:, 0], —, 1, c=cluster.labels\_)

## ⑥ Silhouette Score.

## \* Implementation $\rightarrow$ DBSCAN

### ① Create dataset.

Non-linear dataset

> from sklearn.datasets import  
make\_moons

### ② Make two interleaving half circles.

>  $x, y = \text{make_moons}(n\_samples=250,$   
 $\text{noise}=0.1)$

### ③ plot

> plt.scatter(x[:, 0], x[:, 1])

### ④ Standardization.

### ⑤ Clustering

> dbScan = DBSCAN(eps=0.3)

> dbScan.fit(x\_scaled)

> dbScan.labels\_

### ⑥ Plot

( $\rightarrow c = \text{dbScan.labels}_-$ )