

Clustering

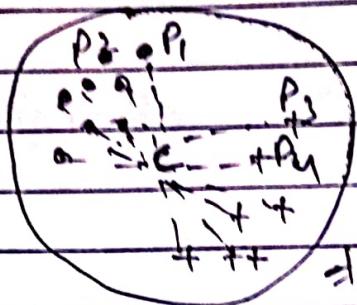
1 Clustering - It is the task of dividing the unlabeled data or data points into different cluster such that similar data point fall in the same cluster.

Application of Cluster:- Social network analysis, medical imaging, image Segmentation and anomaly detection

WCSS \rightarrow Within intra-Cluster Submission of Square

$= \sqrt{(n_1 - p_1)^2 + (n_2 - p_1)^2}$
WCSS is to find the no. of cluster in the fit

①

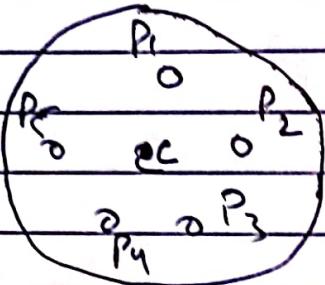


If $K = 1$

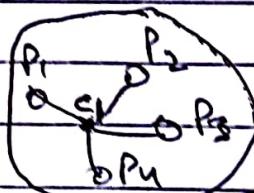
then

SC = Submission of Square

$$SS = (c_1 - p_1)^2 + (c_1 - p_2)^2 + (c_1 - p_3)^2 + (c_1 - p_4)^2 + (c_1 - p_5)^2$$



②



$$\begin{aligned} SS_1 \text{ for cluster } 1 &\Rightarrow (c_1 - p_1)^2 + (c_1 - p_2)^2 \\ &+ (c_1 - p_3)^2 + (c_1 - p_4)^2 \end{aligned}$$

$$\begin{aligned} SS_2 \text{ for } C_2 &\Rightarrow (c_2 - p_5)^2 + (c_3 - p_5)^2 \\ &+ (c_4 - p_5)^2 \end{aligned}$$

$$\text{Inertia} = \text{WCSS} \Rightarrow SS_1 + SS_2$$

$$\text{WCSS} = \sum \sum (c_i - p_i)^2$$

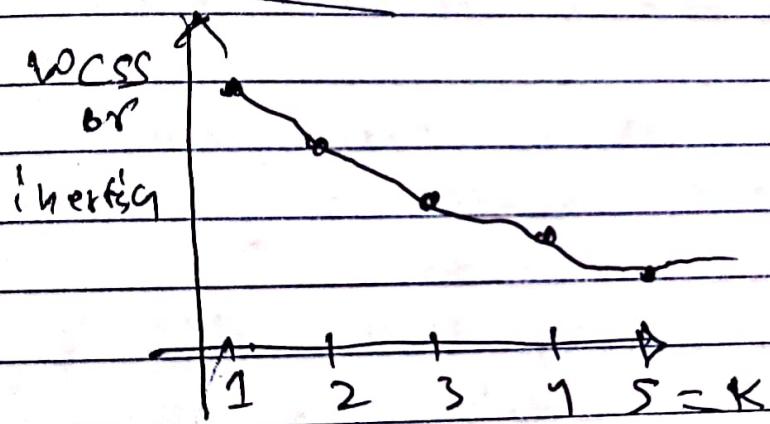
Σ = One Submission for individual cluster

Σ = One Submission for itself.

Q How and when WCSS will be helpful.

Ans If we have $K=3$ then we take $K=4$ then $K=n$. In case of $K=n$ then the WCSS will be 0. So when we have a 1 cluster the distance is going to be maximum. Then if $K=1$ then the distance is going to decrease than more for $K=2$

Elbow graph



Note: for final model we find WCSS for every value of K and we will plot it with the help of Elbow graph and we will check from where there is a no dispersion of K in data and we will select K for a particular model.

- There is a problem to calculate the centroid in the data set in K-Means so to overcome this problem we have $K\text{-Median}++$ to overcome this problem.

Note :- whenever we try to build a cluster we either follow

- Agglomerative
- Divisive (K-means)

Hierarchical clustering

1) Hierarchical clustering - It is a popular method for grouping objects. It creates groups so that objects within a group are similar to each other.

- It is separating data into groups based on some measure of similarity.

Q) What exactly hierarchical do -

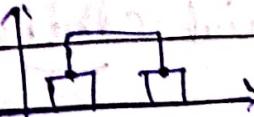
Ans) 1) Every data make one cluster and lie on X-axis of 2D graph.

2) After that point data point find similar or near data point to make a dendrogram

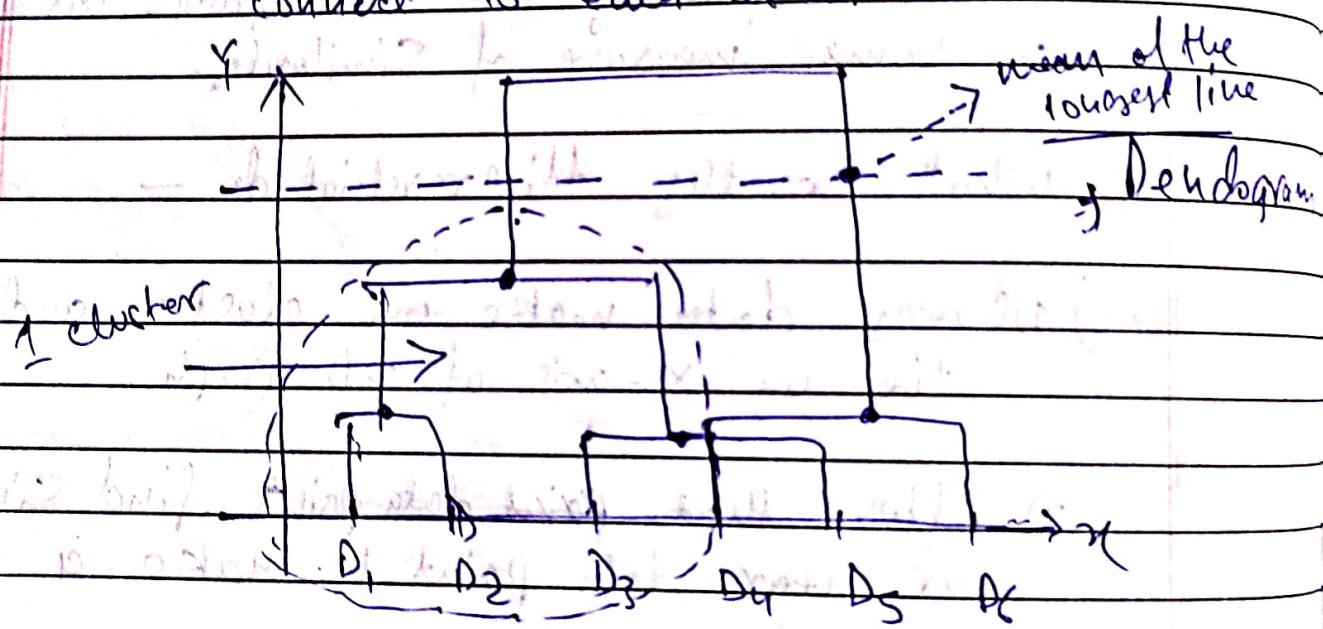
3) It make vertical line based on euclidean distances and join to data point with similarities.

4) After connecting with each data point and make a graph there will be every centroid for every connecting data point.

5) and based on the less distance or the similarities we will connect every centroid.



6) we will fill every point attach or connect to each other.



7) Now we will check which vertical distance is more and after choosing the more vertical line we will find a mean of that long vertical line and we will draw a horizontal line.

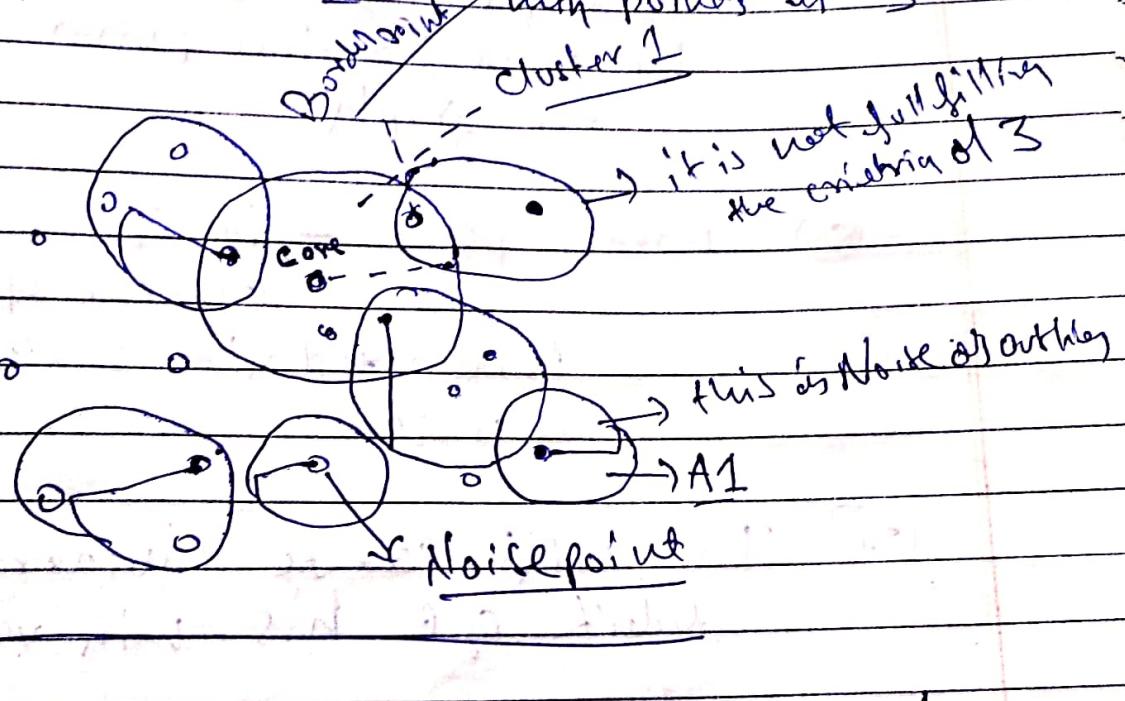
8) And the no. of line of point if connected it will created that much cluster for that dataset.

⑨ This will known as Dendrogram approach or we can say Agglomerative clustering,

DBSCAN Clustering

- 1 DBSCAN :- Density Based Spatial Clustering of Application with Noise.
- 2 DBSCAN is create cluster by -
 - i) Epsilon :- representing a distance / radius.
 - ii) Min-point :- Minimum no. of point requires to make a cluster.
 - iii) Core points :- Satisfy the minimum point.
 - iv) Border-point :- Point is a part of some circle but not forms its own cluster.
 - v) Noise :- outlier and which is not a part of core point as well as a cluster.

3. How to make cluster or group in DBSCAN
if Epsilon is 1 & min points is 3



Ans

- 1) We will choose a random data in the graph and draw a circle whose radius will be 1 and minimum data point will be 3 or greater than 3 and the random point is taken is basically a core point

- 2) Now i will take all the point of the circle of the cluster 1 will be form.

- 3) Now we will do the same thing with the data just inside the circle or cluster 1 but it should have Epsilon 1 & min point 3

- (6) If every point doesn't a part of any core point or cluster - then it will be Noise / outlier points for $\text{in} - A_1$ in the graph.
- (5) Border point - it is a point of cluster of the core point but didn't make it noisy cluster.
- (4) It will create cluster till it didn't get his min value.

Metrics of Cluster

Date _____
Page _____

- 1 How to evaluate the cluster - (Performance Matrix)
 - i) Jaccard Coefficient
 - ii) Entropy
 - iii) Purity
 - iv) Silhouette Coefficient
 - v) Separation
 - vi) Rand Index

2 Rand Index - It is one of the metrics by which we will able to observe or monitor that how good our cluster creation is.

Rand Index = Total Agree

Total Agree +

~~SS~~ → If means that both the point belong to a same cluster for both our algorithm and Ground Truth (ground truth mean that base which i will able to compare)

~~DD~~ → Both point don't belong to the same cluster but both our algorithm and Ground Truth

~~SD~~ →

$S^1 \quad D^0 \rightarrow$ Prediction clustering

	S^1	D^0	
Ground Truth	SS	SD	DS
Prediction	SD	DD	DS

So, Rand Index will be $\frac{SS + DS}{SS + DS + SD + DS}$

* Ground Truth \rightarrow is our Actual Observation or equivalent to what it prepared by us.

3 Jaccard Coefficient: - A total dataset where our cluster is able to

Group or Separate from a ground truth

Some values we are able to get

(Intersection) of cluster = $\frac{SS}{SS + SD + DS}$

$$ss + sd + ds$$

4 Entropy = $-P_i \log(P_i)$ is always try to tell that how random my distribution are.

$$= - \sum P_i \log(P_i)$$

Entropy tell us the randomness in the cluster and we always try to or looking for less randomness into our data distribution. So we can compare in our data which clusters have less randomness.

5 Purity :- Total percentage of data set or a data point clustered directly.

$$\text{Purity of Cluster}_i = P_i = \max(P_{ij})$$

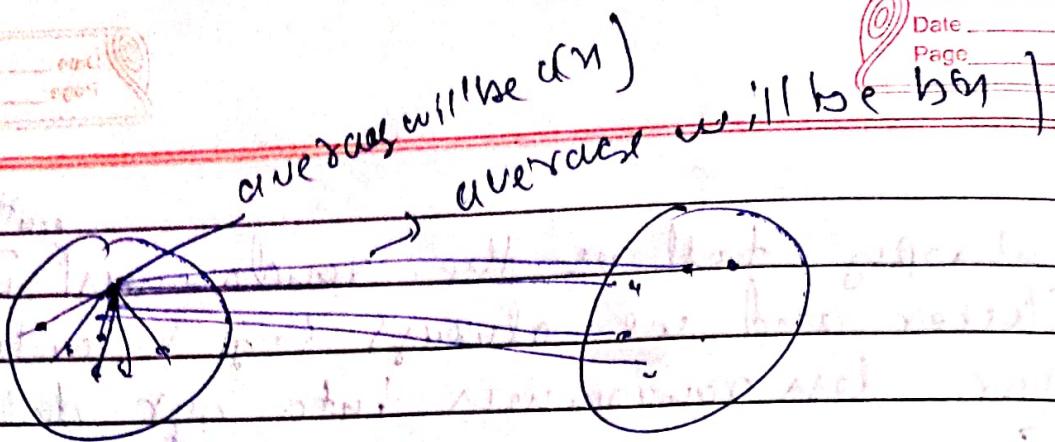
$$\text{Purity of whole cluster} = P(c) = \left(\sum m_i P_i \right) / n$$

6 Silhouette - to find accuracy of the clusters

$$s(n) = \frac{b(n) - a(n)}{\max\{a(n), b(n)\}}$$

$a(n)$ = average distance of n from all the other points in the same cluster.

$b(n)$ = average distance of n from all the points in another cluster.



Note: - In other metric we have to prepare a Ground Truth then only we can able to calculate accuracies or stability of model in terms of clustering.

In case of Silhouette it not required because it following distance based approach to evaluate whether my cluster is good or bad.

$$\text{Silhouette} = \frac{1 - \sum s(n)}{N(n)}$$

PCA

Principal Component Analysis

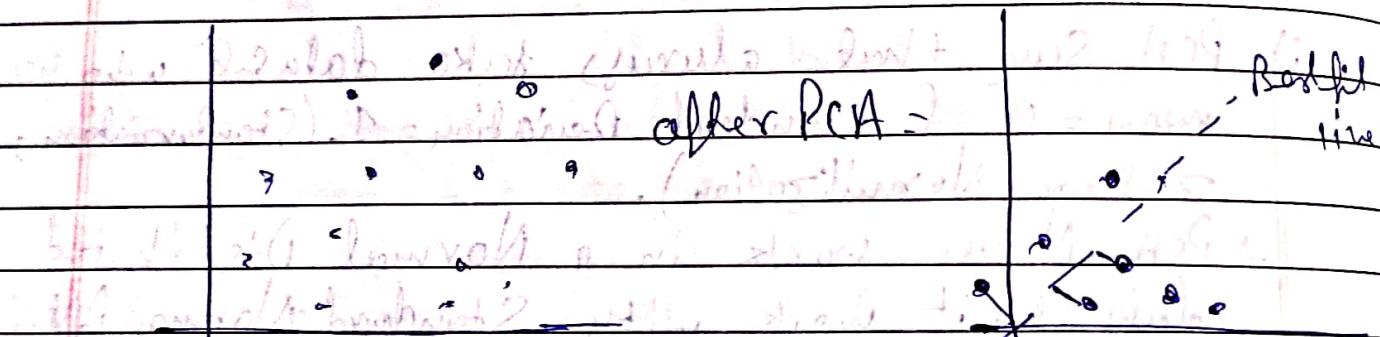
- PCA is not a approach where we can train a model, It is a data pre-processing stage.
- PCA help to reduce a dimension.

1 Step of PCA -

- i) PCA says that always take dataset who has mean = 0 & Standard Deviation = 1. (Standardization, Z-Score Normalization).
- PCA Never work in a Normal Distributed dataset, it work with Standard Normal Distr.
- ii) Draw a best fit line or a Straight line as we do in Linear Regression.
- iii) The line draw is PC 1 which is a Straight line and after that PC 2 will be create which is Perpendicular line of PC 1.
If PC 3 created that it will be perpendicular line of PC 1 & PC 2 and it is z-axis.

Note - In PCA we don't eliminate the columns we reduces the columns.

Note :- PC 1 (Principal Component 1) is the first and most contributing component that other PC 2 & PC 3 ... PC n will have less contribution. Because PC 1 is kind of a straight line which has been drawn base out of the relation. And rest PC line are the perpendicular on top of the first one.

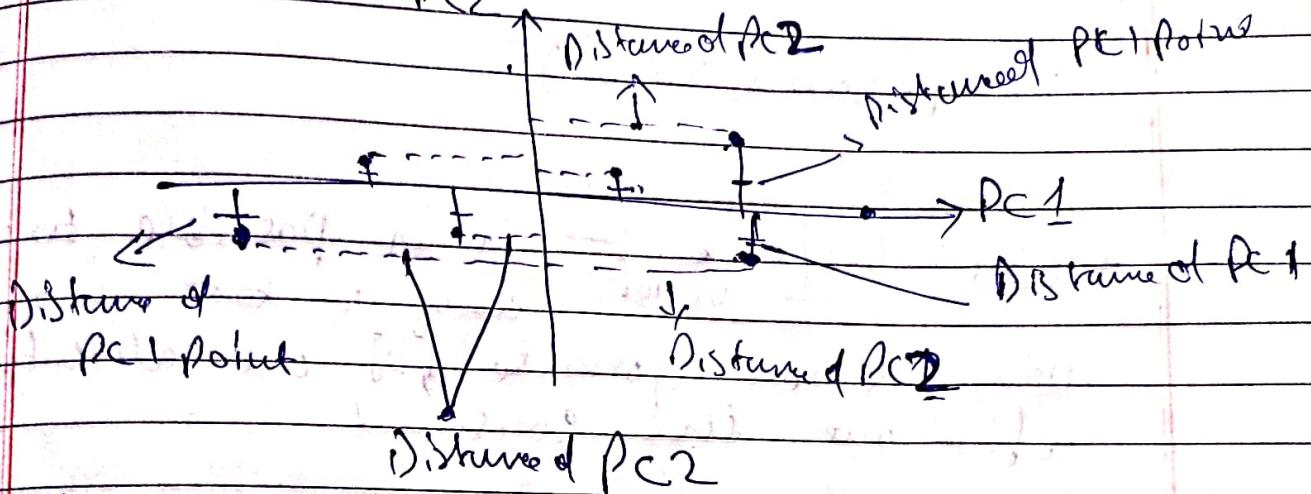


2 How to know that how many PC line to create

Ans EVR (Explained Variance Ratio) :- One PC has the capacity that One single line or P. how much relation this one single line is able to explain is called (EVR).

EVR of PC1 = Distance of PC1 point

Distance of PC1 + Distance of PC2

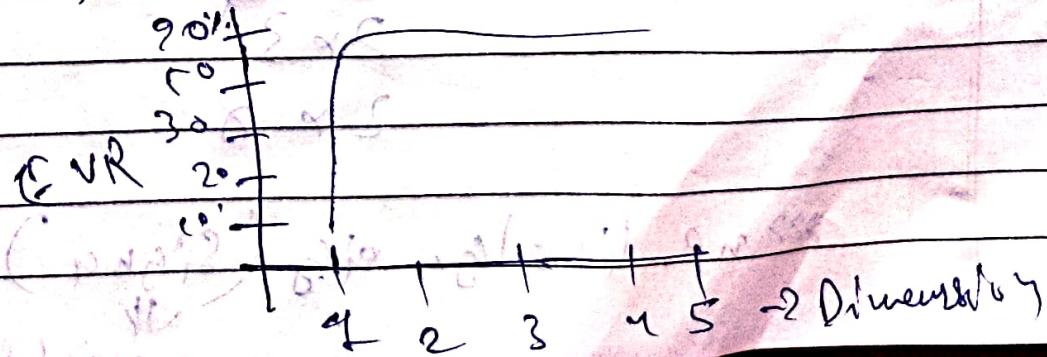


for ex,

$$EVR \text{ of PC1} = \frac{56}{50+5} = 0.91$$

∴ It means ~~mean~~ PC1 line contain 91% of data information.

3 Draw Screen plot to know to know how many dimension to take with with the help of EVR



$$\text{Covariance} = \frac{1}{m-1} \sum_{(c_1, c_2)} (c_1 - \bar{c}_1)(c_2 - \bar{c}_2)$$

?? n.p. cov (c1, c2)

n. eigen value/eigen vector → it is possible to represent any vector in any other vector formate, if we able to find out the constant.

$$T(y = dv)$$

for ex

$$m_1 \quad m_2$$

$$22 \quad 1$$

$$y$$

$$8$$

NA and will P. 3 with 6 rows of 8 where 16

$$14$$

is found of m_1 only need one

$$2 \times 2$$

$$2 \times 3$$

$$2 \times 8$$

?? n.p. linalg.eig (signs) or
variable of cov.

↳ mathematically, $\text{eig}(\text{np.cov}(a_0, a_1))$ - will give

5) \Rightarrow Eigen value always tell the one component $\xrightarrow{\text{PC}}$
 ↳ the other PC.

→ The highest eigen value will be contained
 the PC 1.

Code of PCA -

```

>>> from sklearn.decomposition import PCA
>>> pca = PCA()
>>> principal_component = pca.fit_transform(data)
>>> pca.explained_variance_ratio_
>>> plt.figure()
>>> plt.plot(np.cumsum(pca.explained_variance_ratio_))

```

Note :- here data should be standarization,

PCA - Code

- i) Use Standard Scaling for mean = 0 & $S.D = 1$
- ii) Find out the Covariance Metrics.
- iii) Find the eigenvector & Eigen Value.
- iv) Draw Screen plot and check n-component or no. of PC.

>>> i) Using glass data (iris dataset)

ii) dropping the label column

Code 1: first do Standard Scaling to the data

```
>>> from sklearn.preprocessing import StandardScaler
```

```
>>> scaler = StandardScaler()
```

```
>>> scaler.fit_transform(data)
```

```
>>> scaled_data = scaler.fit_transform(data)
```

Code 2: To change in Data Frame

```
>>> df = pd.DataFrame(scaled_data, columns = data.columns)
```

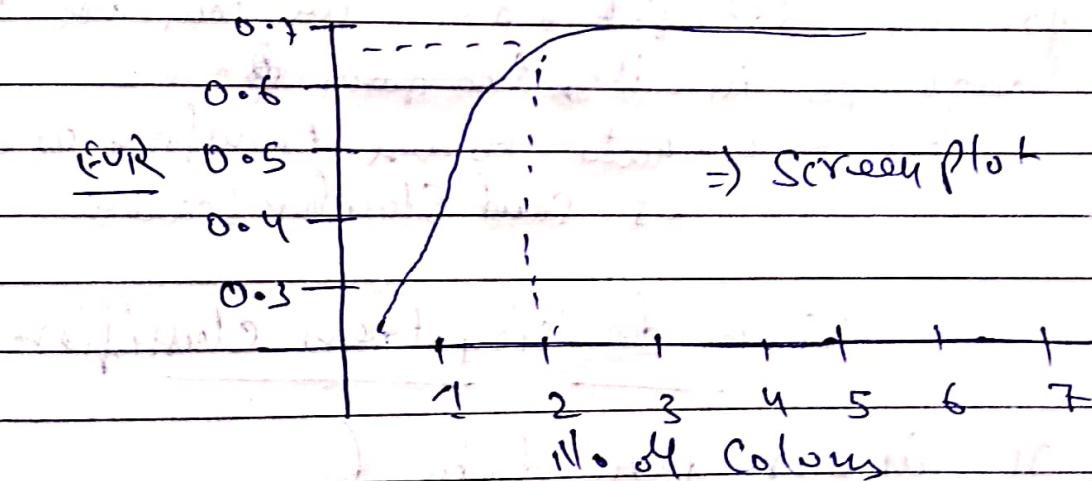
Code 3 >>> df.describe() (# to know the whether mean & S.D = 1 or not)

```
code 4 >>> from sklearn.decomposition import PCA
>>> pca = PCA()
>>> pca.fit_transform(df)
```

Now we have to take the EVR or the no. of PC to take in the PCA with the help of Screen plot with the help EVR

```
code 5 >>> plt.figure()
```

```
>>> plt.plot(np.cumsum(pca.explained_variance_ratio_))
>>> plt.xlabel("No. of columns")
>>> plt.ylabel("EVR")
```



This tells the No. of Principal Component (PC) with the presence of data file.

No. 2 lies - 70% of data So we will take n-component = 2

Code 6 \Rightarrow `pca1 = PCA(n_components=2)`

Code 7 \Rightarrow `new_data = pca1.fit_transform(df)`

Code 8 \Rightarrow ~~to create new data into Data frame~~

\Rightarrow `y = pd.DataFrame(new_data, columns=['PC1', 'PC2'])`

Now we have used PCA to reduce our Dimension or columns for better performance for our data set. Now we can use any algorithm in it's ~~framing~~.

Note - we have to use transform as we have used Standard Scaler.

Using decision tree classifier

\Rightarrow `y` is our PCA data

\Rightarrow ~~y~~ is our $y = \text{data}.\text{Class}$ our label or outcome

Code 9 \Rightarrow `from sklearn.tree import DecisionTreeClassifier`
 \Rightarrow `dt_model = DecisionTreeClassifier()`
 \Rightarrow `dt_model.fit(y, y)`

Note :- Now we can use for predict but the problem is that after PCA our columns has change to 2 and in the original data set there 9 columns. So, unknown person will give 9 columns. And we also uses StandardScaler.

So we have to used ~~PCA~~ of PCA1.transform to reduce the columns and StandardScaler.transform to change the value and ~~vector~~ for the prediction purpose. Showen in the given code in 10

Code 10 :- dtc model. predict (pca1.transform(StandardScaler.
 ([-0.12, 0.34, 0.44, ..., 97])))

Clustering Code

>> Using Mall dataset

KMeans

code1 >> from sklearn.cluster import KMeans

code2 >> wcss = []

>> for i in range(1, 15):

>> [kmean = KMeans(n_clusters = i, init = 'k-means++', random_state = 30)

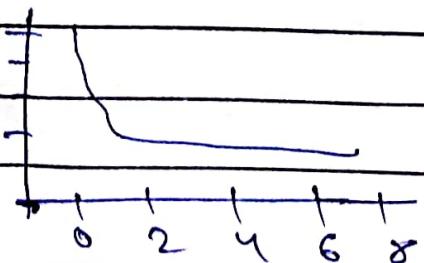
>> kmean.fit(X)

>> wcss.append(kmean.inertia_)

Now we will do elbow graph with the help of WCSS so we could know the no. of cluster should we take.

WCSS = inertia = help to know or find the no. of cluster that we should use

Code 3 >> plt.plot(range(1, 15), wcss)



code 4 >>> Kmean1 = KMeans(n_clusters=5, init="k-means++",
random_state=30)
>>> Kmean1.fit_predict(u)

code 5 >>> we will u a column of predict in our dataset

>>> np["cluster number"] = Kmean1.fit_predict(u)

code 6 >>> Kmean1.predict([25,32])

DBSCAN

>>> from sklearn.cluster import DBSCAN
>>> dbSCAN = DBSCAN(eps=1, min_samples=3)

>>> dbSCAN.fit(u)