

Dt: 3/05/22  
3:30 AM

## \* Un Supervised Machine Learning:

# We don't have any output columns,  
Variable / Feature

↓ (We Group them)

### Clustering :-

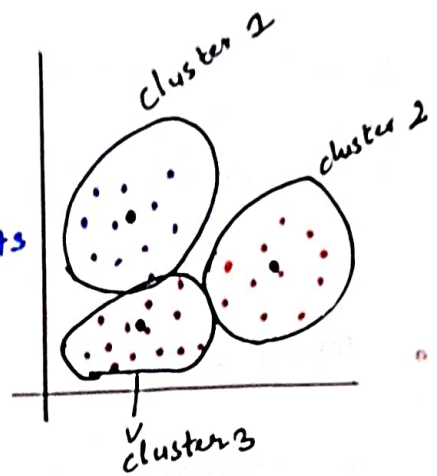
# Grouping of similar objects

1. K-means Clustering
2. Hierarchical clustering
3. DB Scan clustering

# Unlabel data

1. What is a cluster?

A:- A group of objects that are <sup>records.</sup> "similar to other objects" in the cluster, and dissimilar to data points in other clusters.



# groups :- segments

2. Clustering applications?

• Retail / marketing :

- \* Identifying buying patterns of Customers
- \* Recommending new books (or) movies to new Customers.

• Banking :

- \* Fraud detection in credit card use.
- \* Identifying clusters of customers (eg. loyal)

• Insurance :

- \* Fraud detection in claims analysis
- \* Insurance risk of customers.

- Publication :

- \* auto-categorizing news based on their content
- \* Recommending similar news articles

- Medicine :

- \* characterizing patient behaviour

- Biology :

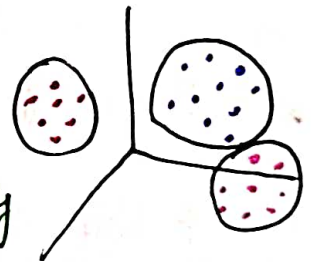
- \* clustering genetic markers to identify family ties.

### 3. Clustering algorithms

- \* Partitioned-based clustering

- Relatively Efficient

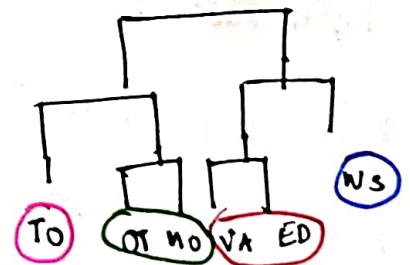
Eg:- k-means, k-median, Fuzzy C-means.



- \* Hierarchical clustering

- Produces trees of clusters

Eg:- Agglomerative, Divisive





## \* Density-based clustering

- Produces arbitrary shaped clusters

eg:- DBSCAN



## 1. K-Means.

### \* Partitioning clustering

- \* K-means divides the data into non-overlapping subsets (clusters) without any cluster internal structure

- \* records within a cluster are very similar

- \* records across different clusters are very different

Ex:- covid zones

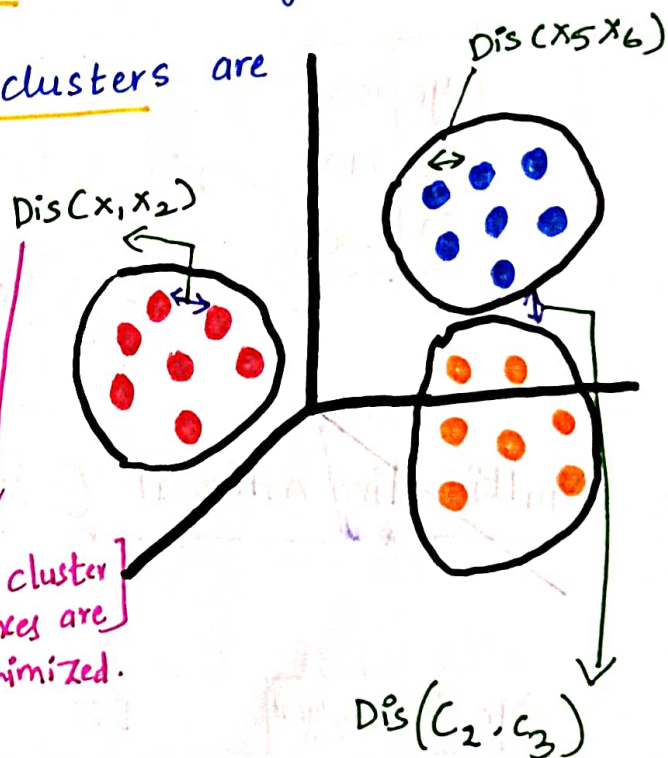
\* Determine The similarity (or) dissimilarity

- intra - within
- inter - outside

⇒ Aim:-

- \* Between two clusters, The distance should be maximum - Better The model

- \* within cluster, The data point should be closely to each other - Better The model



## \* One-dimensionality (similarity / distance)

	Age
1	54
2	50

We, calculate Distance.

$$\sqrt{(54-50)^2} = 4.$$

$$\text{Dis}(x, x_2) = \sqrt{\sum_{i=0}^n (x_{1i} - x_{2i})^2}$$

distance formula.

# 3 dimensional - 3 column:

$$\sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2} \text{ (Euclidean)}$$

## \* two-dimensionality (similarity / distance)

Age	Income
54	190
50	200

$$\text{Dis}(x, x_2) = \sqrt{(x_1 - x_2)^2}$$

$$= \sqrt{(54-50)^2 + (190-200)^2}$$

$$= 10.77$$

## \* Multi-dimensional (similarity / distance)

Age	Income	Educational
54	190	3
50	200	8

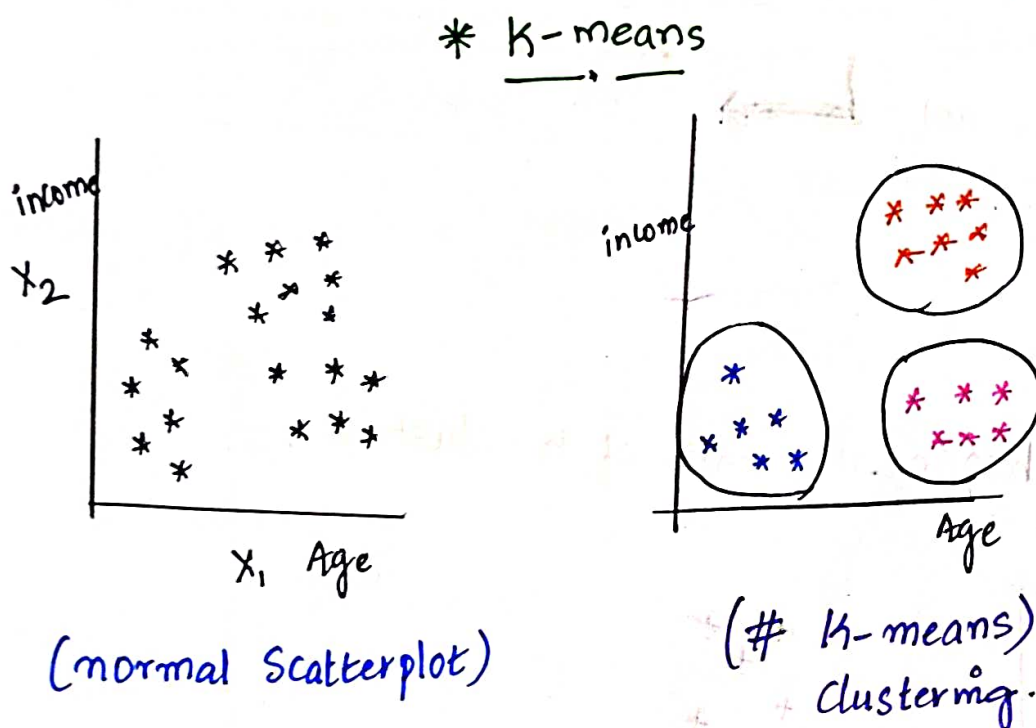
$$\text{Dis}(x, x_2) = \sqrt{(x_{1i} - x_{2i})^2}$$

$$= \sqrt{(54-50)^2 + (190-200)^2 + (3-8)^2}$$

$$= 11.87$$

Q:- what is the difference b/w "K means" & KNN?

A: Both use Distance formulas to group them.  
but, K-means is clustering, KNN- classification.  
Non supervised (only group)  
supervised (O/P) predict



\* Steps to K-means :

Step: 1 (choose The no. of "K" of clusters)

Step: 2 (select at random "K" points, The Centroids)  
(not necessarily from Your dataset)

Step: 3 (assign Each data point to closest Centroid)  
→ (that forms "K" clusters)



↓  
Step: 4 (compute & place new Centroid of Each cluster)

↻  
Step: 5 (Re assign Each data point to the new closet Centroid. if any reassignment took place. go to Step 4. Otherwise go to FIN)

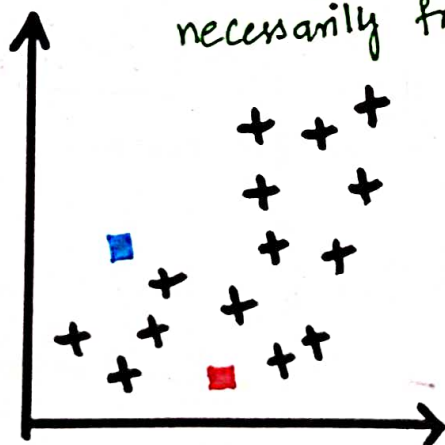
\* K-value is not Fixed...

→ Your Model is Ready

⇒ Step: 1 :- choose the no. of K-clusters : 2

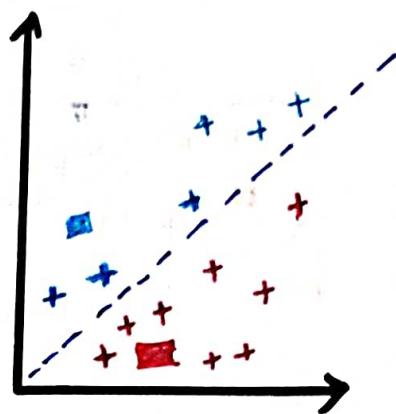
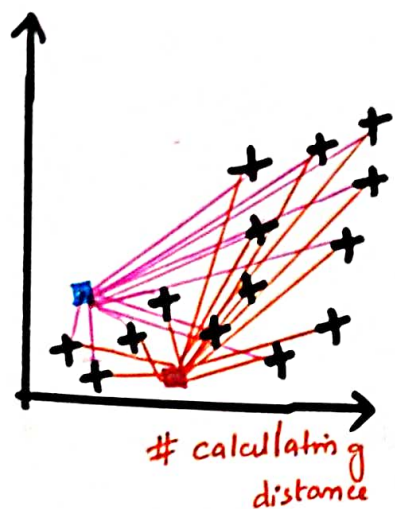


⇒ Step: 2 select at random K-points, The Centroids (not necessarily from your dataset)

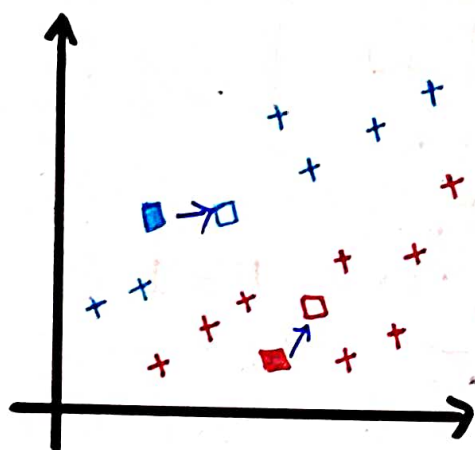


Centroid.  
Based on  
#assumptions

Step:-3 (assign Each data point to The closest Centroid that forms K-clusters.)

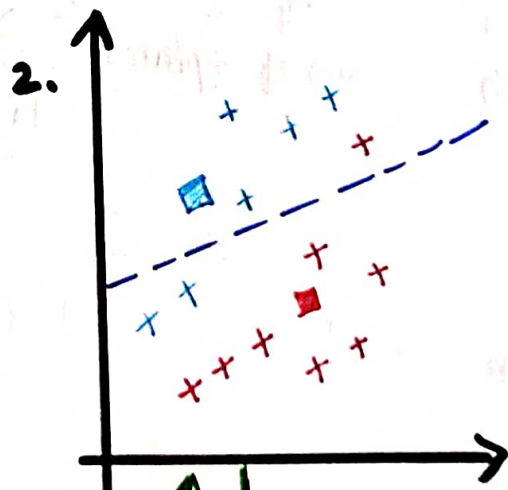
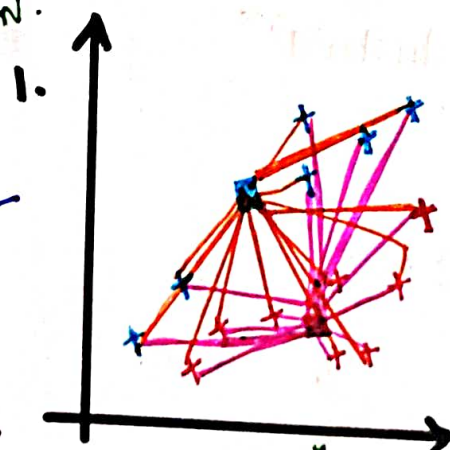
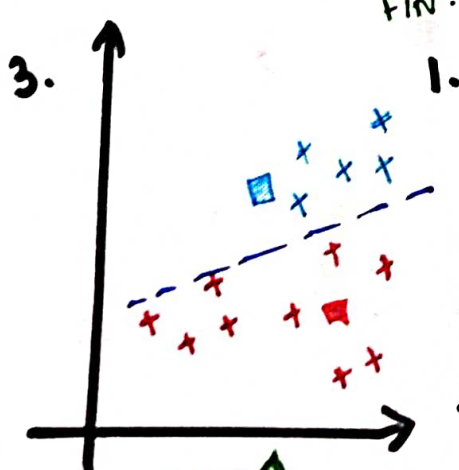


Step:4 Compute and replace The new Centroid of Each other.



Calculate the original Centroid  
 $Avg = \bar{x} \bar{y}$   
 # replace with avg.

Step:5 Reassign each data point to the new closest Centroid.  
 if any reassignment took place, go to Step 4. Otherwise go to FIN.



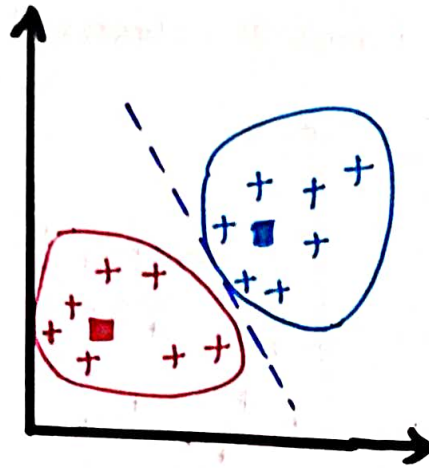
# at least one datapoint

changed also, we have continue process.

# again calculate Centroid, again

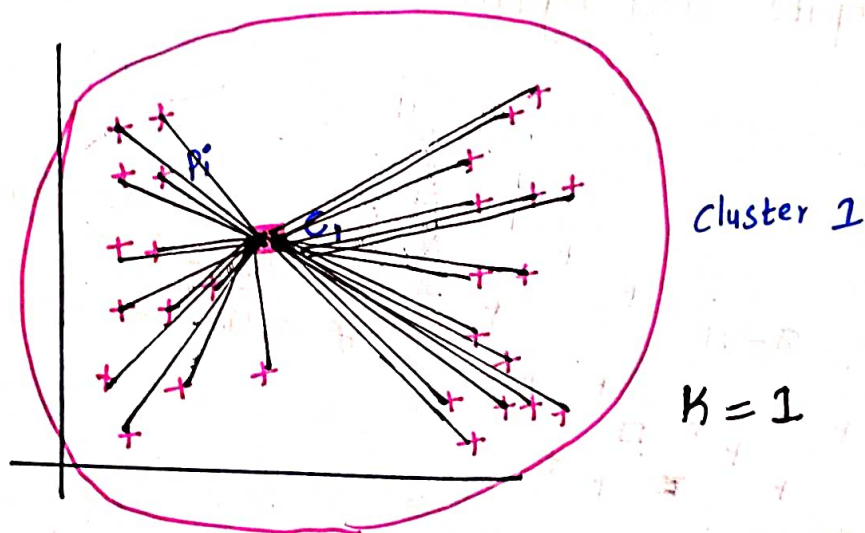


FINAL : Model is Ready.



Que :- How to choose right number of clusters ?

A :-

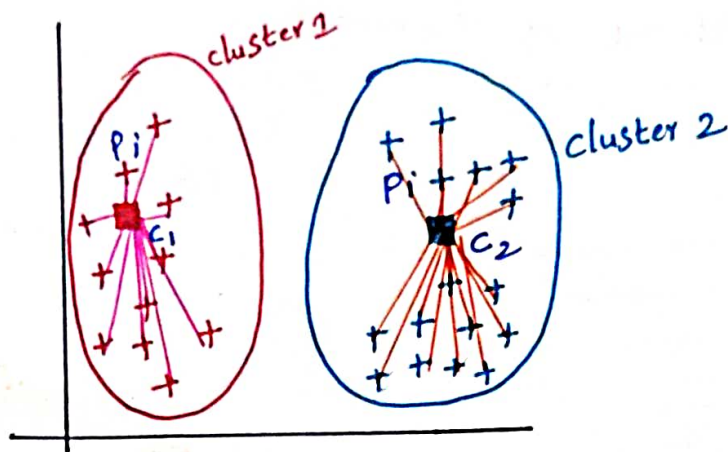


\* WCSS  $\Rightarrow$  (with cluster Sum of Squares)  
 $\Downarrow$   
 with Sum of Squares =  $\sum_{P_i \text{ in cluster } 1} \text{distance}(P_i, C_1)^2$

\* WCSS =  $\text{distance}(P_i, C_1)^2$   
 #  $K=1$

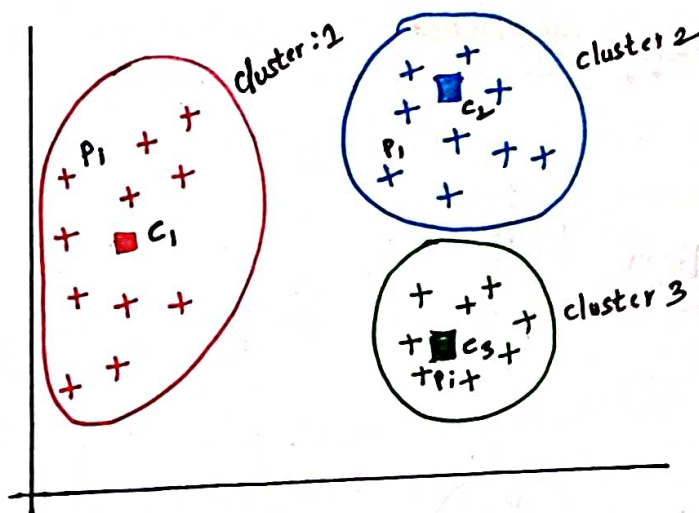
# K = 2

# Calculate SSE



$$* WCSS = \sum_{P_i \text{ in cluster } 1} \text{distance}(P_i, c_1)^2 + \sum_{P_i \text{ in cluster } 2} \text{distance}(P_i, c_2)^2$$

(SSE) + . . . . . (SSE)



# K = 3

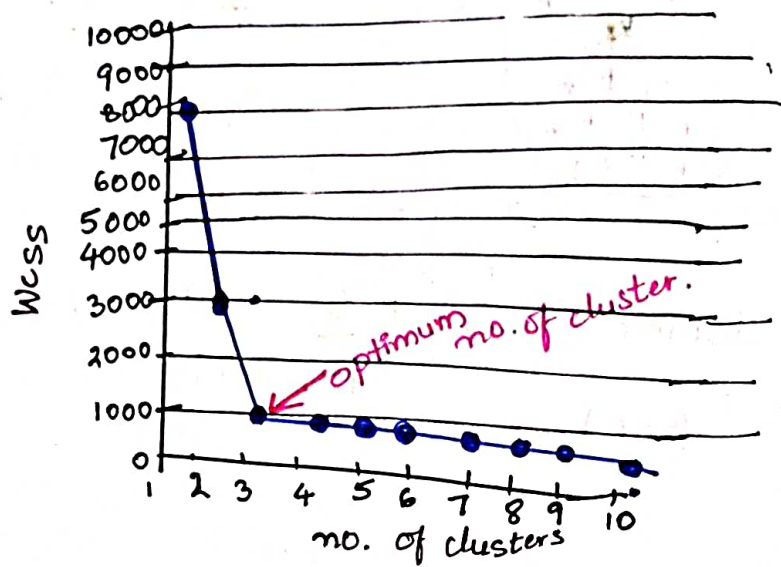
Apply loop.

$$WCSS = \sum_{P_i \text{ in cluster } 1} \text{distance}(P_i, c_1)^2 + \sum_{P_i \text{ in cluster } 2} \text{distance}(P_i, c_2)^2 + \sum_{P_i \text{ in cluster } 3} \text{distance}(P_i, c_3)^2$$

03/05/22

6:00 Am.

## \* choosing The right no. of clusters

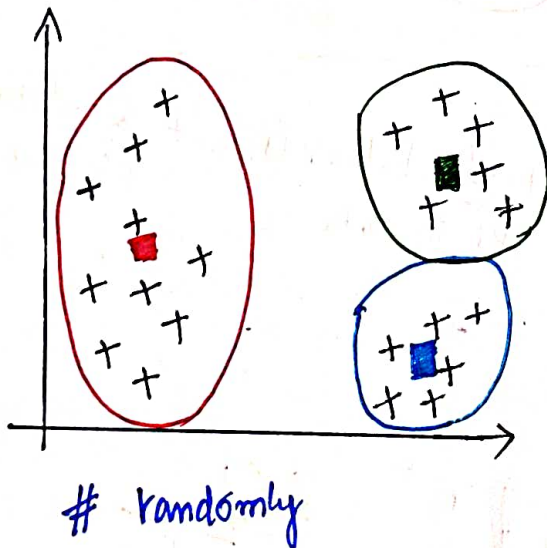


## \* Elbow Technique



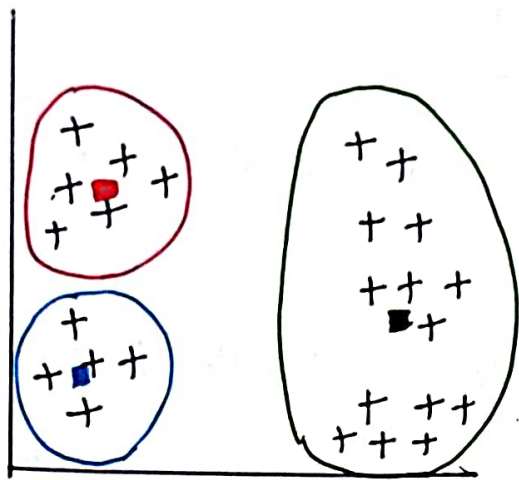
## \* Elbow method

## \* Random initialization trap



Que :- What would happen if we had a bad random initialisation?





# different types of clusters But, Data Set is common.

# Based on random initialization, it is going to change clusters

If you choose wrong initialization, Total problem gets wrong. it is the problem in K-means.

A: Solution for this is K-means ++

⇒ The Random initialisation Trap

One major drawback of K-means clustering is

Random initialisation of Centroids.

The formation of clusters is closely bound by the initial position of a centroid.

The random positioning of Centroids can completely alter clusters and can result in random formation.

\* The Solution is K-means ++

\* K-means ++ is an algorithm that is used to initialise  
 \* The K-means algorithm instead of "arithmetic mean" it  
 Take "Weighted mean" For calculating centroid.

## K-means ++

1. Choose one Centroid uniformly at random from among the data points.
2. For each data point say  $x$ , compute  $D(x)$ , which is distance between  $x$  and nearest Centroid that has already been chosen.
3. Choose one new data point at random as a new Centroid, using weighted probability distribution where a point " $x$ " is chosen with probability proportional to  $D(x)^2$ .
4. Repeat steps 2 and 3 until  $K$  Centres have been chosen.
5. Proceed with Standard K-means clustering.

Ex:-

Age	Edu	Exp
21	B.Tech	0
24	B.Tech	2
22	B.Tech	0
28	M.Tech	3

group

Encoding

grouping same element (or) similar

Multiple feature overall

(calculate distance formula)

CODE :-

```
# df = pd.read_csv("Mall_Customers.csv")
```

```
# df.head()
```

**Out :-**

CustomerID	Genre	Age	Annual income (K\$)	Spending score (1-100)
1	Male	19	15	39
2	Male	21	15	81
3	Female	20	16	6
4	Female	23	16	77
5	Female	21	17	40

```
# df.shape
```

**Out:** (200, 5)

```
# df.info()
```

**Out:** 200 entries, 0 to 199  
column Non null count Dtype

CustomerID

200 non null

int64

Genre

200 non null

object

Age

"

"

Annual income

"

"

Spending score

"

```
# df.isnull().sum()
```

**Out**

0.



X # to understanding, we take only two columns.

#  $x = df.iloc[:, [3, 4]].values$

## MODEL

from sklearn.cluster import KMeans.

# using the elbow method to find the optimal no. of clusters.


#  $WCSS = []$

for  $k$  in range(1, 21):

$kmeans = KMeans(n\_clusters = k, init = "k-means++")$

$kmeans.fit(x)$

$WCSS.append(kmeans.inertia_)$

- # WCSS values is called inertia 

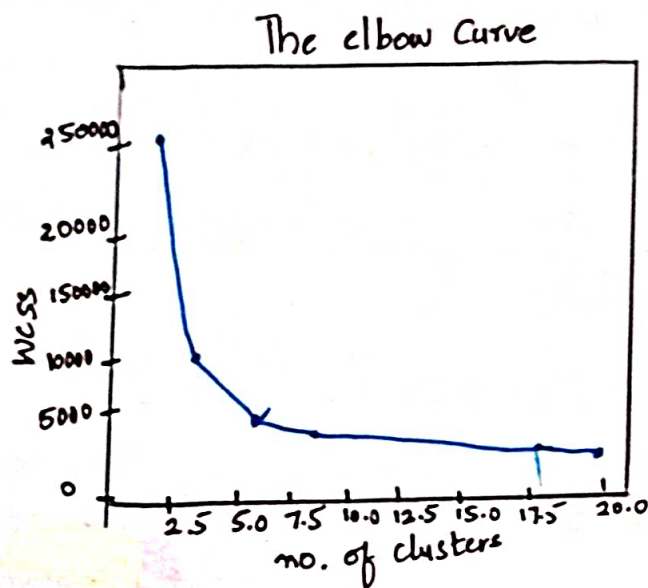
# WCSS

out :-  $[269981.2800 \dots \# k=1$   
 $1811363.5959 \dots \# k=2$   
 $106348.373062,$   
 $73679.789039$

$8504.24137246]$

## # Elbow curve

```
# plt.plot(range(1,21), wcss)
# plt.title('The elbow method')
# plt.xlabel('no. of clusters')
# plt.ylabel('wcss')
# plt.show()
```



fit The K-means Model on the data.

```
# kmeans = KMeans(n_clusters = 5, init = "k-means++",
                  random_state = 42)
```

# predict

```
# y_kmeans = kmeans.fit_predict(x)
```

## Visualising The clusters

# cluster 1

```
# plt.scatter(x[y_kmeans == 0, 0], x[y_kmeans == 0, 1],  
              S=100, c="red", label="cluster 1")
```

# cluster 2

```
# plt.scatter(x[y_kmeans == 1, 0], x[y_kmeans == 1, 1],  
              S=100, c="blue", label="cluster 2")
```

# cluster 3

```
# plt.scatter(x[y_kmeans == 2, 0], x[y_kmeans == 2, 1],  
              S=100, c="green", label="cluster 3")
```

# cluster 4

```
# plt.scatter(x[y_kmeans == 3, 0], x[y_kmeans == 3, 1],  
              S=100, c="cyan", label="cluster 4")
```

# cluster 5

```
# plt.scatter(x[y_kmeans == 4, 0], x[y_kmeans == 4, 1],  
              S=100, c="magenta", label="cluster 5")
```

# Centroid

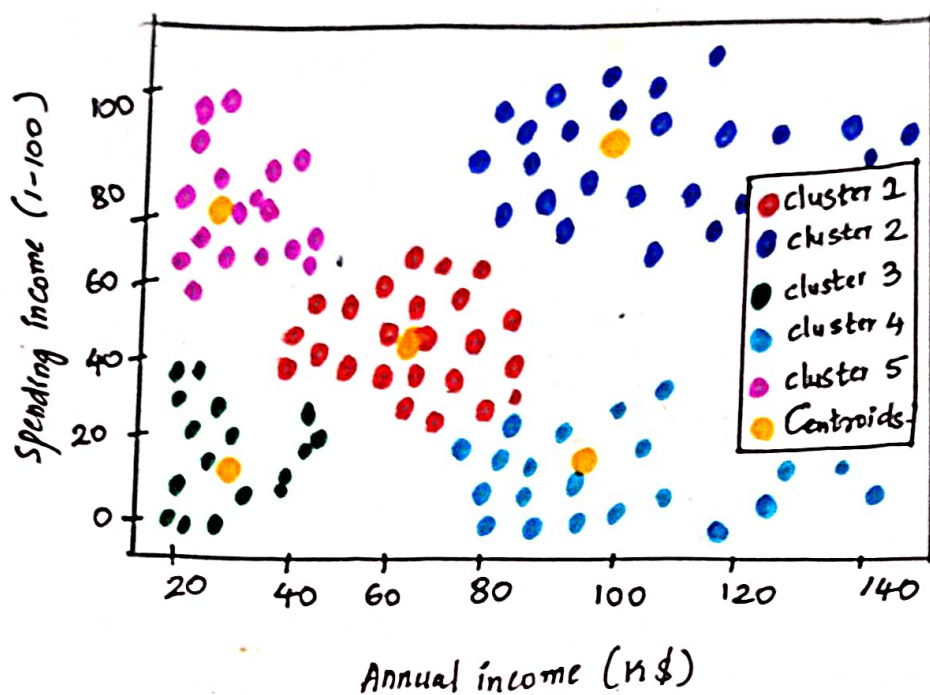
```
# plt.scatter(x[kmeans.cluster_centers_[:, 0],  
              x[kmeans.cluster_centers_[:, 1],  
              S=100, c="yellow", label="centroids")
```

```
# plt.title("clusters of customers")  
# plt.xlabel("Annual income ('k $)")  
# plt.ylabel("Spending score (1-100)")  
# plt.legend()  
# plt.show()
```



Out :-

clusters of customers.



# original data

```
# a = df.iloc[:, 3]
```

```
# b = df.iloc[:, 4]
```

```
# plt.scatter(a, b)
```

Out :

