

① Types of ML:

- i) Supervised ML & ii) Unsupervised ML

i) Supervised :

→ we train the ML using labelled datasets & based on training, the machine predicts the output.

machine train and  $\Rightarrow$  ip & corresponding o/p & then, en predict and  $\Rightarrow$  test cases as per

provide  
ex. 1 datasets  $\Rightarrow$  dog & cat

$\hookrightarrow$  shape, size, color, height etc.

then ip a picture & ask ml to identify & predict o/p.

→ Categories of supervised : —

i) Classification: o/p variable categorical "yes" or "no"

ii) Regression: linear relationship b/w ip & o/p variables

## ① Unsupervised ML :

- ↪ group or categorize the unsorted dataset according to the similarities, pattern, and differences.
- ↪ models are trained with data that is neither labelled and acts on with out supervision.

Ex: basket of fruits & images are unknown to model.

task ML to find patterns & categories  
→ shape, colour, size, height etc.

## ② → Categories of Unsupervised ML :

### i) Clustering:

↪ find inherent groups from the data.

↪ it is way to group objects in to clusters.  
based on similarities, fewer similarities, no-similarity.

↪ algo's used:

~~using~~ k-Means clustering algo

↪ Mean shift algo

↪ principal component analysis.  
etc.

### ii) Association:

↪ find interesting relations among variables within a large dataset.

↪ find dependency of one data item on another data item.

→ Advantages:

- i) used for complicated tasks compared to supervised algs work on unlabelled datasets.
- ii) Unlabelled data getting easier compare to labelled data.

→ disadv:

- i) can be less accurate
- ii) algs not trained with the exact op in prior.

④ K-Means clustering algo:

↳ groups the unlabeled data into different clusters.

↳  $k \rightarrow$  no. of pre-defined clusters.

↳ unlabeled data divided into  $k$ -different clusters in such a way that each dataset belongs only one group. that has similar properties.

↳ it allow us to cluster data into different groups and discover the categories of group.

↳ Centroid-based algo, each cluster associated with a centroid.

aim is minimize the sum of distance b/w data point & their corresponding clusters.

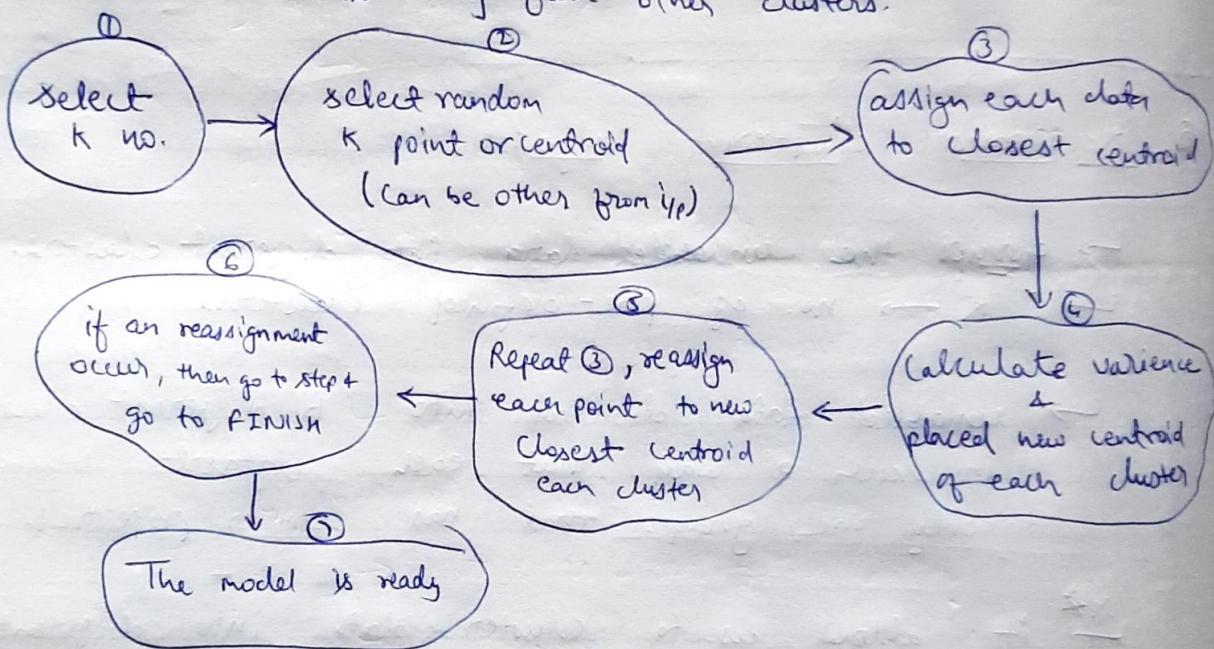
↳  $k$  should be predetermined in this algo.

→ take unlabelled dataset & divided into k-clusters, and repeats process until don't find best cluster.

→ Two tasks perform by k-Means algo:

- (i) Determines the best value for k-center points or centroids by an iterative process.
- (ii) assign each data point to its closest k-center.  
which are near to particular k-center, create a cluster.

Hence, each cluster has datapoints with some commonalities & it is away from other clusters.



0 → choose K (no. of clusters) =

Elbow method:

→ concept WCSS (within cluster sum of square)

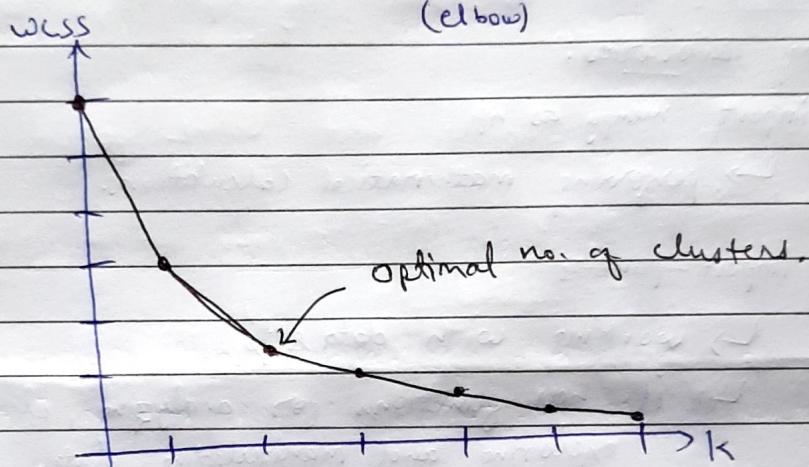
$$WCSS = \sum_{p_i \in C_1} \text{dis}(p_i; c_1)^2 + \sum_{p_i \in C_2} \text{dis}(p_i; c_2)^2 +$$

$$\sum_{p_i \in C_3} \text{dis}(p_i; c_3)^2$$

→ sum of squares of the dis betw each point & its centroid with in cluster 3.  
same for other 2.

To find optimal value of  $k$ .

- ④ i) execute k-means clustering on given dataset for different  $k$  values (range 1-10).  
ii) for each  $k$ , calculate WCSS.  
iii) plot curve WCSS values & no. of clusters.  
iv) sharp point of bend or a point of plot looks like arm. Consider that best value for  $k$



Note: We can choose no. of cluster equal to data points then WCSS becomes zero, that end point of plot.

## -:- Python implementation of K-Means :-

Steps to be:

- i) Data pre-processing
  - (i) finding optimal no. of cluster using elbow method
  - (ii) Training K-Means algo on training dataset
  - (iv) visualizing the clusters.

### ① Data pre-processing:-

#### (i) Importing Libraries:

a) import numpy as np

↳ performing mathematical calculations.

b) import pandas

↳ working with datasets

↳ it has functions for analyzing, learning, exploring and manipulating data.

↳ make readable data.

c) matplotlib.pyplot

↳ serve as visualization utility

↳ graphs, bubble chart, pie chart etc.

#### (ii) Importing data:

we need to deal with huge datasets while analyzing the data, which usually get in csv file.

using: `read_csv( )`

↳ location of file.

## a) information of dataset

↳ `data.info()`

→ Point information about Data frame  
Contains:

no. of columns, no. of rows, Dtype  
and memory uses.

↳ `head()` & `describe`

↳ description of first 5 & last 5 rows  
of data frame.

↳ `isnull()` & `sum()`

↳ checking missing values in data frame  
if any present then give no. of  
vacant or missing values.

## (iii) Extracting independent variables:

↳ we want group based upon spending score  
& spending score data is important for  
mall owner.

↳ we neglect other arguments.

use: `iloc[:, [3, 4]].values`

starting col. last column  
rows      columns  
(empty bcoz we want all rows)

## ② Finding the Optimal number of clusters using elbow method:

- ↳ use the KMeans class of sklearn.cluster library to form clusters
- ↳ we store wcss value in a list for 1-10 clusters.  
then we draw the plot by plotting wcss vs clusters, e
- ↳ sharp bend consider best value for k (no. of clusters).

⑦

```
KMeans(n_clusters=i, init='K-Means++', random_state=42)  
    ↳ no. of cluster (default is 8)  
kmeans.fit(X)  
list.append(kMeans.inertia_)
```

init: { 'K-means++' }

- ↳ Select initial cluster centroids using sampling based on an empirical probability distribution of the points contribution to overall inertia.
- ↳ Making several trials at each sampling step and choosing best centroid.

random state:

- ↳ Determines random no. generation for centroid initialization. use an int to make the randomness.

inertia\_ :

- ↳ Sum of squared distances of samples to their closest cluster center.

Expt. No. \_\_\_\_\_

fit () :

↳ Compute k-mean clustering.

## (3) Training k-mean algo on data:

`kMeans(n_clusters=5, init='k-means++', random_state=42)``y_predict = kmeans.fit_predict(data)`

compute cluster centers & predict  
cluster index for each sample.

## (4) Visualizing the clusters:

↳ To visualize the clusters will use scatter plot using scatter() function of matplotlib

`Scatter(x-axis, y-axis, marker="d", c=" ", label="")`

design for points      colour      label in particular cluster.

X-axis:

↳ we have different values for different clusters

at well →`x[y_predict == 0, 0]`

(that's for  
1<sup>st</sup> cluster)

→ first column from dataset

Y-axis:`x[y_predict == 0, 1]`

for 1<sup>st</sup>  
cluster.

→ 2<sup>nd</sup> column from dataset

↳ coordinates of cluster center

↳ scatter (kMeans.cluster\_centers\_[ :, 0 ], kMeans.cluster\_centers\_[ :, 1 ],  
          s= " ", c= " ", label = "centroids")

↳ legend():  
    ↳ place legend on the axes.

### ⑤ Result:

↳ different clusters with different colour.

↳ five different clusters

- i) rich , spending low
- ii) rich , spending high
- iii) medium, spending medium
- iv) poor , spending low
- v) poor , -- high.

### ⑥ printing data clusterwise

↳ we add new column in dataset then  
give condition for different cluster &  
print for all clusters.