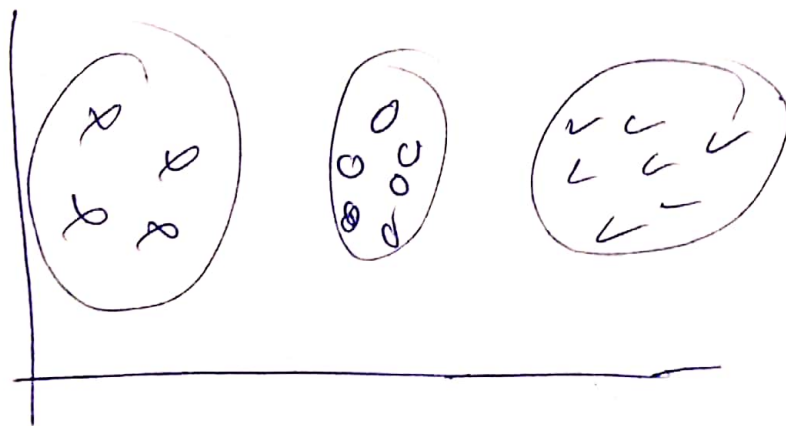# K - Means - Theory

## What is Clustering?

Unsupervised learning in which data is organized into distinct groups, based on similarity they exhibit.



After organizing data points into distinct clusters, we observe similarities within the clusters.

## Why Clustering?

↳ Clustering helps organize unknown observations into groups based on similarity.

↳ We do Clustering when there is no ~~labeled~~ labeled Variable (dependant Column) in data Set

It helps us to find _patterns_ . Widely used for initial analysis of data.

## Applications

* Segmenting Customers based on purchasing behaviour . (:- Customer Segmentation)

* Segmenting medicines based on Constituents, reactions, function & results (:- Medicine)

* In _Psychology_ .

## Why K-means Clustering?

↳ K-means Clustering performs better than Hierarchial Clustering ~~discussed in~~ in large data set.

↳ Because as data size increases Computational time for ~~hi~~ Hierarchial Clustering increases.

## What is K-means Clustering?

↳ An _unsupervised_ learning technique.

↳ In which data is organized into distinct _groups_ having Centroids (Mean values).

↳ K denotes number of clusters or groups.

**Step 1 :-** Choose the value of k - where k indicates number of clusters (we'll discuss later how to choose optimal k)

**Step 2 :-** Initialize (mean) & or (Centroid) of each cluster taking random values.

**Step 3 :-** Calculate [Euclidean distance] of each attribute for each observation from each Centroid.

**Step 4 :-** Based on (nearest distance) from Centroids, [assign] observation to clusters.

**Step 5 :-** After each assignment (recalculate) mean of each attribute across all Clusters.

**Step 6 :-** (Repeat) Step 3, 4 & 5 until [Convergence] of Centroids i.e., Centroids don't change significantly.
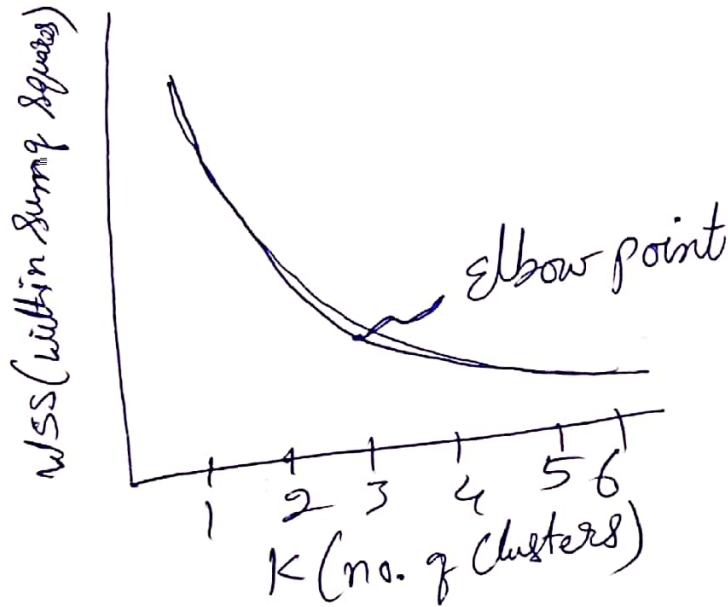
( see the Example in the pdf files )

Choose optimal k

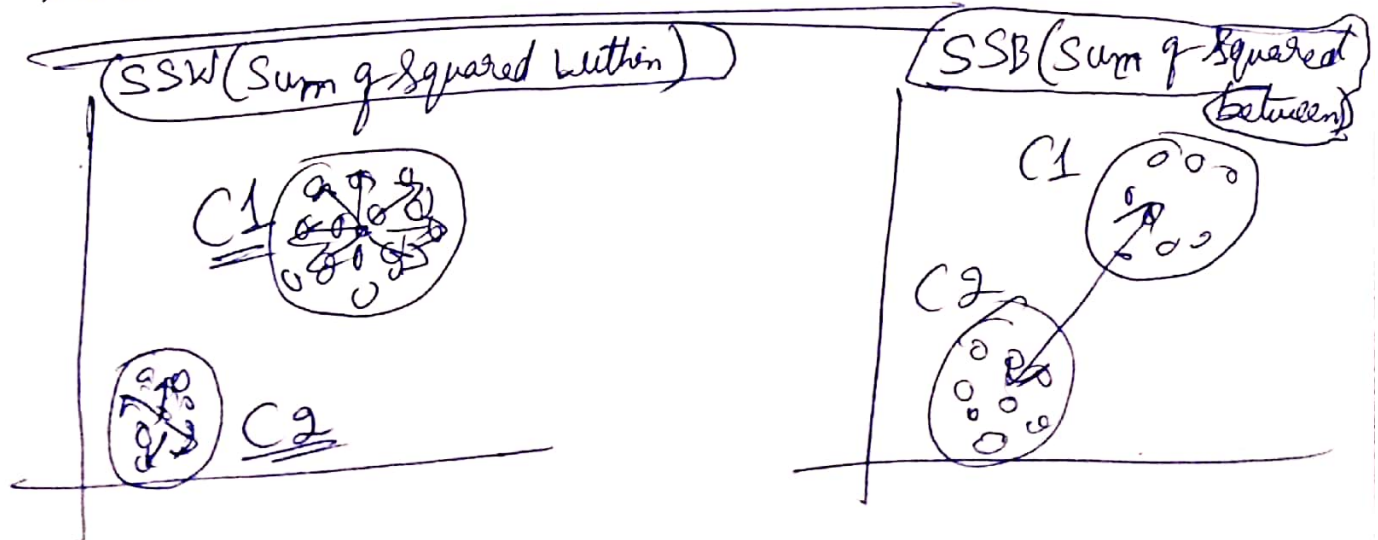↳ Choose optimal K is the key of K-means Clustering because

Improper selection of k may lead to erroneous assignment of observations.

There are a number of methods to choose optimal k such as Elbow method, GAP analysis, Average Silhoultte but the most common is Elbow method.



WSS (within Sum of squares) vs K (no. of Clusters)

Elbow point

Let's know about SSW & SSB

SSW (Sum of Squared Within)

C1
C2

SSB (Sum of Squared between)

C1
C2

Improper selection of K may lead to __erroneous assignment__ of observations.

There are a number of methods to choose optimal K such as __Elbow method__, GAP analysis, Average Silhouette but the most common is
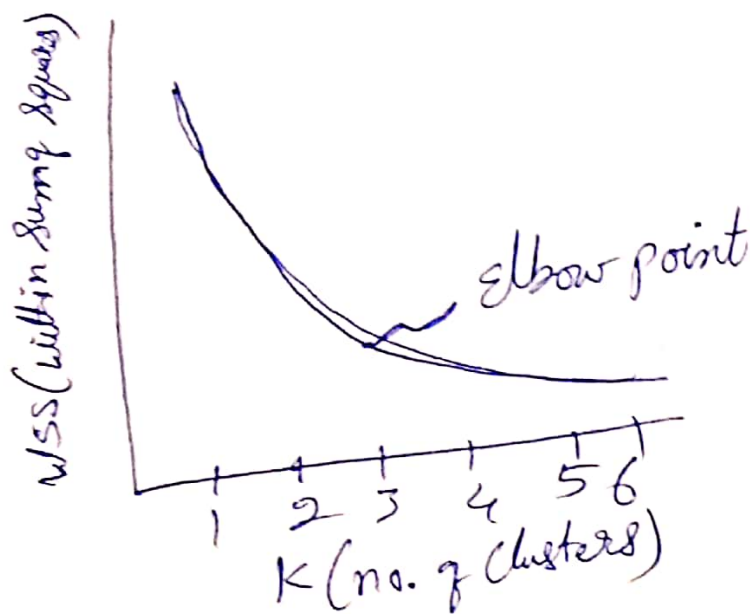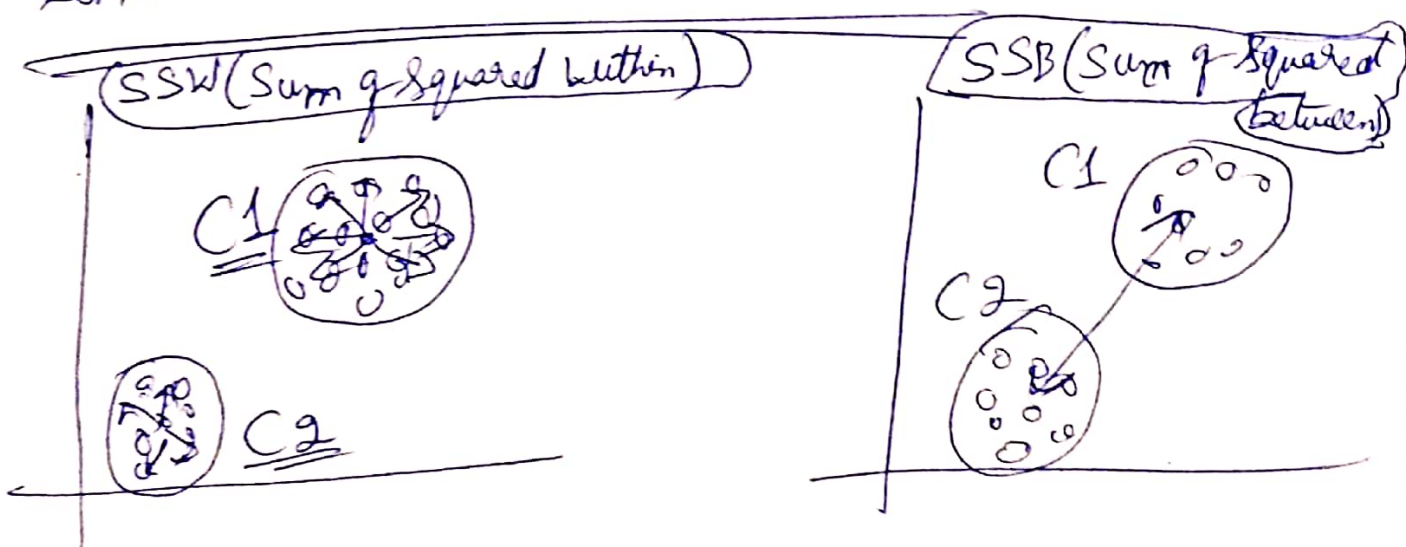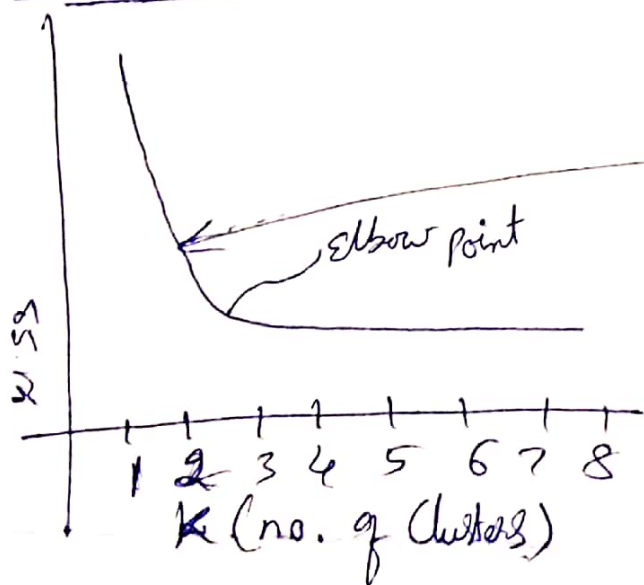
Elbow method.



Let's know about SSW & SSB

SSW (Sum of Squared Within)



SSB (Sum of Squared between)

## For perfect seperation

→ We want SSW to be least.

→ The smaller SSW the more compact data points are in clusters.

→ So that there is no intersection between 2 clusters.

→ SSB should be maximum

→ Max SSB indicates Centroids are distantly & perfectly positioned.

## Choose optimal K



Elbow point

K (no. of clusters)

→ Elbow method is very useful to find optimal K.

→ In this method SSW (within Sum of Square) decrease as number of clusters increases.

→ Elbow method helps us to find optimal K with increasing number of K.

→ We should take that number of K where elbow is spotted.

→ Adding clusters beyond Elbow point doesn't significantly reduce SSW.

# K - Means

① Import packages (Load the data)

②③ Clean Data ⟶ [ ~~Nan~~; fill / drop } numpy,
                    outlier treatment } Pandas

③ EDA ⟶ All visualizations should be
                                    done here ← matplotlib / seaborn + ~~Stat~~ statistics

→ Boxplt
→ Correlation
→ Histogram
→ Pairplt

~~Steps~~ ∴ Since it is an unsupervised
~~④ to~~ learning; we do not have to ① Create
features / labels & ② Split data ⟹ Not applicable

## Step 1

Finding ~~the~~ Clusters with Elbow method

```
SSW = [ ]
cluster_range = range(1,10)
    for i in cluster_range:
        model = KMeans (n_clusters = i, init = 'Kmeans++',
                        max_iter = 300, random_state = 0)
        model.fit(x)
        SSW.append(model.inertia_)
```

→ Print the clusters with SSW (use dictionary)

```
SSW_df = pd.DataFrame({ ....:.. ---, ... })
Print (SSW_df
```

## Plot the clusters

```
plt. figure (figsize = (10,17))
plt. plot (cluster_range, SSW, marker=o, color= " ")
plt. xlabel (" ...... ")
plt. ylabel (" _ _ _ ")
plt. title (" _     ")
plt. show ()
```

( Step 2 )

## Build a K-means model

```
Kmeans = KMeans(n_clusters = 4 , init = "Kmeans++",
                n_init = 10, random_state =42)

Y_kmeans = Kmeans.fit_predict(x)
```

→ ( Fitting the model )

## Step - 3

Visualizing the Clusters

plt. Scatter (X[Y-means == 0, 0], X[Y-means == 0,1], S=100, C=purple,
label = cluster1

$- - - - - \quad = = 1,0], [- - - - = 1,1,]$  $- - cluster 2$

$2, 0][- - - 2,1]$  $- cluster 3$

$3,0)[- . - 3,1]$  $- cluster 4$

plt. Scatter (kmeans.cluster centers _[:,0], kmeans cluster [:, 1]

S=300, marker = "s", C = 'red', label = 'Centroids')

plt. title (: - - - - -)
plt. xlabel (- - - - - .)
plt. ylabel (- - - - -)
plt. legend ()
plt. show .