

## Clustering

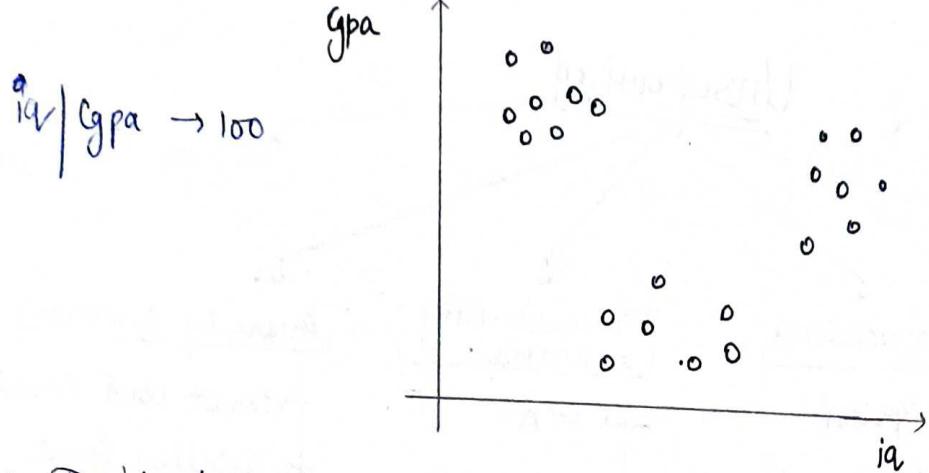
### Application of Clustering

1. Customer Segmentation.
2. Data Analysis
3. Semi Supervised Learning
4. Image Segmentation

### k-Means Geometric Intuition

#### Problem statement

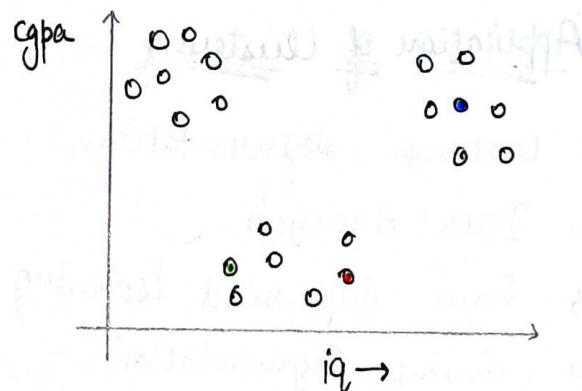
1. Decide k cluster
2. Initialize centroids
3. start Iterating
  - Assign clusters
  - Move centroid
  - Check and stop



1. Decide k cluster → We have to decide the cluster. later we will learn how to decide the cluster but now, we are assuming the value of k is 3.

2. Initialize centroids

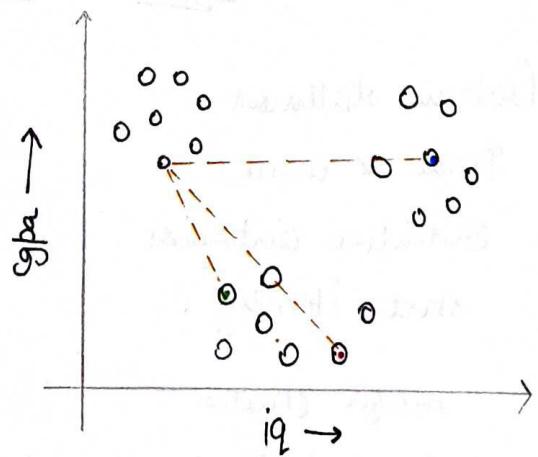
Randomly we can initialize centroids on any 3 points.



3. Starting Iterating → (loop)

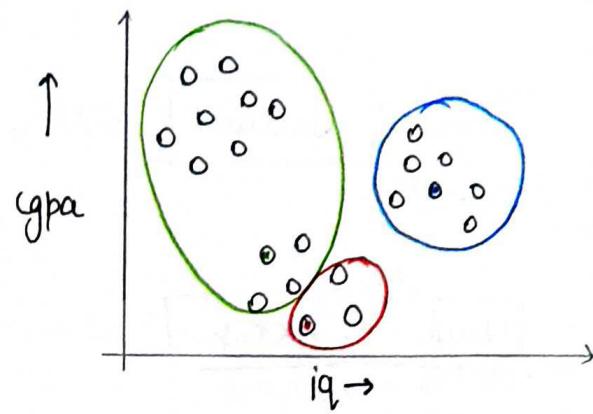
Find euclidean distance between point and centroids.

We have to find euclidean distance of all points.



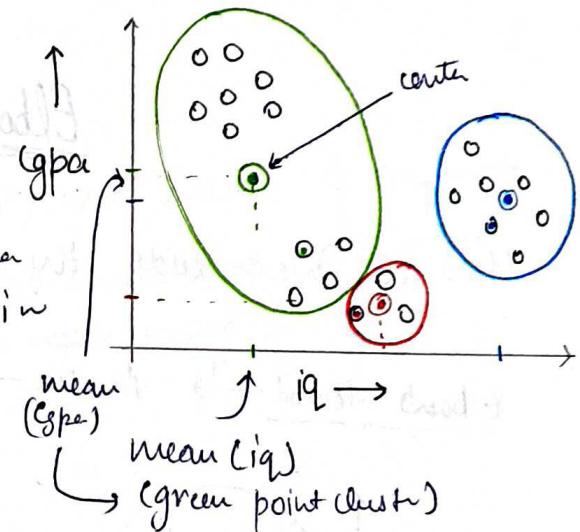
## Assign clusters

Make cluster of those points which are near to the centroid.



## Move centroid

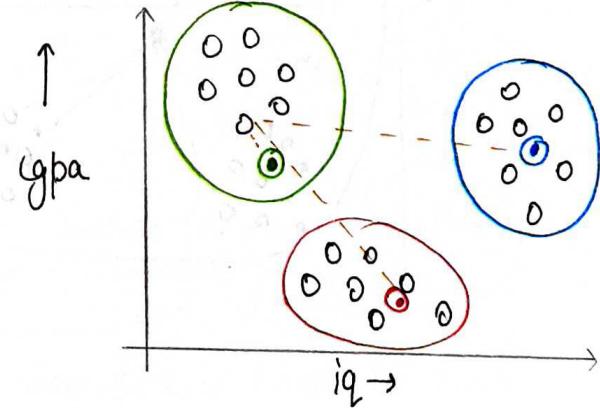
We have to find the mean of point which are present in cluster. example:- Mean  $\rightarrow$  Gpa  
Mean  $\rightarrow$  iq, Gpa and iq  $\rightarrow$  present in green cluster.



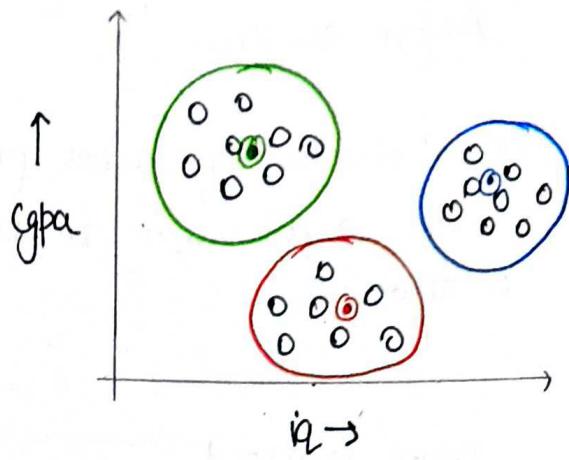
## Check and stop

condition  $\rightarrow$  if previous position of centroid = current position of centroid  
then  $\rightarrow$  stop otherwise algo work in continuously.

Now, Again find euclidian distance from point to centroid and again make clustering.



Move centroid

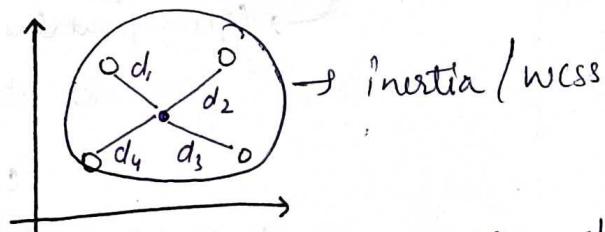


Check and stop

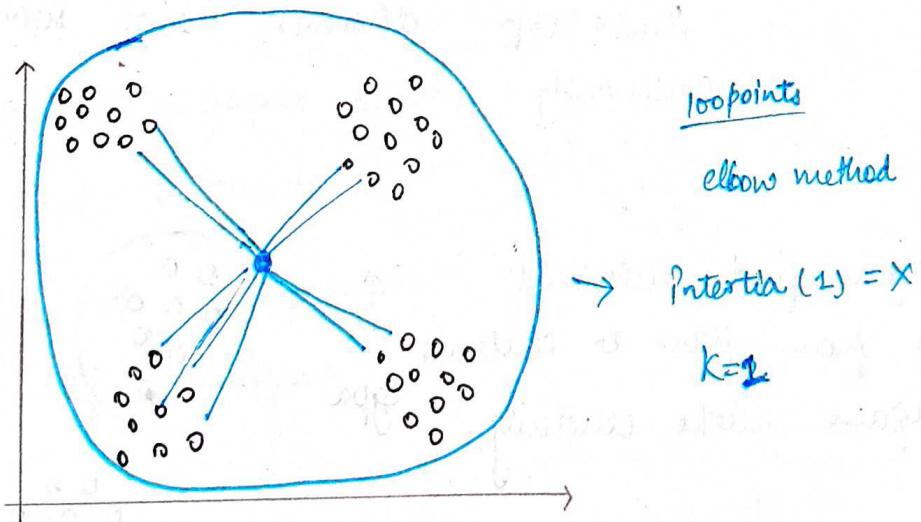
### Elbow Method

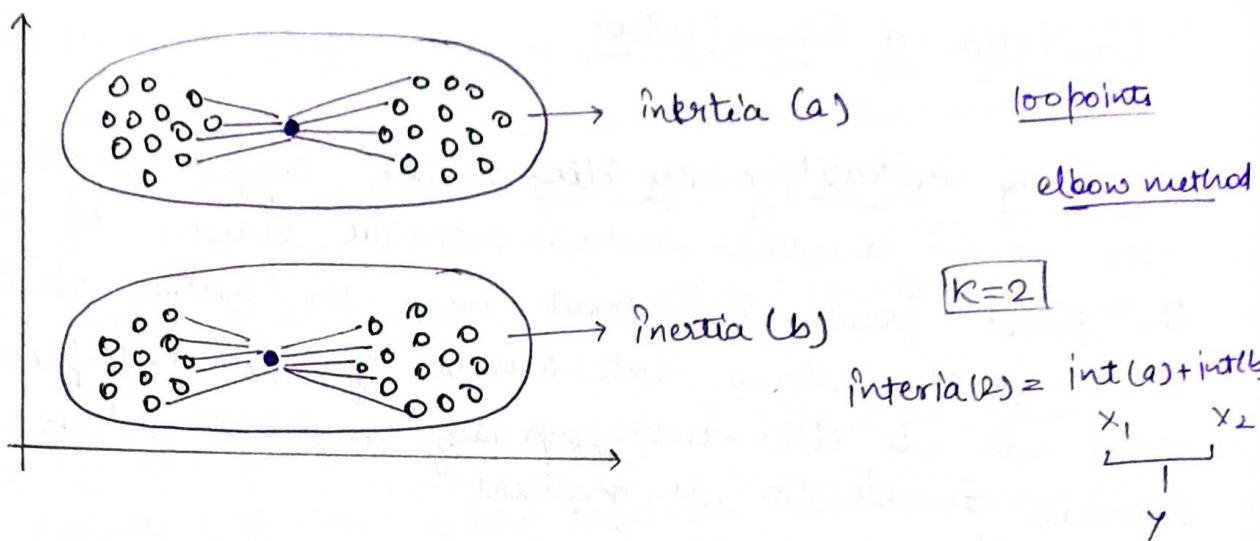
How do we decide the correct value of k?

Elbow Method  $\rightarrow$  Inertia  $\rightarrow$  WCSS  $\rightarrow$  within Cluster sum of square d distance



$$WCSS = d_1^2 + d_2^2 + d_3^2 + d_4^2$$



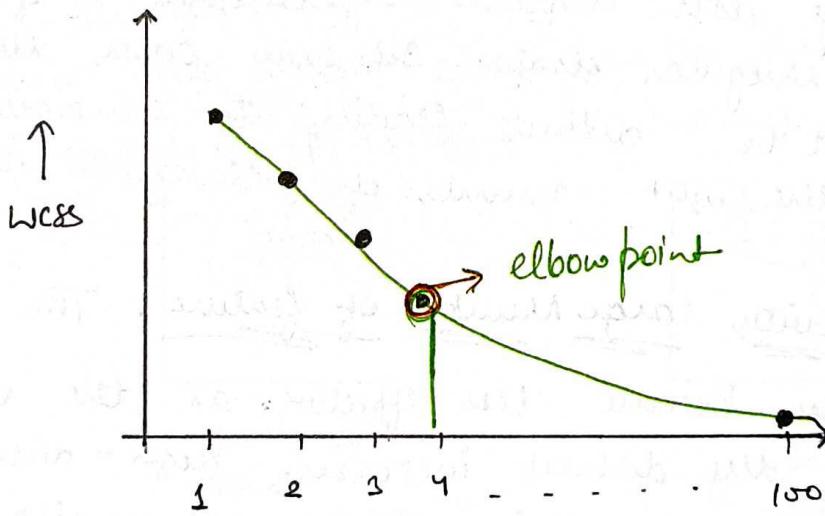


1)  $x > y ?$

2)  $y > x ?$

→ If we have 100 data points and inertia of 100 points is zero because 100 points → 100 centroid ⇒ between point and centroid, so ~~inertia~~ inertia is 0  
 $WCSS = d_1^2 + d_2^2 + d_3^2 + \dots + d_{100}^2 = 0$

$1 > 2 > 3 > 4 > 5 > \dots > 100$   
 ↴ inertia



→ Why we Reduce WCSS?

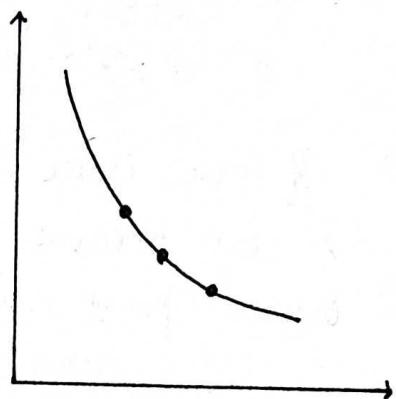
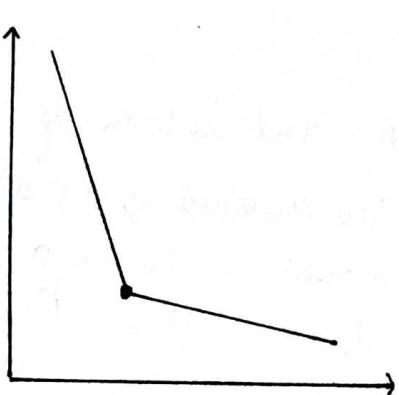
Because WCSS is like Variance →

$$WCSS = \frac{(d_1^2 + d_2^2 + \dots + d_{100}^2)}{n}$$

When WCSS Reduce → Variance also Reduce (not too much reduce) because 100 data → create 100 cluster.

## Limitation of Elbow Method

Subjectivity in identifying the elbow: The biggest challenge with the elbow method is the subjective nature of identifying the "elbow" point. The point where the inertia starts decreasing at a slower rate can be open to interpretation and may not be clear-cut, especially in datasets where the decrease in inertia is gradual.



Not Suitable for All Datasets: The method does not work well if the data is not very clustered or if the clusters have an irregular shape. In such cases, the elbow might not be distinct leading to ambiguity in choosing the right number of clusters.

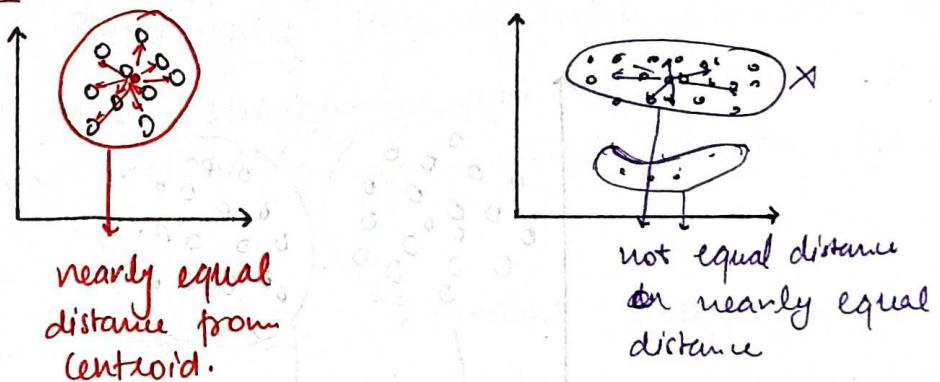
Performance with large Number of feature: The elbow method can become less effective as the number of features in the dataset increases. High-dimensional data can make the identification of a clear elbow more difficult.

Doesn't Consider Cluster Quality: The elbow method focuses solely on the variance within the clusters and does not take into account the quality of the clusters formed. It's possible to choose a  $k$  where clusters are not meaningful or well-separated.

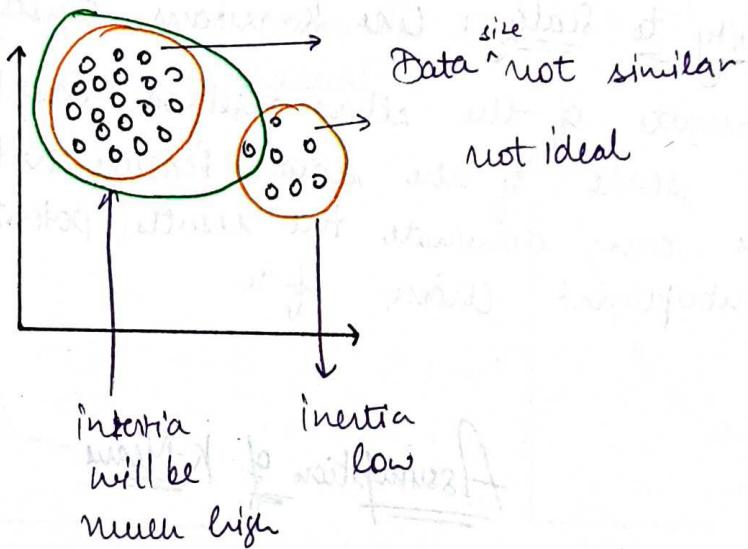
Sensitivity to scaling: Like k-means clustering itself, the results of the elbow method can be sensitive to the scale of the data. Features with larger scales can dominate the results, potentially leading to suboptimal choice of  $k$ .

### Assumption of K-Means

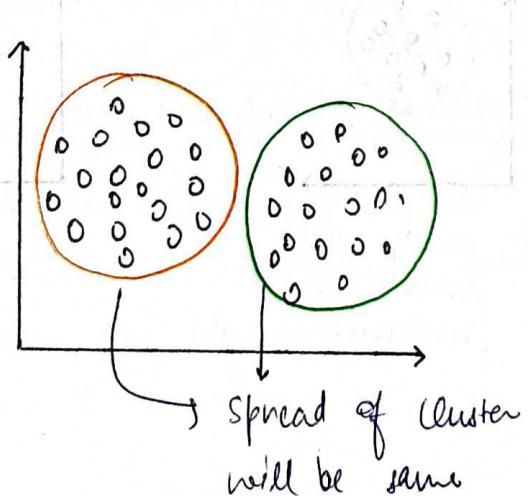
Spherical Cluster Shape: k-means assume that the clusters are spherical and isotropic, meaning they are uniform in all directions. Consequently, the algorithm works best when the actual cluster in the data are circular (in 2D) or spherical in (in higher dimensions).



Similar Cluster Size :- The algo tends to perform better when all cluster are of approximately the same size. If one cluster is much larger than others, k-means might struggle to correctly assign the points to the appropriate cluster.



Equal Variance of Cluster ! K-means assume that all clusters have similar variance. The algo uses the Euclidean distance metric, which can bias the clustering toward clusters with lower variance.



Clusters are well Separated: The algorithm works best when the clusters are well separated from each other. If clusters are overlapping or intertwined, k-means might not be able to distinguish them effectively.

Number of Cluster (k) is predefined: k-means requires the number of clusters ( $k$ ) to be specified in advance. Choosing the right value of  $k$  is crucial, but it is not always straightforward and typically requires domain knowledge or additional methods like the Elbow method or Silhouette analysis.

large n, small k: k-means is generally more efficient and effective when the dataset is large ( $large n$ ) and the number of clusters is small ( $small k$ ).

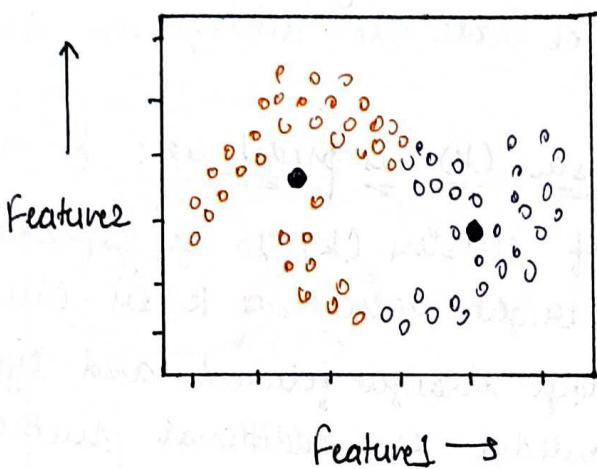
### Limitation of K Means

Number of Clusters: Determining the optimal number of clustering ( $k$ ) is not straightforward and often requires domain knowledge or method like the elbow method.

Require cluster of similar size: k-means require the clusters to be of similar size.

Similar Variance between cluster: k-means require the cluster to be of similar variance.

Assumption of Spherical Cluster: KMeans assume that clusters are spherical and of similar size, which might not be the case in real world data.



Vulnerability to outliers: Outliers can significantly distort the mean value of a cluster, leading to misleading result.

Hard Clustering: Each data point is forced into exactly one cluster, which may not be suitable for all application, especially where data can belong to multiple cluster.

Soft Clustering → Not use in Kmean

Yaussian Mixture Model

↳ data [  $\frac{0.78}{\text{prob of Clust 2.}}$ ,  $\frac{0.62}{\text{prob of Clust 1.}}$  ]  
↳ prob of  
Cluster 1 assign  
data 1

Higher-Dimensional Challenges: In very high-dimensional spaces, the distance between data points can become less meaningful, affecting the performance of k-means.

Sensitive to Scale: The measure is sensitive to the scale of the features. Hence, feature scaling (like standardization) is often recommended before applying k-means.

### Silhouette Score

Cohesion (intra cluster distance)

Definition: Cohesion refers to the degree to which elements within the same cluster are close to each other. It measures how tightly grouped the data points in a cluster are.

Ideal Scenario: High cohesion means that data points in a cluster are similar or near to each other, indicating a good clustering where each cluster is distinct and meaningful.

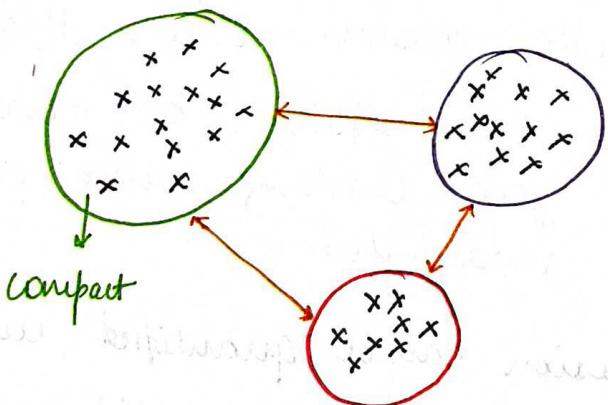
Measurement: Cohesion can be quantified using metrics such as the sum of squared distance of data points from their respective cluster centroids. In k-means this is often referred as inertia or within-cluster sum of squares.

## Separation (inter cluster distance)

Defination: Separation, on the other hand, refers to how distinct or well-parted different clusters are from each other. It measures the extent to which clusters are different or distant from each other.

Ideal Scenario: High separation means that clusters are well-differentiated and far apart, indicating that the algo has done a good job in distinguishing between different groups in the data.

Measurement: Separation can be quantified metrics such as the distance between cluster centroids, or more complex measures like the silhouette score which consider both cohesion and separation.



- 1.) Cohesion = ↓ low
- 2.) Separation = high

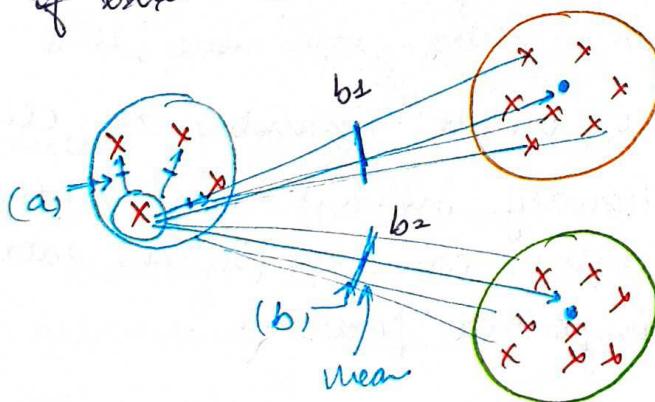
## Silhouette Score

The silhouette score is a measure used to assess the quality of cluster created by a clustering algo. It provides a succinct graphical representation of how well each data point lies within its cluster, which is a combination of both cohesion and separation. The value of the silhouette score range -1 to +1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighbouring clusters.

$$S_{ci} = \frac{b_{ci} - a_{ci}}{\max\{a_{ci}, b_{ci}\}}$$

(a) = mean of distance b/w selected point and all point in particular cluster.

(b) = mean of distance b/w selected point and ~~out~~ point of other cluster.



if  $b_1 < b_2$   
↳ we consider  $b_1$   
if  $b_2 < b_1$   
↳ consider  $b_2$

## Interpretation:

close to +1: Indicates that the data point is far away from the neighbouring clusters.

Close to 0: Indicates that the data point is on or very close to the decision boundary b/w two neighbouring clusters.

Close to -1: Indicates that the data point may have been assigned to the wrong cluster.

## Kmeans Hyperparameter

```
class sklearn.cluster.KMeans(n_clusters=8, init='k-means++',  
n_init=10, max_iter=300, tol=0.0001, verbose=0,  
random_state=None, copy_x=True, algorithm='lloyd')
```

### Number of Clusters (k):

Description: This is the number of cluster you want the algorithm to form, as well as the number of centroids to generate.

Importance: Choosing the right number of cluster is crucial as it significantly influence the clustering results. Too many clusters can overfit the data, while too few can miss important patterns.

### Initialization Method

Description: This parameter specifies the method for initializing the centroids. Common methods include

'Random' (randomly choosing K data points as initial centroids) and 'k-means++' (a smarter way of initializing centroid to improve convergence).

Importance: Good initialization can lead to faster convergence and better clustering. k-means++ is generally preferred over random initialization.

### Number of Initialization Runs (n\_init):

Description: This is the number of times the KMeans algo will be run with different centroid initializations. The final results will be the best output of n-init consecutive runs in term of inertia.

Importance: Multiple initialization can prevent the algo from falling into sub-optimal solutions, but increase computational cost.

### Maximum Number of Iteration (max\_iter):

Description: The maximum number of iterations the algo will run for each initialization.

Importance: A higher number of iteration allows more time for convergence but increase computational time. Usually, KMeans converges well before reaching the maximum number of iteration.

### Tolerance (tol):

Description: This is the tolerance to declare convergence if the centroids do not move significantly (as

defined by this parameter) in consecutive iterations, the algorithm stops.

**Importance:** A smaller tolerance can lead to a more precise solution, but might increase computation time. A large tolerance might speed up the algo but can lead to less precise clustering.

### Algorithm

**Description:** The computational algo to use. Options typically include 'auto', 'full', or 'Elkan'. 'full' corresponds to the classic EM-style algo, while 'elkan' is an optimized variant that is generally more efficient.

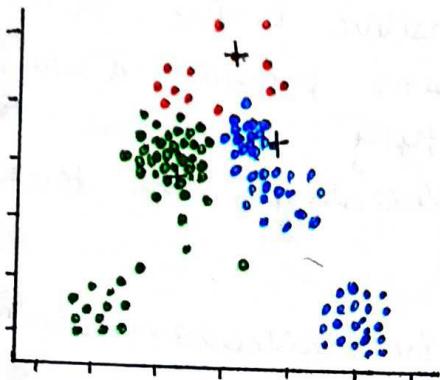
**Importance:** The choice of algo affects the computational efficiency. 'Elkan' is often faster but works only with Euclidean distance.

### Random State:

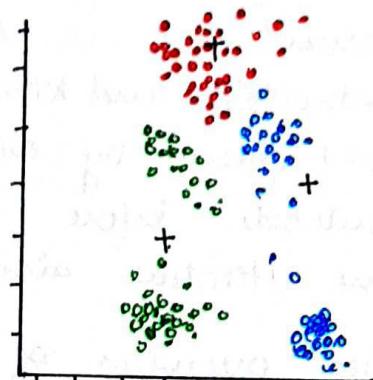
**Description:** Determines the random number generation for centroid initialization, which influences the behaviour of the algo.

**Importance:** Setting a random state ensures reproducibility, as KMeans can produce different results on different runs due to its stochastic nature.

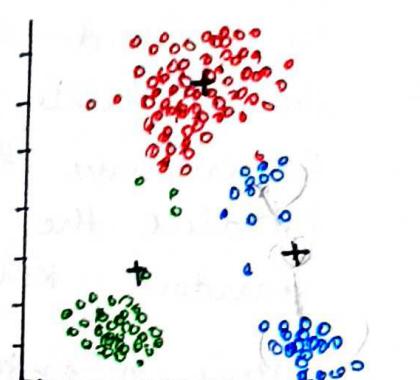
## Kmeans++



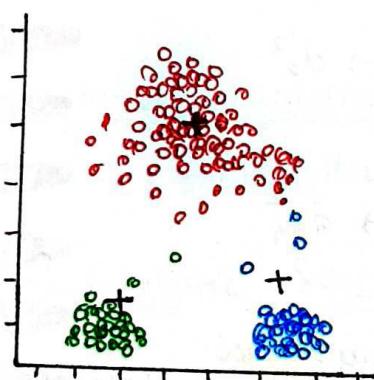
iteration 1



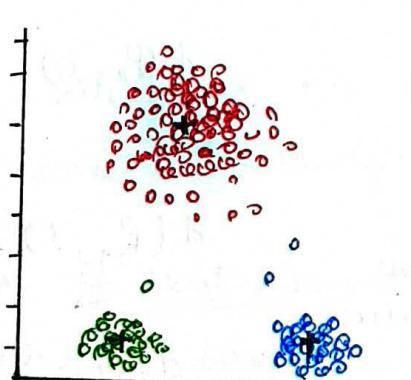
iteration 2



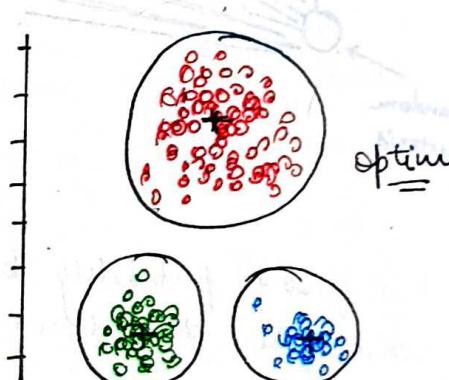
iteration 3



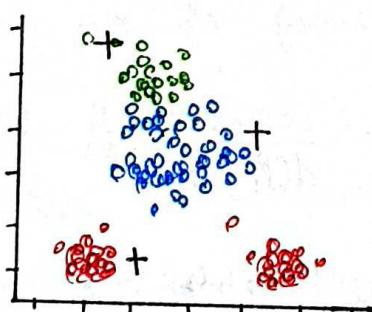
iteration 4



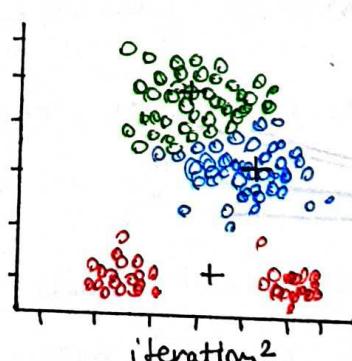
iteration 5



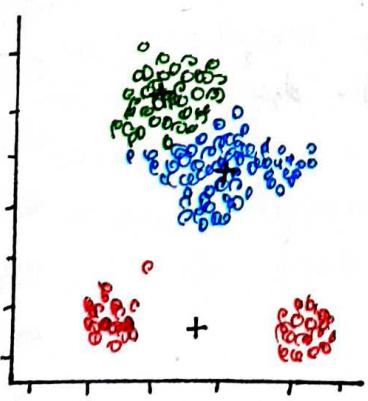
iteration 6



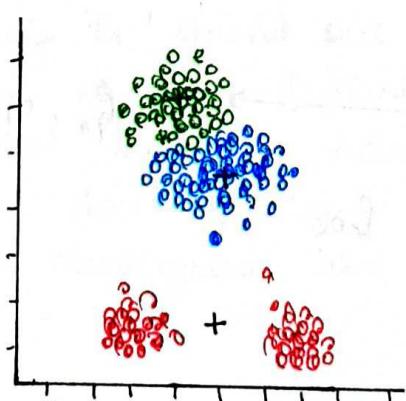
iteration 1



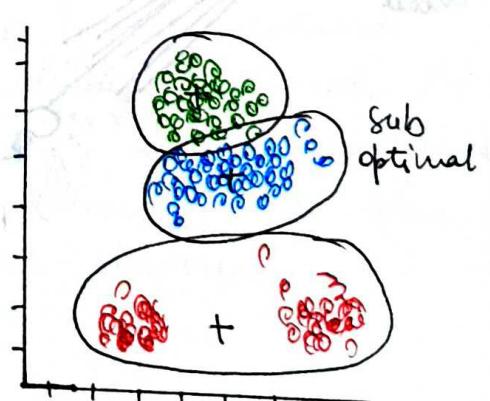
iteration 2



iteration 3



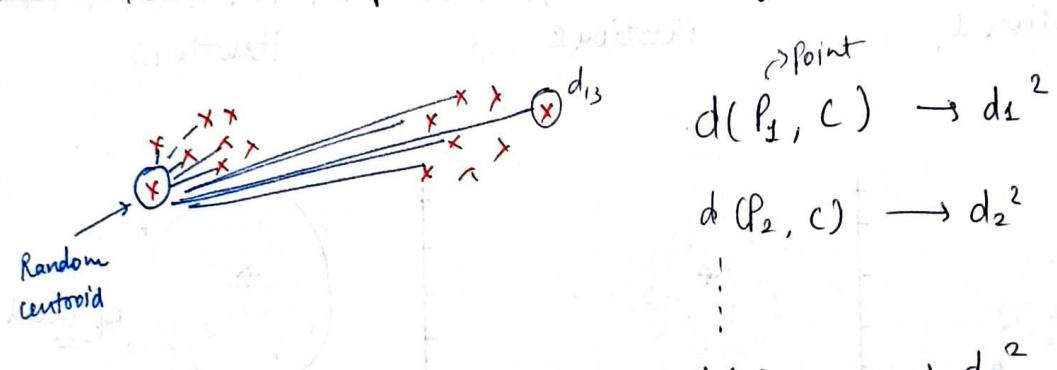
iteration 4



iteration 4

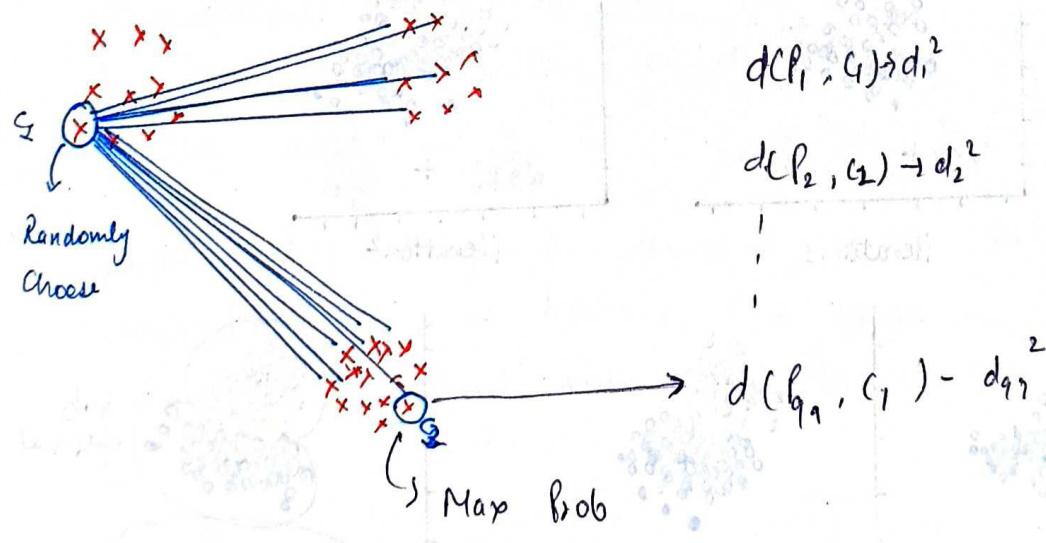
KMeans++ is an algorithm for choosing the initial values (or "seeds") for the KMeans clustering algorithm. The standard KMeans algo is sensitive to the initial starting points (centroids), and KMeans++ provide a way to overcome this problem by specifying a procedure to initialize the centroids before proceeding with the standard KMeans iterative algo.

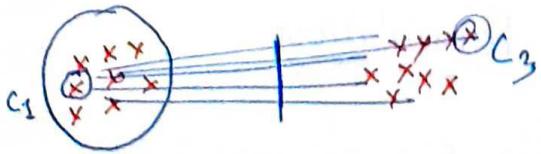
Here's a simplified overview of how KMeans++ works:



\* We find probability because some time data have outliers

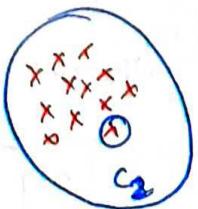
for 3 Cluster





$c_3 \rightarrow$  every  $\rightarrow$  nearest point centroid

if  $c_2$  is nearest then we can find the distance b/w  $c_2$  and every point.



Initial Centroid Selection: The first centroid is chosen uniformly at random from the data points that are being clustered.

Distance Calculation: calculate the distance of each data point from the nearest, previously chosen centroid.

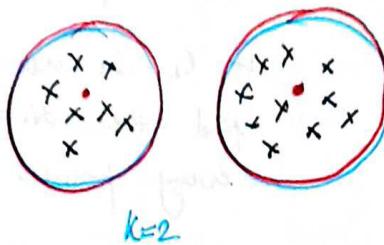
Probabilistic Selection of Next Centroids: Choose the next centroid from the data points with a probability proportional to the square of the distance from the point to its nearest centroid. This step biases the algo to select data points that are far from the existing centroids.

Repeat until k Centroids: Repeat steps 2 and 3 until  $k$  centroids have been chosen.

Proceed with Standard KMeans: Once the initial centroids are chosen, proceed with the standard KMeans clustering algo. This method tends to spread out the initial centroids, which can lead to better clustering results compared to selecting the initial centroids randomly, as the standard KMeans algorithm does. By doing so, KMeans++ can often lead to faster convergence and better clustering.

## Kmeans Mathematical formulation

Kmeans



[Lloyd's Algo] step

↓  
Approximation

$$J = \underset{\mu_1, \mu_2, \dots, \mu_k}{\operatorname{argmin}}$$

$$\left[ \sum_{i=1}^k \left[ \sum_{x \in S_i} [ \|x - \mu_i\|^2 ] \right] \right] \rightarrow \text{Inertia}$$

such that

$$S_i \cap S_j = \emptyset$$

$\mu_1, \mu_2, \dots, \mu_k \rightarrow \text{centroid}$

↳ Not common point  
betw 2 cluster or

k centroid

$S_i \rightarrow \text{current cluster}$

→ we have to find the value of  $\mu_1, \mu_2, \dots, \mu_k$  which reduce the value of inertia.

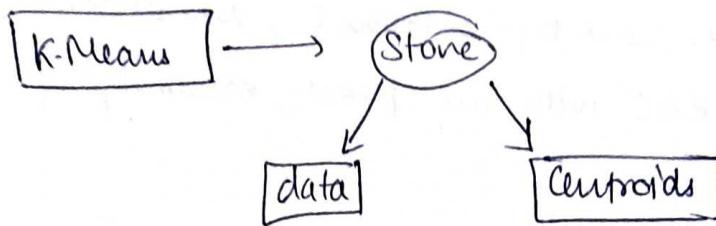
→ This is not but sol<sup>n</sup> because everytime we take new initial ~~and~~ then everytime new solution form through Lloyd's Algo.

KMeans Time Complexity

Kmeans → time complexity

no. of rows      O(n k d i)      Iteration  
 Centroids      dimension

## K-Means Space Complexity



$$O(nd + kd)$$

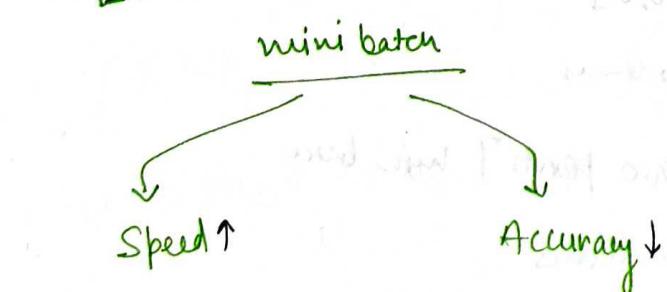
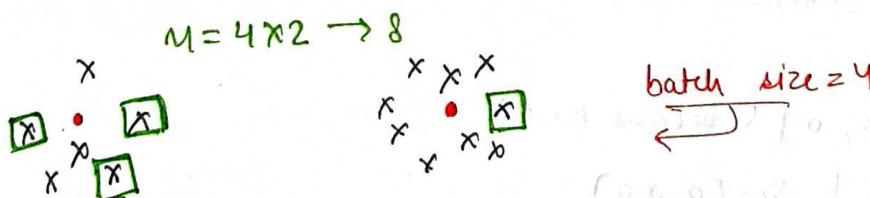
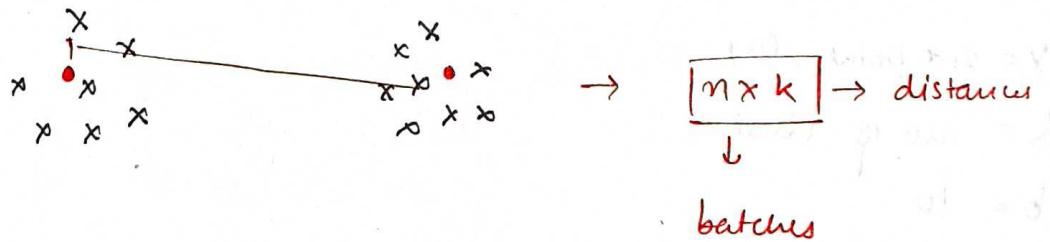
no. of rows  $\leftarrow$  data  
cols  $\leftarrow$  data  
no. of centroids  $\leftarrow$  Centroids  
cols  $\leftarrow$  Centroids

## Mini Batch K-Means

Variation  $\rightarrow$  k-means

b) Mini Batch k-Means

c) Gradient Descent



## Algorithm 1 Mini-batch k-Means

1. Given:  $k$ , mini-batch size  $b$ , iteration  $t$ , data set  $X$
2. Initialize each center with one picked randomly from  $X$
3.  $V \leftarrow 0$
4. for  $i = 1$  to  $t$  do
5.    $M \leftarrow b$  examples picked randomly from  $X$  (mini batch)
6.   for  $x \in M$  do
7.      $d(x, c) \leftarrow f(c, x)$  // Cache the center nearest to  $x$
8.   end for
9.   for  $x \in M$  do
10.      $e \leftarrow d(x, c)$  // Get cached center for this  $x$
11.      $v_{cl} \leftarrow v_{cl} + 1$  // Update per-center counts ~~choose closest initially~~
12.      $n \leftarrow \frac{1}{v_{cl}}$  // Get per-center learning rate
13.      $c \leftarrow (1 - n)e + nx$  // Take gradient step
14.   end for
15. end for



$X = 200$  points ( $n$ )

$k = \text{No. of cluster}$

$b = 10$

$i = 10$  (iteration)

Step 3:  $V = [0, 0]$  (because  $k=2$ )

if  $k=3$   $V = [0, 0, 0]$

Step 4: loop  $\rightarrow 1 - t/10$  times

Step 5: [Randomly 10  $\rightarrow 200$  points] mini batch

Step 6: Nested loop  $\rightarrow 10$  times

$$x_1 \rightarrow c_1 \text{ Step: 7 } d(x, c) = [x, c]$$

$$x_2 \rightarrow c_1$$

$$x_3 \rightarrow c_2$$

:

$$x_{10} \rightarrow c_2$$

Step 8: loop end

Step 9: loop  $\rightarrow$  10 times

Step 10:  $c \leftarrow d[x] \Rightarrow$  we have to find nearest centroid with the help of step 7:  $d(x, c)$ .

Step 11:  $v[c] \leftarrow v[c] + 1 \Rightarrow$  (update value)

$c_i$

$$v = [1, 0]$$

( $\downarrow$  centroid near the point 1).

Step 12:  $\eta \leftarrow \frac{1}{v[c]} = \frac{1}{1} = 1$  (learning rate)

Step 13:  $c_1 = (1 - \eta) c_1 + \eta x_1$

$$c_1 = \underbrace{(1 - 1)}_0 c_1 + 1 x_1$$

$$c_1 = x_1$$

another example

$x_2 \rightarrow c_1 \Rightarrow$  2 points near centroid

$$v_2 = [2, 0]$$

$$\eta = \frac{1}{2} = 0.5$$

$$c_1 = (1 - 0.5)c_1 + 0.5x_2$$