

Table Of Contents

1. Concept of K-Means Clustering
2. Math Intuition Behind K-Means
3. Cluster Building Process
4. Edge Case Scenarios of K-Means
5. Challenges and Improvements in K-Means

Main Requirements for clustering Algorithms

- ① It should be Scalable.
- ② It should be able to deal with attributes of different type.
- ③ It should be able to discover arbitrary shape clusters.
- ④ It should have an inbuilt ability to deal with noise & outliers.
- ⑤ The clusters should not vary with the order of input records.
- ⑥ It should be able to handle data of high dimension.
- ⑦ It should be easy to interpret & use.

⑩ K-Means

old Algorithm

Proposed in 1957 by Stuart Lloyd, and
in 1965 by Edward W. Forgy. It
Sometimes called Lloyd-Forgy Algo.

- K-Means is a clustering approach in which data is grouped into K distinct non overlapping clusters based on their distances from the K centres.
- the value of 'K' needs to be specified first & then the algorithm assigns the points to exactly one cluster.

Theory of K-means

→ Let C_1, C_2, \dots, C_k be the k -clusters.

→ then we can write

$$C_1 \cup C_2 \cup C_3 \cup \dots \cup C_k = \{1, 2, 3, \dots, n\} \rightarrow$$

$$\text{and also, } C_k \cap C_{k'} = \phi \longrightarrow$$

i.e.,
Each data point
has been assigned
to a cluster

this means, clusters are non-overlapping.

→ the idea behind the K-means approach is that within-cluster variation amongst the points should be minimum.

The within-cluster variance is denoted by $\rightarrow \underline{w(C_k)}$

→ Hence, according to statement above, we need to minimise the variance for all clusters mathematically,

WCSS \leftarrow
within cluster
sum of squares

$$\text{minimize}_{C_1, C_2, \dots, C_k} \left\{ \sum_{k=1}^k w(C_k) \right\}$$

→ the next step is to define the criterion for measuring the within cluster variance.

Generally, the criterion is the Euclidean distance between two data points.

$$w(C_k) = \frac{1}{|C_k|} \sum_{i, j \in C_k} \sum_{p=1}^p (x_{ip} - x_{jp})^2$$

$p = (1, 2, \dots, p)$

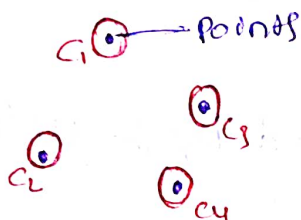
→ the above formula says that we are
 ① calculating the distance b/w all the
 points in a cluster, ② then we are repeating
 it for all the k -clusters (so two summations)

Eg:-

Let's take 4 points

Case ① :-

4 points
 4 clusters



→ now, each point
 acts as a cluster
 and so it is
 centroid of
 itself.

∴ distance b/w
 the point &
 centroid (same point)
 is zero (0)

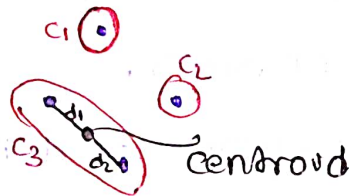
→ And it is
 same for
 all the 4 clusters
 in this.

i.e) $w(c_1) =$
 $0 + 0 + 0 + 0 = 0$
 $c_1 \quad c_2 \quad c_3 \quad c_4$

↓
Summation of
 distance.

Case ② :-

4 points
 3 clusters



→ now we need
 distance b/w
 point & the
 centroid of
 cluster it is
 present in.

$$w(c_3) = c_1 + c_2 + c_3$$

$$w(c_3) = 0 + 0 + (d_1 + d_2)$$

① individual
 sum of distance
 of a cluster

$$\sum_{i=1}^p (x - c)^2$$

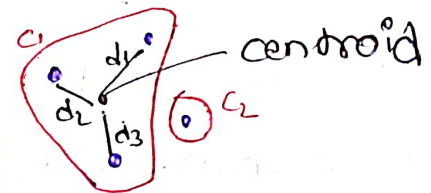
L centroid

② All the
 clusters
 sum

$$\sum_{j=1}^k \sum_{i=1}^p (x - c)^2$$

Case ③ :-

4 points
 2 clusters



$$w(c_4) = c_1 + c_2$$

$$w(c_4) = (d_1 + d_2 + d_3) + 0$$

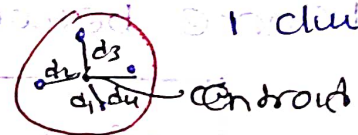
↓

(distance b/w point & centroid) Summation of distance in a cluster.

Summation of the
 individual
 summations of
 distance in cluster

Case ④ :-

4 points
 1 cluster



$$w(c_1) = c_1$$

↓

$$w(c_1) = d_1 + d_2 + d_3 + d_4$$

Goal is to minimize $w(c_1) + w(c_2) + w(c_3) + w(c_4)$ → low variance for all clusters

So, we got

$$w(c_4) = 0$$

$$w(c_3) = d_1 + d_2$$

$$w(c_2) = d'_1 + d'_2 + d'_3$$

$$w(c_1) = d''_1 + d''_2 + d''_3 + d''_4$$

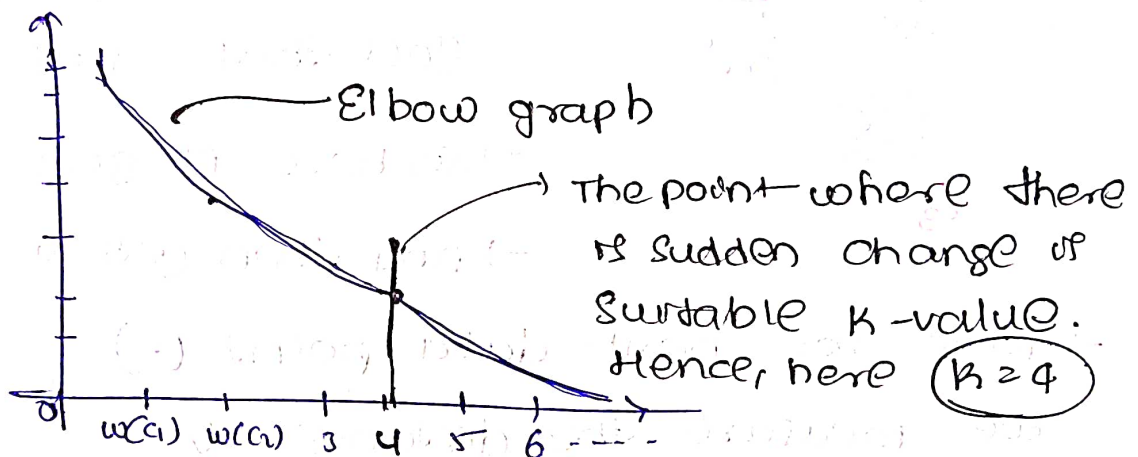
By observation, we can say that

As no. of clusters increase, $w(c_k)$ value decreases & tends to '0'

$$w(c_k) \propto \frac{1}{\text{no. of clusters}(k)}$$

$$w(c_1) > w(c_2) > w(c_3) > \dots > w(c_k)$$

and the graph goes as -



within cluster summation or square ($wcss$)

→ this is how we find the value of 'k' in k-means

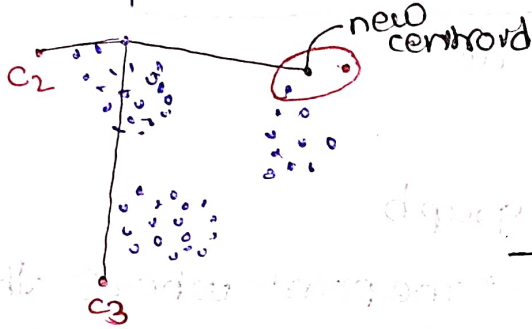
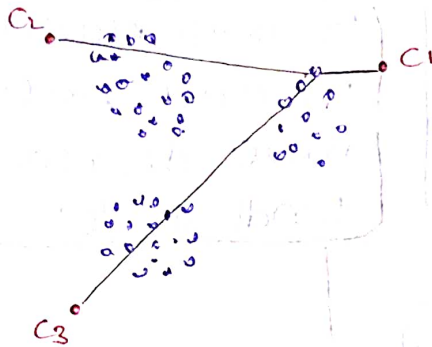
But, here is the problem in k-means,

⊛ After selecting k value, it takes k-points as centroids and starts building the clusters normally.

But thing is where do we place those k-points?

* cluster building process

let's say we got: $K=3$



→ After, Randomly placing the K points,

→ It takes a point from our data (.) and measures the distance to all the K -centroid points (.), And whichever distance is smaller, the point goes to that cluster.

→ So, here it goes to C_1

→ And then gets a new centroid.

→ now for next data point (.) we measure the distance of point to the "new centroids"

→ and in this way, cluster keeps building and centroid gets updated.

→ this process continues until all the data points are finished.



this is how grouping happens.

Edge case scenario

Problem with k-means

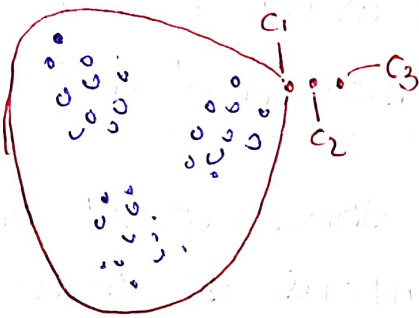
→ In the above example, we took $k=3$ and then it generated 3 random points and build cluster.

→ So, what if all the 3 points are close?

If all are close to only c_1 .

And,

$c_1 = 0$
 $c_2 = 0$ } elements



In this case, all data points would go to cluster ① (c_1)

which is not what we are expecting!!

Algorithm

- ① Find wcss for k-groups possible ($\text{max } k = n$) ^{not elem}
 - ② Build Elbow graph, and understand k-value (point where there's differentiation)
 - ③ Randomly assign k-centres.
 - ④ calculate distance of all points from all k-centres and allocate the points to cluster based on shortest distance to cluster centroid.
- The model's inertia is the mean squared distance b/w each instance & closest centroid. Goal is to have model with low inertia.
- ⑤ once all points assigned to cluster, recompute centroid.

- ⑥ Repeat steps 4 & 5 until the location of the centroids stop changing and the cluster allocation of the points becomes constant.
-

Note:

→ K-means is only recommended for spherical data (or) symmetrical data.

→ For, non-spherical data there won't be fair chance for K centroids to form clusters, and it will become high biased clusters.

* Challenges and Improvements in K-Means

- ① we need to specify the no. of clusters beforehand.
 - ② It is required to run the algorithm multiple times to avoid a suboptimal solution.
 - ③ K-means does not behave very well when the clusters have varying sizes, different densities or non-spherical shapes.
 - ④ An improvement to K-means is proposed in 2003 by Charles Elkan.
 - It accelerates the algorithm by many unnecessary distance calculations which is achieved by exploiting the TRIANGLE INEQUALITY.
- In a triangle with sides a, b, c $\Rightarrow a + b > c$
i.e., the straight line is always shortest

⑤ Another Important Improvement to k-means algorithm called "k-means++", was proposed in 2006 by David Arthur & Sergei.

→ they introduced a smarter initialization step that tends to select centroids that are distant from one another, and this makes the k-means algorithm much less likely to converge to a suboptimal solution.

⑥ Yet another Important variant of the k-means algorithm was proposed in a 2010 paper by David Sculley.

→ Instead of using the full dataset at each iteration, the algorithm is capable of using mini-batches, moving the centroids just slightly at each iteration. This speeds up the algorithm typically by a factor of 3 or 4 and makes it possible to cluster huge datasets that do not fit in memory.

Scikit-Learn implements this algorithm in the "MiniBatchKMeans" class.
