

K-Means Clustering

clustering: it is unsupervised machine learning techniques in which we combine similar type of data into the group.

- * K-mean clustering is very easy to visualize in 2D and 3D but in the higher dimensions (large number of columns) then it is harder to visualize it. we only perform it
- * for the higher dimension it is very effective
- * the steps in k-mean clustering:

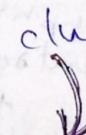
(1.) Decide n-clusters



(2.) Initialize Centroids



(3.) Assign cluster

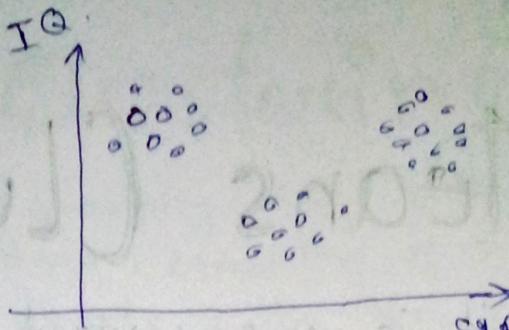


(4.) Move Centroids



(5.) Finish

इस algorithm में programmer को इस algorithm की बताना पड़ता है कि हमें कितने clusters करने हैं।

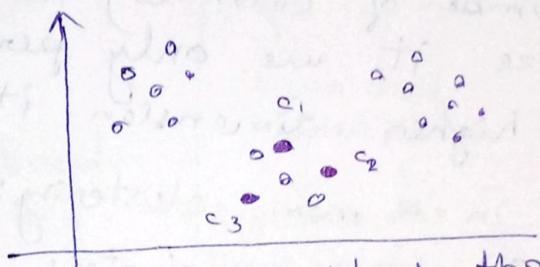


Step-1 let us assume in the above graph there is have to make a three cluster.

$$k = 3$$

Step-2 initialize centroid

algorithm take any three points in the graph and assume it as the centroids



c₁, c₂, c₃ is centroid of three cluster.

Step-3 Assign clusters.

use the euclidean distance for each point on the graph calculate the distance from each centroid. which one is minimum them that point is member of that cluster when centroid lies.

example.

P₁ distance from c₁ = 20

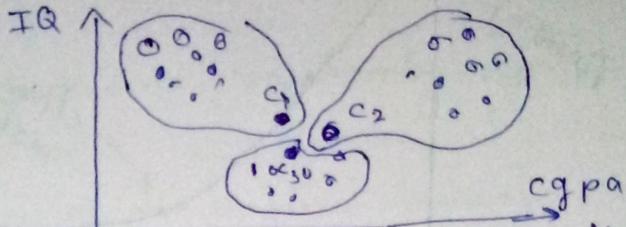
P₁ distance from c₂ = 10

P₁ distance from c₃ = 22

Hence, P₁ is the part of C₂ cluster.

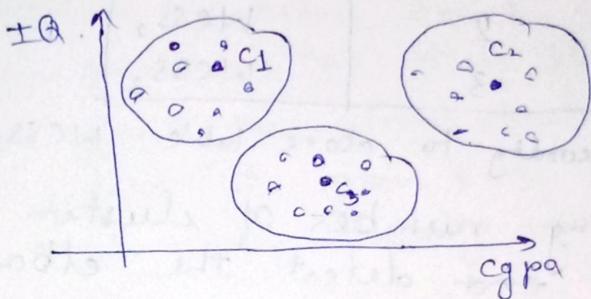
the above approach is apply for each point

Step-4 Move centroid.



according to above graph, cluster is formed.

Now for moving centroid for each cluster.
calculate the intersection point of. mean of
cgpa and mean of IQ for each cluster



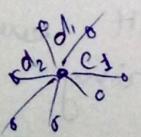
Step-5 Finish.

if the location of each centroid ~~is~~ in
previous and present step is same
then stop it else continue from step-3

How to decide the number of cluster:

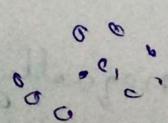
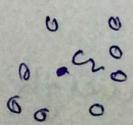
for finding the number of cluster we use
elbow method. in which we plot a graph
that is called elbow curve and has
two axis number of clusters and ~~is~~ inertia
(within cluster sum of square distance (WCSS)).

c_i = centroid.

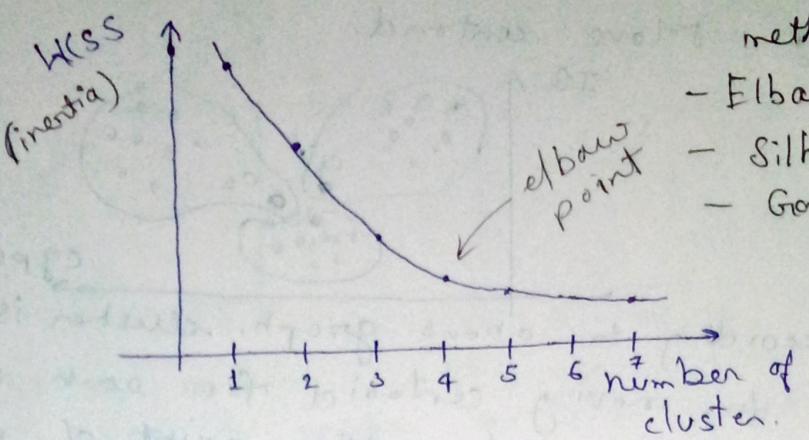


$$WCSS = d_1^2 + d_2^2 + d_3^2 + \dots + d_r^2$$

c_1, c_2 = centroid.



$$WCSS = WCSS_1 + WCSS_2 + \dots + WCSS_n$$



- Elbow method
- Silhouette Score
- Gap Statistic

no. of cluster	WCSS
1	WCSS ₁
2	WCSS ₂
3	WCSS ₃

according to above table $WCSS_1 > WCSS_2 > WCSS_3$

for finding number of cluster we want to select who detect the elbow point in elbow curve.

elbow point = a point in which curve is parallel to x-axis.

Some Drawback of k-mean clustering:

- it is difficult to determine the optimal k for a given dataset
- Different initial partitions can result in different final clusters
- It doesn't work well with clusters of different size and density.
- it is very sensitive to scaling and requires careful pre-processing.
- it does not produce the same result every time
- it is only for spherical clustering
- It requires manual work
- it clusters everything.
- euclidean distance can unequally weigh the factors.

- it gives the local optima of the squared error function.
- choosing the centroids randomly can't give fruitful results.
- As the number of dimensions increases, a distance based similarity measure converges to a constant value between any given examples.

Question:

Solve a numerical problem on k mean clustering
The problem problem has given 15 points.
 we are also given the information that we need to make 3 clusters. it means $k=3$. we will solve this numerical on k-means clustering using

first randomly choose 3 centroids from the given points. let us consider $A_2(2,6)$, $A_7(5,10)$, $A_{15}(6,11)$ as the centroids of the initial clusters.

Centroid-1 = $(2,6)$ is associated with C_1

Centroid-2 = $(5,10)$ is associated with C_2

Centroid-3 = $(6,11)$ is associated with C_3

Now we will find the euclidean distance between each point and the centroids. Based on the minimum distance of each point from the centroids. we will assign the points to a cluster.

point.	1st centroid (2, 6)	2nd centroid (5, 10)	3rd centroid (6, 11)	
(2, 10)	4	3	9.13	c ₂
(2, 6)	0	5	6.4	c ₁
(11, 11)	10.29	6.08	5	c ₃
(6, 9)	5	1.41	2	c ₂
(6, 4)	4.47	6.08	7	c ₁
(1, 2)	4.12	8.94	10.2	c ₁
(5, 10)	5	0	1.41	c ₂
(4, 9)	3.6	1.41	2.8	c ₂
(10, 12)	10	5.38	4.1	c ₃
(7, 5)	5.09	5.38	6.08	c ₁
(9, 11)	8.60	4.12	3	c ₃
(4, 6)	2	4.12	5.38	c ₁
(3, 10)	4.12	2	3.16	c ₂
(3, 8)	2.23	2.82	4.29	c ₁
(6, 11)	6.4	1.41	0	c ₃

calculate
new centroid

$$c_1(x) = \left(\frac{2+6+1+7+4+3}{6}, \frac{6+4+2+5+6+8}{6} \right)$$

$$\Rightarrow (3.83, 5.16)$$

$$c_2(x) = \left(\frac{2+6+5+4+3}{5}, \frac{10+9+10+9+10}{5} \right)$$

$$= (4, 9.6)$$

$$c_3(x) = \left(\frac{11+10+9+6}{4}, \frac{11+12+11+11}{4} \right)$$

$$= (9, 11.25)$$

Now that we have calculated new centroid for each cluster, we will calculate the distance of each data point from the new centroids. Then we will assign the points to clusters based on their distance from the centroids.

point	1st centroid (3.83, 5.16)	2nd centroid (4, 9.6)	3rd centroid (9.11, 2.5)	
(2, 10)	5.169	2.040	7.111	C ₂
(2, 6)	2.013	4.118	8.750	C ₁
(11, 11)	9.291	7.139	2.016	C ₃
(6, 9)	4.403	2.088	3.750	C ₂
(6, 4)	2.461	5.946	7.846	C ₁
(1, 2)	4.249	8.171	12.230	C ₁
(5, 10)	4.972	1.077	4.191	C ₂
(4, 9)	3.837	0.600	5.483	C ₂
(10, 12)	9.204	6.462	1.250	C ₃
(7, 5)	3.171	5.492	6.562	C ₁
(9, 11)	7.792	5.192	0.250	C ₃
(4, 6)	0.850	3.600	7.25	C ₁
(3, 10)	4.904	1.077	6.129	C ₂
(3, 8)	2.950	1.897	6.824	C ₂
(6, 11)	6.223	2.441	3.010	C ₂

Now calculate new centroid.

$$C_1(x) = (4.4, 6)$$

$$C_2(x) = (4.193, 9.571)$$

$$C_3(x) = (10, 11.333)$$

the above iteration is continue till our present centroids and previous centroids are similar.

point	1st centroid (4.46)	2nd centroid (4.143, 9.57)	3rd centroid (10, 11, 33)	Assigned cluster
(2, 10)	5.7	2.1	8.1	C ₂
(2, 6)	2.4	4.1	9.6	C ₁
(11, 11)	9.4	7.0	1.0	C ₃
(6, 9)	4.8	1.9	4.6	C ₂
(6, 4)	2.0	5.8	8.3	C ₁
(1, 2)	3.9	8.1	12.9	C ₃
(5, 10)	5.9	0.9	5.1	C ₂
(4, 9)	4.4	0.5	6.4	C ₂
(10, 12)	9.5	6.3	0.6	C ₃
(7, 5)	3.0	5.3	7.8	C ₁
(4, 6)	1.4	3.5	8.0	C ₁
(3, 10)	5.9	1.2	7.1	C ₂
(3, 8)	3.5	1.9	7.7	C ₂
(6, 11)	6.7	2.3	4.8	C ₂
(9, 11)	8.1	5.0	1.0	C ₃

calculate new centroid = ②

$$C_1(x) = (4, 4, 6)$$

$$C_2(x) = (4.143, 9.57)$$

$$C_3(x) = (10, 11, 33)$$

Hence here we see that the previous and present centroid is same hence we stop there our iteration.

points lies in cluster-1:

(2, 6), (6, 4), (1, 2), (7, 5), (4, 6)

points lies in cluster-2:

(2, 10), (6, 9), (5, 10), (4, 9), (3, 8), (6, 11), (3, 10)

points lies in cluster-3:

(1, 11), (10, 12), (3, 11)

Applications of k-mean clustering:

- Document classification
- Customer segmentation
- Cyber profiling
- Image segmentation
- Fraud detection in banking and insurance.

Advantages of k-means clustering Algo:

- easy to implement
- Scalability
- Convergence
- Generalization
- Choice of centroids

Disadvantages of k-means clustering Algo:

- Deciding the number of clusters
- choice of initial centroids
- effects of outliers
- curse of dimensionality