

# Non-hierarchical clustering

## K-means clustering

when you have large amount of dataset hierarchical clustering is not useful (Agglomerative and division)

when we have large amount of data we will use partition clustering (specifying the cluster how many we want) one of the way of partition clustering is K-means clustering

partition =  $K = 3$

$K$  = no. of clustering in dataset  
means = method to identify the cluster

Note: always confirm  $K$  values from domain expert.

# process of K-means clustering

Step 1 :- set K values

For example I have a dataset

$$A = [1, 2, 4, 7, 5]$$

A	B
1	1
2	1
4	5
7	7
5	7

→ So I have dataset which has two variable A and B which have 5 rows

and select K value as 2

$$\underline{\underline{K = 2}}$$

Step 2 : partition of dataset

we can do partition into three ways

- ① user-specified initial partition
- ② user-specified initial centroids
- ③ random partition

① ② is depend of user input  
Note :- mostly we will use random partition.

step 3

From random partition happens.  
firstly it will check the  $k$  values and  
give centroids that here <sup>no.</sup>  $k = \text{no. centroids}$   
but how it give centroid position?

solution

Let say we have co-ordinates of data  
points like  $(1,1)$   $(2,1)$   $(4,5)$   $(5,7)$   $(7,7)$   
now random partition will happen  
let say cluster

$$A = (1,1) (2,1) (4,5)$$

$$B = (5,7) (7,7)$$

and we will take out the centroid position  
of cluster through this formula.

$$X = \frac{X_1 + X_2 + X_3 + X_n}{\text{no. co-ordinates}}$$

$$Y = \frac{Y_1 + Y_2 + Y_3 + Y_n}{\text{no. co-ordinates}}$$

$(X, Y)$  will be centroid position.

Let compute for A cluster

$$x = \frac{(1+2+4)}{3} = 3.33$$

$$y = \frac{(1+1+5)}{3} = 2.33$$

$(3.33, 2.33)$  will be my centroid position.

For B cluster.

$$x = \frac{(7+5)}{2} = \frac{12}{2} = 6$$

$$y = \frac{(7+7)}{2} = \frac{14}{2} = 7$$

$(6, 7)$  will be my centroid position

step 4

NAME	
DATE	/ /

Now compute Euclidean distance of each record from each centroid and reassigns to closest cluster.

$$\frac{A}{\sqrt{(1-2 \cdot 33)^2 + (1-2 \cdot 33)^2}} = 1.89$$

	A	B
similarity		
1.37		7.81
3.14		7.21
6.60		2.93
5.37		

So whoever in A cluster is higher values than ~~B~~ B cluster will move towards B cluster and B cluster value will move. So A cluster value so in this example 3.14 is less than 2.93 so we will move ~~close~~ it.

So we have move the datapoint into A cluster and find the new centroid values for each cluster and this cluster goes on until we don't reach perfect values which don't overlap each other.

How to select K values.

Step 1:- get the mean values of each cluster.

Step 2:- take the average distance in cluster by giving different K values.

Now plot elbow chart

Step 3:- whoever got the least values of Average distance for example 6 cluster ~~was~~

Step 4:- Now take out the mean point 6 cluster is 3. the 3 values will be our best cluster.

A best cluster

because  
it not rising  
too much  
No decline or  
flatter too  
much.

Y axis

Elbow chart.