

# K-means algorithm

## Problem

For the given dataset, consider the number of cluster as 3. Find the datapoints in these 3 clusters using K-means algorithm.

$$(x, y) = \{(2, 4), (2, 6), (5, 6), (4, 7), (8, 3), (6, 6), (5, 2), (5, 7), (6, 3), (4, 4)\}$$

Use Euclidean distance as a distance function.

Iteration 01: Centroid point } Sometimes let's centroid be given

Centroid  $C_1 \rightarrow (1, 5)$

Centroid  $C_2 \rightarrow (4, 1)$

Centroid  $C_3 \rightarrow (8, 4)$

$C_1, C_2, C_3 \rightarrow \text{clusters}$

Euclidean distance

$$D = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

• Calculate the distance of each point from each of the centroid of three clusters.

Given data points		Distance to Centroid			cluster Number
x	y	<u>C<sub>1</sub></u> (1, 5)	<u>C<sub>2</sub></u> (4, 1)	<u>C<sub>3</sub></u> (8, 4)	
1	2, 4	1.41	3.61	6.00	C <sub>1</sub>
2	2, 6	1.41	5.39	6.32	C <sub>1</sub>
3	5, 6	4.12	5.10	3.61	C <sub>3</sub>
4	4, 7	3.61	6.00	5.00	C <sub>1</sub>
5	8, 3	7.28	4.47	1.00	C <sub>3</sub>
6	6, 6	5.10	5.39	2.83	C <sub>3</sub>
7	5, 2	5.00	1.41	3.61	C <sub>2</sub>
8	5, 7	4.47	6.08	4.24	C <sub>3</sub>
9	6, 3	5.39	2.83	2.24	C <sub>3</sub>
10	4, 4	3.16	3.00	4.00	C <sub>2</sub>



clusters after iteration 1 are

cluster 1  $\rightarrow (2,4), (2,6), (1,7)$

cluster 2  $\rightarrow (5,2), (4,4)$

cluster 3  $\rightarrow (5,6), (8,3), (6,1), (5,7), (6,3)$

Now recompute the new cluster by determining the new centroid.  
 $\rightarrow$  Centroid can be computed by taking the mean of data points in that cluster.

Hence, Iteration 2

Centroid  $\rightarrow C_1 = (2.66, 5.66)$

Centroid  $\rightarrow C_2 = (4.5, 3)$

Centroid  $\rightarrow C_3 = (6, 3)$

$$\frac{2+2+4}{3} = \frac{10}{3} = 2.66$$

$$\frac{1+6+7}{3} = \frac{14}{3} = 5.66$$

Now determine the distance of all data points from new centroid.

x, y	distance to			cluster number
	$(2.66, 5.66)$ $C_1$	$(4.5, 3)$ $C_2$	$(6, 3)$ $C_3$	
2, 4	1.79	2.68	4.12	$C_1$
2, 6	0.74	3.31	4.12	$C_1$
5, 6	2.36	2.04	1.41	$C_3$
1, 7	1.90	1.93	2.83	$C_1$
8, 3	5.87	3.5	2.83	$C_3 \rightarrow$
6, 6	3.56	3.35	1	$C_3$
5, 2	4.54	1.72	3.16	$C_2$
5, 7	2.70	2.13	2.24	$C_3$
6, 3	1.27	1.5	2	$C_2$
4, 1	2.13	1.72	2.14	$C_2$

Observe data point (6, 3) is changed moved from cluster  $C_3$  to  $C_2$ . Hence further centroid need to be calculated.

$C_1 \rightarrow (2,4), (2,6), (1,7)$

$C_2 \rightarrow (5,2), (6,3), (4,4)$

$C_3 \rightarrow (5,6), (8,3), (6,1), (5,7)$





2

### Iteration 3

How centroids will be

Centroid  $c_1 \rightarrow (2.66, 5.66)$

Centroid  $c_2 \rightarrow (5, 3)$

Centroid  $c_3 \rightarrow (6.5, 5)$

					cluster Number
2	4	1.79	2.16	4.12	$c_1$
2	6	6.74	4.24	4.03	$c_1$
5	6	2.38	3.00	1.12	$c_3$
4	7	1.90	4.12	2.50	$c_1$
✓ 8	3	5.97	3.00	3.30	(2)
6	6	3.36	3.11	0.50	$c_3$
5	2	4.34	1.00	3.64	$c_2$
9	7	2.70	4.00	1.80	$c_3$
6	3	4.27	1.00	2.50	$c_2$
4	4	2.13	1.41	2.50	$c_2$

Data point (8,3) is moved from cluster  $c_3$  to  $c_2$ . hence again new centroid need to be calculated.



### Iteration 4:

Centroid  $\rightarrow c_1 = (2.66, 5.66)$

Centroid  $\rightarrow c_2 = (5.25, 3)$

Centroid  $\rightarrow c_3 = (5.33, 6.43)$

x	y	(2.66, 5.66)	(5.25, 3)	(5.33, 6.43)	Cluster Members
2	4	1.74	9.88	2.06	C1
2	6	0.74	4.80	4.385	C1
5	6	2.31	3.09	0.47	C3
4	7	1.90	4.37	4.58	C3
8	3	5.92	3.06	19.4	C2
6	6	3.36	3.01	0.950	C1
5	2	4.34	1.05	4.34	C2
5	7	2.70	4.07	4.8009	C3
6	3	4.26	0.25	3.40	C2
4	4	2.13	2.02	2.68	C2

data point  $(4, 4)$  is moved from  $c_1$  to  $c_2$ .  
 $\therefore$  New centroid need to be recalculated.

### Iteration 5

Centroid  $\rightarrow c_1 = (2.5)$  , Centroid  $c_2 = (5.75, 3)$

Centroid  $\rightarrow c_3 = (5.64)$

x	y	(2.5)	(5.75, 3)	(5.64)	Cluster Members
2	4	1.00	3.88	3.91	C1
2	6	1.00	4.80	3.04	C1
5	6	3.18	3.09	0.50	C3
4	7	2.83	4.37	1.12	C3
8	3	6.32	2.25	4.51	C2
6	6	4.12	3.01	1.12	C3
5	2	4.24	1.25	6.50	C2
5	7	3.61	4.07	0.90	C3
6	3	4.47	0.25	3.64	C2
4	4	2.24	2.02	2.68	C2

Now there is no data point which is moving from one cluster to another.

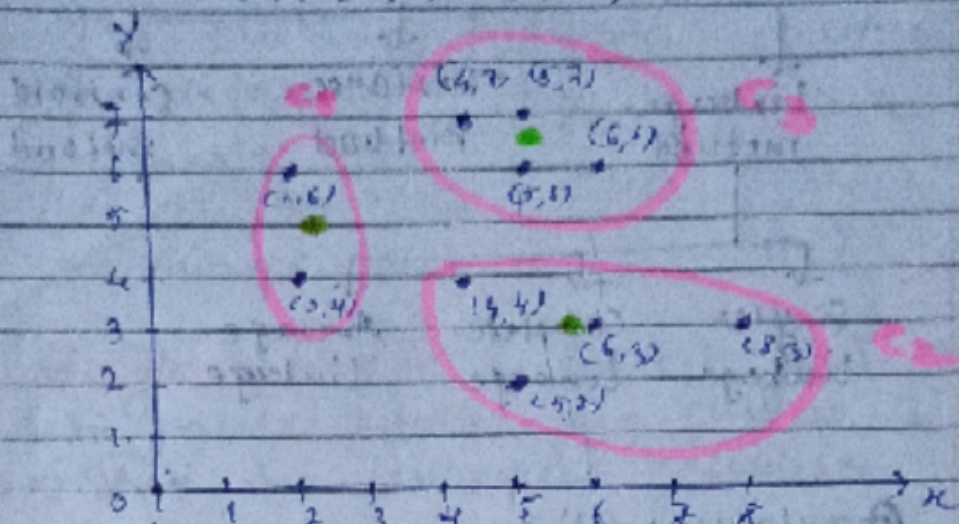


These are the final centroids of clusters where datapoints are

$$c_1 \rightarrow (2, 4), (2, 6)$$

$$c_2 \rightarrow (8, 3), (5, 2), (6, 3), (4, 4)$$

$$c_3 \rightarrow (4, 1), (4, 7), (6, 1), (5, 7)$$



~~K-Means~~

**K-Medoids**

**Manhattan distance**

between  $(x_1, y_1)$  and  $(x_2, y_2)$  is

$$d(P, Q) = ||P - Q||$$

$$= |x_1 - x_2| + |y_1 - y_2|$$

→ It is very similar to K-means algorithm

→ Like K-means and K-medoids are partition based method and try to minimize the distance between data points and centroids

→ In K-medoids - data points are chosen as centroid and use Manhattan distance to calculate the distance.

→ It is more robust to noise and outliers as compared to K-means.