

Example

Problem: Cluster the following eight points (with (x, y) representing locations) into three clusters A1(2, 10) A2(2, 5) A3(8, 4) A4(5, 8) A5(7, 5) A6(6, 4) A7(1, 2) A8(4, 9). Initial cluster centers are: A1(2, 10), A4(5, 8) and A7(1, 2). The distance function between two points $a=(x1, y1)$ and $b=(x2, y2)$ is defined as: $\rho(a, b) = |x2 - x1| + |y2 - y1|$.

Use k-means algorithm to find the three cluster centers after the second iteration.

Solution:

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean 1	Dist Mean 2	Dist Mean 3	Cluster
A1	(2, 10)				
A2	(2, 5)				
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

First we list all points in the first column of the table above. The initial cluster centers – means, are (2, 10), (5, 8) and (1, 2) - chosen randomly. Next, we will calculate the distance from the first point (2, 10) to each of the three means, by using the distance function:

point	mean1
$x1, y1$	$x2, y2$
(2, 10)	(2, 10)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned}
 \rho(\text{point}, \text{mean1}) &= |x2 - x1| + |y2 - y1| \\
 &= |2 - 2| + |10 - 10| \\
 &= 0 + 0 \\
 &= 0
 \end{aligned}$$

point mean2
 $x1, y1$ $x2, y2$
(2, 10) (5, 8)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned} \rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |5 - 2| + |8 - 10| \\ &= 3 + 2 \\ &= 5 \end{aligned}$$

point mean3
 $x1, y1$ $x2, y2$
(2, 10) (1, 2)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned} \rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |1 - 2| + |2 - 10| \\ &= 1 + 8 \\ &= 9 \end{aligned}$$

So, we fill in these values in the table:

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean 1	Dist Mean 2	Dist Mean 3	Cluster
A1	(2, 10)	0	5	9	1
A2	(2, 5)				
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

So, which cluster should the point (2, 10) be placed in? The one, where the point has the shortest distance to the mean – that is mean 1 (cluster 1), since the distance is 0.

Cluster 1
(2, 10)

Cluster 2

Cluster 3

So, we go to the second point (2, 5) and we will calculate the distance to each of the three means, by using the distance function:

point	mean1
$x1, y1$	$x2, y2$
(2, 5)	(2, 10)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned}\rho(\text{point}, \text{mean1}) &= |x2 - x1| + |y2 - y1| \\ &= |2 - 2| + |10 - 5| \\ &= 0 + 5 \\ &= 5\end{aligned}$$

point	mean2
$x1, y1$	$x2, y2$
(2, 5)	(5, 8)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned}\rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |5 - 2| + |8 - 5| \\ &= 3 + 3 \\ &= 6\end{aligned}$$

point	mean3
$x1, y1$	$x2, y2$
(2, 5)	(1, 2)

$$\rho(a, b) = |x2 - x1| + |y2 - y1|$$

$$\begin{aligned}\rho(\text{point}, \text{mean2}) &= |x2 - x1| + |y2 - y1| \\ &= |1 - 2| + |2 - 5| \\ &= 1 + 3\end{aligned}$$

$$= 4$$

So, we fill in these values in the table:

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean 1	Dist Mean 2	Dist Mean 3	Cluster
A1	(2, 10)	0	5	9	1
A2	(2, 5)	5	6	4	3
A3	(8, 4)				
A4	(5, 8)				
A5	(7, 5)				
A6	(6, 4)				
A7	(1, 2)				
A8	(4, 9)				

So, which cluster should the point (2, 5) be placed in? The one, where the point has the shortest distance to the mean – that is mean 3 (cluster 3), since the distance is 0.

Cluster 1
(2, 10)

Cluster 2

Cluster 3
(2, 5)

Analogically, we fill in the rest of the table, and place each point in one of the clusters:

Iteration 1

		(2, 10)	(5, 8)	(1, 2)	
	Point	Dist Mean 1	Dist Mean 2	Dist Mean 3	Cluster
A1	(2, 10)	0	5	9	1
A2	(2, 5)	5	6	4	3
A3	(8, 4)	12	7	9	2
A4	(5, 8)	5	0	10	2
A5	(7, 5)	10	5	9	2
A6	(6, 4)	10	5	7	2
A7	(1, 2)	9	10	0	3
A8	(4, 9)	3	2	10	2

Cluster 1
(2, 10)

Cluster 2
(8, 4)
(5, 8)
(7, 5)
(6, 4)
(4, 9)

Cluster 3
(2, 5)
(1, 2)

Next, we need to re-compute the new cluster centers (means). We do so, by taking the mean of all points in each cluster.

For Cluster 1, we only have one point A1(2, 10), which was the old mean, so the cluster center remains the same.

For Cluster 2, we have $((8+5+7+6+4)/5, (4+8+5+4+9)/5) = (6, 6)$

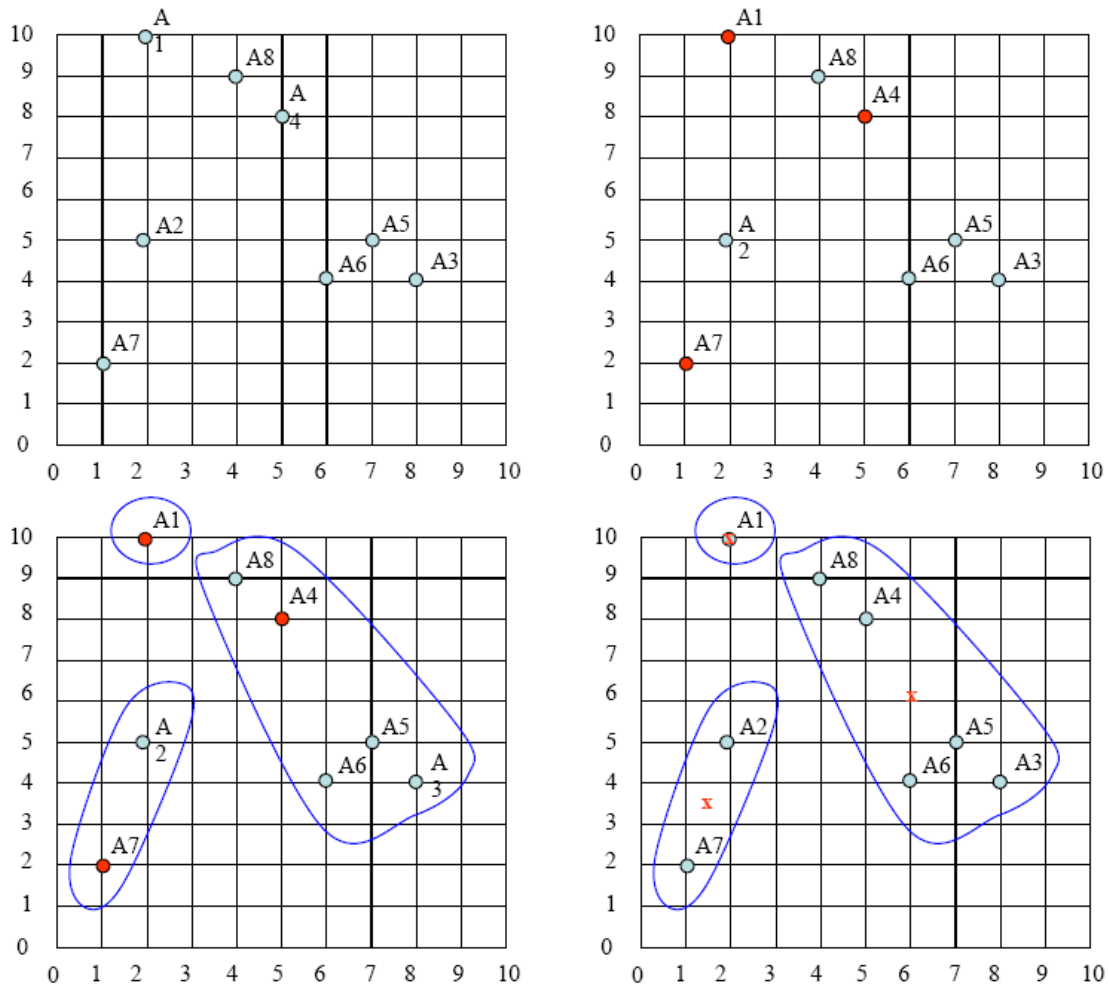
For Cluster 3, we have $((2+1)/2, (5+2)/2) = (1.5, 3.5)$

new clusters: 1: {A1}, 2: {A3, A4, A5, A6, A8}, 3: {A2, A7}

b) centers of the new clusters:

C1= (2, 10), C2= $((8+5+7+6+4)/5, (4+8+5+4+9)/5) = (6, 6)$, C3= $((2+1)/2, (5+2)/2) = (1.5, 3.5)$

c)



The initial cluster centers are shown in red dot. The new cluster centers are shown in red x.

That was Iteration1 (epoch1). Next, we go to Iteration2 (epoch2), Iteration3, and so on until the means do not change anymore.

In Iteration2, we basically repeat the process from Iteration1 this time using the new means we computed.

d)

We would need two more epochs. After the 2nd epoch the results would be:

1: {A1, A8}, 2: {A3, A4, A5, A6}, 3: {A2, A7}

with centers $C1=(3, 9.5)$, $C2=(6.5, 5.25)$ and $C3=(1.5, 3.5)$.

After the 3rd epoch, the results would be:

1: {A1, A4, A8}, 2: {A3, A5, A6}, 3: {A2, A7}

with centers $C1=(3.66, 9)$, $C2=(7, 4.33)$ and $C3=(1.5, 3.5)$.

