

## K-Means Clustering - Manual Mathematical Calculation

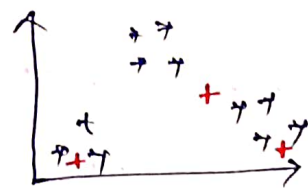
- K Means is an unsupervised ML algo used for clustering.
- Goal is to categorize the data into K-no of clusters.
- It begins by selecting K-points randomly as the initial centroids then assign each data points to the nearest centroid.
- After every point is assigned to a centroid, the centroids are recalculated by taking the mean of all the datapoints that were assigned to each centroid then the datapoints are reassigned to the nearest centroid based on the updated one.
- This process repeats until the cluster no longer change.

(or)

Step 1: Choose the number of clusters (K) you wish to put the data into, let us choose 3 (K).

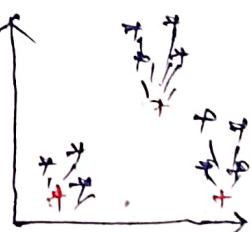


Step 2: Randomly initialize K points, these are called centroids.



Step 3: Identify the points closest to each centroid.

Step 4: Calculate the mean of the points in each cluster and move each centroid to that mean point.



Step 5: Repeat Step 3 & 4 until the centroid value changes the same.

# K Means Clustering - Solved Example

- Datapoints:

$A_1(2, 10)$ ,  $A_2(2, 5)$ ,  $A_3(8, 4)$ ,  $B_1(5, 8)$ ,  $B_2(7, 5)$ ,  $B_3(6, 4)$   
 $C_1(1, 2)$ ,  $C_2(4, 9)$

- Distance function: Euclidean distance

- Suppose initially we assign  $A_1, B_1$  and  $C_1$  as the center of each cluster respectively.

Iteration 1:

Iteration 1		Distance to				Cluster	New Cluster	
Datapoints		$x_1$ $y_1$	$x_2$ $y_2$	1	2			
Dist to Centroid	A1	2	10	0	3.61	8.06	1	
	A2	2	5	5	4.24	3.16	3	
	A3	8	4	8.49	5.00	7.28	2	
	B1	5	8	3.61	0.00	7.21	2	
	B2	7	5	7.07	3.61	6.71	2	
	B3	6	4	7.21	4.12	5.39	2	
	C1	1	2	8.06	7.21	0.00	3	
	C2	4	9	2.24	1.41	7.62	2	

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

Sample calculation

$$\text{Euclidean distance } d(P_1, P_2) = \sqrt{(2-2)^2 + (10-10)^2} = 0$$

New Centroid Calculation for cluster

1st cluster  $A_1(2, 10) - (2, 10)$

2nd cluster

$$\rightarrow \frac{8+5+7+6+4}{5}$$

$$\rightarrow \frac{4+8+5+4+9}{5} = (6, 6)$$

3rd cluster

$$\rightarrow \frac{(2+1)}{2} = 1.5$$

$$\frac{(5+2)}{2} = 3.5$$

$$(1.5, 3.5)$$

$$d_2(P_1, P_2) = \sqrt{(2-2)^2 + (10-5)^2} = \sqrt{(5)^2} = 5$$

$$d_3(P_1, P_2) = \sqrt{(2-8)^2 + (10-4)^2} = \sqrt{(6)^2 + (6)^2} = 8.49$$

$$d_4(P_1, P_2) = \sqrt{(2-5)^2 + (10-8)^2} = \sqrt{(3)^2 + (2)^2} = \sqrt{13} = 3.61$$

$$d_5(P_1, P_2) = \sqrt{(2-7)^2 + (10-5)^2} = \sqrt{(5)^2 + (5)^2} = 7.07$$

$$d_6 = 7.21, d_7 = 8.06, d_8 = 2.24$$

Current Centroids

$A_1(2, 10)$   $B_1(6, 6)$   $C_1(1.5, 3.5)$

Distance to

Data points			2 10	6 6	1.5 3.5	Cluster	New cluster
A1	2	10	0.0	5.66	6.52	1	1
A2	2	5	5.0	4.12	1.58	3	3
A3	8	4	8.49	2.83	6.52	2	2
B1	5	8	3.61	2.84	5.70	2	2
B2	7	5	7.07	1.41	5.70	2	2
B3	6	4	7.21	2.00	4.53	2	2
C1	1	2	8.06	6.40	1.58	3	3
C2	4	9	2.24	3.61	6.04	2	1

New Centroids

$C(3, 9.5)$   $(6.5, 5.25)$   $(1.5, 3.5)$

Distance to

Data points		3 9.5	6.5 5.25	1.5 3.5	Cluster	New cluster
A1	2 10	1.12	6.54	6.52	1	1
A2	2 5	4.61	4.51	1.58	3	3
A3	8 4	7.42	1.95	6.52	2	2
B1	5 8	2.50	3.13	5.70	2	1
B2	7 5	6.02	0.56	5.70	2	2
B3	6 4	6.26	1.35	4.53	2	2
C1	1 2	7.76	6.39	1.58	3	3
C2	4 9	1.12	4.51	6.04	1	1

New Centroids.

$A1(3.67, 9)$   $B1(7, 4.33)$   $C1(1.5, 3.5)$   
Distance to

Datapoints			3.67, 9	7	4.33	1.5	3.5	Cluster	New Cluster
A1	2	10	1.94	7.56	6.52	1		1	1
A2	2	5	4.33	5.04	1.58	3		3	3
A3	8	4	6.62	1.05	6.52	2		2	2
B1	5	8	1.67	4.18	5.70	1		1	1
B2	7	5	5.21	0.67	5.70	2		2	2
B3	6	4	5.82	1.05	4.58	2		2	2
C1	1	2	7.49	6.44	1.58	3		3	3
C2	4	9	0.33	5.55	6.04	1		1	1

Cluster = New Cluster  
Stop the iteration.

$A1, B1, C2 \rightarrow$  belongs to  $\rightarrow$  cluster 1  
 $A3, B2, B3 \rightarrow$  belongs to  $\rightarrow$  cluster 2  
 $A2, C1 \rightarrow$  belongs to  $\rightarrow$  cluster 3