



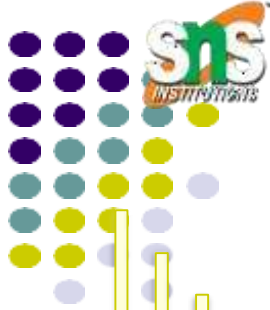
SNS COLLEGE OF TECHNOLOGY

Coimbatore-35

An Autonomous Institution

Accredited by NBA – AICTE and Accredited by NAAC – UGC with ‘A++’ Grade

Approved by AICTE, New Delhi & Affiliated to Anna University, Chennai



DEPARTMENT OF COMPUTER APPLICATIONS

23CAT702 – MACHINE LEARNING

II YEAR III SEM

UNIT III – DISTANCE-BASED MODELS

TOPIC 20 – Clustering around Medoids



What is k-Medoid



K-Medoids (also called Partitioning Around Medoid) algorithm was proposed in 1987 by Kaufman and Rousseeuw. A medoid can be defined as a point in the cluster, whose dissimilarities with all the other points in the cluster are minimum. The dissimilarity of the medoid(C_i) and object(P_i) is calculated by using $E = |P_i - C_i|$

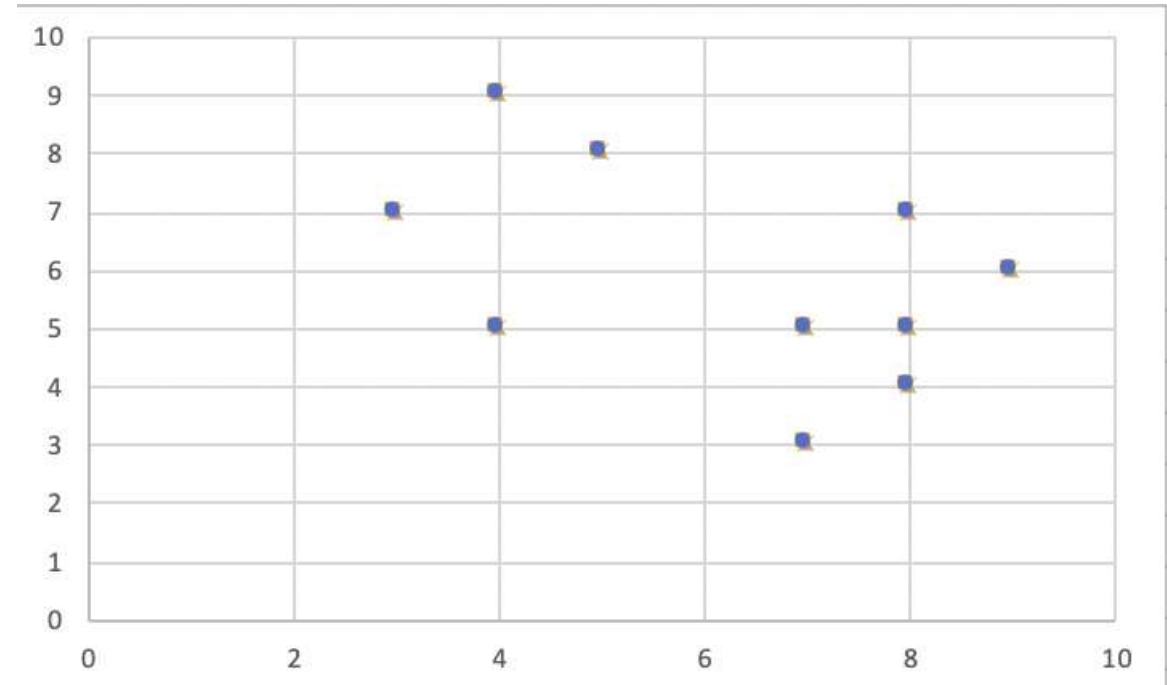
Algorithm

1. Initialize: select k random points out of the n data points as the medoids.
2. Associate each data point to the closest medoid by using any common distance metric methods.
3. While the cost decreases: For each medoid m , for each data o point which is not a medoid:
 1. Swap m and o , associate each data point to the closest medoid, and recompute the cost.
 2. If the total cost is more than that in the previous step, undo the swap.



Example

	X	Y
0	8	7
1	3	7
2	4	9
3	9	6
4	8	5
5	5	8
6	7	3
7	8	4
8	7	5
9	4	5



Step 1: Let the randomly selected 2 medoids, so select $k = 2$, and let **C1** $-(4, 5)$ and **C2** $-(8, 5)$ are the two medoids.



Example...



Step 2: Calculating cost. The dissimilarity of each non-medoid point with the medoids is calculated and tabulated:

	X	Y	Dissimilarity from C1	Dissimilarity from C2
0	8	7	6	2
1	3	7	3	7
2	4	9	4	8
3	9	6	6	2
4	8	5	-	-
5	5	8	4	6
6	7	3	5	3
7	8	4	5	1
8	7	5	3	1
9	4	5	-	-

Here we have used Manhattan distance formula to calculate the distance matrices between medoid and non-medoid points. That formula tell that **Distance = |X1-X2| + |Y1-Y2|**.

Each point is assigned to the cluster of that medoid whose dissimilarity is less. Points 1, 2, and 5 go to cluster C1 and 0, 3, 6, 7, 8 go to cluster C2. The Cost = (3 + 4 + 4) + (3 + 1 + 1 + 2 + 2) = 20



Example...



Step 3: randomly select one non-medoid point and recalculate the cost. Let the randomly selected point be (8, 4). The dissimilarity of each non-medoid point with the medoids – C1 (4, 5) and C2 (8, 4) is calculated and tabulated.

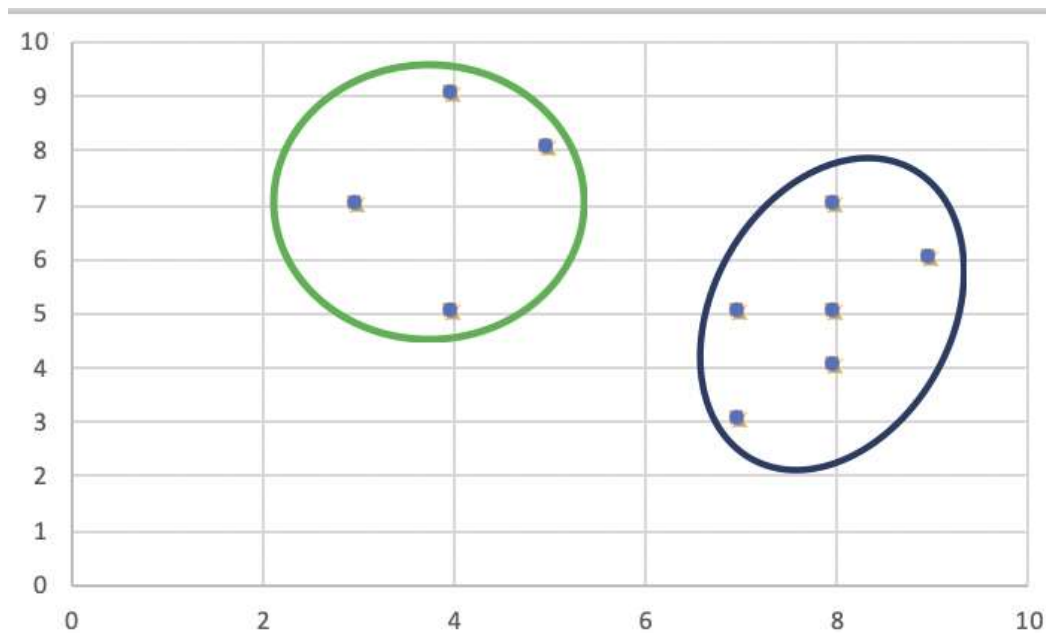
	X	Y	Dissimilarity from C1	Dissimilarity from C2
0	8	7	6	3
1	3	7	3	8
2	4	9	4	9
3	9	6	6	3
4	8	5	4	1
5	5	8	4	7
6	7	3	5	2
7	8	4	-	-
8	7	5	3	2
9	4	5	-	-



Example...



Each point is assigned to that cluster whose dissimilarity is less. So, points 1, 2, and 5 go to cluster C1 and 0, 3, 6, 7, 8 go to cluster C2. The New cost = $(3 + 4 + 4) + (2 + 2 + 1 + 3 + 3) = 22$ Swap Cost = New Cost – Previous Cost = $22 - 20$ and $2 > 0$ As the swap cost is not less than zero, we undo the swap. Hence (4, 5) and (8, 5) are the final medoids.





Advantages:

- It is simple to understand and easy to implement.
- K-Medoid Algorithm is fast and converges in a fixed number of steps.
- PAM is less sensitive to outliers than other partitioning algorithms.





References

1. [Tutorial](#) - Tutorial with introduction of Clustering Algorithms (k-means, fuzzy-c-means, hierarchical, mixture of gaussians) + some interactive demos (java applets).
2. Digital Image Processing and Analysis-byB.Chanda and D.Dutta Majumdar.
3. H. Zha, C. Ding, M. Gu, X. He and H.D. Simon. "Spectral Relaxation for K-means Clustering", Neural Information Processing Systems vol.14 (NIPS 2001). pp. 1057-1064, Vancouver, Canada. Dec. 2001.
4. J. A. Hartigan (1975) "Clustering Algorithms". Wiley.
5. J. A. Hartigan and M. A. Wong (1979) "A K-Means Clustering Algorithm", Applied Statistics, Vol. 28, No. 1, p100-108.
6. [D. Arthur](#), [S. Vassilvitskii](#) (2006): "How Slow is the k-means Method?,"
7. D. Arthur, S. Vassilvitskii: "[k-means++ The Advantages of Careful Seeding](#)" 2007 Symposium on Discrete Algorithms (SODA).
8. www.wikipedia.com



Thank You