

k-medoids

The **k-medoids algorithm** is a **clustering algorithm** related to the **k-means** algorithm and the medoidshift algorithm. Both the k-means and k-medoids algorithms are partitional (breaking the dataset up into groups) and both attempt to minimize the distance between points labeled to be in a cluster and a point designated as the center of that cluster. In contrast to the k-means algorithm, k-medoids chooses datapoints as centers (**medoids** or exemplars) and works with an arbitrary matrix of distances between datapoints instead of l_2 . This method was proposed in 1987^[1] for the work with l_1 norm and other distances.

k-medoid is a classical partitioning technique of clustering that clusters the data set of n objects into k clusters known *a priori*. A useful tool for determining k is the **silhouette**.

It is more robust to noise and outliers as compared to **k-means** because it minimizes a sum of pairwise dissimilarities instead of a sum of squared Euclidean distances.

A **medoid** can be defined as the object of a cluster whose average dissimilarity to all the objects in the cluster is minimal. i.e. it is a most centrally located point in the cluster.

1 Algorithms

The most common realisation of k-medoid clustering is the **Partitioning Around Medoids (PAM)** algorithm. PAM uses a greedy search which may not find the optimum solution, but it is faster than exhaustive search. It works as follows:^[2]

1. Initialize: randomly select (without replacement) k of the n data points as the medoids
2. Associate each data point to the closest medoid.
3. While the cost of the configuration decreases:
 - (a) For each medoid m , for each non-medoid data point o :
 - i. Swap m and o , recompute the cost (sum of distances of points to their medoid)
 - ii. If the total cost of the configuration increased in the previous step, undo the swap

Other algorithms than PAM have been suggested in the literature, including the following **Voronoi iteration** method:^[3]

1. Select initial medoids
2. Iterate while the cost decreases:
 - (a) In each cluster, make the point that minimizes the sum of distances within the cluster the medoid
 - (b) Reassign each point to the cluster defined by the closest medoid determined in the previous step.

2 Demonstration of PAM

Cluster the following data set of ten objects into two clusters i.e. $k = 2$.

Consider a data set of ten objects as follows:

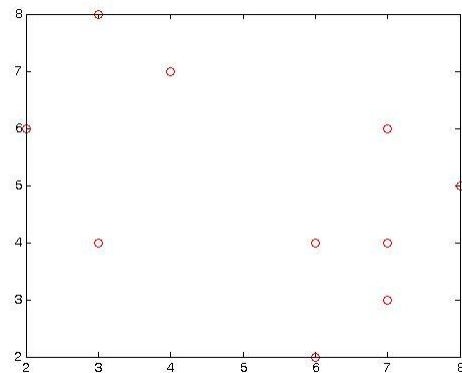


Figure 1.1 – distribution of the data

2.1 Step 1

Initialize k centers.

Let us assume x_2 and x_8 are selected as medoids, so the centers are $c_1 = (3, 4)$ and $c_2 = (7, 4)$

Calculate distances to each center so as to associate each data object to its nearest medoid. Cost is calculated using **Manhattan distance** (Minkowski distance metric with $r = 1$). Costs to the nearest medoid are shown bold in the table.

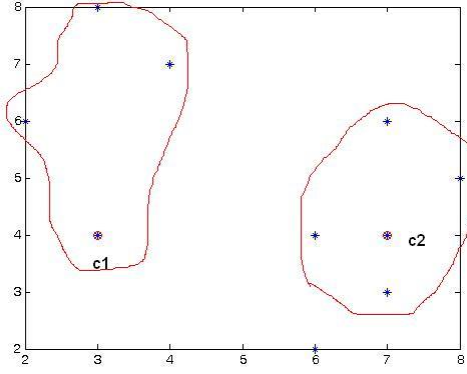


Figure 1.2 – clusters after step 1

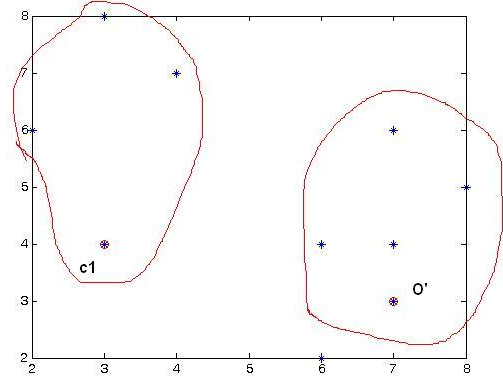


Figure 1.3 – clusters after step 2

Then the clusters become:

$$\text{Cluster}_1 = \{(3,4)(2,6)(3,8)(4,7)\}$$

$$\text{Cluster}_2 = \{(7,4)(6,2)(6,4)(7,3)(8,5)(7,6)\}$$

Since the points (2,6) (3,8) and (4,7) are closer to c_1 hence they form one cluster whilst remaining points form another cluster.

So the total cost involved is 20.

Where cost between any two points is found using formula

$$\text{cost}(x, c) = \sum_{i=1}^d |x_i - c_i|$$

where x is any data object, c is the medoid, and d is the dimension of the object which in this case is 2.

Total cost is the summation of the cost of data object from its medoid in its cluster so here:

$$\begin{aligned} \text{total cost} &= \{\text{cost}((3, 4), (2, 6)) + \text{cost}((3, 4), (3, 8)) + \text{cost}((3, 4), (4, 7))\} \\ &\quad + \{\text{cost}((7, 4), (6, 2)) + \text{cost}((7, 4), (6, 4)) + \text{cost}((7, 4), (7, 3)) \\ &\quad + \text{cost}((7, 4), (8, 5)) + \text{cost}((7, 4), (7, 6))\} \\ &= (3 + 4 + 4) + (3 + 1 + 1 + 2 + 2) \\ &= 20 \end{aligned}$$

2.2 Step 2

Select one of the nonmedoids O'

Let us assume $O' = (7,3)$, i.e. x_7 .

So now the medoids are $c_1(3,4)$ and $O'(7,3)$

If c_1 and O' are new medoids, calculate the total cost involved

By using the formula in the step 1

$$\begin{aligned} \text{total cost} &= 3 + 4 + 4 + 2 + 2 + 1 + 3 + 3 \\ &= 22 \end{aligned}$$

So cost of swapping medoid from c_2 to O' is

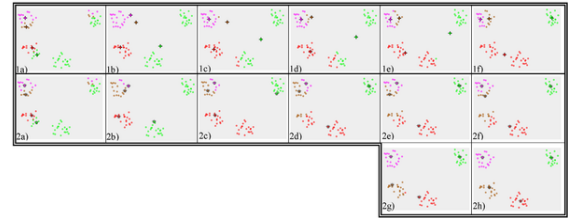


Figure 2. *K-medoids versus k-means.* Figs 2.1a-2.1f present a typical example of the *k-means* convergence to a local minimum. This result of *k-means* clustering contradicts the obvious cluster structure of data set. In this example, *k-medoids* algorithm (Figs 2.2a-2.2h) with the same initial position of medoids (Fig. 2.2a) converges to the obvious cluster structure. The small circles are data points, the four ray stars are centroids (means), the nine ray stars are medoids.^[4]

So moving to O' would be a bad idea, so the previous choice was good. So we try other nonmedoids and found that our first choice was the best. So the configuration does not change and algorithm terminates here (i.e. there is no change in the medoids).

It may happen some data points may shift from one cluster to another cluster depending upon their closeness to medoid.

In some standard situations, *k-medoids* demonstrate better performance than *k-means*. An example is presented in Fig. 2. The most time-consuming part of the *k-medoids* algorithm is the calculation of the distances between objects. If a quadratic preprocessing and storage is applicable, the distances matrix can be precomputed to achieve consequent speed-up. See for example,^[3] where the authors also introduce a heuristic to choose the initial *k* medoids.

3 Software

- **ELKI** includes several k-means variants, including an EM-based k-medoids and the original PAM algorithm.
- **Julia** contains a k-medoid implementation in the JuliaStats clustering package.
- **R** includes variants of k-means in the “flexclust” package and PAM is implemented in the “cluster” package.
- **RapidMiner** has an operator named KMedoids, but it does *not* implement the KMedoids algorithm correctly. Instead, it is a k-means variant, that substitutes the mean with the closest data point (which is not the medoid).
- **MATLAB** implements PAM, CLARA, and two other algorithms to solve the k-medoid clustering problem.

4 References

- [1] Kaufman, L. and Rousseeuw, P.J. (1987), Clustering by means of Medoids, in Statistical Data Analysis Based on the L_1 –Norm and Related Methods, edited by Y. Dodge, North-Holland, 405–416.
- [2] Sergios Theodoridis & Konstantinos Koutroumbas (2006). *Pattern Recognition 3rd ed.* p. 635.
- [3] H.S. Park , C.H. Jun, A simple and fast algorithm for K-medoids clustering, Expert Systems with Applications, 36, (2) (2009), 3336–3341.
- [4] The illustration was prepared with the Java applet, E.M. Mirkes, K-means and K-medoids: applet. University of Leicester, 2011.

5 Text and image sources, contributors, and licenses

5.1 Text

- **K-medoids** *Source:* <https://en.wikipedia.org/wiki/K-medoids?oldid=703337360> *Contributors:* Dina, Giftlite, WorldsApart, Andreas Kaufmann, 3mta3, Tommytao, Mandarax, Qwertyus, Oliekirk, SmackBot, Mauls, Cronholm144, Jonnat, Anthony Bradbury, Alai-bot, Narayanese, Talgalili, Magioladitis, Subail, Justin wen, Shrikantnangare, Turketwh, AlleborgoBot, Vkotor, Tosanjay, Xinunus, BobKawanaka, Bradpitcher~enwiki, Agor153, XLinkBot, Addbot, Yobot, AnomieBOT, Surturpain, Kxx, FrescoBot, Dan Golding, Is-mailari, EmausBot, Combee123, Sgoder, AvicAWB, Chire, Donner60, Pokbot, ClueBot NG, Cheater no1, Gongzhitao, Chimpooop, BG19bot, Wesamaaa, Eracle.adeluca, Khirodkantnaik, HelpUsStopSpam and Anonymous: 63

5.2 Images

- **File:Edit-clear.svg** *Source:* <https://upload.wikimedia.org/wikipedia/en/f/f2/Edit-clear.svg> *License:* Public domain *Contributors:* The Tango! Desktop Project. *Original artist:* The people from the Tango! project. And according to the meta-data in the file, specifically: “Andreas Nilsson, and Jakub Steiner (although minimally).”
- **File:K-means_versus_k-medoids.png** *Source:* https://upload.wikimedia.org/wikipedia/commons/b/b3/K-means_versus_k-medoids.png *License:* CC BY-SA 3.0 *Contributors:* Own work *Original artist:* Agor153
- **File:Kmedoid1.jpg** *Source:* <https://upload.wikimedia.org/wikipedia/commons/e/e1/Kmedoid1.jpg> *License:* Public domain *Contributors:* en:image:Kmedoid1.jpg *Original artist:* en>User:Shrikantnangare
- **File:Kmedoid2.jpg** *Source:* <https://upload.wikimedia.org/wikipedia/commons/8/82/Kmedoid2.jpg> *License:* Public domain *Contributors:* en:image:Kmedoid2.jpg *Original artist:* en>User:Shrikantnangare
- **File:Kmedoid3.jpg** *Source:* <https://upload.wikimedia.org/wikipedia/commons/a/a4/Kmedoid3.jpg> *License:* Public domain *Contributors:* en:image:Kmedoid3.jpg *Original artist:* en>User:Shrikantnangare

5.3 Content license

- Creative Commons Attribution-Share Alike 3.0