

MCOKE: Multi-Cluster Overlapping K-Means Extension Algorithm

Said Baadel, Fadi Thabtah, Joan Lu

Abstract—Clustering involves the partitioning of n objects into k clusters. Many clustering algorithms use hard-partitioning techniques where each object is assigned to one cluster. In this paper we propose an overlapping algorithm MCOKE which allows objects to belong to one or more clusters. The algorithm is different from fuzzy clustering techniques because objects that overlap are assigned a membership value of 1 (one) as opposed to a fuzzy membership degree. The algorithm is also different from other overlapping algorithms that require a similarity threshold be defined *a priori* which can be difficult to determine by novice users.

Keywords—Data mining, k-means, MCOKE, overlapping.

I. INTRODUCTION

DATA clustering, also known as unsupervised classification is a research field widely studied in data mining and machine learning domains due to its applications to segmentation, summarization, learning, and target marketing [1]. Clustering involves the partitioning of a set of objects or data into clusters or subsets such that the objects or data in each subset or cluster contains similar traits based on measured similarities [2] and data from different clusters are dissimilar [13]. There are many techniques that have been explored for clustering processes such as distance-based, probabilistic, and density/grid-based etc. with the distance based being very popular in research fields. Two major types of distance-based algorithms are hierarchical and partitioning. The first type typically represents clusters hierarchically through a dendrogram and using a similarity criterion either splits or merges the partitions to create a tree-like structure [3]. The second type divides data into several initial clusters partitions and iteratively data is assigned to their closest cluster partition or centroid using a dissimilarity criterion. According to [4] hierarchical algorithms are more suitable for small data sets and partitioning algorithms for large data sets.

Algorithms that perform these two types of clustering use mostly crisp or hard-clustering techniques where each object is assigned to a single cluster. Fuzzy clustering techniques allow objects to belong to multiple clusters with different degrees [5] by assigning membership degrees to the objects and assigning the object to the cluster that has the highest degree. Overlapping clustering techniques allows an object to belong to one or more clusters [6]. This has several applications such as dynamic system identification, document

categorization (document belonging to different clusters) etc. among others.

K-means is one of the most frequently used partitioning clustering algorithms and also considered one of the simplest methods for clustering [7]. This study proposes modified K-means to produce overlapping clusters by using the maximum Euclidean distance assigned to all the objects as the global threshold to reassign data to overlapping clusters.

The paper is organized as follows: Section II will give a review of related work on overlapping cluster algorithms. Section III presents MCOKE: Multi-Cluster Overlapping K-Means Extension (MCOKE) algorithm that we propose. After that the experiments and results are discussed. Finally Section V mentions the conclusions and future work.

II. RELATED WORK ON OVERLAPPING CLUSTERING

Most research on overlapping clustering has focused on algorithms that evolve fuzzy partitions of data [8] and based around the Fuzzy C-means and many of its variants [5], [9]. Data objects are assigned membership degrees (values between 0 and 1) to a particular cluster where the total sum of all membership values must add up to 1 hence generating soft partitions. Objects are eventually assigned to clusters that have the highest degree of membership. If the highest degree of membership is not unique, then an object is assigned to an arbitrary cluster that achieves the maximum. By adding a constraint where the data object must belong to a cluster with the highest membership degree, a “1” is imposed on every object in the matrix thus degenerating it to hard-partitioning.

Other overlapping algorithms such as Overlapping Partitioning Cluster (OPC) [10] that does not consider the similarity of objects in the same cluster but rather use a similarity threshold that determines whether the object will belong to a cluster or not. An object that has a distance of less than the threshold can belong to multiple clusters this way. The Overlapping K-means (OKM) [11] has a similar threshold to determine the belonging of an object. The drawback of such algorithms is that the threshold is determined *a priori* and may not be easy to determine the right threshold for different data samples.

The MCOKE algorithm assigns the global threshold to determine the belonging of a data object to a cluster once the K-means algorithm finishes its iterations and picks the maximum distance (*maxdist*) of all objects that were assigned to the clusters. In this case, the threshold *maxdist* is not a priori and can change depending on the data.

Said Baadel is with the Canadian University of Dubai, Dubai, UAE (e-mail: baadel@tud.ac.ae).

Fadi Thabtah is with the Canadian University of Dubai, Dubai, UAE.

Joan Lu is with The University of Huddersfield, UK.

III. MCOKE: MULTI-CLUSTER OVERLAPPING K-MEANS EXTENSION ALGORITHM

The MCOKE algorithm consists of two procedures. The standard K-means clustering that iterates through the data objects in order to attain a distinct partitioning of the data points given a priori number of C clusters by minimizing the distance between the objects and the cluster centroids.

Input: Number of clusters K , A set of data vectors

Output: Membership Matrix

Step One:

- Draw randomly k initial cluster prototypes

Step Two:

- For each data point, compute the centroid it is closest to using Euclidean distance measure
- Assign the data point to the cluster
- Re-compute and update the centroids

Step Three:

- If not converged, repeat Step Two; otherwise return assignment vector, final centroids vector, and maximum Euclidean distance (*maxdist*) used in the assignment

Step Four:

- Draw initial membership matrix MT
- For each data point, compare with final centroid vector distance with *maxdist*
- Add new member in MT if distance is shorter than *maxdist*

Fig. 1 MCOKE Algorithm

The goal of K-means is essentially optimizing the objective function provided as

$$J = \sum_{i=1}^C \sum_{x_i \in \mu_i} d(x_i, v_i) \quad (1)$$

where v_i is the center of cluster μ_i and $d(x_i, v_i)$ is the Euclidean distance between a point x_i and v_i . The objective J will therefore decrease with every iteration until it converges to a local minimum. Hence we know that each object assigned to a particular cluster centroid will have the minimum distance to it compared to other clusters. After the distinct partitioning of the data set, each object is assigned to one cluster. The maximum distance allowed by K-means in assigning the data objects to a cluster is saved as *maxdist* and is used in our algorithm as the global threshold of belonging to a cluster to be used in procedure two below when assigning objects to multi-clusters.

The clustering algorithm starts by comparing the cluster centroids to a random generated initial cluster centers and

iteratively re-computes the cluster centroids to a more sensible location until the centroids do not change.

The second procedure in our algorithm uses the results produced from the K-means algorithm to generate a membership table MT (of dimension $N \times C$) such that $MT(i,j)$ denotes a member of object i to cluster j where $i = 1, \dots, N$ and $j = 1, \dots, C$. Each object in $MT(i,j)$ is assigned a 1 (one) to denote membership to that cluster and a 0 (zero) for non-membership of a cluster.

The algorithm then iterates through the table MT and compares the distance of the objects assigned to their respective clusters with the other final centroids in the table. If the object distance is less than the *maxdist* (used as the threshold for belonging to a cluster) generated from K-means algorithm then that object is also assigned to that cluster centroid and the membership table is updated with a 1 (one). The complexity of the MCOKE is similar to that of K-means procedure of $O(n)$ where n is the number of data points in the dataset.

IV. EXPERIMENTS

We run three experiments to highlight and test the algorithm. Experiment 1 uses a sample data to demonstrate the application of MCOKE as compared to the standard K-means algorithm.

TABLE I
SAMPLE EXPERIMENTAL DATA

Book	Price, Chapter
ENGLISH	40,4
INTRO COMPUTER	37,4
STRATEGY	45,4
E-BUSINESS	34,2
ECON	38,5
STATISTICS	38,6
MATHS	38,3

We set k , the number of clusters to 4 and the following membership matrix is returned where each data set belongs to one cluster. The algorithm labels the clusters as $C0$, $C1$, $C2$, $C3$, and a 1 or 0 is assigned depending on whether they belong to the cluster.

TABLE II
K-MEANS MEMBERSHIP TABLE

Vector	Cluster Centroid	C0	C1	C2	C3
40,4	42,5,4	0	1	0	0
37,4	37,5,3,5	1	0	0	0
45,4	42,5,4	0	1	0	0
34,2	34,2	0	0	0	1
38,5	38,5,5	0	0	1	0
38,6	38,5,5	0	0	1	0
38,3	37,5,3,5	1	0	0	0

We then run the same data set through MCOKE and a membership matrix with the same cluster centroids returned as follows.

TABLE III
IMPLEMENTATION RESULT I OF MCOKE ALGORITHM

Vector	Cluster Centroid	C0	C1	C2	C3
40,4	42,5,4	1	0	1	0
37,4	37,5,3,5	0	0	1	1
45,4	42,5,4	1	0	0	0
34,2	34,2	0	1	0	0
38,5	38,5,5	0	0	1	1
38,6	38,5,5	0	0	0	1
38,3	37,5,3,5	0	0	1	1

Vector	C0	C1	C2
13.74,1.67,2.25,16.4,118,2.6,2.9,0.21,1.62,5.85,0.92,3.2,1060	0	1	1
13.56,1.73,2.46,20.5,116,2.96,2.78,0.2,2.45,6.25,0.98,3.03,1120	0	1	1
13.29,1.97,2.68,16.8,102,3,3.23,0.31,1.66,6,1.07,2.84,1270	0	1	0
13.72,1.43,2.5,16.7,108,3.4,3.67,0.19,2.04,6.8,0.89,2.87,1285	1	1	0

Here we note that while the cluster centroids are maintained in this instance, the membership matrix is updated each time the distance from the vector to the cluster centroid is less than *maxdist*.

We then conduct an experiment on real data that strongly show overlapping clustering from the UCI Machine Learning Repository domain on wine [12]. The data contains 13 attributes {alcohol, malic acid, ash, magnesium,...} that describe the constituents found in three types of wines. Class 1 of the wines has 59 instances depicted below.

The results suggest that out of the 59 instances of class 1 with 13 different attributes, only 10 instances can be uniquely identified with only 1 cluster while the rest overlap to multiple clusters. The test data is run 3 more times returning the following results.

TABLE V
IMPLEMENTATION RESULT III OF MCOKE ALGORITHM ON UCI DATA WINE

No. of Run	Objects in C1	Objects in C1	Objects in C2
1	17	48	42
2	26	52	36
3	42	17	48

A third experiment is conducted on UCI Repository on the Iris Plants dataset [12]. The data is classified into 3 classes of plant species with 50 instances on each class. There are 4 attributes {sepal length, sepal width, petal length, and petal width} on the set. We run the data through MCOKE with different *k* clusters and the following result on Table VI is reported with the number of objects in each cluster.

TABLE VI
IMPLEMENTATION RESULT I OF MCOKE ALGORITHM ON UCI DATA IRIS

K	C1	C2	C3	C4	C5	Mindist	Maxdist	Iter
3	50	98	92	-	-	0.05993	1.66064	5
4	92	50	50	98	-	0.10377	1.6468	7
5	50	58	26	66	45	0.05993	1.23935	10

The table shows the overlapping objects in different cluster sizes and the minimum and maximum Euclidean distance with the iterations taken to achieve the results.

V.CONCLUSION AND FUTURE WORK

In this paper we introduced MCOKE algorithm, an extension of K-means algorithm which is to partition *n* objects into *C* clusters that may overlap with each other. The algorithm differs from other algorithms in that it does not require a similarity threshold to be defined *a priori* which may be difficult to set depending on the data samples but rather uses the maximum distance (*maxdist*) allowed in K-means to assign objects to a given cluster as the global threshold. An initial run of the standard K-means algorithm produces a membership matrix table where each object is assigned to one

cluster with the minimum Euclidean distance from the centroid. At this stage all objects are assigned to at least one cluster hence the objective function of K-means is optimized. A second run is done through the membership table comparing the vectors with the final cluster centroids determined by the K-means algorithm and a new membership is assigned to the dataset object and the table updated if the distance between the two is less than *maxdist*.

Finally, while we acknowledge that more tests need to be done, MCOKE algorithm returns desirable overlapping cluster results each time its run. However, the algorithm suffers the same drawbacks of the standard K-means algorithm in that it only works with numerical data and that the objective function of K-means is designed to optimize while under constrain of assigning the data objects to hard-partition. Future work can be done to address these issues. We also propose extending the algorithm to be able to combine together every multi-cluster (with object mapping) that are very close to each other and assigning them a new cluster name on the fly while still maintaining the original clusters.

REFERENCES

- [1] C.C. Aggarwal, C.K. Reddy. Data Clustering: Algorithms and Applications. CRC Press, 2014.
- [2] A.K. Jain, R.C. Dubes. Algorithms for Clustering Data. Prentice Hall, 1988.
- [3] E. Boundaillier, G. Hebrail. Interactive interpretation of hierarchical clustering. Intell. Data Anal. 1998.
- [4] O.A. Abbas. Comparisons between Data Clustering Algorithms. The International Arab Journal of Information Technology, Vol 5. No. 3. 2008.
- [5] F. Höppner, F. Klawonn, R. Kruse, T. Runkler, Fuzzy Cluster Analysis: Methods for Classification, Data Analysis and Image Recognition, Wiley, 1999.
- [6] B. S. Everitt, S. Landau, M. Leese, "Cluster Analysis", Arnold Publishers, 2001
- [7] A. Jaini. Data Clustering: 50 years beyond k-means. Pattern Recognition Letters, 31(8): pp. 651-666, 2010.
- [8] E.R. Hruschka et. al. A survey of Evolutionary Algorithms for Clustering. IEEE Trans. Vol. 39, pp. 133-155, 2009.
- [9] J. C. Bezdek, Pattern Recognition with Fuzzy Objective Function Algorithm, Plenum Press, 1981.
- [10] Y. Chen, H. Hu. An overlapping Cluster algorithm to provide non-exhaustive clustering. Presented at European Journal of Operational Research. pp. 762-780, 2006
- [11] G. Cleuzious. An extended version of the k-means method for overlapping clustering. IEEE International Conference on Pattern Recognition. 2008
- [12] K. Bache, M. Lichman. UCI Machine Learning Repository (<http://archive.ics.uci.edu/ml>). Irvine, CA: University of California, School of Information and Computer Science. 2013
- [13] N. Abdelhamid, A. Ayesh, F. Thabtah. Phishing detection based Associative Classification data mining. Expert Systems with Applications Journal. Vol. 41 (13). 2014