

Fuzzy Logic based Document Inspection Using TF-IDF

M. Kishore Babu

¹Assistant Professor, Department of CSE,
 Universal College of Engineering & Technology,
 Guntur, AP, India.

Abstract - Clustering is one the main area in data mining literature. There are various algorithms for clustering. There are several clustering approaches available in the literature to cluster the document. But most of the existing clustering techniques suffer from a wide range of limitations. The existing clustering approaches face the issues like practical applicability, very less accuracy, more classification time etc. In recent times, inclusion of fuzzy logic in clustering results in better clustering results. One of the widely used fuzzy logic based clustering is Fuzzy C-Means (FCM) Clustering. In order to further improve the performance of clustering, this thesis uses Modified Fuzzy C-Means (MFCM) Clustering. Before clustering, the documents are ranked using Term Frequency-Inverse Document Frequency (TF-IDF) technique. From the experimental results, it can be observed that the proposed technique results in better clustering results when compared to the existing technique.

Keywords - Fuzzy C-means Clustering, Datasets, and Multi View Point Clustering.

1. Introduction

Clustering is one of the most exciting and important topics in information exploration. The aim of clustering is to find intrinsic components in information, and arrange them into meaningful subgroups for further analysis and analysis. There have been many clustering methods released every year. They can be suggested for very unique analysis areas, and developed using absolutely different methods and methods.

Nevertheless, according to majority of folks [1], more than 50 years after it was presented, the simple algorithm k-means still continues to be as one of the top 10 data mining methods these days.

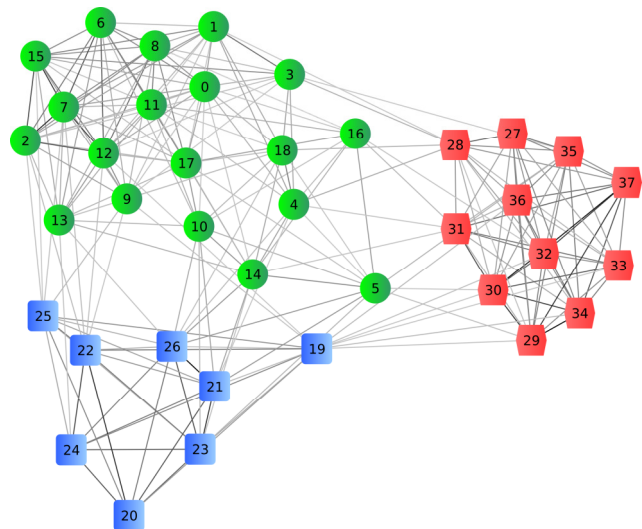


Fig .1. Data clustering in data analysis.

Needless to bring up, k-means has more than a few basic drawbacks, such as sensitiveness to initialization and to cluster dimension, and its efficiency can be more intense than other state-of-the-art methods in many websites. Regardless of that, its convenience, understandability, and scalability are the reasons for its remarkable reputation. An criteria with adequate efficiency and functionality in most of application scenarios could be much better one with better performance in some situations but restricted utilization due to great complexness. While providing affordable outcomes, k-means is quick and simple to combine with other methods in bigger techniques.

A typical strategy to the clustering problem is to treat it as an marketing procedure. An maximum partition is found by improving a particular operate of likeness (or distance)

among information. Generally, there is an implied supposition that the real implied framework of information could be properly described by the likeness system described and embedded in the clustering requirements operate. Hence, effectiveness of clustering methods under this approach depends on the suitability of the likeness evaluate to the information at hand. For example, the unique k-means has sum-of-squared-error purpose operate that uses Euclidean distance. In a very rare and high-dimensional domain like written text records, rounded k-means, which uses cosine similarity (CS) instead of euclidean range as the evaluate, is considered to be more appropriate.

As outlined in the sharp case may not be quickly generalized for unclear clustering. This is due to the fact that in fuzzy methods an example does not are part of a group completely but has restricted account principles in most clusters. More about clustering methods can be discovered in. Clustering considerable quantities of information requires a long time. Further, new unlabeled information places, which will not fit in storage, are becoming available. To group them, either sub testing is required to fit the information in storage or time will be significantly affected by hard drive accesses making clustering an unpleasant choice for information analysis.

Another resource of huge information places is streaming information where you do not shop all the information, but process it and remove it. There are some very huge information places for which a little branded information is available and the relax of the data is unlabeled i.e. for example, computer attack detection. Our first purpose is to obtain a novel technique for measuring likeness between information things in rare and high-dimensional sector, particularly written text records.

2. Background Approach

The likeness between two records d_i and d_j is identified w.r.t. the position between the two factors when looking from the source. To create a new idea of likeness, it is possible to use more than just one referrals factor. We may have a more precise evaluation of how near or remote a couple of factors are, if we look at them from many different opinions. A presumption of group subscriptions has been created before to the evaluate. The two things to be calculated must be in the same group, while the factors from where to set up this measurement must be outside of the group.

$$\begin{aligned} \text{MVS}(d_i, d_j | d_i, d_j \in S_r) \\ &= \frac{1}{n-n_r} \sum_{d_h \in S \setminus S_r} (d_i - d_h)^t (d_j - d_h) \\ &= \frac{1}{n-n_r} \sum_{d_h} \cos(d_i - d_h, d_j - d_h) \|d_i - d_h\| \|d_j - d_h\|. \end{aligned}$$

As shown in the above figure, the likeness between two factors d_i and d_j within group S_r , considered from a factor d_h outside this group, is similar to the item of the cosine of the position between d_i and d_j looking from d_h and the Euclidean ranges from d_h to these two factors. This meaning is in accordance with the supposition that d_h is not in the same group with d_i and d_j . Little sized the ranges $k_{di} \cap d_{hk}$ and $k_{dj} \cap d_{hk}$ are, the greater the opportunity that d_h is actually in the same group with d_i and d_j , and the likeness depending on d_h should also be minute indicate this prospective.

procedure BUILDMVSMATRIX(A)

```

for  $r \leftarrow 1 : c$  do
     $D_{S \setminus S_r} \leftarrow \sum_{d_i \notin S_r} d_i$ 
     $n_{S \setminus S_r} \leftarrow |S \setminus S_r|$ 
end for
for  $i \leftarrow 1 : n$  do
     $r \leftarrow \text{class of } d_i$ 
    for  $j \leftarrow 1 : n$  do
        if  $d_j \in S_r$  then
             $a_{ij} \leftarrow d_i^t d_j - d_i^t \frac{D_{S \setminus S_r}}{n_{S \setminus S_r}} - d_j^t \frac{D_{S \setminus S_r}}{n_{S \setminus S_r}} + 1$ 
        else
             $a_{ij} \leftarrow d_i^t d_j - d_i^t \frac{D_{S \setminus S_r} - d_j}{n_{S \setminus S_r} - 1} - d_j^t \frac{D_{S \setminus S_r} - d_i}{n_{S \setminus S_r} - 1} + 1$ 
        end if
    end for
end for
return  $A = \{a_{ij}\}_{n \times n}$ 
end procedure

```

Fig .2. Process of multi view clustering in real time data sets.

The overall likeness between d_i and d_j is identified by getting regular over all the opinions not that belong to group S_r . It is possible to claim that while most of these opinions are useful, there may be some of them providing deceiving details just like it may occur with the source factor. However, given a huge enough variety of opinions and their wide range, it is affordable to believe that most they will be useful.

3. Proposed Approach

This constraint makes FCM to be exceptionally touchy to clamor. The general rule of the method introduced in this paper is to consolidate the area data into the FCM calculation amid arrangement. Keeping in mind the end

goal to join the spatial connection into FCM calculation, the destination capacity of aforementioned comparison may seems equal and other critical issues will be planned and the proposed calculation focused around regularization term in predefined and different incidents is punished by a regularization term, which is propelled by the above NEM calculation and adjusted focused around the paradigm of FCM calculation. The new target capacity of the PFCM is characterized as takes after:

$$J_{PFCM} = \sum_{k=1}^n \sum_{i=1}^c (u_{ik})^q d^2(x_k, v_i) + \gamma \sum_{k=1}^n \sum_{j=1}^n \sum_{i=1}^c (u_{ik})^q (1 - u_{ij})^q w_{kj}$$

The parameter g ($g \geq 0$) controls the impact of the punishment term. The relative criticalness of the regularizing term is conversely corresponding to the signal-to-noise ratio (SNR) of the picture. Lower SNR would oblige a higher estimation of the parameter g , and the other way around. At the point when $g = 0$, JPFCM measures up to JFCM. The significant contrast between NEM calculation and PFCM calculation is that the punishment term in the NEM is boosted to get the arrangements while in the PFCM it ought to be minimized so as to fulfill the standard of FCM calculation. Moreover, the punishment term in the PFCM calculation has the weighting example to control the level of fluffiness in the ensuing participation capacity in spite of the punishment term in the NEM calculation that is fresh. This new punishment term is minimized when the enrollment esteem for a specific class is vast and the participation values for the same class at neighboring pixels is additionally huge, and the other way around. As it were, it obliges the pixel's participation estimation of a class to be associated with those of the neighboring pixel.

4. Experimental Evaluation

To confirm the key benefits of our suggested techniques, we evaluate their efficiency in tests on document data. The purpose of this area is to evaluate MVSC-IR and MVSC-IV with the current techniques that also use specific likeness actions and requirements features for document clustering. The likeness actions to be compared includes Euclidean range, cosine likeness, and extended Jacquard coefficient.

The information corpora that we used for tests include of 20 standard papers information places. Besides reuters7 and k1b, which have been described in information previously, we included another 18 written text selections

so that the evaluation of the clustering techniques is more thorough and comprehensive.

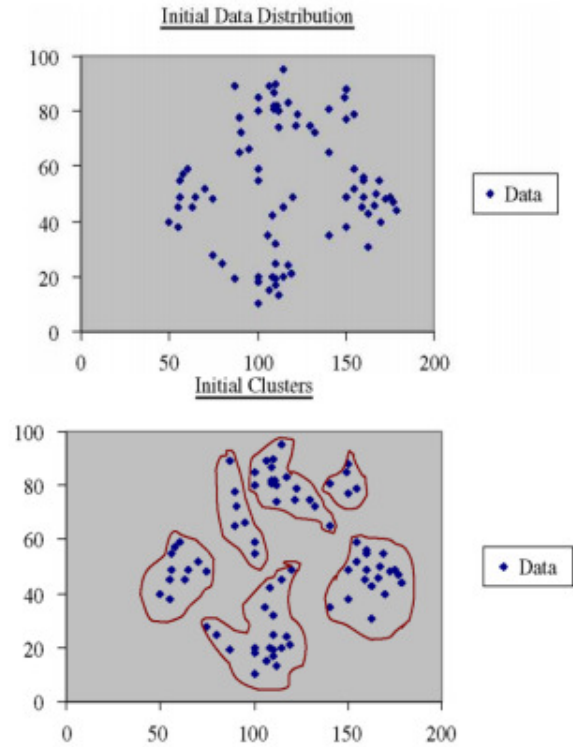


Fig .3. Cluster results with processing application development.

It has been known that requirements function-based partitioned clustering techniques can be delicate to group dimension and balance. In the ingredients of IR, parameter λ which is known as the controlling aspect, $\lambda \in [0, 1]$. To analyze how the dedication of λ could impact MVSCIR's performance, we analyzed MVSC-IR with different values of λ from 0 to 1, with 0.1 step-by-step period. The assessment was done in accordance with the clustering outcomes in NMI, FScore, and Precision, each averaged over all the 20 given information places. Since the evaluation analytics for different data places could be very different from each other, simply taking the common over all the information places would not be very meaningful. Hence, we applied the technique used to convert the analytics into comparative analytics before calculating.

The customized powerful clustering criteria is as follow:

Input: k : variety of groups (for powerful clustering initialize $k=2$)

Fixed variety of groups = yes or no (Boolean).

D: a information set containing n things.

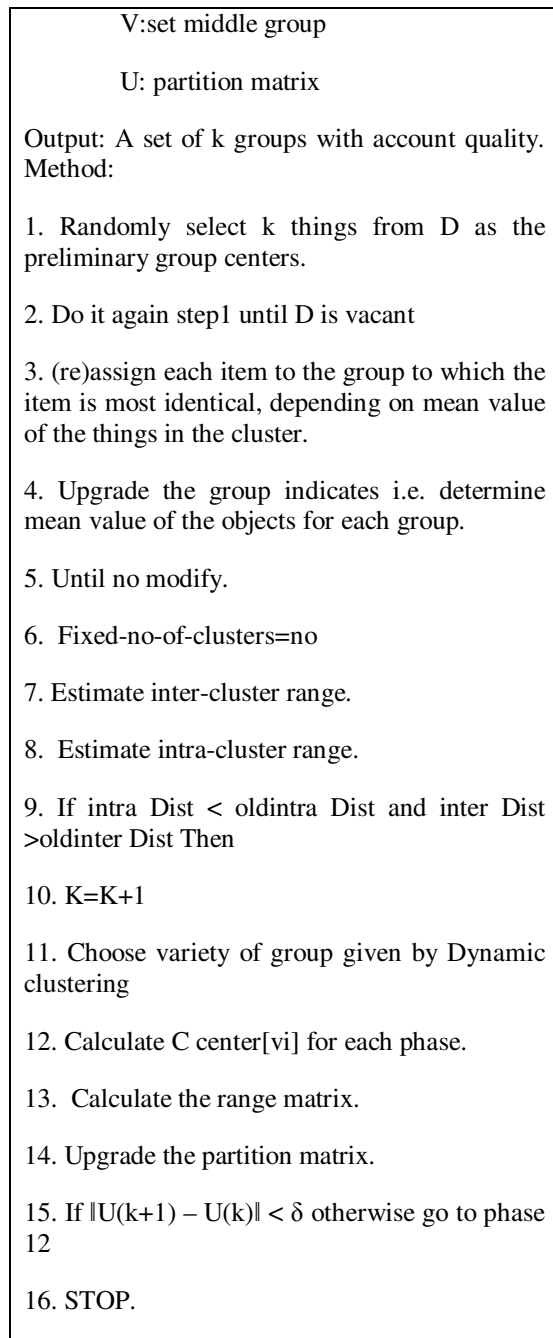


Fig .4. Algorithm for fuzzy c means process.

Powerful Indicates clustering strategy is the new methodology to group the information things into variety of groups, which is unidentified originally. Variety of categories (clusters) is some beneficial integer. The collection is done by measuring the range between item and centroid. Objects are iteratively arranged into the current categories or a new cluster development is done with those things centered up on the threshold restrict. Thus the objective of dynamic clustering is to classify the

information. It could enhance the possibilities of discovering the global optima with cautious choice of preliminary group. In this algorithm information things are saved in additional storage and transferred to primary storage individually. Only the group representatives are saved completely in primary storage to alleviate area restrictions. Therefore, an area need of these criteria is very small, necessary only for the centroid of the categories. This algorithm is non-iterative and therefore it is time need is also little.

5. Conclusion

Clustering decides the connections between information objects in the data source. The things are arranged or arranged based on the key of “maximizing the infraclass similarity and reducing the interclass similarity”. It discovers out something useful from data source. Clustering has its roots in many areas, such as information exploration, research, biology, and device learning etc. Clustering methods can be divided into various types: Dividing methods, Hierarchical methods, Solidity centered methods, Grid-based methods; Design centered methods, Probabilistic methods, and Chart theoretic and Unclear methods. The Powerful mean algorithm are the significant concentrate of this dissertation work. Dynamic mean criteria generate good groups automatically because there is no need to described the number of groups before head but in Powerful mean criteria each data factor can be a participant of one and only one group at a time. In other terms we can say that the sum of account grades of each information point in all groups is similar to one and in all the staying groups its account quality is zero. In our thesis dynamic criteria is customized using fuzzy criteria. By implementing fuzzy criteria over Powerful criteria we can show the account of each information factor in all groups. By applying Unclear criteria over Powerful criteria clustering can be at an extremely quicker rate. It is appropriate to a large amount of information saved in databases. The overall results are significant in displaying that Powerful criteria display membership of each information factor in every groups.

References

- [1] X. Wu, V. Kumar, J.R. Quinlan, J. Ghosh, Q. Yang, H. Motoda, G.J. McLachlan, A. Ng, B. Liu, P.S. Yu, Z.-H. Zhou, M. Steinbach, D.J. Hand, and D. Steinberg, “Top 10 Algorithms in Data Mining,” Knowledge Information Systems, vol. 14, no. 1, pp. 1-37, 2007.
- [2] I. Guyon, U.V. Luxburg, and R.C. Williamson, “Clustering: Science or Art?,” Proc. NIPS Workshop Clustering Theory, 2009.
- [3] I. Dhillon and D. Modha, “Concept Decompositions for Large Sparse Text Data Using Clustering,”

- Machine Learning, vol. 42, nos. 1/2, pp. 143-175, Jan. 2001.
- [4] S. Zhong, "Efficient Online Spherical K-means Clustering," Proc. IEEE Int'l Joint Conf. Neural Networks (IJCNN), pp. 3180-3185, 2005.
- [5] A. Banerjee, S. Merugu, I. Dhillon, and J. Ghosh, "Clustering with Bregman Divergences," J. Machine Learning Research, vol. 6, pp. 1705-1749, Oct. 2005.
- [6] E. Pekalska, A. Harol, R.P.W. Duin, B. Spillmann, and H. Bunke, "Non-Euclidean or Non-Metric Measures Can Be Informative," Structural, Syntactic, and Statistical Pattern Recognition, vol. 4109, pp. 871-880, 2006.
- [7] M. Pelillo, "What Is a Cluster? Perspectives from Game Theory," Proc. NIPS Workshop Clustering Theory, 2009.
- [8] D. Lee and J. Lee, "Dynamic Dissimilarity Measure for Support Based Clustering," IEEE Trans. Knowledge and Data Eng., vol. 22, no. 6, pp. 900-905, June 2010.
- [9] A. Banerjee, I. Dhillon, J. Ghosh, and S. Sra, "Clustering on the Unit Hypersphere Using Von Mises-Fisher Distributions," J. Machine Learning Research, vol. 6, pp. 1345-1382, Sept. 2005.
- [10] Mr.Kamakshaiah K, Dr.R..Seshadri "Water Quality Analysis Using Enhanced K-Means Clustering" International Journal of Advanced Research in Computer Science and Software Engineering, Volume 5, Issue 10 October (2015) ISSN : 2277-128X
- [11] T. Q. Chen and Y. Lu, "Color image segmentation an innovative approach", Pattern recognition, vol. 35, 2002, pp. 395-405.
- [12] Y. Yang, C. Zheng, and P. Lin, "Fuzzy c-means clustering algorithm with a novel penalty term for image segmentation" Optoelectronic review, Vol.13, Issue 4, 2005, pp. 309-315.
- [13] Mr. Kamakshaiah K, Dr.R..Seshadri "Ground Water Quality Assessment Using Data Mining Techniques" International journal of Computer Applications Volume 76-No 15, August 2013, ISSN: 0975-8887.
- [14] Mr.Kamakshaiah K, Dr.R..Seshadri, "Classification Ground Water Process Using PC Based K-Means Clustering ", International Journal of Applied Sciences, Engineering and Management, Vol 3,-No61, November 2014,ISSN: 2320-3439.