

# Review of gravitational search algorithm for Clustering problem

Pawan<sup>1</sup>, Parveen Khanchi<sup>2</sup>

<sup>1</sup>M.Tech Student, Department of ECE, BIMT, Chidana, India

<sup>2</sup>Head of the Department of ECE, BIMT, Chidana, India

**Abstract:** In statistic and data mining, agglomerative clustering is well known for its efficiency in clustering large data sets. The aim is to group data points into clusters such that similar items are lumped together in the same cluster. In general, given a set of objects together with their attributes, the goal is to divide the objects into clusters such that objects lying in one cluster should be as close as possible to each other's (homogeneity) and objects lying in different clusters are further apart from each other. However, there exist some flaws in classical agglomerative clustering algorithm. According to the method, first, the algorithm is sensitive to selecting initial threshold level and on the other hand, the agglomerative clustering is NP hard problem in selecting the optimum threshold level so that maximum F-measure or correct assignment of data to right clusters can be obtained.

In this paper, to solving the agglomerative clustering problem, we provide optimizing threshold level in clustering to decide the number of clusters, which in this algorithm we consider the issue of how to derive an optimization model to the maximum accuracy which is measured in terms of F-measure. We introduce the optimization algorithm named Gravitational Search Algorithm (GSA) to optimize k-means algorithm to guarantee the result of clustering is more accurate than clustering by basic clustering algorithms. F-measure is used to compare the performance of both algorithms.

**Keywords:** Optimized Agglomerative Clustering (OAC and hierarchical clustering..

## I. INTRODUCTION

In general, the clustering techniques are broadly divided into partition and hierarchical methods. K-mean is the most popular partition clustering algorithm which is based on centroid [1 and 2]. The k-mean algorithm takes less number of iterations and time to partition the data set, but not suitable for large data set and hence more inconsistent [3 and 4]. The hierarchical clustering algorithm is either divisive or agglomerative. The divisive method works on the top-down approach [5 and 6]. In this method, all the data objects are arranged within a big single cluster, and the larger cluster is continuously divided into smaller clusters until each cluster consists of a single object. Agglomerative hierarchical clustering technique works on bottom-up approach [7-9]. The technique starts with N clusters and each of which contains exactly one data object. A series of merge operations is then followed that eventually forces all clusters into the same single cluster. Many papers have been reported in [10-15] with drawbacks in the traditional agglomerative hierarchical cluster technique as listed below: high computational complexity to merge the two closest leader clusters, consumption of more time to merge two closest leader clusters by minimum or maximum or average cutting distance, consumption of (N-1) iterations to categorize large data set where N denotes the dataset size, high misclassification errors due to the number of iterations, high complication to predict the outliers on the large data set and consumption high computational complexity  $O(N^3)$  for large dataset.

In [16], Hisashi Koga et.al., have reported a fast approximation algorithm for single linkage method, that reduces the time complexity by rapidly finding the near

clusters to be connected by locality-sensitive hashing (LSH). P.A. Vijaya et.al. [17], have reported another hierarchical clustering algorithm (Leaders-Sub leaders) for large data sets, that uses an incremental clustering principle to generate a hierarchical structure for finding the sub-clusters within each cluster. Antti Honkela et.al. [18], have been suggested two variants of an agglomerative technique for learning a hierarchy of independent variable of group analysis. It resembles hierarchical clustering, but the choice of clusters to merge is based on variation Bayesian model comparison. It also allows determining optimal cutoff points on the hierarchy. In [19], J.A.S. Almeida et.al, have suggested to improve agglomerative hierarchical clustering algorithm based on single linkage. They have reported two analysis tools that involve identification of outliers and normal clusters for the data set. The main drawback in this method is its high computational complexity for merging every leader clusters pair. In [20], a distance hierarchy scheme for both categorical and numerical values has been reported with an integrated scheme for agglomerative hierarchical clustering with mixed data. The integrated approach allows for expressing similarity relationships between categorical values and hence produces better clustering results. Manoranjan Dash et.al. [25], reported a fast Hierarchical Agglomerative Clustering (HAC) based on partially overlapping partitioning (POP). They presented a two-phase algorithm for HAC based on POP. In phase-1, large numbers of small clusters are merged and in-phase-2,

The history of extraction of patterns from data is centuries old. The earlier method which has been used is Bayes' theorem (1700s) and regression analysis (1800s). [1] In the field of computer technology, using

the ever growing power of computers, we develop an essential tool for working with data. Such as, it is being able to work with increasing size of the datasets and complexity. And also an urgent need to further refine the automatic data processing, which has been aided by other discoveries in computer science, means that our ability for data collection storage and manipulation of data has been increased. As definition, Data mining or important part of Knowledge Discovery in Database (KDD), used to discover the most important information throughout the data, is a powerful new technology. Across a myriad variety of fields, data are being collected and of course, there is an urgent need to computational technology which is able to handle the challenges posed by these new types of data sets.

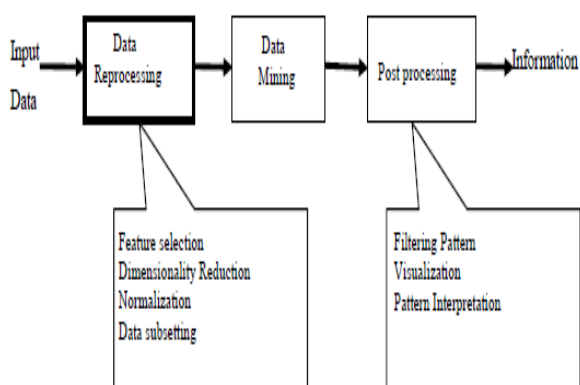


Figure 1.1 The overall steps of the process of Knowledge Discovery in Database (KDD)

There are challenges in traditional data analysis techniques and always new types of datasets. In order to cope with these new challenges, researchers have been developing more efficient and scalable tools that can more easily handle diverse types of data. In particular, data mining draws upon ideas such as:

- 1- Sampling, estimating and hypothesis testing from statistic
- 2- Search algorithms, modelling techniques and learning theories from artificial intelligence, pattern recognition and machine learning.

The field of Data mining grows up in order to extract useful information from the rapidly growing volumes of data. It scours information within the data that queries and reports can't effectively reveal.

This process contains a series of transformation steps, from data pre-processing to data mining results. [1] And also data mining has been adopting from other areas, such as: optimization, evolutionary computing, information theory, signal processing, visualization and information retrieval. [6] Agglomerative clustering is also an important method for data mining clustering which is gaining attention these days.

In Agglomerative clustering algorithm, a cut height parameter is required to determine the dissimilarity threshold at which clusters are allowed to be merged together. This parameter greatly influences the

clustering accuracy, as measured by the Rand index, of the final clusters produced. For instance, using a very high cut height or dissimilarity threshold would result in most data being included in one giant cluster since a weak measure of similarity is enforced during the merging process. So an optimum selection of threshold level is necessary. In our work we work towards optimising this threshold height.

## II. DIFFERENT CLUSTERING

### 2.1 Hierarchical Clustering

The set of given data objects are partitioned in form of a tree like structure or nested clusters in hierarchical clustering. The hierarchical methods can be classified into two types.

- Agglomerative and
- Divisive

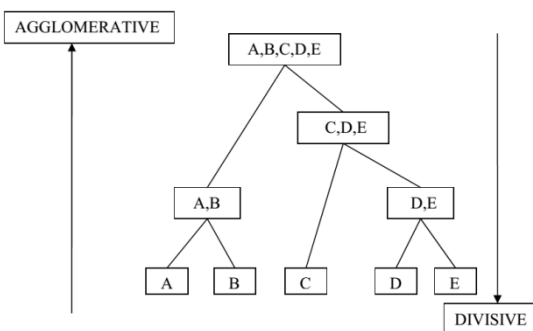


Figure 2.1: Agglomerative and Divisive clustering

In agglomerative method also known as bottom-up approach, each object forms a separate group. It successively merges the groups close to one another by checking the similarity function, until all the groups are merged into one, that's until the top most level of hierarchy is reached or until a termination condition holds. In divisive clustering also known as top-down approach, initially all the objects are grouped into a single cluster which can also be called as parent. In each successive iteration, a cluster is split up into smaller clusters, until eventually each object is in one cluster or until a termination condition holds.

#### 2.1.1 Agglomerative Method

This method begins by treating each object as an individual cluster and then proceeds by merging two nearest clusters. The distance between any two clusters  $m$  and  $n$  is defined by a metric  $D_{m,n}$ . Metrics can be single-link, complete-link and group average etc. A general class of metrics was given by Lance and Williams [21]. If  $D_{k,i}$  be the distance between cluster  $k$  and the union of cluster  $i$  and cluster  $j$ , then:

$$D_{k,ij} = \alpha_i D_{k,i} + \alpha_j D_{k,j} + \beta D_{i,j} + \gamma |D_{k,i} - D_{k,j}|$$

The agglomerative method is as follows:

- Consider each object to be an atomic cluster. The  $(n \times n)$  distance matrix represents the distance between all possible pairs of clusters.
- Find the smallest element in the matrix. This corresponds to the pair of clusters that are most similar. Merge these two clusters, say  $m$  and  $n$ , together.
- Measure the distances between the newly formed cluster and the other remaining clusters using a distance function. Delete the row and column of  $m$  and overwrite row and column of cluster  $n$  with the new values.
- If the current number of clusters is more than  $k$  then go to step 2; otherwise stop.

The merging process can continue until all the objects are in one cluster. The advantages of hierarchical methods are that they are easy to implement computationally. They are able to tackle larger datasets than the  $k$ -medoids method and we can run the algorithm without providing the input  $k$  (the number of clusters to be formed). The drawbacks of agglomerative method are:

- The algorithm has  $O(n^3)$  time complexity. Even though the order of the distance matrix decreases with each iteration, the cost of Step2 on iteration  $k$  is  $O((n - k)^2)$ , and we are guaranteed  $(n - k)$  iterations before we get to  $k$ ;
- The clusters produced are heavily dependent on the metric  $D_{i, j}$ . Different metrics can produce different clusters. For instance, the complete-link metric tends to produce spherical clusters, whereas the single-link metric produces elongated clusters [21].

## 2.2 Problem Description

The agglomerative clustering is the unsupervised approach which makes cluster of same kind of data to make data analysis easier. This clustering algorithm is described in previous chapter in section 3.1. We are targeting the bottom approach of agglomerative clustering in which every node starts with its own cluster and based on similarity value, other nodes are combined with its cluster. This similarity measure is calculated in terms of Euclidean distance between nodes. These steps are repeated for all nodes and clusters are dependent upon the minimum distance between nodes. Once all nodes are clustered depending upon their distance, number of clusters are selected using a threshold level. Classical agglomerative clustering approach selects the threshold level to divide the cluster tree into specific clusters numbers by using Euclidean distance. The accuracy of this clustering algorithm depends upon the threshold level. For example if actual cluster number is 4 but threshold level divides the cluster tree into 5 clusters then false

clustering will be high. So to improve the accuracy and F-measure, it should be set at optimal position. To fulfill this purpose evolutionary optimization algorithms have been used which uses minimum distance concept to cluster. Particle swarm optimization, bacterial foraging optimization, genetic algorithm etc are used earlier for data clustering purpose. These evolutionary techniques can be categorized as global optimization and local optimization techniques. As all such kind of algorithms look for local minimum position for which cost function has minimum value but local optimization algorithms like genetic algorithms (GA), ant colony optimization (ACO), particle swarm optimization (PSO) etc, sometimes jumps over the local minimum point whereas global optimization techniques like gravitational search algorithm (GSA) which came into existence in 2013 has no such issue, it checks for all iteration values for local minimum point. In case of multi objective functions, global optimization algorithms perform well. But these suffers from a drawback of speed. Iteration speeds of global techniques are less than local. So these take long time to process.

So in our work we have used a gravitational search algorithm technique for data clustering which used agglomerative clustering objective function for data clustering.

## IV. CONCLUSION

In this work, the problem was to solve the agglomerative clustering problem by introducing a clustering technique – Gravitational search Algorithm in agglomerative clustering which is an optimization algorithm to tune the cutoff level of clustering tree. The problem in clustering as we notice is because of the cut off level in agglomerative method which is based on the Euclidean distance between nodes. Sometimes we have a poor clustering (some clusters don't have any member). The goal is clustering in the best behaviour, which should be to group similar data points as much as possible. But with classical agglomerative clustering this rarely is the case.

To optimize the clustering, we propose an algorithm. In this data clustering method concept of maximizing the F measure between every cluster head and other data points. Gravitational search algorithm (GSA) is used. Cut off position optimized by GSA serves as input to the clustering algorithm. Position of level is initialized randomly in this also but later on it changes the position as per tuning method of respective optimization technique. Comparison of results with classical agglomerative algorithm is done in terms of F-measure. Standard deviation comparison if it, as in figure 3.4 shows improvement by proposed algorithm.

## REFERENCES

- [1]. Kriti Srivastava, R. Shah, D. Valia, and H. Swaminarayan, "Data Mining Using Hierarchical Agglomerative Clustering Algorithm in Distributed Cloud Computing Environment", *International Journal of Computer Theory and Engineering*, Vol. 5, No. 3, June 2013
- [2]. Marek Lipczak, Evangelos Milios, "Agglomerative Genetic Algorithm for Clustering in Social Networks", Genetic and Evolutionary Computation Conference, GECCO 2009, Proceedings, Montreal, Québec, Canada, July 8-12, 2009.
- [3]. Singaravelu.S, A.Sherin and S.Savitha, "Agglomerative Fuzzy K-Means Clustering Algorithm", *A Journal of Nehru Arts and Science College (NASC)*, Vol 1 (2013)
- [4]. R. Krishnamoorthy and S. SreedharKumar, "New optimized agglomerative clustering algorithm using multilevel threshold for finding optimum number of clusters on large data set," *Emerging Trends in Science, Engineering and Technology (INCOSET), 2012 International Conference on*, Tiruchirappalli, Tamilnadu, India, 2012, pp. 121-129.
- [5]. Jake M Drew and Tyler Moore, "Optimized combined – clustering methods for finding replicated criminal websites", *EURASIP Journal on Information Security* (2014).
- [6]. Sanjay Tiwari, Mahinder Kumar Rao, "Optimization In Association Rule Mining Using Distance Weight Vector And Genetic Algorithm" *International Journal of Advanced Technology & Engineering Research (IJATER)*, Volume 4, Issue 1, Jan. 2014.
- [7]. Poonam Schrawat, Manju, "Association Rule Mining Using Particle Swarm Optimization", *International Journal of Innovations & Advancement in Computer Science*, Volume 2, Issue 1 January 2014
- [8]. R.Jensi and G.Wiselin Jiji, "Hybrid Data Clustering Approach Using K-Means And Flower Pollination Algorithm", *Advanced Computational Intelligence: An International Journal (ACII)*, Vol.2, No.2, April 2015
- [9]. Khalid Raza, "Clustering analysis of cancerous microarray data", *Journal of Chemical and Pharmaceutical Research*, 2014, 6(9)
- [10]. P. Ramachandran, N.Girija, "Early Detection and Prevention of Cancer using Data Mining Techniques", *International Journal of Computer Applications*, Volume 97– No.13, July 2014.
- [11]. Hlaudi Daniel Masethe, Mosima Anna Masethe, "Prediction of Heart Disease using Classification Algorithms", *Proceedings of the World Congress on Engineering and Computer Science 2014 Vol II*
- [12]. P.Kalyani, "Medical Data Set Analysis – A Enhanced Clustering Approach" *International Journal of Latest Research in Science and Technology*, Volume 3, Issue 1: Page No.102-105 ,January-February 2014
- [13]. Ibrahim M. El-Hasnony, Hazem M. El Bakry, Ahmed A. Saleh, "Data Mining Techniques for Medical Applications: A Survey", *Mathematical Methods in Science and Mechanics*, 2014
- [14]. Sundararajan S, Dr. Karthikeyan S, "An Hybrid Technique for Data Clustering Using Genetic Algorithm with Particle Swarm Optimization", *International Journal of Advanced Research in Computer Science and Software Engineering*, Volume 4, Issue 12, December 2014
- [15]. Sundararajan S., And Karthikeyan S, "An Efficient Hybrid Approach For Data Clustering Using Dynamic K-Means Algorithm And Firefly Algorithm", *ARNP Journal of Engineering And Applied Sciences*, Vol. 9, No. 8, August 2014
- [16]. Sandeep Rana, Sanjay Jasola, Rajesh Kumar, "A hybrid sequential approach for data clustering using K-Means and particle swarm optimization algorithm", *International Journal of Engineering, Science and Technology*, Vol. 2, No. 6, 2010
- [17]. Sandeep U. Mane, Pankaj G. Gaikwad, "Hybrid Particle Swarm Optimization (HPSO) for Data Clustering", *International Journal of Computer Applications (0975 8887)* Volume 97 - No. 19, July 2014
- [18]. T. Niknam, M. Nayeripour and B.Bahmani Firouzi, "Application of a New Hybrid optimization Algorithm on Cluster Analysis", *International Journal of Computer, Electrical, Automation, Control and Information Engineering* Vol:2, No:10, 2008
- [19]. Amin Rostami and Maryam Lashkari, "Extended Pso Algorithm For Improvement Problems K-Means Clustering Algorithm", *International Journal of Managing Information Technology (IJMIT)* Vol.6, No.3, August 2014
- [20]. M. Bhanu Sridhar1, Y. Srinivas2, M. H. M. Krishna Prasad, "Software Reuse in Cardiology Related Medical Database Using K-Means Clustering Technique", *Journal of Software Engineering and Applications*, 2012.