

An Application of Genetic Algorithm with Iterative Chromosomes for Image Clustering Problems

Mr.R.Balakrishnan¹, Mr.U.Karthick Kumar²

¹Assistant Professor & Head
Department of MCA & Software Systems
VLB Janaki Ammal College of Arts and Science
Coimbatore, Tamil Nadu, India -641042

²Assistant Professor
Department of MCA & Software Systems
VLB Janaki Ammal College of Arts and Science
Coimbatore, Tamil Nadu, India -641042

Abstract

Many heuristic algorithms have been applied to the clustering problem, which is known to be NP Hard. This paper represents a Genetic Algorithm for clustering on image data. Genetic algorithms have been used in a wide variety of fields to perform clustering, however, the technique normally has a long running time in terms of input set size. This paper proposes an efficient genetic algorithm for clustering on very large data sets, especially on image data sets. In this study, a heuristic method based on Genetic Algorithms (GA) is adopted to automatically determine the number of cluster centroids during unsupervised classification. Efficient time techniques are used as a performance measure for clustering on image data. This paper compares Genetic algorithm with K-Means algorithm for clustering on image data.

KeyWords: *Data Mining, Clustering, K-means algorithm, Genetic algorithm, Image data, Heuristic method.*

1. Introduction

Data mining techniques are the result of a long process of research and product development [20]. This evolution began when business data was first stored on computers, continued with improvements in data access, and more recently, generated technologies that allow users to navigate through their data in real time. Data mining takes this evolutionary process beyond respective data access and navigation to prospective and proactive information delivery. Data mining is ready for application in the business community because it is supported by three technologies such as Massive data collection, Powerful multiprocessor computers and Data mining algorithms [11]. Data mining deals with large databases that impose on clustering analysis additional severe computational requirements. Cluster analysis is mainly conducted on computers to deal with very large scale and complex datasets. With the development of computer based techniques, clustering has been widely used in data mining, ranging from web mining, image processing, machine

learning, artificial intelligence, pattern recognition, social network analysis, bioinformatics, geography, geology, biology, psychology, sociology, customers behavior analysis, marketing to e-business and other fields[14].

The task of grouping data points into clusters of "similar" items are a form of unsupervised learning that has applications in many fields. For instance, current techniques used for machine vision require processing of digital information obtained from pixels [2]. A very important step in this digital information processing is to group the data in some fashion so that patterns can be recognized. Clustering can be used for this task. In the medical field, clustering of data can be used to determine if a drug provides greater benefits to a certain group of patients. Grouping of information is used in the engineering field to determine what factors lead to the failure of a component in a system [14]. And in marketing, data clustering can give a clearer picture of how to focus an advertising campaign to the proper audience. The concept of clustering has been around for a long time. It has several applications, particularly in the context of information retrieval and in organizing web resources. The main purpose of clustering is to locate information and in the present day context, to locate most relevant electronic resources [3]. The research in clustering eventually led to automatic indexing as well as to retrieve electronic records. Clustering method in which we make group of objects that are somehow similar in characteristics.

The ultimate aim of the clustering is to provide a grouping of similar records. A clustering is essentially a set of such clusters, usually containing all objects in the data set. Additionally, it may specify the relationship of the clusters to each other, for example a hierarchy of clusters embedded in each other. Clustering can be roughly distinguished by hard clustering in which each object belongs to a cluster or not and soft clustering in which each

object belongs to each cluster to a certain degree e.g. a likelihood of belonging to the cluster. This paper will discuss the use of Genetic Algorithms (GA) for the task of clustering data. The running time for most clustering becomes quite large as the size of the input data set increases. Here, the application of GA algorithm for clustering on very large data sets, such as image data sets, will be addressed. The paper is organized as follows. In section II we explain related work, in section III, we review the clustering techniques. In Section IV-we describe clustering problems and in section V detail Genetic algorithms for clustering on very large data sets. In section VI describe performance analysis and conclusions describes in section VII.

2. Related work

Clustering has become a widely studied problem in a variety of application domains, such as in data mining and knowledge discovery [1], [2] statistical data analysis [3], [4] data classification and compression [6], medical image processing [5] and bioinformatics [6]. Several algorithms have been proposed in the literature for clustering [7], [8]. A. L. Abul is explained about Cluster Validity Analysis Using Sub sampling [10]. The objective of all clustering algorithms is to divide a set of data points into subsets so that the objects within a subset are similar to each other and objects that are in different subsets have diverse qualities [11], [12], [13]. Bradley and Fayyad have proposed an algorithm for refining the initial cluster centers. Not only are the true clusters found more often, but the clustering algorithm also iterates fewer times [15]. Some clustering methods improve performance by reducing the distance calculations. For example, Judd *et al.* proposed a parallel clustering algorithm P-CLUSTER [16] which uses three pruning techniques.

The K-Means algorithm [7] is well known for its efficiency in clustering large data sets. Fuzzy versions of the k-means algorithm have been reported in Ruspini [9] and Bezdek, where each pattern is allowed to have membership functions to all clusters rather than having a distinct membership to exactly one cluster. Kanungo *et al.* [17] proposed a filtering algorithm which begins by storing the data points in a $k-d$ tree. For each node of the tree maintaining a set of candidate centers, they will be pruned, or filtered as they are propagated to the node's children. Kanungo *et al.*'s implementation of the filtering algorithm is more robust than Alsabti's, because Alsabti's method relies on a less effective pruning mechanism based on computing the minimum and maximum distances to each cell.

Ming-Chuan Hung, explained about an Efficient k -Means Clustering Algorithm Using Simple Partitioning [20]. The Genetic algorithm described in [18] uses a multi step procedure. The authors refer to this procedure as a semi

supervised form of learning. In [19] a GA is used to solve the clustering problem for a data set of geographical data. Similarly, Yan-He Chen [22] describes about Genetic algorithm for Aerial image clustering. Bandyopadhyay, S., and Maulik describes Genetic algorithm for clustering using application of image classification by automatic evolution [24]. Rothlauf, F explained uses of genetic algorithm and evolutionary algorithms[23]. Yang, G., Reinstein define new genetic algorithm for optimisation problem[26].

3. Clustering Analysis

Cluster Analysis groups data objects based on information found in the data that describes objects and its relationship. The main goal of clustering is that a object within a group be similar to one another and different from the objects in other groups. The similarity within a group and the difference between the groups defines the clustering.

3.1 Types of Clustering

Cluster and clustering have different meaning. Cluster is a collection of objects and collection of clusters is referred as clustering and it has various types which have been mentioned below:

Hierarchical Cluster: Set of nested cluster which is organised as a tree. Each node in the tree refers clusters while parent is a main cluster and its children are subclusters and the root is a cluster containing all the objects.

Partitional Clustering: Set of data objects which is unnested non – overlapping subsets called clusters such that each data object have exactly one subset.

Exclusivive Clustering: It assign each object to a single cluster.

Overlapping Cluster: Object in a data set may belongs to more than one group. It also called non-exclusive clustering which is used an object is between two or more clusters and assigned to any of these clusters.

Fuzzy or Soft Clustering: Every objects belongs to every cluster and hold value in between 0 or 1. Sum of value of object is 1 then it is absolutely belongs to that group and sum of value of object is 0 then it is absolutely not belongs to that group.

Complete Clustering: Assign all objects to cluster.

Partial Clustering: Not assign all objects to cluster when some objects in data sets may not belongs to group.

3.2 Types of Cluster

Cluster groups define valid set of groups. There are several different cluster group are in practice such as

Well Seperated Cluster: A cluster is a set of objects which each object is very similar to other objects within a cluster and not similar to objects in other clusters.

Prototype Based Cluster: Prototype of a object is called centroid. A cluster is a set of objects in which each object is closer to prototype than to the prototype of other cluster.

Graph Based Cluster: If data is represented as a graph then group of object are connected to one another but no connection to outside group objects.

Density Based Cluster: If a cluster is a dense region of objects then it is seperated from other and it surrounded by a region of low density.

Conceptual Cluster: A cluster as a set of objects that share some property example: it share some closest centroid value.

3.2 Cluster Techniques

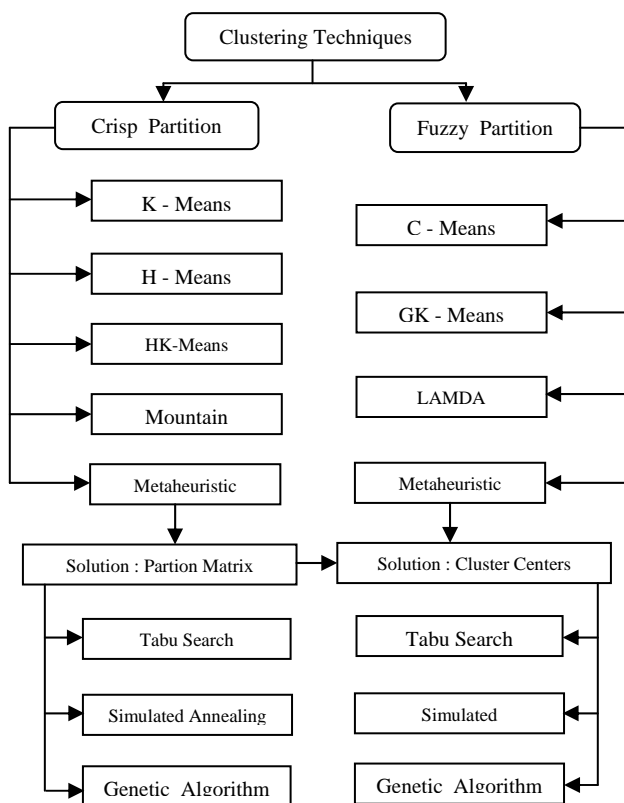


Figure. 1 Clustering Techniques

Bustos and Sellier categorise a number of lustering techniques (**Figure 1**) and show the K-means and the C-means to be amongst the most popular [25]. Cluster analysis has become an important technique in exploratory data analysis, pattern recognition, machine learning, neural computing, and other engineering. The clustering aims at identifying and extracting significant groups in underlying data. The four main classes of clustering algorithms are partitioning methods, hierarchical methods, density based clustering and grid-based clustering.

In the field of clustering, K-means algorithm is the most popularly used algorithm to find a partition that minimizes mean square error (MSE) measure. Although K-means is an extensively useful clustering algorithm, it suffers from several drawbacks. The objective function of the K-means is not convex and hence it may contain local minima. Consequently, while minimizing the objective function, there is possibility of getting stuck at local minima (also at local maxima and saddle point)[25]. The performance of the K-means algorithm depends on the initial choice of the cluster centers. Besides, the Euclidean norm is sensitive to noise or outliers. Hence K-means algorithm should be affected by noise and outliers. In addition to the K-means algorithm, several algorithms, such as Genetic Algorithm (GA) have been used for clustering.

4. Clustering Problem

The diversity of applications for clustering has lead to many problem definitions. The objective of all clustering algorithms is to divide a set of data points into subsets so that the objects within a subset are similar to each other and objects that are in different subsets have diverse qualities [11], [13], [20]. The fact that there are many different methods used to quantify the similarity and diversity of data points leads to the many different variations of the problem. For our comparison, we defined the clustering problem as the task of dividing an input data set into a desired number of subgroups so that the Euclidean distance between each data point and its corresponding cluster center is minimized. This is a very common method of defining the clustering problem. Total distances of each point to its cluster center is known as the total distance measurement of the clustering and is calculated as shown in (1).

$$F = \sum_{k=1}^n \sum_{x \in C_k} \sum_{a=1}^A (x_a - mk_a)^2 \quad (1)$$

In this formula n is the number of clusters, x represents a data point, Ck represents cluster k, mk represents the mean of the cluster k, and A is the total number of attributes for a data point. This formula simply calculates the Euclidean distance of each point in the input data set to its cluster center. It minimizes the total distance measurement of a clustering leads to an optimal cluster solution. This definition, like all clustering definitions, requires finding

an optimal collection of subsets for a group of data points. This class of problem is known to be NP-Hard. Work has been done to develop both approximate and exact solution algorithms for solving various clustering problems [23] but the solutions appear to be impractical, as either the number of data points in the input set or the number of clusters desired becomes large. As a result, there have been a wide variety of heuristic algorithms developed for this clustering problem. These algorithms do not guarantee any quality in the solutions they find but they do run in polynomial time with respect to the number of objects in the input data set and the number of desired clusters.

5. GENETIC ALGORITHMS

Genetic Algorithms (GA), first proposed by John Holland in the 1975s, are a category of EC that use concepts derived from evolution. Proper application of a GA finds a balance between exploration and exploitation of a given optimization problem's search space. First, a population of chromosomes is created and initialized. These chromosomes each contain a collection of genes and each gene has a value. A single chromosome is an encoded version of a solution to the problem that the GA is attempting to optimize. The GA performs exploration or exploitation of the problem's search space by evolving the population of chromosomes through a series of generations. During each generation of the GA, parent chromosomes are selected from the population. These parent chromosomes are combined to form children chromosomes and then the child chromosomes are mutated. In a generational type GA, an entirely new population for each generation is formed by creating multiple child chromosomes. For a steady state GA, the child chromosomes are used to replace members of the current population but a new population is not formed during each generation.

A very important step in the GA is the selection of parents for the next generation of chromosomes. In order to provide a guided search, which is appropriate for the given optimization problem, the selection of parents needs to be based on the quality of the solution that their chromosomes represent. A property called fitness is used to quantify the quality of a given solution and a fitness function is used to calculate the fitness value of each chromosome in a given population before parent selection is made. A variety of different selection methods are used by GA but they all use the principle that higher fit chromosomes are more likely to be chosen as parents. This fitness selection provides the GA direction for the search of an optimization problem's search space.

5.1 A Genetic Algorithm On Image Data

GA has been successfully implemented for various clustering problems using different chromosome encoding schemes and fitness functions. A GA performs clustering on an input set of data objects so that supervised learning can be applied to predict class labels in the second step. The input for the GA is a set of data objects that have both numeric and label attributes and a desired number of clusters. The goal of the GA is to produce clusters of data objects that minimize cluster dispersion and are as pure as possible in relation to the label attributes. The GA uses a two component fitness function where the first component measures within cluster variance using a distance metric and the second component measures the similarity of the labeled attributes of the data objects.

A very large input data set can be preprocessed to make a representative set that can be used by the algorithm for better time and space efficiency. In GA implement two alternate preprocessing methods for clustering algorithm such as

- The first Preprocessing method used random sampling to obtain a data set with fewer points. This reduced data set was then used in evaluating the fitness of the chromosomes.
- The second preprocessing method used summarization of the input data set and is based on the work presented in reference [21].

For this method, a grid is first constructed and then the input data set is applied to this grid. A single point location and corresponding weight is calculated for each region defined by the grid. The location of the representative point is chosen as the mean value of all the points in the region and the weight of the representative point is equal to the number of points that it replaces.

5.2 Pseudocode for Genetic Algorithm

Here, we take n number of chromosomes for n number of centroids, m number of Samples S_j and S_j^i is j^{th} sample assigned to i^{th} chromosome. We take maximum iteration to find best fitness solution (F). In search space at each iteration it simultaneously found objects on label as L_a and input attribute as I_a . First chromosome is consider as a parent chromosome and in each iteration it build child chromosomes. Below are GA rule to find cluster construction.

Genetic Algorithm is more efficient because final solution is evaluated by frequency of fitness value and samples matching value.

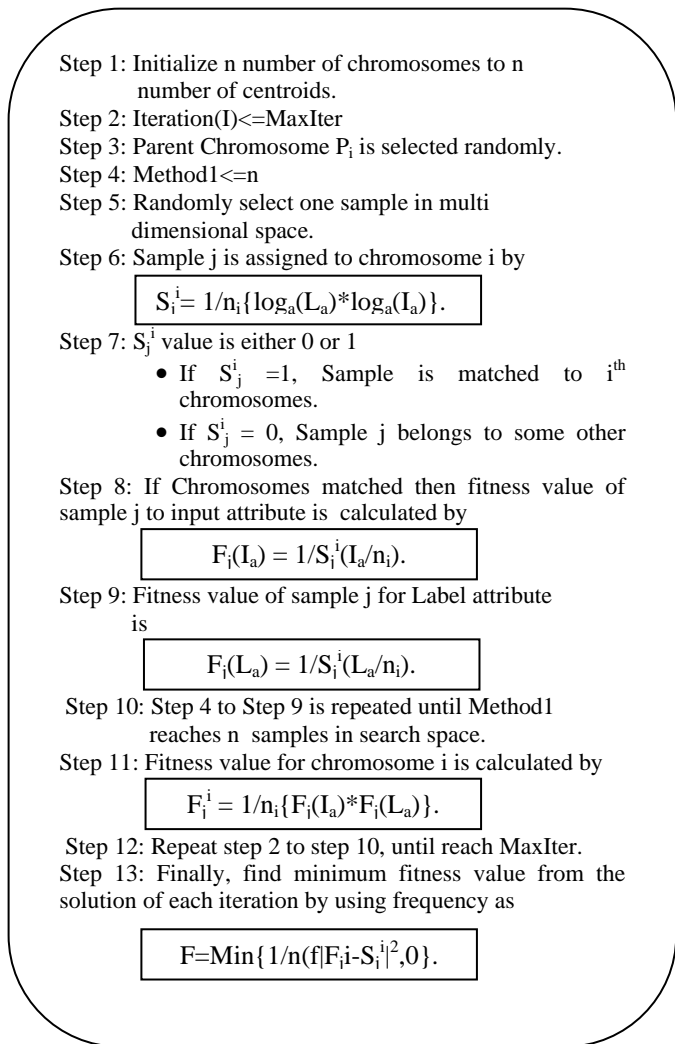


Figure. 2 Pseudocode for Genetic Algorithm

6. Performance Analysis

Clustering is an important task with applications in many fields. Heuristic algorithms are used for this task in an attempt to provide acceptable results, both in terms of solution quality and running time, because all of the non-trivial clustering problem variations are NP-Hard. In this paper comparison are done by using preprocessing techniques and performance measure for these algorithms are done on the basis of time. In K-means algorithm at different runs it produces poor results when the initial centroids are choosing randomly. It is important to realize that the choice of the initial centroid has a huge effect on the final result. K-means algorithm for multiple runs on large data sets is not work and it take more time to complete. For clustering on very large data sets, such as image data sets, the size of the associated databases makes it necessary to modify traditional GA because of their slow running times and combinations of input and label attributes.

Here, comparisons is done on two data sets first one is Artificial data sets created manually and it contains groups of color points. Genetic and K – means algorithm are applied to this data sets and clustering group of data by different colors. Second real group of image data sets is taken from google image. Two data sets contain six numerical attributes with values between 0 and 255. Each data set contains n points with the points centered around k cluster centers or chromosomes. The k cluster centers are first allocated by randomly and each of the six attributes values from a range of 0 to 255 is uniformly allocated in k- means algorithm and randomly allocated in Genetic algorithm. To calculate the minimum distance(D) between two cluster points is defined by D/r, where r is a variable used to define cluster. This process is repeated until it reach data points and clusters. Below are some sample images are done by experiments.

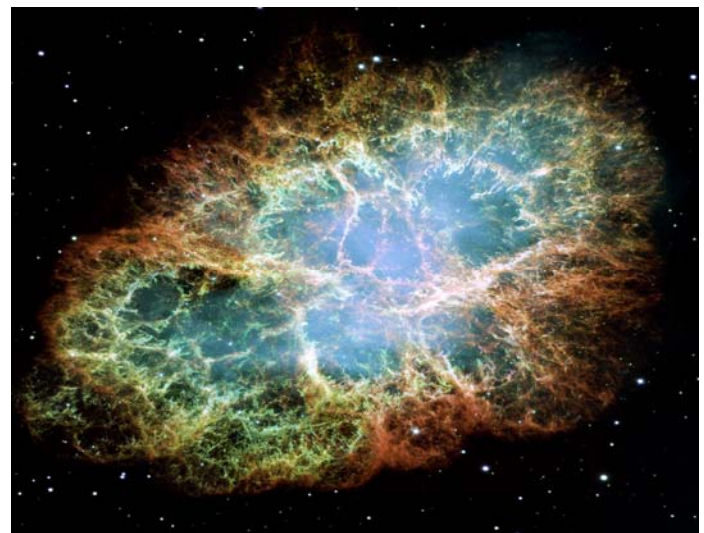


Figure. 3 Original Image taken from Google.

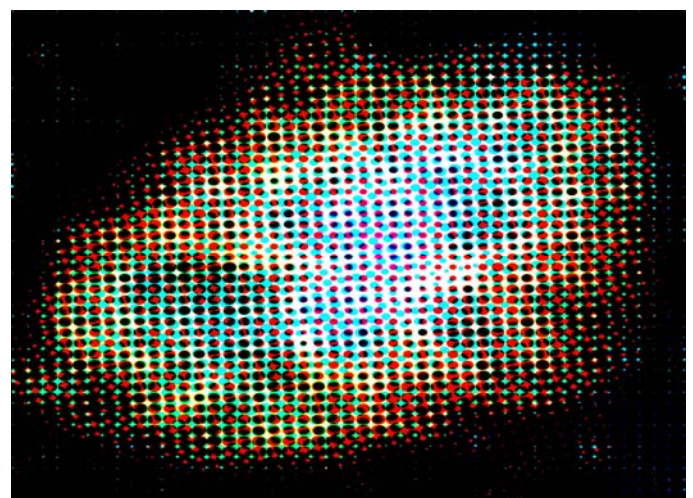


Figure. 4 Using GA to group red space.



Figure. 5 Group red space using GA view in black and white.

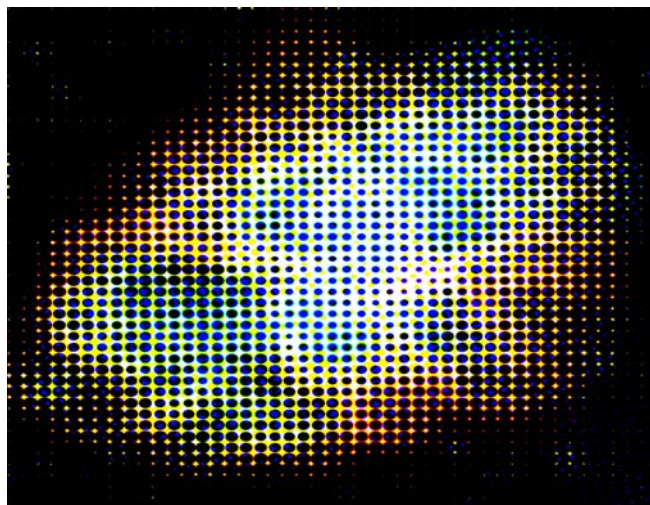


Figure. 8 Using GA to group blue space.

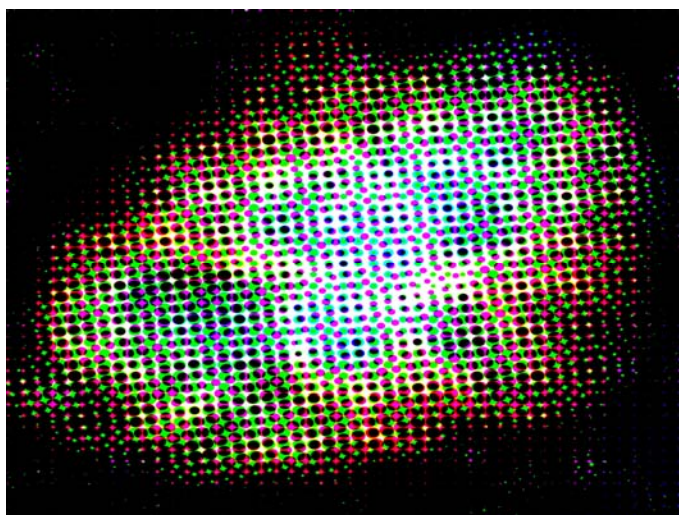


Figure. 6 Using GA to group green space.

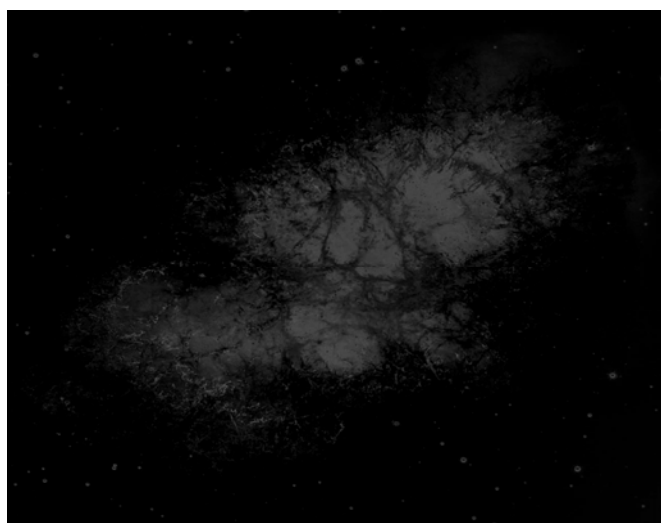


Figure. 9 Group blue space using GA view in black and white.

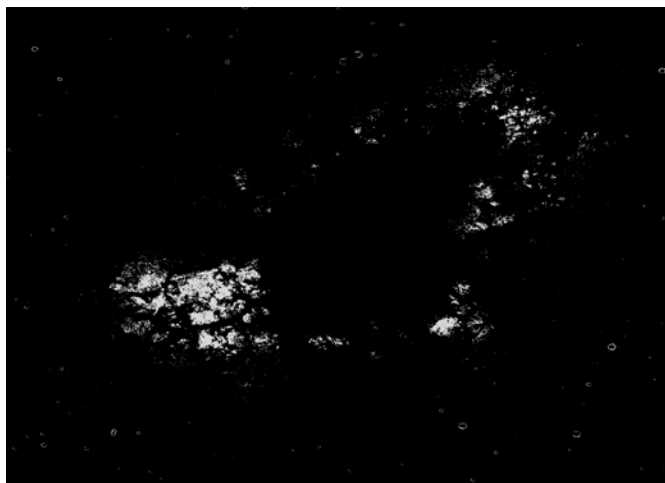


Figure. 7 Group green space using GA view in black and white.

Table 1: Results On Artificial Data Set

<i>Input set</i>	<i>Algorithm</i>	<i>Running Time(sec)</i>	<i>Distance Measurement</i>
3 Centroids 1000 points	GA	230	5.23×10^4
	K-means	219	7.89×10^4
5 Centroids 10000 points	GA	457	4.56×10^5
	K-means	512	2.34×10^6
7 centroids 20000 points	GA	890	8.27×10^8
	K-means	2023	1.10×10^{10}

Table 2: Results On Real Image Data Set

Input set	Algorithm	Running Time(sec)	Distance Measurement
3 Centroids 500 images	GA	153	3.87×10^5
	K-means	127	7.56×10^5
5 Centroids 1000 images	GA	296	8.90×10^5
	K-means	313	2.23×10^6
7 centroids 2000 images	GA	789	9.27×10^7
	K-means	1567	3.43×10^9

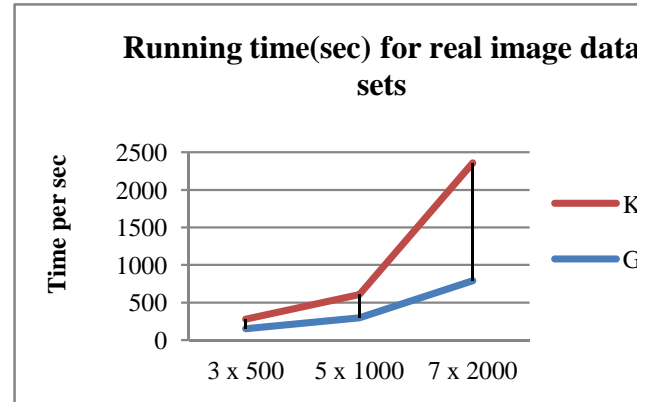


Figure. 12 Performance Analysis for running time using Artificial Data Sets

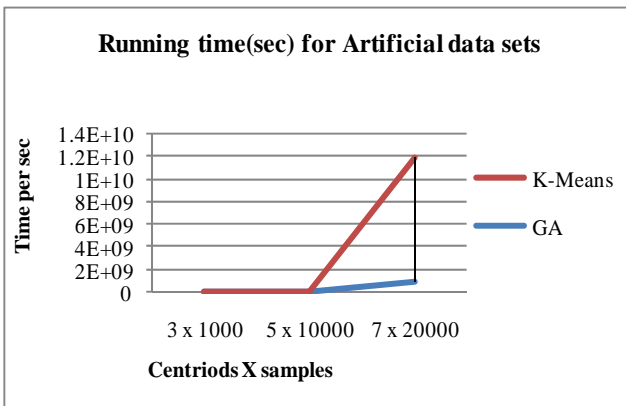


Figure. 10 Performance Analysis for running time using Artificial Data Sets

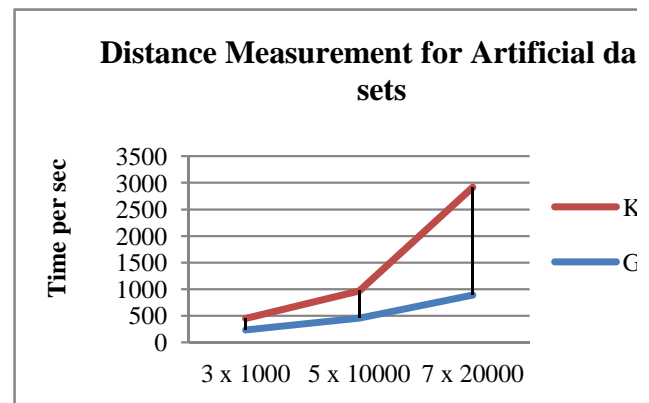


Figure. 13 Performance Analysis for distance measurement using Artificial Data Sets

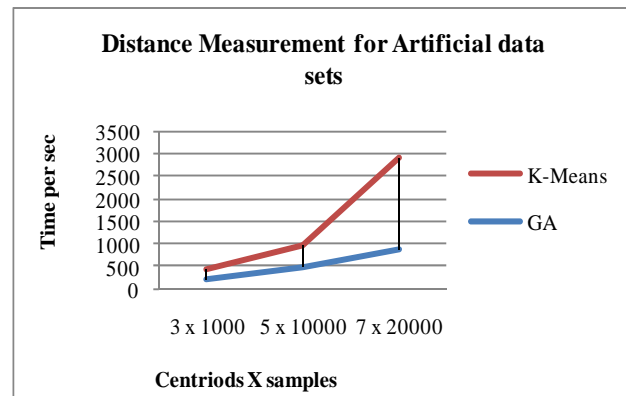


Figure. 11 Performance Analysis for distance measurement using Artificial Data Sets

7. Conclusion

The clustering performance of the two algorithms has been evaluated using real world and artificial data sets. The satisfactory results have demonstrated the effectiveness of the two algorithms in discovering structures in data. The scalability tests have shown that the two algorithms are efficient when clustering. But for large complex data sets in terms of number of records and the number of clusters GA shows better result as compared to k-means algorithm. The summary preprocessing method that prevent the creation of representative points for regions that contained less than a certain minimum threshold of points. This refinement removes the negative effect that outlier points have on the clustering quality. It makes the GA run faster than k-means algorithm because there would be fewer points in the processed data set.

Reference

- [1]. Z. Huang, "Extensions to the k-means algorithm for clustering large data sets with categorical values," Data

- [2]. Mining and Knowledge Discovery, Vol. 2, 1998, pp. 283-304.
- [3]. J. R. Wen, J. Y. Nie, and H. J. Zhang, "Query clustering using user logs," ACM Transactions on Information Systems, Vol. 20, 2002, pp. 59-81.
- [4]. J. Banfield and A. Raftery, "Model-based gaussian and non-gaussian clustering," Biometrics, Vol. 49, 1993, pp. 15-34.
- [5]. J. L. Bentley, "Multidimensional binary search trees used for associative searching," Communications of the ACM, Vol. 18, 1975, pp. 509-517.
- [6]. D.A. Clausi, "K-means iterative fisher unsupervised clustering algorithm applied to image texture segmentation," Pattern Recognition, Vol. 35, 2002, pp. 1959-1972.
- [7]. F. X. Wu, W. J. Zhang, and A. L. Kusalik, "Determination of the minimum samples size in micro array experiments to cluster genes using K-means clustering," in Proceedings of 3rd IEEE Symposium on Bioinformatics and Bioengineering, 2003, pp. 401-406.
- [8]. K. Alsabti, S. Ranka, and V. Singh, "An efficient k-means clustering algorithm," in Proceedings of 1st Workshop on High performance Data Mining, 1998.
- [9]. R. C. Dubes and A. K. Jain, "Algorithms for Clustering Data", Prentice Hall, 1988.
- [10]. E. R. Ruspini, "A new approach to clustering," Inform. Contr., vol. 19, pp. 22-32, 1969.
- [11]. L. Abul, R. Alhaji, F. Polat and K. Barker "Cluster Validity Analysis Using Sub sampling," in proceedings of IEEE International Conference on Systems, Man, and Cybernetics, Washington DC, Oct. 2003 Volume 2: pp. 1435-1440.
- [12]. J. Grabmeier and A. Rudolph, "Techniques of cluster algorithms in data mining," Data Mining and Knowledge Discover, 6, 2002, pp. 303-360.
- [13]. L. O Hall, I. B. Ozyurt, J. C. Bezdek, "Clustering with a genetically optimized approach," IEEE Transactions on Evolutionary Computation, 3(2), 1999, pp. 103-112.
- [14]. K. Jain, M. N. Murty, and P. J. Flynn, "Data clustering: a review," ACM Computing Surveys, 31(3), 1999, pp. 264-323.
- [15]. P. Berkhin, "A Survey of Clustering Data Mining Techniques" Kogan, Jacob; Nicholas, Charles; Teboulle, Marc (Eds.) Grouping Multidimensional Data, Springer Press (2006) 25-72.
- [16]. P. S. Bradley and U. M. Fayyad, "Refining initial points for k-means clustering," in Proceedings of 15th International Conference on Machine Learning, 1998, pp. 91-99.
- [17]. D. Judd, P. McKinley, and A. Jain, "Large-scale parallel data clustering," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 20, 1998, pp. 871-876.
- [18]. T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: analysis and implementation," IEEE Transactions on Pattern Analysis and Machine Intelligence, Vol. 24, 2002, pp. 881-892.
- [19]. Demiriz, K. P. Bennett, and M. J. Embrechts, "Semi-supervised clustering using genetic algorithms," R.P.I. Math Report No. 9901, Rensselaer Polytechnic Institute, 1999.
- [20]. M. Painho and F. Bação, "Using genetic algorithms in clustering problems," in Proceedings of GeoComputation Conference, 2000.
- [21]. J. Han and M. Kamber, "Data Mining: Concepts and Techniques," Morgan Kaufmann Publishers, 2000.
- [22]. W. DuMouchel, C. Volinsky, T. Johnson, C. Cortes, and D. Pregibon, "Squashing flat files flatter," in Proceedings of the ACM SIGKDD Conference on Knowledge Discovery and Data Mining, 1999, pp. 6-15.
- [23]. Ya-Wei Ho; Chih-Hung Wu; Chih-Chin Lai, "Aerial image clustering using genetic algorithm," IEEE transactions on Pattern Analysis and Machine Intelligence, Vol. 24, 2009.
- [24]. P. K. Agarwal and C. M. Procopiuc, "Exact and approximation algorithms for clustering," in Proceedings of the ninth annual ACM SIAM symposium on Discrete algorithms, 1998, pp. 658-667.
- [25]. Venkatesh Katari, Suresh Chandra Satapathy, JVR Murthy, PVGD Prasad Reddy, "Hybridized Improved Genetic Algorithm with Variable Length Chromosome for Image Clustering", IJCSNS International Journal of Computer Science and Network Security, VOL.7 No.11, November 2007.
- [26]. Yang, G., Reinstein, L.E., Pai, S., Xu, Z., Carroll, and D.L., "A new genetic algorithm technique in optimization of prostate implants". Medical Physics, 35(5), pp.104-112.



R. Balakrishnan MSc., M.Phil., Phd.,

He is a Phd Research Scholar in Manonmaniam Sundaranar University, Tirunelveli, Tamilnadu, India. He is working as a Assistant Professor and Head of the department in VLB Janaki Ammal Arts and Science College, Coimbatore. He has 12 years of experience in teaching line 6 years of experience in research. He conducted International, National Conference and he presented paper in International, National Conference and Journals His Interest areas are Data Mining, Image Processing, Current research project Genetic Algorithm using , Image Processing.



U. Karthick Kumar MSc., DCA., MCA., M.Phil.,

He is a Post Graduate with M.Phil from Bharathiar University, Coimbatore, Tamilnadu, India. He is working as a Assistant Professor in VLB Janaki Ammal Arts and Science College, Coimbatore. He has three years of experience in research. He presented paper in International, National Conference and Journals. His Interest areas are Grid Computing, Data Mining, Image Processing, Mobile Computing and Data Structures. Current research project Ant colony, PSO, Priority Based Pheromone Algorithm, Fair Scheduling, Sensor Network, Genetic Algorithm.