

## Optimal Data Analysis Based on Unsupervised Learning New Projection K-Means Initialization Clustering Algorithm and EMCA

Jyoti Kurmi, M. Tech. Scholar, Department of CSE, SVCST, RGPV, Bhopal, India; jyotikurmi11@gmail.com  
Amit Thakur, Asst. Professor in Department of CSE, SVCST, RGPV, Bhopal, India; amit.svcst@gmail.com

**ABSTRACT:** - Data mining could be a method of extracting desired and helpful data from the pool of information. Clusterin data processing is that the grouping of information points with some common similarity. Cluster is a vital aspect of information mining. It simply clusters the information sets into given no. of clusters. Various no. of ways are used for the information cluster among that K-suggests that is that the most generally used cluster formula. During this paper, we've briefed within the kind of a review work done by completely different researcher's victimization K-means cluster formula. As a partition primarily based cluster algorithmic program, K-Means is wide employed in several areas for the options of its efficiency and simply understood. However, it's documented that the K-Means algorithmic program could get suboptimal solutions, looking at the selection of the initial cluster centers. During this paper, they propose a projection-based K-Means low-level formatting formula. The planned formula initial uses standard mathematician kernel density estimation techniques to search out the extremely density information areas in one dimension. Then the projection step is to iteratively use density estimation from the lower variance dimensions to the upper variance ones till all the scale square measure computed. Experiments on actual datasets show that our technique will get similar results compared with different standard ways with fewer computation tasks.

**KEYWORDS:** Data mining, Data sets, clustering, clustering method, K-means clustering, unsupervised learning.

### I. INTRODUCTION

Data mining consists of extract, transform, and load dealings information onto the information warehouse system; data processing includes anomaly detection, association rule learning, classification, regression, summarization, and cluster. Data processing is one of the most vital analysis fields that are because of the expansion of each component and package technologies that has imposed organizations to depend heavily on these technologies. Data processing ideas and strategies may be applied in varied fields like promoting, medicine, property, client relationship management, engineering, web mining, etc. varied cluster algorithms consistent with totally different techniques are designed and applied to numerous data processing issues with success. During this paper, bunch analysis is completed by exploitation easy k mean cluster and changed k mean cluster. Standardization and classification is a very important preprocessing step in to standardize the values of all variables from dynamic vary into specific

vary. Cluster analysis is a kind data processing technique that is used to search out information segmentation and pattern info. By cluster, the information individuals get the information distribution, observe the character of every cluster, and build additional study on explicit clusters. The aim of cluster analysis is that the objects during a cluster ought to be kind of like each {other} and totally different from the objects in other groups. Bunch is way higher once there's larger similarity at intervals a gaggle and larger the distinction between the groups. Thus we will say that information has got to be used with the rule to extract helpful data from it. Varied bunch algorithms consistent with totally different techniques are designed and applied to numerous data processing issues with success. Describes the varied data processing techniques that permit extracting unknown relationships among the information things from massive data assortment that are helpful for deciding. The wide-spread use of distributed data systems ends up in the development of huge information collections in business, science and on the Web [1]. These information collections contain a wealth of data, that but have to be discovered. Businesses will learn from their dealings information additional regarding the behavior of their customers and thus will improve their business by exploiting this information. Science will get from data-based information (e.g. satellite data) new insights on analysis queries. Internet usage data is analyzed and exploited to optimize data access. Therefore data processing generates novel, unknown interpretations of information [2]. In recent years, there's a tremendous increase in the usage of the web. The usage of the web generates millions of information. This information is gaining its size because the year passes. The information is generated at a record rate on a daily basis. To research that information and cluster into a cluster is a tedious task. The problem additionally lies in storing and retrieving of information. The analysis of those information points into a totally different cluster is additionally a difficult task. Researchers have calculable that the quantity of data within the world doubles every twenty months. But data cannot be used directly. Its really worth is expected by extracting data helpful for call support. In most areas, information analysis was historically a manual method. Once the dimension of information manipulation and exploration goes on the far side human capabilities, individuals search for computing technologies to modify the method [3].

### Clustering Method

Clustering could be a method of grouping information objects into disjointed clusters in order that the

information within the same cluster are similar, however, information happiness to take issuing completely different cluster differ. A cluster is collections of information objects that are like each {other} area unit in the same cluster and dissimilar to the objects are in other clusters. The demand for organizing the sharp increasing information and learning valuable data from information, that makes agglomeration techniques are widely applied in several application areas like AI, biology, client relationship management, information compression, data processing, data retrieval, image process, machine learning, marketing, medicine, pattern recognition, psychology, statistics and then on. Cluster analysis could be a tool that's accustomed observes the characteristics of clusters and to target a selected cluster for any analysis. Agglomeration is unattended learning and doesn't place confidence in predefined categories. In agglomeration, we tend to live the unsimilarity between objects by activity the space between every combination of objects. These measures embrace the Euclidian [4].

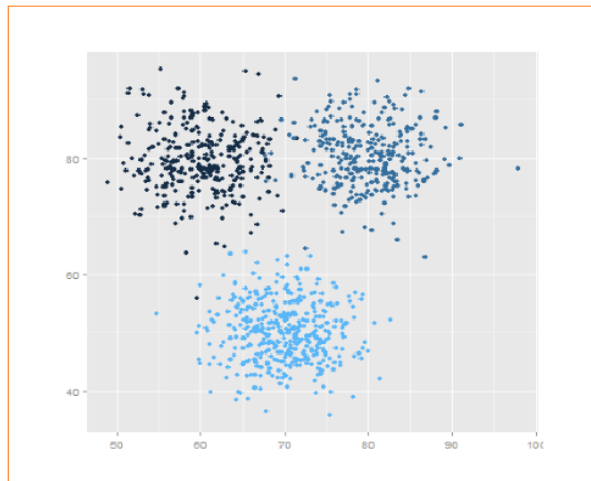


Fig.1 three clusters in clustering

## II. RELATED WORK

Shafeeq et al. [5] present a changed K-means algorithmic rule to boost the cluster quality and to mend the optimum range of clusters. As an input range of clusters (K) given to the K-means algorithmic rule by the user. However, within the sensible state of affairs, it's terribly tough to fix the number of clusters prior to. The strategy projected during this paper works for each of the cases i.e. for a celebrated range of clusters prior to likewise as an unknown range of clusters. The user has the flexibleness either to mend the range the number of clusters or input the minimum number of clusters needed. The new cluster centers are computed by the algorithmic rule by incrementing the cluster counters by one in every iteration until it satisfies the validity of cluster quality. This algorithmic rule can overcome this drawback by finding the optimum range of clusters on the run.

Soumi Ghosh et al. [6] present a comparative discussion of 2 cluster algorithms particularly the center of mass-based mostly K-Means and representative object-based FCM (Fuzzy C-Means) cluster algorithms. This discussion is on the premise of performance analysis of the potency of cluster output by applying these algorithms.

F.A. Ramadan et al. [7] propose an economical increased k-means algorithmic rule to beat issues in existing k-means. Original means is known because of its ease, simplicity, speed of convergence and flexibility to thin information In spite of its large number of benefits, it suffers from sure disadvantages. These drawbacks are the formatting of centroids, problem to converge to native minimum i.e updating of centroids until a native minimum isn't fount & execution of recurrent whereas loops. All these issues are handled by the projected k-means cluster algorithmic rule. The enhanced algorithmic rule first assigns datasets to its highest center of mass and so work out distance with different centroids. In the next step, the 2 distances are compared and if the new distance tiny is little than the previous distance then the data point is touched to a new cluster otherwise it is small then it's allotted to the same cluster. This method can save a great deal of your time and improve the potency. This algorithmic rule uses 2 new functions. The first one is distance () perform that's accustomed work out distance between every data point and its nearest cluster head. The second is distance new () perform accustomed work out distance between data points and different remaining clusters. The experimental result shows that the improved k-means algorithmic rule is far quick and economical than the first k-means.

Binu Thomas et al. [8] gave a comparative analysis between the k-means cluster algorithmic program and the fuzzy cluster algorithmic program. In this paper, the research worker additionally discusses the benefits and limitations of fuzzy c-means algorithmic programs means may be a partial primarily based cluster algorithmic program whereas Fuzzy means is not a partial based cluster algorithmic program. Fuzzy C means principally works in 2 methods. In the initial method, cluster centers are calculated and in second the information points are assigned to the calculated cluster center with the assistance of Euclidian distance. This method is nearly just like typical k-means with a touch distinction. In fuzzy means algorithmic program membership worth starting from zero to one is assigned to knowledge item in cluster. 0 membership indicates that the information purpose isn't a member of cluster whereas one indicates the degree to that information represents a cluster. The problem round-faced by the fuzzy c-means algorithmic program is that the ad of membership worth of information points in every cluster is restricted to one. Algorithm conjointly face drawback in addressing outliers. On the opposite hand comparison

with k-means shows that the fuzzy algorithmic program is economical in getting hidden patterns and information from natural data with outlier points.

Kohei Arai et al. [9] have projected hierarchical k-means which mixes k-means and hierarchical algorithmic program. The strategy executes k-means for a few mounted ranges of times then applies for the hierarchical algorithmic program on centroids obtained as a result of executions of k-means. The centroids, therefore, obtained from hierarchical algorithmic programs are then used as initial centroids for Kmeans. However, authors have recommended that their technique works higher (in terms of speed) as compared to ancient k-means for advanced cluster task (large numbers of information set and lots of dimensional attributes)

T. Gonzalez et al.[10] technique picks up an initial center of mass arbitrarily and also the remaining centroids are selected because the information that has the best minimum-distance to the antecedently designated centroid. This technique was originally developed as a 2-approximation to k-center cluster drawback.

Ismail Bin Mohamad et al. [11] applying standardization before cluster ends up in higher quality, economical and correct cluster results. The author has experimented on min-max, z-score and decimal scaling techniques and complete that among the 3 techniques, z-score provides the best result for infectious diseases dataset with improved accuracy over ancient k-means. But the author has commented that the choice of standardization technique ought to be tired in accordance with the character of the chosen dataset.

Rakesh Kumar et al. [12] planned a ranking mechanism that uses the many ratings of a review and calculates the mixture score of the product. The ranking of varied products is completed by suggests that of their reviews rating through rank selection methodology. The planned product-ranking approach victimization reviews rating establish the highest list of products and facilitate the client in selecting the simplest product. During this framework, the collected information is preprocessed and transformed for feature choice. When omitting the unimportant options the classification method train the information set to induce the ultimate model. Currently, the ranking approach picks the highest k-products. The planned approach considerably reduces the user time in choosing the proper product

Utkarsh Gupta et al. [13] planned a unique recommender system that supported a hierarchical cluster formula. The Item specific or user-specific data is classified into a collection of clusters victimization hierarchical cluster formula referred to as Chameleon. Following this, a legal system is employed to predict the rating of a selected item given by users. The method started with the set of

users with their options, supported that cluster is completed victimization ranked cluster formula. Then for a given item and a user, the mapping is completed to predict rating The prediction is completed by mapping a user into a selected cluster then the selected theme is applied for all users present therein cluster for the particular item. The performance of Chameleon based mostly recommender system is evaluated by comparing it with an existing technique supported K-means cluster formula. The results showed that Chameleon based mostly primarily recommender system reduce errors considerably as compared to K-means based Recommender System. The dataset used could be a motion-picture show rating dataset with a sample of 80k ratings with data regarding users and things. The variety of users is 943, with feature set (age, gender, occupation, pin code). The quantity of things is 1682 with feature set (release year, motion-picture show type). The planned approach is best than the prevailing K-Means based mostly approach in terms of low Mean Absolute Error.

Joy deep Das et al. [14] present a Recommender System supported information cluster techniques. This approach affects the quantifiability drawback related to the advice task. Totally different vote systems' algorithms are wont to mix opinions from multiple users for recommending things of interest to the new user. During this work, the author used the DBSCAN cluster formula for a bunch the users. Betting on the cluster to that the item belongs vote algorithms suggest things to the user. The thought behind this approach is "clusters -then apply voting" that partitions the users of the RS into teams then apply the advice formula one by one to every cluster. The planned system recommends an item to a user of a cluster supported rating statistics of the opposite users of that cluster. This approach avoids computations over the whole information, limits it to the targeted information and reduces the period of time of the formula. The formula is tested on the Netflix prize dataset. Netflix with 17770 rating files specified one per picture is taken into account. The picture rating file consists of the rating data with the attribute set (movie id, year of unfairness, title, average rating, genre) given by the shoppers to its picture. The rating of every picture given by all the shoppers is employed to calculate a median rating. The system recommends, per the user's preference of picture genres. For choosing the foremost well-liked things in an exceedingly cluster, a vote primarily based formula is applied one by one to the clusters.

H. Altay Guvenir et.al [15] have planned a brand new classification formula VFI5 and have applied to a drawback of medical diagnosis of erythematic squalors. There are several authors WHO have used a medical specialty dataset from UCI (the University of CA at Irvine) ranging from his work wherever he applied his new developed formula VFI5. This represents an idea description by a group of feature intervals. The

classification of a brand new instance is predicated on a vote among the classification created by the values of every feature one by one. All training examples are processed quickly. The VF15 formula constructs intervals for every feature from the training examples. For every interval, one price and therefore the votes of every category therein interval are maintained. Thus, an interval could represent many categories by sorting the vote for every category. This formula has obtained ninety six.25% of classification accuracy.

**III.MATLAB Tool**

It is simulating on mat lab 2013a and for this work, we use Intel 2.4 GHz Machine. Mat lab setup may be a high-level language and technical calculate and interactive surroundings for algorithmic program development, information visualization, records analysis, and numeric computation Mat lab may be a software system program that permits you to try to information manipulation and visualization, calculations, mathematics, and programming. It is wont to do terribly easy moreover as terribly subtle tasks. Image procedure function provides a comprehensive set of reference normal methods and graphical tools for image process, analysis, visualization, and algorithmic program development. You'll be able to perform image improvement, image declaring, feature detection, noise decrease, image segmentation, abstraction transformations, and image registration. Several functions within the tool case are multithreaded to require good things about multicourse and digital computer computers. The Performance analysis of MATLAB version R2013a i.e. used for this thesis simulation results of the image process provides processor optimized libraries for quick execution and image computation. It uses its JIT (just in time) compilation technology to produce execution speeds that rival ancient programming languages. It may also more advantage of multi-core and digital computer computers, MATLAB give several multi-rib algebra and numerical perform. These functions mechanically execute on multiple procedure thread during a single MATLAB session, enabling them to execute quicker on multicourse computers.

**IV RESULT ANALYSIS**

Data mining using a k-means clustering-based cluster center. Find minimum error of medical dataset analysis and the best possible solution. The results of the graph support IDCbreastcancer\_dataset Experimenting with an analysis result support the use of IDCbreastcancer\_dataset and the first set of random values. KMICA is an additional error rate compared to EMCA and the analysis of the time of use of KMICA and the square measurement EMCA takes average time. EMCA is better compared to KMICA due to the redundancy of information, but EMCA is minimal redundancy. The IDCbreastcancer\_dataset analysis is illustrated in Figure: 5.2 below. Experimenting with the

analysis of the results supports the use of Instantyeast\_Dataset and the first defined random values. KMICA is an additional error rate compared to EMCA and the analysis of the time of use of KMICA and the square measurement EMCA takes average time. EMCA is better compared to KMICA due to the redundancy of information is more however EMCA is minimal redundancy. The instantyeast\_Dataset analysis is illustrated in Figure: 5.3 below.



Fig.2 results analysis between KMICA & EMCA

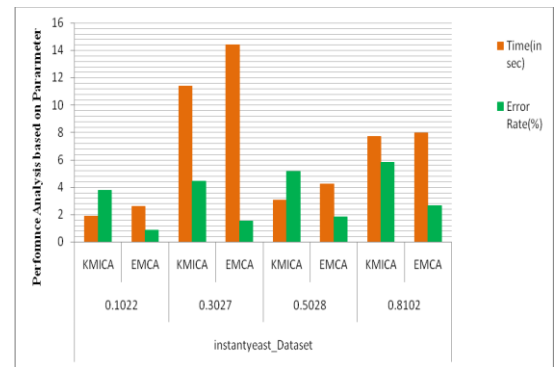


Fig.3 results analysis between KMICA & EMCA

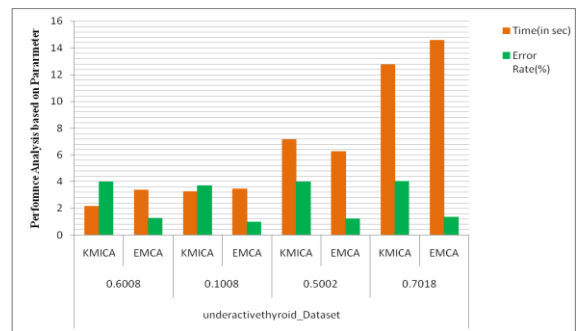


Fig4 results analysis between KMICA & EMCA

The results of the three-way experiment support the underactive use of the thyroid dataset and the first defined random values. KMICA is an additional error rate compared to EMCA and the analysis of the time of use of KMICA and the square measurement EMCA takes average time. EMCA is better compared to KMICA due to the redundancy of information is more however EMCA is

minimal redundancy. The underactivethyroid\_Dataset analysis is illustrated in Figure: 5.4 below.

## V. CONCLUSION

Clustering algorithms share some vital common problems that require being self-addressed to create them eminently. Some problems are thus present that they're not even specific to unsupervised learning and may be thought of as an area of an overall data processing framework. Different problems are resolved in sure algorithms we conferred. In fact, several algorithms were specifically designed to handle a number of these problems and k-means are concentrated on these problems which may be self-addressed in the next analysis. In this paper varied machine learning techniques employed in varied health care sectors are mentioned. Because of the big volume and complexness of information, there's a desire to method these data. For this medical data processing will facilitate to arrange some strategies for identification and deciding activates. In this, we've got centered on the utilization of machine learning techniques for medicine classification in past years. Data mining using a k-means clustering-based cluster center. Find the minimum error of medical dataset analysis. KMICA is an additional error rate compared to EMCA and the analysis of the time of use of KMICA and the square measurement EMCA takes average time. EMCA is better compared to KMICA due to the redundancy of information is more however EMCA is minimal redundancy.

## REFERENCES

- [1]. S. Anupama Kumar and M. N. Vijayalakshmi "Relevance of data mining techniques in identification sector", International Journal of Machine Learning and Computing, Volume 3, Issue 1, February 2013.
- [2]. Fayyad, U. Data Mining and Knowledge Discovery: Making Sense Out of Data. IEEE Expert, v. 11, no. 5, pp. 20-25, October 1996
- [3]. E. A. Khadem, E. F. Nezhad, M. Sharifi, "Data Mining: Methods & Utilities", Researcher2013; 5(12):47-59. (ISSN: 1553-9865).
- [4]. D. Napoleon, P. Ganga Lakshmi "An Efficient K-Means Clustering Algorithm for Reducing Time Complexity using Uniform Distribution Data Point" IEEE 2010, pp. 42-45.
- [5]. Shafeeq, A., Harsha, K., iDynamic Clustering of Data with Modified K-Means Algorithm, International Conference on Information and Computer Networks, vol. 27,2012
- [6]. Soumi Ghosh, Sanjay Kumar Dubey, Comparative Analysis of K-Means and Fuzzy C-Means Algorithms, International Journal of Advanced Computer Science and Applications, Vol. 4, No.4, 2013.
- [7]. [7] FAHIM, SALEM A.M, TURKEY F.A, RAMADAN M.A "An efficient enhanced k-means clustering algorithm" Journal of Zhejiang University SCIENCE A ISSN 1009-3095 (Print); ISSN 1862-1775 (Online)
- [8]. Raju G, Binu Thomas, Sonam Tobgay and Th. Shanta Kumar "Fuzzy Clustering Methods in Data Mining: A comparative Case Analysis" 2008 International Conference on advanced computer theory and engineering,2008 IEEE
- [9]. Kohei Arai, Ali Ridho Barakbah, "Hierarchical K-means: an algorithm for centroids initialization for K-means ", Reports of the Faculty of Science and Engineering, Saga University, Vol. 36, No.1, 2007.
- [10]. T. Gonzalez, "Clustering to minimize the maximum intercluster distance". Theoretical Computer Science, Vol. 38, pp. 293-306, 1985.
- [11]. Ismail Bin Mohamad, Dauda Usman, "Standardization and Its Effects on K-Means Clustering Algorithm", Research Journal of Applied Sciences, Engineering and Technology, Vol. 6, 2013.
- [12]. Rakesh Kumar, Aditi Sharan, Payal Biswas, "Framework for Ranking Products Using Ranked Voting Method", 2016, Second International Conference on Computational Intelligence & Communication Technology© 2016 IEEE DOI 10.1109/CICT.2016.138,
- [13]. Utkarsh Gupta and Dr. Nagamma Patil, "Recommender System Based on Hierarchical Clustering Algorithm Chameleon", 2015 IEEE International Advance Computing Conference (IACC), 978-1-4799-8047-5/15/\$31.00 c\_2015 IEEE.
- [14]. Joydeep Das, Partha Mukherjee, Subhashis Majumder, and Prosenjit Gupta, "Clustering-Based Recommender System Using Principles of Voting Theory", 2014 International Conference on Contemporary Computing and Informatics (IC3I) @2014 IEEE.
- [15]. Güvenir, H., Demiröz, G., &İlter, N. "Learning differential diagnosis of erythematous-squamous diseases using voting feature intervals". Artificial Intelligence in Medicine, 13(3) 147-165, 1998.