

# Comparative Data Analysis based K-Medoid Method and Effective K-Medoid Algorithm

Kapil Patidar<sup>1</sup>, Prof. Ankita Porwal<sup>2</sup>

Department of Computer Science and Engineering

Maxim Institute of Technology, RGPV, India

kapilpatidar9981@gmail.com, ankitaporwal19@gmail.com

**Abstract:-** Bunching examination is an unmistakable undertaking that tries to distinguish homogeneous gatherings of items in light of the estimations of their properties. The different grouping methods considered are displayed with a plan to upgrade them. K-medoid grouping calculations are broadly utilized for some down to earth applications. Unique K-medoid calculation select introductory centroid and medoids haphazardly that influence the nature of the subsequent bunches and some of the time it creates insecure what's more, exhausts groups which are negligible. The first k-implies calculation is computationally costly and requires time relative to the result of the quantity of information things, number of bunches and the quantity of emphases. The Efficient K-medoid calculation gave better outcomes just when the underlying segment was near the last arrangement. A few endeavors have been accounted for to take care of the group introduction issue. Creator proposes a technique that refines the underlying point prone to be near the methods of the joint likelihood thickness of the information. Unique K-medoid calculation select starting centroid and medoids haphazardly that influence the nature of the subsequent bunches and here and there it produces flimsy and purge bunches which are good for nothing. The first k-implies calculation is computationally costly and requires time corresponding to the result of the quantity of information things, number of groups and the quantity of cycles. Proficient k-Medoid bunching calculation has the precision higher than the k-medoid calculation. The new approach for the K-medoid calculation wipes out the lack of existing medoid. It initially ascertains the underlying centroid k according to prerequisites of clients and afterward gives better, powerful and great group without scarifying exactness. It creates stable bunches to enhance precision. It additionally lessens the mean square mistake and enhances the nature of grouping. As per our trial comes about, the proficient k-Medoid grouping calculation has the precision higher than the first one.

**Keywords:-** data mining, data warehouse, Clusters, k-Medoid, Principal Direction Divisive Partitioning, Hierarchical Clustering, Dataset.

## I. Introduction

Astounding advances in data find, planning power, data transmission, and limit capacities are engaging relationship to organize their diverse databases into data circulation focuses. Data warehousing is described as a method of united data organization and recuperation. Data warehousing, comparable to data mining, is a for the most part new term

notwithstanding the way that the thought itself has been around for an impressive period of time. Data warehousing addresses a flawless vision of keeping up a central vault of all progressive data. Centralization of data is anticipated that would enlarge customer get to and examination. Passionate creative advances are making this vision a reality for a few associations and similarly as shocking advances in data examination writing computer programs are allowing customers to get to this data straight forwardly. The data examination writing computer programs is the thing that sponsorships data mining. Frequently, the data to be mined is at first isolated from an endeavor data dissemination focus into a data mining database or data shop. There is some honest to goodness advantage if your data is currently some bit of a data stockroom. The issues of decontaminating data for a data appropriation focus and for data mining are on a very basic level the same. If the data has as of presently been cleansed for a data circulation focus, at that point it without a doubt won't require additionally cleaning in order to be mined. Moreover, we will have adequately kept an eye on extensive segments of the issues of data union and put set up help systems. The data mining database might be a true blue rather than a physical subset of our data conveyance focus, gave that the data stockroom DBMS can reinforce the additional advantage solicitations of data mining.

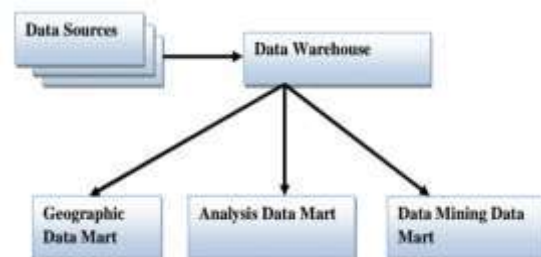


Figure 1 data warehouse and relation with other streams

Information mining data store expelled from a data stockroom a data appropriation focus is not a need for data mining. Setting up a broad data stockroom that joins data from various sources, decides data respectability issues, and weights the data into an inquiry database can be a gigantic task, now and again taking years and costing an immense number of dollars. You could, regardless, mine data from at least one operational or esteem based databases Information Sources Information Warehouse Geographic Information Mart Examination Data Mart Data Mining Data Store by simply removing it into a read-just database. These new database limits as a kind of data shop.

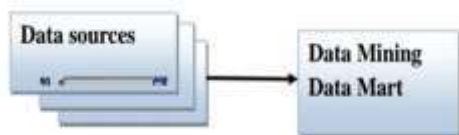


Figure 2 data warehouse and data mart

**How Does Data Mining Work:** While considerable scale information advancement has been creating separate trade and logical structures, data mining gives the association between the two? Data mining programming dismembers associations and cases in set away trade data considering open-completed customer request. A couple of sorts of logical writing computer programs are available: real, machine learning, and neural frameworks. Generally, any of four sorts of associations are searched for:

**Classes:** Stored data is used to discover data in predestined social events. For example, a diner system could mine customer purchase data to center when customers visit and what they normally organize. This information could be used to extend action by having step by step specials.

**Clusters:** Data things are gathered by associations or customer slants. For occasion, data can be mined to recognize business pieces or customer affinities.

**Associations:** Data can be mined to recognize affiliations. The ale diaper outline is an instance of familiar mining.

**Sequential cases:** Data is mined to imagine direct cases and examples. For occasion, an outside rigging retailer could predict the likelihood of a rucksack being gotten in perspective of a customer's purchase of snoozing packs and climbing shoes.

**Data mining involves five essential segments**

- I. Extract, change, and weight trade data onto the data conveyance focus structure.
- II. Store and manage the data in a multidimensional database structure.
- III. Give data access to business examiners and information development specialists.
- IV. Break down the data by application programming.
- V. Introduce the data in an accommodating association, diagram and table.

**Area Technology:** Every one of the investigations are performed on a 3-GHz Pentium PC machine with 1 GB or Higher, essential memory, running on Microsoft Windows/NT. Each one of the tasks is made in Microsoft Visual studio .net (VB or ASP or C#). Notice that we don't direct differentiation our aggregate number of runtime and those in some circulated reports running on the RISC workstations in light of the way that particular machine designs may differ fantastically on the incomparable runtime for similar figure. Or maybe, we realize their figuring's to the best of our understanding in light of the dispersed reports on the same machine and break down in a similar running condition. On the off chance that it's not all that much inconvenience in like manner observe that run time used here means the total execution time, that is, the period amidst information and yield, as opposed to CPU time measured in the tests in some written

work. I feel that run time is a more broad measure since it requires the total running venture eaten up as the measure of cost; however CPU time considers only the CPU's cost resource.

The application has no window based GUI.

- I. The application will work only for VB .net (VB or ASP or C#).
- II. The application relies upon Boolean connection rules
- III. This application is work for 30 things not more than that

The recognizing in order to group technique portrayed in this paper begins disconnected event groups from the earliest starting point plan of all around passed on events. Next, each event gathering is taken care of progressively by sub-parceling the social affair into a starting "best" figure of consistent bundles where the amount of gatherings is picked moderate while so far guaranteeing that the most outrageous compass of each bundle is not mishandled. At long last, the social occasion is taken care of using a best down k-medoid computation to find the best course of action of cluster delegates for the get-together. The rest of the pages of this paper may depict the method in banter ask. To begin with we will depict the k-medoid progression tolerating an event gathering has been found and early on courses of action of gatherings and representatives have been distributed. Next we will portray little settle ups to deal with the phenomenal circumstances where low quality packs are confined and events are doled out to raunchy gatherings. By then we will return to depict the procedure for confining the starting best figure of groups and their specialists given a subjective event bundle. Finally, we will depict the system for encircling event get-togethers given the starting scattering of all around geo-referenced events.

**II. Literature Survey**

**M. N. Vijayalakshmi et al. [1]** outline that different information mining methods like arrangement, bunching are apply on the understudy's information base. This examination can be utilized to empower the student and showing group increment the execution. These strategies can likewise be consolidated with other particular revelation model to build the limit of the model. In this paper clarify the numerous systems of information mining as indicated by Instructive information to outline another condition Result of this paper is that training framework can upgraded their execution by utilizing information mining systems. In this paper demonstrates that each technique has its own key range in which it performs precise. Bharat Chaudhari, Manan Parikh [2] speaks to relative investigation of grouping calculations utilizing weka tools. Clustering is a procedure in which information is partitioned into various groups concurring their usefulness. Information of one bunch is diverse to another group however inside that bunch information is homogenous. In this paper they think about the execution of bunching calculation in term of class savvy bunch building capacity of calculation. The results of this paper are that k mean is superior to other grouping algorithm (Hierarchical Clustering

calculation, Density based grouping calculation) however is create quality when I utilize substantial measure of information.

**Sharaf Ansari et al. [3]** .management framework. Understudy execution in college courses give a diagram of grouping calculations utilized as a part of information mining. They speak to a vital part in our life since I require much data (information) also, we realize that information mining is a procedure to extricating information and perceive the examples. In this paper they give an outline of some grouping investigation strategies, for example,

**Narendra Sharma et al.[4]** speaks to the correlation between different grouping calculation utilizing weka apparatus .There are different devices in information mining which are utilized to examination the information. They enable the clients to investigation the information in various measurement or points, sort it, and abridge the connections recognized. Weka is additionally an information mining device which is utilized for investigation the information. The fundamental goal is to demonstrate the examination of the diverse distinctive grouping calculations of weka and to discover which calculation will be most appropriate for the clients. Each calculation has their own particular significance and we utilize them on the conduct of the information, yet on the premise of this exploration I found that k-implies grouping calculation is least complex calculation when contrasted with different calculations.

**T. Saravanan et al [5]**.The creator in writing uses Principal Component Analysis for estimation diminishment and to find beginning bunch centers.

**N. Aggarwal et al. [6]**.In in any case data set is pre-taken care of by changing all data qualities to positive space and after that data is arranged and isolated into k proportionate sets and subsequently focus estimation of each set is taken as starting focus point.

**B. Barisi et al. [7]**. In writing a dynamic answer for k – Means is recommended that computation is laid out with pre-processor using layout authenticity record that thusly chooses the appropriate number of packs that assemble the capability for gathering in light of present circumstances.

**S. Na et al. [8]** .In a framework is proposed to make estimation self-sufficient of number of cycles that goes without from enrolling detachment of each data point to group centers on and on, saving running time and diminishing computational multifaceted nature.

**A. Shafeeq et al.[9]**. In the writing component infers calculation is proposed to improve the group quality and propelling the amount of gatherings. The customer has the versatility either to settle the amount of bundles or information the base number of gatherings required. In the past case it works same as k-Means calculation. In the last case the calculation figures the new gathering centers by enlarging the cluster incorporate by one every accentuation until the point that it satisfies the authenticity of gathering quality.

**Osama Abu Abba et al.[10]**.This paper is planned to study and look at changed information bunching calculations. The calculations under scrutiny are: k-implies calculation, various leveled grouping calculation, self-organizing maps calculation and desire boost grouping calculation. All these calculations are contrasted agreeing with the accompanying variables: size of dataset, number of groups, kind of dataset and sort of programming utilized. A few conclusions that are extricated have a place with the execution, quality, and precision of the bunching calculations.

**Manish Verma et al. [11]**. This paper surveys six sorts of bunching systems k-Means Clustering, Hierarchical Grouping, DB Scan bunching, Density Based Clustering, Optics, EM Algorithm. These grouping strategies are executed and dissected utilizing a bunching instrument WEKA. Execution of the 6 strategies are exhibited and thought about.

**Tajunisha et al.[12]**. In this paper, I have separated the execution of our proposed procedure with the present works. In our proposed strategy, I have used Principal Part Analysis (PCA) for estimation diminishing and to find the starting centroid for k-suggests. Next I have used heuristics approach to manage reduce the amount of detachment calculation to select the data point to gathering. By differentiating the results on iris data set, it was discovered that the results obtained by the proposed method are more fruitful than the present framework.

**D.Napoleon et al. [13]**. K-implies grouping calculation frequently does not work splendidly for high estimation, thus, to improve the capability, apply PCA on one of a kind data set and get a decreased dataset containing possibly uncorrelated factors. In this paper principal part examination and direct change is used for dimensionality decreasing and initial centroid is prepared, at that point it is associated with KMeans grouping calculation.

**Kehar Singh et al. [14]**. K-implies is extraordinarily outstanding on the grounds that it is sensibly essential and is computationally snappy and memory capable however there are distinctive sorts of constrains in k suggests estimation that makes extraction genuinely troublesome. In this paper I am discussing these limitations and how these obstructions will be removed.

**N. S. Chandolika et al.[15]**.This paper assesses execution to two well-known characterization calculations for assault order. Bayes net and J48 calculation are examined the key contemplations are to use data burrowing strategies gainfully for interference assault grouping.

### III. SIMULATION AND RESULTS ANALYSIS

**(a)Simulation:** This section contains data of apparatuses utilized while actualizing the proposed approach and some other customary strategies. The trials are performed on a 2-GHz Pentium PC machine with 512 megabytes fundamental memory, running on Microsoft Windows/NT. Every one of the projects is composed in Microsoft Visual studio .net(C#

7.0).And likewise K-Medoid calculations and Efficient K-Medoid calculations. This part additionally introduces the examination diagram by taking diverse parameters.

The tests are performed on a 2-GHz Pentium PC machine with 512 megabytes primary memory, running on Microsoft Windows/NT. Every one of the projects is composed in Microsoft Visual studio .net(C# 7.0). Notice that we don't straightforwardly contrast our supreme number of runtime and those in some distributed reports running on the RISC workstations on the grounds that diverse machine architectures may vary incredibly on the outright runtime for the same calculations. Rather, I actualize their calculations to the best of our insight in view of the distributed reports on the same machine and look at in the same running environment. If you don't mind likewise take note of that run time utilized here means the aggregate execution time, that is, the period in the middle of information and yield, rather than CPU time measured in the trials in some writing. I feel that run time is a more extensive measure since it takes the aggregate running time devoured as the measure of expense, while CPU time considers just the expense of the CPU asset. The examinations are sought after on both engineered and genuine information sets. The engineered information sets which I utilized for our examinations were created utilizing the methodology. I allude peruses to it for more subtle elements on the era of information sets. I report test results on two engineered information sets. In this information set, the normal exchange size and normal maximal possibly visit thing set size are situated to 4 and 5, individually, while the quantity of exchanges in the dataset is situated to 100 K. It is a meager dataset. The regular thing sets are short and not various. The second engineered information set we utilized is. The normal exchange size and normal maximal possibly visit thing set size are situated to 20 and 25, individually. There exist exponentially various regular thing sets in this information set when the bolster edge goes down. There are likewise really long successive thing sets and additionally a substantial number of short incessant things set in it. It contains bounteous blends of short and long regular thing sets.

**(b) Results Analysis:** Comparison On The Basis Of Varying Number OF Bunches and Records: I perform Comparative investigation between Existing calculation, K-Medoids and The Proposed Proficient K-Medoids calculation on the premise of two parameters i. Number of Records and Number of Clusters with Execution Time (in milliseconds).

**Cluster Based:** In the Cluster based similar examination. I perform Comparison amongst existing and proposed calculation with various Number of Clusters of Employee information - set and Execution Time (in milliseconds).K-Medoid and Efficient K-Medoid Algorithm with Number of Cluster and Execution Time.

**Portrayal:** Above figure indicates examination between K-medoid and Efficient K-Medoid Calculation.

As chart demonstrate that when number of groups is less, Efficient K-Medoid Calculation sets aside little opportunity to execute than the K-medoid. If the amount of gatherings is more than it is again real that Efficient K-Medoid Algorithm takes little time to execute than the K-medoid.

Number of clusters	Time taken to execute (In millisecond) K-Medoid Algorithm	Time taken to execute (In millisecond) Efficient K-Medoid Algorithm
2	21026	16343
3	42126	26922
4	66414	48533
5	76223	49342

Table 5.5.1 Comparison between K-Medoid and Efficient K-Medoid Algorithm with Number of Cluster and Execution Time

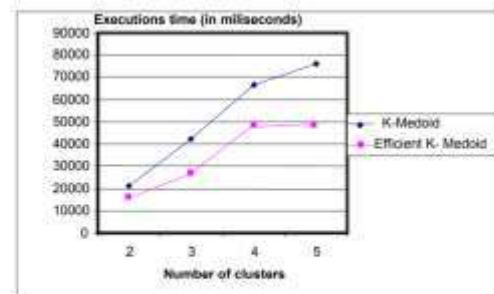


Figure 5.5.1 Graph Represent Number of Clusters and Execution Time for K-Medoids Algorithm and Efficient K-Medoid Algorithm

At the particular number of the records the execution time taken by Efficient K-Medoid Algorithm takes roughly little time than K-medoid.

Number of Records	Time taken to execute (In millisecond) K-Medoid Algorithms	Time taken to execute (In millisecond) Efficient K-Medoid Algorithm
300	95672	59735
400	123272	89332
500	139826	106243
600	170231	128338

Table 5.5.2 Comparison between K-Medoids Algorithm and Efficient K-Medoids Algorithm with Number of Records and Execution Time.

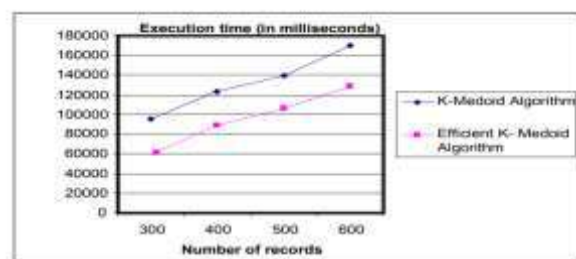


Figure 5.5.2 Graph Represent the Number of Records and Execution Time for K-Medoid Algorithm and Efficient K-Medoid Algorithm

**Record Based:** In the record based relative examination I perform Comparison amongst existing and proposed calculation with various Number of Records of Employee information set and Execution Time (in milliseconds). Diagram Represent the

Number of Records and Execution Time for K-Medoid Algorithm and Efficient K-Medoid Algorithm.

**Depiction:** Above figure indicates correlation between K-medoid and Efficient K-Medoid calculation. As graph show that when number of records is less or progressively, Efficient K-Medoid sets aside little opportunity to execute than the K-medoid.

## VII. CONCLUSION

The examination work predominantly manages execution time; it is watched that the correlation between k-medoid and productive k-medoid calculation demonstrates that when number of bunch is slight or more productive k-medoid calculation sets aside not as much opportunity to execute than the k-medoid calculation and same for the record case. According to the numerical examination occurs, the proposed framework is an effective grouping methodology. It can be associated with an extensive variety of kind of collection issues or joined with some other information digging strategies for getting more empowering comes about for applications. Future scope K-implies calculations, (Principal Direction Divisive Partitioning (PDDP) calculation. In new Approach of traditional parcel based grouping calculation, the estimation of k (the number of sought groups) is given as information parameter, paying little respect to dispersion of the information focuses. It is ideal to build up some measurable strategies to figure the estimation of k, depending on the information dissemination.

## REFERENCE

- [1]. A. Kumar and M. N. Vijayalakshmi Zhao, Kaidi and Liu, Bing, Tirpark, Thomas M. and Weimin, Xiao, "A Visual Data Mining Framework for Convenient Identification of Useful Knowledge", ICDM '05 Proceedings of the Fifth IEEE International Conference on Data Mining, vol.-1, no.-1, pp. - 530-537, Dec 2005.
- [2]. B. Chaudhari<sup>1</sup>, Manan Parikh<sup>2</sup>" A Comparative Study of clustering algorithms Using weka tools" International Journal of Application or Innovation in Engineering & Management (IJAIEM) Volume 1, Issue 2, October 2012.
- [3]. S.Ansari<sup>1</sup>, Sailendra Chetlur<sup>2</sup>, Srikanth Prabhu<sup>3</sup>, N. Gopalakrishna Kini<sup>4</sup>, Govardhan Hegde<sup>5</sup>, Yusuf Hyder<sup>6</sup>" An overview of clustering algorithms used in data mining" International Journal of Emerging Technology and Advanced Engineering Website: www.ijetae.com (ISSN 2250- 2459, ISO 9001:2008 Certified Journal, Volume 3, Issue 12, December 2013).
- [4]. N. Sharma , Aman Bajpai and Mr. Ratnesh Litoriya" Comparison the various clustering algorithms of weka tools" International Journal of Emerging Technology and Advanced Engineering Website: www.ijetae.com (ISSN 2250-2459, Volume 2, Issue 5, May 2012).
- [5]. T. Saravanan, "Performance Analysis of k-means with different initialization methods for high dimensional data" International Journal of Artificial Intelligence & Applications (IJAIA), Vol.1, No.4, October 2010.
- [6]. N. Aggarwal and Kriti Aggarwal, "A Mid- point based k -mean Clustering Algorithm for Data Mining". International Journal on Computer Science and Engineering (IJCSE) 2012.
- [7]. B. Barisi Baridam," More work on k-means Clustering algorithm: The Dimensionality Problem". International Journal of Computer Applications (0975 - 8887) Volume 44- No.2, April 2012.
- [8]. S. Na, Li Xumin, Guan Yong "Research on K-means clustering algorithm". Proc of Third International symposium on Intelligent Information Technology and Security Informatics, IEEE 2010.
- [9]. A. Shafeeq and Hareesha "Dynamic clustering of data with modified K-mean algorithm", Proc. International Conference on Information and Computer Networks (ICICN 2012) IPCSIT vol. 27 (2012) © (2012) IACSIT Press, Singapore 2012.
- [10].M. Verma, Mauli Srivastava, Neha Chack, Atul Kumar Diswar, Nidhi Gupta," A Comparative Study of Various Clustering Algorithms in Data Mining," International Journal of Engineering Research and Applications (IJERA) ISSN: 2248-9622 www.ijera.com, vol. 2, Issue 3, pp.1379-1384,May-Jun. 2012.
- [11].T. and Saravanan, "Performance analysis of k-means with different initialization methods for high dimensional datasets, "International Journal of Artificial Intelligence & Applications (IJAIA), vol. 1, no.4, pp.44-52,Oct. 2010.
- [12].N. ,S. Pavalakodi, "A New Method for Dimensionality Reduction using K-Means Clustering Algorithm for High Dimensional Data Set," International Journal of Computer Applications (0975- 8887),vol. 13, no.7, pp.41-46, Jan 2011.
- [13].K. Singh, Dimple Malik and Naveen Sharma, "Evolving limitations in K-means algorithm in data mining and their removal, "IJCEM International Journal of Computational Engineering &Management, vol. 12, pp.105-109,Apr. 2011.
- [14].C. V. D. Nandavadekar, "Comparative Analysis of Two Algorithms for Intrusion Attack Classification Using KDD CUP Dataset, "International Journal of Computer Science and Engineering (IJCSE), vol.1, pp.81-88,Aug 2012.
- [15].C. M. Bishop and Michael E. Tipping, "A Hierarchical Latent Variable Model For Data Visualization," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 20, no. 3, pp.281-293, Mar. 1998.