# K-Means and K-Medoids Algorithms Comparision on TB Data

Ashwini L P
B.E., (M.Tech), Department of CS&E,
Adhichunchanagiri Institute of Technology,
Chikkamagaluru

Sunitha M R
B.E., M.Tech., (Ph.D)., Department of CS&E,
Adhichunchanagiri Institute of Technology,
Chikkamagaluru

*Abstract*:- **In identifying and forming the group's dynamic data clustering is a challenge. Undirected knowledge discovery is the outcome of unsupervised learning usually. By victimization the similarity measures the cluster detection algorithmic rule finds the clusters of information that are nearly same to 1 another. Deciding the appropriate algorithmic rule to bring sensible clump result's a problem. Identifying the dataset provided parameters information and attributes information of the results changes in both K- medoids method and K- means method. The proposed method presents a comparison of each algorithmic rules specializing in the result of dataset on every algorithm**

*Keywords: K-medoids algorithm, Clustering method, K-means algorithm.*

## I. INTRODUCTION

The main objective of cluster objects is to search out objects that demonstrate homogeneity in the attributes and group them into one cluster, dissimilar attributes are grouped into another cluster based on different potentials. Clustering method area unit sometimes unattended learning techniques since they look after info while not noted or pre-classified labels, that's that the foremost distinction between cluster and classification. cluster algorithms unit sometimes classified based on Partitioning, that partitions the record set into initial partitions therefore utilize repetitive relocation technique to reinforce the value of clusters and terminates once it converges; Hierarchical-based, that decomposes the dataset into a hierarchy and will be either agglomerate or factious, the previous uses a bottom-up approach that starts with one object and builds clusters by adding similar objects whereas the latter uses a top-down approach that starts with one cluster and divides the dataset into smaller clusters; a Grid which is based on density builds clusters so that a minimum of particular threshold like neighbors for each different members inside each clusters and conjointly based on Model, which builds a cluster models first and finally teams the points of data which is according to the model. The k-means method and k-medoids method comprise partition based mostly techniques.

Choosing the suitable approach to research knowledge may even be tough while not correct understanding regarding the knowledge, the particular domain and additionally the expected result. during this paper, the comparative analysis among two connected approaches is reviewed. clump technique is one amongst powerful tools in data processing wont to discover information. Grouping knowledge according specific parameters eventually varied supported their knowledge and downside domain. There's no specific account all knowledge varieties. This analysis targeted on two techniques that area unit K-Means algorithmic rule and K-Medoids algorithmic rule. Cluster analysis, to boot referred to as unattended knowledge classification, could be an important subject in data processing [1]. Its aim is to partition a group of patterns into clusters of comparable knowledge points. K-Medoids being planned to urge every centers to be one amongst the aim itself. Even K-Medoids algorithmic rule is similar to K-means; variety of the info could also be appropriate to cluster by exploitation K-Medoids algorithmic rule. Many domains information being utilized to research these two techniques to avoid data specific bias.

## II. LITERATURE REVIEW

### DATA PREPROCESSING

Once the data is taken or loaded the data must be preprocessed. The available data is unstructured and need to be processed for further use. The data contains many unfilled or missing values, special characters and noise in it. The special characters and noise is removed manually by looking the data. In this paper the data is cleaned by filling the missing values or unfilled attributes values. The missing attributes values are filled by taking the mean of all the attributes values in that corresponding column or row. If more number of attributes values is missing then entire row or column is deleted, so that it reduces the data. There are many preprocessing techniques only two are used in this paper. Clean and reduce data preprocessing techniques are applied in our work.

### K MEANS CLUSTERING

K-means method (Mac Queen, 1967) is one in all the most unsupervised learning calculations that make out

Special Issue - 2017

International Journal of Engineering Research & Technology (IJERT)
ISSN: 2278-0181
NCETAIT - 2017 Conference Proceedings

the quality grouping draw back. The system takes once a transparent and simple approach to characterize a given informational index through associate unequivocal assortment of teams (accept k bunches) mounted from the sooner. K-centroids are characterized for every bunch. These centroids have to be compelled to be set in associate passing crafty strategy as a results of entirely extraordinary space causes entirely distinctive outcome. During this method, the upper determination is to position them the foremost extreme total as potential inaccessible from one another. The subsequent stride is to need every reason satisfaction to a given informational index and partner it to the closest center of mass. At the purpose once no style is incomplete, the initial step is finished associated an early gathering age is finished. Currently we'd like to re-figure k new centroids as centers of the teams ensuing from the past stride. Once the k new centroids are collected, a substitution excluding have to be compelled to be completed between indistinguishable informational index focuses and moreover the nearest new center of mass. A circle has been made. As a consequences of this circle we have a tendency to might see that the k centroids revision their space regular until no additional changes territory unit done. At last, this equation goes for limiting Partner in Nursing target work, amid this case a square blunder work. The target operation is

$$J = \sum_{j=1}^{K} \sum_{j=1}^{n} ||x_i \ (j) - Cx_j||^2$$

Where ∥ xi (j) – Cj ∥ two could be a chosen distance live between a knowledge purpose xi (j) and therefore the cluster center cj, is associate degree indicator of the gap of the n knowledge points from their individual cluster centers.

The formula consists of the subsequent steps:

1. Place K points into the area delineate by the objects that are being clustered. These points represent initial cluster centroids.

2. Assign each object to the cluster that has the closest center of mass.

3. Once all objects are assigned, cipher the positions of the K centroids. Repeat steps two and three till the centroids now not move. This produces a separation of the objects into teams from with the metric to be reduced are often calculated.

*K-MEDOIDS CLUSTERING*

K-medoids could be a cluster rule associated with the k-means rule and it is called medoid shift rule, each k-means clustering and k-medoids clustering algorithms are partitioned (breaking the dataset up into groups) and every decide to minimize sq. error, the gap between points tagged to be in associate passing cluster and some extent elect as a result of the middle of that cluster. In distinction to the k-means algorithm, k-medoids chooses information points as centers (medoids or exemplars). K-medoids is in addition a partitioning technique of cluster that clusters the information set of n objects into k clusters with k far-famed a priori. It

may be plenty of strong to noise and outliers as compared to k-means clustering method as a result of it reduce a complete of general attempt wise different similarities as an alternative of a complete square geometrician distance. Several of the distinction operate is implausibly affluent, but in an application. This tends to use the square Euclidian value. The center of medoid as a finite information that could be for a data purpose from any data set, overall average distinction to all information points is lowest such that it is the main centre set purpose among different data set.

The most general understanding of k-medoid bunch is that the "Partitioning around Medoids" (PAM) algorithmic rule and it is listed below:
1. Initialization step: Initially select a k information point out of n information points as a center point.
2. Allocation step: Group each data to the nearest center value (medoid).
3. Update step: For each center value (medoid) m and each information purpose of related to m swap m and o and calculate the full worth of the configuration. Medoid o with the lowest worth of the configuration is selected.
4. Repeat step 2 and step 3 until there is no alteration among the allocation.

## III. METHODOLOGY

The outline overview of the system is called the architecture of the system. It is the model which illustrates the characteristics of the modules, sub modules associated with the main system. It explains the flow of the overall system working, the levels in which the data flow, the modules and the processes which should be executed in what order. It should be able to categorize both visible and non-visible activities of the system. Figure 1, shows the proposed system architecture. The system architecture diagram of the proposed method mainly contains 3 modules. The three different modules are namely Data Collection Approach, Data preprocessing, cluster formation using data mining techniques and result Analysis. Using the health care data is more advantageous it helps the health care professionals and researchers to gain knowledge about the diseases and also the accuracy of the diseases increasing year by year.

**Data collection:** Data collection is the approach in which the health care data is collected from the specified hospitals or by some websites. The information of websites data is unstructured and fuzzy and hence it is subjected to the pre-processing techniques.

**Pre-processing techniques:** It includes clean and reducing the data. The data is cleaned by removing any noise present in the data. The unfilled value are filled by taking the mean of all the values and if more values are missing then the entire tuple is deleted, this is how the data is reduced. The xsl sheet data is converted into arff style data. These are the pre-processing steps we undertake for the data set sample.

**Clustering:** Clustering can be viewed as the most imperative unsupervised learning issue; in this way, as each

**Special Issue - 2017**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCETAIT - 2017 Conference Proceedings**

other issue of this type, it manages to find a structure in an collection of unlabeled information. The other meaning of bunching could be "the procedure of sorting out items into gatherings whose individuals are comparative somehow". A group is in this manner a gathering of articles which are

"comparative" amongst them and are "unique" to the items having a place with different bunches. In clustering first we cluster the data by using k-means clustering algorithm, then again we cluster data using k-medoids (PAM) algorithm.
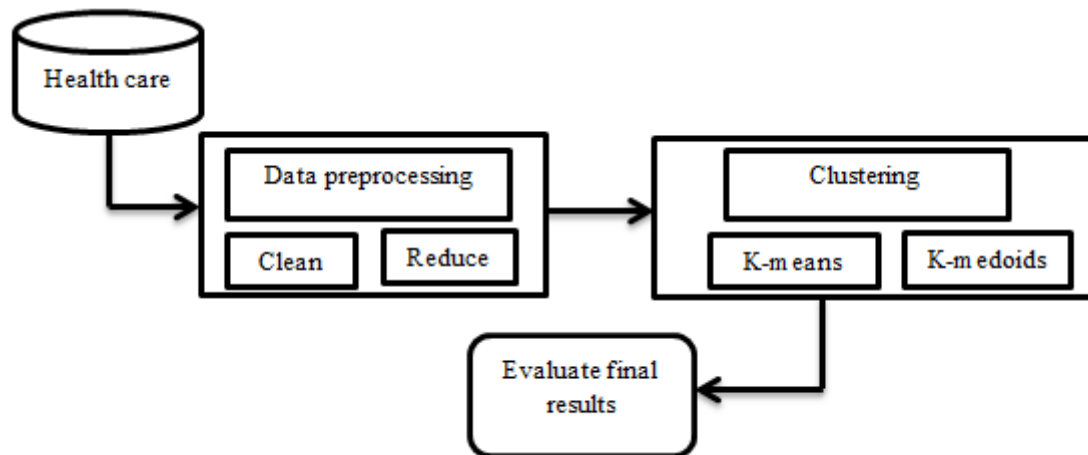


Figure 1: Architecture diagram for clustering TB data sample

## IV. PERFORMANCE COMPARISION OF STRATEGIES FOR CHOOSING INITIAL MEDOID

Performance of unvarying cluster algorithmic rule depends up on selection of 'cluster centers' in every step. An experiments square measure conducted on T.B. information set sample. T.B. (TB) dataset is shown in below Figure a pair of. A world wide causes of death due to Tuberculosis (TB) is found one in every of the ten people. T.B. (TB) is caused by microorganisms (Mycobacterium tuberculosis) that almost all typically have an effect on the lungs. T.B. can be cured and prevented. TB is unfolded from one person to another person through air. once individuals with respiratory organ sneeze or spit, TB cough, they spread the germs of TB into an air. If someone has inhale solely many of these germs are to be infected.

Regarding tierce of a world population had latent TB problem, which suggests individuals are infected by TB microorganism however aren't (yet) sick with the sickness and can't transmit the sickness. Individuals infected with TB microorganism have a tenth period of time danger of having sick with TB. On the other hand, persons with high immune power like individuals having an HIV, deficiency disease or polygenic disorder, or peoples using tobacco, have a higher danger for falling sick. Once someone develops an active TB sickness, the symptoms of these disease are (weight loss, cough, night sweats, or fever) could also be delicate for several days. This will cause due of days in seeking a care,

and leads to transform of an microorganism to another. Individuals with an active TB problem will be infecting to 10–15 peoples through shut in contacts over a courses of a year. While not correct handling of disease, leads to forty five percent of HIV-negative individuals with TB on this the average and a nearly all the people having HIV-positive with a TB problem will die.

In 2015, 10.4 million people fell ill with TB and 1.8 million died from the disease (including 0.4 million among people with HIV). Over 95% of TB deaths occur in low- and middle-income countries. Six countries account for 60% of the total, with India leading the count, followed by Indonesia, China, Nigeria, Pakistan and South Africa.

**Attributes**
Country –country name
ISO2- International standard organization (ISO) code alpha 2(country code)
ISO3-country code alpha 3
ISO_numeric-Numeric country code
Year
e_most_exc_tbhiv_100k-Mathematical model for estimating TB incidence
e_inc_tbhiv_100k_hi-HIV positive per 100000 population high bound
e_mort_exc_tbhiv_num- Number of deaths from TB
e_pop_num - Estimated total population number
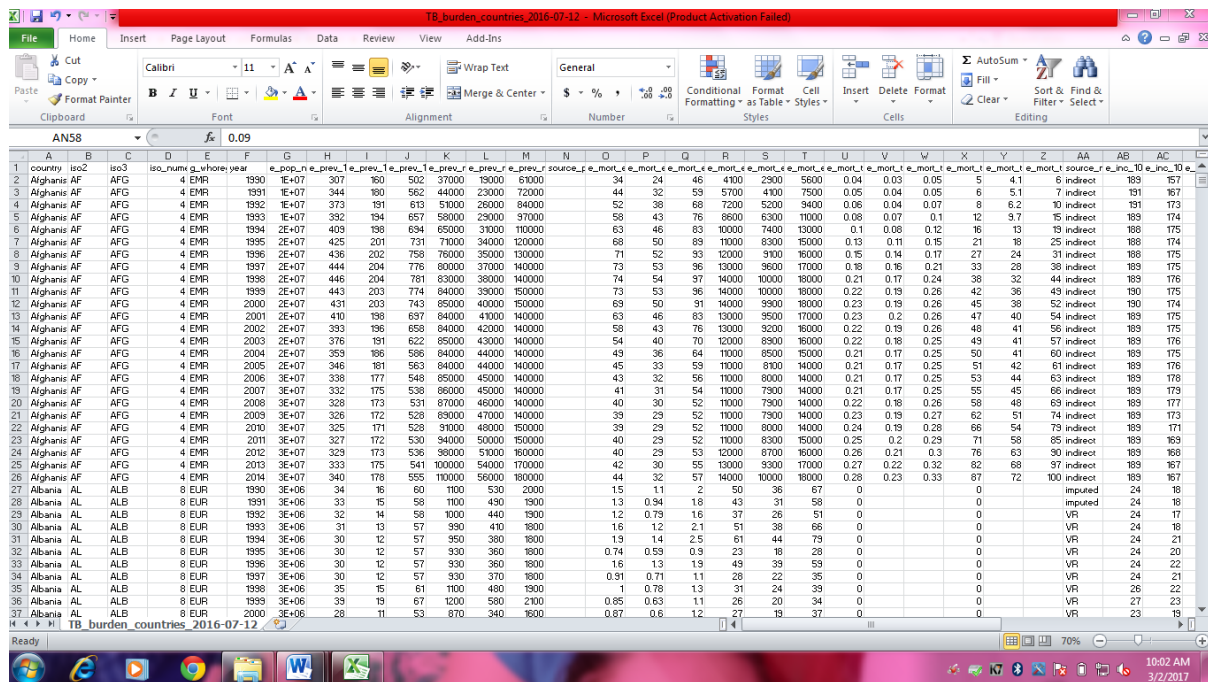new_ep_f014- female extra pulmonary aged 0-14 years

**Special Issue - 2017**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCETAIT - 2017 Conference Proceedings**

Figure 2: Tuberculosis dataset sample

## V. COMPARISON BETWEEN K-MEAN CLUSTERS AND K-MEDIOD ALGORITHMIC RULE WITH A VARIETY OF CLUSTER NUMBER AND ITS EXECUTION TIME

Table 1: For different number of clusters and time taken for execution using K-Mean method and K-Medoid Algorithms

| Number of clusters | Execution Time (in milliseconds) K-means Algorithm | Execution Time (in milliseconds) K-medoids Algorithm |
|---|---|---|
| 2 | 23123 | 20456 |
| 3 | 42245 | 41275 |
| 4 | 68632 | 69452 |
| 5 | 76893 | 75649 |

The Figure 3 shows a comparison among K-means and a K-medoids algorithm. Since graph shows the variety of a cluster is a smaller amount of time takes to execute than a K-mean method. If a quantity of a clusters is over its once it is extra accurate that K-medoid takes a lesser time to compute 15345 records than a K-mean cluster method. At a foremost variety of a cluster enhances than time taken for execution by a K-medoid method is a smaller amount than the K- mean.
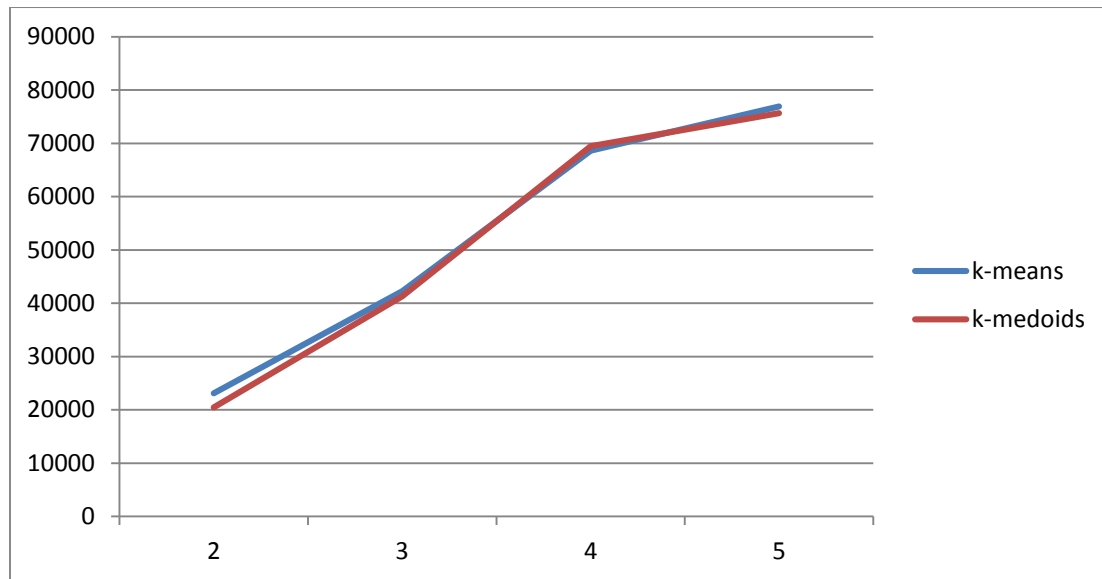
**Special Issue - 2017**

**International Journal of Engineering Research & Technology (IJERT)**
**ISSN: 2278-0181**
**NCETAIT - 2017 Conference Proceedings**

Figure 3: Comparison of a K-means and K-medoid algorithm

## VI. CONCLUSION

In this paper comparison of both K-means and K-medoids algorithmic rules are discussed. K-medoid technique has higher performance than K-means agglomeration and this takes a lesser execution time than K-mean method. Conjointly numerous strategies for choosing original clusters are conferred and compared. The K-Medoid agglomeration may be a higher for agglomeration performance and its execution time will be a smaller amount. As a result, even if the strategy of choosing first clusters delineate within the planned method sufficient for the use of once allowing for each performance and also for a computation time.

## REFERENCES

[1]  S. A. Raut, et.al, "Bioinformatics: Trends of Gene Expression and Analysis," proceedings of International Conference On Bioinformatics and Biomedical Technology, 16-18 April 2010,Chengdu, China.

[2]  S. A. Raut, et.al. "Gene Expression Analysis- A Review of large datasets," vol.4, Issue 1, Journal of Computer Science and Engineering, November 2010.

[3]  Xiong, H., et.al. "K-Means clustering versus validation measures: A Data distribution perspective", IEEE Trans. Syst., Man,C ybernet. PartB,39:318331.http://www.ncbi.nlm.nih.gov/pubmed/1909553, 2009.

[4]  S. Ray, et.al, "Purpose of number of clusters in k-means clustering and application in a color image segmentation," In Proceedings of the 4th International Meeting on Advances in Pattern Recognition and Digital Techniques, pp.137-143. 2014.

[5]  G. Sheikholeslami et.al, "Wave-Cluster: A Multi-Resolution Clustering Approach for Very Large Spatial Databases," Proc. 24th Int. Conf. on Very Large Data Bases. New York, pp. 428-439,2008.

[6]  R. Sibson, "SLINK: An optimally efficient algorithm for the single link cluster method," The Comp. Journal, 16(1), 2015, pp. 30-34.

[7]  T. Zhang, et.al, "BIRCH: An Efficient Data of Clustering Method for Very Large Databases," Proc. ACM SIGMODInt. Conf. on

[8]  Zhang Y, et.al "An efficient Clustering algorithm", In Proceedings of Second International Conference on Machine Learning and Cyber netics, November 2013.