

# Disease Prediction Model Using Clustering Classification Techniques for Diabetes Mode in Data Mining: Review Paper

Peddineni Kalpana<sup>1</sup>, Dr. B. V. V. Siva Prasad<sup>2</sup>

<sup>1</sup>Research Scholar, Career Point University

<sup>2</sup>Research Supervisor, Career Point University

**Abstract** - People are the most composite living beings on this globe. It is difficult to imagine how billions of minuscule parts, everyone with its own personality, cooperate in an arranged way for the benefit of the all-out being. A framework is an association of different organs masterminded together with the goal that they can complete complex capacities for the body. Body capacities are the physiological or mental elements of body frameworks. Endurance is the body's most significant desire, and it relies upon the body's looking after homeostasis. Homeostasis is a circumstance of relative consistency of body's inner climate, and it relies upon the body's doing numerous activities in an organized way constantly. Its significant activities or capacities are reacting to changes in the body's current circumstance, trading materials between the climate and cells, utilizing food sources, and incorporating the entirety of the body's assorted activities. In the event that there is any change in the homeostasis illnesses set in [1]. A sickness is an unusual circumstance influencing any piece of the body. Infections are generally grouped into transferable and non-transmittable illness.

**General Terms**—Medical data mining, clustering, rule-based classification using M-tree, K-means, Weka.

**Index Terms** – K-means clustering, Categorical data, rule-based classification, M-tree, Pima Indian Diabetics.

## I. INTRODUCTION

The information mining functionalities are utilized to determine the sort of examples to be found in the information mining task. The information mining functionalities fundamentally incorporate affiliation rule mining, arrangement, forecast and grouping. Affiliation investigation is utilized for finding fascinating relations between factors with regards to

huge information bases, which in given as rules to client. Order predicts the class marks. Expectation is utilized to get to the worth of a property that a given example is probably going to have. Bunching is the way toward gathering the information into classes or groups so that articles inside a group have high closeness in contrast with each other yet are unlike items in different bunches. Grouping is managed learning calculations in stands out from bunching, which are unaided learning calculation [1]. Arrangement is a regulated model, which maps or orders an information thing into one of a few predefined classes. Information arrangement is a two-venture measure. In the initial step, a model is constructed portraying a foreordained arrangement of information classes or ideas. The most widely recognized grouping information mining procedures are Case-Based Reasoning, M tree, Back spread neural organization, Radial premise neural organization, Bayesian arrangement, rough set Approach, Fuzzy Set Approaches, and K-closest neighbor classifiers. In this paper a fell K-implies bunching and M tree has been utilized to sort diabetic's patients. Writing review of grouping of diabetic informational collection is advised in segment 2. For culmination M - tree and K-mean bunching have been momentarily clarified in segment 3 and 4. Preprocessing of diabetic informational index and working of fell K-implies bunching and M tree classifier is clarified in area 5, trailed by results and end in segment 6 and 7 individually.

## II. MEDICAL DATA MINING

Information Mining for Healthcare Management (DMHM) is a promising field where scientists from both scholarly community and industry have recognized the capability of its contact on improved medical care by designing examples and patterns in a lot of complex information produced by medical services exchanges. Information mining likewise assists with finding appealing business bits of knowledge to help settle on business choices that can influence cost productivity but keep a top notch of care. Medical services the board has gotten thousand arrangements of consideration in current occasions and use of information mining methods to this field is ahead expanding ubiquity [15].

The information rich nature of the medical services area has made it an ideal air, where information on information mining ought to likewise must be expanded further for the expanding need. However, the theoretical idea of implied medical services information has brought about the under-use of a particularly major segment of the overall medical services conveyance system [3]. There are numerous calculations for these issues, however they are not severe and exact. Patient Reported Outcomes (PROs) in clinical examinations have continuously improved in recurrence, for their significance in assessing treatments and maturing treatment plans. Because of the phenomenal development pace of Health Care information, which is being gathered, stacked and put away through the World Wide Web and got to electronically in practically all fields of human undertaking, there is an imperative requirement for modern apparatuses and strategies that can switch amazingly enormous numerous information.

### III. RELATED WORK ON DIABETIC DATA SET

#### CLASSIFICATION TECHNIQUES: MACHINE LEARNING TECHNIQUES

Diabetics: -Disease description

Diabetes mellitus diabetes is an infection wherein the body cannot deliver or incapable to appropriately utilize and store glucose (a type of sugar). Glucose upholds in the circulation system causing one's blood glucose or "sugar" to ascend excessively high. There are two significant sorts of diabetes. In type 1 diabetes, the body totally quits creating any insulin, a

chemical that empowers the body to utilize glucose found in food sources for fuel.

Individuals with type 1 diabetes should take day by day insulin infusions to endure. This type of diabetes typically creates in kids or youthful grown-ups yet can happen at whatever stage in life.

Diabetes: Type 2 (likewise called grown-up beginning or non-insulin-subordinate) diabetes results when the body does not deliver sufficient insulin as well as cannot utilize insulin appropriately (insulin obstruction). This type of diabetes normally happens in individuals who are more than 40, overweight, and have a family background of diabetes, albeit today it is progressively happening in more youthful individuals, especially teenagers [2], [3].

World Health Organization (WHO) report had shown a stamped expansion in the quantity of diabetics and this pattern is required to fill in the following years and years. In the International Diabetes Federation Conference 2003 held in Paris, India was named, as "Diabetes Capital of the World," as of around 190 million diabetics around the world, in excess of 33 million are Indians. The overall figure is relied upon to ascend to 330 million, 52 million of them Indians by 2025, generally because of populace development, maturing, urbanization, unfortunate dietary patterns, and an inactive way of life. Ineffectively oversaw diabetes can prompt a large group of long-haul complexities among these are coronary episodes, strokes, visual impairment, kidney disappointment, vein illness.

### IV. LITERATURE REVIEW OF CLASSIFICATION OF DIABETIC DATASET

A great deal of examination work has been done on different clinical informational collections including Pima Indian diabetes dataset. Grouping exactness accomplished for Pima Indian diabetes dataset utilizing 22 distinct classifiers is given in [4] and utilizing 43 unique classifiers is given in [5]. The exhibition of proposed fell model utilizing k-means and M tree is contrasted and [4] and [5]. The consequences of [5] and [4] are appeared in Table 1 and Table 2 separately. The precision of the majority of these classifiers is in the scope of 66.6% to 77.7%. Mixture K-means and M tree [6]

accomplished the grouping precision of 92.38% utilizing 10 overlay cross approvals for constant information. Further fell learning framework dependent on Generalized Discriminate investigation (GDA) and Least Square Support Vector Machine (LS\_SVM), showed precision of 82.05% for determination of Pima dataset [7]. Further creators have accomplished arrangement precision of 72.88 % utilizing ANN, 78.21% utilizing DT\_ANN where M tree C4.5 is utilized to distinguish pertinent highlights and given as contribution to ANN [8], 79.50% utilizing Cascaded GA\_CFS\_ANN, important component recognized by Genetic calculation with Correlation based element determination is given as contribution to ANN [9], 77.71% utilizing GA improved ANN, 84.10% utilizing GA enhanced ANN with significant highlights distinguished by M tree and 84.71% with GA advanced ANN with applicable highlights distinguished by GA\_CFS[10]. Creators have accomplished a precision of 96.68% for diabetic dataset utilizing fell k-means and K-closest neighbor [11].

#### M-TREE: PROPOSED ALGORITHM

M-tree represents a supervised approach to classification. A M tree is a simple tree structure where non-terminal nodes represent tests on one or more attributes and terminal nodes reflect M outcomes. The basic M-tree induction algorithm was enhanced The WEKA classifier package has its own version of known as J4.8. Information gain and gain ratio measures are used by as splitting criterion, respectively. The summary of M tree algorithm is given.

1. Choose an attribute that best differentiates the output attribute values.
2. Create a separate tree branch for each value of the chosen attribute.
3. Divide the instances into subgroups so as to reflect the attribute values of the chosen node.
4. For each subgroup, terminate the attribute selection process if:
  - (a) The members of a subgroup have the same value for the output attribute, terminate the attribute selection process for the current path and label the branch on the current path with the specified value.
  - (b) The subgroup contains a single node or no further distinguishing attributes can be determined. As in (a),

label the branch with the output value seen by the majority of remaining instances.

© For each subgroup created in (iii) that has not been labeled as terminal, repeat the above process.

#### K-MEANS CLUSTERING:

K-implies [15] is one of the easiest solo learning calculations and follows dividing strategy for grouping. K-implies calculation takes the information boundary, k as number of groups and segments a dataset of n objects into k bunches, so the subsequent objects of one group are unlike that of other bunch and like objects of a similar bunch. In k-implies calculations starts with arbitrarily chose k items, addressing the k introductory group place or mean. Next each item is doled out to one the group dependent on the closeness of the article with bunch focus. To relegate the item to the nearest focus, a nearness measure specifically Euclidean distance is utilized that evaluates the idea of nearest. After every one of the articles are disseminated to k bunches, the new k group places are found by taking the mean of objects of k bunches separately. The interaction is rehashed till there is no adjustment of k bunch communities. K-implies calculation targets limiting a target work to be specific amount of squared mistake (SSE). SSE is characterized as Where E is amount of the square mistake of items with bunch implies for k group. p is the article have a place with a bunch  $C_i$  and  $m_i$  is the mean of group  $C_i$ . The time intricacy of K-implies is  $O(t*k*n)$  where t is the quantity of cycles, k is number of bunches and n is the absolute number of records in dataset. K-implies dividing calculation: (Input is k is the quantity of groups, D is input informational collection. Yield is k groups).

1. Randomly pick k items from D as the underlying group places.
2. Rehash
3. Dole out each item from D to one of k bunches to which the article is most comparative dependent on the mean worth of the articles in the group.
4. Update the group implies by taking the mean worth of the items for every one of k bunch.
5. Until no adjustment of group implies/min blunder E is reached.

K-MEANS AND M-TREE: COMPARISON ALGORITHM

Data preprocessing: -

The PIMA diabetic information base comprises of two classes in the informational collection (for example Tried positive, Tested Negative) each having 8

highlights: Number of times pregnant, Plasma glucose fixation a 2 hours in an oral glucose resistance test, Diastolic pulse (mm Hg), Triceps skin overlay thickness (mm), 2-Hour serum insulin (mu U/ml), Body mass file (weight in kg/ (stature in m)<sup>2</sup>), Diabetes family capacity and Age (years). The information is profited from UCI Machine Learning the information preparing procedures, when applied before mining, can generously improve the general nature of the examples mined as well as the time needed for the real mining. Information preprocessing is a critical advance in the information disclosure measure since quality Ms should be founded on quality information. An aggregate of 768 cases is accessible in PIDD. 5 patients had a glucose of 0, 11 patients had a weight record of 0, 28 others had a diastolic pulse of 0, 192 others had skin crease thickness readings of 0, and 140 others had serum insulin levels of 0. In the wake of erasing these cases there were 392 cases with no missing qualities (130 tried positive cases and 262 tried negative) [16].

V.WORKING OF PROPOSED METHOD

In the main phase of proposed model, basic K-implies grouping (with k = 2) of Weka apparatus, is applied to 392 diabetics patient's examples as gotten in area 5.1. The wrongly ordered examples are wiped out to get last 299 examples.

As a piece of preprocessing the persistent information is changed over to downright frame by rough width of the ideal stretches, in light of the assessment of clinical specialists as demonstrated in table 3. At long last, in the subsequent stage, the effectively characterized tests from first stage and the downright information is given as contribution to M tree C4.5 (weka J4.8). The information is parceled utilizing (a) 60-40 proportion dividing strategy (preparing test) and (b) 10-overlay cross approval. For fulfillment not many of the exhibition measurements have been

talked about. Genuine positive (TP) compares to the quantity of positive models accurately anticipated by the classifier. Bogus negative (FN) compares to the quantity of positive models wrongly anticipated as negative by the classifier. Bogus positive (FP) relates to the quantity of negative models wrongly anticipated as certain by the classifier. Genuine negative (TN) relates to the quantity of negative models effectively anticipated by the classifier. The genuine positive rate (TP rate) or affectability is the negligible portion of positive models anticipated accurately by the model.  $TP\ Rate = TP / (TP + FN)$ . The bogus positive rate (FP rate) or Specificity the small part of negative models anticipated as a positive class.  $FP\ Rate = FP / (TN + FP)$ . Accuracy is the negligible portion of records that really ends up being positive in the gathering the classifier has pronounced as a positive class.  $Accuracy = TP / (TP + FP)$ . Review is the small amount of positive models effectively anticipated by the classifier.  $Review = TP / (TP + FN)$ . F-measure is utilized to inspect the tradeoff among review and exactness.  $Measure = 2 * TP / (2 * TP + FP + FN)$ .

VI.RESULT ANALYSIS AND COMPARISON

Cluster instances	%	accuracy
M-TREE	69.0	8.990
K-MEANS	255	33.2031

VII.RESULT ANALYSIS

The result analysis is based on two parameters on the basis of data set records considered. The comparison analysis on k-means and M-tree algorithm performed as given follows:

1. Accuracy of k-means and M-tree proposed algorithm on disease data set.
2. Goodness fit ratio of prosed algorithm and K-means algorithm result determined.

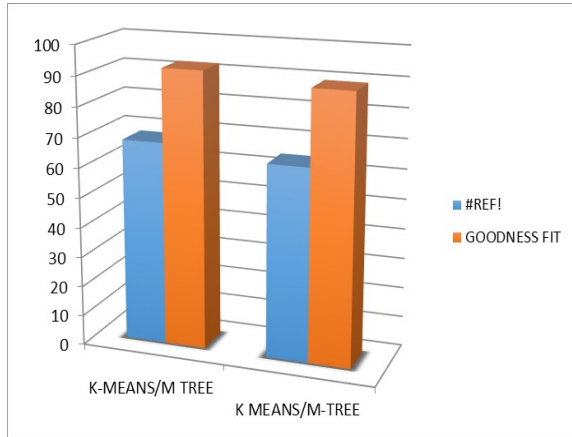


Fig 1.0 Result Comparison Between K-means and Mtree Algorithm.

### VIII.CONCLUSION WORK

The exhibition of characterization calculation relies upon the nature of information. The K-implies grouping is utilized to distinguish and dispose of erroneously ordered cases. Further the nonstop information is changed over to clear cut information by counseling clinical master's recommendation. The accurately arranged case by k-implies is utilized as contribution to M tree after transformation of constant information to straight out information. The proposed fell shows improved arrangement of 92.33% for PIMA diabetic dataset utilizing 60-40 % preparing testing parceling technique with preprocessed information. Further outcomes showed that the exhibition of fell model with absolute information created relatively less number of rules which are not difficult to decipher contrasted with rules produced by M tree with natural information. The arrangement correctness's got by the proposed fell K\_ implies grouping and M tree classifier is probably the best outcome contrasted and the aftereffects of M tree revealed in the writing. The M-tree calculation proposed model of diabetes is exceptionally precision of positive trial of diabetes patient.

### REFERENCE

[1] G.Kesavaraj, Dr.S.Sukumaran, "A Study on Classification Techniques in Data Mining", 4th ICCNT, IEEE 2013, Tiruchengode, India.

[2] Karthika Jayprakash, Nidhi Kargathra, Pranay Jagtap, Suraj Shridhar, Archana Gupta, "Comparison of Classification Techniques for Heart Health Analysis System", International Journal of Computer Sciences and Engineering, Vol. 4, No.2, pp.92-95, 2016.

[3] Divya Tomar and Sonali Agarwal, "A survey on Data Mining approaches for Healthcare", International Journal of Bioscience and Bio-Technology Vol.5, No.5, pp. 241-266, 2013.

[4] Tanvian and, rekha pal and Sanjay kumar dubey, "Data mining in healthcare informatics: Techniques and applications", 3rd International Conference on Computing for Sustainable Global Development, IEEE 2016.

[5] B. Sunil Srinivas, Dr. A. Govardhan, Dr. C. Sunil Kumar, "Data Mining Issues and Challenges in Healthcare Domain", International Journal of Engineering Research & Technology, Vol. 3, No. 1, January 2014.

[6] Shanta kumar B. Patil, Y. S. Kumaraswamy, "Intelligent and Effective Heart Attack Prediction System Using Data Mining and Artificial Neural Network", European Journal of Scientific Research, Vol. 31, No.4, pp. 642-656, 2009.

[7] Polat K., Gunes S., Aslan A., "A cascade learning system for classification of diabetes disease: Generalized discriminate analysis and least square support vector machine", Expert systems with applications, Vol. 34, No. 1, pp. 214-221, 2008.

[8] Asha Gowda Karegowda, A.S. Manjunath, M.A. Jayaram, "Application of Genetic Algorithm Optimized Neural Network Connection Weights for Medical Diagnosis of Pima Indians Diabetes", International Journal on Soft Computing, Vol.2, No.2, pp. 15-23, 2011.

[9] VeenaVijayan V and Aswathy Ravikumar, "Study of Data Mining Algorithms for Prediction and Diagnosis of Diabetes Mellitus", International Journal of Computer Applications, Vol. 95, No. 17, pp. 12-16, June 2014.

[10] M.A.NisharaBanu, B Gomathy, "DISEASE PREDICTING SYSTEM USING DATA MINING TECHNIQUES", International Journal of Technical Research and Applications, Vol. 1, No. 5, pp. 41-45, Nov-Dec 2013.

- [11] JyotiSoni, Ujma Ansari, Dipesh Sharma, SunitaSoni, “Predictive data mining for medical diagnosis: an overview of heart disease prediction”, International Journal of Computer Science and Engineering, vol. 17, No. 8, pp. 43-48, March 2011.
- [12] SitiFarhanah, BtJaafar, DannawatyMohd Ali, “Diabetes mellitus forecast using artificial neural networks”, Asian conference of paramedical research proceedings, September 2005, Kuala Lumpur, Malaysia.
- [13] K.Vijaya Lakshmi, Prof.M.Padmavathamma, “Modeling an Expert System for Diagnosis of Gestational Diabetes Mellitus Based On Risk Factors”, IOSR Journal of Computer Engineering (IOSRJCE), Vol. 8, No. 3, pp. 29-32, Jan. – Feb. 2013.
- [14] Ms. A J. Chamatkar et al, “Importance of Data Mining with Different Types of Data Applications and Challenging Areas”, International Journal of Engineering Research and Applications, Vol. 4, No. 5, pp.38-41, May 2014.
- [15] Kritika Yadav, Mahesh Parmar, “Review Paper on Data Mining and its Techniques of Mahatma Gandhi National Rural Employment Guarantee Act”, International Journal of Computer Science and Engineering (IJCSE), Vol. 5, No. 4, pp. 68-73, April 2017.
- [16] Richa Sharma, Dr. Shailendra Narayan Singh, Dr. Sujata Khatri, “Medical Data Mining Using Different Classification and Clustering Techniques: A Critical Survey”, Second International Conference on Computational Intelligence & Communication Technology, IEEE 2016