

Data Warehouse and Mining Tools and Technology in Medical Science

Samir Suman¹, Dr. Prakash Kumar Pathak², Ms. Yamini Sharma³

¹M.Tech (CSE) Student from MD University, Rohtak, Haryana, India

²Associate Professor CSE Department, WCTM, Gurugram

³Assistant Professor CSE Department, WCTM, Gurugram

ABSTRACT

The amount of data in electronic and real world is constantly on the rise. Therefore, extracting useful knowledge from the total available data is very important and time consuming task. Data mining has various techniques for extracting valuable information or knowledge from data. These techniques are applicable for all data that are collected in all fields of science. Several research investigations are published about applications of data mining in various fields of sciences such as defense, banking, insurances, education, telecommunications, medicine and etc. This investigation attempts to provide a comprehensive survey about applications of data mining techniques in breast cancer diagnosis, treatment & prognosis till now. Further, the main challenges in these area is presented in this investigation. Since several research studies currently are going on in this issues, therefore, it is necessary to have a complete survey about all researches which are completed up to now, along with the results of those studies and important challenges which are currently exist in this area for helping young researchers and presenting to them the main problems that are still exist in this area.

In this study we have conclude the techniques , breast cancer and Chronic and kidney Disease.. The present study is focused on the usage of classification techniques in the field of medical science and bioinformatics. The goal of data mining application is to turn that data are facts, numbers, or text which can be processed by a computer into knowledge or information. The main purpose of data mining application in healthcare systems is to develop an automated tool for identifying and disseminating relevant healthcare information. Breast cancer is one of the leading cancers for women in developed countries including India. It is the second most common cause of cancer death in women. The high incidence of breast cancer in women has increased significantly in the last years. Chronic- Kidney-Disease prediction using weka data mining tool and its usage for classification in the field of medical bioinformatics. It firstly classifies dataset and then determines which algorithm performs better for diagnosis and prediction of Chronic- Kidney-Disease. Finally, the

existing data mining techniques with data mining algorithms and its application tools which are more valuable for healthcare services are discussed in detail

Keywords: Data mining, medical data, cancer diagnosis, cancer treatment, cancer prognosis, risk factor.

I INTRODUCTION

Nowadays in all fields of sciences including genetics, education, earth science, agriculture and medicine the amount of data is increasing dramatically. Analyzing this huge amount of data to extract the novel and usable information or knowledge is very complicated and time consuming task. Data mining techniques are useful for this matter. Moreover, the advancement of the healthcare database management systems creates a huge number of medical databases. Creating knowledge and management of large amounts of heterogeneous data has become a major field of research, namely data mining. Data Mining is a process of identifying novel, potentially useful, valid and ultimately understandable patterns in data [1]. Data mining techniques can be classified into both unsupervised and supervised learning techniques. Unsupervised learning technique is not guided by variable and does not create a hypothesis before analysis. In the present study, we have focused on the usage of classification techniques in the field of medical science and bioinformatics. Classification is the most commonly applied data mining technique, and employs a set of pre-classified examples to develop a model that can classify the population of records at large.

Generally, in the medical world, there are two phases for making the decisions. These two phases are.

(1) *Differential Diagnosis (DD)*: in this phase, all information of patients including their medical history, symptoms of disease, results of various testing such as blood testing and etc. are perceived by doctors as the input data. These data are processed by doctors based on their medical knowledge for disease diagnosis. Sometimes several diseases have some similar symptoms, therefore, medical doctors must be assign arbitrary weights to each one of inputs and make patterns, match these patterns with the patterns of various diseases and finally select the closest match and diagnosis the exact disease.

(2) *Final or Provisional Diagnosis (FD)*: in this phase, the preliminary recommendations and treatments would be start according to the identified disease. In this step, a physician with medical knowledge and his/her logic, continues checkups and records the results of continually perceives or tests, and decides the final treatments and prognosis.

Data mining has various techniques (such as: Classification, Clustering, Regression, Association Rules and etc.) and algorithms (such as: Decision Trees, Genetic Algorithm, Nearest Neighbor method etc.) for analyzing the huge amount of raw or multi-dimensional data. In the other words, data mining has capabilities for intelligent data analysis to extract hidden knowledge from large databases of medical or clinical data that are collected from medical centers or hospitals. These knowledge provide useful

information to improve decision support, prevention, diagnosis and treatment in medical world. Further, data mining has ability to identify association rules or establish relationships between various features such as: patient's personal data, disease symptoms and etc.

This investigation attempts to represent the results of several research works which are published related to data mining applications in prediction, diagnosis or treatment of breast cancers.

This paper is organized in four sections. Section 2nd includes some basic concepts related to this paper. Section 3rd presents data mining applications or usages in early diagnosis, treatment and prognosis of various cancers. Section 4th concludes this paper and presents our future works.

II METHODOLOGY

Data Mining Process

In the KDD process, the data mining methods are for extracting patterns from data. The patterns that can be discovered depend upon the data mining tasks applied. Generally, there are two types of data mining tasks: *descriptive data mining tasks* that describe the general properties of the existing data, and *predictive data mining tasks* that attempt to do predictions based on available data. Data mining can be done on data which are in quantitative, textual, or multimedia forms.

Data mining applications can use different kind of parameters to examine the data. They include association (patterns where one event is connected to another event), sequence or path analysis (patterns where one event leads to another event), classification (identification of new patterns with predefined targets) and clustering (grouping of identical or similar objects). Data mining involves some of the following key steps[3]-

- (1) *Problem definition*: The first step is to identify goals. Based on the defined goal, the correct series of tools can be applied to the data to build the corresponding behavioural model.
- (2) *Data exploration*: If the quality of data is not suitable for an accurate model then recommendations on future data collection and storage strategies can be made at this. For analysis, all data needs to be consolidated so that it can be treated consistently.
- (3) *Data preparation*: The purpose of this step is to clean and transform the data so that missing and invalid values are treated and all known valid values are made consistent for more robust analysis.

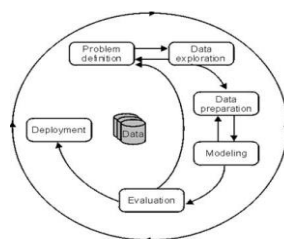


Fig.2. Data Mining Process Representation

(4) *Modelling*: Based on the data and the desired outcomes, a data mining algorithm or combination of algorithms is selected for analysis. These algorithms include classical techniques such as statistics, neighbourhoods and clustering but also next generation techniques such as decision trees, networks and rule based algorithms. The specific algorithm is selected based on the particular objective to be achieved and the quality of the data to be analysed.

(5) *Evaluation and Deployment*: Based on the results of the data mining algorithms, an analysis is conducted to determine key conclusions from the analysis and create a series of recommendations for consideration.

III DATA MINING CLASSIFICATION METHODS

The data mining consists of various methods. Different methods serve different purposes, each method offering its own advantages and disadvantages. However, most data mining methods commonly used for this review are of classification category as the applied prediction techniques assign patients to either a "benign" group that is non- cancerous or a "malignant" group that is cancerous and generate rules for the same. Hence, the breast cancer diagnostic problems are basically in the scope of the widely discussed classification problems.

In data mining, classification is one of the most important task. It maps the data in to predefined targets. It is a supervised learning as targets are predefined. The aim of the classification is to build a classifier based on some cases with some attributes to describe the objects or one attribute to describe the group of the objects. Then, the classifier is used to predict the group attributes of new cases from the domain based on the values of other attributes. The commonly used methods for data mining classification tasks can be classified into the following groups[4].

3.1. Decision Trees (DT's)

A decision tree is a tree where each non-terminal node represents a test or decision on the considered data item. Choice of a certain branch depends upon the outcome of the test. To classify a particular data item, we start at the root node and follow the assertions down until we reach a terminal node (or leaf). A decision is made when a terminal node is approached. Decision trees can also be interpreted as a special form of a rule set, characterized by their hierarchical organization of rules.

3.2. Support Vector Machine(SVM)

Support vector machine (SVM) is an algorithm that attempts to find a linear separator (hyper -plane) between the data points of two classes in multidimensional space. SVMs are well suited to dealing with interactions among features and redundant features.

3.3. Genetic Algorithms (GAs) / Evolutionary Programming (EP)

Genetic algorithms and evolutionary programming are algorithmic optimization strategies that are inspired by the principles observed in natural evolution. Of a collection of potential problem solutions that compete with each other, the best solutions are selected and combined with each other. In doing so, one expects that the overall goodness of the solution set will become better and better, similar to the process of evolution of a population of organisms. Genetic algorithms and evolutionary programming are used in data mining to formulate hypotheses about dependencies between variables, in the form of association rules or some other internal formalism.

3.4. Fuzzy Sets

Fuzzy sets form a key methodology for representing and processing uncertainty. Uncertainty arises in many forms in today's databases: imprecision, non-specificity, inconsistency, vagueness, etc. Fuzzy sets exploit uncertainty in an attempt to make system complexity manageable. As such, fuzzy sets constitute a powerful approach to deal not only with incomplete, noisy or imprecise data, but may also be helpful in developing uncertain models of the data that provide smarter and smoother performance than traditional systems.

3.5. Neural Networks

Neural networks (NN) are those systems modeled based on the human brain working. As the human brain consists of millions of neurons that are interconnected by synapses, a neural network is a set of connected input/output units in which each connection has a weight associated with it. The network learns in the learning phase by adjusting the weights so as to be able to predict the correct class label of the input.

3.6. Rough Sets

A rough set is determined by a lower and upper bound of a set. Every member of the lower bound is a certain member of the set. Every non-member of the upper bound is a certain non-member of the set. The upper bound of a rough set is the union between the lower bound and the so-called boundary region. A member of the boundary region is possibly (but not certainly) a member of the set. Therefore, rough sets may be viewed as with a three-valued membership function (yes, no, perhaps). Rough sets are a mathematical concept dealing with uncertainty in data. They are usually combined with other methods such as rule induction, classification,

In order to carry out experiments and implementations WEKA is used as the data mining tool for the users to classify the accuracy on the basis of datasets by applying different algorithmic approaches in the field of bioinformatics. In this work we have used the data mining techniques to predict the survivability of Chronic-Kidney disease through classification of different algorithms accuracy.

Explorer: The explorer interface has several panels like pre-process, classify, cluster, associate, select attribute and visualize. But in this interface our main focus is on the Classification Panel.



Experimenter: This interface provides facility for systematic comparison of different algorithms on basis of given datasets. Each algorithm runs 10 times and then the accuracy gets reported.

PRELIMINARY Classification is a supervised learning technique. It maps the data into predefined groups. It is used to develop a model that can classify the population of records at large level. Classification algorithm requires classes to be defined based on the data attribute value. It describes these classes according to the characteristics of the data that is already known to belong to the classes. The classifier training algorithm uses these pre-defined examples to determine the set of parameters required for proper discrimination. In Classification, training examples are used to learn a model that can classify the data samples into known classes. The Classification process involves following steps:

- Create training data set.
- Identify class attribute and classes.
- Identify useful attributes for classification (Relevance analysis).
- Learn a model using training examples in Training set.
- Use the model to classify the unknown data samples.

IV CLASSIFIERS

Used In this work six classification algorithms have been used for classification task to study their classification accuracy and performance over the Chronic-Kidney-Disease data set. The classifiers in Weka have been categorized into different groups such as Bayes, Functions, Lazy, Rules, Tree based classifiers etc. A good mix of algorithms has been chosen from these groups which are used in distributed data mining. They include Naive Bayes (from Bayes), Multilayer Perceptron, SVM, J48, Conjunctive rule and Decision Table. The following sections explain a brief about each of these algorithms.

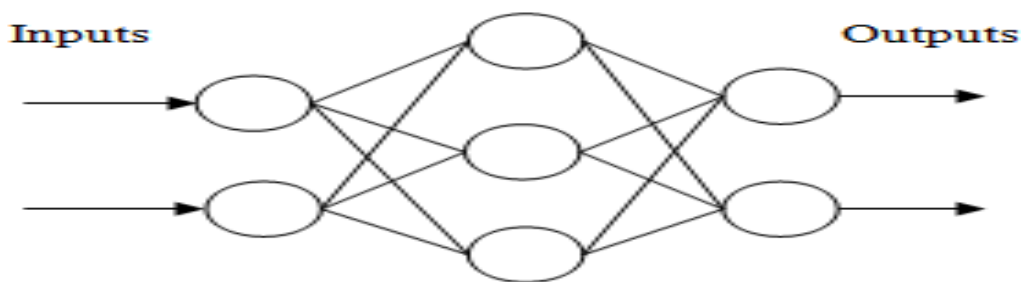
4.1 Naïve Bayes Classifier

It is one of the fastest statistical classifier algorithm works on probability of all attribute contained in data sample individually and then classifies them accurately. It is used to predict class membership probabilities i.e. probability about the tuple that belongs to the particular class or not. Bayesian classification is based on Bayes theorem. Abstractly, naive Bayes is a conditional probability model: given a problem instance to be classified, represented by a vector $X = (x_1, x_2, \dots, x_n)$ representing some n features (independent variables), it assigns to this instance probabilities $p(C_k | x_1, \dots, x_n)$ for each of k possible outcomes or classes. The problem with the above formulation is that if the number of features n is large or if a feature can take on a large number of values, then basing such a model on probability tables is infeasible. We therefore reformulate the model to make it more tractable. Using Bayes' theorem, the conditional probability can be decomposed as $p(C_k | X) = \prod (X|C_k)(X)$ In other

words, using Bayesian probability terminology, the above equation can be written as $\text{Posterior} = \text{prior} \times \text{likelihoodevidence}$

4.2 Multilayer Perceptron

It is the most popular network architecture in today's world. Each unit performs a biased weighted sum of their inputs and pass this activation level through a transfer function to produce their output. The units are arranged in a layered feed forward topology. The network has a simple input-output model, with the weights and thresholds. Such networks can model functions of almost arbitrary complexity, with the number of layers, and the number of units in each layer, determining the function complexity. The important issues in Multilayer Perceptron are the design specification of the number of hidden layers and the number of units in these layers. Multilayer Perceptron is a nonlinear classifier based on the Perceptron. A Multilayer Perceptron (MLP) is a back propagation neural network with one or more layers between input and output layer. The following diagram illustrates a perceptron network with three layers.



4.3 Support Vector Machine (SVM)

Support Vector Machine (SVM) is based on the concept of decision planes that define decision boundaries. A decision plane is one that separates between a set of objects having different class memberships. The standard SVM takes a set of input data and predicts, for each given input, which of two possible classes comprises the input, making the SVM a non-probabilistic binary linear classifier.

4.4 J48

J48 classifier is a simple C4.5 decision tree for classification. It is supervised method of classification. It creates a small binary tree. It is univariate decision tree. It is an extension of ID3 algorithm. In this classifier Divide and Conquer approach is used to classify the data. It divides the data into range based on the attribute value for that value that are found in training sample. As this approach is range based and univariate [11], it does not perform better than multivariate approach. As this is decision tree approach it is very much useful in predicting the values. J48 accuracy of correctly classified instance is much more than that of the other algorithms which are univariate in nature [10].

4.5 Conjunctive Rule

It is a decision-making rule in which the intending buyer assigns least values for a number of factors and discards any result which does not meet the bare minimum value on all of the factors i.e. a superior performance on one factor cannot recompenset for deficit on another. Conjunctive rule uses the AND logical relation to correlate stimulus attributes. Conjunctive rule is a simple well interpretable 2-class classifier.

V DECISION TABLE

A decision table is a predictive modeling tool that performs classification. It incorporates an inducer (an algorithm for generating decision table models), and a visualizer. Unlike the evidence model, the Decision Table model does not assume that the attributes are independent. A decision table is a hierarchical breakdown of the data, with two attributes at each level of the hierarchy. The Decision Table inducer identifies the most important attributes (columns) for classifying the data, and the accompanying visualizer displays the resulting model graphically. It summarizes the dataset with a decision table which contains the same number of attributes as the original dataset. Decision Table employs the wrapper method to find a good subset of attributes for inclusion in the table. By eliminating attributes that contribute little or nothing to a model of the dataset, the algorithm reduces the likelihood of over-fitting and creates a smaller and condensed decision table.

VI CHARACTERISTICS REQUIRED FOR CLASSIFICATION

Algorithm In this work, we have focused on the following three measures namely correctly classified instances, incorrectly classified instances, and accuracy. (i) Correctly classified instance: These are the instances which are correctly classified by any classification algorithm. Percentage of correctly classified instances is called as accuracy. (ii) Incorrectly classified instances: These instances are not correctly classified by the algorithm. Sometimes it is observed that the data which is incorrectly classified may contain inconsistent data, noisy data or data out of scope. (iii) Accuracy: Accuracy is how a measured value is closed to the true value. The general formula is given below: $Accuracy = \frac{Tp+Tn}{P+N}$ (1) where, Tp indicates True positive, Tn indicates True negative, P indicates total positive, N indicates total negative. And $P = Tp + Fp$, $N = Fp + Tn$. In classification system, the algorithm with highest accuracy will be selected for the prediction. Accuracy of the algorithm varies according to the dataset used. So before using the algorithms for prediction system, we must check the accuracy of the algorithm. So it will reduce the cost of doing trial and error of using algorithms in the prediction system.

6.1 Performance Evaluation

10-fold cross validation technique is used to evaluate the performance of classification methods, Data set is randomly sub divided into ten equal sized partitions. Among the partitions nine of them are used as training set and the remaining one is used as a test set. Evaluation of performance is compared using Mean absolute error, Rootmean squared error, Receiver Operating Characteristic (ROC) Area and Kappa statistics. Large test sets gives a good assessment of the classifier's performance and small training sets which result in a poor classifier.

6.2 Kappa Statistics

Kappa Statistics measure degree of agreement between two sets of categorized data. Kappa result varies between 0 to 1 intervals. Higher the value of Kappa means stronger the agreement. Kappa is a normalized value of agreement for chance of agreement. $K = \frac{PA - (E)1 - P(E)}{1 - P(E)}$ Where P(A) = percentage of agreement P(E) = chance of agreement. If K =1 agreement is perfect between the classifier and ground truth. If K=0 indicates there is a chance of agreement.

6.3 Mean Absolute Error (MAE)

The mean absolute error (MAE) is a quantity used to measure predictions of the eventual outcomes. The mean absolute error is given by $MAE = \frac{1}{n} \sum_{i=1}^n |f_i - y_i|$ The mean absolute error is an average of the absolute errors $e_i = |f_i - y_i|$, where f_i = prediction, y_i = true value.

6.4 Root Mean Squared Error (RMSE)

Root mean squared error is the square root of the mean of the squares of the values. It squares the errors before they are averaged [18] and RMSE gives a relatively high weight to large errors. The RMSE E_i of an individual program i is evaluated by the equation:

$$E_i = \sqrt{\frac{1}{n} \sum_{j=1}^n \left(\frac{P(i,j) - T_j}{T_j} \right)^2}$$

where, P(i,j) = the value predicted by the individual program i = fitness case T_j = the target value for fitness case j .

6.5 Receiver Operating Characteristic (ROC) Area

ROC Area is defined as area under the ROC curve which is the probability of randomly chosen positive instance that is ranked above randomly chosen negative one. Receiver Operating Characteristic represents test performance guide for classifications accuracy of diagnostic test based on: excellent (0.90-1), good (0.80-0.90), fair (0.70-0.80), poor (0.60-0.70), fail (0.50-0.60).

VII DISCUSSION AND RESULT

We have conducted two experiments based on the dataset with all above discussed classification algorithms; first without using feature selection and second with using Genetic search feature selection. First the results of the classification algorithms based on parameters such as accuracy of classification, kappa statistics, MAE, RMSE, model building time, model testing time, and ROC are shown in the following Table 2, where model building time and model testing time are generated by WEKA Tool itself during classification.



Attributes selection

First of all, we have to find the correlated attributes for finding the hidden pattern for the problem stated. The WEKA data miner tool has supported many in built learning algorithms for correlated attributes. There are many filtered tools for this analysis but we have selected one among them by trial.[5]

Totally there are 520 records of data base which have been created in Excel 2007 and saved in the format of CSV (Comma Separated Value format) that converted to the WEKA accepted of ARFF by using command line premier of WEKA.

The records of data base consist of 15 attributes, fromwhich 10 attributes were selected based on attribute selection in explorer mode of WEKA 3.6.4

CLASSIFICATION OF ATTRIBUTES

S.NO.	Attributes	Data Type
01.	Name	Text
02.	Age	Numeric(Integer)
03.	Education	Text
04.	Sex	Character
05.	Fluoride Level	Numeric(Real)
06.	Profession	Text
07.	Praganancy status	Boolean
08.	Drinking water	Text
09.	Duration	Numeric(Integer/Real)
10.	Known status of fluoride	Boolean
11.	Neck Pain	Numeric(Binary)
12.	Joint Pain	Numeric(Binary)
13.	Body Pain	Numeric(Binary)
14.	Foot Neck Pain	Numeric(Binary)
15.	Disease Level	Text

We have chosen Symmetrical random filter tester for attribute selection in WEKA attribute selector. It listed 14 selected attributes, but from which we have taken only 8 attributes. The other attributes are omitted for the convenience of analysis of finding impaction among peoples in the district

S.NO.	Attributes	Data Type
01.	Age	Numeric(Integer)
02.	Education	Text
03.	Fluoride Level	Numeric(Real)
04.	Drinking water	Text
05.	Duration	Numeric(Integer/Real)
06.	Neck Pain	Numeric(Binary)
07.	Joint Pain	Numeric(Binary)
08.	Body Pain	Numeric(Binary)
09.	Foot Neck Pain	Numeric(Binary)
10.	Disease Level	Text

K-Means Method

The k-Means algorithm takes the input parameter, k, and partitions a set of n objects into k clusters so that the resulting intracluster similarity is high but the intercluster similarity is low. Cluster similarity is measured in regard to the mean value of the objects in a cluster, which can be viewed the cluster's centroid or center of gravity.

The k-Means algorithm proceeds as follows

First, it randomly selects k of the objects, each of which initially represents a cluster mean or center. For each of the remaining objects, an object is assigned to the cluster to which it is the most similar, based on the distance between the object and the cluster mean. It then computes the new mean for each cluster. This process iterated until the criterion function converges. Typically, the square-error criterion is used, defined as [2] [3] [4] $E = \sum_i |p - m_i|^2$

Where E is the sum of the square error for all objects in the data set; p is the point in space representing a given object; and m_i is the mean of cluster C_i . In other words, for each object in each cluster, the distance from the object to its cluster center is squared, and the distances are summed. This criterion tries to make the resulting k clusters as compact and as separate as possible

K-Means algorithm

Input;

= k: the number of clusters,

D: a data set containing n objects

Output: A set of k clusters.

Method:

arbitrarily choose k objects from from D as the initial cluster centers;

(2) (re)assign each object to the cluster to which the object is the most similar, based on the mean value of the objects in the cluster;

(3) Update the cluster means, i.e., calculate the mean value of the objects for each cluster;

(4) until no change;

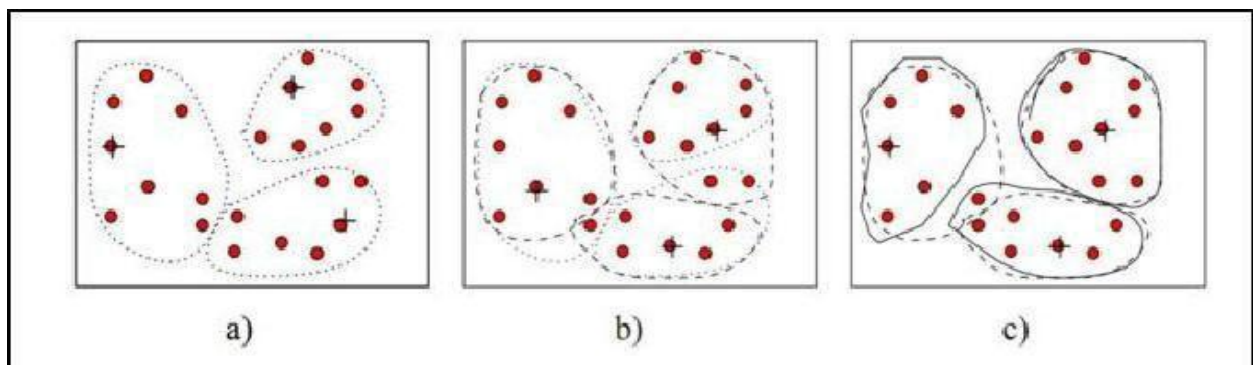
Suppose that there is a set of objects located in space as depicted in the rectangle. Let $k = 3$; i.e., the user would like the objects to be partitioned into three clusters.

According to the algorithm above we arbitrarily choose three objects as the three initial cluster centers, where cluster centers are marked by a "+". Each objects is distributed to a cluster based on the cluster center to which it is the nearest. Such a distribution forms encircled by dotted curves.

Next, the cluster centers are updated. That is the mean value of each cluster which is recalculated based on the current objects in the cluster. Using the new cluster centers, the objects are redistributed to the clusters based on which cluster center is the nearest. Such a redistribution forms new encircled by dashed curves. This process iterates. The process of iteratively reassigning objects to clusters to improve the partitioning is referred to as iterative relocation. Eventually, no redistribution of the objects in any cluster occurs, and so the process terminates. The resulting cluster is returned by the clustering process.

K-Means in WEKA

The learning algorithm k-Means in WEKA 3.6.4 accepts the training data base in the format of ARFF.



It accepts the nominal data and binary sets. So our attributes selected in nominal and binary formats naturally. So there is no need of preprocessing for further process. We have trained the training data by using the 10 Fold Cross Validated testing which used our trained data set as one third of the data for training and remaining for testing.

After training and testing this gives the following results.

Euclidean distance

K-means cluster analysis supports various data types such as Quantitative, binary, nominal or ordinal, but do not support categorical data. Cluster analysis is based on measuring similarity between objects by computing the distance between each pair. There are a number of methods for computing distance in a multidimensional environment. Distance is a well understood concept that has a number of simple properties. Distance is always positive Distance from point x to itself is always zero Distance from point x to point y cannot be greater than the sum of the distance from x to some other point z and distance from z to y. Distance from x to y is always the same as from y to x. It is possible to assign weights to all attributes indicating their importance. There are number of distance measures such as Euclidean distance, Manhattan distance and Chebychev distance. But in this analysis Weka tool used Euclidean distance. Euclidean distance of the difference vector is most commonly used to compute distances and has an intuitive appeal but the largest valued attribute may dominate the distance. It is therefore essential that the attributes are properly scaled. Let the distance between two points x and y be $D(x,y) = \sqrt{\sum(x_i - y_i)^2}$

Clustering of Disease Symptoms

The collected disease symptoms such as Neck pain, Joint pain, Body pain, Foot Neck as raw data, supplied to kmeans method is being carried out in weka using Euclidean distance method to measure cluster centroids. The result is obtained in iteration 12 after clustered. The centroid cluster points are measured based on the diseases symptoms and the water they are drinking. Based on the diseases symptoms in raw data the kmeans clustered two main clustering units. From the confusion matrix above we came to know that the district mainly impacted by skeletal osteoporosis.

VIII CONCLUSION

The main objective of this chapter is to predict chronic kidney disease. We have used six algorithms i.e. Naive Bayes, Multilayer Perceptron, SVM, J48, Conjunctive Rule and Decision Table for our experiments. These algorithms are implemented using WEKA data mining tool to analyze accuracy which is obtained after running these algorithms in the output window. These algorithms have been compared with classification accuracy to each other on the basis of correctly classified instances, time taken to build model, time taken to test the model, mean absolute error, Kappa statistics and ROC Area. In the experiments Multilayer perceptron algorithm gives better classification accuracy and prediction performance to predict chronic kidney disease (CKD) using relevant dataset available at UCI machine learning repository.

Data mining applied in health care domain, by which the people get beneficial for their lives. As the analog of this research we found out that the meaningful hidden pattern from the real data set collected the people impacted in Krishnagiri district is by drinking high fluoride content of potable water. By which we can easily know that the people do not get awareness among themselves about the fluoride impaction. If it continues in this

way, it may lead to some primary health hazards like Kidney failure, mental disability, Thyroid deficiency and Heart disease.

REFERENCES

- [1] Jain, M. Murty, and . Flynn, “Data clustering: A review,” A M Computing Surveys, vol. 31, no. 3, pp. 264–323, 1999.
- [2] Jiawei Han and MichelineKamber – Data mining concepts and Techniques. - Second Edition –Morgan Kaufmann Publishers
- [3] ArunK.Pujari –Datamining Techniques – University Press.
- [4] Introduction to Datamining with case studies - G.K.Gupta PHI. Fuzzy Models for Social Scientists - W.B.VasanthaKandasamy (e-book :<http://mit.iitm.ac.in>)
- [5] BerryMjLinoff G Data mining [10] Professionals statement calling for an Techniques: for Marketing, Sales and Customer support USA.Wiley,1997.
- [6] Weka3.6.4 data miner manual.
- [7] Water Quality for Better Health – TWAD Released Water book.
- [8] Data mining Learning models and Algorithms for medical applications – White paper - PlamenaAndreeva, Maya Dimibova, Petra Radeve
- [9] Elementary Fuzzy Matrix Theory and End to water Fluoridation – Conference Report (www.fluoridealert.org)
- [11] Analysis of Liver Disorder Using Data mining algorithms - Global Journal of computer science and Technology 1.10 issue 14 (ver1.0) November 2010 page 48.
- [12] The WEKA Data Mining Software: An Update, Peter Reutemann, Ian H. Witten, Pentaho Corporation, Department of Computer Science