



# Comparison of Different Classification Techniques Using WEKA for Diabetic Diagnosis

Arka Haldar<sup>1</sup>, G.Prudhvi Raj<sup>2</sup>, S.V.S.S Lakshmi<sup>3</sup>

P.G Student, Dept. of Computer Science and Systems Engineering, Andhra University, Visakhapatnam, India<sup>1</sup>

U.G Student, Dept. of Computer Science & Engineering, ANITS, Visakhapatnam, India<sup>2</sup>

Assistant Professor, Dept. of. Computer Science & Engineering, ANITS, Visakhapatnam, India<sup>3</sup>

**ABSTRACT:** Medical professionals need a reliable prediction methodology to diagnose factors influencing diabetes. There are large quantities of information about patients and their medical conditions. Generally, data mining (sometimes called data or knowledge discovery) is the process of analyzing data from different perspectives and summarizing it into useful information. Data mining software is one of a number of analytical tools for analyzing data. It allows users to analyze data from many different dimensions or angles, categorize it, and summarize the relationships identified. Weka is a data mining tools. It contains many machine learning algorithms. It provides the facility to classify our data through various algorithms. Classification is an important data mining technique with broad applications. It classifies data of various kinds. Classification is used in every field of our life. Classification is used to classify each item in a set of data into one of predefined set of classes or groups. In this paper we are studying the various Classification algorithms. The thesis main aims to show the comparison of different classification algorithms using Waikato Environment for Knowledge Analysis or in short, WEKA and find out which algorithm is most suitable for user working on haematological data of diabetic patients. To use proposed model, new Doctor or patients can predict the factors influencing diabetes. The best algorithm based on the diabetic hematological data is Naive Bayes with an accuracy of 76.3021% and the total time taken to build the model is at 0.06 seconds. Naive Bayes classifier has the lowest average error at 29.71% compared to others.

**KEYWORDS:** Data Mining, J48 Decision tree, Zero R, Naïve Bayes.

## I. INTRODUCTION

Data mining technique is a process of discovering pattern of data. The patterns discovered must be meaningful in that they lead to some advantage. The overall goal of the data mining process is to extract information from a data set and transform it into an understandable data in order to aid user decision making [9]. Data mining is being used in several applications like banking, insurance, hospital and Health informatics. In case of health informatics, Data mining plays a vital role in helping physicians to identify effective treatments, and Patients to receive better and more affordable health services. In hematology laboratory, it has become a powerful tool in managing uncountable laboratory information in order to seek knowledge that is underlying or within any given information.

Comparison of Different Classification Techniques Using WEKA for diabetic hematological data is a challenging and interesting task in medical research area. To find out which classification algorithms is better it is very difficult to compare different classification algorithms in different dataset. Our dissertation concerns with which algorithm, which is capable to Diagnose Diabetic Hematological data accurately as well as quickly. With this purpose to perform a better approach, we divide this problem of Diabetic Hematology Data comments into two phases: Data Collection, Classification algorithm, . We proceed in the following ways to achieve our purpose successfully.



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 1, January 2018

- We are going to collect diabetic hematological data from datasets.
- We are going to apply diabetic hematological data in WEKA then find three classification algorithms performance.
- We studied various journals and articles regarding performance evaluation of Data Mining algorithms on various different tools, some of them are described here.
- There are related works using data mining techniques to diagnose several types of diseases and phenomena, such as Automated Diagnosis of Thalassemia Based on Data Mining Classifiers, etc. And many others tried to find their own formula. This paper presents an investigation for thalassemia existence by using data mining classifiers depending on CBC. They do that but they say MCV is the main feature. They should need use Hemoglobin is the main feature to classified thalassemia.
- K.Rajesh et al in their paper "Application of Data Mining Methods and Techniques for Diabetes Diagnosis." they provide a comparative analysis of different algorithms. This project aims for mining the relationship in diabetes data for efficient classification. But they need proposed a model that can diagnose diabetes dataset.
- Satish Kumar David et al in his research paper "Comparative Analysis of Data Mining Tools and Classification Techniques using WEKA in Medical Bioinformatics." Studied the performance of Tree Random Forest, J48 decision tree, Bayes Naïve Bayes and Lazy.IBK. In this paper, they compared algorithms based on their accuracy, learning time and error rate. They observed that there is a direct relationship between execution time in building the tree model and the volume of data records, while there is also an indirect relationship between execution time in building the model and the attribute size of the data sets. Through experiment, they conclude that Bayesian algorithms have better classification accuracy over and above compared algorithms.
- Salvitha et al in their article "Evaluating Performance of Data Mining Classification Algorithm in Weka". They provide performance of different dataset use data mining classification. The main aim of this paper is to judge the performance of different data mining classification algorithms on various datasets.
- Nookala et al in their article "Performance Analysis and Evaluation of Different Data Mining Algorithms used for Cancer Classification." In this study, they have made a comprehensive comparative analysis of 14 different classification algorithms and their performance has been evaluated by using 3 different cancer data sets. The results indicate that none of the classifiers outperformed all others in terms of the accuracy when applied on all the 3 data sets. Most of the algorithms performed better as the size of the data set is increased. They recommend the users not to stick to a particular classification method and should evaluate different classification algorithms and select the better algorithm.
- Vaithyanathan et al in their paper "comparison of different classification techniques using different datasets". They used three dataset from benchmark data set (UCI) and they used four classifier algorithms J48, Multilayer Perception, Bayes Net, and Naïve Bayes Update. This work has been carried out to make a performance evaluation above algorithms.
- Tiwari et al in their research paper "Performance analysis of Data mining algorithms in Weka". The aim of their paper is to judge the accuracy of different data mining algorithms on various data sets.
- Bin Othman et al "Comparison of different classification techniques using WEKA for breast cancer". In this paper they present the comparison of different classification techniques using Waikato Environment for Knowledge Analysis or in short, WEKA. The aim of their paper is to investigate the performance of different classification or clustering methods for a set of large data. The algorithm or methods tested are Bayes Network, Radial Basis Function, Pruned Tree, Single Conjunctive Rule Learner and Nearest Neighbors Algorithm. The best algorithm based on the breast cancer data is Bayes network classifier with an accuracy of 89.71% and the total time taken to build the model is at 0.19 seconds. Bayes network classifier has the lowest average error at 0.2140 compared to others.
- All the previous works tried to makes a model to diagnosis diosis, and most of them just try to use one data

Mining technique they consider it the best one without any comparison with the other techniques in the domain. In this study, I will used more than one classifier to get most significance one, and make a model that can easily diagnosis hematological data comments.

The main contributions of the thesis are summarized follow:

- J48 based on decision tree algorithm has been achieved to classify different types of hematological data comment.
- Naïve Bayes algorithm has been obtained for high probability of hematological data comment
- 
- Multilayer perception algorithm has been obtained mathematical or computational model for information processing based on a connectionist approach.
- A comparison with different classification techniques has made with optimal features to show which method is appropriate for hematological data.



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 1, January 2018

## II. MATERIAL AND METHODS

We have used the popular, open-source data mining tool Weka (version 3.8) for this analysis. A data set has been used and the performance of a comprehensive set of classification algorithms (classifiers) has been analyzed. The analysis has been performed on a HP Windows 10 system with Intel® Core™ i7 CPU, 2.30 GHz Processor and 8.00 GB RAM.

The diabetic hematological parameter are composed of Number of times pregnant, Plasma glucose concentration a 2 hours in an oral glucose tolerance test, Diastolic blood pressure (mm Hg), Triceps skin fold thickness(mm),2-Hour serum insulin ( $\mu$  U/ml), Body mass index (weight in kg/(height in m)<sup>2</sup>),Diabetes pedigree function ,Age (years),Class variable (0 or 1) This data has been provided by National Institute of Diabetes and Digestive and Kidney Diseases

### 2.1 Dataset and Preprocessing

The experimental dataset consists of 768 samples. Its attributes represents the features as Number of times pregnant, Plasma glucose concentration a 2 hours in an oral glucose tolerance test, Diastolic blood pressure (mm Hg), Triceps skin fold thickness (mm),2-Hour serum insulin ( $\mu$  U/ml), Body mass index (weight in kg/(height in m)<sup>2</sup>),Diabetes pedigree function, Age (years),Class variable (0 or 1)

Table 1: CBC Test Features

Term	Normal Value
Number of times pregnant	N. A
Plasma glucose concentration	70-130mg/dl
Diastolic blood pressure	80-90
Triceps skin fold thickness	N A
2- Hour serum insulin	2.6-24.9
Body mass index	18.5-24.9
Diabetes pedigree function	N A
AGE	N. A
Class variable	0



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 1, January 2018

## 2.2 Classification Methods

Three candidate classifiers are considered in this study: Decision Tree (J48), Naïve Bayes, and ZeroR.

### 2.2.1 J48 Algorithm

J48 algorithm is called as optimized implementation of the C4.5 or improved version of the C4.5. The output given by J48 is the Decision tree. A Decision tree is same as that of the tree structure having different nodes, such as root node, intermediate nodes and leaf node. Each node in the tree contains a decision and that decision leads to our result as name is decision tree. Decision tree divide the input space of a data set into mutually exclusive areas, where each area having a label, a value or an action to describe or elaborate its data points. Splitting criterion is used in decision tree to calculate which attribute is the best to split that portion tree of the training data that reaches a particular node. Practical elucidation can be seen in Fig 3.

### 2.2.2 Zero R Algorithm

ZeroR is the simplest classification method which relies on the target and ignores all predictors. ZeroR classifier simply predicts the majority category (class). Although there is no predictability power in ZeroR, it is useful for determining a baseline performance as a benchmark for other classification methods. Practical elucidation can be seen in Fig 4.

### 2.2.3 Naive Bayes

Naive Bayes implements the probabilistic Naïve Bayes classifier. Naïve Bayes Simple uses the normal distribution to model numeric attributes. Naïve Bayes can use kernel density estimators, which develop performance if the normality assumption is grossly correct; it can also handle numeric attributes using supervised discretization. Naïve Bayes Updateable is an incremental version that processes one request at a time. It can use a kernel estimator but not discretization. Practical elucidation can be seen in Fig 3.

## III.RESULTS AND DISCUSSION

The results and a detailed classification accuracy analysis emphasizing on the classification errors will be presented in following Sections. Three experiments were conducted in each type: the first one is to measure the performance of the decision tree classifier; the second one is to measure the performance of the naïve bayes classifier, the third one to measure the performance of the Zero R.

### 3.1 Experiments with full features

In these experiments we used the whole records attributes of each sample. The decision tree classifier gives a result with general accuracy of 97.16%, the naïve bayes classifier gives a result with general accuracy of 70.28%, and finally the zero R gives a result with general accuracy of 86.55% as shown in Fig.2, Table3.

# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 1, January 2018

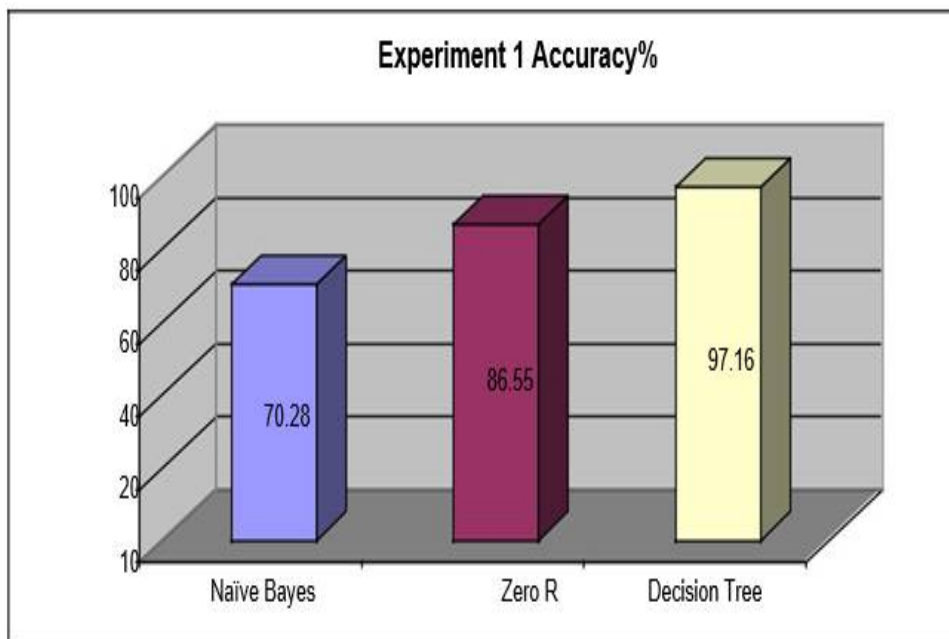


Figure 2: experiment 1 classifiers accuracy values

Table 3: Simulation Result of Each Algorithm for Exprement1

Algorithm (Total instances,425)	Correctly Classified Instances %	Incorrectly Classified Instances %	Time Taken (seconds)	Kappa statistic
J48 Decision tree	73.8281	26.1719	0.28	0.4164
Zero R	65.1042	34.8958	0	0
Naïve Bayes	76.3021	23.6979	0.06	0.4664

Based on the above Fig. and Table, we can clearly see that the highest accuracy is 76.3021% and the lowest accuracy is 65.1042%. We can say that Naive Bayes is better.

# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijirccce.com](http://www.ijirccce.com)

Vol. 6, Issue 1, January 2018

## J48 Decision Tree

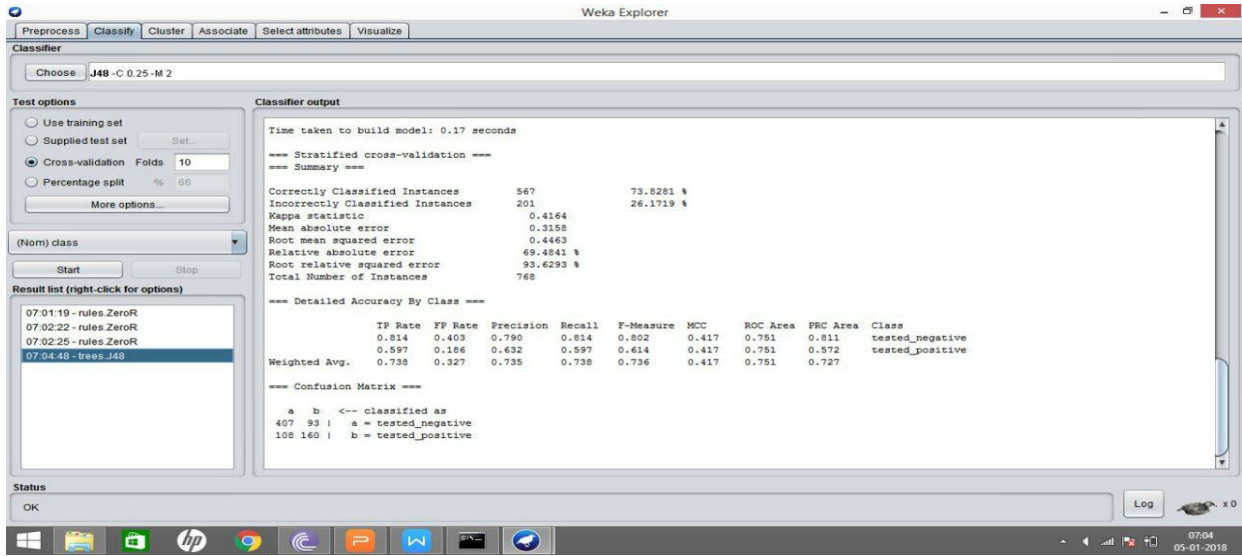


Figure 3: Practical elucidation of J 48 Decision Tree in WEKA software.

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a dataset into smaller and smaller subsets while at the same time an associated decision tree is incrementally developed. The final result is a tree with decision nodes and leaf nodes. A decision node (e.g., Outlook) has two or more branches (e.g., Sunny, Overcast and Rainy). Leaf node (e.g., Play) represents a classification or decision. The top most decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

## ZeroR

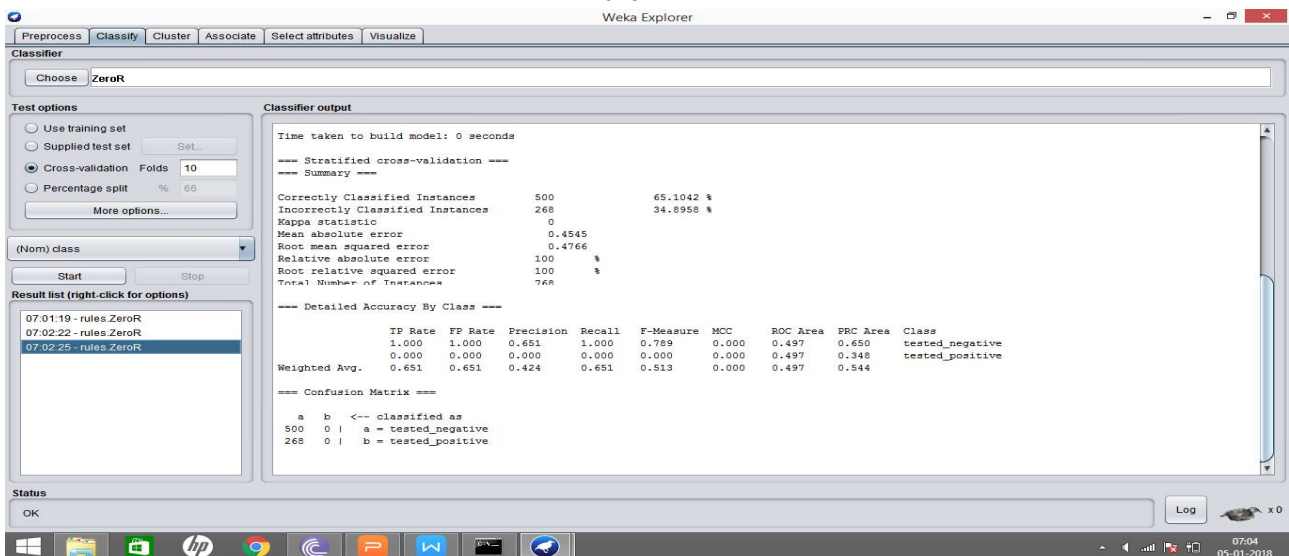


Figure 4: Practical elucidation of ZeroR in WEKA software.

# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 1, January 2018

ZeroR is the simplest classification method which relies on the target and ignores all predictors. ZeroR classifier simply predicts the majority category (class). Although there is no predictability power in ZeroR, it is useful for determining a baseline performance as a benchmark for other classification methods. Basically used to construct a frequency table for the target and select its most frequent value.

## Naïve Bayes

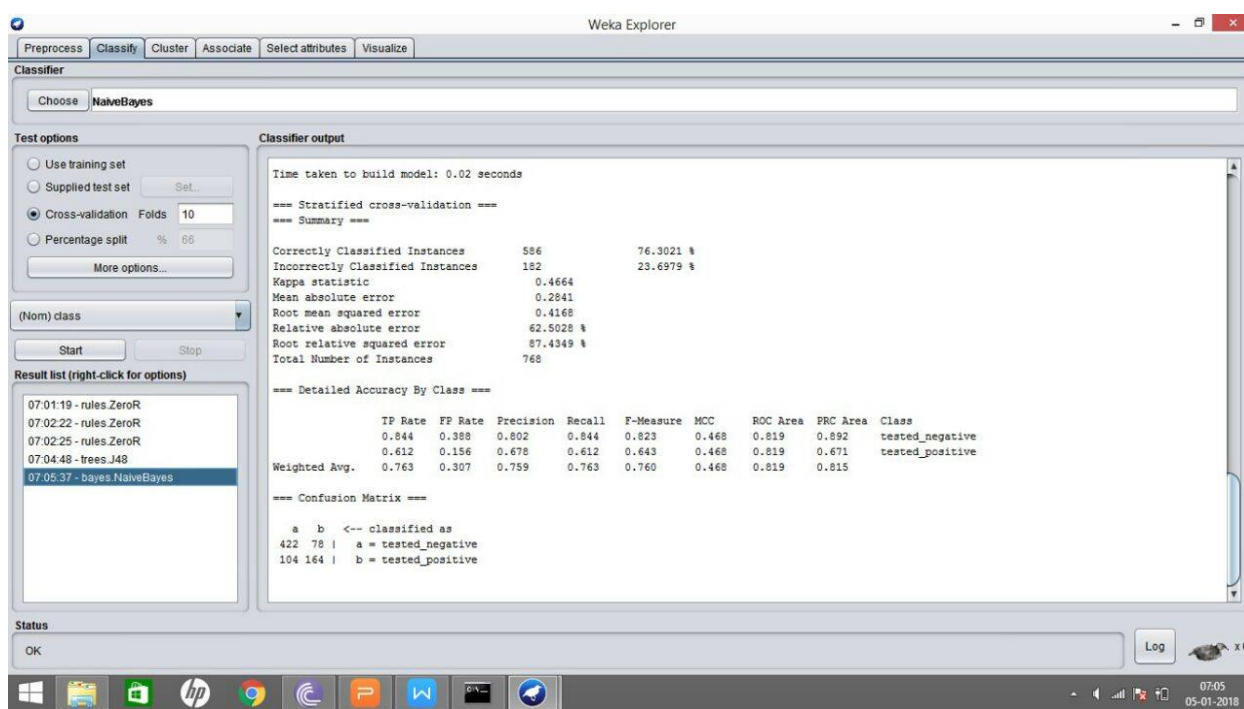


Figure 5: Practical elucidation of Naïve Bayes in WEKA software.

The Naive Bayesian classifier is based on Bayes' theorem with independence assumptions between predictors. A Naive Bayesian model is easy to build, with no complicated iterative parameter estimation which makes it particularly useful for very large datasets. Despite its simplicity, the Naïve Bayesian classifier often does surprisingly well and is widely used because it often outperforms more sophisticated classification methods.

## IV. CONCLUSION

As a conclusion, we have met our objective which is to evaluate and investigate three selected classification algorithms based on Weka. The best algorithm based on the diabetic data is Naive Bayes with an accuracy of 76.3021% and the total time taken to build the model is at 0.06 seconds. Naive Bayes classifier has the lowest average error at 29.71% compared to others. These results suggest that among the machine learning algorithm tested, Naive Bayes classifier has the potential to significantly improve the conventional classification methods for use in medical or in general, bioinformatics field.



# International Journal of Innovative Research in Computer and Communication Engineering

(A High Impact Factor, Monthly, Peer Reviewed Journal)

Website: [www.ijircce.com](http://www.ijircce.com)

Vol. 6, Issue 1, January 2018

## REFERENCES

- [1] Vaithyanathan, V., K. Rajeswari, Kapil Tajane, and Rahul Pitale. "Comparison of Different Classification Techniques Using Different Datasets." *Vol.6, no. 2 (2013)*.
- [2] Sharma, Narendra, Aman Bajpai, and Mr Ratnesh Litoriya. "Comparison the various clustering algorithms of weka tools." *Volume2, no.5 (2012)*.
- [3] Salvithal, Nikhil N., and R. B. Kulkarni. "Evaluating Performance of Data Mining Classification Algorithm in Weka." *Vol 2., no. 10 (2013)*.
- [4] Khan, S. A., J. H. Epstein, K. J. Olival, M. M. Hassan, M. B. Hossain, K. B. M. A. Rahman, M. F. Elahi et al. "Hematology and serum chemistry reference values of stray dogs in Bangladesh." *Vol. 1: 13-20 (2011)*.
- [5] Zhang, Wenjing, Donglai Ma, and Wei Yao. "Medical Diagnosis Data Mining Based on Improved Apriori Algorithm." *Journal of Networks 9, no. 5 (2014): 1339-1345*.
- [6] Nookala, Gopala Krishna Murthy, Bharath Kumar Pottumuthu, Nagaraju Orsu, and Suresh B. Mudunuri. "Performance analysis and evaluation of different data mining algorithms used for cancer classification." *International Journal of Advanced Research in Artificial Intelligence (IJARAI) 2, no. 5 (2013)*.
- [7] Tiwari, Mahendra, Manu Bhai Jha, and OmPrakash Yadav. "Performance analysis of Data Mining algorithms in Weka." *IOSR Journal of Computer Engineering (IOSRJCE) ISSN (2012): 2278-0661, Vol.6, Iss.3*.
- [8] Kaushik H. Raviya, Biren Gajjar "Performance Evaluation of Different Data Mining Classification Algorithm Using WEKA" *Vol. 2, Issue. 1. (2013)*.
- [9] Saichanma, Sarawut, Sucha Chulsomlee, Nonthaya Thangrua, Pornsuri Pongsuchart, and Duangmanee Sanmun. "The Observation Report of Red Blood Cell Morphology in Thailand Teenager by Using Data Mining Technique." *Advances in hematology 2014 (2014)*.
- [10] bin Othman, Mohd Fauzi, and Thomas Moh Shan Yau. "Comparison of different classification techniques using WEKA for breast cancer." *3rd Kuala Lumpur International Conference on Biomedical Engineering 2006. Springer Berlin Heidelberg, 2007*.
- [11] Elshami, E. H., & Alhalees, A. M. (2012). *Automated Diagnosis of Thalassemia Based on Data Mining Classifiers. In The International Conference on Informatics and Applications (ICIA2012) (pp. 440-445). The Society of Digital Information and Wireless Communication*.
- [12] Pankaj saxena & sushma lehri, International Journal of Computer & Communication Technology ISSN (PRINT): "Analysis of various clustering algorithms of data mining on health informatics" *Vol. 4, Issue. 2. (2013)*,
- [13] Ms S. Vijayarani , Ms M. Muthulakshmi, International Journal of Advanced Research in Computer and Communication Engineering: "Comparative Analysis of Bayes and Lazy Classification Algorithms". *Vol.2, Issue. 8, (2013)*.
- [14] Rajesh, K., and V. Sangeetha. "Application of data mining methods and techniques for diabetes diagnosis." *International Journal of Engineering and Innovative Technology (IJEIT) Volume. 2, Issue. 3 (2012)*.
- [15] David, Satish Kumar, Amr TM Saeb, and Khalid Al Rubeaan. "Comparative Analysis of Data Mining Tools and Classification Techniques using WEKA in Medical Bioinformatics." *Computer Engineering and Intelligent Systems 4, no. 13 (2013): 28-38*.