# COMPARISION AND COMBINATION OF MINING TECHNIQUES FOR GENE ANALYSIS TO IDENTIFY DENGUE

Dave Kaveri Atulbhai

*M.E. Student, Computer Engineering, Gujarat Technical University, Gujarat, India*

## ABSTRACT

*Data Mining is focused on deriving knowledge from multiple data formats using intelligent analytics techniques. Data mining is a field having various techniques that converts the raw data into useful information in various research fields. There are two primary goals for data mining prediction and description. Prediction involves fields or variables in the data sets to predict unknown or future values of other diseases possibilities. On the other hand description involves finding of pattern describing the data that can be present in knowledge base provided for disease prediction. We can predict diseases like hepatitis, Lung cancer liver disorder, breast cancer or heart diseases, diabetes etc, it helps in finding the patterns to decide future trends in medical field. We can use in finding the knowledge and prediction, detection of Dengue as well.*

**Keywords :-** *Data Mining, Raw Data, Information, Knowledge Discovery, Prediction, Description, Medical Analysis, Dengue.*

## 1. INTRODUCTION:

Information technology has generated large amount of data-base and huge amount of data in various research fields. To research in knowledge mining has given rise to store data and manipulate previously stored data for further decision making process. Dengue is a threatening disease caused by female mosquitoes. It is typically found in widespread hot regions. From long periods of time, Experts are trying to find out some of features on Dengue disease so that they can rightly categorize patients because different patients require different types of treatments [2]. There are various data mining techniques available with suitable dependent on domain application. By using data mining we can examine large amount of routine samples collected in disease prediction. Best results are achieved by balancing knowledge of experts for describing the problem and goals with search capabilities. Hospitals must also want to minimize cost of clinical test. It can be achieved by employing appropriate computer based information and decision sport system. Here, data mining plays an important role to give many results faster and accurate by using various algorithms [8].

## 2. DATA MINING PROCESS AND MEDICAL SCIENCE:

Data mining is the process of selecting, exploring and modeling large amounts of data in order to discover unknown patterns or relationships which provide a clear and useful result to the data analyst [8].

### 2.1 Data mining involves the following Steps:

1. Problem definition: The first step is to identify goals.

2.  Data exploration: 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 for more robust analysis.
4.  Modeling: Based on the data and the desired outcomes, a data mining algorithm or combination of algorithms is selected for analysis. The specific algorithm is selected based on the particular objective to be achieved and the quality of the data to be analyzed.
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.
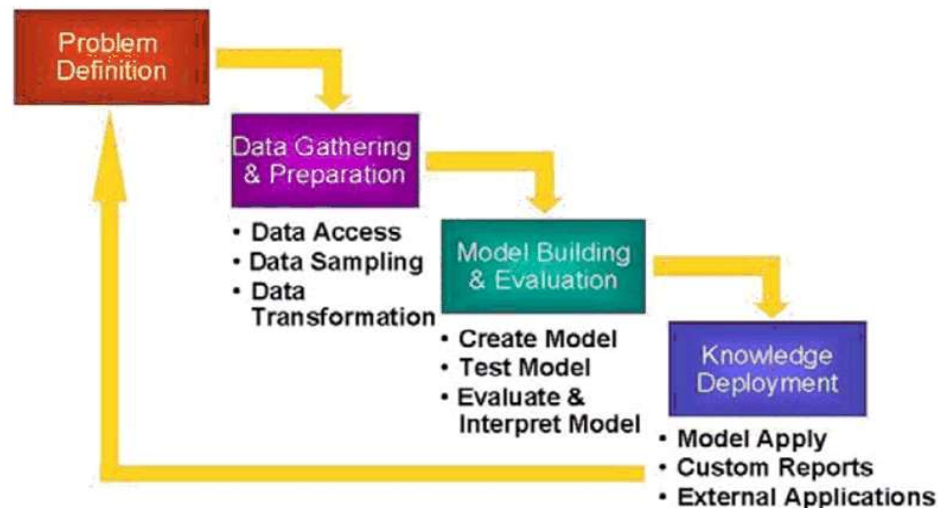


**Fig -1**: Data Mining Process

**2.2 Data mining for Medical Science:**

Data mining is a relatively new field of research whose major objective is discovering knowledge from large amounts of data. In medical and health care areas a large amount of data is becoming available due to regulations and the availability of computers and technology. On the one hand, practitioners are expected to use all this data in their work but, at the same time, such a large amount of data cannot be processed by humans in a short time to make diagnosis, prognosis and treatment schedules. A major objective is to evaluate data mining tools in medical and health care applications to develop a tool that can help make timely and accurate decisions.

The Healthcare industry is among the most information intensive industries. Medical information, knowledge and data keep growing on a daily basis. It has been estimated that an acute care hospital may generate five terabytes of data a year [7]. The ability to use these data to extract useful information for quality healthcare is crucial. Medical informatics plays a very important role in the use of clinical data. In such discoveries pattern recognition is important for the diagnosis of new diseases and the study of different patterns found when classification of data takes place.

The application of artificial intelligence in healthcare is relatively new. Data mining can be applied to the medical databases, which will predict or classify the data with a reasonable accuracy. For a good prediction or classification the learning algorithms must be provided with a good training set from which patterns and rules can be extracted to help classify the testing dataset. We can predict diseases like hepatitis, Lung cancer liver disorder, breast cancer or heart diseases, diabetes, arbovirus disease etc, it helps in finding the patterns to decide future trends in medical field.

## 3.PREDICTION OF AN AIRBOVIRUS – DENGUE DISEASE:

Dengue infection is an epidemic disease typically found in tropical region. Symptoms of the disease show rapid and violent for patients in a short time. The World Health Organization (WHO) classifies the dengue infection as Dengue Fever (DF) and Dengue Hemorrhagic Fever (DHF). Symptoms of DHF are divided into 4 types. The problem might be happen when an expert misdiagnoses dengue infection. For Example, an expert diagnosed a patient as non dengue or DF even if a patient was a DHF patient. That might be the cause of dead if patient did not receive treatment. Therefore, we selected data mining approach to solve this problem [2].

For properly categorizing our dataset, different classification techniques are used. These techniques are Naïve Bayesian, REP Tree, Random tree, J48, SMO, SVM, Decision Tree Approach, and Spatial Data Analysis as well.

In this survey study, [1] Represents Comparative Analysis of Machine learning techniques for classification of Arbovirus - Research efforts have reported with increasing confirmation that the support vector machines (SVM) have greater accurate diagnosis ability. The accuracy of svm outperforms the Naïve bayes technique. As shown in the experiment results, support vector machine has the highest classification precision most of the time. However support vector machine is very time consuming because of more parameters, demands more computation time. [2] Represents Data Mining of Dengue Infection Using Decision Tree in which each dataset consists of more than 400 attributes. To accomplish the knowledge discovery task, we consider employing decision tree as a data mining tool. We propose a set of meaningful attributes from the temporal data. Our experiments are divided into 4 parts. The first two experimental results show the useful knowledge to classify dengue infection from 2 different datasets respectively. Another objective of this research is to detect the day of defervescence of fever which is called day0. At the end we obtained very low accuracy in day-4 as we found that the tree is over fit. The experimental results shown that the decision tree approach did not suit this task thus we think we should to select a new classification approach in the future works. [3] Represents the classification techniques to determine the population of Dengue fever infected cases in Jhelum district and in surrounding areas geographically. So, we can compare performance of different classification techniques. Objective of this study also includes the comparison of different classification algorithms with the help of graphs, based on our dataset. We have implemented all the techniques by using weka tool and all the procedure of implementation is within it. At the end, After analysis of our dataset with each technique we are paralleling them in the conclusion. When we have done the comparison among all of them we concluded that naïve Bayes Technique is greatest among all others. As the accuracy of Naive bayes is 92% which was biggest of all. Naïve Bayes is the best also for the aim that it gives the probability and efficiency while Random Tree and REP Tree do not give us probability. It has concluded that NB and J48 are the top performance classifier techniques by way that, they has achieved an accuracy of 92% and 88%, takes fewer time to run and shows ROC area=0.815, and had smallest error rate. We can find the following comparison chart below.[4] Represents implementation the automated system to reorganize the Dengue fever using the Microscope blood image report. Microscope image report as input and signals are filtered and the feature characteristics are extracted, the features are fed to the neural networks. Classification is carried using and the Back Propagation Network (BPN).Which gives the 98% accurate result with a short period. [5] Represents prediction of chronic kidney disease. Study is focused on the usage of classification techniques in the field of medical science and bioinformatics. The main focus of this paper is Chronic-Kidney-Disease prediction using weka data mining tool and its usage for classification in the field of medical bioinformatics. 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, K appa statistics and ROC Area. In the experiments Multilayer perception algorithm gives better classification accuracy and prediction performance to predict chronic kidney disease (CKD) using relevant dataset available at UCI machine learning repository. [6] Represents a study of different data mining techniques that can be employed in automated HCV infection prediction systems. The system extracts hidden knowledge from a historical HCV patient's database. The models are trained and validated against a test dataset. Three classification techniques which are Naïve Bayes, decision trees and neural networks applied on three HCV database of different size, then determine Performance accuracy and effectiveness of each technique using Lift Chart. Performance analysis of these classifiers is summarized in Table.

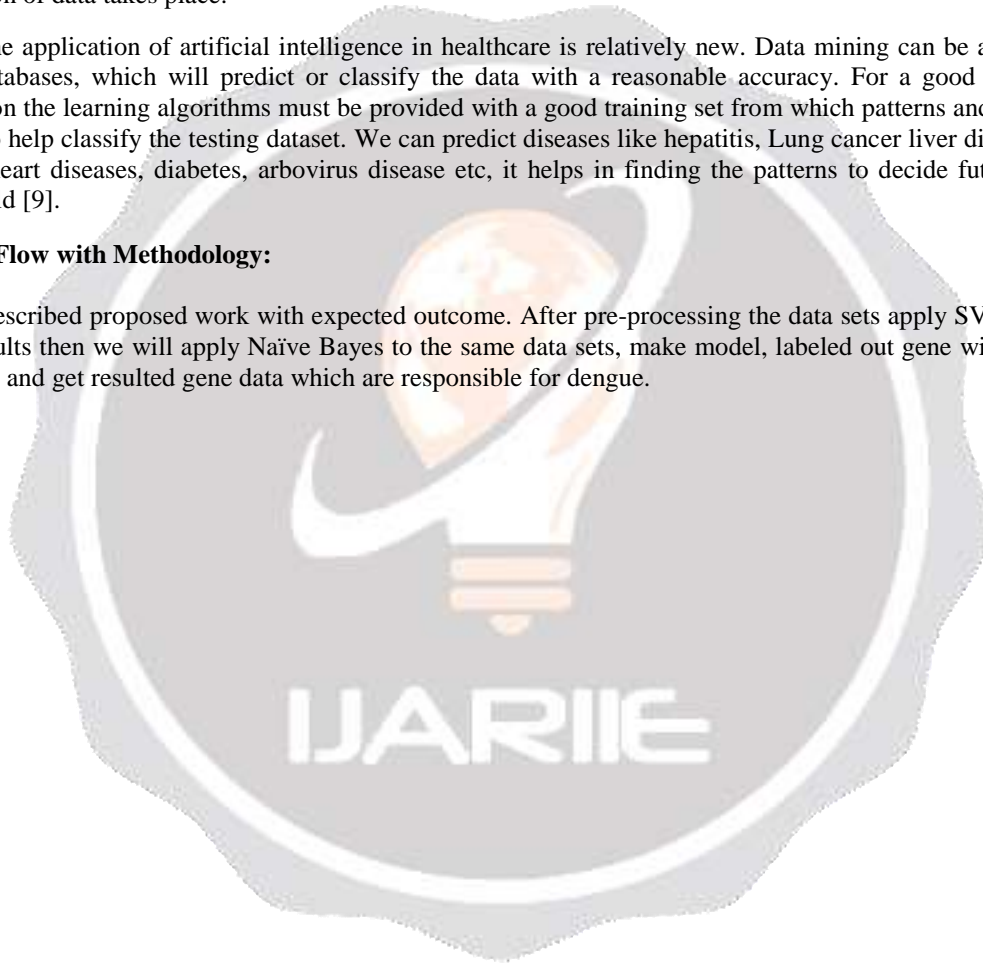## 4. METHODOLOGY, EXPECTED OUTCOMES AND IMPLEMENTATION STRATEGIES:

Data mining is a relatively new field of research whose major objective is discovering knowledge from large amounts of data. In medical and health care areas a large amount of data is becoming available due to regulations and the availability of computers and technology. On the one hand, practitioners are expected to use all this data in their work but, at the same time, such a large amount of data cannot be processed by humans in a short time to make diagnosis, prognosis and treatment schedules. A major objective is to evaluate data mining tools in medical and health care applications to develop a tool that can help make timely and accurate decisions [10].

The Healthcare industry is among the most information intensive industries. Medical information, knowledge and data keep growing on a daily basis. It has been estimated that an acue care hospital may generate five terabytes of data a year [10]. The ability to use these data to extract useful information for quality healthcare is crucial. Medical informatics plays a very important role in the use of clinical data. In such discoveries pattern recognition is important for the diagnosis of new diseases and the study of different patterns found when classification of data takes place.

The application of artificial intelligence in healthcare is relatively new. Data mining can be applied to the medical databases, which will predict or classify the data with a reasonable accuracy. For a good prediction or classification the learning algorithms must be provided with a good training set from which patterns and rules can be extracted to help classify the testing dataset. We can predict diseases like hepatitis, Lung cancer liver disorder, breast cancer or heart diseases, diabetes, arbovirus disease etc, it helps in finding the patterns to decide future trends in medical field [9].

**4.1 Work Flow with Methodology:**

Here, we described proposed work with expected outcome. After pre-processing the data sets apply SVM algorithm and get results then we will apply Naïve Bayes to the same data sets, make model, labeled out gene with respective expressions and get resulted gene data which are responsible for dengue.
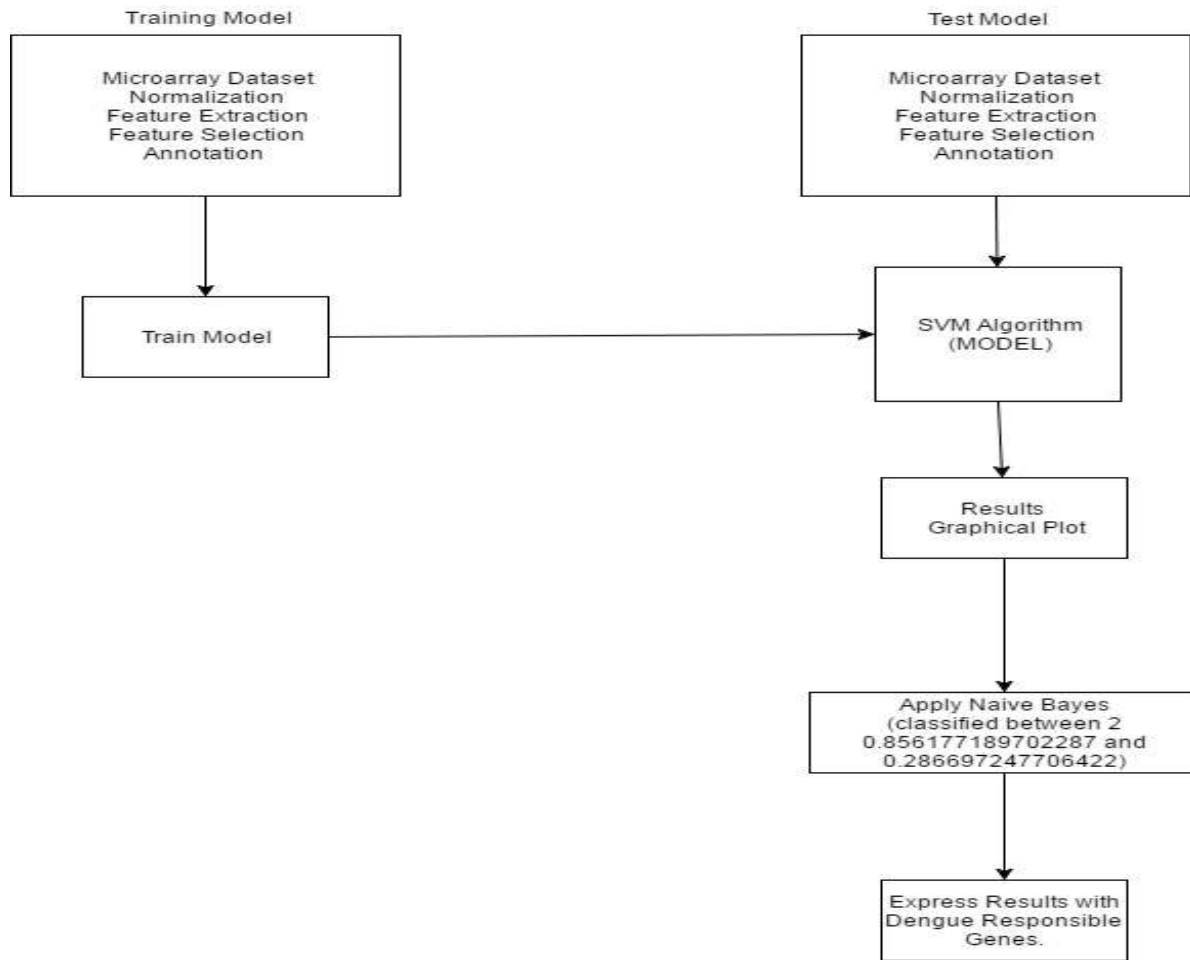
**Fig – 4.1** Flow Diagram

**4.2 Expected Outcome:**

Classified highly down regulated gene which are responsible for dengue**.**

**4.3 Implementation Strategies:**

| Input Dataset | Reports of Patients |
|---|---|
| Total Dataset | 758 MB to 1256 MB |
| Software | R Studio |
| Library | R Library |
| Programming Language | R Language by R Studio |

**4.4 IMPLEMENTATION AND RESULT:**

**Let's consider following work flow Algorithm for mentioned system:**

1.  Start with Micro Array data set.

#Plot Histogram and Boxplot before preprocessing.

**Fig – 4.1 (A):** Before Normalization

2. Apply Normalization with RMA.

3. Express data in log2 expression as data not in a feasible range.
   #Plot Histogram and Boxplot after Normalization.

**Fig – 4.1 (B):**   After Normalization

4. Feature Extraction.
   #Apply eBayes Algorithm for p.value, Log FC value, Adj p.value calculations.

**Fig – 4.1 (C)**: Feature Extraction

5.   Feature Selection.
      #values with p.value < 0 used to select differently expressed or say significant genes.
      #Log FC calculated bay p.value analysis.

*If ( logFC < 0)*
*Then negative logFC value which represents logarithmic foldness of down regulation*
*Else if ( logFC > 0)*
*Then positive logFC value which represents logarithmic foldness of up regulation*

Note: LogFC = 0 consider as default bu LIMMA Package Library.

#Genes which shows negative logFC value denotes the downregulation of their corresponding gene function

#which means these genes are infected during disease.

#This table shows the predicted negative logFC values

**Fig – 4.1 (D):** Feature Selection

6.  Apply Annotation.
    #generated data set having Database & CDF (Chip Description Format), we consider Database as Database having gene name and description as well.
    #to convert in readable format we are using Annotation.



**Fig – 4.1 (E):** Annotation

7.  Check Duplication & Redundancy.

8.  Create Train Model.

9.  Apply pre-processed test dataset as input to machine learning algorithm.

10. SVM Algorithm with e1071 Package in R.Studio.

11. Analyze Generated results and graphical plot.

**Fig – 4.1 (F):** SVM Prediction Result Plot.

12. Interchange the row and column of csv file with header names - Name and Result.

13. Apply Naive Bayes.

14. Result classified between 2 values 0.856177189702287 and 0.286697247706422.

15. Give the value of correctly interpreted down regulated genes.
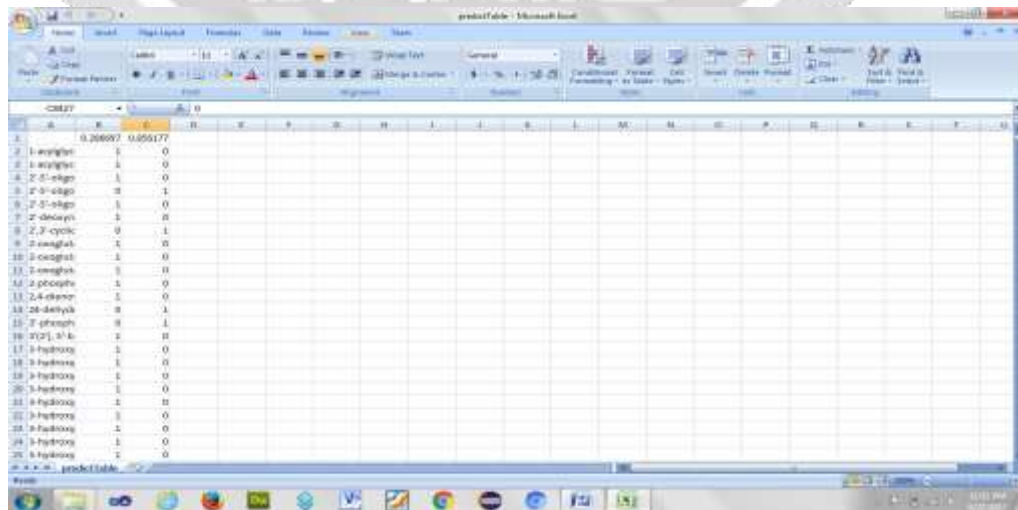    #Express Results with Dengue Responsible Genes.



**Fig – 4.1 (G):** Prediction Result.

## 5. COMPARISION & COMBINATION:

The best accuracy is achieved by svm - 0.9042. The good performance of SVM indicates that they can possibly be adjusted to improve specific ad hoc methods for prediction of susceptibility to complex diseases such as Dengue. The total accuracy for 10-fold and 100 cross-fold validation on training data for svm was 90.3 in both cases. We show that sensitivity of naIve Bayes is more than svm in many classification tasks in bioinfonnatics and related scientific fields. The accuracy of svm outperforms the Naive bayes technique. As shown in the experiment results, support vector machine has the highest classification precision most of the time .However support vector machine is very time consuming because of more parameters, demands more computation time. As well we have calculated that SVM having more risk rate compare to naïve bayes & efficiency of naïve bayes is little good than all the methods for bio medical gene analysis.

| *Measures* | *SVM (%)* | *Naïve Bayes (%)* |
|---|---|---|
| Sensitivity | 47.23 | 98.77 |
| Efficiency | 72.77 | 94.05 |
| Accuracy | 90.18 | 78.05 |
| Risk Rate | 83.43 | 15.98 |

**Table – 5.1 (A):** Comparisons

In this mentioned model we have combined most suitable methods for medical data sets – SVM & Naïve Bayes. As per comparative analysis we can justify that SVM is little more accurate with high risk rate while Naïve Bayes having less accuracy with more efficiency and lowest risk rate. So, we classified gene with more attribute expressions by SVM then apply Naïve Bayes to classified gene in respected range to get dengue infected data from gene data sets. We can increase Efficiency of this model by 29.48% with 70.52% Accuracy.

## 6. CONCLUSION:

Firstly, there exists a wide class of algorithms and techniques for information extraction and knowledge discovery in medical science. Best results are achieved by balancing knowledge of experts for describing the problem and goals with search capabilities. Hospitals must also want to minimize cost of clinical test. It can be achieved by employing appropriate computer based information and decision sport system. Here, data mining plays an important role to give many results faster and accurate by using various algorithms.By analyzing the different techniques in mentioned researches we can say that more accurate method for medical science is Naïve Bayes as classification methods is more useful in medical science and disease prediction. But, we can combine different methods to get accurate knowledge discovery. Previously studies are along with comparative analysis only.

This paper presents a generic survey of the Dengue Disease Prediction using different data mining techniques. All techniques we explained are widely using in medical science. Data set for dengue prediction is DNA microarray data which have information of gene's expression responsible for dengue virus' successful attack. This data is available on GEO dataset NCBI. For this project .cel format is used as mostly microarray chip machines saved information in this format. Before doing any analysis, Normalization of these data is necessary. Hence, after analyzing all the researches we can conclude that combination of different methods can help to acquire more accurate results.

We apply RMA to normalize microarray data. Then, features selection by using Wrapper Method. We will improve result of detection and diagnose relevant data using 2 way methods  –  SVM for classification on gene expressions and Naïve Bayes to classified dataset for high level of classification and accurate and efficient results about dengue detection. We can increase Efficiency of this model by 29.48% with 70.52% Accuracy. In future we will predict survivability of patients on basis of gene expressions and their value changes.

## 4. REFFERENCES:

[1] Shameem A. Fathima and NisarHundewale, Senior Member, "Comparative Analysis of Machine learning Techniques for classification of Arbovirus", Proceedings of the IEEE-EMBS International Conference on Biomedical and Health Informatics (BHI 2012) Hong Kong and Shenzhen, China, 2-7 Jan 2012.

[2] Kamran Shaukat1,NayyerMasood,Sundas Mehreen1 and UlyaAzmeen ," Dengue Fever Prediction: A Data Mining Problem", Dar et al.J Data Mining Genomics Proteomics 2015.

[3] N.Subitha and Dr.A.Padmapriya "Diagnosis for Dengue Fever Using Spatial Data Mining",International Journal of Computer Trends and Technology (IJCTT) ,August 2013.

[4] Daranee, Pratap Suriyaphol and Nuanwan," Data Mining of Dengue Infection Using Decision Tree",July 2015.

[5] cLambodar Jena, Narendra Ku. Kamila," Distributed Data Mining Classification Algorithms for Prediction of Chronic- Kidney-Disease", International Journal of Emerging Research in Management &Technology, November 2015.

[6] Ahmed A.A.Radwan and TarekAbd El-Hafeez," An Analysis of hepatitis C virus prediction using different data mining techniques",December 2013.Reecha P. Yadav, Vinuchackravarthy Senthamilarasuand Krishnan Kutty, Sunita P. Ugale,".

[7] Huang, H. et al. "Business rule extraction from legacy code", Proceedings of 20th International Conference on Computer Software and Applications, IEEE COMPSAC'96, 1996, pp.162-167

[8] Aarti Sharma,Rahul Sharma,Vivek Kr. Sharma,Vishal Shrivatava, "Application of Data Mining – A Survey Paper", Department of CS & IT, A.C.E.I.T., Jaipur International Journal of Computer Science and Information Technologies- 2014.

[9] Arun George Eapen, "Application of Data mining in Medical Applications", University of Waterloo, Waterloo, Ontario, Canada, 2004.

[10] Irshad Ullah, "Data Mining Algorithms And Medical Sciences", International Journal of Computer Science & Information Technology (IJCSIT), Vol 2, No 6, December 2010