# Comparative Study of Data Mining Classification Algorithms in Heart Disease Prediction

Durgadevi. A[1], Prof. Saravanapriya.K[2]

[1,2]PG and Research Department of Computer Applications Sacred Heart College, Tirupattur, Vellore

*Abstract:* This research paper intends to provide an accuracy of the best algorithm from the classification algorithms. The main objective of this research work is to predict more accurately the presence of heart disease in patients. The dataset used here is from the UCI Machine Learning Repository based Hungarian-14-heart disease with 294 instances. In this research paper, we use weka tool with five classifiers like Naïve Bayes, Logistic function, RBF Network function, Decision Table rule, SMO function and their performance on the diagnosis has been compared. From that Naïve Bayes provides 86% accurate result in 0 seconds and RBF Network provides 86% of accuracy in 0.17 seconds. The research result shows both the Naïve Bayes outperforms with a time duration of 0 seconds to build the model.

*Keywords:* CVD - Cardio Vascular Disease, Machine learning algorithm, Supervised and Unsupervised algorithms, Logistic function Algorithm, RBF - Radial Basis Function Network, SMO function -Sequential Minimal Optimization Algorithms.

## 1. INTRODUCTION

Various data mining algorithms and techniques like Classification, Clustering, Artificial Intelligence, Regression, Neural Networks, Association Rules, Genetic Algorithm, Decision Trees, Nearest Neighbour method etc., are used for knowledge discovery from databases. In this paper we were discussed about the five different classifiers used to predict the heart disease with 14 attributes.

*Classification algorithms:* Classification is the most commonly applied data mining technique, which needs a set of pre-classified examples to develop a model that can classify the population of records at huge. Fraud detection and credit risk applications are particularly well suited to such kind of analysis. This attitude commonly employs decision tree or neural network-based classification algorithms. Mainly data classification process involves learning and classification. In learning the training data are analysed by classification algorithm. During classification test data are used to estimate the accuracy of the classification rules. In the case of accuracy is acceptable the rules can be applied to the new data tuples. For a fraud detection type of application, this would comprise complete records of both fraudulent and valid activities determined on a record-by-record basis. To determine the set of parameters required for appropriate discrimination the classifier-training algorithm uses these pre-classified examples. Then the algorithm encodes these parameters into a model called a classifier.

In *Supervised learning* the computer is presented inputs and their desired outputs, and the goal is to learn a general rule that maps inputs to outputs. In *Unsupervised learning* there is no labels are given to the learning algorithm, leaving it on its own to find structure in its input. Unsupervised learning can be a goal in itself or a means towards an end. *Bayes' Theorem* finds the probability of an event occurring given the probability of another event that has already occurred. If B represents the dependent event and A represents the prior event, Bayes' theorem can be stated as follows.

Prob (B given A) = Prob(A and B)/Prob(A)

***RBF Network*** is a three-layer network, namely the input, the output and the hidden layer, where each hidden unit in a hidden layer implements a radial activated function. The main advantages of RBF's over feed-forward networks are its accuracy and shorter computational time. ***Sequential minimal optimization*** is an algorithm for solving the quadratic programming problem that arises during the training of support vector machines. ***Logistic function*** predicts the probability of an outcome that can only have two values. The prediction is based on the use of one or several predictors. **Decision table** is the rule which is used to make the decision. Decision tables may be defined in a variety of different ways.

## 2. DATASET DESCRIPTION

In this paper we have used Hungarian-14-heart disease dataset with 14 attributes. Using the WEKA tool to analyze the data and it provides the result with the accuracy of the algorithm.

**Table 2.1 Hangarian-14-heart-disease attributes**

| NO | Attribute Name | Description |
|---|---|---|
| 1 | Age | Age in Year |
| 2 | Sex | (value 1: Male; value 0 : Female) |
| 3 | Cp | value 1: typical type 1 angina, value 2:typical type angina, value 3: non-angina pain; value 4:asymptomatic) |
| 4 | Trestbps | mm Hg on admission to the hospital |
| 5 | Chol | Serum Cholesterol (mg/dl) |
| 6 | Fbs | Fasting Blood Sugar (value 1: > 120 mg/dl; value 0: < 120 mg/dl) |
| 7 | Restecg | Resting electrographic results (value 0: normal; value 1: 1 having ST-T wave abnormality; value 2: showing probable or definite left ventricular hypertrophy) |
| 8 | Thalach | (value 3: normal; value 6: fixed defect; value 7:reversible defect) |
| 9 | Exang | exercise induced angina (value 1: yes; value 0: no) |
| 10 | Oldpeak | ST depression induced by exercise relative to Rest |
| 11 | Slope | the slope of the peak exercise ST segment (value 1: Un sloping; value 2: flat; value 3: down sloping) |
| 12 | Ca | number of major vessels colored by floursopy (value 0– 3) |
| 13 | Thal | value 3: normal; value 6: fixed defect; value 7: reversible defect |
| 14 | Num | Diagnosis of heart disease. |

**Table 2.2 accuracy of different classification algorithms**

| Algorithm to be tested | Correctly classified instances | Incorrectly classified instances | Time taken to build model in seconds |
|---|---|---|---|
| Naïve Bayes | 86% | 14% | 0 sec |
| RBFNetwork | 86% | 14% | 0.17sec |
| SMO | 84% | 16% | 0.12sec |
| Logistic | 85% | 15% | 0.31sec |
| Decision Table | 79% | 21 % | 0.17sec |

The above table shows Hangarian-14-heart-disease dataset with 14 attributes are used to predict the heart disease in this study.

## 3. PERFOMANCE COMPARISON OF CLASSIFICATION ALGORITHMS

In the above table shows the accuracy with the correctly classified and incorrectly classified instances and time taken to build the model.

All the algorithm have same limitation a time when the dataset is very large. Using the split 66.0% train data test option is used to split the data.



**Figure 3.1 Screenshot for accuracy of Naïve bayes with 86%**



**Figure 3.2 Screenshot for accuracy of RBFNetwork function with 86%**

**Figure 3.3 Screen shot of report viewer (CVD)**

In figure 3.1 and 3.2 shows the naïve bayes and the RBFNetwork function results.Both the naïve bayes and the RBFNetwork function will provides the 86% of accuracy. The Naïve Bayes algorithm provides 0.7009 kappa statistics and the RBFNetwork function provides 0.6983 kappa statistics. The Naïve Bayes theorem provides the highest kappa statistics. Figure 3.3 shows the CVD report.

## 4.  CONCLUSION AND FUTURE WORK

Naïve Bayes, Logistic function, RBF Network, Decision Table, SMO function algorithms are used for testing and the testing results show that the Naïve Bayes algorithm outperforms with 86% accuracy in 0 seconds. A comparative study is applied to determine the most effective techniques that are capable for the detection of heart valve disease with a high accuracy. To implement the combination of classification techniques to improve the performance of the algorithms.

## REFERENCES

[1]  Anita Devi, Abhishek Misal, "A Survey on Classifiers Used in Heart Valve Disease Detection" ,International Journal of Advanced Research in Electrical, Electronics and Instrumentation Engineering Vol. 2, Issue 1, January 2013

[2]  K.R. Lakshmi , M.Veera Krishna and S.PremKumar,"Performance comparison of Data Mining Techniques for Predicting of Heart Disease Survivability", International Journal of Scientific and Research Publications, Volume 3, Issue 6, June 2013 1 ISSN 2250-3153

[3]  B.Venkatalakshmi, M.V Shivsankar ,"Heart Disease Diagnosis Using Predictive Data mining" International Journal of Innovative Research in Science, Engineering and Technology Volume 3, Special Issue 3, March 2014 International Conference on Innovations in Engineering and Technology (ICIET'14)

[4] Sona Baby, Ariya T.K, "A survey paper of data mining in medical diagnosis", international journal of research in IJRCCT computer and communication technology.

[5] K.Srinivas, Dr. G. Raghavendra Rao and and Dr. A. Govardhan, "survey on prediction of heart morbidity using data mining techniques",International Journal of Data Mining & Knowledge Management Process (IJDKP) Vol.1, No.3, May 2011

[6] Beant Kaur, Williamjeet Singh," Review on Heart Disease Prediction System using Data Mining Techniques" International Journal on Recent and Innovation Trends in Computing and Communication ISSN: 2321-8169 Volume: 2 Issue: 10

[7] Prof. Gondkar Mayura D.1, Prof. Pawar Suvarna E, " A Survey On Data Mining Techniques To Find Out Type Of Heart Attack", IOSR Journal of Computer Engineering (IOSR-JCE) e-ISSN: 2278-0661, p-ISSN: 2278-8727Volume 16, Issue 1, Ver. V (Jan. 2014), PP 01-05

[8] Vikas Chaurasia,etal,Carib.j.SciTech, "Early Prediction of Heart Diseases Using Data Mining Techniques" , ,2013,Vol.1,208-217

[9] Aqueel Ahmed, Shaikh Abdul Hannan, "Data Mining Techniques to Find Out Heart Diseases: An Overview", International Journal of Innovative Technology and Exploring Engineering (IJITEE) ISSN: 2278-3075, Volume-1, Issue-4, September 2012

[10] Dr.Hari Ganesh S and Gajenthiran M, "Comparative study of Data Mining Approaches for prediction Heart Diseases", IOSR Journal of Engineering (IOSRJEN) www.iosrjen.org ISSN (e): 2250-3021, ISSN (p): 2278-8719 Vol. 04, Issue 07 (July. 2014), ||V3|| PP 36-39

[11] Devendra Shah, Ketan Kadam, Shubham Shinde , Akash Doiphode, "Heart Disease Prediction: A Data Mining Aspect", International Journal of Engineering Technology, Management and Applied Sciences

[12] Zarna Parekh, Avaniba Parmar, "Survey Paper on Early Diagnosis Ofcardio-Vascular Disease Using Data Mining And Neural Network", IJSRD - International Journal for Scientific Research & Development| Vol. 2, Issue 10, 2014 ISSN (online): 2321-0613

[13] Jyoti Soni, Ujma Ansari, Dipesh Sharma and Sunita Soni, "Predictive Data Mining for Medical Diagnosis: An Overview of Heart Disease Prediction", International Journal of Computer Applications (0975 – 8887) Volume 17–No.8, March 2011

[14] Vikas Chaurasia and Saurabh Pal, "Data Mining Approach to Detect Heart Dieses", International Journal of Advanced Computer Science and Information Technology (IJACSIT) Vol. 2, No. 4, 2013, Page: 56-66, ISSN: 2296-1739.