

Design of a Predictive Model for Congenital Heart Disease Using Neural Networks

¹B. Vijayalakshmi and ²R. Sugumar

¹Department of Computer Science, Bharathiar University, Coimbatore, India

²Department of CSE, Velammal Institute of Technology, Chennai, India

Abstract: The proportion of deaths caused by heart disease is the highest in south India (25 per cent) and lowest - 12 per cent - in the central region of India. The prediction of heart disease survivability has been a challenging research problem for many researchers. Since the early dates of the related research, much advancement has been recorded in several related fields. Therefore, the main objective of this manuscript is to report on a research work where we took advantage of those available technological advancements to develop prediction models for heart disease survivability. We used three popular data mining algorithms CART (Classification and Regression Tree), ID3 (Iterative Dichotomized 3) and decision table (DT) extracted from a decision tree or rule-based classifier to develop the prediction models using a large dataset.

Key words: Back propagation • Data mining • Heart disease • Multilayer perceptron neural network • Neural Network

INTRODUCTION

Heart disease is the major cause of deaths. The World Health Organization (WHO) has estimated that 12 million deaths occur worldwide, every year due to the Heart diseases. In 2008, 17.3 million people died due to Heart Disease. Over 80% of deaths in world are because of Heart disease. WHO estimated by 2030, almost 23.6 million people will die due to Heart disease. Predication should be done to reduce risk of Heart disease. Diagnosis is usually based on signs, symptoms and physical examination of a patient. Almost all the doctors are predicting heart disease by learning and experience. The diagnosis of disease is a difficult and tedious task in medical field. Predicting Heart disease from various factors or symptoms is a multi-layered issue which may lead to false presumptions and unpredictable effects. Congenital heart disease is due to Cardiac malformations and are the most common birth defect in humans, affecting nearly 1% of all live births [1] and 1.35 million infants per year worldwide [2]. This number is probably even an underestimation, given that mild defects can be clinically unremarkable for decades. Furthermore, CHD are identified in about 10% of stillbirths and thus account for a substantial number of fetal deaths [3, 4]. The reported incidence of CHD varies substantially between different regions of the world, with the highest rate in Asia (0.93%) and lower rates in Europe (0.82%) and North America

(0.69%). The observed differences might be attributed to genetic, environmental as well as socioeconomic factors (e.g. parental consanguinity) and/or differences in healthcare and referral systems [2, 4]. Advances in cardiovascular medicine and surgery have lead to significant.

Review of Literature

Heart Diseases

Congenital Heart Disease: Congenital heart disease refers to a problem with the heart's structure and function due to abnormal heart development before birth. Congenital means present at birth [5].

Heart Muscle Disease (Cardiomyopathy):

Cardiomyopathy is a chronic disease of the heart muscle (myocardium), in which the muscle is abnormally enlarged, thickened and/or stiffened. The weakened heart muscle loses the ability to pump blood effectively, resulting in irregular heartbeats (arrhythmias) and possibly even heart failure [6].

Dilated Cardiomyopathy: The chambers of the heart are dilated (enlarged) because the heart muscle is weakened and cannot pump effectively. There are many causes, the most common being myocardial ischemia (not enough oxygen supplied to the heart muscle) due to coronary artery disease [6].

Restrictive Cardiomyopathy: Cardiomyopathy is an ongoing disease process that damages the muscle wall of the lower chambers of the heart. Restrictive cardiomyopathy is a form of cardiomyopathy in which the walls of the heart become rigid [5].

Hypertrophic Cardiomyopathy: Cardiomyopathy is an ongoing disease process that damages the muscle wall of the lower chambers of the heart. It is a form of cardiomyopathy in which the walls of the heart's chambers thicken abnormally. Hypertrophic cardiomyopathy is also referred as idiopathic hypertrophic sub aortic stenosis and asymmetrical septal hypertrophy [5].

Heart Disease Prediction: In [7], the neural network approach is used for analyzing the heart disease dataset. Applying feed forward neural network model and back propagation learning algorithm with variable learning rate and momentum the heart disease database are trained by the neural network. The input layer contains 13 neurons to represent 13 attributes. It consists of 4 class labels namely normal person, first stroke, second stroke and end of life. The output layer consists of two neurons to represent these four classes. Some of the neural networks are constructed with and without hidden layer that is single and multilayer networks are trained. The dataset was collected from Cleveland database [8]. This dataset classifies the person into normal and abnormal person based on heart diseases. The dataset consists of 414 instances, 13 attributes and a class attribute. Both test and training data are used for performance analysis. In a trained network, the test data is given as the input. With the adjusted weights, the output of the net is calculated. From the experimental results, the author concluded that efficiency of the classification process is increased by applying parallel approach which is adopted in the training phase. In future this work will be enhanced by applying genetic algorithm using neural networks.

In [9], the data mining classification techniques namely RIPPER classifier, decision tree, Artificial Neural Networks and Support vector machine are used for predicting cardiovascular heart disease. The performance factors used for comparing these techniques are sensitivity, accuracy, specificity, error rate, true positive rate and false positive rate. To measure the unbiased estimate of prediction models the author used 10 fold cross validation method. This model was developed by using data mining classification tool weka version 3.6. It contains 14 attributes and 303 instances. From an experiment, the results are compared. Error rates for

RIPPER, Artificial Neural Networks, Support vector machine and Decision tree are 0.2756, 0.2248, 0.1588 and 0.2755 respectively. The accuracy of RIPPER, Artificial Neural Networks, Support Vector Machine and Decision tree are 81.08%, 80.06%, 84.12% and 79.05% respectively. When compared to four classification models, the Support Vector Machine has given least error rate and highest accuracy. The author concluded that the Support Vector Machine is the best technique for predicting the cardiovascular disease. In future, in order to improve the efficiency of the classification techniques by creating meta models.

In [10], the heart attack symptoms are predicted using biomedical data mining techniques. The author used data classification which is based on supervised machine learning algorithms. For data classification the Tanagra tool is used. Using entropy based cross validations and partitioned techniques, the data is evaluated and the results are compared. The algorithms used in these techniques are K -nearest neighbors, K-means and Mean Clustering Algorithm (EMC) is the extension of the K-mean algorithm for clustering process which reduces the number of iterations. As a result the author analyzed that the mean clustering algorithm performs well when compared to other algorithms. To run the data the time taken is very fast and it gives the result of accuracy about 82.90%. Further this work will be enhanced by applying unsupervised machine learning algorithm.

In [11] the fuzzy expert system is designed for heart disease diagnosis with reduced number of attributes. The author finds that how genetic algorithm and fuzzy logic combine together for efficient and cost effective diagnosis of heart disease. The genetic algorithm and two models of fuzzy system Mamdani and Takagi-sugeno were used to find the cost. The dataset were taken in Cleveland clinic foundation dataset [6]. The input field is a set of all the selected features and the output of the system is to get a value '1' or '0' that indicates the presence or absence of disease. It is further enhanced by using art classifiers like Decision tree, Naïve bayes, Classification via clustering and SVM classifier.

Algorithms Used for Model Building and Performance Measures: An attempt was made to construct prediction model using Decision Tree, Neural Networks. After constructing the models, performance of each models were evaluated and also their performances were compared to each other. In this section the algorithms used to build the models and matrices used for performance measure and comparison are discussed in detail.

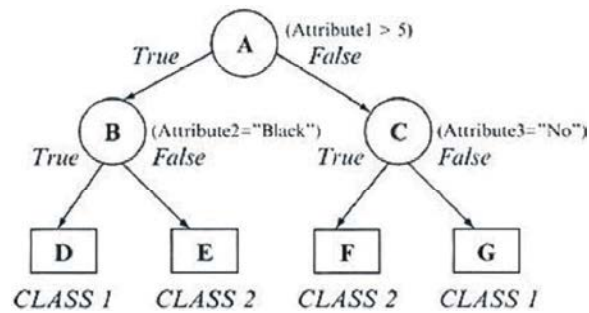


Fig. 3.1: A simple Decision Tree

Decision Trees: Han and Kamber (2006) defined decision tree as a flowchart like tree structure, where each internal node (non-leaf node) denotes a test on an attribute, each branch represents an outcome of the test and each leaf node (or terminal node) holds a class label. The topmost node in a tree is the root node.

As the construction of decision tree classifiers does not require any domain knowledge or parameter setting and therefore is appropriate for exploratory knowledge discovery they have become popular. Decision trees can handle high dimensional data. Their representation of acquired knowledge in tree form. Their representation of acquired knowledge in tree form is intuitive and generally easy to assimilate by humans. The learning and classification steps of decision tree induction are simple and fast. In general, decision tree classifiers have good accuracy. However, successful use may depend on the data at hand.

Decision tree induction algorithms have been used for classification in many application areas, such as medicine, manufacturing and production, financial analysis, astronomy and molecular biology (Han and Kamber, 2006).

3.1.1 J48 Classifier Algorithm: During the late 1970s and early 1980s, J. Ross Quinlan, a researcher in machine learning, developed a decision tree algorithm known as ID3 (Iterative Dichotomiser). This work expanded on earlier work on concept learning systems, Quinlan later presented C4.5 (a successor of ID3), which became a benchmark to which newer supervised learning algorithms are often compared. C4.5 algorithm is an improvement of IDE3 algorithm. It is based on Hunt's algorithm and also like IDE3, it is serially implemented. Pruning takes place in C4.5 by replacing the internal node with a leaf node thereby reducing the error rate. Unlike IDE3, C4.5 accepts both continuous and categorical attributes in building the decision tree. It has an enhanced method of tree pruning that reduces misclassification errors due noise or

too-much details in the training data set. Like IDE3 the data is sorted at every node of the tree in order to determine the best splitting attribute. It uses gain ratio impurity method to evaluate the splitting attribute (Quinlan, 1993). The C4.5 algorithm uses the concept of information gain or entropy reduction to select the optimal split. Suppose that we have a variable X whose k possible values have probabilities P_1, P_2, \dots, P_k the smallest number of bits, on average per symbol, needed to transmit a stream of symbols representing the values of X observed called the entropy of X and is defined as:

$$H(x) = \sum_j P_j \log_2 (P_j) \quad (\text{Formula 3.1})$$

For an event with probability p , the average amount of information in bits required to transmit the result is $-\log_2(P)$. For example, the result of a fair coin toss, with probability 0.5, can be transmitted using $\log_2(0.5) = 1$ bit, which is a zero or 1, depending on the result of the toss. For 41 variables with several outcomes, we simply use a weighted sum of the $P(j)$'s, with weights equal to the outcome probabilities, resulting in the formula 3.1. C4.5 uses this concept of entropy as follows. Suppose that we have a candidate split S , which partitions the training dataset T into several subsets, T_1, T_2, \dots, T_k .

J48 examines the normalized information gain (difference in entropy) that results from choosing an attribute for splitting the data. To make the decision, the attribute with the highest normalized information gain is used. Then the algorithm recurs on the smaller subsets. The splitting procedure stops if all instances in a subset belong to the same class. Then a leaf node is created in the decision tree telling to choose that class. But it can also happen that none of the features give any information gain. In this case J48 creates a decision node higher up in the tree using the expected value of the class. J48 can handle both continuous and discrete attributes, training data with missing attribute values and attributes with differing costs. Further it provides an option for pruning trees after creation.

Neural Networks: Neural networks are network structures consisting of a number of nodes connected through directional links. Each node represents a processing unit and the links between nodes specify the causal relationship between connected nodes (Kantardzic, 2003).

According to Two Crows Corporation (2005), a neural network (Figure 3.2) starts with an input layer, where each node corresponds to a predictor variable.

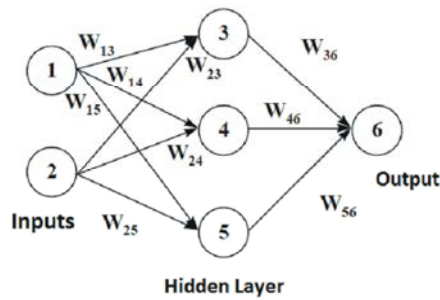


Fig. 3.2: A Neural Network Architecture

These input nodes are connected to a number of nodes in a hidden layer. Each input node is connected to every node in the hidden layer. The nodes in the hidden layer may be connected to nodes in another hidden layer, or to an output layer. The output layer consists of one or more response variables.

After the input layer, each node takes in a set of inputs, multiplies them by a connection weight W_{xy} (e.g., the weight from node 1 to 3 is W_{13} - Figure 3.2), adds them together, applies a function (called the activation or squashing function) to them and passes the output to the node(s) in the next layer. For example, the value passed to node 4 is:

Activation function applied to $([W_{14} * \text{value of node 1}] + [W_{24} * \text{value of node 2}])$

Each node may be viewed as a predictor variable (nodes 1 and 2 in this example) or as a combination of predictor variables (nodes 3 through 6). Node 6 is a non-linear combination of the values of nodes 1 and 2, because of the activation function on the summed values at the hidden nodes. In fact, if there is a linear activation function but no hidden layer, neural nets are equivalent to a linear regression; and with certain non-linear activation functions, neural nets are equivalent to logistic regression.

The connection weights (W 's) are the unknown parameters which are estimated by a training method. Originally, the most common training method was backpropagation; newer methods include conjugate gradient, quasi-Newton, Levenberg-Marquardt and genetic algorithms. Each training method has a set of parameters that control various aspects of training such as avoiding local optima or adjusting the speed of conversion (Two Crows Corporation, 2005).

The architecture (or topology) of a neural network is the number of nodes and hidden layers and how they are connected. In designing a neural network, either the user

or the software must choose the number of hidden nodes and hidden layers, the activation function and limits on the weights. While there are some general guidelines, you may have to experiment with these parameters.

Multilayer Perception: As Giudici (2003) stated that the multilayer perceptron is a neural network which is the most used architecture for predictive data mining. It is a feedforward network with possibly several hidden layers, one input layer and one output layer, totally interconnected. It can be considered as a highly nonlinear generalization of the linear regression model when the output variables are quantitative, or of the logistic regression model when the output variables are qualitative. The network is feedforward if the processing propagates from the input side to the output side unanimously, without any loops or feedbacks. In a layered representation of the feedforward neural network, there are no links between nodes in the same layer; outputs of nodes in a specific layer are always connected as inputs to nodes in succeeding layers. This representation is preferred because of its modularity, i.e., nodes in the same layer have the same functionality or generate the same level of abstraction about input vectors (Kantardzic, 2003).

Proposed Method: In the present study a most frequently used Back Propagation Neural Network Model is used to perform the Congenital Heart Disease Diagnosis classification based on the signs, symptoms and physical evaluation of a patient which are given in the below table.

Backpropagation Neural Network: Back propagation Neural Network is a multilayer Feed Forward Neural Network Model [12] which contains one input layer, one output layer and one or more hidden layers. As the name implies the input layer receives signals from the external nodes and transmits these signals to other layers without performing any computations at that layer. The output layer receives the signals from an input layer through a weighted connection links, performs computations at that layer and produces output of the network. The hidden layer of a network receives signals from an input layer through a weighted connection links, performs the computation and transmits these results as signals to the output layer through a weighted connection links w_{11} , w_{12} , \dots w_{pm} . The hidden layer of a network is neither an input layer nor an output layer instead it acts as input or output layer based on situation. The architecture of a Back propagation Neural Network Model [13] is shown in

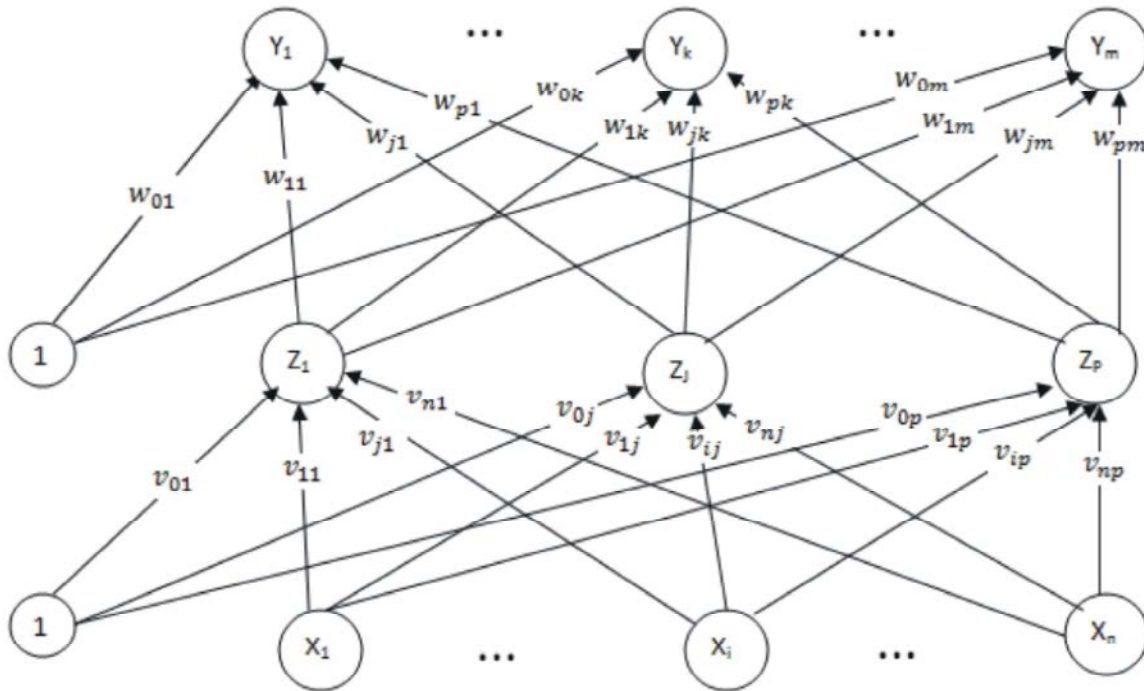


Fig. 4.1: Architecture of a Back propagation Neural Network Model

Figure 4.1. Neural Network is trained by using Generalized Delta Rule also called as Back propagation Rule [14]. Training a Back propagation Network involves three stages. First stage is the Feed Forward phase of input training pattern, in which each of the hidden nodes Z_j receives the input signals x_1, x_2, \dots, x_n from input nodes X_i through connection weights $v_{11}, v_{12}, \dots, v_{np}$, computes the net input as the sum of the products of the input signal and weights i.e $Z_i n_j = v_{0j} + \sum x_i v_{ij}$ applies an activation function to produce the response and finally sends these signals to the output nodes Y_k . Similar to hidden nodes, each of the output nodes Y_k receives the signals z_1, z_2, \dots, z_p from the hidden nodes Z_j through the connection weights $w_{11}, w_{12}, \dots, w_{pn}$, calculates the net input as $y_i n_k = w_{0k} + \sum z_j w_{jk}$, applies an activation function to produce the output of the network. An activation function that is used at both hidden and the output layer is a sigmoid function i.e $f(x) = 1/(1 + \exp(-x))$. Once the response of the net is calculated then each output node compares its response with the target values t_k to determine the associated error. Based on these error the factor δ_k is determined as $\delta_k = (t_k - y_k) * f'(y_i n_k)$. δ_k is used to distribute the errors at Y_k back to all the hidden nodes that are connected to Y_k . Similarly factor δ_j is computed as $\delta_j = \delta_i n_j * f'(z_i n_j)$ for each hidden node Z_j to propagate errors back to input layer. Once all the δ factors are determined and propagated to appropriate layers, then the weights

are adjusted simultaneously. The weight adjustment from hidden node to output node is based on the δ_k as $\Delta w_{jk} = \alpha \delta_k z_j$ and the weights between the hidden node to input node is based on the factor δ_j as $\Delta v_{ij} = \alpha \delta_j x_i$. Therefore the new weights between the input nodes and hidden nodes are $v_{ij}(\text{new}) = v_{ij}(\text{old}) + \Delta v_{ij}$, $v_{0j}(\text{new}) = v_{0j}(\text{old}) + \Delta v_{0j}$ and the new weights between the hidden nodes and output nodes are $w_{jk}(\text{new}) = w_{jk}(\text{old}) + \Delta w_{jk}$, $w_{0k}(\text{new}) = w_{0k}(\text{old}) + \Delta w_{0k}$. This process is continued until the stopping condition. The stopping condition may be the minimization of mean squared error (MSE) value or the number of epochs it has reached. The least mean squared error (MSE) can be computed as $\text{MSE} = 0.5 \sum (t_k - y_k)^2$.

Since the Neural Network solutions will not depends on algorithmic solution instead it depends on examples of the previous cases it gives more accurate results than the human diagnosis.

Confusion Matrix: In classification problems, the primary source of performance measurements is a confusion matrix (Coincidence matrix, classification matrix or a contingency table). Given m classes, a confusion matrix is a table of at least size m by m (Olson and Delen, 2008).

If the instance is positive and it is classified as positive, it is counted as a true positive (TP); if it is classified as negative, it is counted as a false negative (FN). If the instance is negative and it is classified as

		Predicted class	
		C ₁	C ₂
Actual class	C ₁	true positives	false negatives
	C ₂	false positives	true negatives

Fig. 4.2: A simple Confusion Matrix

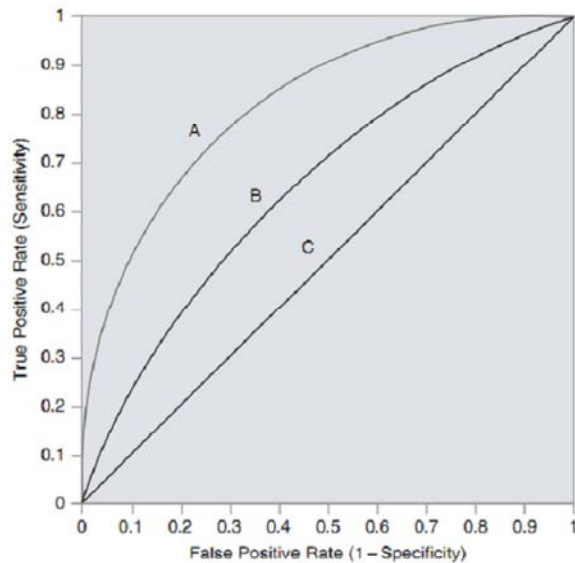


Fig. 4.3: A Sample ROC curves

negative, it is counted as a true negative (TN); if it is classified as positive, it is counted as a false positive (FP). These terms are useful when analyzing a classifier's ability.

Figure 4.2 shows a confusion matrix for a two-class classification problem. The numbers along the diagonal from upper-left to lower-right represent the correct decisions made and the numbers outside this diagonal represent the errors. The equations of most commonly used metrics that can be calculated from the coincidence matrix are discussed below.

The overall accuracy of a classifier is estimated by dividing the total correctly classified positives and negatives by the total number of samples. Other performance measures, such as recall (sensitivity), specificity and F-measure are also used for calculating other aggregated performance measures (e.g., area under the ROC curves).

The accuracy of a classifier on a given test set is the percentage of test set tuples that are correctly classified by the classifier.

In order to plot an ROC curve for a given classification model, true positive (TP) rate is plotted on

the Y axis and false positive (FP) rate is plotted on the X axis. As Han and Kamber (2006) described the process of drawing Roc curve we start at the bottom left-hand corner (where the true positive rate and false-positive rate are both 0), we check the actual class label of the tuple at the top of the list. If we have a true positive (that is, a positive tuple that was correctly classified), then on the ROC curve, we move up and plot a point. If, instead, the tuple really belongs to the 'no' class, we have a false positive. On the ROC curve, we move right and plot a point. This process is repeated for each of the test tuples, each time moving up on the curve for a true positive or toward the right for a false positive.

Area under the Curve (AUC) is a portion of the area of the unit square and its value will always be between 0 and 1.0. A perfect accuracy gets a value of 1.0. The diagonal line $y = x$ represents the strategy of randomly guessing a class. For example, if a classifier randomly guesses the positive class half the time (much like flipping a coin), it can be expected to get half the positives and half the negatives correct; this yields the point (0.5; 0.5) in ROC space, which in turn translates into area under the ROC curve value of 0.5. No classifier that has any classification power should have an AUC less than 0.5. For example, In Figure 3.4 classification performances of three classifiers (A, B and C) are shown in a single ROC graph. Since the AUC is the commonly used metric for performance comparison of prediction models, one can easily tell that the best performing classifier (out of the three that is being compared to each other) is A, followed by B. The classifier C is not showing any predictive power; staying at the same level as random chance.

Experimentation and Results: The experiment is carried out on a available database for heart disease. The dataset contains total 573 records. The dataset is divided into 2 sets training (303 records) and testing set (270 records). A data mining tool Weka 3.6.6 is used for experiment. Parameters used for experiment are listed below.

Patient ID: Patient Identification number.

Diagnosis: Value 1: $\leq 50\%$ (no heart disease)
Value 0: $> 50\%$ (has heart disease)

The remaining parameters are listed out in tabular form as below.

Table 1: Switzerland Data Set Clinic Foundation

Name	Type	Description
Age	Continuous	Age in years
Sex	Discrete	1 = male 0 = female
cp	Discrete	Chest pain type: 1 = typical angina 2 = atypical angina 3 = non-anginal pa 4 = asymptomatic
Trestbps	Continuous	Resting blood pressure (in mm Hg)
Chol	Continuous	S erum cholesterol in mg/dl
Fbs	Discrete	Fasting blood sugar > 120 mg/dl: 1 = true 0 = false
Restecg	Discrete	Resting electrocardiographic results: 0 = normal 1 = having ST-T wave abnormality 2 = showing probable or definite left ventricular hypertrophy by Estes' criteria
Thalach	Continuous	Maximum heart rate achieved
Exang	Discrete	Exercise induced angina: 1 = yes 0 = no
Slope	Discrete	The slope of the peak exercise segment: 1 = up sloping 2 = flat 3 = down sloping
Diagnosis	Discrete	Diagnosis classes: 0 = healthy 1 = possible heart disease

Table 2: Shows the experimental result. We have carried out some experiments in order to evaluate the performance and usefulness of different classification algorithms for predicting heart patients

Evaluation Criteria	Classifiers			
	CART	ID3	Decision Table (DT)	Predictive algorithm
Timing to build model (in Sec)	0.23	0.02	0.03	0.02
Correctly classified instances	253	221	250	252
Incorrectly classified instances	50	75	53	49
Accuracy (%)	83.49%	72.93%	82.5%	96%

CONCLUSION

An average of 241 instances out of total 303 instances is found to be correctly classified with highest score of 253 instances compared to 221 instances, which is the lowest score. The total time required to build the model is also a crucial parameter in comparing the classification algorithm.

In this simple experiment, from Table 2, we can say that a ID3 and DT requires the shortest time which is around 0.02 and 0.03 seconds consecutive with compared

to CART which requires the longest model building time which is around 0.23 seconds. The empirical results show that we can produce short but accurate prediction list for the heart patients by applying the predictive models to the records of incoming new patients. This study will also work to identify those patients who needed special attention.

REFERENCES

- Hoffman, J.I.E. and S. Kaplan, 2002. The incidence of congenital heart disease, J. Am. Coll. Cardiol., 39: 1890-1900.
- Van Der Linde, D., 2011. Birth prevalence of congenital heart disease worldwide: a systematic review and meta-analysis, J. Am. Coll. Cardiol., 58: 2241-2247.
- Hoffman, J.I., 1995. Incidence of congenital heart disease: II, Prenatal Incidence. Pediatr Cordiol, 16: 155-165.
- Fahed, A.C., B.D. Gelb, J.G. Seidman and C.E. Seidman, 2013. Genetics of congenital heart disease: the glass half empty. Circ. Res., 112: 707-720.
- Vijayarani, S. and S. Sudha, 2012. A Study of Heart Disease Prediction in Data Mining, IRACST - International Journal of Computer Science and Information Technology & Security (IJCSITS, 2(5): 2249-9555.
- www.webmd.com/heart-disease/guide/heart-diseasesymptoms-types.
- Rani, K. Usha, 2011. Analysis of Heart Diseases Dataset Using Neural Network Approach, (IJDKP, 1(5).
- Cleveland and heart disease dataset, sci2s.ugr.es/keel/dataset.php?cod=57.
- Karabulut Esra Mahsereci and Turgay İbrikçi, 2011. Effective Diagnosis of Coronary Artery Disease Using the Rotation Forest Ensemble Method, June 2011, Accepted: 30 August 2011 / Published online: 13 September 2011, Springer Science Business Media, LLC.
- Jaya Rama krishniah V.V., D.V. Chandra Sekar, K. Ramch and H. Rao, 2012. Predicting the Heart Attack Symptoms using Biomedical Data Mining Techniques, ISSN – 2278-1080, The International Journal of Computer Science & Applications (TIJCSA), 1(3).
- Ephzibah, E.P., 2011. A Hybrid Genetic-Fuzzy Expert System for Effective Heart Disease Diagnosis, D.C. Wyld et al. (Eds.), ACITY 2011, CCIS, Springer-Verlag Berlin Heidelberg 2011, 198: 115-121.

12. Rimer Michael, E., 2007. Improving Neural Network Classification Training, Ph.D dissertation, submitted to Dept. of Computer Science, Brigham young University.
13. Fausett Laurene, Fundamental of Neural Networks, 3rd Edition, Pearson Education.
14. Rajasekaran, S. and G.A. Vijayalakshmi Pai, Neural Networks, Fuzzy logic and Genetic Algorithms synthesis and Applications, Eastern Economy Edition.