RESEARCH ARTICLE                                                                      OPEN ACCESS

# Hybridization of Improved K-Means and Artificial Neural Network for Heart Disease Prediction

Shubhada Bhalerao [1], Dr. Baisa Gunjal [2]
ME Student, Associate Professor & HOD
Department of Computer Engineering
Department of Information Technology
Amrutvahini COE, Sangamner,
Maharashtra - India.

## ABSTRACT

The healthcare field is composed of a variety of data, but the task of getting efficient knowledge and hidden information from those data is very difficult and time-consuming by the traditional method. Data mining technique is supportive for extracting valuable knowledge as well as hidden information from heath care system. Recently, abundant software's, tools and various algorithms have been studied by the researchers for developing efficient medical decision support systems. Moreover, latest algorithms and innovative tools are continuously added day by day for making existing work more efficient. Last few years, heart disease is the major reason of death all over the world. Detection of heart disease is one of the vital issues and many researchers are developing intelligent medical decision support systems to get better the ability of the physicians. In heart disease diagnosis and treatment, single data mining techniques are giving the satisfactory level of accuracy. But to increase the accuracy level we can use hybrid data mining techniques. In this paper, Improved K-means and ANN are combined to achieve an efficient result in heart disease diagnosis.

*Keywords: -* Classification, Data mining, Heart disease prediction, Improved K-mean, ANN, Hybridization.

## I. INTRODUCTION

Medical informatics is a department where there is an activity of designing, developing, adoption and use of IT-based innovations in health care field. In short it is a junction of the medical field and information technology which provides considerable improvements in a perfection of healthcare activities and its effectiveness. The health-care field is composed of a variety of data, but the task of getting efficient knowledge and gaining hidden information from those data is very difficult and time-consuming by the traditional method. So, to maintain a quality of care and effectiveness as well as to gain knowledge from hidden information there a need of applying computer science technique on available information. Data mining holds great potential to explore the hidden patterns in huge information that can be used for clinical diagnosis. Data mining allow health systems to use data systematically and do the analysis for identifying inefficiencies, best practices that improve care and reduce costs.

Today's lifestyle invites the various kind of disease in the human body. In Spite of all disease, a large number of populations are suffering from multiple kinds of heart disease and count of a patient dying due to heart disease is increasing day by day. So there is need of tool or technique which helps in early detection of heart disease with higher accuracy. It is a challenging task to analyze heart disease just with the help of patient's report mostly; doctors take the decision on their experience and knowledge. But unfortunately, due to the complex processes and different symptoms and pathological tests, the correct diagnosis of heart diseases is a difficult task which causes a delay in the proper treatment. Therefore, there is a requirement to develop the prediction systems for heart disease which can help the medical experts in the early and accurate diagnosis of heart disease. This thesis proposes method planned for developing new hybrid algorithms which are designed for providing automatic computerize analysis and decision support system for heart disease diagnosis. Fig.1 shows a complete scenario of identify the Heart disease patients [2]. Firstly identify the disease database for diagnosis. After that find out the different data mining techniques for improvement of efficiency and accuracy of the prediction. To reduce the cost and for getting best outcome hybrid data mining technique can be used.
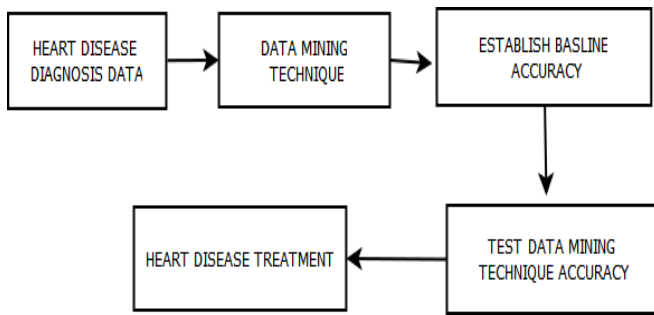
Fig. 1 Complete Scenario Of Identify The Heart Disease Patients.

## II. RELATED WORK

Carlos Ordonez [3] has worked on a prediction of heart disease with the help of Association rules. The author used a simple mapping algorithm. This algorithm continuously treats attributes as numerical or categorical. In this medical records are converted into a transaction format. To mine the constrained association rules an improved algorithm is used. A mapping table is prepared and attribute values are mapped to items. The decision tree is used for mining data because they automatically Split numerical values [3]. The split point chosen by the Decision tree are of little use only. Clustering is used to get a global understanding of data

Chaitrali S. Dangare [4] has proposed heart disease prediction system based on neural network. The HDPS system predicts the probability of having a Heart disease. For prediction, the system uses gender, blood pressure, cholesterol, sugar level, age like 13 clinical parameters. For better accuracy two more parameters are added i.e. Obesity and smoking. From the results, it has been noticed that neural network predicts heart disease with nearly 100% accuracy.

M.Akhil Jabbar, B.L.Deekshatulu, Priti Chandra [5] put forward a new algorithm in which classification is done using KNN with Genetic Algorithm. To provide optimal solution genetic algorithms perform a global search on complex large and multimodal dataset. From the results the author made conclusion that hybridizing GA with KNN performs well and give great accuracy in heart disease prediction.

Ankita Dewan [1] proposed an efficient algorithm for heart disease prediction in this proposed work the author have made hybridization of genetic algorithm and back propagation. They have observed that neural network is best among all the classification techniques for a non-linear data. BP algorithm is the best classifier of Artificial Neural Network which is a common method of training. In this, the primary system output is compared to the expected output, and the system is adjusted until the difference between the two is minimized. But it has a drawback of being stuck in local minima.

Milan Kumari et al. [6] proposed a system for prediction of cardiovascular disease dataset with the help of different data mining algorithms, like Support Vector Machine, Artificial neural networks (ANNs), Decision Tree, and RIPPER classifier. The author did performance analysis of these algorithms through several statistical analysis factors such as sensitivity, specificity, accuracy, error rate, True Positive Rate and False Positive Rate. The accuracy percentages are 81.08%, 80.06%, 79.05% and 84.12% for RIPPER, ANN, Decision Tree and SVM respectively. While the results of error rate are 2.756, 0.2755, 0.2248 and 0.1588 for RIPPER, Decision Tree, ANN and SVM respectively. The author concluded that out of these four classification models SVM is giving better prediction result for cardiovascular disease with least error rate and highest accuracy.

Rajkumar, A. and G.S. Reena et al. [7] studied the diagnosis of heart disease dataset with Tanagra tool .It is used to compare the performance accuracy of data mining algorithms for prediction of heart disease. In this author have studied the various algorithms such as Naive Bayes, k-nn and Decision tree. In their research, Naive Bayes algorithm is giving best compact time for processing dataset and showed better performance in accuracy prediction. The overall system gave 52.33% accuracy.

Gudadhe et al. developed a system by using Support Vector Machine and multilayer Perceptron neural network architecture to classify heart disease. The database of heart disease is divided into two classes which show a presence of heart disease or absence of heart disease with an accuracy of 80.41%, this division is done using support vector machine. Whereas the artificial neural network classifies the heart disease data into 5 classes with an accuracy of 97.5% [8].

M. Durairaj, V. Ranjani has observed a detection of heart disease using classification technique which gives 60% accuracy [9].

Javad Kojuri et al. (2015) was enrolled a total of 935 cardiac patients with chest pain and no diagnostic electrocardiogram (ECG). The study is done by using two types of data: nominal (clinical data) and quantitative (ECG findings). To classify data into two groups different neural networks are used that are radial basis function (RBF) and multi-layer perceptron (MLP). Result describe that the RBF neural network had an accuracy of 83% with ECG findings and an accuracy of 78% with clinical features [10].

Srinivas et al. [11] studied the diagnosis of cardiovascular disease (CVD) with the help of Neural Network, Decision Trees, Naïve Bayes data mining technique. The system gives the accuracy of 82.5%.

A decision support system for analysis of Congenital Heart Disease has been proposed by Vanisree K et al. [12]. The core

of the proposed system is based on Back propagation Neural Network (multi-layered Feed Forward Neural Network). The attribute set used in this work are the signs, symptoms and the results of physical evaluation of a patient. The proposed system achieved an accuracy of 90%.

Preeti Gupta et al. (2014) proposed the prediction of heart disease diagnosis with UCI dataset for enhancement of the accuracy of the results by using Artificial Neural Network optimized with Genetic Algorithm. In this study neural network was optimized with Genetic Algorithm for the accuracy enhancement. The MATLAB GUI feature is used for designing purpose. The proposed method achieved an accuracy of 97.83 [13].

Resul Das and Ibrahim Turkoglu [14] used a methodology of SAS software for diagnosing of heart disease. He used a neural network developed a system. In this creation of new models is done by combining the posterior probabilities or predicted values from multiple predecessor models. The 89.01% classification accuracy was obtained from the experiments done on Cleveland heart disease database. The specificity and sensitivity value obtained are 95.91% and 80.95% respectively.

Alizadehsani et al diagnosis the CAD via stenosis of LAD vessel by using C4.5 Classifier data mining technique he has also studied the diagnosis of CAD via stenosis of LCX vessel by using the KNN. The accuracy achieved by KNN and C4.5 is 61.39%, 74.20 % respectively [15].

## III. PROPOSED SYSTEM

The system will help to predict heart disease using machine learning hybrid approach. The User first goes from the process of authentication and after that, the user will enter the details like age, gender, blood pressure, etc. From the report .After providing the data the user will receive output like whether he requires diagnosis or not. All the system is developed by doing hybridization of Improved K-mean and ANN algorithm. The whole system is a combination of Training phase and Prediction phase .10 Parameters used as an input from dataset namely Age, Sex, Blood pressure, chest pain, blood sugar etc. The user will load the data set first then parsing is done after all this clustering is carried out by using an Improved K-Mean algorithm. All this clustered data is provided as input to ANN algorithm which is a combination of back propagation and feeds forward to prepare trained data set. In Prediction phase, the user will again load the data set entry and the system will detect the heart disease with the help of trained data .The accuracy of system is calculated by using a result of prediction phase.

## IV. PROBLEM DEFINITION

The problem is to determine as: - Medical Misdiagnoses are a serious risk to our healthcare profession. Due to this, people will develop a fear to visit a hospital for treatment. The system can put an end to medical misdiagnosis by developing a system where in the heart disease can be predicted using machine learning hybrid approach.
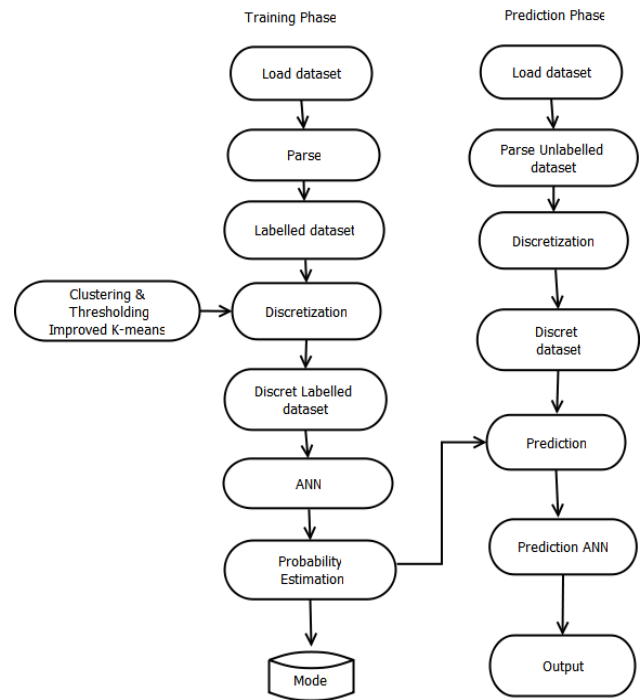
## V. SYSTEM DESIGN

### A. System Workflow



Fig 2: System Workflow

Fig 2 shows workflow of the whole system which consists of Improved K-Means and ANN.

### B. System Modules

1) **Client**

a) **Input Features:** Client first registers and login than enters his features like age, gender, blood pressure, etc.

b) **Predict Output:** Receiving output with suggested medicine name and displayed on a screen.

2) **Server**

a) **ANN:** Receives normalize data and apply ANN for the training and save this train data.

**b) Prediction:** Receive current input features and apply prediction.

3) **Admin**

**a) Load data set:** Admin load dataset for the training purpose.

**b) Clustering:** Applying k-means algorithm for the input features.

**c) Export to Server:** Sending normalizes data to the server for the training.

## C. Proposed Algorithm

1. Start

2. Load Data set

3. Perform Parsing

4. Accept label data set

5. Apply IMPROVED K-Means algorithm for Clustering

6. Perform Discretization

7. Apply ANN for Training

8. Repeat the steps from 2 to 6 for prediction

9. Prediction of disease

10. Show result

11. Stop

# VI. USER INTERFACE FRAMEWORK

Fig.3 and Fig.4 show the user interface of the system in which Fig. 4 shows the testing multiple entries with accuracy percentage.
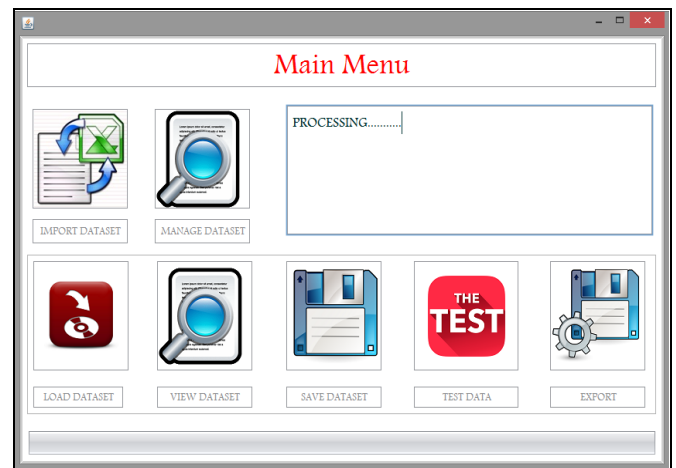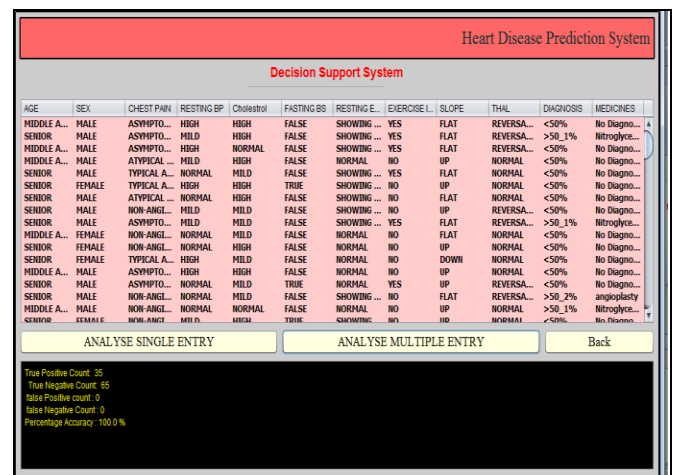


Fig.3 Main Menu



Fig.4 Test Dataset

# VII. EXPERIMENTAL RESULT

## A. Data Source

The experimental results of the heart disease prediction system using Improved K-means and ANN Approach are explained in this section. In this work the clinical data related to Heart diseases are considered.Cleveland dataset is used which is taken from Data mining repository of the University of California, Irvine (UCI) [16]. Cleveland dataset shows the classification of a person into the normal and abnormal person with respect to the presence of heart diseases. The dataset consists of total 303 records and it is divided into 2 sets training and testing set. The dataset consists of 13 attributes (inputs), and 4 classes (outputs). Out of that 13 attribute (input), this system is developed using 10 attributes Tables I illustrate the

representation of the attributes [16] and Fig.5 shows sample dataset.

TABLE I

CLEVELAND HEART DISEASE DATASET ATTRIBUTES

| Attribute | Description | Range |
|---|---|---|
| **Age** | Age in years | Continuous |
| **Sex** | (1 = male; 0 = female) | 0,1 |
| **Cp** | chest pain type<br>  -- Value 1: typical angina<br>  -- Value 2: atypical angina<br>  -- Value 3: non-anginal pain<br>  -- Value 4: asymptomatic | 1,2,3,4 |
| **trestbps** | resting blood pressure (in mm Hg on admission to the hospital) | Continuous |
| **Chol** | serum cholesterol in mg/dl | Continuous |
| **Fbs** | (Fasting blood sugar .120mg/dl )<br>(1=true; 0=false) | 0,1 |
| **Restecg** | Resting electrocardiographic results:<br> -- Value 0: normal<br> -- Value 1: having ST-T wave abnormality (T wave inversions and/or ST elevation or depression of > 0.05 mV)<br>-- Value 2: showing probable or definite left ventricular hypertrophy by Estes' criteria. | 0,1,2 |
| **Exang** | Exercise induced angina(1=yes;0=no) | 0,1 |
| **Slope** | the slope of the peak exercise ST segment –Value 1: up sloping -- Value 2: flat | 0,1,2 |
| | --Value3:downslopig | |
| **Thal** | 3 = normal; 6 = fixed defect; 7 = reversible defect | 3,6,7 |
| **Num** | Class(0=healthy, 1/2/3/4=have heart disease) | 0,1,2,3,4 |

| Age | sex | cp | trestbps | chol | fbs | restecg | thalach | exang | oldpeak | slope | ca | thal | num |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 63 | 1 | 1 | 145 | 233 | 1 | 2 | 150 | 0 | 2.3 | 3 | 0 | 6 | 0 |
| 67 | 1 | 4 | 160 | 286 | 0 | 2 | 108 | 1 | 1.5 | 2 | 3 | 3 | 2 |
| 67 | 1 | 4 | 120 | 229 | 0 | 2 | 129 | 1 | 2.6 | 2 | 2 | 7 | 1 |
| 37 | 1 | 3 | 130 | 250 | 0 | 2 | 187 | 0 | 3.5 | 3 | 0 | 3 | 0 |
| 41 | 0 | 2 | 130 | 204 | 0 | 2 | 172 | 0 | 1.4 | 1 | 0 | 3 | 0 |
| 56 | 1 | 2 | 120 | 236 | 0 | 2 | 178 | 0 | 0.8 | 1 | 0 | 3 | 0 |
| 62 | 0 | 4 | 140 | 268 | 0 | 2 | 160 | 0 | 3.6 | 3 | 2 | 3 | 3 |
| 57 | 0 | 4 | 120 | 354 | 0 | 0 | 163 | 1 | 0.6 | 1 | 0 | 3 | 0 |
| 63 | 1 | 4 | 130 | 254 | 0 | 2 | 147 | 0 | 1.4 | 2 | 1 | 7 | 2 |
| 53 | 1 | 4 | 140 | 203 | 1 | 2 | 155 | 1 | 3.1 | 3 | 0 | 7 | 1 |
| 57 | 1 | 4 | 140 | 192 | 0 | 0 | 148 | 0 | 0.4 | 2 | 0 | 6 | 0 |
| 56 | 0 | 2 | 140 | 294 | 0 | 2 | 153 | 0 | 1.3 | 2 | 0 | 3 | 0 |
| 56 | 1 | 3 | 130 | 256 | 1 | 2 | 142 | 1 | 0.6 | 2 | 1 | 6 | 2 |
| 44 | 1 | 2 | 120 | 263 | 0 | 0 | 173 | 0 | 0 | 1 | 0 | 7 | 0 |

Fig 5: Sample Dataset

### B. Performance Metrics

The Improved K-mean –ANN Approach for Heart Disease Prediction is evaluated to compute the accuracy. TABLE II shows Confusion matrix. A confusion matrix is obtained to calculate the accuracy of classification. A confusion matrix shows how many instances have been assigned to each class. We are going to use two classes in our experiment, and therefore, we have a 2x2 confusion matrix.

TABLE III

CONFUSION MATRIX

| | Normal | Disease |
|---|---|---|
| **Normal** | TP | TN |
| **Disease** | FN | FP |

Accuracy is determined as follows:

$$Accuracy = \frac{TP+TN}{TOTAL}$$

Class T=Normal

Class F=Disease

TP=Got a positive result for the normal patient.

TN= Got a negative result for the normal patient.

FP= Got a positive result for disease patient

FN= Got a negative result for Disease patient

### C. Experimental Results

After testing the system 100 times with different 100 datasets this prediction system is giving average accuracy percentage 96.74. The system is tested with 100 permutation combination of normal and disease entries some sample calculations is shown in Table III. Table IV. And Fig.6 shows accuracy comparison. This accuracy comparison shows that with fewer attribute our system is giving good result.

TABLE IIIII
SAMPLE CALCULATION

| Sr No | Name of dataset | Total Entries | TP | TN | FP | FN | Accuracy Percentage |
|---|---|---|---|---|---|---|---|
| 1 | T1 | 100 | 23 | 70 | 7 | 0 | 93 |
| 2 | T2 | 100 | 30 | 66 | 4 | 0 | 96 |
| 3 | T3 | 100 | 35 | 64 | 2 | 0 | 98 |
| 4 | T4 | 100 | 24 | 75 | 1 | 0 | 99 |
| 5 | T5 | 100 | 26 | 74 | 0 | 0 | 100 |

TABLE IVII
ACCURACY COMPARISON

| Sr. No | Author | Description | Accuracy |
|---|---|---|---|
| 1 | Milan Kumari[6] | Studied cardiovascular disease using SVM, ANN, Decision tree | SVM=84.12% ANN=80.06% Decision tree=79.05 |
| 2 | Vanisree K[12] | Diagnosis the Congenital Heart Disease based on BPNN | 90% |
| 3 | Resul Das and Ibrahim Turkoglu [14] | Diagnosis heart disease using NN | 89% |
| 4 | Hongmei [17], | Developed support system for diagnosis of five major heart disease using MLP | 63.2 – 82.9% |
| 5 | K.Wankhade [18] | Used SVM and MLPNN to a diagnosis of heart disease. | 97.5% |
| 1 | Bhuvaneswari Amma N.G[19] | Predict the risk of cardiovascular disease using Genetic based neural network | 94.17% |
| 2 | Our System | Diagnosis heart angiographic disease status using improved K-means and ANN | 96.74% |

## VIII. CONCLUSIONS AND FUTURE WORK

The overall objective of this system is to predict more accurately the presence of heart disease using machine learning and data mining techniques. In this paper, we have presented a new approach that based on improved K-Means and ANN to model heart disease prediction. To validate the system, we have tested 100 data sets which show that our approach is giving a better result for classification. This prediction model is very helpful for doctors to diagnosis heart diseases efficiently with fewer attributes. In future, we can use the same hybrid technique to diagnosis other multiple diseases like–Cancer prediction, HIV prediction etc. to improve the prediction results.
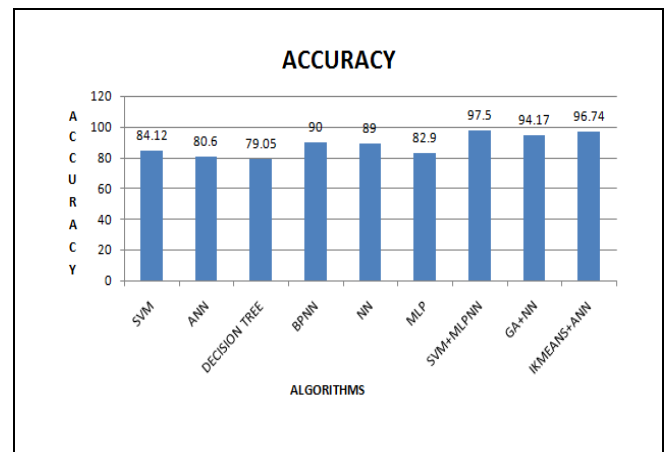


Fig.6 Accuracy Comparison

## ACKNOWLEDGMENT

## REFERENCES

[1]. Ankita Dewan, Meghna Sharma," Prediction of Heart Disease Using a Hybrid Technique in Data Mining Classification", 2nd International Conference on Computing for Sustainable Global Development IEEE, pp 704-706, 2015.

[2]. Ankur Makwana, Jaymin Patel," Decision Support System for Heart Disease Prediction using Data Mining Classification Techniques", International Journal of Computer Applications (0975 8887),Volume 117 - No. 22, pp 1-5 ,May 2015.

[3]. Carlos Ordonez, Edward Omincenski and Levien de Braal ,"Mining Constraint Association Rules to Predict Heart Disease", Proceeding of 2001, IEEE International Conference of Data Mining, IEEE Computer Society, ISBN-0-7695-1119-8, 2001, pp: 433-440.

[4]. Chaitrali S. Dangare, Sulabha Apte, "A Data Mining Approach for Prediction of Heart Disease Using Neural Networks", International Journal of Computer Engineering and Technology (IJCET), ISSN 0976 –6367(Print), ISSN 0976 – 6375(Online) Volume 3, Issue 3, October-December (2012), © IAEME

[5]. M Akhil Jabbar, BL Deekshatulu, Priti Chandra," Heart disease classification using nearest neighbor classifier with feature subset selection", Anale. Seria Informatica, 11, 2013.

[6]. Milan Kumari, Sunila Godara, "Comparative Study of Data Mining Classification Methods in Cardiovascular Disease Prediction", IJCST Vol. (2), Issue (2), June 2011.

[7]. Rajkumar, A. and G.S. Reena, Diagnosis of Heart Disease Using Data mining Algorithm. Global Journal of Computer Science and Technology, 2010. Vol. 10 (Issue 10).

[8]. Shadab Adam Pattekari and Asma Parveen," PREDICTION SYSTEM FOR HEART DISEASE USING NAIVE BAYES", International Journal of Advanced Computer and Mathematical Sciences ISSN 2230-9624, Vol 3, Issue 3, 2012, pp 290-294.

[9]. M. Durairaj, V. Ranjani ,"Data Mining Applications In Healthcare Sector: A Study," International Journal Of Scientific & Technology Research Volume 2, Issue 10, October 2013

[10]. Javad Kojuri, Reza Boostani, PooyanDehghani, FarzadNowroozipour, Nasrin Saki, " Prediction of acute myocardial infarction with artificial neural networks in patients with nondiagnostic electrocardiogram ",Journal of Cardiovascular Disease Research ,Vol 6 Issue 2 ,pp-51-59, Apr-Jun 2015

[11]. B.Srinivasa Rao, Dr. K. Nageswara Rao, Dr. S. P. SETTY, " An Approach for Heart Disease Detection by Enhancing Training Phase of Neural NetworkUsing Hybrid Algorithm", Advance Computing Conference (IACC), 2014 IEEE International,pp- 1211 – 1220, 21-22 Feb. 2014

[12]. Vanisree K, Jyothi Singaraju, "Decision Support System for Congenital Heart Disease Diagnosis based on Signs and Symptoms using Neural Networks", International Journal of Computer Applications (0975 8887) Volume 19 No.6, April 2011.

[13]. Preeti Gupta,BikrampalKaur,"Accuracy Enhancement of Artificial Neural Network using Genetic Algorithm",*International Journal of Computer Applications (0975 – 8887) Volume 103 – No 13, October 2014*

[14]. R. Das, I. Turkoglu, A. Sengur, Effective Diagnosis of Heart Disease through Neural Network Ensemble, "Expert Systems with Applications" vol. 36 issue 4, pp. 7675-7680, May (2009). [Available]: 10.1016/j.eswa.2008.09.013.

[15]. Roohallah Alizadehsani,[1] Jafar Habibi,[1] Zahra Alizadeh Sani,[2,*] Hoda Mashayekhi,[1] Reihane Boghrati,[1] Asma Ghandeharioun,[1] Fahime Khozeimeh,[3] and Fariba Alizadeh-Sani[3] ," Diagnosing Coronary Artery Disease via Data Mining Algorithms by Considering Laboratory and Echocardiography Features", s

Cardiovascular Med. 2013 Aug; 2(3): 133–139.Published online 2013 Jul 31

[16]. Lichman, M. (2013). UCI Machine Learning Repository [http://archive.ics.uci.edu/ml]. Irvine, CA: University of California, School of Information and Computer Science.

[17]. S. Prabhat Panday, N. Godara, "Decision Support System for Cardiovascular Heart Disease Diagnosis using Improved Multilayer Perceptron," *International Journal of Computer Applications* (0975 – 8887) Vol. 45– No.8, May (2012).

[18]. Gudadhe, K. Wankhade, and S. Dongre, "Decision support system for heart disease based on support vector machine and artificial neural network", In proceedings of IEEE International Conference on Computer and Communication Technology (ICCCT), pp. 741–745, November 2010.

[19]. Bhuvaneswari Amma N.G," Cardiovascular Disease Prediction System using Genetic Algorithm and Neural Network". International Conference on Computing, Communication and Applications, pp 1-5, Feb 2012

[20]. S.B.Bhalerao, DR. B.L.Gunjal, "Survey Of Heart Disease Prediction Based On Data Mining Algorithms", International Journal Of Advance Research And Innovative Ideas In Education, ISSN (O)-2395-4396 Vol-2 Issue-2 pp-856-860, 2016.