

## **Heart Disease Diagnosis Using Genetic Algorithms And Fuzzy Inference System With Multiple Fuzzy Labels**

**E. P. Ephzibah**

*School of Information Technology and Engineering,  
VIT University, Vellore, Tamilnadu, India,  
ep.ephzibah@vit.ac.in*

**T. Santhanam**

*D.G. Vaishnav College,  
Arumbakkam, Chennai, Tamilnadu, India.*

### **Abstract**

The focus of the work is to analyze the existing medical data and to find out the relevant features for disease classification and prediction using genetic algorithm and fuzzy logic techniques. Genetic algorithm is a stochastic searching technique that helps in identifying the relevant features. As the solution is obtained in parallel, the search becomes faster and simpler. Identifying the best fitness function is another important task. Appropriate genetic operators like genetic crossover, mutation and selection have been assigned suitable values for better performance. Classification is a task that evaluates the proposed work. The work takes into consideration the fuzzy classification system that can handle uncertain data. The work has introduced the concept of providing fuzzy labels ranging from 3 to 11 with a step value of 2. The fuzzy region for the attributes in the dataset is split into equal segments. Each segment is labeled as “1”, “2”, “3” and so on. As the number of labels increases the accuracy also increases. The proposed work after evaluation is found to produce an accuracy of 90.0% with a sensitivity of 100% and specificity of 81.25%. Increasing the number of labels to the next level doesn't improve the performance but reaches a saturation point beyond which it doesn't show any improvement in the performance. Hence in the proposed work is stopped after 11 labels. It is observed that the proposed work has outperformed some of the existing research works in the literature.

**Keywords:** Genetic algorithms, feature selection, fuzzy labels, rule generation, classification.

## **Introduction**

Coronary heart disease abbreviated as CHD is a major cause of death in many countries including India. The disease is caused by the deposition of fatty materials on the walls of the arteries that enter the heart. When the deposit becomes thicker that increases the risk rate for heart attacks where the heart doesn't get enough blood and oxygen supply. The reasons for the disease are hereditary, diabetics, high blood pressure, metabolic syndrome, kidney disease etc.[1] Due to the increase in the number of heart disease patients there are many cardiac hospitals established all over the world. Heart patients are allowed to undergo many tests. Using the test results, the doctors are able to diagnose whether the patient is really suffering from the disease or not. The proposed work will definitely help the doctors to perform the task more effectively and accurately using computing tools and techniques.

Computing tools exist for various disease diagnoses. "The advent of digital whole-slide scanner has been a revolution in imaging technology for histopathology" according to Metin Gurcan, Ph.D., an associate professor of Biomedical Informatics at The Ohio State University Medical Center [2]. A digital phototrichogram is a computerized automated system to define the hair growth and is commonly used by dermatologists. A fully computerized system called Trichoscan is available to the dermatologists for a thorough study of hair growth and to solve the issues related to it [3]. Another tool called Positron Emission Tomography (PET) is combined with a Computed Tomography (CT) scan to show images of the body tissues. Together, PET and CT scans create powerful images of the anatomy and biological functions of the body to reveal disease states that allow physicians to better view the structural details, location and changes in body tissues [4]. The above mentioned are a few among the hundreds of diagnosing tools in medical sector.

Heart disease diagnosing systems are coming in practice for accurate diagnosis and early treatment. The work focuses on the soft computing techniques like genetic algorithms and fuzzy logic for effective diagnosis of heart disease in patients. However, there are many research papers existing for the prediction of this disease. The proposed work has been proved to be more effective and accurate than the existing ones.

Genetic algorithm is an adaptive searching algorithm that provides solution to any optimization problem. It is a best robust evolutionary algorithm that operates on parallel solutions [5]. Genetic algorithm helps to escape from local minima and also handles multidimensional objective functions. The individuals in a genetic algorithm are characterized by chromosomes which are collection of genes. The gene decides the characteristic of the chromosomes for the solution. Every chromosome in a population is evaluated using fitness functions. A fitness function is an objective function that decides the fitness of the chromosome for its survival. The solution is obtained after several iterations of evaluations and also by applying the operators like crossover and mutation. The iterations are called as generations. Genetic algorithms are chosen for the work as they are capable of handling solutions of large space.

Feature selection is an important task when dealing with large volume of data. It is true that the prediction accuracy can be improved even after using a subset of features. A fitness function or objective function that helps to identify the important

features for prediction is used. This helps in improving the performance of the classification system. This method is cheaper, safer and helps us to easily understand the domain. In general features selection algorithms are categorized into two types as: Filters and Wrappers. The filters do not rely on any particular classification algorithm. Instead the features are selected based on the statistical properties of the individual features. The wrappers do rely on the classification system of the model and hence can improve the accuracy as and when iterated. The work in our approach focuses on the both of these methods in different situations.

In any real world application there is uncertainty that exists. Especially when taking the medical data into consideration uncertain information may confuse the physicians preventing them to take correct decisions. This kind of problems can be effectively handled using fuzzy systems as they are capable of handling uncertain data. Fuzzy system involve many phases like fuzzification, understanding the knowledge base containing the IF...THEN rules, the fuzzy inference system that maps the input to the corresponding output using the rules in the knowledge base, finally the defuzzification phase. The fuzzification and defuzzification phases are used to convert the crisp data into fuzzy data and fuzzy data into crisp output respectively. The fuzzy data can be interpreted using the linguistic labels that can help us to understand the rules easily.

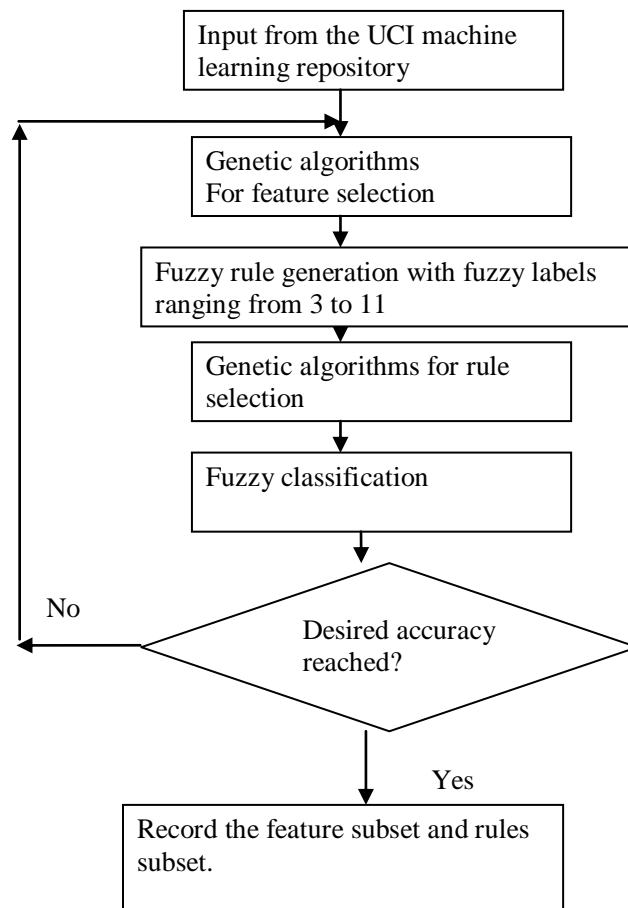
## **Literature Survey**

The related work in the literature has been analyzed under two divisions like disease diagnosis in general, soft computing techniques like genetic algorithms and fuzzy logic as a classifier for heart disease data. The analysis is based on the articles published in the past five years. According to the authors of the paper [6], the Adaptive Neuro fuzzy Inference System (ANFIS) is suitable for classification in the quantitative analysis of the intramuscular EMG signals among the other neural network classifiers. It has been also proved that the features selection techniques like autoregressive and discrete wavelet transform along with ANFIS has produced a classification accuracy of 95%. In paper [7] the authors have proposed a decision support system for malaria disease diagnosis using a fuzzy approach that has produced an accuracy of 80%. In [8] the authors have explored a computer-aided system that helps the physicians in diagnosing Alzheimer's disease. The authors have taken the Eigen image decomposition of single photon emission computed tomography images. The support vector machine has classified the data and has produced an accuracy of 92.78% accuracy. In [9] a computer aided detection system has been designed for accurate diagnosis of human brain tumor through MRI images. The computing techniques that were used are feedback pulse coupled neural network, discrete wavelet transform and principal component analysis for segmentation, feature extraction and reducing the dimensionality of the wavelet respectively. Feed-forward back propagation neural network has been used as the classifier. All the computing tools and techniques have been applied on the data and produced an accuracy of 99% with a total of 101 images.

In paper [10] the authors have used the computational fluid dynamics method to understand the blood flow in human body. The authors have proposed this method for identifying the risk factors for the coronary artery disease development. They have given an accuracy of 84.3% over 130 samples using fractional flow reserve derived from CT angiography. In paper [11] the authors have applied rough sets with logistic regression on heart disease data taken from the UCI machine learning repository and have produced an accuracy of 83.93%. In paper [12] the authors have proposed an algorithm that helps the physicians to detect the coronary heart disease using rose angina questionnaire. The model of classification has proved to be 53% and 89% in terms of sensitivity and specificity respectively. Authors in paper [13] have used ensemble based methods for predicting the heart disease with neural network classifier. Ensemble method is a combination of the predicted values using multiple predecessor models. Their proposed method has given an accuracy of 89.01% as the classification accuracy. In paper [14] the authors have taken a hybrid method of combining the artificial neural network and fuzzy neural network for the diagnosis of heart disease data taken from the UCI machine learning repository and have found the classification accuracy to be 86.8% using a k fold cross validation method. The paper [15] focuses on classification of heart disease in patients using artificial immune recognition system with fuzzy resource allocation mechanism. The data taken from the UCI repository has been preprocessed using a new weighting scheme using k-nearest neighbor method. The classification accuracy was to be 87%.

### **Materials and Methods**

The proposed system has series of steps which are to be done in a sequential manner. The sequence is depicted pictorially in Fig1.

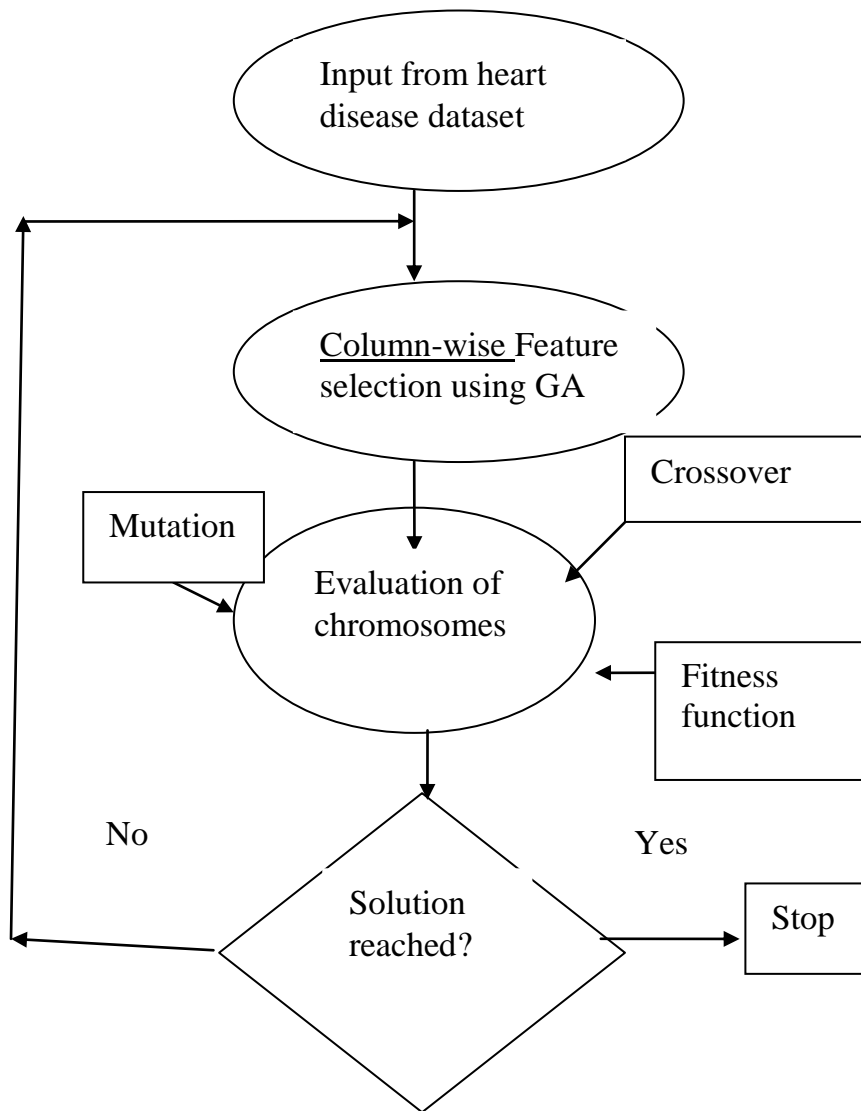


**Figure 1:** Proposed Work

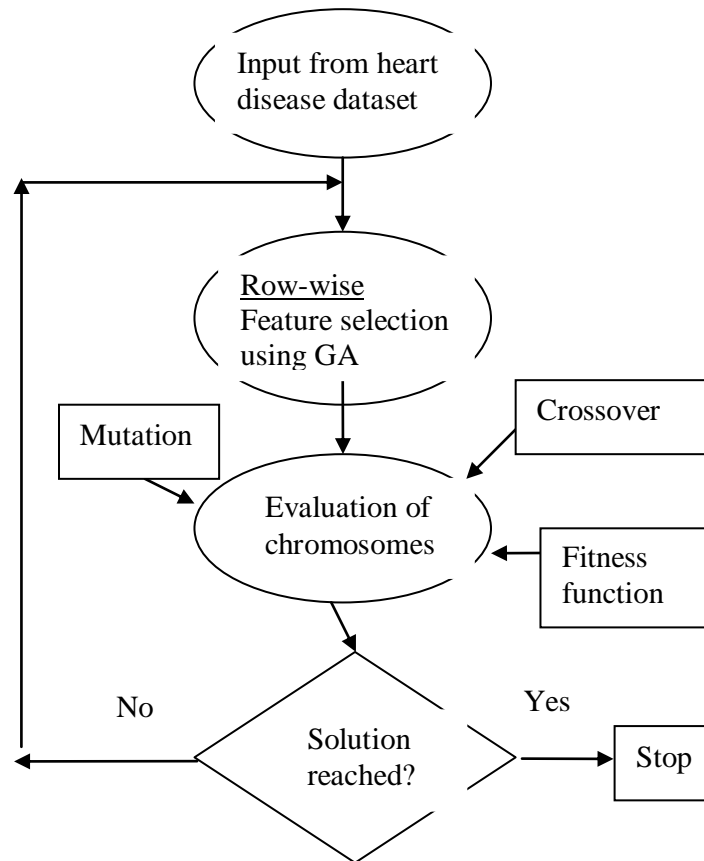
### **Genetic algorithms**

Genetic algorithms are used in the proposed work for a stochastic search that provides optimal solution to the problem. The algorithm has been used twice in the work for two different purposes. One for the selection of features (attributes) subset from the original set of features and the other for selection of rules based on the evaluation criteria using the fitness function. Hence they are called as column wise selection for attributes and row wise selection for rules.

The pictorial representation of the functioning of genetic algorithms in our work is as given in figures 2 and 3:



**Figure 2:** Selecting Attributes Using Genetic Algorithms



**Figure 3:** Selecting Rules For Fuzzy System Using Genetic Algorithm

### Genetic Search

Genetic algorithms perform a probabilistic search. The search takes into account the chromosomes in the population based on their fitness function value. This value helps to identify the survivability of the chromosomes into the next generation. The chromosome is evaluated and selected with probabilistic phenomena that the best chromosomes are not always selected and also the worst are not necessarily excluded. Random picking of the chromosomes is done with respect to the crossover and mutation points.

### Genetic Operators

Genetic operators are namely, the selection, the reproduction, the crossover and the mutation. The search is manageable and the process becomes optimal with the help of these operators. Without these operators the search would get trapped in local minima or may result in inefficient performance. The operators deal with the population in each and every generation and produce a progress in the solution space.

Selection is the process of selecting the chromosomes from the population (parent chromosomes) that have a better fitness value. They are treated as child chromosomes for the next generation. Selection can promote the solution to the optimal solution.

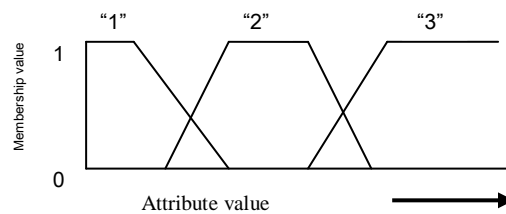
The selected chromosomes are taken for reproduction. The reproduction phase has the operators like crossover and mutation. Two parents are selected and a crossover point is fixed. The genes in those two chromosomes are interchanged among the chromosomes. Thus providing a crossover operation in which a chromosome becomes much better than other chromosomes. In mutation a single parent is selected. Based on the mutation point a gene in the chromosome is flipped. This again provides a better child chromosome for the next generation. The action takes place based on the mutation rate.

### Genetic Feature Selection

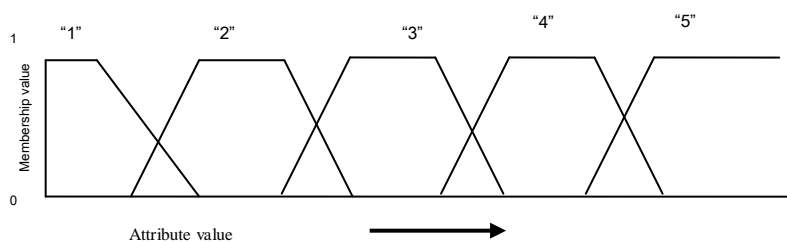
Genetic algorithm plays the role of feature selector in this work. It is to select a subset of size  $d$  ( $d < n$ ) that leads to the smallest classification error, where  $n$  is the total number of features. It is an optimization problem that searches for the feature subsets in the solution space. The two objectives in this type of search are to minimize the number of features as well as to minimize the classification error rate.

### Fuzzy Labels and Regions

Fuzzy labels deal with numeric values. The labels are called the linguistic variables that can help anyone to easily understand the mapping between the input value and the output label. The ranges of the label vary based on the type of attribute. The attribute value is divided into different labels as '1', '2', and '3'. The figure 4 depicts the split in every case of labeling.

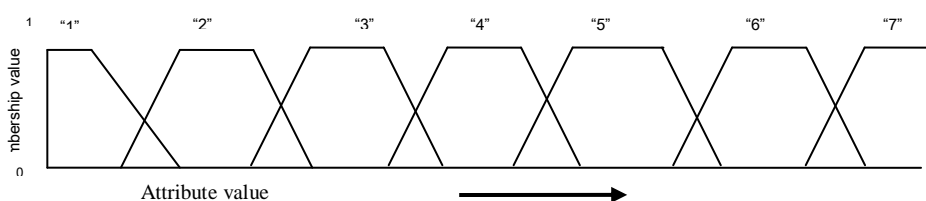


**Figure 4:** Fuzzy labels for the attribute values



**Figure 5:** Fuzzy logic with 5 labels for the attribute values





**Figure 6:** Fuzzy logic with 7 labels for the attribute values

Figure 6: Fuzzy logic with 7 labels for the attribute values

In Figure 5 and Figure 6, the labels are split into 5 and 7 divisions respectively. Thus the fuzzy region has been split based on the number of labels providing more importance to every region.

The region for the fuzzy labels can change depending on the attribute value. The following table 1 provides the division in regions and labels for every attribute in the dataset.

**Table 1:** Attribute names and their corresponding regions and labels in the fuzzy logic

Attribute number	Attribute name	Regions	Labels
1	Age	29-45 45-61 61-77	"1" "2" "3"
2	Sex	1-male 0-female	"1" "2"
3	Cp	1-typical 2-atypical 3-non-anginal 4-asymptotic	"1" "2" "3" "4"
4	Trestbps	94-129 129-164 164-200	"1" "2" "3"
5	Chol	126-272 272-418 418-564	"1" "2" "3"
6	Fbs	1-true 2-false	"1" "2"
7	Restecg	1-normal 2-abnormality 3-hypertrophy	"1" "2" "3"
8	Thalach	71-114 114-158	"1" "2"

		158-202	“3”
9	Exang	1-yes 2-no	“1” “2”
10	Oldpeak	0-2 2-4 4-6	“1” “2” “3”
11	Slope	1-upsloping 2-flat 3-downsloping	“1” “2” “3”
12	Ca	1-none 2-mild 3-moderate 4-severe	“1” “2” “3” “4”
13	Thal	3- normal 6- fixed defect 7- reversable defect	“1” “2” “3”

### Fuzzy Rule Generation

Fuzzy rules play an important role in fuzzy inference system. These rules provide the understanding about the classifier as they contain an antecedent and consequent part. Rule generation is done based on the data available in the dataset. The rules are derived from the training data taken separately from the dataset. The generated rules undergo a process where a rule subset is selected using genetic algorithms. Feature selection and rule selection are performed using effective fitness functions that can upgrade the performance of the proposed model. Hence the best set of attributes and rules together form the model for the heart disease diagnosis system in the proposed work.

### Fuzzy Classification

Fuzzy classification is one of the applications of fuzzy theory. Every instance in the dataset can have membership in many different classes to different degrees. It is to be noted that all the membership values for an instance sum to 1. The AND operator combines the input to produce the output. Among the various methods for defuzzification, the chosen method is the centroid method. This method is the most prevalent and physically appealing of all the defuzzification methods [16]. This method is also called as center of area or center of gravity method as it returns the center of area under the curve. Among the other defuzzification methods like bisector, largest of maximum, smallest of maximum and middle of maximum this method has been referred and used by many researchers [17] [18] [19][20].

### Stratified 10 Fold Cross Validation

The standard method for system evaluation is the stratified k fold cross validation method. This evaluation method splits the data into 10 folds or subsets of equal size. Each subset will get a chance as a train set and a test set. The subsets are stratified

before the cross validation. The evaluated estimates are taken and their average is calculated as the final estimate. Thus stratification reduces the estimate’s variance.

### **Results and Discussion**

The program was written in MATLAB R2010a and executed. The genetic algorithm has been implemented and features were selected. For the selected set of features the rules were generated for fuzzy inference system. Genetic algorithms were again used for rule selection. The rules obtained from the training data were reduced to a limited number. The dataset contains 297 instances with 13 input attributes and a class label. The data has been divided into 10 folds based on the stratified approach. Tables 2-6 provide the detailed information about the performance measures like accuracy, specificity and sensitivity. The program executed for every fold is recorded. The average of all the folds gives the final accuracy, specificity and sensitivity. Figure 7 depicts the performance measure graphically. This provides the variation and the improvement in the accuracy as the number of labels increases. It is not true that increasing the number of labels will always improve the performance. At a certain level of splitting the fuzzy region it reaches a saturation point in the performance. Hence the proposed work is stopped after 11 labels. Table 7 shows the performance measures of the proposed work with labels like 3, 5, 7, 9, 11 and 13. It is also observed that the performance reaches a saturation point after 11 labels.

**Table 2:** Performance metrics for 3 labels

<b>SNo.</b>	<b>No.of Folds</b>	<b>Accuracy</b>	<b>Specificity</b>	<b>Sensitivity</b>
1	Fold1	90	81.25	100
2	Fold2	86.6667	87.5	85.7143
3	Fold3	83.3333	78.9474	90.9091
4	Fold4	80.3333	92.85	68.75
5	Fold5	86.6667	94.444	75.0
6	Fold6	80.0	88.8889	66.6667
7	Fold7	83.3333	64.2857	100
8	Fold8	90	81.25	100
9	Fold9	83.333	68.75	100
10	Fold10	90	81.25	100
<b>Average</b>		<b>85.36663</b>	<b>81.9416</b>	<b>88.70401</b>

**Table 3:** Performance metrics for 5 labels

<b>SNo.</b>	<b>No.of Folds</b>	<b>Accuracy</b>	<b>Specificity</b>	<b>Sensitivity</b>
1	Fold1	90	81.25	100
2	Fold2	90	81.25	100
3	Fold3	83.3333	78.9474	90.9091

4	Fold4	83.3333	78.9474	90.9091
5	Fold5	86.6667	94.444	75
6	Fold6	86.6667	88.8889	66.6667
7	Fold7	83.3333	64.2857	100
8	Fold8	90	81.25	100
9	Fold9	83.333	68.75	100
10	Fold10	90	81.25	100
<b>Average</b>		<b>86.66663</b>	<b>79.92634</b>	<b>92.34849</b>

**Table 4:** Performance metrics for 7 labels

SNo.	No.of Folds	Accuracy	Specificity	Sensitivity
1	Fold1	90	81.25	100
2	Fold2	90	81.25	100
3	Fold3	83.3333	78.9474	90.9091
4	Fold4	83.3333	78.9474	90.9091
5	Fold5	86.6667	94.444	75
6	Fold6	86.6667	88.8889	66.6667
7	Fold7	86.6667	94.444	75
8	Fold8	90	81.25	100
9	Fold9	86.6667	94.444	75
10	Fold10	90	81.25	100
<b>Average</b>		<b>87.33334</b>	<b>85.51157</b>	<b>87.34849</b>

**Table 5:** Performance metrics for 9 labels

SNo.	No.of Folds	Accuracy	Specificity	Sensitivity
1	Fold1	90	81.25	100
2	Fold2	90	81.25	100
3	Fold3	86.6667	94.444	75
4	Fold4	86.6667	94.444	75
5	Fold5	86.6667	94.444	75
6	Fold6	86.6667	88.8889	66.6667
7	Fold7	86.6667	94.444	75
8	Fold8	90	81.25	100
9	Fold9	86.6667	94.444	75

10	Fold10	90	81.25	100
<b>Average</b>		<b>88.00002</b>	<b>88.61089</b>	<b>84.16667</b>

**Table 6:** Performance measure for 11 labels

SNo.	No.of Folds	Accuracy	Specificity	Sensitivity
1	Fold1	90	81.25	100
2	Fold2	90	81.25	100
3	Fold3	90	81.25	100
4	Fold4	90	81.25	100
5	Fold5	90	81.25	100
6	Fold6	90	81.25	100
7	Fold7	90	81.25	100
8	Fold8	90	81.25	100
9	Fold9	90	81.25	100
10	Fold10	90	81.25	100
<b>Average</b>		<b>90</b>	<b>81.25</b>	<b>100</b>

**Table 7:** Accuracy details regarding the various labels

SNo.	No.of labels	Specificity	Sensitivity	Accuracy
1	3	81.94	88.7	85.4
2	5	79.93	92.35	86.66
3	7	85.51	87.35	87.33
4	9	88.61	84.17	88.00
5	11	81.25	100	90

**Figure 7.** Performance measures for the proposed method

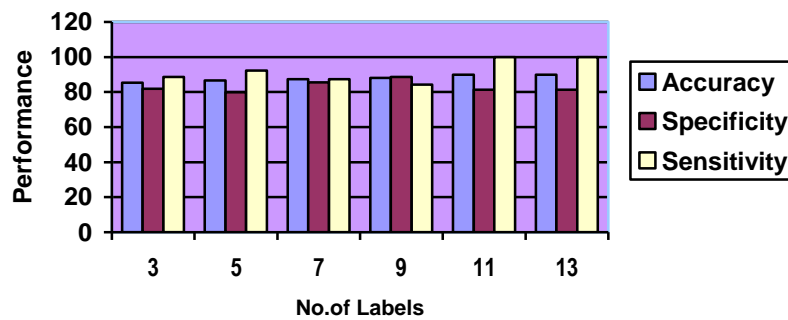


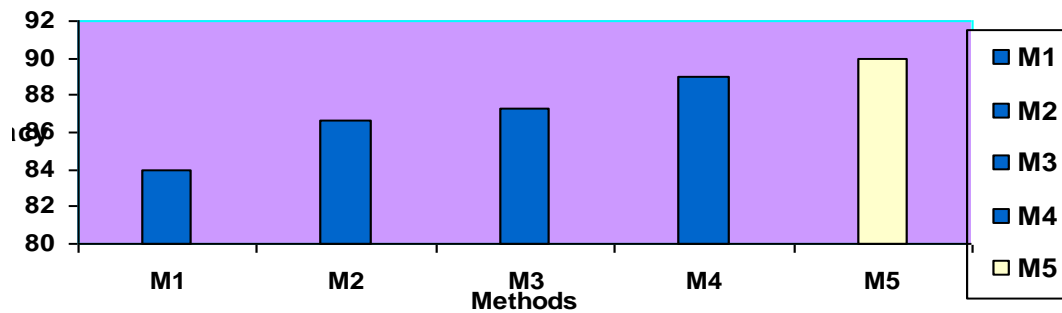
Table 8: lists the existing methods in the literature and compares the accuracy of those methods with the proposed method. This is pictorially depicted in figure 8.

**Table 8:** Comparison of the proposed work with the existing works

SNo.	Existing Methods	Accuracy
M1	Rough sets with logistic regression [11]	83.93%
M2	Artificial neural network and fuzzy neural network [14]	86.66%
M3	Artificial immune recognition system with fuzzy resource allocation mechanism.[15]	87.33%
M4	Ensemble methods [13]	89.01%
M5	Proposed Method **	90.0%

Note :\*\* Genetic algorithm and fuzzy inference system.

**Figure 8. Performance Comparison with existing methods**



## Conclusion

Disease diagnosis requires a strong truth based prediction as it deals with life. It is a daily practice for the physicians to examine the patient and carefully diagnose the disease. The proposed work helps the medical practitioners to correctly examine and diagnose the disease using some of the intelligent soft computing techniques like genetic algorithms, fuzzy inference system using a stratified approach for validation. The propose work uses the genetic algorithm a heuristic searching technique for both attributes selection and rule election. The fuzzy region of the fuzzy classifier has been split into 3, 5, 7, 9 and 11 labels. Dividing the data into different fuzzy regions and labeling them into linguistic fuzzy labels adds more weight to the system. Even though splitting the fuzzy region into different labels improves the performance of the classifier, there is a saturation point after which there is no improvement in the performance. The work can be more beneficial and provide better accuracy when compared with the existing systems.

## References

- [1]. Longo, D.L., Kasper, J.L., Jameson, A.S., Fauci, S.L., Hauser, J., Loscalzo, Harrison's, 2012, "Principles of Internal Medicine", McGraw-Hill.
- [2]. [www.osc.edu/press/computer\\_assisted\\_diagnosis\\_tools\\_to\\_aid\\_pathologist\\_s](http://www.osc.edu/press/computer_assisted_diagnosis_tools_to_aid_pathologist_s).
- [3]. Hillmann, K., Blume-Peytavi, U., 2009, "Diagnosis of hair disorders", PubMed, *Semin Cutan Med Surg.*, pp.33-38.
- [4]. <http://www.radiologyinfo.org>
- [5]. Nguyen, T., Khosravi, A., Creighton, D., Nahavandi, S., 2015, "Classification of Healthcare Data using Genetic Fuzzy Logic System and Wavelets", *Expert Systems with Applications* (42), pp.2184-2197.
- [6]. Abdulhamit Subasi, 2012, "Classification of EMG signals using combined features and soft computing techniques". *Applied soft computing*, Volume 12, pp.2188-2198.
- [7]. Faith-Michael, E., Uzoka, Joseph Osuji, Okure Obot, 2011, "Clinical decision support system (DSS) in the diagnosis of malaria: A case comparison of two soft computing methodologies", *Expert Systems with Applications* 38, pp.1537–1553
- [8]. Ignacio Alvarez Illán, Juan Manuel Górriz, Javier Ramírez, Elmar, W. Lang, Diego Salas-Gonzalez, Carlos, G., Puntonet, 2012, "Bilateral symmetry aspects in computer-aided Alzheimer's disease diagnosis by single-photon emission-computed tomography imaging", *Artificial Intelligence in Medicine* 56, pp.191–198.
- [9]. Muhammad Aziz Rahman, Nicola Spurrier, Mohammad Afzal Mahmood, Mahmudur Rahman, Sohel Reza Choudhury, Stephen Leeder, 2014, "Computer-aided diagnosis of human brain tumor through MRI: A survey and a new algorithm", *Expert Systems with Applications* 41, pp.5526–5545.
- [10]. Zhonghua Sun, Lei Xu., 2014, "Computational fluid dynamics in coronary artery disease", *Computerized Medical Imaging and Graphics*, Article in Press.
- [11]. Yuehjen E. Shao, Chia-Ding Hou, Chih-Chou Chiu., 2014, "Intelligent modeling schemes for heart disease classification", *Applied Soft Computing* (14), pp.47–52.
- [12]. Muhammad Aziz Rahman, Nicola Spurrier, Mohammad Afzal Mahmood, Mahmudur Rahman, Sohel Reza Choudhury, Stephen Leeder, 2013, "Rose Angina Questionnaire: Validation with cardiologists' diagnoses to detect coronary heart disease in Bangladesh", *Indian Heart Journal*. Volume 65, pp.30–39.
- [13]. Resul Das, Ibrahim Turkoglu, Abdulkadir Sengur., 2009, "Effective diagnosis of heart disease through neural networks ensembles", *Expert Systems with Applications* 36, pp.7675–7680.

- [14]. Humar Kahramanli, Novruz Allahverdi., 2008, “Design of a hybrid system for the diabetes and heart diseases” *Expert Systems with Applications*, 35, pp.82–89.
- [15]. Kemal Polat, Seral Sahan, Salih Gu˘nes., 2007, “Automatic detection of heart disease using an artificial immune recognition system (AIRS) with fuzzy resource allocation mechanism and k-nn (nearest neighbour) based weighting preprocessing”, *Expert Systems with Applications* 32, pp.625–631.
- [16]. Lee, C., 1990, “Fuzzy logic in control systems: fuzzy logic controllers—part I and II”, *IEEE Trans. Syst. Man. Cybernet.* 20, pp.404-435.
- [17]. Feilong Liu., 2008, “An efficient centroid type-reduction strategy for general type-2 fuzzy logic system”, *Information Sciences* 178 pp.2224–2236.
- [18]. O.W. Samuel, M.O. Omisore, B.A. Ojokoh., 2013, “A web based decision support system driven by fuzzy logic for the diagnosis of typhoid fever”, *Expert Systems with Applications* 40, pp.4164–4171.
- [19]. Debabrata Pal, K.M. Mandana, Sarbajit Pal, Debranjana Sarkar, Chandan Chakraborty., 2012, “Fuzzy expert system approach for coronary artery disease screening using clinical parameters”, *Knowledge-Based Systems* 36, pp.162–174.
- [20]. Stavros Lekkas, Ludmil Mikhailov., 2010, “Evolving fuzzy medical diagnosis of Pima Indians diabetes and of dermatological diseases”, *Artificial Intelligence in Medicine* 50, pp.117–126.