

# A hybrid cost-sensitive ensemble for heart disease prediction

Zhenya Qi

Tianjin University <https://orcid.org/0000-0003-4938-9785>

Zuoru Zhang (✉ [zhangzuoru@tju.edu.cn](mailto:zhangzuoru@tju.edu.cn))

<https://orcid.org/0000-0003-2118-6548>

---

## Research article

**Keywords:** cost-sensitive, ensemble, heart disease

**Posted Date:** February 7th, 2020

**DOI:** <https://doi.org/10.21203/rs.2.22946/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on February 25th, 2021. See the published version at <https://doi.org/10.1186/s12911-021-01436-7>.

## RESEARCH

# A hybrid cost-sensitive ensemble for heart disease prediction

Zhenya Qi<sup>1†</sup> and Zuoru Zhang<sup>2\*</sup>

\*Correspondence:

[zhangzuoru@tju.edu.cn](mailto:zhangzuoru@tju.edu.cn)

<sup>2</sup>School of Mathematical Science, Hebei Normal University, Yuhua District, 050024 Shijiazhuang, PR China

Full list of author information is available at the end of the article

<sup>†</sup>Equal contributor

## Abstract

Heart disease is the primary cause of morbidity and mortality in the world. It includes numerous problems and symptoms. The diagnosis of heart disease is difficult because there are too many factors to analyze. What's more, the misclassification cost could be very high. In this paper, we firstly propose a cost-sensitive ensemble model to improve the accuracy of diagnosis and reduce the misclassification cost. The proposed model contains five heterogeneous classifiers: random forest, logistic regression, support vector machine, extreme learning machine and k-nearest neighbor. Then, experiments are done on three datasets from UCI machine learning repository. The highest classification accuracy of 91.74%, highest G-mean of 90.55%, highest precision of 96.11%, highest recall of 89.61% and lowest misclassification cost of 30.32% are achieved by the proposed model according to ten-fold cross validation. The results demonstrate that the performance of the proposed model is superior to those of previously reported classification techniques.

**Keywords:** cost-sensitive; ensemble; heart disease

## 1 Introduction

Heart disease is any disorder that influences the heart's ability to function normally [1]. As the leading cause of death, heart disease is responsible for nearly 30% of the global deaths annually [2]. In China, it is estimated that 290 million people are suffering from heart disease, and the rate of death caused by heart disease is more than 40% [3]. According to The European Society of Cardiology (ESC), nearly half of the heart disease patients die within initial two years [4]. Therefore, accurate diagnosis of heart disease in early stages is of great importance in improving security of heart [5].

However, as it's associated with numerous symptoms and various pathologic features such as diabetes, smoking and high blood pressure, the diagnosis of heart disease remains a huge problem for less experienced physicians [6]. In order to detect heart disease, several diagnostic methods have been developed, Coronary angiography (CA) and Electrocardiography (ECG) are the most widely used among them, but they both have serious defects. ECG may fail to detect the symptoms of heart disease in its record [7] while CA is invasive, costly and needs highly-trained operators [8].

Computer-aided diagnostic methods based on machine learning predictive models can overcome these difficulties. Such methods are noninvasive and provide proper and objective diagnoses, and hence can reduce the suffering of patients [9]. Various

machine learning predictive models such as random forest (RF) [10], logistic regression (LR) [11], support vector machine (SVM) [12], extreme learning machine (ELM) [13] and k-nearest neighbor (KNN) [14] have been developed and widely used as classifiers to assist doctors in diagnosing heart disease. Dogan et al. [15] built a RF classification model for symptomatic heart disease. The clinical characteristics of the 1545 and 142 subjects were used for training and testing respectively, and the classification accuracy was 78%. Detrano et al. [16] proposed a LR classifier for heart disease classification and obtained an accuracy of 77% in 3 patient test groups. Gokulnath and Shantharajah [17] proposed a classification model based on genetic algorithm (GA) and SVM, obtaining an accuracy of 88.34% on Cleveland heart disease dataset. Subbulakshmi et al. [18] performed a detailed analysis of different activation functions of ELM using Statlog heart disease dataset. The results indicated that ELM achieved an accuracy of 87.5%, higher than other methods. Duch et al. [19] used KNN classifier to predict heart disease on Cleveland heart disease dataset and achieved an accuracy of 85.6%, superior to other machine learning techniques.

It is realized that no single model exists that is superior for all classification problems, because different machine learning algorithms consider datasets with different features in different aspects [20]. One way to overcome the limitations of a single classifier is to use an ensemble model. An ensemble model is the combination of multiple sets of classifiers, it outperforms the individual classifiers because the variance of error estimation is reduced [21]. In recent years, many ensemble approaches have been proposed to improve the performance of heart disease diagnosis systems. For instance, Das et al. [22] proposed a neural networks ensemble and obtained 89.01% classification accuracy from the experiments made on the data taken from Cleveland heart disease dataset. Bashir et al. [23] employed the ensemble of five heterogeneous classifiers on five heart disease datasets. The proposed ensemble classifier achieved the high diagnosis accuracy of 87.37%. Khened et al. [24] presented an ensemble system based on deep fully convolutional neural network (FCN) and achieved a maximum classification accuracy of 100% on Automated Cardiac Diagnosis Challenge (ACDC-2017) dataset. Therefore, we use an ensemble classifier to predict the presence or absence of heart disease in present study.

From the previous studies, it is observed that traditional medical decision support systems usually focused only on the maximization of classification accuracy without taking the unequal misclassification costs between different categories into consideration. However, in the field of medical decision making, it is often the minority class that is of higher importance [25]. Further, the cost associated with missing a patient (false negative) is much higher than that of mislabeling a healthy instance (false positive) [26]. Therefore, traditional classifiers inevitably result in a defective decision support system. In order to overcome this limitation, in this paper we combine the classification results of individual classifiers in a cost-sensitive way so that classifiers that help reduce the costs gain more weights in the final decision.

The contribution of the proposed research is to design a cost-sensitive ensemble classification model based on machine learning algorithms. Five individual classifiers including SVM, ELM, KNN, LR and RF are used for diagnosis of heart disease. Relief feature selection algorithm is used to select the most important features that

have great influence on target predicted value. In order to evaluate the performance of ensemble model, various performance evaluation metrics such as classification accuracy, misclassification cost (MC), G-mean, precision, recall and receiver optimistic curves (ROC) are used. In addition, data preprocessing techniques are applied to the heart disease datasets. The main contributions of the proposed research are as follows:

(1) The proposed ensemble model is a novel combination of heterogeneous classifiers which had outstanding performance in previous studies. Besides, their characteristics are quite complementary.

(2) We have used the accuracy and MC to combine the results of individual classifiers. The proposed ensemble model not only focuses on high classification accuracy, but also concerns the costs patients have to pay for misclassification.

(3) Compared with five individual classifiers and previous studies, the proposed ensemble model has achieved excellent classification results. The best performance of the proposed model are as follows: accuracy of 91.74%, MC of 30.32%, G-mean of 90.55%, precision of 96.11% and recall of 89.61%.

The rest of the paper is organized as follows. Section 2 offers brief background information concerning Relief algorithm and each individual classifier. Section 3 presents the framework of the proposed cost-sensitive ensemble diagnosis model in detail. Section 4 describes the experimental results and compares the ensemble model with individual classifiers and previous methods. Finally, the conclusions and directions for future works are summarized in Section 5.

## 2 Backgrounds and Preliminaries

### 2.1 Relief Feature Selection Algorithm

Relief is a kind of famous filter feature selection algorithm which adopts a relevant statistics to measure the importance of the feature. This statistics can be seen as the weight of each feature. Top  $k$  features of bigger weights are selected. Therefore, the key is to determine the relevant statistics [27].

Assume  $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$  is a dataset.  $x_i$  is an input feature vector and  $y_i$  is a class label corresponding to  $x_i$ . First, select a sample  $x_i$  randomly. Then, Relief attempts to find out its nearest sample  $x_{i,nh}$  from samples of its same class and nearest sample  $x_{i,nm}$  from samples of its different class using the same techniques as in KNN,  $x_{i,nh}$  is called "near-hit",  $x_{i,nm}$  is called "near-miss". Next, update the weight of a feature  $A$  in  $W$  as described in Algorithm 1 [28, 29]. Repeat the random sampling steps for  $m$  times and get the average value of  $W[A]$ ,  $W[A]$  is the weight of feature  $A$ .

#### RELIEF Algorithm

**Require:** for each training instance, a vector of feature values and the class value

$n \leftarrow$  number of training instances

$a \leftarrow$  number of features

**Parameter:**  $m \leftarrow$  number of random training instances out of  $n$  used to update  $W$

Initialize all feature weights  $W[A] := 0.0$

**For:**  $i := 1$  to  $m$  do

Randomly select a target instance  $R_i$

find a nearest hit  $H$  and nearest miss  $M$  (instances)

**For:**  $A := 1$  to  $a$  do

$W[A] := W[A] - diff(A, R_i, H)/m + diff(A, R_i, M)/m$

**End For**

**End For**

**Return** the weight vector  $W$  of feature scores that compute the quality of features

Algorithm 1: Pseudocode of the Relief algorithm

In Algorithm 1,  $diff(x_a^j, x_b^j)$  depends on the type of feature  $j$ . For discrete feature  $j$ :

$$diff(x_a^j, x_b^j) = \begin{cases} 0, & x_a^j = x_b^j \\ 1, & otherwise, \end{cases}$$

for continuous feature  $j$ :

$$diff(x_a^j, x_b^j) = |x_a^j - x_b^j|.$$

Repeatedly operate for  $n$  times, then average the weights of each feature. Finally, choose the top  $k$  features for classification.

## 2.2 Machine Learning Classifiers

Machine learning classification algorithms are used to distinguish heart disease patients from healthy people. Five popular classifiers and their theoretical backgrounds are discussed briefly in this paper.

### 2.2.1 Random Forest

RF is a machine learning algorithm based on the ensemble of decision trees [30]. In traditional decision tree methods such as C4.5 and C5.0, all the features are used for generating the decision tree. In contrast, RF builds multiple decision trees and chooses the random subspaces of the features for each of them. Then, the votes of trees are aggregated and the class with the most votes is the prediction result [31].

As an excellent classification model, RF can successfully reduce the overfitting and calculate the nonlinear and interactive effects of variables. Besides, the training of each tree are done separately, so it could be done in parallel, which reduced the training time needed. Finally, combining the prediction result of each tree could reduce the variance and improve the accuracy of the predictions. There are many studies showing the performance superiority of RF over other machine learning methods [32–34]. In present study, the number of decision trees to build the RF is 50.

### 2.2.2 Logistic Regression

LR is a generalized linear regression model [35]. Therefore, it is similar with multiple linear regression in many aspects. Usually, LR is used for binary classification problems where the predictive variable  $y \in [0, 1]$ , 0 is negative class and 1 is positive class. But it can also be used for multi-classification.

In order to distinguish heart disease patients from healthy people, a hypothesis  $h(\theta) = \theta^T X$  is proposed. The threshold of classifier output is  $h_\theta(x) = 0.5$ , which is to say, if the value of hypothesis  $h_\theta(x) \geq 0.5$ , it will predict  $y = 1$  which means that the person is a heart disease patient, otherwise the person is healthy. Hence, the prediction is done.

The sigmoid function of LR can be written as:

$$h_\theta(x) = \frac{1}{1 + e^{-z}},$$

where  $z = \theta^T X$ .

The cost function of LR can be written as:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(y_i, y'_i),$$

where  $m$  is the number of instances to be predicted,  $y_i$  is the real class label of the  $i$ th instance, and  $y'_i$  is the predicted class label of the  $i$ th instance.

$$\text{cost}(y_i, y'_i) = \begin{cases} 0, & y_i = y'_i \\ 1, & \text{otherwise.} \end{cases}$$

### 2.2.3 Support Vector Machine

Invented by Cortes and Vapnik [36], SVM is a supervised machine learning algorithm which has been widely used for classification problems [26, 37, 38]. The output of SVM is in the form of two classes in a binary classification problem, making it a non-probabilistic binary classifier [39]. SVM tries to find a linear maximum margin hyperplane that separates the instances.

Assume the hyperplane is  $w^T x + b = 0$ , where  $w$  is a dimensional coefficient vector, which is normal to the hyperplane of the surface,  $b$  is offset value from the origin, and  $x$  is dataset values. Obviously, the hyperplane is determined by  $w$  and  $b$ . The data points nearest to the hyperplane are called support vectors. In the linear case,  $w$  can be solved by introducing Lagrangian multiplier  $\alpha_i$ . The solution of  $w$  can be written as:

$$w = \sum_{i=1}^m \alpha_i y_i x_i,$$

where  $m$  is the number of support vectors and  $y_i$  are target labels to  $x$ . The linear discriminant function can be written as:

$$g(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i x_i^T x + b\right),$$

$\text{sgn}$  is the sign function that calculates the sign of a number,  $\text{sgn}(x) = -1$  if  $x < 0$ ,  $\text{sgn}(x) = 0$  if  $x = 0$ ,  $\text{sgn}(x) = 1$  if  $x > 0$ . The nonlinear separation of data set is performed by using a kernel function. The discriminant function can be written as:

$$g(x) = \text{sgn}\left(\sum_{i=1}^m \alpha_i y_i K(x_i, x) + b\right),$$

where  $K(x_i, x)$  is the kernel function. The polynomial kernel function is used in present study.

### 2.2.4 Extreme Learning Machine

ELM was first proposed by Huang et al. [40]. Similar to a single layer feed-forward neural network(SLFNN), ELM is also a simple neural network with a single hidden

layer. However, unlike a traditional SLFNN, the hidden layer weights and bias of ELM are randomized and need not to tune, and the output layer weights of ELM are analytically determined through simple generalized inverse operations [40, 41].

Considering the gradient-based learning algorithms suffer from slow training speed and poor generalization performance, ELM has overcome the two limitations by randomly initiating hidden layer parameters without iteratively tuning and tedious iterative process [42]. Then, ELM can obtain the global optimal solution at extremely fast learning speed[43]. With these advantages, ELM has gained popularity in many application fields [44–46].

### 2.2.5 *K-Nearest Neighbor*

KNN a supervised classification algorithm. Its procedure is as follows: when a new case is given, first search the database to find the  $k$  historical cases which are closest to the new case, namely  $k$ -nearest neighbors, and then these neighbors vote on the class label of the new case. If a class has the most nearest neighbors, the new case is determined to belong to the class [47]. The following formula is used to calculate the distance between two cases [48]:

$$d(x_i, x_j) = \sum_{q \in Q} w_q (x_{iq} - x_{jq})^2 + \sum_{c \in C} w_c L_c(x_{ic}, x_{jc}),$$

where  $Q$  is the set of quantitative features and  $C$  is the set of categorical features,  $L_c$  is an  $M \times M$  symmetric matrix,  $w_q$  is the weight of feature  $q$  and  $w_c$  is the weight of feature  $c$ . In present study, the weights of features are gained by Relief feature selection algorithm, and the number of  $k$  is 5.

As a representative of lazy learning algorithms, KNN has no explicit training process. Therefore, it doesn't have to spend time training the model. Besides, KNN can be easily understood and implemented. These advantages have made KNN widely used in various tasks [49, 50].

## 3 Proposed Framework

The proposed classification system consists of four main components: (1) preprocessing of data, (2) feature selection using Relief algorithm, (3) training of individual classifiers, and (4) prediction result generation of the ensemble classifier. A flow chart of the proposed system is shown in Figure 1. The main components of the system are described in the following subsections.

### 3.1 Data Preprocessing

The aim of data preprocessing is to obtain data from different heart disease data repositories and then process them in the appropriate format for the subsequent analysis [51]. The preprocessing phase involves missing-value imputation and data normalization.

#### 3.1.1 *Missing-value Imputation*

Missing data in medical data sets must be handled carefully because they have a serious effect on the experimental results. Usually, researchers choose to replace the missing values with the mean/mode of the attribute depending on its type

[23]. Mokeddem [51] used weighted KNN to calculate the missing values. In present study, features with missing values more than 50% of all instances are removed, then group mean instead of simple mean are used to substitute remaining missing values, as Bashir et al did in their study [38]. For example, if the case with a missing value is a patient, the mean value for patients is calculated and inserted in place of the missing value. In this way the class label is taken into consideration, thus the information offered by the dataset could be fully utilized.

### 3.1.2 Data Normalization

Before feature selection, the continuous features are normalized to ensure that they have the mean 0 and variance 1, thus the effects of different quantitative units are eliminated.

## 3.2 Feature Selection

The features are selected by the Relief algorithm and the obtained result is a feature rank. A higher ranking means that the feature has stronger distinguishing quality and a higher weight [52]. Features whose relevant statistics are higher than the predefined threshold could be selected. Besides, we can also select the top  $k$  features in the feature rank. In present study, the rejected features are not used for subsequent modules and analysis.

## 3.3 Training of Individual Classifiers

In the training phase, the dataset is randomly split into training set and test set. Then the training set is used to train individual classifiers. In present study, ten-fold cross validation are used to generate the training set and test set, which means that 90% data are used for training and 10% data are used for testing purpose.

## 3.4 Prediction Result Generation

The classification accuracy and misclassification cost (MC) of each classifier are taken into account during the process of generating the final prediction result. In present study, in order to compare the misclassification costs for the different classifiers conveniently, the value of the correct classification cost is set as 0, and the MC is split into two scenarios. In the first scenario, healthy people are diagnosed with heart disease, resulting in unnecessary and costly treatment. In the second scenario, heart disease patients are told that they are healthy, as a result they may miss the best time for treatment, which may cause the disease to deteriorate or even death. The cost matrix is presented in Table 1. Considering the different costs people have to pay for misclassification, we set  $cost_1 = 5$  and  $cost_2 = 1$ . Afterwards, an index  $E$  is constructed to evaluate the performance of each classifier:

$$E_i = \frac{Accuracy_i + 1 - \frac{MC_i}{cost_1 + cost_2}}{2},$$

where  $Accuracy_i$  represents the accuracy and  $MC_i$  represents the MC of  $i$ th classifier during the training phase (the formula to calculate the MC is presented in Section 4.2).  $E_i$  stands for the efficiency of  $i$ th classifier to improve the accuracy



and reduce the MC simultaneously. The weights of individual classifiers are based on  $E_i$  and they are calculated as:

$$w_i = \frac{E_i}{\sum_{i=1}^n E_i},$$

where  $n$  is the number of classifiers. Finally, the instances of the test set are imported into each classifier, and the outputs of ensemble classifier are the labels with the highest weighted vote[53].

## 4 Experimental Results

This section involves the exhibition of experimental results on different heart disease datasets. First, important features selected by Relief algorithm are reported. Then, the performance of individual classifiers and the ensemble classifier are showed. Finally, a comparison is made between the performance of the proposed ensemble model and those of previous studies. The experiment is implemented on MATLAB 2018a platform, and the performance parameters of the executing host were Win 10, Inter (R) 1.80 GHz Core (TM) i5-8250U, X64, and 16 GB (RAM).

### 4.1 Datasets Description

Three different datasets are used in the proposed research, they are Statlog, Cleveland and Hungarian heart disease datasets from UCI machine learning repository. Statlog dataset consists of 270 instances, Cleveland dataset consists of 303 instances and Hungarian dataset consists of 294 instances. The three datasets share the same feature set. The details of feature information are presented in Table 2.

### 4.2 Performance Evaluation Metrics

Various performance metrics are used to evaluate the performance of the classifiers in this study. First, the confusion matrix is introduced in Table 3. In the matrix, the classification result of a two-class problem is divided into four parts: true positive ( TP ), true negative ( TN ), false positive ( FP ) and false negative ( FN ). Based on these error measures, accuracy, MC, G-mean, precision and recall are used to evaluate the performance of different classifiers. G-mean is used because it evaluates the degree of inductive bias which considers both positive and negative accuracy [54]. The metrics are calculated as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\%, \quad (1)$$

$$MC = \frac{FP \times cost_2 + FN \times cost_1}{TP + TN + FP + FN} \times 100\%, \quad (2)$$

$$G - mean = \sqrt{\frac{TP}{TP + FN} \times \frac{TN}{TN + FP}} \times 100\%, \quad (3)$$

$$Precision = \frac{TP}{TP + FP} \times 100\%, \quad (4)$$

$$Recall = \frac{TP}{TP + FN} \times 100\%. \quad (5)$$

Ten-fold cross validation is used to obtain the final results. The ensemble model runs on each test set and processes each instance individually. The evaluation metrics of the ten folds are averaged to verify the superiority of the proposed ensemble classifier.

#### 4.3 Performance on Statlog dataset

Figure 2 shows the the weights of different features on Statlog dataset. Experiments are performed with different numbers of selected features, and the performance of classifiers on 7 features is found to be the best. Therefore, only the performance of classifiers on 7 features is reported. The most important 7 features selected by Relief are given in Table 4.

Table 5 and Figure 3 indicate the comparison of performance evaluation metrics for the proposed ensemble with individual classifiers on Statlog dataset. The proposed ensemble has obtained the highest accuracy of 91.74%, the highest G-mean of 90.55% and the lowest MC of 30.32%. KNN has obtained the highest recall of 89.20%. RF has obtained the highest precision of 92.84%. Overall, the proposed ensemble classifier performs best, and LR performs worst on Statlog dataset. The ROC curve and the area under the curve(AUC) of the proposed ensemble classifier are shown in Fig 4.

#### 4.4 Performance on Cleveland dataset

Figure 5 shows the weights of different features on Cleveland dataset. Experiments are performed with different numbers of selected features, but the performance of classifiers on 7 features is the best, so only the performance of classifiers on 7 features is reported. The most important 7 features selected by Relief are given in Table 6.

Table 7 and Figure 6 indicate the comparison of performance evaluation metrics for the proposed ensemble with individual classifiers on Cleveland dataset. It is clear from the results that the proposed ensemble algorithm has obtained the highest accuracy of 91.26%, the highest G-mean of 90.55% and the highest precision of 95.20%. KNN has obtained the highest recall of 87.89% and the lowest MC of 37.20%. Overall, the proposed ensemble classifier performs best, KNN takes the second place and LR performs worst on Cleveland dataset. The ROC curve and the AUC of the proposed ensemble classifier are shown in Fig 7.

#### 4.5 Performance on Hungarian dataset

For Hungarian dataset, Slope, Ca and Thal are deleted during the process of missing-value imputation because these features have missing values more than 50% of all instances. Therefore, only ten features can enter into the feature selection

module. Figure 8 shows the weights of different features on Hungarian dataset. Experiments are performed with different numbers of selected features, but the performance of classifiers on 5 features is the best, so only the performance of classifiers on 5 features is reported. The most important 5 features selected by Relief are given in Table 8.

Table 9 and Figure 9 indicate the comparison of performance evaluation metrics for the proposed ensemble with individual classifiers on Hungarian dataset. The proposed ensemble algorithm has obtained the highest accuracy of 90.79%, the highest precision of 96.11%, the highest recall of 89.61% and the lowest MC of 39.77%. RF has obtained the highest G-mean of 92.65%. Overall, the proposed ensemble classifier performs best, RF takes the second place and LR performs worst. The ROC curve and the AUC of the proposed ensemble classifier are shown in Fig 10.

#### 4.6 Comparison of the Results with Other Studies

Table 10 shows the classification accuracies of our study and previous methods. The results show that our proposed method obtains superior and promising results in classifying heart disease patients. In addition, most previous studies did not take different kinds of misclassification costs into consideration, and the limitation is conquered in present study. Thus, we believe that the proposed model can be beneficial in aiding physicians in making better decisions.

## 5 Conclusions and Future Works

In this study, a cost-sensitive ensemble model based on five different classifiers is presented to assist the diagnosis of heart disease. The Statlog heart disease dataset, Cleveland heart disease dataset and Hungarian heart disease dataset are selected to test the model. The performance of classifiers are presented using different parameters such as accuracy, MC, G-mean, precision and recall. The significant results achieved by the proposed ensemble model are as follows:

- (1) The highest accuracy achieved by the ensemble model is 91.74%, which is a significant improvement for accurately diagnosing the heart disease patients. The individual classifiers have not achieved such accuracy.
- (2) The highest G-mean of 90.55%, highest precision of 96.11%, highest recall of 89.61% and lowest MC of 30.32% are achieved by the ensemble model, which indicate that the proposed ensemble model has reduced the MC, thus relieved the pain patients had to suffer during the process of diagnosis.

Kononenko [55] applied various machine learning techniques and compared the performance on eight medical datasets using five different parameters: performance, transparency, explanation, reduction, and missing data handling. While individual classifiers have shortcomings on some of these aspects, the ensemble model is able to overcome their deficiencies. For example, RF can generate explicit rules for decision making, and the basic idea of KNN is "to solve new problems by identifying and reusing previous similar cases based on the heuristic principle that similar problems have a high likelihood of having similar solutions" [56], which is easily understood by physicians. On the other hand, LR, SVM and ELM are more like a "black box", and physicians are willing to accept a "black box" classifier only when it outperforms a very large margin all other classifiers, including the physicians themselves,

but such situation is highly improbable [55]. In addition, KNN is a lazy evaluation method while the other four are eager evaluation methods. Eager algorithm generates frequent itemset rules from a given data set and predicts a class for test instance based on multicriteria approach from selected frequent itemset rules [23]. If no matching is found, default prediction (i.e., the most frequent class in data set) is performed, which may not be correct. In contrast, lazy algorithm uses a richer hypothesis space, it makes judgment according to a small proportion of the instances in the database, thus overcomes the limitation of eager algorithms. However, lazy algorithm uses more time for prediction, as multicriteria matching is performed for each instance in data set [57], while eager algorithm is able to generate the prediction results at a very fast speed after the training phase. From the above discussion, it can be concluded that the selected classifiers complement each other very well. In any scenario where one classifier has some limitations, the other classifier overcome them. As a result, better performance is achieved.

Moreover, the present study takes MC into consideration and tries to reduce it. Most traditional algorithms focus only on the classification accuracy, ignoring the cost patients have to pay for misclassification. But the diagnostic mistakes are of higher importance in the medical field, and the price of a false negative instance is clearly much higher than that of a false positive one. Aiming at this problem, the present study has adopted a new method to combine the prediction results of heterogeneous classifiers and significantly reduced the MC.

There are limitations within this study that constrain the improvement of the results. For instance, all the selected classifiers used supervised learning algorithms. Therefore, they often made similar mistakes during the prediction phase, which reduced the accuracy of the ensemble model. Besides, the fitness functions used by classifiers in training phase still aimed only on improving the accuracy, thus the individual classifiers were not cost-sensitive. In the future, new fitness functions and learning algorithms, such as unsupervised and semi-supervised algorithms can be incorporated into the proposed ensemble classifier to improve its performance.

## Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Funding

This study was not funded.

Competing interests

The authors declare that they have no competing interests.

Author's contributions

Qi and Zhang designed research, performed research, analyzed data, and wrote the paper.

## Acknowledgements

The authors acknowledge the editor and anonymous reviewers for their supportive works and insightful comments.

## Availability of data and materials

The data used in this study is available in UCI Machine Learning Repository.

## Author details

<sup>1</sup>College of Management and Economics, Tianjin University, Nankai District, 300072 Tianjin, PR China. <sup>2</sup>School of Mathematical Science, Hebei Normal University, Yuhua District, 050024 Shijiazhuang, PR China.

## References

- Heart disease. <http://health.allrefer.com/health/heart-disease-info.html> Accessed:17.04.06
- World Heart Federation Report. <http://www.world-heart-federation.org/Accessed:01.12.16>
- for Cardiovascular Diseases, N.C.: The epidemic of heart disease. Encyclopedia of China Publishing House (2019)
- Lopez-Sendon, J.: The heart failure epidemic. *Medicographia* **33**(2), 363–369 (2011)
- Amato, F., Lopez, A., Pena-Mendez, E.M., Vanhara, P., Hampl, A., Havel, J.: Artificial neural networks in medical diagnosis. *Journal of Applied Biomedicine* **11**(2), 47–58 (2013)
- Xu, M., Shen, J.: Information sharing system for heart disease emergence treatment based on an information fusion model. *Industrial Engineering Journal* **12**(4), 61–66 (2009)
- Giri, D., Acharya, U.R., Martis, R.J., Sree, S.V., Lim, T.C., Thajudin Ahamed, V.I., Suri, J.S.: Automated diagnosis of coronary artery disease affected patients using lda, pca, ica and discrete wavelet transform. *Knowledge-Based Systems* **37**(2), 274–282 (2013)
- Safdar, S., Zafar, S., Zafar, N., Khan, N.F.: Machine learning based decision support systems (dss) for heart disease diagnosis: a review. *Artificial Intelligence Review* **2017**, 1–27 (2017)
- U Rajendra, A., Oliver, F., Vinitha, S., Swapna, . G., Roshan Joy, M., Nahrizul Adib, K., Suri, J.S.: Linear and nonlinear analysis of normal and cad-affected heart rate signals. *Computer Methods & Programs in Biomedicine* **113**(1), 55–68 (2014)
- Mejia, O.A.V., Antunes, M.J., Goncharov, M., Dallan, L.R.P., Veronese, E., Lapenna, G.A., Lisboa, L.A.F., Dallan, L.A.O., Brandao, C.M.A., Zubelli, J., Tarasoutchi, F., Pomerantzeff, P.M.A., Jatene, F.B.: Predictive performance of six mortality risk scores and the development of a novel model in a prospective cohort of patients undergoing valve surgery secondary to rheumatic fever. *PLoS ONE* **2018**, 1–14 (2018)
- Lukacs Krogager, M., Skals, R.K., Appel, E.V.R., Schnurr, T.M., Engelbrechtsen, L., Have, C.T., Pedersen, O., Engstrom, T., Roden, D.M., Gislason, G., Poulsen, H.E., Kober, L., Stender, S., Hansen, T., Grarup, N., Andersson, C., Torp-Pedersen, C., Weeke, P.E.: Hypertension genetic risk score is associated with burden of coronary heart disease among patients referred for coronary angiography. *PLoS One* **13**(12), 1–17 (2018)
- Tomar, D., Agarwal, S.: Feature selection based least square twin support vector machine for diagnosis of heart disease. *International Journal of Bio-Science and Bio-Technology* **6**, 69–82 (2014)
- Subbulakshmi, C.V., Deepa, S.N.: Medical dataset classification: A machine learning paradigm integrating particle swarm optimization with extreme learning machine classifier. *The scientific world journal* **2015**, 1–12 (2015)
- Jabbar, M.A., Deekshatulu, Chandra, P.: Heart disease classification using nearest neighbor classifier with feature subset selection. *Computer Science & Telecommunications* **2**, 47–54 (2013)
- Dogan, M.V., Grumbach, I.M., Michaelson, J.J., Philibert, R.A.: Integrated genetic and epigenetic prediction of coronary heart disease in the framingham heart study. *Plos One* **13**(1), 1–18 (2018)
- Detrano, R., Janosi, A., Steinbrunn, W., Pfisterer, M., Schmid, J.J., Sandhu, S., Guppy, K.H., Lee, S., Froelicher, V.: International application of a new probability algorithm for the diagnosis of coronary artery disease. *American Journal of Cardiology* **64**(5), 304–310 (1989)
- Gokulnath, C.B., Shantharajah, S.P.: An optimized feature selection based on genetic approach and support vector machine for heart disease. *Cluster Computing* (4), 1–11 (2018)
- Subbulakshmi, C.V., Deepa, S.N., Malathi, N.: Extreme learning machine for two category data classification. In: *IEEE International Conference on Advanced Communication Control & Computing Technologies* (2012)
- Duch, W., Adamczak, R., K., G.: A new methodology of extraction, optimization and application of crisp and fuzzy logical rules. *IEEE Transactions on Neural Networks* **12**(2), 277–306 (2001)
- Yingsang, L.O., Fujita, H., Pai, T.: Prediction of coronary artery disease based on ensemble learning approaches and co-expressed observations. *Journal of Mechanics in Medicine & Biology* **16**(01), 1–10 (2016)
- Eom, J.H., Kim, S.C., Zhang, B.T.: Aptacdss-e: A classifier ensemble-based clinical decision support system for cardiovascular disease level prediction. *Expert Systems with Applications* **34**(4), 2465–2479 (2008)
- Das, R., Turkoglu, I., Sengur, A.: Effective diagnosis of heart disease through neural networks ensembles. *Expert Systems with Applications* **36**(4), 7675–7680 (2009)
- Bashir, S., Qamar, U., Khan, F.H.: A multicriteria weighted vote-based classifier ensemble for heart disease prediction. *Computational Intelligence* **32**(4), 615–645 (2016)
- Khened, M., Kollerathu, V.A., Krishnamurthi, G.: Fully convolutional multi-scale residual densenets for cardiac segmentation and automated cardiac diagnosis using ensemble of classifiers. *Medical Image Analysis* **51**, 21–45 (2018)
- Krawczyk, B., Schaefer, G., Wozniak, M.: A hybrid cost-sensitive ensemble for imbalanced breast thermogram classification. *Artificial Intelligence in Medicine* **65**(3), 219–227 (2015)
- Liu, N., Shen, J., Xu, M., Gan, D., Qi, E.S.: Improved cost-sensitive support vector machine classifier for breast cancer diagnosis. *Mathematical Problems in Engineering* **4**, 1–13 (2018)

27. Wei, Z., Junjie, C.: Relief feature selection and parameter optimization for support vector machine based on mixed kernel function. *International Journal of Performability Engineering* **14**(2), 280–289 (2018)
28. Ul Haq, A., Jian Ping, L., Memon, M.H., Nazir, S., Sun, R.: A hybrid intelligent system framework for the prediction of heart disease using machine learning algorithms. *Mobile Information Systems* **2018**, 1–21 (2018)
29. Urbanowicz, R.J., Meeke, M., Lacava, W., Olson, R.S., Moore, J.H.: Relief-based feature selection: Introduction and review. *Journal of Biomedical Informatics* **85**, 189–203 (2018)
30. Breiman, L.: Random forest. *Machine Learning* **45**, 5–32 (2001)
31. Hajjalian, H., Toma, C.: Network anomaly detection by means of machine learning: Random forest approach with apache spark. *Informatica Economica* **22**(4), 89–98 (2018)
32. Eccel, E., Ghielmi, L., Granitto, P., Barbiero, R., Grazzini, F., Cesari, D.: Prediction of minimum temperatures in an alpine region by linear and non-linear post-processing of meteorological models. *Nonlinear Processes in Geophysics* **14**(3), 211–222 (2007)
33. Whitrow, C., Hand, D.J., Juszczak, P., Weston, D., Adams, N.M.: Transaction aggregation as a strategy for credit card fraud detection. *Data Mining & Knowledge Discovery* **18**(1), 30–55 (2009)
34. Kaya, G.O.: A hybrid method based on empirical mode decomposition and random forest regression for wind power forecasting. *Journal of Multiple-Valued Logic & Soft Computing* **31**(1/2), 123–137 (2018)
35. Larsen, K., Petersen, J.H., Budtz-Jorgensen, E., Endahl, L.: Interpreting parameters in the logistic regression model with random effects. *Biometrics* **56**(3), 909–914 (2015)
36. Cortes, C., Vapnik, V.: Support-vector networks. *Machine Learning* **20**(3), 273–297 (1995)
37. Davari, D.A., Khadem, S.E., Asl, B.M.: Automated diagnosis of coronary artery disease (cad) patients using optimized svm. *Computer Methods & Programs in Biomedicine* **138**, 117–126 (2017)
38. Bashir, S., Qamar, U., Khan, F.H.: Bagmoov: A novel ensemble for heart disease prediction bootstrap aggregation with multi-objective optimized voting. *Australas Phys Eng Sci Med* **38**(2), 305–323 (2015)
39. Ghumbre, S., Patil, C., Ghatol, A.: Heart disease diagnosis using support vector machine. In: *International Conference on Computer Science and Information Technology (ICCSIT)*, Pattaya, Thailand (2011)
40. Huang, G.B., Zhu, Q.Y., Siew, C.K.: Extreme learning machine: Theory and applications. *Neurocomputing* **70**(1), 489–501 (2006)
41. Huang, G.B., Wang, D.H., Lan, Y.: Extreme learning machines: a survey. *International Journal of Machine Learning & Cybernetics* **2**(2), 107–122 (2011)
42. Wu, D., Qu, Z.S., Guo, F.J., Zhu, X.L., Wan, Q.: Hybrid intelligent deep kernel incremental extreme learning machine based on differential evolution and multiple population grey wolf optimization methods. *AUTOMATIKA* **60**(1), 48–57 (2019)
43. Huang, G.B., Wang, D.: Advances in extreme learning machines (elm2010). *Neurocomputing* **128**(16), 1–3 (2014)
44. Wentao, Z., Pan, L., Qiang, L., Dan, L., Xinwang, L.: Selecting the optimal hidden layer of extreme learning machine using multiple kernel learning. *KSII Transactions on Internet & Information Systems* **12**(12), 5765–5781 (2018)
45. Nahato, K.B., Nehemiah, K.H., Kannan, A.: Hybrid approach using fuzzy sets and extreme learning machine for classifying clinical datasets. *Informatics in Medicine Unlocked* **2**, 1–11 (2016)
46. Raghuvanshi, B.S., Shukla, S.: Class-specific kernelized extreme learning machine for binary class imbalance learning. *Applied Soft Computing* **73**, 1026–1038 (2018)
47. Wang, X., Li, H., Zhang, Q., Wang, R.: Predicting subcellular localization of apoptosis proteins combining go features of homologous proteins and distance weighted knn classifier. *BioMed Research International* **2016**(2), 1–8 (2016)
48. Uguroglu, S., Carbonell, J., Doyle, M., Biederman, R.: Cost-sensitive risk stratification in the diagnosis of heart disease. In: *Twenty-sixth Aaai Conference on Artificial Intelligence* (2012)
49. Qiao, S., Yan, B., Jing, L.: Ensemble learning for protein multiplex subcellular localization prediction based on weighted knn with different features. *Applied Intelligence* (1), 1–12 (2017)
50. Cai, Z., Gu, J., Wen, C., Dong, Z., Huang, C., Hui, H., Tong, C., Li, J., Chen, H.: An intelligent parkinson's disease diagnostic system based on a chaotic bacterial foraging optimization enhanced fuzzy knn approach. *Computational & Mathematical Methods in Medicine* **2018**(3), 1–24 (2018)
51. Mokeddem, S.A.: A fuzzy classification model for myocardial infarction risk assessment. *Applied Intelligence* (12), 1–18 (2017)
52. Zhang, L.X., Wang, J.X., Zhao, Y.N., Yang, Z.H.: A novel hybrid feature selection algorithm: using relief estimation for ga-wrapper search. In: *International Conference on Machine Learning & Cybernetics* (2004)
53. Saha, S., Ekbal, A.: Combining multiple classifiers using vote based classifier ensemble technique for named entity recognition. *Data & Knowledge Engineering* **85**(8), 15–39 (2013)
54. Chong Zhang, H.L., Kay Chen Tan, Hong, G.S.: A cost-sensitive deep belief network for imbalanced classification. *arXiv* (2018)
55. Kononenko, I.: Machine learning for medical diagnosis: history, state of the art and perspective. *Artificial Intelligence in Medicine* **23**(1), 89–109 (2001)
56. Ahmed, M.U., Begum, S., Olsson, E., Ning, X., Funk, P.: Case-based reasoning for medical and industrial decision support systems. Springer (2010)
57. Houeland, T.G., Aamodt, A.: An efficient hybrid classification algorithm - an example from palliative care. Springer **6679**, 197–204 (2011)
58. Kahramanli, H., Allahverdi, N.: Design of a hybrid system for the diabetes and heart diseases. *Expert Systems with Applications* **35**(1-2), 82–89 (2008)
59. Lee, S.H.: Feature selection based on the center of gravity of bswfms using newfm. *Engineering Applications of Artificial Intelligence* **45**, 482–487 (2015)
60. Karegowda, A.G., Jayaram, M.A., Manjunath, A.S.: Feature subset selection problem using wrapper approach in supervised learning. *International Journal of Computer Applications* **1**(7), 13–17 (2010)

## Figures

[scale=1]flow

**Figure 1** Flowchart of the proposed ensemble classifier

[height=7cm,width=14cm]statlog

**Figure 2** Feature weights on Statlog dataset

[height=8cm,width=14cm]sta

**Figure 3** Performance comparison on Statlog dataset

[height=8cm,width=14cm]roc2

**Figure 4** ROC curve of proposed ensemble classifier on Statlog dataset



[scale=0.8]cleveland

**Figure 5** Feature weights on Cleveland dataset

[height=8cm,width=14cm]cleve

**Figure 6** Performance comparison on Cleveland dataset

[height=8cm,width=12cm]roc1

**Figure 7** ROC curve of proposed ensemble classifier on Cleveland dataset

[height=8cm,width=12cm]hungarian

**Figure 8** Feature weights on Hungarian dataset

[height=8cm,width=14cm]hun

**Figure 9** Performance comparison on Hungarian dataset

[height=8cm,width=12cm]roc3

**Figure 10** ROC curve of proposed ensemble classifier on Hungarian dataset

## Tables

**Table 1** The cost matrix used by the classifiers

Predicted	Reality	
	sick	healthy
sick	0	$cost_2$
healthy	$cost_1$	0

**Table 2** Features of heart disease datasets

Feature	Description	Value
Age	Age in years	Continuous value
Sex	Gender	1 : male;0 : female
Cp	Chest Pain Type	1: typical angina 2: atypical angina 3: non-anginal pain 4: asymptomatic
Trestbps	Resting Blood Sugar	Continuous value in mm hg
Chol	Serum Cholestorol	Continuous value in mm/dl
Fbs	Fasting Blood Sugar	0 :< 120mg/dl 1 :> 120mg/dl 0 : normal
Restecg	Resting ECG Results	1 : having ST-T wave abnormality 2 : probable or definite left ventricular hypertrophy
Thalach	Maximum heart rate achieved	Continuous value
Exang	Exercise induced angina	0 :no 1 :yes
Oldpeak	ST depression induced by exercise relative to rest	Continuous value
Slope	Slope of the peak exercise ST segment	1 =upsloping 2 =flat 3 =downsloping
Ca	Number of major vessels colored by flourosopy	0,1,2,3
Thal	Heart beat	3 :normal 6 :fixed defect 7 :reversable defect
Num	Predicted Class	0, 1

**Table 3** The confusion matrix

Predicted	Actual	
	Positive	Negative
Positive	TP	FP
Negative	FN	TN

**Table 4** Features selected by Relief and their ranking on Statlog dataset

Ranking	Feature	Weight
1	Cp	0.18
2	Ca	0.1581
3	Sex	0.1522
4	Thal	0.1323
5	Exang	0.093
6	Slope	0.09
7	Restecg	0.0815

**Table 5** Experimental results on Statlog dataset

Statlog dataset	Accuracy(%)	MC(%)	G-mean(%)	Precision(%)	Recall(%)
RF	90.69	34.67	90.09	92.84	85.84
LR	84.60	54.00	83.80	85.99	78.41
SVM	90.18	35.39	89.61	91.79	85.76
ELM	90.13	36.48	89.50	92.25	85.18
KNN	88.99	33.32	88.92	86.96	89.20
Proposed ensemble	91.74	30.32	90.21	92.21	86.57

**Table 6** Features selected by Relief and their ranking on Cleveland dataset

Ranking	Feature	Weight
1	Cp	0.2361
2	Ca	0.1883
3	Thal	0.1106
4	Sex	0.1056
5	Slope	0.1043
6	Exang	0.0931
7	Restecg	0.05688

**Table 7** Experimental results on Cleveland dataset

Cleveland dataset	Accuracy(%)	MC(%)	G-mean(%)	Precision(%)	Recall(%)
RF	89.58	41.67	89.84	94.32	85.13
LR	81.93	76.69	81.84	91.68	71.92
SVM	88.86	41.83	89.08	92.64	85.43
ELM	88.92	42.70	89.17	93.19	84.96
KNN	88.11	37.20	87.91	89.47	87.89
Proposed ensemble	91.26	39.90	90.55	95.20	85.70

**Table 8** Features selected by Relief and their ranking on Hungarian dataset

Ranking	Feature	Weight
1	Cp	0.3648
2	Exang	0.1743
3	Sex	0.1429
4	Oldpeak	0.1314
5	Trestbps	0.0624

**Table 9** Experimental results on Hungarian dataset

Hungarian dataset	Accuracy(%)	MC(%)	G-mean(%)	Precision(%)	Recall(%)
RF	88.31	55.98	92.65	95.78	87.21
LR	75.86	100.07	78.01	89.32	73.47
SVM	81.56	85.72	85.79	94.16	83.44
ELM	83.16	75.04	85.38	92.17	81.81
KNN	84.12	71.25	86.39	93.02	83.85
Proposed ensemble	90.79	39.77	89.27	96.11	89.61

**Table 10** Comparison of our results with those of other studies

Author	Method	Classification accuracy(%)
Our study	Ensemble classifier	91.74
Tomar and Agarwal[12]	Feature selection-based LSTSVM	85.59
Kahramanli and Allahverdi[58]	Hybrid neural network	86.8
Subbulakshmi et al[18]	ELM	87.5
Lee[59]	Graphical characteristics of BSWFM combined with Euclidean distance	87.4
Das et al.[22]	Neural networks ensemble	89.01
Karegowda et al[60]	GA + Naive Bayes	85.87
Duch et al[19]	KNN	85.6

LSTSVM: Least Square Twin Support Vector Machine; BSWFM: bounded sum of weighted fuzzy membership functions; GA: genetic algorithm.

# Figures

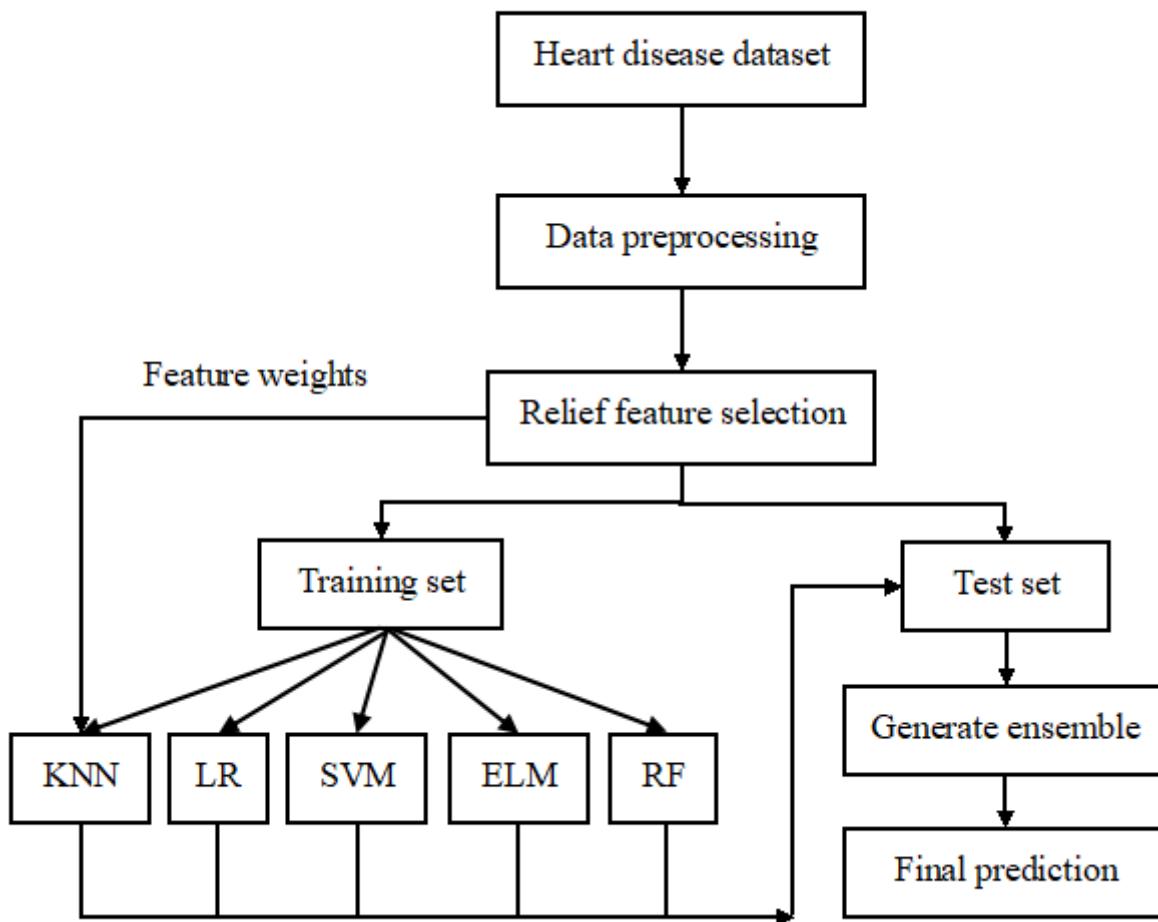
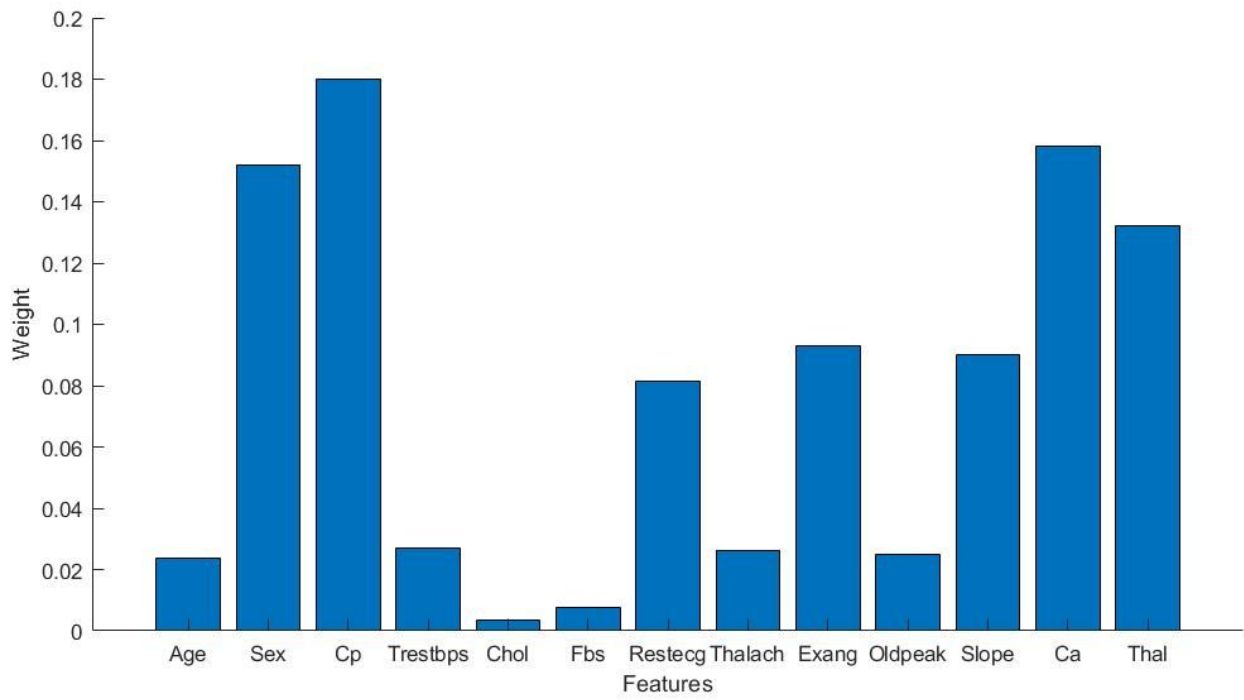


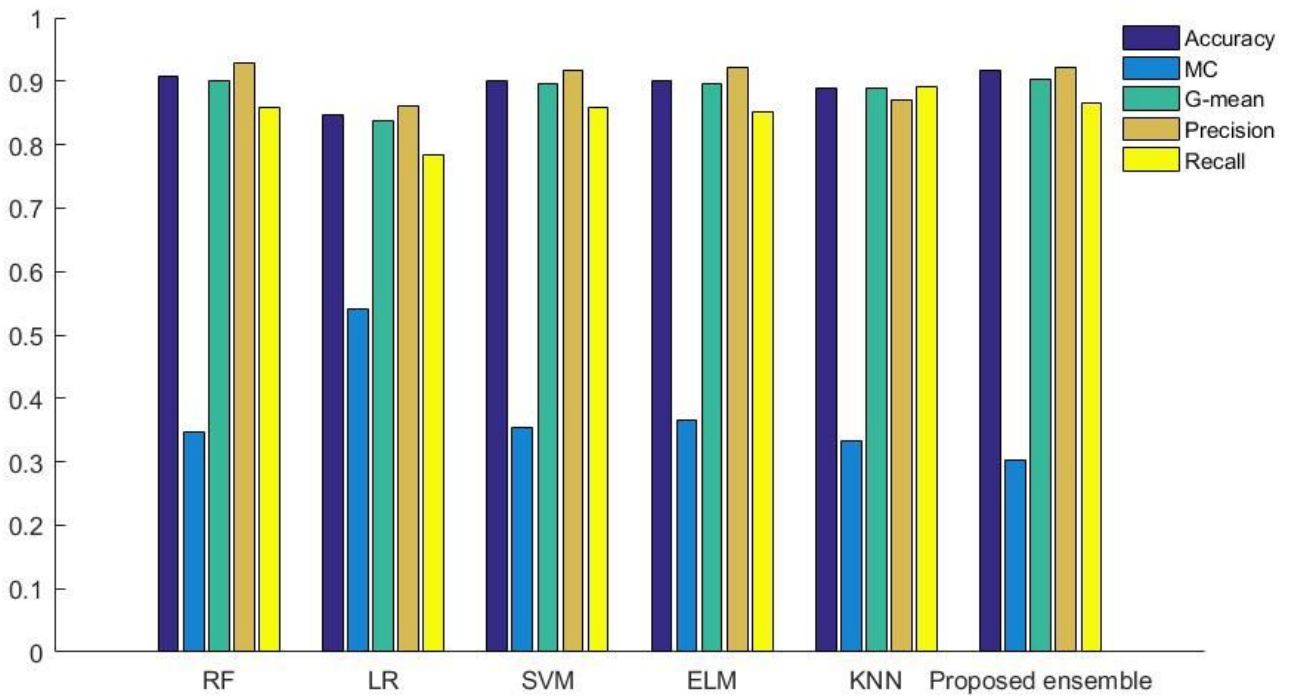
Figure 1

Flowchart of the proposed ensemble classifier



**Figure 2**

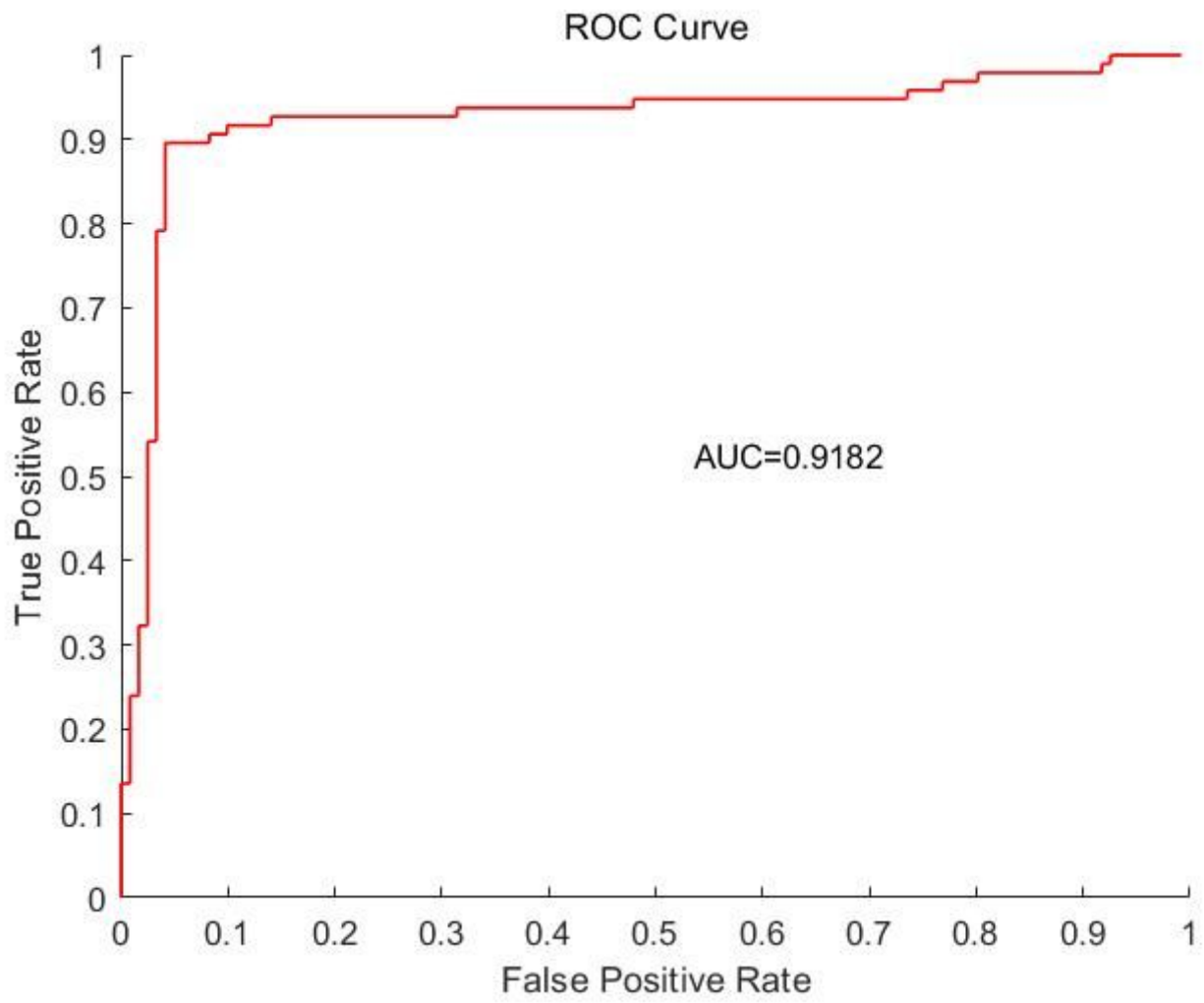
Feature weights on Statlog dataset



**Figure 3**

Performance comparison on Statlog dataset





**Figure 4**

ROC curve of proposed ensemble classifier on Statlog dataset

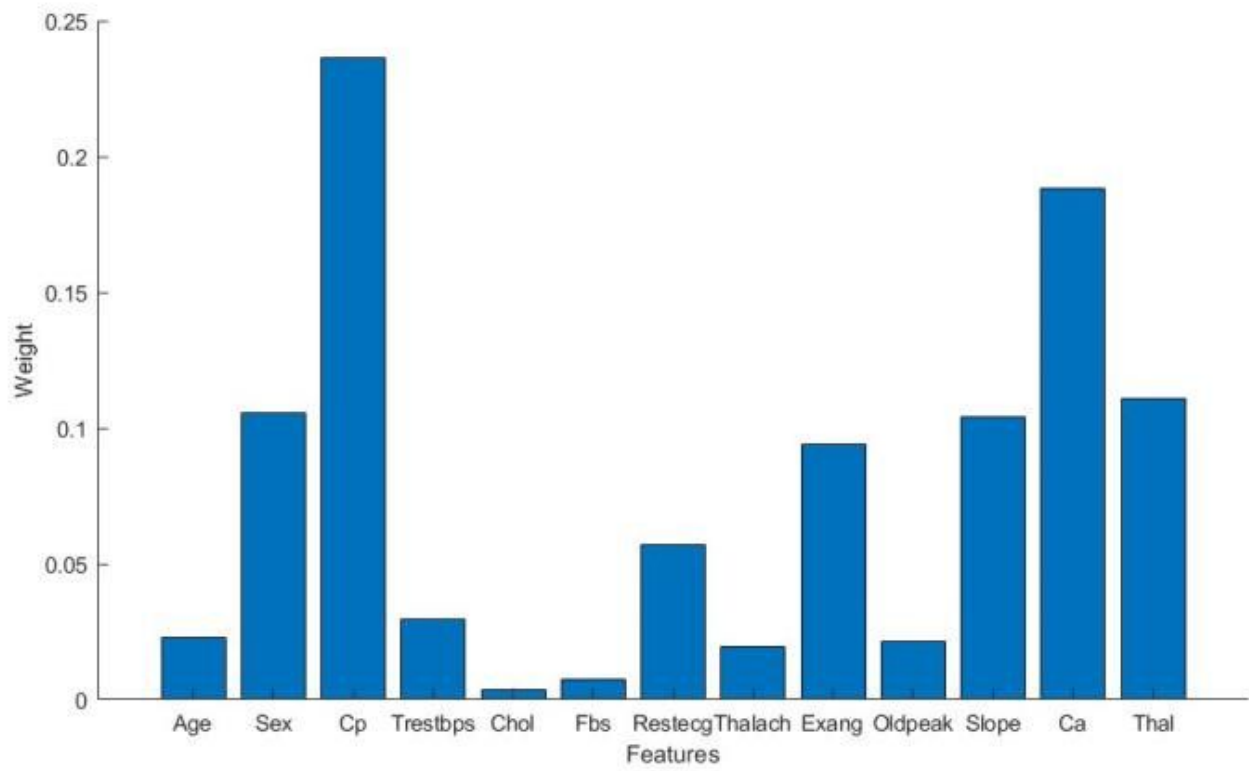


Figure 5

Feature weights on Cleveland dataset

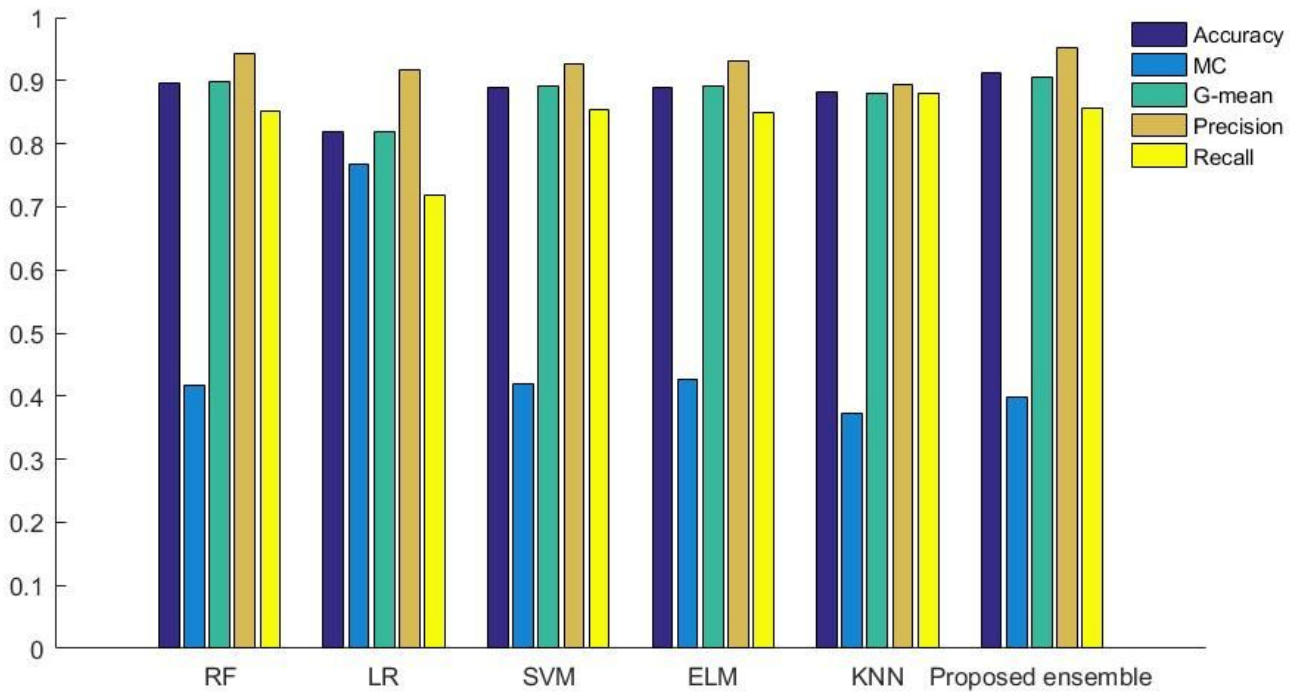


Figure 6

Performance comparison on Cleveland dataset

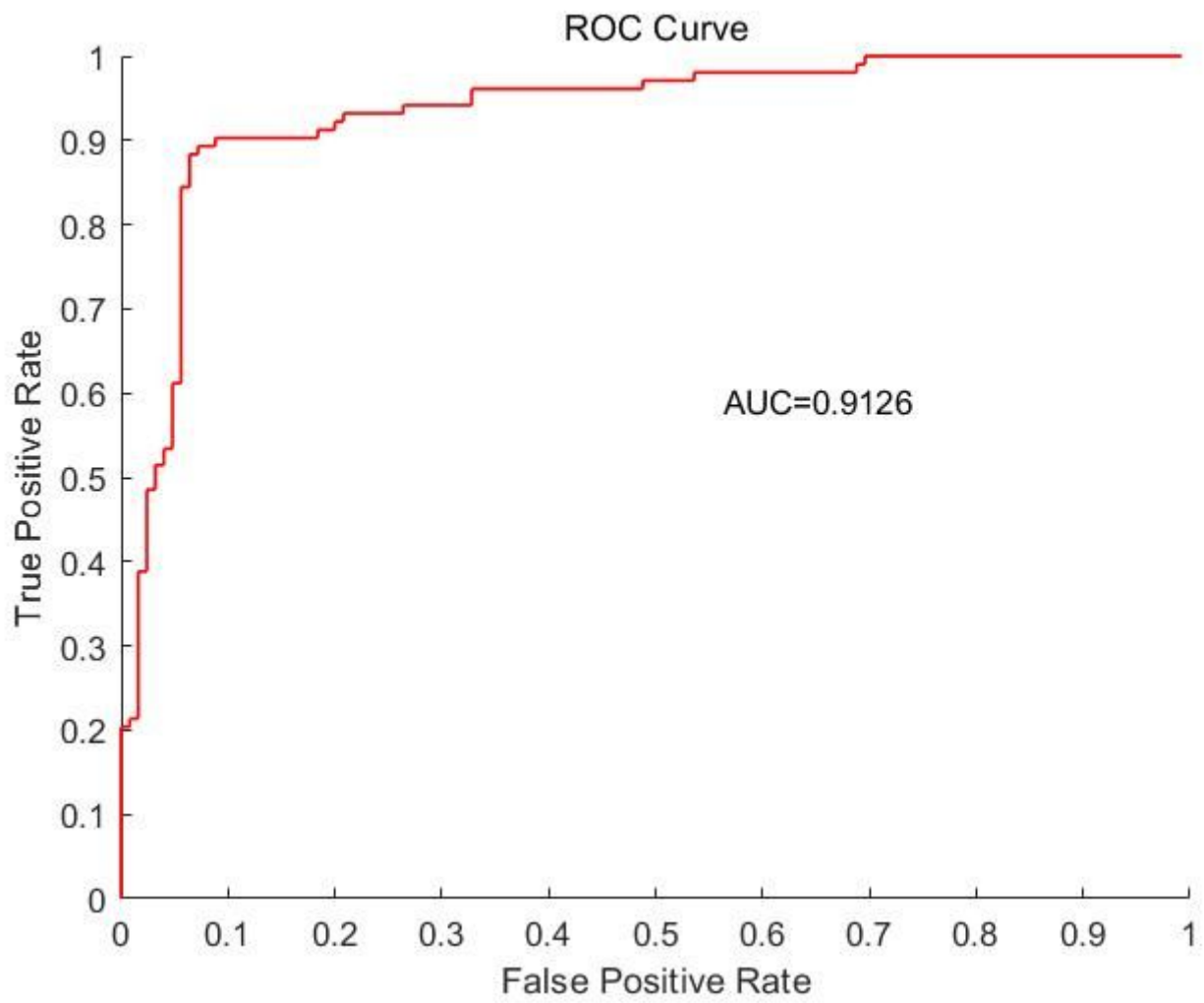


Figure 7

ROC curve of proposed ensemble classifier on Cleveland dataset

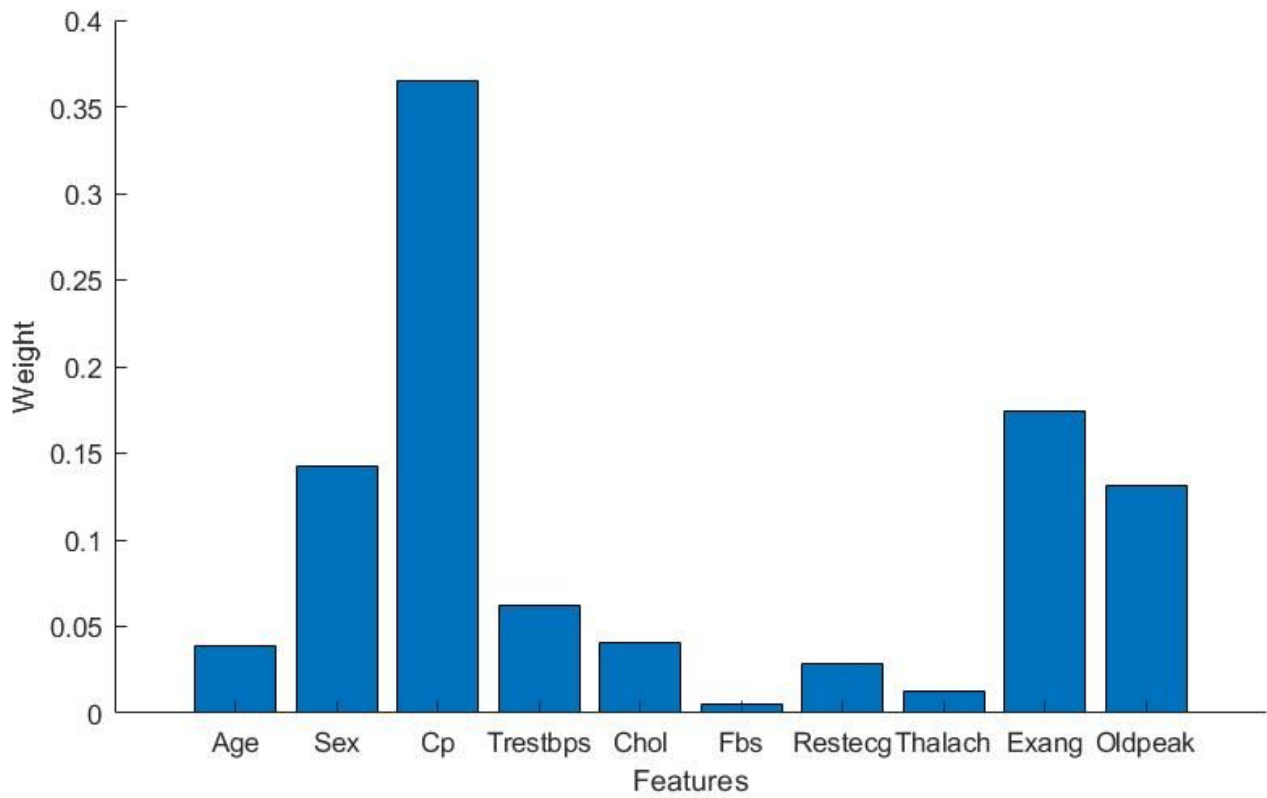


Figure 8

Feature weights on Hungarian dataset

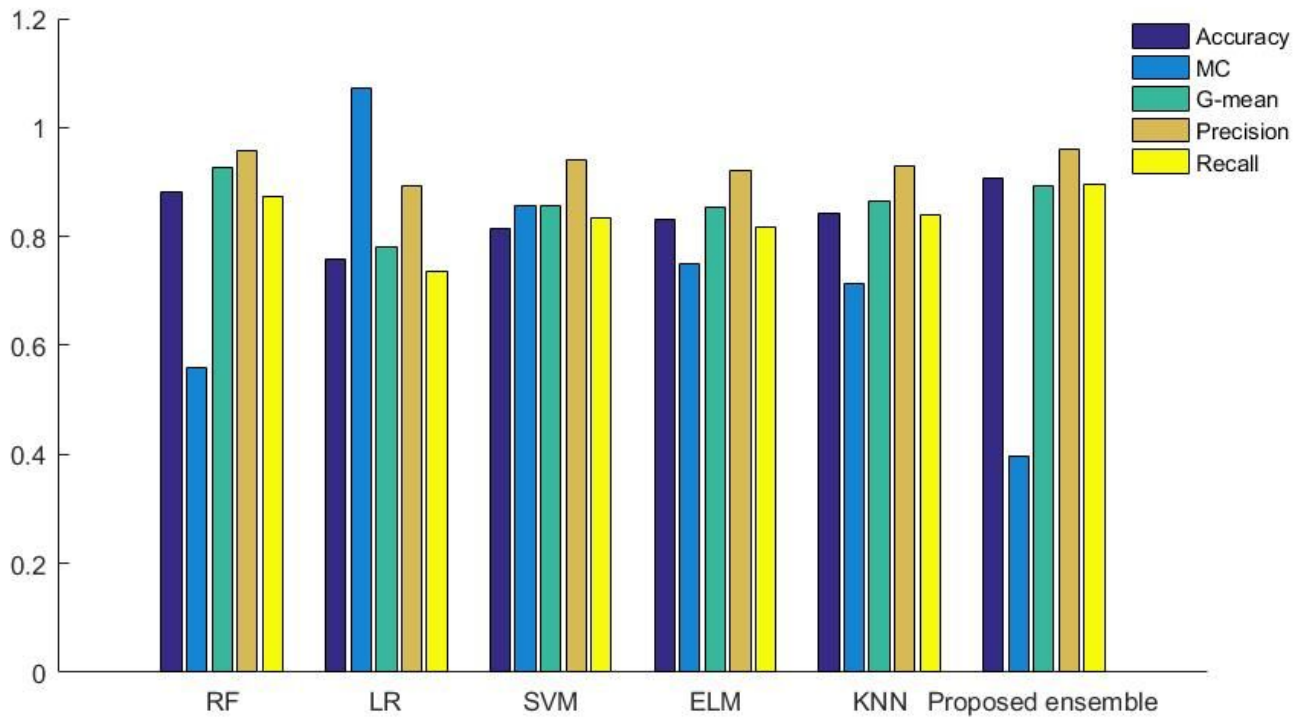


Figure 9

Performance comparison on Hungarian dataset

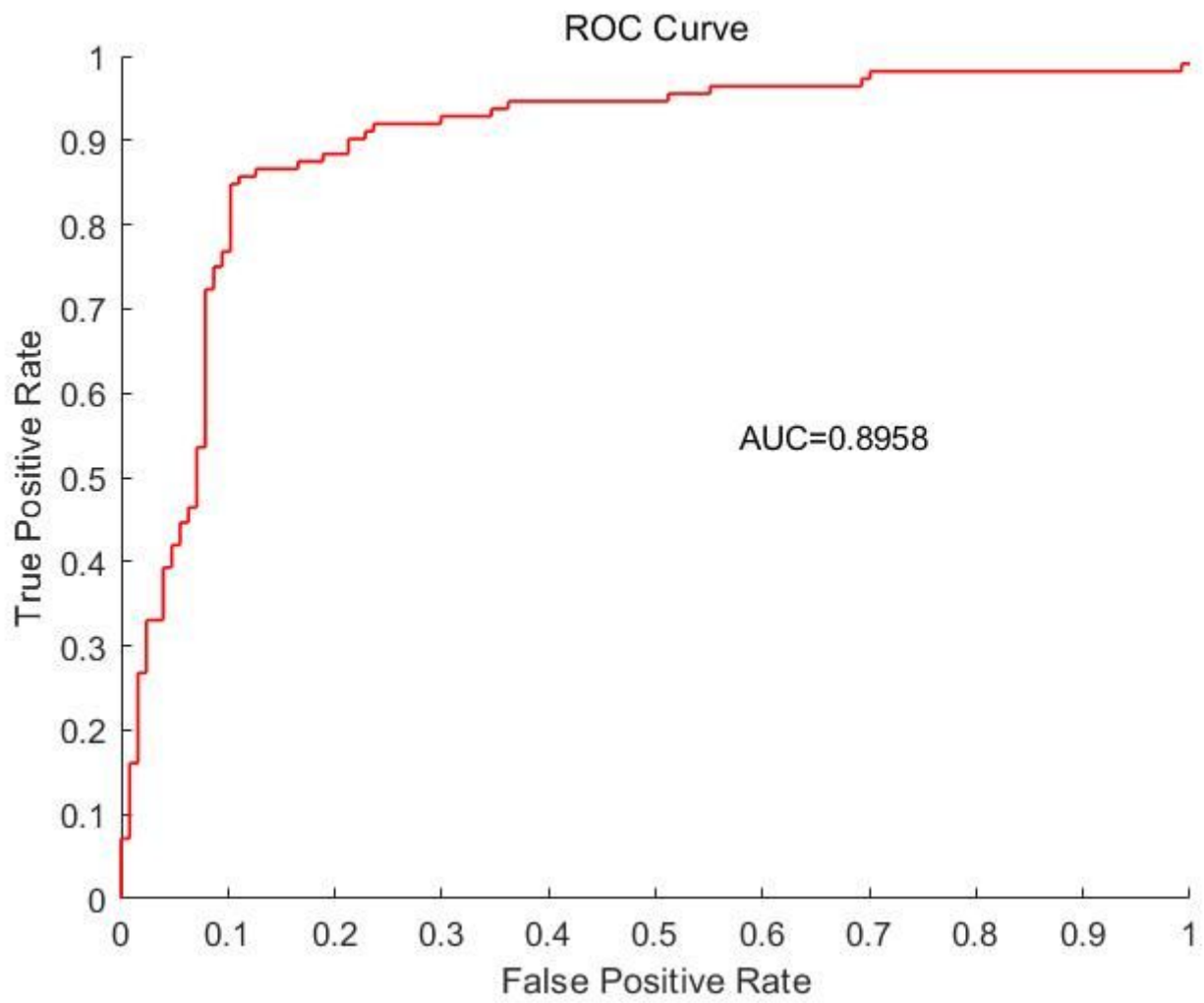


Figure 10

ROC curve of proposed ensemble classifier on Hungarian dataset