**WILEY** | Hindawi

*Research Article*

# Improving the Accuracy for Analyzing Heart Diseases Prediction Based on the Ensemble Method

**Xiao-Yan Gao,[1] Abdelmegeid Amin Ali,[2] Hassan Shaban Hassan,[2] and Eman M. Anwar** (iD) [3,4]

[1]*School of Mathematics and Statistics, Yulin University, Yulin 719000, China*
[2]*Faculty of Computers and Information, Department of Computer Science, Minia University, Minya, Egypt*
[3]*Faculty of Computers and Information, Department of Information System, Minia University, Minya, Egypt*
[4]*Al-Obour High Institute for Management and Informatics, Obour, Egypt*

Correspondence should be addressed to Eman M. Anwar; emyanwar616@gmail.com

Heart disease is the deadliest disease and one of leading causes of death worldwide. Machine learning is playing an essential role in the medical side. In this paper, ensemble learning methods are used to enhance the performance of predicting heart disease. Two features of extraction methods: linear discriminant analysis (LDA) and principal component analysis (PCA), are used to select essential features from the dataset. The comparison between machine learning algorithms and ensemble learning methods is applied to selected features. The different methods are used to evaluate models: accuracy, recall, precision, F-measure, and ROC. The results show the bagging ensemble learning method with decision tree has achieved the best performance.

## 1. Introduction

Nowadays, the cardiac disease is one of the most critical problems relating to human safety. The treatment of heart problems has recently been stated in a study that has received huge attention in the medical system worldwide. Cardiac diseases are one of the most principal causes of death worldwide. On median, 17.7 million deaths result from heart disease which counts for about 31% throughout the world in 2016, according to World Health Organization (WHO) [1]. The cardiac cases number, as the focus of this study, shows that 82% of the cases are from low and middle countries, 17 million are under 70 years of age and prone to noninfectious diseases, 6.7 million are affected by stroke, and 7.4 million people are suffering from heart disease (WHO, 2016) [2]. In the US and other developed countries, about half of all deaths are caused by heart disease; also, one-third of all people's deaths worldwide are related to heart disease. Cardiac disease affects not just people's health but the economies and costs of countries as well. The most common cardiac disorders are those of microvascular origin, primarily cardiac disorders and stroke. After several years of exposure to unhealthy lifestyles, cardiovascular disease clinically presents itself in early stages of life, as well as at an old age. The main cardiac medical conditions include overweight, diabetes, family history, smoking, and high cholesterol [3].

To examine the cardiac disease mischance, the particular issues which need to be discussed are those related to the behaviors. Furthermore, patients will undergo extensive examinations, such as blood pressure, glucose, vital signs, chest pain, electrocardiograms, maximum heart rate, and elevated levels of sugar, but the bright side may be that successful treatment is feasible if the disease is easily and early detected and anticipated, but treatment for all of these cardiac patients is depending on clinical studies, the patient history, and the responses to questions by the patient [4]. All of these techniques (history analysis, physical examination research, and medical professional evaluates) often cause inaccurate diagnosis and mechanical failure besides delaying the diagnosis tests. In addition, it is also more expensive and computation intensive, and it takes a lot of time for evaluations to be carried out [5].

Determining the probability of having cardiac disease manually is hard to depend on as risk factors. Recently, to solve difficult issues, a range of data mining techniques and

machine learning techniques are built [6, 7]. Still, more advanced machine learning will assist us to identify patterns and their useful knowledge. While it has several uses in the medical field, machine learning is mainly utilized to forecast the heart disease. In order to diagnose diseases, many researchers have been interested in utilizing machine learning because it helps minimize diagnostic time and demonstrates accuracy and effectiveness. Using machine learning techniques, as a matter of fact, several diseases can be identified, but heart diagnosis is the main objective of this article since heart disease is the leading cause of death nowadays and since successful heart disease diagnosis is highly helpful in saving lives [8].

Machine learning (ML) plays a significant role in disease predicting [9]. It predicts whether the patient has a particular disease type or not based on an efficient learning technique [7–10]. In this paper, we are utilizing supervised learning techniques for predicting the early stage of heart disease. Ensemble algorithms and several algorithms such as a k-nearest neighbor (KNN), support vector machine (SVM), decision tree (DT), Naive Bayes (NB), and random forest (RF) are used to classify whether the people tested belong to the class of heart disease or healthy people. Furthermore, two techniques for feature extraction, linear discriminant analysis (LDA) and principal component analysis (PCA), are used to select essential features from the dataset.

The rest of this paper is structured as follows: Section 2 describes the literature review of the current research proposed in this field. Section 3 describes the proposed architecture and methodology. In Section 4, experimental results and the comparison between classification techniques are presented. Finally, Section 5 describes the conclusion of the paper.

## 2. Literature Review

There are many literature contributions to heart disease diagnoses using data mining and machine learning techniques [11]. Reddyet al. [12] used RF, SVM, NB, NN, and KNN with multiple feature selection such as correlation matrix, recursive feature elimination (RFE), and learning vector quantization (LVQ) model to classify the cardiac disease into normal or abnormal. The results show that RF accomplished the optimal performance. Atallah and Al-Mousa [13] utilized stochastic gradient descent (SGD), KNN, RF, logistic regression (LR), and voting ensemble learning to predict cardiac diseases. The voting ensemble learning model has achieved the best accuracy of 90%. Pillaiet al. [14] used a recurrent neural network (RNN), a genetic algorithm, and K-mean to predict heart diseases. RNN has achieved the highest accuracy, and K-mean has achieved the lowest accuracy. Kannan and Vasanthi [15] used four machine learning algorithms: LR, RF, SVM, and stochastic gradient boosting (SGB) to predict heart diseases. The model prediction showed that LR has a best accuracy of 86.5%. Raza [16] applied an ensemble learning model, multilayer perceptron, LR, and NB to classify heart diseases. The result shows that ensemble learning has improved the prediction performance of cardiac disease compared to other algorithms. Oo and Win [17] used feature subset selection (CFS) with sequential minimal optimization (SMO) to predict heart diseases. The result shows that the CFS-SMO algorithm has achieved the best accuracy 86.96%. Nalluri et al. [18] used two techniques (XGBoost and LR) to improve heart disease prediction. The result showed that LR with an accuracy of 85.68% was better than XGBoost, which achieved an accuracy of 84.46%. Bhat et al. [19] proposed a model that is a combination of multilayer perceptron network (MLP) with a backpropagation algorithm to diagnose heart disease. The result shows that the proposed model has reduced error and an improved accuracy of 80.99%. Abushariah et al. utilized [20] ANN and adaptive neuro-fuzzy inference system (ANFIS) to predict cardiac disease. ANN has an obtained optimal accuracy of 87.04%, but ANFIS has achieved the lowest accuracy of 75.93%. Hasanet al. [21] utilized MLP with backpropagation and SVM to classify heart disease. The result showed that MLP achieved the highest accuracy of 98%. Chen et al. [22] used ANN with multiple features to diagnose cardiac disease. The results showed that ANN achieved the best accuracy of 80%. Sonawane and Patil [23] used vector quantization algorithm neural network to predict heart disease. Sapra et al. [24] utilized two datasets (Z-Alizadeh Sani and Cleveland heart disease dataset) that were trained by six machine learning algorithms (LR, deep learning (DL), DT, RF, SVM, and ensemble learning (gradient boosted tree)) to classify cardiac diseases. The results showed that gradient boosted tree achieved the best accuracy of 84% compared to other algorithms. Haq et al. [25] used seven machine learning algorithms: LR, ANN, KNN, NB, SVM, DT, and RF with three feature selections: minimal-redundancy-maximal-relevance (mRMR), Relief, and Shrinkage and Selection Operator (LASSO) to predict heart disease. LR with Relief achieved the highest accuracy of 89% compared to other techniques.

## 3. The Proposed System of Predicting Heart Disease

The objective of the proposed system technique is to use ensemble techniques to improve the performance of predicting heart disease. Figure 1 describes the architecture of the proposed system. It is structured into six stages, including data collection, data preprocessing, feature selection, data splitting, training models, and evaluating models.

The steps of the proposed approach are explained in detail as follows.

*3.1. Data Collection.* The heart disease dataset [26] is utilized for training and evaluating models. It consists of 1025 records, 13 features, and one target column. The target column includes two classes: 1 indicates heart diseases, and 0 indicates nonheart disease. Table 1 describes the details of the features.

*3.2. Data Preprocessing.* The features are scaled to be in the interval [0, 1]. It is worth noting that missing values are deleted from the dataset.
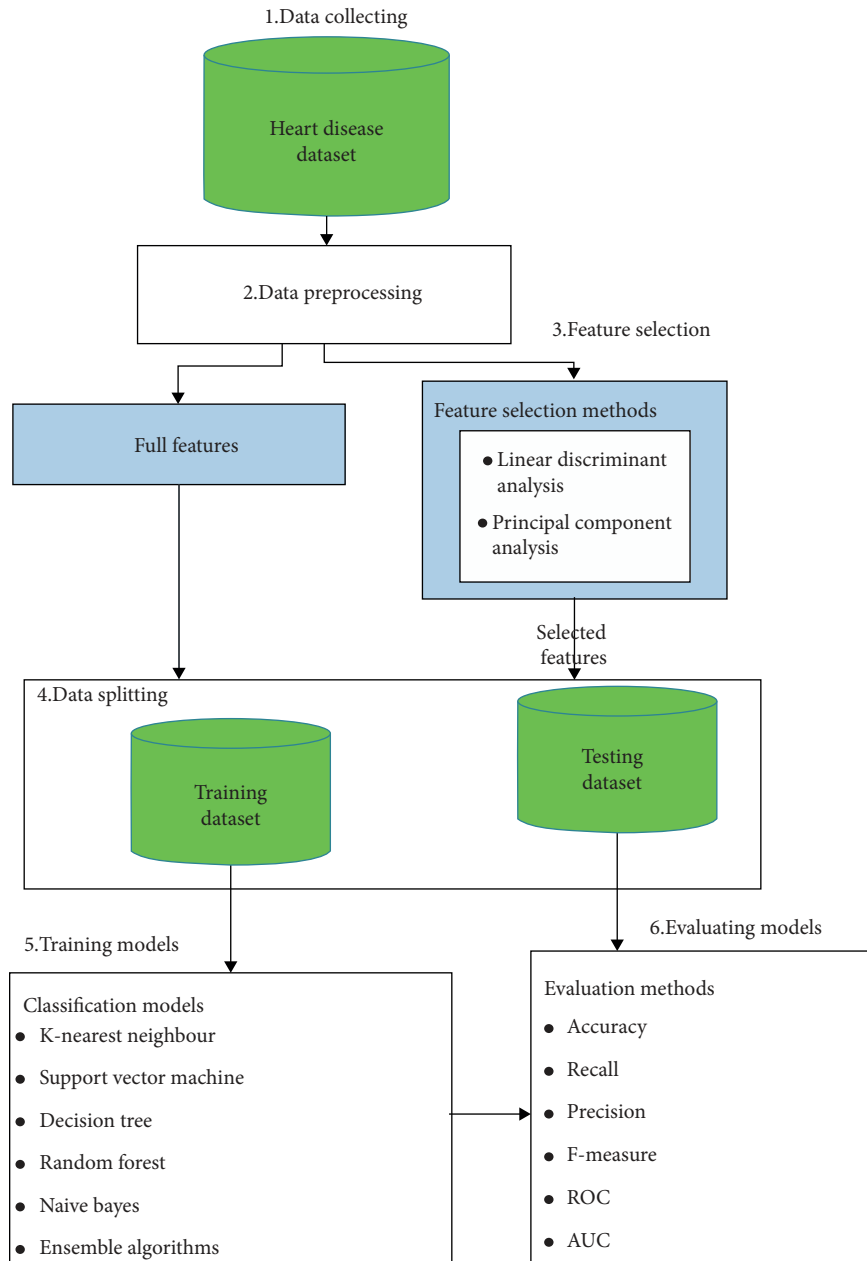
FIGURE 1: The structure of proposed system for prediction heart disease.

*3.3. Feature Extraction (FE).* The extraction of the best features is a crucial phase because irrelevant features often affect the classification efficiency of the machine learning classifier. In this phase, linear discriminant analysis (LDA) [27] and principal component analysis (PCA) [28, 29] are used to select essential features from the dataset.

*3.4. Data Splitting.* In this step, the heart disease dataset is divided into a 75% training set and a 25% as the testing set. The training set is utilized for training the models, and the testing set is utilized to evaluate the models. Also, ninefold cross-validation is utilized in the training set.

*3.5. Training Models.* Different types of machine learning algorithms: KNN, DT, RF, and NB are applied to classify heart disease. Also, two types of ensemble techniques: boosting and bagging are applied to classify heart disease:

(1) KNN is a nonparametric technique of lazy learning to enable the prediction of the new sample classification. It is utilized in several groups. It can be utilized in both the forecast problems of regression and classification. However, it is often utilized in classification when it applies to industrial problems as it fairs across all criteria examined when assessing a technique's functionality, but it is utilized mostly because of its ease of understanding and lower computation time [8–25, 27–30].

TABLE 1: Heart disease dataset descriptions.

| No. | Features | Descriptions |
|---|---|---|
| 1 | Age | Age of patient (years) |
| 2 | Sex | 1: male, 0: female |
| 3 | Chest pain (CP) | CP types<br>1 = typical angina<br>2 = atypical angina<br>3 = nonangina pain<br>4 = asymptomatic |
| 4 | RestBP | Resting blood pressure |
| 5 | Chol | Serum cholesterol in mg/dl |
| 6 | FBS | Fasting blood sugar larger 120 mg/dl (1 true) |
| 7 | RestECG | Resting electrocardiographic result |
| 8 | Thalach | Maximum heart rate accomplished |
| 9 | Exang | Exercise-induce angina (1 yes) |
| 10 | Oldpeak | ST depression induce: exercise relative to rest |
| 11 | CA | Number of major vessels (0–3) |
| 12 | Slope | Slope of peak exercise ST |
| 13 | Thal | No explanation provided, but probably thalassemia |
| 14 | Num | Diagnosis of cardiac disease:<br>1: yes<br>0: no |

(2) DT is a structure of a tree that functions on the condition's principle. It is accurate and has powerful algorithms that are utilized for predictive modeling. In particular, it has allocated internal nodes, branches, and a terminal node to include them. Every internal node carries a "test" on features, and branches carry the test conclusion, and the class label is meant for each leaf node. It is utilized both for classifications and regression [31].

(3) RF has called random decision forests to perform a ML role that can be utilized for problems with classification and regression. They function by constructing a different number of DT classifiers or regressors, and the output is obtained by enhancing all DT's output to settle a single outcome [32].

(4) NB is a family of fundamental probabilistic classifiers that focuses on applying the Bayes theorem with clear assumptions of (naive) independence between the attributes. It is extremely scalable, requiring several linear parameters for various parameters (features/predictors) in a learning problem [33].

(5) Ensemble techniques are methods that can be utilized to enhance the performance of a classifier. It is an effective classification method that combines a weak classifier with a strong classifier to improve the weak learner's efficiency [34]. The ensemble technique is used in the proposed technique to enhance the accuracy of various algorithms for diagnosing heart disease. Compared to an individual classification, the purpose of combining multiple algorithms is to obtain better performance. Figure 2 explains how the ensemble approach is utilized to enhance heart disease diagnosis.

There two types of ensemble techniques: boosting and bagging.

(a) Boosting means producing a model sequence that aims to correct the errors that have arisen in the models. The dataset is split into different subsets in detail [35]. The classification algorithm is then trained on a sample to create a series of average efficiency models as shown in pseudocode of boost algorithm, where B is the number of base hypotheses and $e$ is exp $1/e = 0.368$. Consequently, based on the previous model's elements not properly classified, new samples are produced. Then, by combining the weak models, the ensemble method increases its efficiency. The pseudocode for boosting is provided in Algorithm 1.

(b) Bagging: it refers to taking a replacement training set with multiple subsets and training a model for each subset [35]. The average of the forecast values of the submodels together are as stated by the final performance forecast. A voting procedure for each classification model is then performed as shown in pseudocode of bagging algorithm. Consequently, the classification outcome is determined based on the majority of the average values. The pseudocode for bagging is provided in Algorithm 2.

3.6. Evaluating Models. Evaluation of the proposed model is performed focusing on some criteria, namely, accuracy, recall, precision, F-score, ROC, and AUC.

Accuracy is one of the most important performance metrics for classification. It is defined as the proportion between the correct classification and the total sample, as shown in the following equation:

$$Acc = \frac{(TP + TN)}{(TP + TN + FP + FN)}. \tag{1}$$

Recall is the small portion of sufficient instances over the overall quantity of applicable instances which have been recovered. The recall equation is shown as follows:
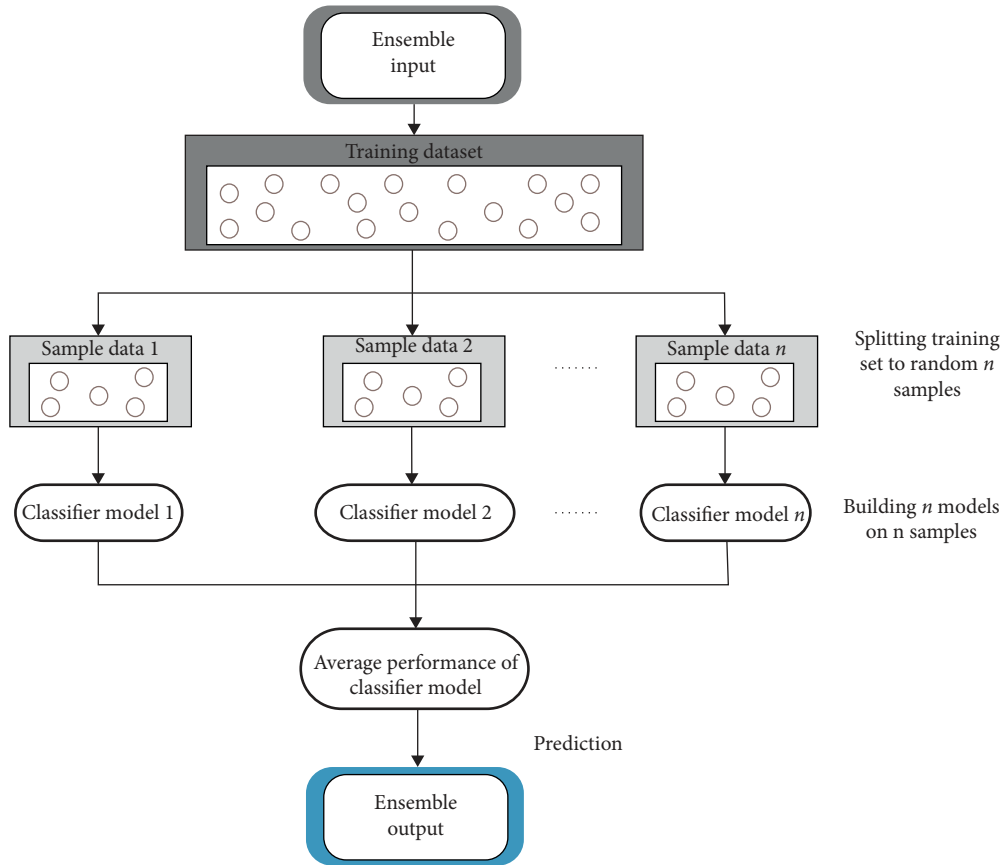
Figure 2: Building an ensemble learning prediction of heart disease.

$$recall = \frac{TP}{(TP + FN)}. \tag{2}$$

Precision is identified as follows:

$$precision = \frac{TP}{(TP + FP)}. \tag{3}$$

The F-measure is often referred to as the F1-score as follows, and it measures the mean value of precision and recall:

$$F - measure = \frac{(2 * precision * recall)}{(precision + recall)}. \tag{4}$$

The receiver operating characteristic curve (ROC) is a graph illustrating the efficiency of a classification algorithm at all classification thresholds. Two parameters are shown in this curve: true positive and false positive. The area under the curve (AUC) is the indicator of a classifier's ability to differentiate among classes and is utilized as a ROC curve description. The greater the AUC is, the greater the model's efficiency is in differentiating between the positive and negative groups.

## 4. Experimental Results

This section includes a discussion of the experimental results of classification algorithms.

### 4.1. Experimental Setup.
The experimental results have been implemented using Python. They have also been executed using Intel (R) Core i7 CPU and 8 GB of memory.

### 4.2. The Result of Applying Feature Selection Methods

#### 4.2.1. Selected Features by PCA.
Table 2 shows the score of all features extracted by PCA based on the variance of the projected features to determine the most important features. As shown in Figure 3, the number of important features was extracted by PCA equal six ($n = 6$) features. CP feature has the best score, and it is the most important feature for predicting cardiac disease.

#### 4.2.2. Selected Features by LDA.
Table 3 shows the rank of all features extracted by LDA based on the distance between features to determine the most important features. As shown in Figure 4, the number of important features was extracted by LDA equal six ($n = 6$) features. CP and CA features have the highest scores, and they are the most important feature for the prediction of cardiac disease.

### 4.3. Results of Applying the Machine Learning (ML) Algorithms to Selected Features

#### 4.3.1. Selected Features by PCA.
Table 4 shows that DT is the best performance with 98.3% accuracy, 98.7% recall,

```
Input: training number of samples M, classifier C, number iteration N
Output: result E
Training:
    Normalize weights and make the total weight is w
    Mi = sample from M
    Ci = training classifier on Mi by C
    e_i = 1/w ∑ weight (Xi)
    Bi = e_i/1 − e_i
        Weight (X_i) = weight (X_i) B_i, for all Xi where C_i (X_i) = y_i
End for
E = avg ∑ log (1/B_i)
        Ci (Xi) = y
```

ALGORITHM 1: The pseudocode of boost algorithm.
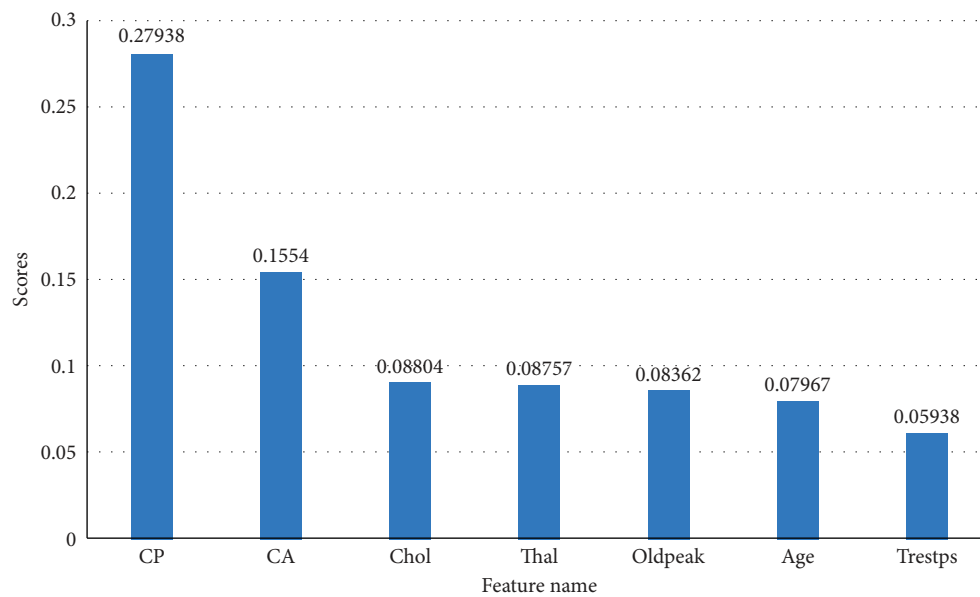


FIGURE 3: Important features extracted by PCA.

```
Input: training number of samples M, classifier C, number iteration N.
Output: result E.
Training:
        For i = 1 to N
            Mi = bootstrap sample from M
            Ci = training classifier on Mi by C
        End for
        E = avg ∑ Ci
        Ci (Xi) = y
```

ALGORITHM 2: The pseudocode of bagging algorithm.

98% AUC, and 98% precision, while the worst performance was achieved by NB: 83.7% of accuracy, 88% of recall, 81.9% of precision, 85% of F-measure, and 92% of AUC. For the KNN, we applied experiments with various $k = 1, 2, 3, 5,$ and 9. The optimal value is $k = 1$ that achieved the highest performance, an accuracy of 0.98%, 97% recall, 99% precision, and 98% AUC. NB is 83.7% classification

accuracy, 88% recall, and 81.9% precision. SVM recorded an accuracy of 84.7%, 88% recall, 83% precision, and 91% AUC. RF is 97.9% accuracy, 98% recall, 98% AUC, and 97.5% precision. The DT performance with the PCA FE algorithm outperforms the other five classification algorithms, and KNN is the second important classification algorithm.
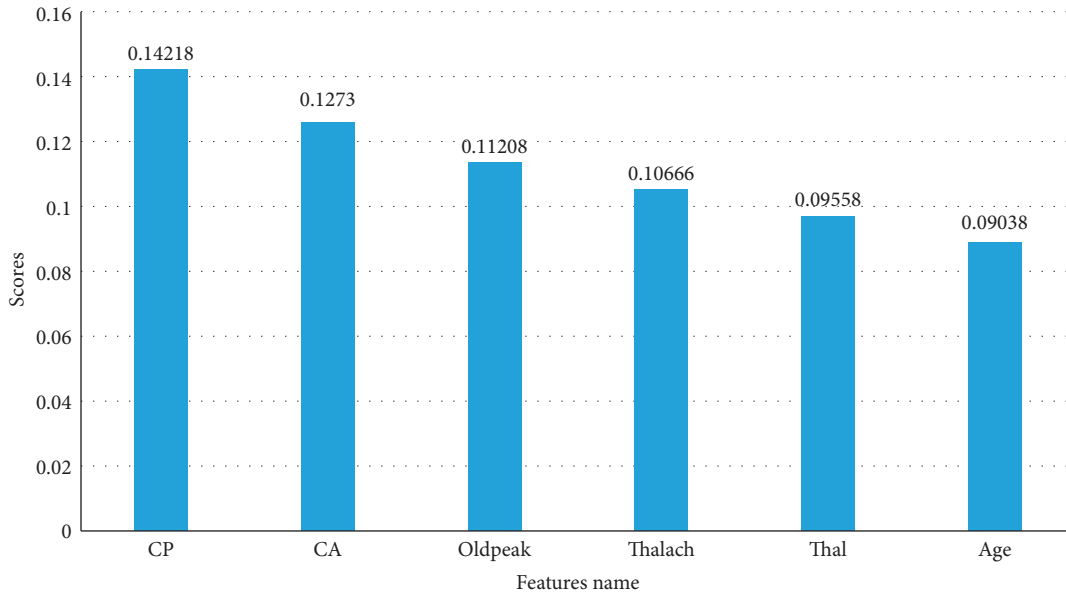
Figure 4: Important features extracted by LDA.

Table 2: The score of all features by PCA FE.

| Features | Score |
| --- | --- |
| Age | 0.07967 |
| Sex | 0.03749 |
| CP | 0.27938 |
| Trestbps | 0.05938 |
| Chol | 0.08804 |
| Fps | 0.01043 |
| Restecg | 0.01873 |
| Thalach | 0.03671 |
| Exang | 0.03001 |
| Oldpeak | 0.08362 |
| Slope | 0.03358 |
| CA | 0.15540 |
| Thal | 0.08757 |

Table 3: Score of all features by LDA.

| Features | Score |
| --- | --- |
| Age | 0.09038 |
| Sex | 0.03614 |
| CP | 0.14218 |
| Trestbps | 0.06957 |
| Chol | 0.08395 |
| Fps | 0.00905 |
| Restecg | 0.01809 |
| Thalach | 0.10666 |
| Exang | 0.06321 |
| Oldpeak | 0.11208 |
| Slope | 0.04581 |
| CA | 0.12730 |
| Thal | 0.09558 |

*4.3.2. Selected Features by LDA.* According to Table 5, it is obvious that the DT, KNN, and RF have the highest performance of accuracy, recall, precision, and F-measure which were 98.4%, 98.5%, 98%, and 98%, respectively. In KNN, we fed different $K = 1, 3, 7, 9,$ and 13. Once again, the optimal value is $k = 1$ that achieved 98.4% accuracy. SVM reported 87% accuracy, 93% recall, and 84% precision. NB reported 86.9% accuracy, 93% recall, and 83% precision. The RF achieved 98.4% classification accuracy, 98% recall, and 98% precision. The worst performance accuracy was achieved by NB and SVM, which have 86.9% and 87%, respectively.

### 4.4. The Result of Applying the Bagging Technique to Selected Features

*4.4.1. Selected Features by PCA.* In the experiment, extracted features by the PCA FE technique were checked on the bagging ensemble learning algorithm with five machine learning algorithms with 9-fold cross-validation methods.

The six important features are utilized. The classification performance was good on 6 features important.

In Table 6, DT achieved the best performance with an accuracy of 98.6%, 99% recall, 99.6% AUC, and 97.8% precision. KNN is the second important classification algorithm that has 97.9% accuracy. The worst performance accuracy was NB, which obtained 83.7%. SVM achieved 85% accuracy, 88.7% recall, 83.5% precision, and 92% AUC. NB has 83.7% classification accuracy, 88% recall, and 82% precision.

*4.4.2. Selected Features by LDA.* In the experiment, bagging ensemble learning algorithm with five machine learning algorithms is applied to selected features by the LDA.

Table 7 shows that the DT and KNN achieved the highest performance with accuracy, recall, AUC, and precision, which are 98.1%, 98.5%, 98.6%, and 98%, respectively. The worst performance is achieved by NB. RF achieved 93.8% accuracy, 94% recall, 98.4% AUC, and 94% precision. RF is the third important classification algorithm that has 93.8% accuracy.

TABLE 4: Result of ML for selected features by PCA.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-Measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN | 98 | 99 | 97 | 98 | 97.9 |
| SVM | 84.7 | 88 | 83 | 85 | 91 |
| DT | 98.3 | 98.7 | 98 | 98 | 98 |
| RF | 97.9 | 98 | 97.5 | 98 | 98 |
| NB | 83.7 | 88 | 81.9 | 85 | 92 |

TABLE 5: Result of ML for selected features by LDA.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN | 98.4 | 98.5 | 98 | 98 | 98 |
| SVM | 87 | 93 | 84 | 88 | 93 |
| DT | 98.4 | 98.5 | 98 | 98 | 98 |
| RF | 98.4 | 98.5 | 98 | 98 | 98.6 |
| NB | 86.9 | 93 | 83.8 | 88 | 93 |

### 4.5. The Result of Applying Boosting Technique to Selected Features

*4.5.1. Selected Features by PCA.* In the experiment, boosting ensemble learning algorithm with five machine learning algorithms are applied to selected features by the PCA.

Table 8 shows RF has achieved the highest accuracy at 98.3%, and SVM has achieved the second-highest accuracy at 98%. The worst accuracy has been performed by SVM at 83%. DT obtained 98.8% recall, 98% AUC, and 97.6% precision. RF obtained 98.7% recall, 99.8% AUC, and 98% precision. The optimal result for KNN when $k = 1$ is 97.8% accuracy.

*4.5.2. Selected Features by LDA.* In the experiment, extracted features by the LDA FE technique were checked on boosting ensemble learning algorithm with five machine learning algorithms.

In Table 9, RF has reported the best performance among other algorithms with an accuracy of 98.2%, 98.5% recall, 98.2% AUC, and 98% precision. In contrast, SVM has registered the lowest performance with 85% accuracy, 94.9% recall, 79.9% precision, and 89.2% AUC.

The optimal result of KNN when $k = 1$ is 98.1% accuracy. The classification accuracy of NB is 86.7%, 90% recall, and 85% precision. DT achieved 98.1% accuracy, 98.5% recall, 98.1% AUC, and 98% precision. DT and KNN are the second important classification algorithms that have 98.1% accuracy.

Table 10 shows the comparison of the results of the proposed model (the bagging ensemble learning method with decision tree) with various other state-of-the-art algorithms. It is obvious from Table 10 that a state-of-the-art algorithm's optimal performance achieved an accuracy of 89.5% [36]. On the other side, the proposed model performance has achieved 98.6% accuracy. So it is clear that the proposed model outperforms other competitors [18, 20, 24, 36, 37] significantly.

TABLE 6: Result of applying the bagging technique to selected features by PCA.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN | 97.9 | 99 | 97 | 98 | 97.9 |
| SVM | 85 | 88.7 | 83.5 | 85.9 | 92 |
| DT | 98.6 | 99 | 97.8 | 98.5 | 99.6 |
| RF | 96.2 | 96.7 | 96 | 96.4 | 99.5 |
| NB | 83.7 | 88.2 | 82 | 84.8 | 92.3 |

TABLE 7: Classification performance of applying the bagging technique to selected features by LDA.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN | 98.1 | 98.5 | 98 | 98.2 | 98.6 |
| SVM | 87 | 93 | 84 | 88 | 93 |
| DT | 98.1 | 98.5 | 98 | 98 | 98.6 |
| RF | 93.8 | 94 | 94 | 94 | 98.4 |
| NB | 86.9 | 93 | 83.6 | 88 | 93.2 |

TABLE 8: Classification performance of applying boosting technique to selected features by PCA.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN | 97.8 | 99 | 97 | 97.9 | 97.9 |
| SVM | 83 | 91.9 | 78.9 | 84.8 | 89.3 |
| DT | 98 | 98.8 | 97.6 | 98.1 | 98 |
| RF | 98.3 | 98.7 | 98 | 98.4 | 99.8 |
| NB | 83.2 | 87 | 81.9 | 84.3 | 92.6 |

TABLE 9: Classification performance of applying boosting technique to selected features by LAD.

| Techniques | Accuracy (%) | Recall (%) | Precision (%) | F-measure (%) | AUC (%) |
|---|---|---|---|---|---|
| KNN ($K = 1$) | 98.1 | 98.5 | 98 | 98.2 | 98.5 |
| SVM | 85 | 94.9 | 79.9 | 86.7 | 89.2 |
| DT | 98.1 | 98.5 | 98 | 98.2 | 98.1 |
| RF | 98.2 | 98.5 | 98 | 98.2 | 98.2 |
| NB | 86.7 | 90 | 85 | 87.6 | 93 |

TABLE 10: Comparison of the results of proposed model with various other state-of-the-art algorithms [18, 20, 24, 36, 38].

| Compared algorithms | Accuracy (%) |
|---|---|
| KNN and NB [36] | 89.5 |
| Boosting [37] | 85.2 |
| XGBoost and LR [18] | 85.68 |
| ANFSI and ANN [20] | 87.04 |
| Ensemble learning (gradient boosted) [24] | 84 |
| Bagging ensemble learning method with decision tree | **98.6** |

## 5. Conclusion

In this paper, we developed the proposed system to predict heart disease. Ensemble methods (boosting and bagging) with feature extraction algorithms (PCA and LDA) are used to improve predicting heart disease performance. The feature extraction algorithms are used to extract essential features from the Cleveland heart disease dataset. Comparison between ensemble methods (boosting and bagging) and five classifiers (KNN, SVM, NB, DT, and RF) is applied to selected features. The experimental results showed that the bagging ensemble learning algorithm with DT and PCA feature extraction method had achieved the best performance.

## Data Availability

The heart disease dataset used to support the findings of this study are available at https://www.kaggle.com/johnsmith88/heart-disease-dataset.

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## References

[1] M. Sanz, A. Marco del Castillo, S. Jepsen et al., "Periodontitis and cardiovascular diseases: consensus report," *Journal of Clinical Periodontology*, vol. 47, no. 3, pp. 268–288, 2020.

[2] World Health Organization. http://www.who.int/cardiovascular diseases/en. 2019.

[3] J. Nahar, T. Imam, K. S. Tickle, and Y.-P. P. Chen, "Association rule mining to detect factors which contribute to heart disease in males and females," *Expert Systems with Applications*, vol. 40, no. 4, pp. 1086–1093, 2013.

[4] S. N. Blair, "Commentary on Wang Y et al. "An Overview of Non-exercise Estimated Cardiorespiratory Fitness: estimation Equations, Cross-Validation and Application"" *Journal of Science in Sport and Exercise*, vol. 1, no. 1, pp. 94-95, 2019.

[5] K. Vanisree and J. Singaraju, "Decision support system for congenital heart disease diagnosis based on signs and symptomsusing neural networks," *International Journal of Computer Application*, vol. 19, pp. 6–12, 2011.

[6] A.-H. Abdel-Aty, H. Kadry, M. Zidan et al., "A quantum classification algorithm for classification incomplete patterns based on entanglement measure," *Journal of Intelligent and Fuzzy Systems*, vol. 38, no. 9, pp. 1–8, 2020.

[7] A. Sagheer, M. Zidan, and M. M. Abdelsamea, "A novel autonomous perceptron model for pattern classification applications," *Entropy*, vol. 21, no. 8, p. 763, 2019.

[8] M. Aljanabi, H. Qutqut, and M. Hijjawi, "Machine learning classification techniques for heart disease prediction: a review," *International Journal of Engineering and Technology*, vol. 7, pp. 5373–5379, 2018.

[9] T. Obasi and M. O. Shafiq, "Towards comparing and using machine learning techniques for detecting and predicting heartattack and diseases," in *Proceedings of the 2019 IEEE International Conference on Big Data (Big Data)*, pp. 2393–2402, Los Angeles, CA, USA, April 2019.

[10] A. A. Ali, H. S. Hassan, and E. M. Anwar, "Heart diseases diagnosis based on a novel convolution neural network and gaterecurrent unit technique," in *Proceedings of the 2020 12th International Conference on Electrical Engineering (ICEENG)*, vol. 145–150, Cairo, Egypt, July 2020.

[11] V. Pham, Q. De Hemptinne, J.-M. Grinda et al., "Giant coronary aneurysms, from diagnosis to treatment: a literature review," *Archives of Cardiovascular Diseases*, vol. 113, pp. 59–69, 2020.

[12] N. S. C. Reddy, S. S. Nee, L. Z. Min, and C. X. Ying, "Classification and feature selection approaches by machine learningtechniques: heart disease prediction," *International Journal of Innovative Computing*, vol. 9, 2019.

[13] R. Atallah and A. Al-Mousa, "Heart disease detection using machine learning majority voting ensemble method," in *Proceedings of the 2019 2nd International Conference on New Trends in Computing Sciences (ICTCS)*, pp. 1–6, Amman, Jordan, October 2019.

[14] N. S. R. Pillai, K. K. Bee, and J. Kiruthika, "Prediction of heart disease using rnn algorithm," *International Research Journal of Engineering and Technology*, vol. 5, 2019.

[15] R. Kannan and V. Vasanthi, "Machine learning algorithms with roc curve for predicting and diagnosing the heart disease," *InSoft Computing and Medical Bioinformatics*, pp. 63–72, 2019.

[16] K. Raza, "Improving the prediction accuracy of heart disease with ensemble learning and majority voting rule," *InU-Healthcare Monitoring Systems*, pp. 179–196, 2019.

[17] A. N. Oo and K. T. Win: Feature Selection Based Sequential Minimal Optimization (Smo) Classifier for Heart Diseaseclassification.

[18] S. Nalluri, R. V. Saraswathi, S. Ramasubbareddy, K. Govinda, and E. Swetha, "Chronic heart disease prediction using datamining techniques," *InData Engineering and Communication Technology*, pp. 903–912, 2020.

[19] R. Bhat, S. Chawande, and S. Chadda, "Prediction of test for heart disease diagnosis using artificial neural network," *Indian Journal of Applied Research*, vol. 9, 2019.

[20] M. A. Abushariah, A. A. Alqudah, O. Y. Adwan et al., "Automatic heart disease diagnosis system based onartificial neural network (ann) and adaptive neuro-fuzzy inference systems (anfis) approaches," *Journal of Software Engineering and Applications*, vol. 7, p. 1055, 2014.

[21] T. T. Hasan, M. H. Jasim, and I. A. Hashim, "Heart disease diagnosis system based on multi-layer perceptron neural networkand support vector machine," *International Journal of Current Engineering and Technology*, vol. 77, pp. 2277–4106, 2017.

[22] A. H. Chen, S.-Y. Huang, P.-S. Hong, C.-H. Cheng, and E.-J. Lin, "Hdps: heart disease prediction system," *In 2011 computing in Cardiology*, vol. 557–560, 2011.

[23] J. S. Sonawane and D. Patil, "Prediction of heart disease using learning vector quantization algorithm," in *Proceedings of the 2014 Conferenceon IT in Business, Industry and Government (CSIBIG)*, vol. 1–5, Indore, India, March 2014.

[24] L. Sapra, J. K. Sandhu, and N. Goyal, "Intelligent method for detection of coronary artery disease with ensemble approach," *Advances in Communication and Computational Technology*, vol. 1033–1042, 2021.

[25] A. U. Haq, J. P. Li, M. H. Memon, S. Nazir, and R. Sun, "A hybrid intelligent system framework for the prediction of heartdisease using machine learning algorithms," *Mobile*

*Information Systems*, vol. 2018, Article ID 3860146, 21 pages, 2018.

[26] https://www.kaggle.com/johnsmith88/heart-diseasedataset.

[27] C. Ricciardi, A. S. Valente, K. Edmund et al., "Linear discriminant analysis and principal component analysis to predict coronary artery disease," *Health Informatics Journal*, vol. 26, no. 3, pp. 2181–2192, 2020.

[28] A. K. Garate-Escamilla, A. H. E. Hassani, and E. Andres, "Classification models for heart disease prediction using featureselection and pca," *Informatics Medicine Unlocked*, vol. 19, Article ID 100330, 2020.

[29] A. A. Ali, H. S. Hassan, and E. M. Anwar, "Improve the accuracy of heart disease predictions using machine learning andfeature selection techniques," in *Proceedings of the International Conference on Machine Learning, Image Processing, Network Security and Data Sciences*, pp. 214–228, Assam, India, April 2020.

[30] M. Yaqoob, F. Iqbal, and S. Zahir, "Comparing predictive performance of k-nearest neighbors and support vector machinefor predicting ischemic heart disease," *Journal of Advanced Scientific Research*, vol. 1, 2020.

[31] Q. Fan, Z. Wang, D. Li, D. Gao, and H. Zha, "Entropy-based fuzzy support vector machine for imbalanced datasets," *Knowledge-Based System*, vol. 115, pp. 87–99, 2017.

[32] D. C. Yadav and S. Pal, "Prediction of heart disease using feature selection and random forest ensemble method," *International Journal for Pharmaceutical Research Scholars*, vol. 12, pp. 56–66, 2020.

[33] S. S. Yadav, S. M. Jadhav, S. Nagrale, and N. Patil, "Application of machine learning for the detection of heart disease," in *2020 2nd International Conference on Innovative Mechanisms for Industry Applications (ICIMIA)*, vol. 165–172, Bangalore, India, March 2020.

[34] H. B. F. David: Impact of Ensemble Learning Algorithms towards Accurate Heart Disease Prediction.

[35] C. Zhang and Y. Ma, "Ensemble Machine Learning;," *Methods and Applications*, Springer Science & Business Media, Berlin, Germany, 2012.

[36] E. Z. Ferdousy, M. M. Islam, and M. A. Matin, "Combination of naive bayes classifier and k-nearest neighbor (cnk) in theclassification based predictive models," *Computer and Information Science*, vol. 6, no. 3, 2013.

[37] K. H. Miao, J. H. Miao, and G. J. Miao, "Diagnosing coronary heart disease using ensemble machine learning," *International Journal of Advanced Computer Science and Applications*, vol. 7, pp. 30–39, 2016.

[38] S. Pouriyeh, S. Vahid, G. Sannino et al., "A comprehensive investigation and comparison of machine learning techniques in the domain of heartdisease," in *Proceedings of the 2017 IEEE Symposium on Computers and Communications (ISCC)*, pp. 204–207, Heraklion, Greece, July 2017.