

The Role of Data Pre-Processing Techniques In Improving Machine Learning Accuracy For Predicting Coronary Heart Disease

Osamah Sami Yousef Elsheikh Fadi Almasalha

Faculty of Information Technology, Applied Science Private University, Amman 11931, Jordan

Abstract: These days, in light of the rapid developments, people work day and night to live at a good level. This often causes them to not pay much attention to a healthy lifestyle, such as what they eat or even what physical activities they do. These people are often the most likely to suffer from coronary heart disease. The heart is a small organ responsible for pumping oxygen-rich blood to the rest of the human body through the coronary arteries. Accordingly, any blockage or narrowing in one of these coronary arteries may cause blood not to be pumped to the heart and from it to the rest of the body, and thus cause what is known as heart attacks. From here, the importance of early prediction of coronary heart disease has emerged, as it can help these people change their lifestyle and eating habits to become healthier and thus prevent coronary heart disease and avoid death. In this paper, we will work to improve the accuracy of machine learning techniques in predicting coronary heart disease using feature processing techniques. Feature processing is a technique used to improve the efficiency of a machine learning model by improving the quality of the feature. The popular Framingham Heart Study dataset was used for validation purposes. The results of the research paper indicate that the use of feature processing techniques had a role in improving the predictive accuracy of poorly efficient classifiers, and shows satisfactory performance in determining the risk of coronary heart disease. For example, the Decision Tree classifier led to a predictive accuracy of coronary heart disease of 91.39% with an increase of 1.39% over the previous work, the Random Forest classifier led to a predictive accuracy of 92.80% with an increase of 2.7% over the previous work, the KNN classifier led to a predictive accuracy of 92.68% with an increase of 3.64% over the previous work, the Multilayer Perceptron Neural Network (MLP) classifier led to a predictive accuracy of 92.64% with an increase of 2.68% over the previous work, and the Naïve Bayes classifier led to a predictive accuracy of 90.56% with an increase of 0.66% over the previous work.

Keywords: Coronary heart disease, heart, machine learning, feature processing, classification technique.

1 Introduction

The heart is one of the most important organs in the human body. It is a small, muscular pumping organ responsible for supplying other organs in the body with oxygen and other important nutrients [1]. This means that a person's life depends on the efficiency of heart function. Therefore, if the heart does not function well, other organs also cannot function well [2].

People, in light of the difficult economic conditions, seek to secure their basic needs by working long hours daily. This lifestyle often does not take into account the diet and health of these people to ensure their safety [3]. This type often leads to a risk of diseases such as diabetes, high cholesterol and blood pressure at an early age, and all of these diseases, if not controlled, can lead to coronary heart disease [3].

Heart disease is a term that refers to any problem that can affect the heart and blood vessels [1]. Usually when people say that they have heart disease, they are referring to

most dangerous and common diseases in the world [5]. coronary heart disease (CHD) which, according to the National Heart, Lung, and Blood Institute ranks among the

In coronary heart disease, a complete or partial blockage of the coronary arteries usually occurs due to blood clotting or the accumulation of fatty plaques on the walls, which leads to the inability of the heart to get enough oxygen [4] and thus it is difficult for the heart to function as efficiently as required.

There are two risk factors for coronary heart disease. The first type is stable and cannot be changed, such as age, gender and family history, while the other type depends on lifestyle such as diabetes, smoking, high cholesterol, high blood pressure, high body mass index, and low exercise [6]. However, the second type of risk factors can usually be controlled, according to experts, by changing our lifestyle and diet, and using certain medications if needed.

In recent years, artificial intelligence techniques have been used extensively in the medical fields in order to improve the efficiency of disease diagnosis/classification in its early stages [7] [8]. Among those techniques stand out machine learning techniques, which are a set of statistical models that help the machine learn from past data [9]. In spite of this, it is often difficult to deal with patient data for diagnosis in the early stages due to reasons such as data

volume, missing values and noise in the data. But machine learning techniques and their capabilities have helped process such data [10].

Also, it is noticeable regarding data features that they may be incomplete and huge. The range of some data features is small while the range is large for other data features. The type of data features is combined between categorical and numerical, all of this will affect the accuracy of machine learning techniques in diagnosing and classifying diseases in their early stages, including coronary heart disease. We will therefore solve these problems by using different techniques to manipulate the features under the so-called feature processing techniques and thus improve the accuracy of machine learning techniques in early prediction of the disease [11].

This research paper is organized as follows: The second section is a review of some relevant work. The third section presents the methodology for this research paper. The fourth section is for presenting, evaluating and discussing the results of the research paper. The fifth section is for conclusion and future work.

2 Related Work

Recently, there has been an increase in the number of papers dealing with the use of machine learning techniques in predicting serious diseases that may affect people's lives, including coronary heart disease.

In [12], the researchers applied a logistic regression technique on the Framingham Heart Study dataset to predict the ten-year risk of coronary heart disease. The researchers used 65% of the dataset for the training set. The accuracy obtained was 84.8%.

The researchers in [13] had a contribution by implementing four machine learning algorithms, namely support vector machine (SVM), neural network, XGBoost, and random forest to predict the ten-year risk of coronary heart disease. The researchers also used the Framingham Heart Study dataset to validate the results. The accuracy obtained was 84.8% for support vector machine, 85.4% for neural network, 86.99% for XGBoost, and 84.9% for random forest.

Also, the researchers in [4] contributed to the literature of this field by using boosting adaptive algorithm on four datasets, namely (UCI Cleveland, UCI Switzerland, UCI Long Beach, and UCI Hungarian) to diagnose coronary heart disease. This approach obtained accuracy (97.16% and 80.14% for Cleveland, 98.63% and 89.12% for Hungarian, 93.15% and 77.78% for Long Beach, 100% and 96.72% for Switzerland) for training and testing set respectively.

In [14], the researchers applied three machine learning algorithms, namely support vector machine, neural network, and Hybrid-SVM on the Framingham Heart Study dataset to predict the ten-year risk of heart attack. The accuracy obtained was 86.03% for support vector machine,

84.7% for neural network, and 94% for Hybrid-SVM.

However, these results were better for some of the machine learning techniques used than those used for [13].

In [15], the researchers applied six algorithms, namely decision tree, boosted decision tree, random forest, support vector machine, neural network, and logistic regression on the Framingham Heart Study dataset to predict the ten-year risk of coronary heart disease. The data was divided into 80% training and 20% testing. The researchers used R Studio and Rapid-Miner in their work. The researchers used three techniques to deal with missing values. The first technique is to ignore missing values, and obtained accuracy of 85% for the decision tree, 63% for the boosted decision tree, and 63% for logistic regression. All this while using the Rapid-Miner tool. Whereas, the R studio tool enabled the researchers to obtain the accuracy of 84% for the decision Tree, 85% for the boosted decision tree, and 84% for logistic regression. Analysis of complete case is the second technique used, as the Rapid-Miner tool enabled the researchers to obtain accuracy of 54% for the decision tree, 64% for the boosted decision tree, 65% for the random forest, 69% for the support vector machine, 69% for the neural network, and 68% for logistic regression. R studio tool obtained accuracy 67%, 81%, 79%, 69%, 67%, and 68% for the decision tree, boosted decision tree, random forest, support vector machine, neural network, and logistic regression respectively. The final technique is to be replaced with the average, and the accuracy obtained while using the Rapid-Miner tool was 62% for the decision tree, 62% for the boosted decision tree, 63% for the random forest, 68% for the support vector machine, 68% for the neural network, and 67% for logistic regression. Whereas, the R Studio tool enabled the researchers to obtain an accuracy of 84% for the decision tree, 84% for the boosted decision tree, 78% for the random forest, 68% for the support vector machine, 71% for the neural network, and 66% for logistic regression.

However, other researchers such as those in [16] applied only one algorithm which is the logistic regression on the Framingham Heart Study dataset to predict the ten-year risk of coronary heart disease. This approach obtained better accuracy of 86.6% than ever.

In [17], the researchers applied the same previous method of logistic regression to the Framingham Heart Study dataset to predict a heart attack. This approach obtained an accuracy of 87%.

Other researchers such as those in [18] applied the neural network algorithm to real data from patient of Paris H^otel-Dieu University Hospital to diagnose coronary heart disease. Their approach used a different number of input factors (6 to 14). The approach obtained 63% for features (age, diabetes, hypertension, obesity, smoking, family anamnesis of CHD), 76% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD), 77% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD, glycaemia, cholesterol total), 81% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD, TG, cholesterol 0.81 69 79 total, HDL, LDL, glycaemia),

83% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD, carotid plaque), 87% for features (diabetes, hyper tension, obesity, smoking, family anamnesis of CHD, PWV index), 91% for features (diabetes, hypertension, obesity, smoking, family anamnesis of CHD, carotid plaque, PWV index), 93% for features (diabetes, hyper- tension, obesity, smoking, family anamnesis of CHD, TG, cholesterol, HDL, 0.93 80 92 LDL, glycaemia, carotid plaque, PWV index), 77% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD, glycaemia, 0.77 53 87 cholesterol total, cGFR), and 77% for features (age, sex, diabetes, hypertension, obesity, smoking, family anamnesis of CHD, glycaemia, cholesterol total, left ventricular hypertrophy).

Those in [19] applied the deep belief algorithm to the KNHANES-6 dataset to predict the risk of coronary heart disease and obtained an accuracy of 82%. However, the researchers applied the genetic algorithm to improve the deep belief network and the obtained accuracy was 74%.

In [20], the researchers applied a logistic regression and neural network to the KNHANES-VI dataset to predict the risk of coronary heart disease. However, this approach obtained accuracy 86.11% for the logistic regression and 87.04% for the neural network. The researchers used a distinct correlation analysis to improve the accuracy of the neural network to become 87.63%.

In other research such as [21], the researchers applied naïve bayes, KNN, random forest, decision tree, SVM, logistic regression, and the ensemble classification approach to the NHANES and Framingham Heart Study dataset, to monitor the risk of chronic diseases. For the NHANES dataset, the decision tree algorithm obtained an accuracy of 97.6%, 96.5% for the ensemble approach, 80.8% for the KNN, 96.4% for logistic regression, 95.7% for naïve bayes, 98.5% for random forest, 95.4% for SVM. Whereas, the results for Framingham Heart Study dataset were as follows: The decision tree obtained an accuracy of 90%, 89.3% for the ensemble approach, 90.1% for the KNN, 90% for the logistic regression, 89.9% for naïve bayes, 90.1% for random forest, and 90.2% for SVM.

Similarly, the researchers of [22] applied naïve bayes, KNN, random forest, decision tree, SVM, logistic regression, neural network, and the ensemble classification approach to the NHANES and Framingham Heart Study dataset to predict Cardiovascular disease. For the NHANES dataset, the decision tree algorithm obtained an accuracy of 97.6%, 96.5% for the ensemble approach, 80.8% for the KNN, 96.4% for logistic regression, 95.7% for naïve bayes, 98.5% for random forest, 95.4% for SVM, 98.8% neural network. Whereas, the results for Framingham Heart Study dataset were as follows: The decision tree obtained an accuracy of 90%, 89.3% for the ensemble approach, 90.1% for the KNN, 90% for logistic regression, 89.9% for naïve bayes, 90.1% for random forest, 90.2% for SVM, and 89% for neural network.

Despite this and many other researches, the field is still open for researchers to conduct their experiments in order

to improve the accuracy of the machine learning techniques for predicting diseases that pose a risk to human life, including coronary heart disease.

3 Research Approach

It is unfortunate to hear that there is an increase in the number of patients diagnosed with coronary heart disease (angina or heart attack) day after day. High blood pressure, high cholesterol, uncontrolled diabetes, smoking, and a diagnosis of cardiovascular impairment and other risks, all increase the chance of diagnosis with coronary heart disease in the future. Therefore, we need an accurate system that helps the patient protect him/herself from the risk of coronary heart disease, relying in this on the patient's demographic information, medical history, medical examination, behavior, and laboratory examination.

Many researchers have developed machine learning models using different classification algorithms such as decision tree, naïve bayes, SVM, KNN, and neural network. Most of these models were utilizing the Cleveland Heart Diseases dataset to predict coronary heart diseases, but few were using the Framingham Study dataset. This paper uses the Framingham Study dataset to validate the resulting model since it includes features for most of the potential risk factors for coronary heart disease and some of these features are not found in the most common dataset of heart disease namely, Cleveland Heart Disease dataset. In this paper, five machine learning classification algorithms were used such as decision tree, naïve bayes, neural network, random forest, and KNN. These five algorithms used the Framingham Heart Study dataset with two events for target (output) features to predict coronary heart disease, as a number of different feature processing techniques will be used to improve the accuracy of machine learning models for predicting coronary heart disease.

3.1 Dataset

The Framingham Heart Study dataset is the first long-term epidemiological study concerned with the possible causes of cardiovascular disease that began in 1948 in Framingham, Massachusetts. The Framingham Heart Study dataset identified the prospective risk factors of cardiovascular diseases and their effects [23]. Our dataset is an educational dataset taken from the National Heart, Lung, and Blood institute that did this study [23].

The dataset contains 19 input features divided into demographic features(Age, Gender), behavioral features(Current Smoker, Cigarettes Per Day, Body Mass Index), medical history features(Prevalent Coronary Heart Diseases, Prevalent Angina Pectoris, Prevalent Myocardial Infarction, Prevalent Stroke, Prevalent Hypertensive, Use Blood Pressure Drugs, Diabetes), medical examination features(Systolic Blood Pressure, Diastolic Blood Pressure, Heart Rate) and laboratory testing features(Glucose, High-Density Lipoprotein, Low-Density Lipoprotein, Total Cholesterol), and two features for outputs (Angina Pectoris,

Myocardial Infarction).

3.2 Feature Processing

Feature processing is a group of techniques that are applied on features to improve the quality of the feature, such as handling missing values, convert the type of feature and many other techniques [11].

Impute Missing Values By KNN

KNN works for missing values by calculating the distance or similarity to find the most similar case in the dataset and changing the missing value with it [24], by applying (1):

$$\text{Dist}(X, Y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

Min Max Normalization

This method converts each numerical feature value to a new one based on the minimum and maximum values of the feature [25], by applying (2):

$$\bar{x} = \frac{x - \text{Min}}{\text{Max} - \text{Min}} \quad (2)$$

Z-Score Standardization

This method converts each numerical feature value to a new value based on the standard deviation and Mean of the feature [25], by applying (3):

$$\bar{x} = \frac{x - \mu}{\sigma} \quad (3)$$

Standard Deviation

One Hot Encoding

One Hot Encoding splits the categorical feature into a separate number of features depending on the number of the cases in the original categorical feature, and gives 0 for absence and 1 for presence in each new feature [26].

Ordinal Encoding

In this technique, every case in the categorical feature is converted into an integer value [26].

Equal Width Discretization

Here is an easy method to sort numeric feature values and divide the range of sort values into pre-defined equal width bins [27] by applying (4) and (5):

$$x = \frac{x - \text{Min}}{\text{Max} - \text{Min}} \quad (4)$$

$$\text{Boundaries} = \text{Min} + (i * W) \quad (5)$$

Equal Frequency Discretization

In this method, first sort the values in ascending order. Divide the range of sort values into predetermined number of equal frequency bins by applying $\frac{x - \text{Min}}{\text{Max} - \text{Min}}$ each bin has the same number of values [27].

3.3 Classification Algorithms

ID3 Decision Tree

features are set as a root node after which the features are split by finding the entropy that uses the data coherence scale; The entropy value is between 0 and 1 [8]. The entropy can be found by applying (6) and (7).

$$\text{Entropy}(F) = - \sum_{\#=1}^c p_{\#} \log_2 p_{\#} \quad (6)$$

$$\text{Gain}(F, A) = \text{Entropy}(F) - \sum_{i=1}^c \frac{|A_i|}{|F|} \text{Entropy}(A_i) \quad (7)$$

The feature with the highest value is identified as the root node of the tree [8].

Random Forest

Random forest is a classification algorithm [28] works by generating multiple decision trees from the dataset [28]. Features are randomly selected from the training set for building trees in random forest [28]. After building each decision tree and finding the outcome for each tree, apply a majority vote to determine the final outcome of the random forest [28]. In the process of building each decision tree, the randomization is applied to find the split node value.

K-Nearest Neighbors

KNN Neighbors is a supervised machine learning algorithm used to predict and classify unknown data from known data by measuring the distance between them [29]. The distance scale is used to measure the distance between two numerical values [30]. The distance can be calculated by applying (8):

$$\text{Cosine}(x, y) = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2} \sqrt{\sum_{i=1}^n y_i^2}} \quad (8)$$

Multilayer Perceptron (MLP) Neural Network

The structure of an artificial neural network is the same as the human brain [31]. Multilayer perceptron (MLP) contains more than one layer (input layer, hidden layer (s), output layer) [1]. First, in the neural network before starting training from the data set, the weight value (w) is randomly determined [1]. Then, the neurons begin to learn from the training set [1].

The activation function applies some gradient processing to the input data to find the neural network output [32]. Sigmoid is a non-linear activation function commonly used in feedforward neural networks [32].

$$F(x) = \frac{1}{1 + e^{-x}} \quad (9)$$

Each decision tree contains a root node, a leaf node, an inner node and branches. In the ID3 decision tree, all the

Back Propagation algorithm is commonly used to train a multilayer perceptron neural network. In the first step of this algorithm, a comparison is made between the prediction output (\hat{Y}) and the actual output (Y) to find the error between them, and this error returns to the neural network and the weight changes based on this error, and the numerical weight changes until the value of \hat{Y} becomes closer to (Y) [1].

Naïve Bayes

is a statistical classification algorithm that works on the basis of Bayes' theory, and Naïve Bayes assumes that each feature is separate, and each variable is distinct in

prediction and occurrence [3]. Naïve Bayes uses the prior probability of Bayes theorem to calculate the likelihood of the relationship between each feature in the test data with each target, the target with the highest probability is selected as the result of the model [33].

Stratified KFold Cross validation

Cross validation is a static method used to test an algorithm by dividing the data set into a training used to train the model and the test used to evaluate the model's performance [34]. In cross-validation, every point has the same chance of being used in the test [34]. In kfold, the dataset is evenly divided into k number of fields [34]. Stratified KFold means that each fold has the same class naming distribution in the original dataset [35]. For each iteration, one test folds and others are used for training [34].

4 Experimental Result And Discussion

In this paper, we used five machine learning classification techniques to predict two primary CHD events, namely, angina pectoris (528 yes, 2735 no) and myocardial infarction (308 yes, 2955 no).

4.1 Tool

RapidMiner is a data science software platform developed by the company of the same name that provides an integrated environment for data preparation, machine learning, deep learning, text mining, and predictive analytics [36]. In machine learning, RapidMiner can be used for feature processing, dataset segmentation, model training, model testing, network research, and performance evaluation [36].

4.2 Performance Evaluation

Performance evaluation is a group of equations used to measure the effectiveness of the classifier or the model [2]. Below is the definition of some essential terms used in the equations of performance evaluation as shown in Table 1:

True Positive (TP)

The person is healthy and also predict as healthy [2].

False Positive (FP):

The person is healthy, but predict as sick [2]

True Negative (TN)

The person is sick and predict as sick [2]

False Negative (FN)

The person is sick, but predict as healthy [2]

Table 1 Confusion matrix

| | Negative (Actual) | Positive (Actual) |
|--------------------|-------------------|-------------------|
| Negative (Predict) | TN | FN |
| Positive (Predict) | FP | TP |

4.3 Performance Metrics

Accuracy (Acc)

Accuracy is an evaluation metric of the total number of predictions the model or the classifier gets right [37]. The accuracy can be calculated by applying (10).

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN} \quad (10)$$

Precision

Precision is an evaluation metric that assesses how accurately a model predicts positive labels [37]. Precision is the percentage of your results which are relevant. The precision can be calculated by applying (11).

$$\text{Precision} = \frac{TP}{TP + FP} \quad (11)$$

F-Measure

F-Measure refers to the mean of consistency between Precision and Recall [37]. The F-Measure can be calculated by applying (12).

$$F - \text{Measure} = 2 * \left(\frac{\text{Recall} * \text{Precision}}{\text{Recall} + \text{precision}} \right) \quad (12)$$

Recall

Recall is a measure of a true positive rate. In other words, the rate at which a healthy person is diagnosed or predicted to be healthy [37]. The recall can be calculated by applying (13).

$$\text{Recall} = \frac{TP}{TP + FN} \quad (13)$$

Specificity

Specificity is a true negative rate measure. In other words, the rate at which a person is diagnosed or predicted to have the disease [37]. The specificity can be calculated by applying (14).

$$\text{Specificity} = \frac{TN}{TN + FP} \quad (14)$$

4.4 Evaluation Result

Accuracy without Feature Processing

As noted in Table 2, the results obtained for calculating the accuracy of predicting coronary heart disease were not comparable without the use of feature processing techniques, as will be shown later in this research paper.

Table 2 Accuracy without feature processing

| Algorithms | Accuracy (%) |
|---------------|--------------|
| Decision Tree | 87.19 |
| Random Forest | 92.68 |
| MLP | 90.56 |
| KNN | 90.50 |
| Naive Bayes | 89 |

Decision Tree

Table 3 Decision Tree confusion matrix

| | No (Actual) | Yes (Actual) |
|---------------|-------------|--------------|
| No (Predict) | 2879 | 205 |
| Yes (Predict) | 76 | 103 |

Table 4 Decision Tree performance metrics

| Accuracy (%) | Precision (%) | F-Measure (%) | Recall (%) | Specificity (%) |
|--------------|---------------|---------------|------------|-----------------|
| 91.39 | 93.36 | 95.35 | 97.43 | 33.50 |

Random Forest

Table 5 Random Forest confusion matrix

| | No (Actual) | Yes (Actual) |
|---------------|-------------|--------------|
| No (Predict) | 2921 | 201 |
| Yes (Predict) | 107 | 34 |

Table 6 Random Forest performance metrics

| Accuracy (%) | Precision (%) | F-Measure (%) | Recall (%) | Specificity (%) |
|--------------|---------------|---------------|------------|-----------------|
| 92.80 | 93.57 | 96.13 | 98.85 | 34.80 |

MLP

Table 7 MLP Neural Network confusion matrix

| | No (Actual) | Yes (Actual) |
|---------------|-------------|--------------|
| No (Predict) | 208 | 2923 |
| Yes (Predict) | 100 | 32 |

Table 8 MLP Neural Network performance metrics

| Accuracy (%) | Precision (%) | F-Measure (%) | Recall (%) | Specificity (%) |
|--------------|---------------|---------------|------------|-----------------|
| 92.64 | 93.36 | 96.06 | 98.92 | 32.51 |

K-Nearest Neighbors

Table 9 K-Nearest Neighbors' confusion matrix

| | No (Actual) | Yes (Actual) |
|---------------|-------------|--------------|
| No (Predict) | 214 | 2930 |
| Yes (Predict) | 94 | 25 |

Table 10 K-Nearest Neighbors' performance metrics

| Accuracy (%) | Precision (%) | F-Measure (%) | Recall (%) | Specificity (%) |
|--------------|---------------|---------------|------------|-----------------|
| 92.68 | 93.20 | 96.08 | 99.15 | 30.53 |

Table 7 Naïve Bayes confusion matrix

| | No (Actual) | Yes (Actual) |
|---------------|-------------|--------------|
| No (Predict) | 2666 | 239 |
| Yes (Predict) | 69 | 289 |

Table 8 Naïve Bayes performance metrics

| Accuracy (%) | Precision (%) | F-Measure (%) | Recall (%) | Specificity (%) |
|--------------|---------------|---------------|------------|-----------------|
| 90.56 | 91.79 | 94.54 | 97.48 | 54.77 |

Accuracy Comparison

Table 9 Comparison between previous work and proposed work in terms of accuracy

| Al go ri thms | Previou s | Propo sed | Dat as et |
|---------------|---------------|-----------|-----------------------|
| Ev en t | | | |
| Decision Tree | 90 [21, 22] | 91.39 | Myocardial Infraction |
| Random Forest | 90.1 [21, 22] | 92.80 | Myocardial Infraction |
| MLP | 89 [22] | 92.64 | Myocardial Infraction |
| | 90.1 [21, 22] | 92.68 | Myocardial Infraction |
| Naïve Bayes | 89.9 [21, 22] | 90.56 | Angina Pectoris |

4.5 Discussion

In this research paper, we used a set of machine learning techniques to predict two events of coronary heart disease namely, Angina Pectoris (528 Yes, 2735 No), and Myocardial Infarction (308 Yes, 2955 No). Despite the previous researchers used many pre-processing techniques, the results obtained from this work were very encouraging compared to other studies that use the same data set to calculate accuracy as shown in Table 13.

It is noted that the techniques that have been used to improve the accuracy of machine learning models or classifiers in predicting coronary heart disease have proven effective and thus have achieved better results than previous research. For example, [21] and [22] used the same data set and obtained by applying the decision tree algorithm a predictive accuracy of 90% to predict coronary heart disease (CHD), while this research paper obtained an accuracy of 91.39%, with a positive increase of 1.39% as shown in Table 3 and Table 4. Also, this research paper and through the application of the random forest algorithm obtained a predictive accuracy of CHD 92.80%, shown in Table 5 and Table 6, which is higher than the result obtained in the decision tree algorithm in this research

paper on the one hand, and on the other hand, higher and better than the results obtained by [21] and [22] and that

Naïve Bayes

was 90.10%, with a positive increase of 2.7%. As for the use of the MLP algorithm in predicting CHD, researchers in [22] obtained an accuracy of predicting the disease 89%,

while this research paper obtained a better accuracy of 92.64%, with a positive increase of 3.64% shown in Table 7 and Table 8. Regarding the use of the KNN algorithm, researchers in [21] and [22] obtained a prediction accuracy of 90.10%, which is less than the prediction accuracy of the disease obtained in this research paper, which is 92.68%, which was applied to calculate the missing values and equal width discretization, with a positive increase of 2.58% as shown in Table 9 and Table 10. The application of the Naïve Bayes in this research paper obtained a predictive accuracy of coronary heart disease 90.56% as shown in Table 11 and Table 12, which is better than the predictive accuracy of 89.90% obtained in [21].

Although the results obtained in predicting coronary heart disease in terms of accuracy were not as significant as it should be, it may contribute to an increase in the number of cases with the correct diagnosis of the disease and at the same time reduce the number of cases that are incorrectly diagnosed with coronary heart disease and thus save lives.

5 Conclusion and Future Work

The heart is among the most important organs of the human body, as any problem with it can damage other important organs in the body, such as the brain. All doctors around the world warn of the sharp increase in the number of heart patients, being a serious disease that may lead to serious complications such as heart failure and cardiac arrest, both of which often lead to death if not diagnosed early.

In this paper, the researchers contributed to improving the accuracy of machine learning classification models in predicting two primary coronary heart disease events, namely, angina pectoris and myocardial infarction through the use of a number of feature processing techniques such as normalization, standardization, and discretization. For the purpose of validating the results obtained, the data set of the Framingham Heart Study was used with two main events (angina pectoris and myocardial infarction (heart attack)), due to its containment and after consulting with cardiologists about the most common factors causing coronary heart disease.

After using feature processing techniques on the dataset used, the accuracy of machine learning algorithms for predicting coronary heart disease improved unevenly. For example, the improvement in accuracy prediction of CHD was 4.2% when using the ID3 decision tree algorithm, 0.14% when using the random forest algorithm, 3.18% when using the KNN algorithm, 2.08% when using the MLP algorithm, and 1.36 when using the Naive Bayes algorithm as shown in Table 2 and Table 13. However, the best prediction accuracy obtained for the ID3 decision tree algorithm is at 91.39% when we applied the equal width discretization method. Whereas, the random forest algorithm achieved a prediction accuracy of 92.80% when we applied the equal width discretization and applied normalization methods. The MLP algorithm achieved an improvement in accuracy prediction by 92.64% when using one of the hot encoding

techniques. 92.68% represents the predictive accuracy obtained with the KNN algorithm when we applied the ordinal coding and standardization techniques. However, However, all of the predicted values obtained were in the case of a myocardial infarction event. Whereas, the value obtained from Naive Bayes algorithm was 90.65% in the case of angina pectoris and when we applied equal frequency discretization. The results obtained confirm the importance of using feature processing techniques in improving the accuracy performance of machine learning algorithms for predicting coronary heart disease compared to previous published research with the same objectives.

In the end, the presence of a correlation between some serious diseases such as the occurrence of stroke, high blood pressure, cardiovascular disease and coronary heart disease leads us in the future to predict such diseases and the effect of each of them on the occurrence of coronary heart disease on the one hand, and on the other hand the effect of the occurrence of coronary heart disease, on these diseases, to prevent death. This is because the patient in such cases does not have enough time to go to the doctor to see him and save his life.

Acknowledgements

The authors are grateful to the Applied Science Private University, Amman, Jordan, for the financial support granted to this research paper.

References

- [1] J. S. Sonawane and D. R. Patil, "Prediction of heart disease using multilayer perceptron neural network," in International Conference on Information Communication and Embedded Systems (ICICES2014), pp. 1–6, Feb 2014.
- [2] Sahoo, S., Subudhi, A., Dash, M. et al. Automatic Classification of Cardiac Arrhythmias Based on Hybrid Features and Decision Tree Algorithm. *Int. J. Autom. Comput.* 17, 551–561 (2020). <https://doi.org/10.1007/s11633-019-1219-2>.
- [3] S. A. Pattekari and A. Parveen, "Prediction system for heart disease using naive bayes," International Journal of Advanced Computer and Mathematical Sciences, vol. 3, no. 3, pp. 290–294, 2012.
- [4] K. H. Miao, J. H. Miao, and G. J. Miao, "Diagnosing coronary heart disease using ensemble machine learning," in (IJACSA) International Journal of Advanced Computer Science and Applications, vol. 7, pp. 30–39, 2016.
- [5] Know the differences cardiovascular disease, heart disease, coronary heart disease, national heart, lung, and blood institute, [Online], Available: http://www.nhlbi.nih.gov/sites/default/files/medi_a/docs/Fact_Sheet_Know_Diff_DesIgn_508.pdf, August 1, 2020.
- [6] R. Das, I. Turkoglu, and A. Sengur, "Effective diagnosis of heart disease through neural networks ensembles," *Expert Systems with Applications*, vol. 36, no. 4, pp. 7675 – 7680, 2009.
- [7] S. I. Ansarullah, P. K. Sharma, A. Wahid, and M. M. Kirmani, "Heart disease prediction system using data mining techniques: A study," *International Research Journal of Engineering and Technology* <http://www.xtgcydzjs.com>

(IRJET), vol. 3, no. 8, pp. 1375–1381, 2016.

[8] S. Bashir, U. Qamar, F. H. Khan, and M. Y. Javed, “An efficient rule-based classification of diabetes using id3, c4.5, and cart ensembles,” in *2014 12th International Conference on Frontiers of Information Technology*, pp. 226–231, Dec 2014.

[9] T. Panch, P. Szolovits, and R. Atun, “Artificial intelligence, machine learning and health systems,” *Journal of Global Health*, vol. 8, Oct. 2018.

[10] G. D. Magoulas and A. Prentza, “Machine learning in medical applications,” in *Machine Learning and Its Applications*, pp. 300–307, Springer Berlin Heidelberg, 2001.

[11] feature processing, amazon, [Online], Available: https://docs.aws.amazon.com/machine-learning/1/at/es/t/d/feature-processing_sing_ht_ml_, July 30, 2020.

[12] A. Gupta and V. Khathuria, “Framingham heart study,” *International Journal on Future Revolution in Computer Science & Communication Engineering*, vol. 4, no. 11, pp. 55–58, 2018.

[13] K. Nagendra and D. M. Ussenaiah, “Analysis of classification algorithms on heart diseases data using association rule mining,” *International Journal of Computational Engineering Research (IJCER)*, vol. 08, no. 6, pp. 39–46, 2018.

[14] K. V. Nagendra and D. Ussenaiah, “Support vector machine and neural network classification improved by bagging,” *International Journal on Future Revolution in Computer Science & Communication Engineering*, vol. 4, no. 2, pp. 125–130, 2018.

[15] J.-J. Beunza, E. Puertas, E. García-Ovejero, G. Villalba, E. Condes, G. Koleva, C. Hurtado, and M. F. Landecho, “Comparison of machine learning algorithms for clinical event prediction (risk of coronary heart disease),” *Journal of Biomedical Informatics*, vol. 97, p. 103257, Sept. 2019.

[16] A. S. T. Nishadi, “Predicting heart diseases in logistic regression of machine learning algorithms by python jupyterlab,” *International Journal of Advanced Research and Publications*, vol. 3, pp. 69–74, August 2019.

[17] A. Bhardwaj, A. Kundra, B. Gandhi, S. Kumar, A. Rehalia, and M. Gupta, “Prediction of heart attack using machine learning,” *IITM Journal of Management and IT*, vol. 10, pp. 20–24, 2019.

[18] A. Vallée, A. Cinaud, V. Blachier, H. Lelong, M. E. Safar, and J. Blacher, “Coronary heart disease diagnosis by artificial neural networks including aortic pulse wave velocity index and clinical parameters,” *Journal of Hypertension*, vol. 37, pp. 1682–1688, Aug. 2019.

[19] K. Lim, B. M. Lee, U. Kang, and Y. Lee, “An optimized DBN-based coronary heart disease risk prediction,” *International Journal of Computers Communications & Control*, vol. 13, pp. 492–502, July 2018.

[20] J. K. Kim and S. Kang, “Neural network-based coronary heart disease risk prediction using feature correlation analysis,” *Journal of Healthcare Engineering*, vol. 2017, pp. 1–13, 2017.

[21] N. S. Rajliwall, G. Chetty, and R. Davey, “Chronic disease risk monitoring based on an innovative predictive modelling framework,” in *2017 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1–8, 2017.

[22] N. S. Rajliwall, R. Davey, and G. Chetty, “Machine learning based models for cardiovascular risk prediction,” in *2018 International Conference on Machine Learning and Data Engineering (iCMLDE)*, pp. 142–148, 2018.

[23] Teaching Datasets, national heart, lung, and blood institute, [Online], Available: <https://bio lin cc.nih bi nih .g ov/teachi ng/>, July 30,

202
0.

[24] Ms.R.Malarvizhi and D. S. Thanamani, "K-nearest neighbor in missing data imputation," *International Journal of Engineering Research and Development*, vol. 5, pp. 5–7, 2012.

[25] C.Saranya and G.Manikandan, "A study on normalization techniques for privacy preserving data mining," *International Journal of Engineering and Technology (IJET)*, vol. 5, pp. 2701–2704, 2013. [26] K. Potdar, T. S., and C. D., "A comparative study of categorical variable encoding techniques for neural network classifiers," *International Journal of Computer Applications*, vol. 175, pp. 7–9, Oct.

201
7.

[27] H. Liu, F. Hussain, C. L. Tan, and M. Dash, "Discretization: An enabling technique," *Data Mining and Knowledge Discovery*, vol. 6, no. 4, pp. 393–423, 2002.

[28] C. B. C. Latha and S. C. Jeeva, "Improving the accuracy of prediction of heart disease risk based on ensemble classification techniques," *Informatics in Medicine Unlocked*, vol. 16, p. 100203,

201
9.

[29] A. H. Khaleel, G. A. Al-Suhail, and B. M. Hussan, "A weighted voting of k-nearest neighbor algorithm for diabetes mellitus," *International Journal of Computer Science and Mobile Computing*, vol. 6, no. 1, pp. 43–51, 2017.

[30] N. Ali, D. Neagu, and P. Trundle, "Evaluation of k-nearest neighbour classifier performance for heterogeneous data sets," *SN Applied Sciences*, vol. 1, Nov. 2019.

[31] Y. Zhang, Z. Lin, Y. Kang, R. Ning, and Y. Meng, "A feed-forward neural network model for the accurate prediction of diabetes mellitus," *INTERNATIONAL JOURNAL OF SCIENTIFIC & TECHNOLOGY RESEARCH*, vol. 7, no. 8, pp. 151–155, 2018.

[32] C. Nwankpa, W. Ijomah, A. Gachagan, and S. Marshall, "Activation functions: Comparison of trends in practice and research for deep learning," *CoRR*, vol. abs/1811.03378, 2018.

[33] Smith, A., Gu, F. & Ball, A.D. An Approach to Reducing Input Parameter Volume for Fault Classifiers. *Int. J. Autom. Comput.* 16, 199–212 (2019). <https://doi.org/10.1007/s11633-018-1162-7>.

[34] P. Refaeilzadeh, L. Tang, and H. Liu, *Cross-Validation*, pp. 532–538. Boston, MA: Springer US, 2009.

[35] R. Kohavi, "A study of cross-validation and bootstrap for accuracy estimation and model selection," in *Proceedings of the 14th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI'95*, (San Francisco, CA, USA), pp. 1137–1143, Morgan Kaufmann Publishers Inc., 1995.

[36] A. Kori, "Comparative study of data classifiers using rapidminer," *International Journal of Engineering Development and Research*, vol.

5, pp. 1041–1043, 2017.

[37] H. M and S. M.N, "A review on evaluation metrics for data classification evaluations," *International Journal of Data Mining & Knowledge Management Process*, vol. 5, pp. 01–11, Mar.

2015



Osmah Sami holds a bachelor's degree in software engineering from Isra'a University, Jordan in 2015, and a master's degree in computer science from the Applied Science Private University, Jordan in 2020. His research interests include machine learning, artificial intelligence, Internet of Things, big data analysis, and cloud computing

E-mail: osamah.sami@asu.edu.jo

ORCID iD: 0000-0001-5619-2687



Yousef Elsheikh is an associate professor of Information Technology at the Applied Science Private University. He holds a PhD in Informatics from the University of Bradford in the United Kingdom and MSc in Information Technology from the University of the West of England in the United Kingdom. He is currently

working as the head of the Software Engineering department at the Applied Science University. His research interests include conceptual modeling, e-business applications, information systems engineering, knowledge based representations, Ontologies and issues in software engineering such as UI and UX.

E-mail: y_elsheikh@asu.edu.jo



Fadi Almasalha An Associate Professor at the Faculty of Information Technology at Applied Science Private University in Amman, Jordan. He received his M.S. in computer Science from New York Institute of Technology, in 2005 and Ph.D. in Computer Science from University of Illinois at Chicago,

in 2011. In fall of 2011, he joined the Department of Computer Science at the Applied Science University as an Assistant professor. Dr. Fadi Almasalha received his Associate rank on 2016, during his appointment as the head of computer science department. Dr. Fadi has published more than 15 technical papers, journals and book chapters in refereed conferences and journals in the areas of multimedia systems, data mining, and cryptography.

E-mail: f_masalha@asu.edu.jo