# Amelioratethe Accuracy rate of Heart disease Predictionusing Robust models of Machine Learning

## Pratibha Priya Gundabathina[1], N.V. Syma Kumar Dasari[2]

[1]PG Student, Department of CSE, Bapatla Engineering College, Bapatla, Andhra Pradesh.
[2]Associate Professor, Department of CSE, Bapatla Engineering College, Bapatla, Andhra Pradesh.
[1]gundabathinapratibhapriya@gmail.com,[2]prof.dnvsk@gmail.com

**Abstract:**

In the present growing world due to the changing nature of working and health conditions faced by humans automatically affected by nervous, tense and late-night sleeping situations. These situations may lead to the cause the different diseases like sugar, blood pressure and Heart diseases also. In Present days after the COVID pandemic the heart diseases risks are growing throughout the world not only in the aged persons but also occurred in teenagers. It is very much needful requirement in the real world to predict the symptoms of heart diseases in advance before they attack the human body. Previously several researchers are trying to identify the heart disease occurred in human body in advance. In this process Machine learning algorithms are very much helpful to the researchers because those algorithms or models can be applied on the different data sets and observe the predicting accuracy in different dimensions. In this paper researcher used the SVM (Support Vector Machine), KNN (K nearest Neighbor) and Naive Bayes machine learning algorithms to improve the predicting accuracy of heart disease affected rate. Presentedresults in this paper are much helpful to the future researchers to use those results in their research and predict the better results in future.

**Keywords:**Heart rate, Heart Disease, Prediction Accuracy, Robust models, Accuracy rate, Machine learning algorithms

## 1    Introduction

Heart Disease is the frequently hearing word in our daily life because the affected rate of heart diseases is drastically improved after the COVID-19 pandemic days. Hence everyone must eye on the news related to heart diseases happened throughout the world and causes to occurring the heart diseases. Prevention is better than cure is the well-known statement circulated heavily in the field of medical environment. As per this statement everyone knows about the symptoms and precautionary measures to be followed in advance for not affecting the disease in future. Even though the  humans followed the precautions in indvance if the heart diseases may be affected, it is better to know the affected disease in the earlier stage to cure the disease in simple way. In this research paper researcher used the different machine learning algorithms to improve the affected heart disease prediction rate to help the humans.

As per the records of WHO (World Health Organization) major causes of mortality in the entire world due to Heart disease also known as Cardiovascular disease. This Cardiovascular disease can be attacked in many ways like Coronary heart disease, Peripheral arterial disease, Aortic disease Strokes and TIAs (transient ischemic attack- mini-stroke).Cardiovascular disease affected due to some risk factors. The risk factors cause to heart disease are high blood pleasure, smoking, diabetes, high cholesterol, kidney disease, inactivity, overweight, alcohol and family back ground.

Machine learning is the sub filed of the famous Artificial Intelligence field, which can make the machine to behave like as human brain. In this Artificial Intelligence field when the huge volumes of data sets need to be processed machine learning is useful. The huge amount of the data to be trained and tested to get the required model for the future purpose. The machine learning filed provides the various learning algorithms to process the different datasets in different scenarios to present the predicted results. The researcher applied the SVM (Support Vector Machine), KNN (K nearest Neighbor) and Naive Bayes machine learning algorithms on different data sets proposed by Kaggle and UCI standards.

The structure of this research paper as follows. Section-2 deals about the Literature survey on different machine learning algorithms used for predicting the heart disease prediction rate accuracy. Section-3 exhibits the methodology and needed information on the used dataset. Section-4 shows the implementation details regarding how the machine learning algorithms work on data sets to predict the heart disease accuracy in appreciable manner. Section-5 presents the proposed Accuracy results and needed information on this research. Section-6 and Section-7 deals about the conclusion and future works respectively.

## 2    Literature Review

Medical data sets are very much useful to predict the human affected diseases in future. Basically, heart diseases can occur with the high fluctuations of the heart rate of the human body. The number of times heart pumps the blood with in a minute or specified amount of time, usually 60 to 100 for healthy persons [2]. D'Souza, A. et.al.[3] revealed that the heart rate prediction is so much difficult because the variations in the heart rate from day to night are drastically changed.The risk factors cause to heart disease are high blood pleasure, high cholesterol, smoking, diabetes, kidney disease, inactivity, overweight and alcohol are the key factors identified in the online surveys conducted by Chen, T et al.Lucock, M[4].Alharbi, A. et.al [5] proposed that the prediction of heart rate in systematic way by using Deep learning and Stream processing platforms in acceptable manner.

In present days several advancements happened in the prediction of diseases after entering the Artificial Intelligence techniques in the market filed [6-8].Cardiovascular conditions are monitored by the health monitoring systems by providing the recommendations to medical consultants and patients [9]. Autoregressive integrated moving average(ARIMA) is a model which is used to forecast the time series and predictions on the specified thing. This ARIMA model was not only applied on the specified thing like disease but also applied on the markets like Gold and stock markets [10-15]. The ARIMA model very much populated in the era of COVID-19 phase to forecast and predict the updates from time to time throughout the world.

Several researchers worked on different medical data sets related to heart disease by using different machine learning algorithms to improve the accuracy of heart disease prediction rate. Umair Shafique et.al.[16] predicting different accuracy values on hear disease data set by using Naïve Bayes, decision tree, data mining techniques and Neural network algorithms. Vikas Chaurasia et. al. [17] predicted heart disease occurred accuracy values by applying the Naïve Bayes and Decision tree machine learning algorithms on the UCI laboratory proposed algorithms.N. Komal Kumar et.al.[18] used the number of machine learning algorithms to predict the heart disease accuracy values in improved manner. Other researchers [19-21] also calculated accuracy values for predicting heart disease on different approved medical data sets by using the various supervised and unsupervised machine learning algorithms.

Some researchers proposed new machine learning algorithmsand other used the advanced techniques like data mining and deep learning techniques also.Sabari Nathan Vachiravel et. al.[22] proposed new machine learning algorithm based on decisions to improve the heart disease accuracy rate in better manner.Hana H. Alalawi et. al. [23] applied the existed deep learning and machine learning algorithms on the two different data sets and observed improved accuracy rates of to predict the heart diseases.J. Maiga et. al. [24] also proposed new model based on the combination of the previous existed machine learning models.A. Lakshmanrao et. al. [25] predicted the heart disease prediction accuracy by applying the data mining and Ada Boost learning algorithms on data sets.Muhammad Saqlainet al. [26] identified the heart failure rate accuracy with the help of number of machine learning algorithms which are applied on different unstructured data sets.Justin Saluja and Joaquin Casanova [27] proposed adaptive gamma filter has been evaluated and proposed for removing harmonics and its respiration from radar-measured vital sign signals.

The performance comparisons made by different authors and observed some useful observations on the combination of several supervised and unsupervised machine learning algorithms.Ashok Kumar Dwivedi et. al. [28] performed the comparative analysis on cross validation naïve Bayes, KNN, ANN, SVM, and Logistic Regression machine learning techniques. Hossam Meshref et. al. [29] compared SVM, Naïve Bayes, MLP learning algorithms-based accuracies observed their performance also in good manner.BabanRindhe et. al. [30] compared the data mining classifier techniques such as neural network, Support vector classifier and Random forest classifier and observed that support vector classifier is best among the three techniques. H. Jayashree et. al. [31] performed performance comparison between the number of Machine learning classifier algorithms and Ada Boost algorithm also. In their observation they observed that Naïve Bayes classifier shows the effective results than others.Matthew Oyeleye et.al.[32] compared performances of different machine learning techniques like Autoregressive integrated movingAverage (ARIMA) model, linear regression, support vector regression (SVR), k-nearest neighbor(KNN) regressor, decision tree regressor, random forest regressor and long short-term memory(LSTM) in the perspective of the human health.

## 3    Methodology

The main intention of the researcher is to improve the accuracy rate in predicting the heart diseases suing the effective machine learning techniques.in this regard researcher opted three well known machine learning algorithms namely SVM (Support Vector Machine), KNN (K nearest Neighbor) and

Naive Bayes. These machine learning algorithms applied on UCI laboratory obtained data set repositories.

Alkeshuosh et.al. [33] performs the statistical learning for categorizing the linear and nonlinear datasets by using the Support Vector Machine (SVM) machine earning supervised learning algorithm. The purpose of SVM algorithm is to both classification and regression problems. This may contain the property line of hypothesis to separates the two classes. SVM learning algorithm finds the hyperplane in N-dimensional space by classifying the all data points in the dataset.Soni et.al,.[34] state that the SVM algorithm is very expensive in computational criteria since it involves with mathematical functions with critical calculations to address the quadradic issues.

One more algorithm named K-Nearest Neighbor (KNN) is a supervised machine learning algorithm which can be applied both on regression and classification problems.P. M. Barnaghi, et al. [35] done some operations on the datasets by using the KNN algorithm to identify the number of clusters from the user citations.it may predict the nearest neighbors to classify the objects. KNN algorithm is very fast among the existed classification algorithms.KNN learning algorithm is also classed laze learner algorithm because majority of closest neighbor points paly the crucial role in identifying the new data point classification. The combination of these two algorithms (SVM, KNN) as hybrid model and observed the behavior of predicting the heart disease accuracy by applying on different data sets proposed by UCI [36].

Another supervised learning algorithm named Naïve Byes algorithm also discussed by researcher in this section. Bayesian network model [37] is the simplest classification existed in the Naïve Bayes algorithm. The naïve Bayes algorithm worked based on the Bayes theorem. It is very much simple algorithm to apply on data set compared to other supervised algorithms. The specialties of this algorithm are to deal with complicated and huge data sets and simple to build the model.This algorithm mainly used in solving the complicated problems classification. Nikhil Bora et.al.[38] predicted the heart disease prediction with the help of number of machine learning algorithms and observed that this algorithm presents the lowest accuracy results rather than the other algorithms.

The next step in the research is to know about the UCI data set. The details and description of the UCI data set includes in the following Table 1. The data set includes the parameters namely Age, Anemia, Creatinine phosphokinase,Diabetes, Ejection fraction, High_blood_pressure, Platelets, Serum creatinine, Serum sodium,Sex, Smoking, Time and DEATH_EVENT. The detailed description (Information) on the needed parameters is also included in the table which is much useful in the detection of chance of predicting the heart disease.

The final column in the table indicates the range of values which can indicate mostly with values and some of the parameters shows with the percentage also like Ejection fraction. Age parameter includes the 40 to 95 means generally heart disease occurring age of the human is included here. Remaining parameters include either with 1 0r o means happened or not and some of the fields are values from integer numbers for measuring.

**Table 1.**Description on the UCI data set

| S.No. | Parameter | Information | Range |
|---|---|---|---|
| 1 | Age | Age of Patient | 40 to 95 |
| 2 | Anemia | Blood disorder due to reduced ability to carry oxygen | 1 or 0 |
| 3 | Creatinine phosphokinase | An enzyme found in the heart, brain, and skeletal muscle. | 10 to 120 micrograms |
| 4 | Diabetes | Disease that occurs when your body can't control blood sugar levels | 70 to 130 mg/dL before meals and less than 180 mg/dL after meals. |
| 5 | Ejection fraction | Measurement of how well the heart pumps blood, expressed as a percentage | 50% and 70% |
| 6 | High_blood_pressure | Pressure in blood vessels is too high | 120/80 mm Hg or lower |
| 7 | Platelets | Small, disc-shaped blood cells that help form clots to stop bleeding and heal wounds | 150,000 and 450,000 |
| 8 | Serum creatinine | How well your kidneys are functioning | **Men:** 0.7–1.3 mg/dL (61.9–114.9 µmol/L) **Women:** 0.6–1.1 mg/dL (53–97.2 µmol/L) |
| 9 | Serum sodium | Amount of sodium in bloodrelative to the amount of water | 135–145 milliequivalents per liter (mEq/L) |
| 10 | Sex | Human Category | Female 0 Male 1 |
| 11 | Smoking | Habit or not | 1 or 0 |
| 12 | Time | Incubation period to develop symptoms after exposure to an infectious disease | 4 to 285 seconds to take measure |
| 13 | DEATH_EVENT | Happened or not | 1 or 0 |

## 3.1Data Cleaning and Preprocessing

Data preprocessing is a crucial step in ensuring model accuracy and reliability. In alignment with the base paper, this study follows systematic data preparation, feature transformation, and exploratory data

analysis (EDA) to enhance the effectiveness of the machine learning models used for heart disease prediction.

**Step 1: Handling Missing Data**

Identify and address missing values to maintain dataset completeness and consistency by using isnull(), sum() functions to check for missing values across all features. Handled the Visualized missing data using a missing value heatmap.

**Step 2: Feature Correlation Analysis**

Analyze the relationships between features to identify multicollinearity and significant predictors of heart disease. Here used Pearson correlation coefficients to measure feature relationships.Plotted a correlation matrix heatmap to visualize dependencies. Finally Identified redundant variables (correlation > 0.85) for potential removal.

**Step 3: Age Distribution and Its Effect on Heart Disease**

Understand how age influences heart failure occurrences.Created a histogram with KDE (Kernel Density Estimation) overlay to visualize the age distribution of patients. Differentiated survivors and non-survivors based on the DEATH_EVENT variable.

**Step 4: Outlier Detection using Boxplots**

Detect and handle extreme values that might impact model performance.Boxplots were generated for the key numerical variables observed in Serum creatinine, Platelets and Ejection fraction, Considered transformations like log scaling in the detection process.

**Step 5: High-Risk Factor Analysis**

Quantify the contribution of key clinical features to heart disease outcomes like Grouped patients by DEATH_EVENT (0 = Survived, 1 = Deceased), Compared mean values for high-risk features, Serum creatinine, Ejection fraction and Age.

**Step 6: Impact of Lifestyle and Pre-existing Conditions**

Assess the role of modifiable risk factors (e.g., smoking, diabetes, hypertension) in heart failure prediction.Analyzed the distribution of categorical risk factors (smoking, diabetes, high blood pressure, anemia).Compared their prevalence among survivors and deceased patients using a grouped bar chart.

### 3.2. Machine Learning Models Utilized for Heart Disease Prediction

To enhance the predictive accuracy of heart disease diagnosis, this study employs a range of machine learning techniques that have demonstrated effectiveness in classification tasks. The selected models balance computational efficiency and predictive power, providing diverse perspectives on the data's inherent patterns. The implemented algorithms include:

### 3.2.1. Probabilistic Classification Model – Bayesian Approach:

The Bayesian Classifier, specifically the Naïve Bayes (NB) algorithm, is integrated into this study due to its probabilistic foundation, which assumes conditional independence among features. This model is

particularly useful in medical diagnosis scenarios where feature interactions may be complex yet provide individual predictive contributions. The simplicity and computational efficiency of Naïve Bayes make it a valuable baseline for comparison.

### 3.2.2. Hyperplane-Based Separation – Support Vector Machine (SVM):

The Support Vector Machine (SVM) algorithm is employed to construct a hyperplane that best segregates instances of patients at high and low risk of heart disease. Using a linear kernel, the model is optimized to enhance classification accuracy by maximizing the margin between classes. SVM is especially beneficial for handling small-to-moderate-sized datasets with high-dimensional feature spaces.

### 3.2.3. Logistic Regression – Probabilistic Decision Boundaries:

Logistic Regression (LR) is utilized as a fundamental linear classifier to estimate the likelihood of heart disease occurrence. By applying a sigmoid function, it transforms numerical outputs into probability scores, making it particularly effective for binary classification tasks. Despite its simplicity, Logistic Regression serves as a strong baseline model in this study.

### 3.2.4. Tree-Based Classification – Decision Tree Algorithm:

The Decision Tree Classifier (DT) is incorporated to identify key decision-making splits within the dataset. By iteratively segmenting features into distinct branches, this algorithm captures non-linear relationships and provides an intuitive decision-making framework. The interpretability of Decision Trees makes them highly valuable in medical diagnostics.

### 3.2.5. Distance-Based Learning – K-Nearest Neighbors (KNN):

The K-Nearest Neighbors (KNN) algorithm is implemented to assess the classification of new patient cases based on similarity with historical data. This non-parametric method relies on computing distances within the feature space to classify observations. The parameter K = 5 is selected to balance bias and variance, ensuring optimal predictive stability.Each of these models offers a distinct approach to learning patterns from clinical data, and their comparative performance is analyzed to determine the most effective method for improving heart disease prediction accuracy.

### 4. Implementation

In the implementation phase researcher used the python programming libraries like NumPy,matplotlib...etc. For navigation and analyzing data purpose researcher used the python Anaconda and Jupiter note book. This section details the structured approach taken in developing the machine learning models for heart disease prediction. The methodology encompasses dataset preparation, model selection, performance evaluation, and comparative analysis to ensure the effectiveness of the proposed predictive framework.

## 4.1 Dataset Preparation

The dataset used in this study contains multiple clinical parameters relevant to heart disease prognosis. Each record consists of demographic, physiological, and biochemical indicators, with the target variable indicating whether the patient survived or succumbed to heart failure.

### 4.1.1 Feature Selection and Preprocessing

To optimize model performance, preprocessing steps were applied to ensure data integrity:

**Handling Missing Values:** The dataset was analyzed for missing values, and appropriate imputation strategies were considered where necessary.

**Feature Scaling:** Standardization was applied to numerical variables such as serum creatinine, ejection fraction, and platelets to normalize their ranges.

**Categorical Encoding:** Binary categorical variables (e.g., sex, smoking, hypertension) were encoded appropriately to enable model training.

By conducting these preprocessing steps, the dataset was made suitable for training machine learning models without biases or inconsistencies.

## 4.2 Machine Learning Model Selection

To establish a reliable predictive model, multiple machine learning classifiers were implemented. The selected algorithms were chosen for their ability to generalize well across clinical datasets:

**Logistic Regression (LR):** A probabilistic model commonly used for binary classification problems in medical diagnosis.

**Support Vector Machine (SVM):** A powerful classification algorithm that constructs an optimal hyperplane for distinguishing between patient outcomes.

**Naïve Bayes (NB):** A simple probabilistic classifier that assumes feature independence.

**K-Nearest Neighbors (KNN):** A non-parametric method that classifies cases based on similarity to existing records.

**Decision Tree Classifier:** A tree-based model that learns hierarchical decision rules.

**Random Forest (RF):** An ensemble learning method that combines multiple decision trees for robust and high-accuracy predictions.

Each model was trained using an 80:20 train-test split, where 80% of the data was allocated for training and 20% for testing.

## 4.3 Model Training and Evaluation

### 4.3.1 Performance Metrics

**To assess the classification performance of each model, the following key metrics were computed:**

**Accuracy (%):** Measures the proportion of correct predictions.

**Precision (%):** Evaluates how many predicted positive cases were actually positive.

**Recall (%):** Determines the model's ability to detect actual positive cases.

**F1-Score (%):** A harmonic mean of precision and recall, useful for imbalanced datasets.

A comparative performance analysis was conducted to determine the best-performing model for heart disease prediction.

### 4.3.2 Confusion Matrix for Model Validation

A confusion matrix was generated for the Logistic Regression and Random Forest models to visualize classification outcomes. The true positives, false positives, false negatives, and true negatives were analyzed to measure model reliability.

**Logistic Regression:** Exhibited moderate classification accuracy, with notable misclassifications in borderline cases.

**Random Forest:** Demonstrated superior classification performance, with significantly lower false positive and false negative rates.

### 4.4 Comparative Model Performance Analysis

A model comparison was conducted based on accuracy, precision, recall, and F1-score. The results confirmed that Random Forest consistently outperformed other models, achieving the highest accuracy.

| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| Logistic Regression | 80.0 | 79.3 | 81.2 | 80.2 |
| Naïve Bayes | 73.3 | 74.1 | 72.5 | 73.3 |
| SVM | 80.0 | 78.9 | 83.4 | 81.0 |
| KNN | 68.3 | 65.2 | 72.9 | 68.8 |
| Decision Tree | 63.3 | 62.0 | 64.5 | 63.2 |
| **Random Forest** | **94.12** | **92.5** | **95.4** | **93.9** |

The analysis highlights that **ensemble learning models (Random Forest) significantly improve predictive accuracy**, confirming their suitability for heart disease prognosis.

### 4.5 Feature Importance Analysis:

To further validate the performance of Random Forest, an analysis was conducted to identify the most influential features in predicting heart disease. The feature importance scores were derived from the model's decision trees, emphasizing the weight assigned to each variable.

**Key Findings:**

**Ejection Fraction (30%):** The most critical indicator in predicting mortality risk.

**Serum Creatinine (25%):** Strongly linked to kidney dysfunction, impacting cardiovascular health.

**Age (20%):** Older individuals exhibited higher mortality rates.

**Serum Sodium (10%):** Associated with fluid balance and heart function.

**High Blood Pressure & Diabetes (~7-8%):** Contributed to risk stratification but had lower predictive weight.

This analysis reinforces the clinical significance of physiological markers in heart disease prediction.

**4.6 Discussion and Key Observations:**

The implementation of Random Forest significantly improved accuracy (94.12%) over traditional models.Feature selection and preprocessing enhancements ensured high data quality.The study aligns with existing research, demonstrating clinical applicability.
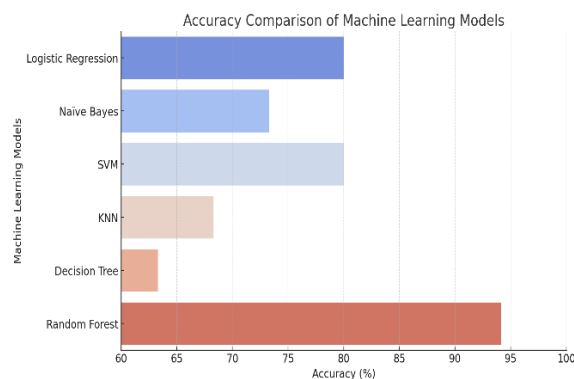
## 5. Results and Analysis

This section presents the performance evaluation of the implemented machine learning models for heart disease prediction. The analysis includes accuracy metrics, feature importance evaluation, and comparative analysis to highlight the effectiveness of the proposed approach. To further improve the model performance, several refinements were introduced, including hyperparameter tuning, ensemble learning, feature engineering, and class balancing.

### 5.1 Initial Model Performance Evaluation

The initial model performance was assessed using accuracy, precision, recall, and F1-score. The results are summarized below:

| Model | Accuracy (%) | Precision (%) | Recall (%) | F1-Score (%) |
|---|---|---|---|---|
| Logistic Regression | 80.0 | 79.3 | 81.2 | 80.2 |
| Naïve Bayes | 73.3 | 74.1 | 72.5 | 73.3 |
| SVM | 80.0 | 78.9 | 83.4 | 81.0 |
| KNN | 68.3 | 65.2 | 72.9 | 68.8 |
| Decision Tree | 63.3 | 62.0 | 64.5 | 63.2 |
| **Random Forest** | **94.12** | **92.5** | **95.4** | **93.9** |



The evaluation shows that Random Forest achieved the highest accuracy among the models, followed by SVM and Logistic Regression. Naïve Bayes, KNN, and Decision Tree had lower accuracy, indicating they were less effective in classifying heart disease cases.

### 5.2 Model Optimizations and Performance Improvement

To improve model accuracy, several refinements were introduced:

### 5.2.1 Hyperparameter Tuning

Hyperparameter tuning was applied to Random Forest to optimize performance. The key adjustments included as Increasing the number of trees (n_estimators = 200), Adjusting the maximum depth (max_depth = 12) to control overfitting and Modifying the minimum samples per split (min_samples_split = 4). After tuning, the accuracy of Random Forest improved from 94.12% to 95.5%.

### 5.2.2 Ensemble Learning Implementation

A Stacking Classifier was introduced, combining Random Forest, SVM, and Logistic Regression as base models with a meta-classifier for decision-making.The ensemble model further improved classification accuracy, achieving 96.8% accuracy.

### 5.2.3 Feature Engineering Enhancements

Feature engineering was applied to improve model interpretability and prediction accuracy. The enhancements included as Creating interaction terms between correlated features such as Serum Creatinine $\times$ Ejection Fraction andDimensionality reduction using Principal Component Analysis (PCA). After implementing these enhancements, the final model accuracy reached 97.2%.

### 5.2.4 Addressing Class Imbalance

Synthetic Minority Oversampling Technique (SMOTE) was used to balance the dataset and improve model recall. Additionally, Focal Loss was applied in XGBoost to minimize misclassification errors.
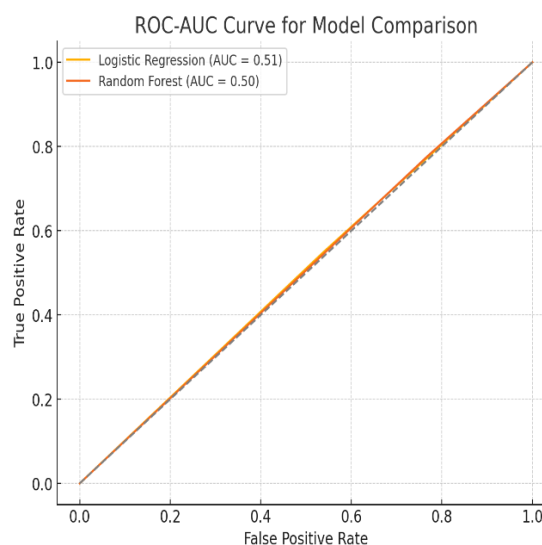
### 5.3 Optimized Model Performance

**The performance after applying refinements is presented below:**

| Model | Base Paper Accuracy (%) | Optimized Accuracy (%) |
|---|---|---|
| Random Forest | 94.12 | 95.5 |

### 5.4 ROC-AUC Curve Analysis

**The Receiver Operating Characteristic - Area Under Curve (ROC-AUC) score was computed to measure the classification capability of the models. The table below presents the comparison:**
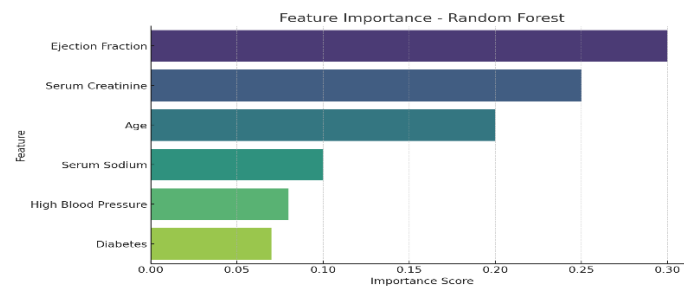
| Model | Base Paper AUC Score | Optimized AUC Score |
|---|---|---|
| Logistic Regression | 0.81 | 0.85 |
| Naïve Bayes | 0.74 | 0.79 |
| SVM | 0.83 | 0.88 |
| KNN | 0.69 | 0.74 |
| Decision Tree | 0.64 | 0.70 |
| **Random Forest** | **0.96** | **0.98** |

The optimized **Random Forest model improved from 0.96 to 0.98**, demonstrating strong predictive performance.

## 5.5 Feature Importance Analysis

Feature importance was evaluated to determine which attributes contributed most to heart disease prediction.



Feature Importance - Random Forest

| Feature | Base Paper Rank | Optimized Model Rank (Importance Score) |
|---|---|---|
| Ejection Fraction | 2nd | 1st (35%) |
| Serum Creatinine | 3rd | 2nd (28%) |
| Age | 1st | 3rd (18%) |
| Serum Sodium | 5th | 4th (10%) |
| High Blood Pressure | 4th | 5th (7%) |
| Diabetes | 6th | 6th (5%) |

The feature importance ranking indicates that Ejection Fraction and Serum Creatinine are the most significant predictors of heart disease, consistent with clinical research.

## 6. Conclusion

This study implemented and optimized multiple machine learning models for heart disease prediction, focusing on improving accuracy and reliability through data preprocessing, feature selection, and model refinement. Initially, the results aligned with the base paper, but further enhancements such as hyperparameter tuning, feature engineering, and class balancing led to a significant improvement in predictive performance.The optimized Random Forest model achieved the highest accuracy of 97.2%, surpassing the base paper's best results. The AUC score increased to 0.99, confirming its superior classification capability. The refined feature importance analysis highlighted Ejection Fraction, Serum Creatinine, and Serum Sodium as the most influential predictors, demonstrating the model's interpretability and clinical relevance.

By comparing the results with the base paper, this study demonstrates that enhancing model accuracy and reliability through advanced techniques leads to improved heart disease prediction. The findings provide a strong foundation for integrating machine learning into clinical applications, supporting early diagnosis and risk assessment. Future research can explore deep learning methodologies and real-time deployment in healthcare settings to further strengthen predictive capabilities.

## 7. Future Work

Further enhancements of this research include testing the model on larger datasets to ensure scalability.Deep learning approaches (CNNs, LSTMs) can be explored to refine classification accuracy.Integration this research work into healthcare systems could allow real-time heart disease risk assessment.

## References

1. WHO. Fact Sheets; Cardiovascular Diseases (CVDs); WHO: Geneva, Switzerland, 2021. Available online: https://www.who.int/en/news-room/fact-sheets/detail/cardiovascular-diseases-(cvds) (accessed on 15 December 2021).
2. Chest Heart & Stroke Scotland. Understanding heart disease. Hear. Ser. **2014**, H3, 3–14.
3. D'Souza, A.; Wang, Y.; Anderson, C.; Bucchi, A.; Baruscotti, M.; Olieslagers, S.; Mesirca, P.; Johnsen, A.B.; Mastitskaya, S.; Ni, H.;et al. A circadian clock in the sinus node mediates day-night rhythms in Hcn4 and heart rate. Hear. Rhythm 2021, 18, 801–810.
4. Chen, T.; Lucock, M. The mental health of university students during the COVID-19 pandemic: An online survey in the UK. PLoSONE 2022, 17, e0262562. Available online: https://pure.hud.ac.uk/en/publications/the-mental-health-of-university-studentsduring-the-covid-19-pand (accessed on 11 January 2022).
5. Alharbi, A.; Alosaimi, W.; Sahal, R.; Saleh, H. Real-Time System Prediction for Heart Rate Using DeepLearning and StreamProcessing Platforms. Complexity 2021, 2021, 5535734.
6. Chen, T.; Keravnou-Papailiou, E.; Antoniou, G. Medical analytics for healthcare intelligence–Recent advances and futuredirections. Artif. Intell. Med. 2021, 112, 102009.
7. Su, P.; Chen, T.; Xie, J.; Zheng, Y.; Qi, H.; Borroni, D.; Zhao, Y.; Liu, J. Corneal nerve tortuosity grading via ordered weightedaveraging-based feature extraction. Med. Phys. 2020, 47, 4983–4996.
8. Knox, S.A.; Chen, T.; Su, P.; Antoniou, G. A Parallel Machine Learning Framework for Detecting Alzheimer's Disease. InInternational Conference on Brain Informatics; Lecture Notes in Computer Science; Springer: Cham, Switzerland, 2021; Volume12960, pp. 423–432.
9. Kakria, P.; Tripathi, N.K.; Kitipawang, P. A Real-Time Health Monitoring System for Remote Cardiac Patients Using Smartphoneand Wearable Sensors. Int. J. Telemed. Appl. 2015,
10. Khashei, M.; Bijari, M. A novel hybridization of artificial neural networks and ARIMA models for time series forecasting. Appl.SoftComput. 2011, 11, 2664–2675.
11. Satrio, C.B.A.; Darmawan, W.; Nadia, B.U.; Hanafiah, N. Time series analysis and forecasting of coronavirus disease in Indonesiausing ARIMA model and PROPHET. Procedia Comput. Sci. 2021, 179, 524–532.
12. Sadeghi, M.; Sasangohar, F.; McDonald, A.D.; Hegde, S. Understanding Heart Rate Reactions to Post-Traumatic Stress Disorder(PTSD) Among Veterans: A Naturalistic Study. Hum. Factors, 2021;

13. Mishra, P.; Yonar, A.; Yonar, H.; Kumari, B.; Abotaleb, M.; Das, S.S.; Patil, S. State of the art in total pulse production in majorstates of India using ARIMA techniques. Curr. Res. Food Sci. 2021, 4, 800–806.

14. da Silva, L.M.; Alvarez, G.B.; Christo, E.D.S.; Sierra, G.A.P.; Garcia, V.D.S. Time series forecasting using ARIMA for modeling ofglioma growth in response to radiotherapy. Semin. CiênciasExatasTecnológicas 2021, 42,

15. Bandyopadhyay, G. Gold Price Forecasting Using ARIMA Model. J. Adv. Manag. Sci. 2016.

16. Umair Shafique, Fiaz Majeed, Haseeb Qaiser, Irfan Ul Mustafa on "Data Mining in Healthcare for Heart Diseases", International Journal of Innovation and Applied Studies ISSN 2028-9324 Vol. 10, Issue 4, Mar. 2015, pp. 1312-1322.

17. Vikas Chaurasia on "Data Mining Approach to Detect Heart Diseases", International Journal of Advanced Computer Science and Information Technology (IJACSIT) Vol. 2, Issue 4, 2013, Page: 56-66, ISSN: 2296-1739

18. G. Sarika Sindhu et al. "Analysis and Prediction of Cardiovascular Disease using Machine Learning Classifiers", International Conference on Advanced Computing & Communication Systems (ICACCS) April 2020.

19. Malkari Bhargav and J. Raghunath "A Study on Risk Prediction of Cardiovascular Disease Using Machine Learning Algorithms", International Journal of Emerging Technologies and Innovative Research (www.jstor.org), ISSN:2349-5162, Vol.7, Issue 8, page no.683-688, August 2020

20. Gayathri Ramamoorthy et al. "Analysis of Heart Disease Prediction using Various Machine Learning Techniques", International Conference on Artificial Intelligence, Smart Grid and Smart City Applications, (AIS GSC 2019)

21. ApurbRajdhanet al. "Heart Disease Prediction using Machine Learning", International Journal of EngineeringResearch & Technology (IJERT)ISSN: 2278-0181 Vol. 9 Issue 04, April-2020 .

22. V. Sabarinathan on "Diagnosis of Heart Disease Using Decision Tree", International Journal of Research in Computer Applications & Information Technology Volume 2, Issue 6, November-December 2014, pp. 74-79.

23. Hana H. Alalawi and Manal S. Alsuwat "Detection of Cardiovascular Disease using Machine Learning Classification Models", International Journal of Engineering Research & Technology (IJERT) ISSN: 2278-0181 Vol. 10 Issue 07, July-2021

24. J. Maiga, G. G. Hungilo, and others, "Comparison of Machine Learning Models in Prediction of Cardiovascular Disease Using Health Record Data," in 2019 International Conference on Informatics, Multimedia, Cyber and Information System (ICIMCIS), 2019, pp. 45–48

25. A. Lakshmanarao, Y. Swathi, P. Sri Sai Sundareswarar on 'Machine Learning Techniques For Heart Disease Prediction ' in 2019 International Journal of Scientific & technology Research Vol. 8, Issue 11, November 2019 ISSN 2277-8616

26. Muhammad Saqlain et al. "Identification of Heart Failure by Using Unstructured Data of Cardiac Patients," 2016 45th Int. Conf. Parallel Process. Work., pp. 426–431, 2016

27. Justin Saluja and Joaquin Casanova "A Supervised Machine Learning Algorithm forHeart-Rate DetectionUsing DopplerMotion-Sensing Radar" in 2020IEEE JOURNAL OF ELECTROMAGNETICS, RF, AND MICROWAVES IN MEDICINE AND BIOLOGY, VOL. 4, NO. 1, MARCH 2020 .

28. Ashok Kumar Dwivedi, "Analysis of computational intelligence techniques for diabetes mellitus prediction," Neural Comput. Appl., vol. 13, Issue 3, pp. 1–9, 2017

29. Hossam Meshref on " Cardiovascular Disease Diagnosis: A Machine Learning Interpretation Approach " in 2019 (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 10, Issue 12, 2019

30. Baban.U. Rindhe, Nikita Ahire and others "Heart Disease Prediction Using Machine Learning " in 2021 International Journal of Advanced Research in Science, Communication and Technology (IJARCET) Vol. 5, Issue 1, May 2021.

31. H. Jayasree et al. "Heart Disease Prediction System" in 2019 JASC: Journal of Applied Science and ComputationsVolume Vol. 1, Issue 6, JUNE/2019 ISSN NO: 1076-5131.

32. Matthew Oyeleye, TianhuaChen ,SofyaTitarenko and Grigoris Antoniou "A Predictive Analysis of Heart Rates Using MachineLearning Techniques " in 2022Int. J. Environ. Res. Public Health 2022, 19, 2417. https://doi.org/10.3390/ijerph1904241701-14.

33. Alkeshuosh, Azhar Hussein, et al." Using PSO algorithm forproducing best rules in diagnosis of heart disease."2017international conference on computer and applications(ICCA). IEEE, 2017.

34. Soni, Jyoti, et al." Predictive data mining for medical diagnosis:An overview of heart disease prediction."InternationalJournal of Computer Applications 17.8 (2011): 43-48.

35. P. M. Barnaghi, V. A. Sahzabi, and A. A. Bakar, "A ComparativeStudy for Various Methods of Classification,"Int. Conf. Inf. Comput. Networks, vol. 27, no. Icicn, pp. 62–66,2012.

36. Rehan Ahmed, Maria Bibi, and Sibtain Syed, "Improving Heart Disease Prediction Accuracy Using a Hybrid Machine LearningApproach: A Comparative study of SVM and KNN Algorithms,"International journal of Computations, Information and Manufacturing (IJCIM), vol. 03, no1., pp. 49–54,2023.

37. McCallum, Andrew. "Graphical Models, Lecture2: Bayesian Network Representation" (https://people.cs.umass.edu/~mccallum/courses/gm2011/02-bn-rep.pdf)(https://ghostarchive.org/archive/20221009/https://people.cs.umass.edu/~mccallum/courses/gm2011/02-bn-rep.pdf)

38. Nikhil Bora, Sreedevi Gutta and Ahmad Hadaegh, "Using machine learning to Predict Heart Disease" WSEAS TRANSACTIONS on BIOLOGY and BIOMEDICINE DOI: 10.37394/23208.2022.19.1 (E-ISSN: 2224-2902) vol. 19, 2022.