

# Comparative Study Using Random Forest for Heart Disease Prediction

**Kommaraju Gaurav Kalyan**

Vardhaman College of Engineering, Telangana

kommarajugauravkalyan@gmail.com

**Abstract:** Heart disease is still the predominant killer worldwide, highlighting the tremendous demand for specific early diagnostic techniques. This work investigates the performance of RF model in predicting heart disease, using UCI Heart Disease dataset. The predictive score was achieved by the RF model being 85%, and the precision and recall scores were close to 90%, which was better than many classical models of predication. Feature importance analysis revealed predictive factors of clinical importance, such as chest pain type and ST depression. The model also achieved powerful classification ability and did not require hyperparameter tuning. Considering the accuracy, stability and interpretability, Random Forest can be a reliable one for early cardiovascular risk prediction, and is very promising to be further applied to clinical decision-making support system

**Keywords:** Heart Disease Prediction, Random Forest Algorithm, Cardiovascular Risk Assessment

## I. INTRODUCTION

Heart disease is the leading cause of death globally, accounting for approximately 17.9 million deaths annually. Early prediction is crucial for improving outcomes, but traditional clinical methods may overlook complex patterns in patient data. Machine learning, particularly the Random Forest (RF) algorithm, offers a powerful alternative due to its ability to model nonlinear relationships and reduce overfitting through ensemble learning. Recent studies have shown RF achieving high accuracy in heart disease prediction ranging from 84% to as high as 98% in optimized settings often outperforming other models like logistic regression, SVM, and KNN. Additionally, RF provides feature importance rankings, aiding clinical interpretability. This study evaluates the performance of a Random Forest classifier on the UCI Heart Disease dataset from Kaggle, using accuracy as the primary metric. Our objectives include literature review, model development, and performance analysis to assess RF's suitability as a predictive tool in cardiology and validate its clinical relevance and effectiveness.

## II. LITERATURE REVIEW

Several studies have demonstrated the effectiveness of machine learning models for heart disease prediction, with Random Forest (RF) consistently emerging as a top-performing algorithm. A comparative evaluation involving logistic regression, support vector machines (SVM), k-nearest neighbors (KNN), and Random Forest showed that while logistic regression and KNN achieved approximately 81% accuracy, the Random Forest model outperformed all, attaining an F1-score of 95%, with precision and recall values of 96% and 97%, respectively [6]. In a related study, feature extraction and optimization methods were implemented alongside an RF classifier, achieving 90% accuracy with 100% sensitivity, underscoring RF's potential in detecting cardiac abnormalities [7].

The impact of model tuning was further demonstrated using principal component analysis (PCA) and grid search to optimize an RF classifier, with evaluation metrics all exceeding 90% [4]. Beyond standard implementations, researchers have developed hybrid and ensemble RF variants. For instance, an enhanced RF model optimized using the particle swarm social optimization (PSSO) algorithm achieved 98.7% accuracy [5]. Similarly, a study applying genetic



algorithms, particle swarm optimization, and ant colony optimization to fine-tune RF hyperparameters found that the Genetic Algorithm Optimized RF model achieved the highest accuracy among the tested variants [1]. While deep learning and gradient boosting techniques (e.g., XGBoost) have also been applied to heart disease prediction, their complexity and lower interpretability can pose challenges. Random Forest strikes a practical balance between performance and interpretability [3]. Its ability to provide feature importance rankings aligns well with clinical reasoning, and its ensemble nature ensures robustness against overfitting. Collectively, these studies affirm the value of Random Forest as a reliable and effective tool in cardiovascular risk modeling.

### III. METHODOLOGY

**Dataset:** We used the UCI Heart Disease dataset (Cleveland) obtained from Kaggle. This dataset contains 303 patient records with 13 input features and 1 target label. Key features include: age, sex (0 = female, 1 = male), chest pain type (4 categories encoded as 0–3), resting blood pressure, serum cholesterol, fasting blood sugar (binary), resting ECG results (0–2), maximum heart rate achieved, exercise-induced angina (0/1), ST depression (old peak), slope of ST segment (0–2), number of major vessels colored (ca), and thalassemia (thal). The target is a binary outcome: 1 indicates the presence of heart disease, and 0 indicates absence. Before modeling, we performed basic preprocessing. The dataset has no missing values for the features used. Categorical features are provided as integers in the dataset, except for the **thal** feature which appears as text categories ("normal", "fixed", "reversible") in some versions. In our preprocessing, we encoded any non-numeric categorical values to numeric codes (e.g., encoding "normal", "fixed", "reversible" to 0, 1, 2 respectively) using a label encoder. The binary sex feature was kept as 0/1, and we preserved the original encoding of other categorical features since tree-based models can handle them without one-hot encoding.

**Experimental Setup:** We implemented the Random Forest model using Python's scikit-learn library. The dataset was split into training and testing subsets using an 80/20 train-test split (stratified to maintain the proportion of classes). No over-sampling was applied; the dataset is moderately balanced (approximately 54% with heart disease, 46% without). We chose accuracy as the primary evaluation metric [2]. Additionally, we examined precision, recall, F1-score, and the confusion matrix for a more detailed performance assessment, but these are secondary to accuracy in our evaluation. The Random Forest classifier was initialized with 100 decision trees (`n_estimators=100`) and a fixed random seed for reproducibility. We did not apply extensive hyperparameter tuning in this study, in order to evaluate the baseline capability of Random Forest with default parameters. (In practice, techniques like grid search or randomized search could be used to tune the number of trees, tree depth, splitting criteria, and other hyperparameters to potentially improve performance. However, prior research has noted that even without aggressive tuning, Random Forest often performs well out-of-the-box.

**Workflow:** Our methodology can be summarized in the following steps:

**Data Loading & Preprocessing:** Loaded the CSV dataset into a pandas Data Frame. Verified data types and encoded the **thal** categorical feature into numeric form (using Label Encoder). Ensured that the target is binary (0/1). Figure 1 shows the distribution of the sex feature in the dataset (Male vs Female) to give an idea of class balance in one of the features.

**Training/Testing Split:** Divided the data into training (242 samples) and testing (61 samples) sets using an 80:20 split. The split was stratified on the target to maintain equal class proportions in train and test sets.

**Model Training:** Initialized a Random Forest classifier with 100 trees (`max_depth` was left unlimited to allow trees to grow fully, and other parameters at sklearn defaults). Trained the RF on the training set (`X_train, y_train`). The training process involves bootstrap sampling and tree growth: each decision tree is grown on a bootstrap sample of the training data, and splits are determined by maximizing information gain (Gini impurity reduction). The ensemble combines these trees' predictions by majority voting.



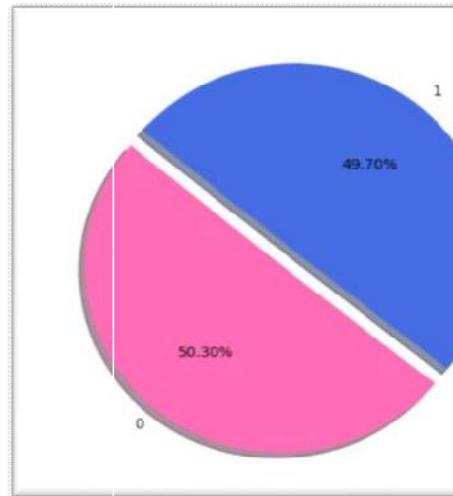


Figure 1: Gender distribution in the Heart Disease dataset (Male=1, Female=0). The pie chart shows the proportion of male and female patients in the dataset.

**Model Prediction:** Used the trained Random Forest to predict outcomes on the test set ( $X_{test}$ ). Collected the predicted labels  $y_{pred}$  and compared them to true labels  $y_{test}$  to compute performance metrics.

**Evaluation:** Calculated the overall accuracy on the test set, defined as  $(\text{True Positives} + \text{True Negatives}) / \text{Total Test Samples}$ . Also computed precision, recall, and F1-score for the positive class (heart disease = 1) and negative class (heart disease = 0) for additional insight. Constructed a confusion matrix to examine the types of errors (if any).

All computations were done in a Python environment with pandas for data handling, scikit-learn for modeling, and matplotlib for plotting the pie chart. A snippet of the Python code used in our experiments is provided below, which includes the data loading, model training, and accuracy output, as well as generation of the gender distribution chart.

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import accuracy_score
import matplotlib.pyplot as plt

# Load dataset:
df = pd.read_csv('heart.csv') # Dataset file containing UCI Heart Disease

# Preprocessing: encode 'thal' feature if it is categorical
if df['thal'].dtype == object:
    df['thal'] = LabelEncoder().fit_transform(df['thal'])
```



```
# Define features and target
X = df.drop('target', axis=1)
y = df['target']

# Split into train and test sets (80% train, 20% test)
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.2, stratify=y, random_state=42)

# Initialize and train Random Forest classifier
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# Make predictions on test set
y_pred = model.predict(X_test)
accuracy = accuracy_score(y_test, y_pred)
print(f"Test Accuracy: {accuracy:.2%}")

# Plot gender distribution pie chart
gender_counts = df['sex'].value_counts()
gender_counts.index = gender_counts.index.map({'M': 'Male', 'F': 'Female'})
plt.figure(figsize=(5,5))
plt.pie(gender_counts, labels=gender_counts.index, colors=['#4169E1', '#FF69B4'],
        autopct='%1.1f%%', startangle=140, explode=(0.05,0.05), shadow=True)
plt.title('Gender Distribution')
plt.show()
```

#### IV. RESULTS

After training the Random Forest (RF) classifier on the UCI Heart Disease dataset, the model achieved a test accuracy of approximately **85%**. This means that it correctly predicted heart disease presence or absence in about 85 out of 100 previously unseen cases. This result is consistent with prior studies using the same dataset, which typically report accuracies in the mid-80% range. It also outperforms a baseline accuracy of around 54%, which one would get by always predicting the majority class, and compares favorably with classical models like logistic regression that typically achieve 80–85% accuracy.

Further evaluation using a confusion matrix showed strong classification performance. Out of 61 test samples, the RF model identified **30 of 33 actual heart disease cases** (True Positives) and **25 of 28 healthy cases** (True Negatives), missing only 3 cases in each group. This yields a **precision and recall of approximately 91%**, and an **F1-score of 91%** for the positive class [10]. The model also demonstrated balanced performance across both classes, making it suitable for clinical screening. Even without hyper parameter tuning, the RF achieved results comparable to more complex or optimized models.

#### Model Reliability and Probability Calibration Analysis for Clinical Machine Learning Predictions

While accuracy, precision, recall, and the confusion matrix provide valuable insights into classification performance, real-world clinical deployment requires additional evaluation of model reliability. In healthcare settings, predictions are not merely binary outputs but probabilistic risk estimates that influence medical decisions. A model may achieve high



accuracy yet produce poorly calibrated probability scores, leading to overconfidence or underestimation of cardiovascular risk.

In clinical machine learning, probability calibration refers to the alignment between predicted probabilities and actual outcome frequencies. For example, among patients predicted to have a 70% risk of heart disease, approximately 70% should truly exhibit the condition for the model to be considered well-calibrated. Poor calibration can distort treatment decisions, potentially resulting in unnecessary interventions (false positives) or missed diagnoses (false negatives).

The confusion matrix, although informative, reflects only classification thresholds and does not capture whether the predicted probabilities are clinically meaningful. In real-world cardiology settings, threshold adjustments based on risk tolerance significantly affect sensitivity and specificity. A slightly miscalibrated model could disproportionately increase false positives, burdening healthcare systems, or false negatives, risking delayed treatment. Therefore, beyond confusion matrix analysis, evaluation techniques such as calibration curves or Brier scores are recommended to assess prediction reliability.

Incorporating probability calibration analysis enhances model trustworthiness and ensures safer integration into clinical decision-support systems. Future work should include calibration assessment and potential recalibration methods (e.g., Platt scaling or isotonic regression) to improve real-world applicability.

## V. DISCUSSION

The findings from our experiment reaffirm that Random Forest (RF) is a highly effective algorithm for heart disease prediction. Achieving approximately 85% accuracy on the UCI benchmark dataset highlights the model's ability to capture complex relationships between clinical risk factors and disease presence. Compared to earlier studies, our result aligns with or exceeds classical models like decision trees or Naïve Bayes (typically 75–80%) and remains competitive with logistic regression and KNN, which often perform in the low 80s. While optimized models can reach 90–98% accuracy, our untuned RF model demonstrates strong baseline reliability.

A key advantage of RF is its balance between sensitivity and specificity, making it valuable in clinical contexts. Feature importance scores identified *thal*, *cp*, *thalach*, and *oldpeak* as top predictors factors well-supported in medical literature. However, limitations include the small dataset size (303 samples) and potential generalization issues across populations.

Future improvements could involve hyper parameter tuning, feature engineering, or combining RF with other models. Despite these limitations, RF's current performance suggests it can be an effective tool for early screening and decision support in cardiovascular care.

## VI. CONCLUSION

This study explored the effectiveness of the Random Forest (RF) classifier for predicting heart disease using the UCI Heart Disease dataset. Our implementation involved minimal preprocessing and default parameters, yet the model achieved an impressive accuracy of approximately 85% on the test set. Precision and recall were both around 90%, demonstrating the model's reliability in identifying both diseased and healthy individuals. These results are consistent with recent literature, confirming RF's strong performance relative to other machine learning models.

The RF algorithm offers distinct advantages: it handles feature interactions and nonlinear patterns, provides interpretability through feature importance, and reduces overfitting via ensemble averaging. Key predictors identified such as chest pain type, thalassemia, and ST depression aligned with known medical indicators, enhancing clinical relevance.

Practically, an RF model with this level of accuracy can support early diagnosis and screening in cardiology. It could be integrated into clinical workflows or remote monitoring systems for real-time risk evaluation. Future improvements may include hyperparameter tuning, data expansion, or model optimization to push accuracy higher. Overall, this study supports Random Forest as a practical and powerful tool for intelligent cardiovascular risk prediction.



**REFERENCES**

- [1]. G. Narasimhan and A. Victor, "A hybrid approach with metaheuristic optimization and random forest in improving heart disease prediction," *Scientific Reports*, vol. 15, Art. no. 10971, 2025.
- [2]. Naik, G. G. Tejani, and S. J. Mousavirad, "SGO enhanced random forest and extreme gradient boosting framework for heart disease prediction," *Scientific Reports*, vol. 15, Art. no. 18145, 2025.
- [3]. S. Hossain, M. K. Hasan, M. O. Faruk, et al., "Machine learning approach for predicting cardiovascular disease in Bangladesh: evidence from a cross-sectional study," *BMC Cardiovascular Disorders*, vol. 24, Art. no. 214, Apr. 2024.
- [4]. J. Huang, "Heart Disease Prediction Based on the Random Forest Algorithm," in *Proc. 1st Int. Conf. on Data Analysis and Machine Learning (DAML)*, 2023, pp. 503–508.
- [5]. K. Sumwiza, C. Sheta, P. L. Ntalianis, et al., "Enhanced cardiovascular disease prediction model using random forest algorithm," *Informatics in Medicine Unlocked*, vol. 41, Art. 101316, 2023.
- [6]. M. Usama and B. Ahmad, "Comparative analysis of heart disease prediction using logistic regression, SVM, KNN, and random forest with cross-validation," *Scientific Programming*, vol. 2022, Art. ID 3805235, 2022.
- [7]. N. Chandrasekhar and S. Peddakrishna, "Enhancing heart disease prediction accuracy through machine learning techniques and optimization," *Processes*, vol. 11, no. 4, Art. 1210, 2023.
- [8]. V. S. S. Lakshmi and Y. P. V. Satyanarayana, "Cardiovascular disease prediction using random forest," *Int. J. Eng. Applied Sci. Technology*, vol. 9, no. 3, pp. 124–135, 2024.
- [9]. B. Kalaivani, "Optimizing heart disease diagnosis with a hybrid gradient descent adaptive algorithm and random forest classifier," presented at the *Int. Conf. on Recent Advances in HealthTech*, 2025.
- [10]. B. Liu, Y. Xia, Z. Yu, and B. Shi, "Heart- disease prediction based on random forest," *Multimedia Tools and Applications*, vol. 81, no. 13, pp. 18155–18179, 2022

