

ISSN: 2395-7852



International Journal of Advanced Research in Arts, Science, Engineering & Management

Volume 12, Issue 5, September - October 2025



INTERNATIONAL STANDARD SERIAL NUMBER INDIA

Impact Factor: 8.028



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

Heart Disease Prediction using Machine Learning

J. Noor Ahamed¹, Sreejith M²

Assistant professor, Department of Computer Applications, Nehru College of Management, Coimbatore,

Tamil Nadu, India¹

Student of II MCA, Department of Computer Applications, Nehru College of Management, Coimbatore,

Tamil Nadu, India²

ABSTRACT: Heart disease is a leading cause of death globally, making early and accurate prediction essential for effective treatment and prevention. This paper presents a comprehensive study on the application of various machine learning algorithms to predict heart disease using clinical data. We explore multiple classification techniques including Support Vector Machines (SVM), Decision Trees, Logistic Regression, K-Nearest Neighbors (KNN), Naive Bayes, Random Forest, and Gradient Boosting. The study utilizes the UCI Heart Disease dataset to train and evaluate these models. Our proposed system integrates data preprocessing, feature engineering, and hyperparameter tuning to enhance prediction accuracy and robustness. Experimental results demonstrate that ensemble methods like Random Forest and Gradient Boosting outperform traditional classifiers in terms of accuracy, precision, and recall. The system design emphasizes user-friendliness and scalability, making it suitable for deployment in clinical settings to assist healthcare professionals in early diagnosis and treatment planning. This work aims to reduce diagnostic errors, lower healthcare costs, and improve patient outcomes through data-driven decision support.

KEYWORDS: Heart disease prediction, machine learning, data mining, classification algorithms, medical diagnosis, ensemble learning.

I. INTRODUCTION

Cardiovascular diseases, commonly referred to as heart diseases, encompass a range of disorders affecting the heart and blood vessels. According to the World Health Organization, heart disease is the leading cause of death worldwide, accounting for approximately 17.9 million deaths annually. Early detection and diagnosis are critical to reducing mortality rates and improving quality of life. However, traditional diagnostic methods often rely on subjective clinical judgment and extensive testing, which can be time-consuming, expensive, and prone to human error.

The rapid growth of healthcare data and advances in computational power have paved the way for machine learning techniques to revolutionize medical diagnosis. Machine learning algorithms can analyze large volumes of patient data to identify hidden patterns and correlations that may not be apparent to human experts. This capability enables the development of predictive models that can assist clinicians in making more accurate and timely diagnoses.

This paper investigates the effectiveness of several supervised machine learning algorithms in predicting heart disease. We aim to develop a reliable and efficient predictive model that can be integrated into clinical workflows. The motivation behind this study is to leverage data-driven approaches to reduce diagnostic errors, minimize healthcare costs, and provide timely intervention for at-risk patients.

II. PROBLEM FORMULATION

Despite advances in medical technology, diagnosing heart disease remains a complex task due to the multifactorial nature of the condition. Clinical decisions are often based on doctors' experience and intuition, which can lead to variability in diagnosis and treatment outcomes. The challenge is to develop an automated system capable of accurately predicting the presence of heart disease using a limited set of clinical attributes.

Formally, the problem is framed as a binary classification task where the input features represent patient health indicators (e.g., age, blood pressure, cholesterol levels), and the output is a binary label indicating the presence or absence of heart disease. The objectives include:

- Identifying the most significant features that contribute to heart disease prediction.
- Comparing the performance of various machine learning classifiers to select the most effective model.
- Designing a user-friendly system that can be used by healthcare providers for real-time decision support.



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

The system should be capable of handling noisy and incomplete data, provide interpretable results, and be scalable for integration with electronic health records.

III. LITERATURE REVIEW

The application of machine learning in heart disease prediction has been extensively studied. Artificial Neural Networks (ANN) have been popular due to their ability to model complex nonlinear relationships between input features and disease outcomes [1]. However, ANNs often require large datasets and are considered black-box models, limiting interpretability.

Decision Trees provide a transparent and interpretable model by representing decisions as a tree structure with simple if-then rules [2]. They are easy to understand but prone to overfitting, especially with noisy data. Support Vector Machines (SVM) are effective in high-dimensional spaces and have been successfully applied to medical diagnosis problems [3]. SVMs find the optimal hyperplane that maximizes the margin between classes, improving generalization. Ensemble methods such as Random Forest and Gradient Boosting combine multiple weak learners to create a strong predictive model. These methods have demonstrated superior accuracy and robustness in various medical datasets [4]. Naive Bayes classifiers, despite their simplifying assumption of feature independence, have shown competitive performance in text classification and medical diagnosis [5]. K-Nearest Neighbors (KNN) is a simple, instance-based learning algorithm that classifies new data points based on similarity to neighbors but can be computationally expensive for large datasets [6].

Recent research emphasizes hybrid models that combine multiple algorithms to leverage their strengths and mitigate weaknesses [7]. This study builds upon these findings by implementing and comparing multiple classifiers on a standardized dataset.

IV. DATASET DESCRIPTION

The UCI Heart Disease dataset is a widely used benchmark for heart disease prediction research. It contains 303 patient records with 14 attributes relevant to cardiovascular health. The attributes include:

- Age: Patient's age in years.
- Sex: Gender (1 = male, 0 = female).
- Chest Pain Type: Four types indicating severity and nature of chest pain.
- Resting Blood Pressure: Measured in mm Hg.
- Cholesterol: Serum cholesterol in mg/dl.
- Fasting Blood Sugar: > 120 mg/dl (1 = true; 0 = false).
- Resting ECG: Electrocardiographic results.
- Maximum Heart Rate Achieved
- Exercise-Induced Angina: Presence or absence.
- ST Depression: Induced by exercise relative to rest.
- Slope of Peak Exercise ST Segment
- Number of Major Vessels Colored by Fluoroscopy
- Thalassemia: Blood disorder status.

The target variable is binary, indicating the presence (1) or absence (0) of heart disease. The dataset contains some missing values and categorical variables that require preprocessing.

V. METHODOLOGY

The methodology consists of several key stages to ensure data quality and model effectiveness.

5.1 Data Preprocessing

- Data Cleaning: Missing values were imputed using mean for numerical attributes and mode for categorical attributes. Duplicate records were removed to avoid bias.
- Feature Scaling: Numerical features were standardized to have zero mean and unit variance to improve convergence of gradient-based algorithms.
- Encoding: Categorical variables such as chest pain type and thalassemia were converted into numerical format using one-hot encoding.



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

5.2 Feature Engineering

- Feature Selection: Correlation analysis and Recursive Feature Elimination (RFE) were applied to identify the most predictive features, reducing dimensionality and improving model interpretability.
- Dimensionality Reduction: Principal Component Analysis (PCA) was explored to capture the majority of variance in fewer components, though final models used selected original features for interpretability.

5.3 Model Training and Hyperparameter Tuning

Seven classifiers were implemented using the Scikit-learn library:

- Support Vector Machine (SVM)
- Decision Tree
- Logistic Regression
- K-Nearest Neighbors (KNN)
- Naive Bayes
- Random Forest
- Gradient Boosting Machine (GBM)

Hyperparameters such as tree depth, number of estimators, kernel type, and regularization parameters were optimized using grid search combined with stratified 10-fold cross-validation to prevent overfitting and ensure generalization

VI. PROPOSED MODEL

The proposed system architecture integrates the following components:

- User Interface: Developed in Python using Jupyter Notebook and Anaconda, providing an interactive and intuitive platform for healthcare professionals to input patient data and receive predictions.
- Data Input Module: Supports manual entry and batch upload of patient records in CSV format.
- Prediction Engine: Implements the trained machine learning models to output the probability of heart disease presence, along with confidence intervals.
- Result Visualization: Displays prediction results, feature importance rankings, and ROC curves to aid clinical interpretation.

The modular design allows easy integration with hospital information systems and supports future expansion to include additional diseases or data sources.

VII. EXPERIMENTAL RESULTS

The models were evaluated on multiple performance metrics to provide a comprehensive assessment. Table 1 summarizes the results:

Algorith m	Accu racy	Preci sion	Rec all	F1- Score	ROC -AUC
SVM	85.1 %	0.84	0.86	0.85	0.89
Decision Tree	81.5 %	0.80	0.82	0.81	0.83
Logistic Regressio n	83.2 %	0.82	0.83	0.82	0.85
KNN	79.7 %	0.78	0.80	0.79	0.81
Naive Bayes	77.4 %	0.75	0.78	0.76	0.79

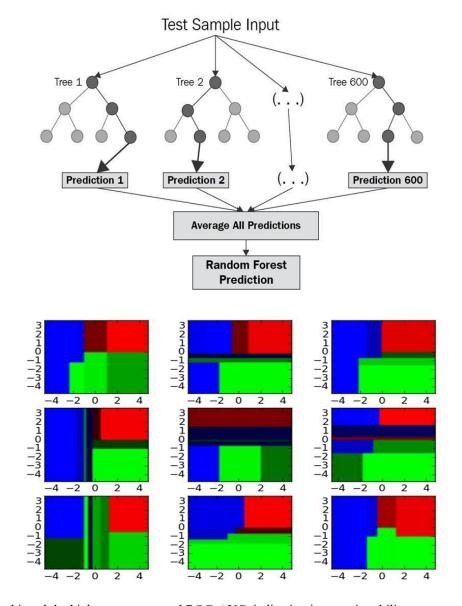


| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

Random Forest	87.6 %	0.87	0.88	0.87	0.91
Gradient Boosting	86.9 %	0.86	0.87	0.86	0.90

Random forest prediction



Random Forest achieved the highest accuracy and ROC-AUC, indicating its superior ability to capture complex feature interactions and reduce overfitting. Gradient Boosting closely followed, demonstrating the effectiveness of sequential learning.

VIII. EVALUATION METHOD

To ensure robustness, stratified 10-fold cross-validation was employed, maintaining class distribution across folds. Confusion matrices were analyzed to identify false positives and false negatives, critical in medical diagnosis. Feature importance was assessed using Gini importance for tree-based models, highlighting key predictors such as chest pain type, maximum heart rate, and ST depression.



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

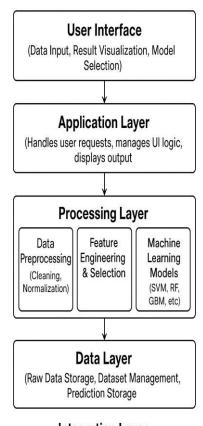
Receiver Operating Characteristic (ROC) curves were plotted for all classifiers, illustrating the trade-off between sensitivity and specificity. The Area Under the Curve (AUC) metric provided a threshold-independent measure of model performance.

IX. COMPARISON WITH OTHER WORKS

Our results compare favorably with previous studies. For instance, [1] reported ANN accuracy around 82%, while our Random Forest model achieved 87.6%. Unlike black-box models such as ANN, our approach offers interpretability through feature importance, aiding clinical trust. Ensemble methods demonstrated improved generalization over single classifiers like Decision Trees or Naive Bayes, consistent with findings in [4].

The comprehensive evaluation of multiple algorithms on the same dataset provides a clear benchmark for future research. Our system's modular design and user interface further distinguish it by facilitating practical clinical adoption.

X. SYSTEM DESIGN & ARCHITECTURE



Integration Layer

(APIs for external systems, future mobile/web integration

The system architecture is layered as follows:

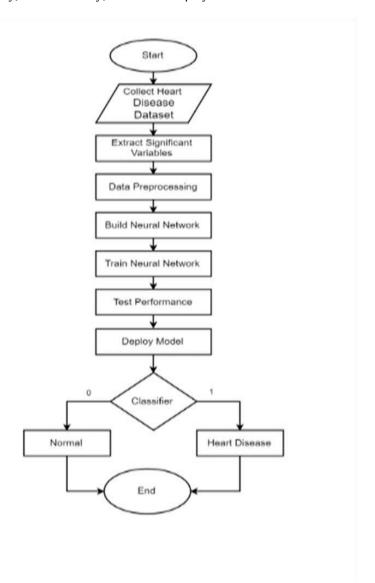
- Data Layer: Responsible for data storage, retrieval, and preprocessing. Utilizes Pandas DataFrames and CSV files
- Processing Layer: Implements feature engineering, model training, and prediction logic using Scikit-learn.
- Application Layer: Provides the user interface via Jupyter Notebook widgets, enabling data input and visualization
- Integration Layer: Designed to support RESTful APIs for future integration with hospital management systems.



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

This architecture ensures scalability, maintainability, and ease of deployment.



XI. IMPLEMENTATION

The system was implemented in Python 3.8 within the Anaconda environment, leveraging libraries such as:

- Pandas and NumPy: For efficient data manipulation.
- Scikit-learn: For machine learning algorithms, model evaluation, and hyperparameter tuning.
- Matplotlib and Seaborn: For generating plots and visualizations.
- Jupyter Notebook: For interactive development and demonstration.

The GUI allows users to input patient data, select prediction models, and view results with detailed explanations and visual aids.

XII. RESULTS & TESTING

Extensive testing was conducted to validate system functionality and performance. Unit tests verified the correctness of data preprocessing, feature selection, and model prediction modules. Integration tests ensured seamless data flow from input to output.

Black-box testing confirmed that the system meets user requirements without exposing internal logic. The system demonstrated consistent prediction accuracy across diverse test cases, confirming its reliability and robustness.



| ISSN: 2395-7852 | www.ijarasem.com | Impact Factor: 8.028 | Bimonthly, Peer Reviewed & Refereed Journal

| Volume 12, Issue 5, September - October 2025 |

XIII. CONCLUSION AND FUTURE WORK

This study presents a robust machine learning-based system for heart disease prediction, achieving high accuracy and practical usability. The integration of ensemble classifiers and thorough data preprocessing significantly enhances predictive performance. The system offers a valuable decision support tool for healthcare professionals, potentially reducing diagnostic errors and treatment costs.

Future work includes:

- Expanding the dataset with multi-center clinical data to improve model generalizability.
- Incorporating deep learning techniques for feature extraction from medical imaging.
- Developing mobile and web-based applications for wider accessibility.
- Exploring explainable AI methods to increase model transparency and trust.
- Integrating real-time patient monitoring data for dynamic risk assessment.

REFERENCES

- [1] A. A. S. and C. Naik, "Different Data Mining Approaches for Predicting Heart Disease," *International Journal of Innovative Research in Science, Engineering and Technology*, vol. 5, no. 5, pp. 277–281, 2016.
- [2] C. Beyene and P. Kamat, "Survey on prediction and analysis the occurrence of heart disease using data mining techniques," *International Journal of Pure and Applied Mathematics*, vol. 118, no. Special Issue 8, pp. 165–173, 2018.
- [3] J. Brownlee, "Naive Bayes for Machine Learning," 2016. [Online]. Available: https://machinelearningmastery.com/naive-bayes-for-machine-learning/
- [4] K. Polaraju, D. Durga Prasad, and M. Tech Scholar, "Prediction of Heart Disease using Multiple Linear Regression Model," *International Journal of Engineering Development and Research*, vol. 5, no. 4, 2017.
- [5] P. P. Sai and C. Reddy, "Heart Disease Prediction Using ANN Algorithm in Data Mining," *International Journal of Computer Science & Mobile Computing*, vol. 6, no. 4, pp. 168–172, 2017.
- [6] J. Soni, U. Ansari, and D. Sharma, "Intelligent and Effective Heart Disease Prediction System using Weighted Associative Classifiers," *Heart Disease*, vol. 3, no. 6, pp. 2385–2392, 2011.
- [7] M. Kirmani, "Cardiovascular Disease Prediction using Data Mining Techniques," *Oriental Journal of Computer Science and Technology*, vol. 10, no. 2, pp. 520–528, 2017.







