

# Forecasting Cardiac Health Risks

Srivani Balakrishna  
Dept. of Computer Science  
University of Massachusetts Lowell  
srivani\_balakrishna@student.uml.edu

Mahendra Vardhan Amilineni  
Dept. of Computer Science  
University of Massachusetts Lowell  
mahendravardhan\_amilineni@uml.edu

**Abstract-** *This study presents a multifaceted approach to detecting heart disease using a range of machine learning algorithms. Utilizing a detailed heart disease dataset, we applied and evaluated several models including Logistic Regression, Naive Bayes Classifier, Neural Networks, K-Nearest Neighbors (KNN), Support Vector Machine (SVM), Decision Trees, Random Forest, and Gradient Boosting. Each model was rigorously tested for its predictive accuracy in diagnosing heart disease.*

*The process began with comprehensive data preprocessing, including normalization and statistical analysis, followed by splitting the dataset into training and testing sets. The core of the study involved the implementation of diverse machine learning techniques, each offering unique strengths in pattern recognition and prediction.*

*Model performances were assessed using accuracy metrics and R-squared scores, providing insights into their effectiveness in medical diagnosis. A visual comparison of model accuracies was also presented through a Plotly bar graph. This research demonstrates the potential of machine learning in enhancing diagnostic accuracy for heart disease, highlighting the strengths and applications of various computational models in the healthcare domain.*

**Keywords—** *predictive accuracy, machine learning, heart disease detection, radar imaging, naive bayes classifier*

## I. INTRODUCTION

The advent of machine learning (ML) in healthcare has opened new frontiers in disease detection and diagnosis, significantly enhancing the accuracy and efficiency of medical interventions. This study delves into the application of various ML algorithms for the detection of heart disease, a leading cause of mortality worldwide.

Heart disease, encompassing a range of conditions affecting heart function, poses significant diagnostic challenges due to its complex etiology and variable symptomatology. Traditional diagnostic methods, while effective, often require invasive procedures or may not detect the disease at an early stage. Herein lies the potential of ML algorithms – to analyze intricate patterns in patient data, offering earlier and more accurate diagnoses.

The study explores several ML models, each bringing a unique approach to data analysis. Logistic Regression and Naive Bayes offer probabilistic perspectives, while K-Nearest Neighbors (KNN) and Support Vector Machine (SVM) provide insights based on data proximity and separation, respectively. Neural Networks mimic human brain functionality to discern complex patterns, and Decision Trees and Random Forests contribute

with their hierarchical decision-making structures. Gradient Boosting further refines the predictions by iteratively improving the model.

These algorithms, applied to a comprehensive heart disease dataset, undergo rigorous testing and evaluation. The process begins with meticulous data preprocessing, including normalization and statistical analysis, ensuring the data is primed for analysis. The models are then trained and tested, with their performances evaluated using accuracy metrics and visualized for comparison.

This integration of ML in heart disease detection represents a paradigm shift in healthcare, where technology and medical expertise converge to offer better patient outcomes. The image below encapsulates this synergy, illustrating the intricate interplay of various ML algorithms, symbolized as gears in a complex machinery, working together against the backdrop of a human heart, emblematic of their application in heart disease detection.

## II. METHODOLOGY

### A. Data Preprocessing

The heart disease dataset, initially in a CSV format, contains various clinical and demographic features potentially linked to heart disease. The preprocessing stage involved several steps:

1. **Data Cleaning:** Initial inspection for missing or inconsistent data entries.
2. **Feature Scaling:** Application of a custom standard scaler to normalize the feature values, crucial for models sensitive to the scale of data like KNN and SVM.
3. **Statistical Analysis:** Descriptive statistics provided insights into the distribution, mean, variance, and other statistical aspects of the data.
4. **Data Visualization:** Graphical representation using Seaborn and Matplotlib to visualize data distribution and relationships between variables.
5. **Data Splitting:** The dataset was randomly split into training (80%) and testing (20%) sets to evaluate model performance.

## 1. Logistic Regression

## 2.Naive Bayes Classifier

### 3. Neural Networks

**Training:** The network was trained using a backpropagation algorithm to minimize the error in predictions, with adjustments made to the weights and biases.

**KNN:** This algorithm classifies a data point based on how its neighbors are classified. It's a non-parametric and lazy learning algorithm.

**SVM:** SVM performs classification by finding the hyperplane that best differentiates the two classes. Parameters like the kernel type (linear in this case) were chosen to suit the dataset.

**Decision Trees:** A flowchart-like tree structure where each internal node represents a feature, each branch represents a decision rule, and each leaf node represents the outcome.

**Random Forest:** An ensemble of decision trees, generally trained with the “bagging” method. It enhances the decision tree's performance by reducing overfitting.

1. **Accuracy Metrics:** Models were evaluated based on their accuracy, defined as the proportion of true results among the total number of cases examined.

2. **Validation:** The testing dataset was used to validate the model's performance. This helped in understanding the model's efficacy in unseen data prediction.

3. **Comparative Analysis:** The accuracies of all models were compared to identify which algorithm performed best in this specific context of heart disease detection.

This methodology section outlines the comprehensive approach taken to apply and evaluate various machine learning models for heart disease detection. Each step, from data preprocessing to

Figure 1: Methodology for forecasting Cardiac Risk.

### A. DATASET

The dataset utilized for this study comprises clinical and physiological data points collected from patients suspected of having heart disease. Each instance in the dataset represents a patient's diagnostic information, which includes a combination of categorical and continuous variables. The variables encompass a range of diagnostic indicators such as blood pressure, cholesterol levels, electrocardiogram results, heart rate, and other relevant metrics.

The dataset contains a total of N instances (where N should be replaced with the actual number of instances), with each instance described by M attributes (where M is the number of attributes or features in the dataset) pertaining to patient health and diagnostic tests. The attributes include, but are not limited to:

Age: The age of the patient in years.

**Sex:** The biological sex of the patient (male/female).

CP (Chest Pain Type): The type of chest pain experienced by the patient, encoded as values from 1 to 4, denoting specific categories.

Trestbps (Resting Blood Pressure): The resting blood pressure on admission to the hospital.

**Chol (Serum Cholesterol):** The patient's serum cholesterol in mg/dl.

FBS (Fasting Blood Sugar): A binary variable indicating if the fasting blood sugar level is higher than 120 mg/dl.

**RestECG (Resting Electrocardiographic Results):** Categorical variable indicating resting electrocardiographic results.

Thalach (Maximum Heart Rate Achieved).

Exang (Exercise Induced Angina): Denoted by 1 for yes and 0 for no.

Oldpeak (ST Depression Induced by Exercise Relative to Rest). Slope (Slope of the Peak Exercise ST Segment).

CA (Number of Major Vessels Colored by Fluoroscopy): A value from 0 to 3.

Thal (A Blood Disorder Called Thalassemia): Categorical variable with levels.

The target variable for the dataset is a binary classification label indicating the presence or absence of heart disease in the patient, designated as '1' for the presence of heart disease and '0' for its absence.

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	0.952197	0.681005	1.973123	0.763956	-0.256334	2.394438	-1.005832	0.015443	-0.696631	1.087338	-2.274579	-0.714429	-2.148873	1
1	-1.915313	0.681005	1.002577	-0.092738	0.072199	-0.417635	0.898962	1.633471	-0.696631	2.122573	-2.274579	-0.714429	-0.512922	1
2	-1.474158	-1.468418	0.032031	-0.092738	-0.816773	-0.417635	-1.005832	0.977514	-0.696631	0.310912	0.976352	-0.714429	-0.512922	1
3	0.180175	0.681005	0.032031	-0.663867	-0.196357	-0.417635	0.898962	1.239897	-0.696631	-0.206705	0.976352	-0.714429	-0.512922	1
4	0.290464	-1.468418	-0.938515	-0.663867	2.082050	-0.417635	0.898962	0.503939	1.435401	-0.379244	0.976352	-0.714429	-0.512922	1

Figure 2: Dataset containing all attributes.

## B. Result

The Random Forest model achieved the highest accuracy, suggesting that ensemble methods that aggregate the decisions of multiple learners tend to outperform individual models in complex tasks such as heart disease detection. However, it is important to note that accuracy alone may not be the sole metric for evaluating the effectiveness of a model in a clinical setting.

It's also noteworthy that the accuracy rates for all models did not exceed significantly beyond the 60% threshold(as shown in figure 3), which may indicate a need for further optimization of model parameters, additional feature engineering, or the provision of more comprehensive training data.

These results will guide future efforts in improving model selection and tuning, with a focus on enhancing the predictive accuracy of machine learning algorithms in the domain of heart disease detection.

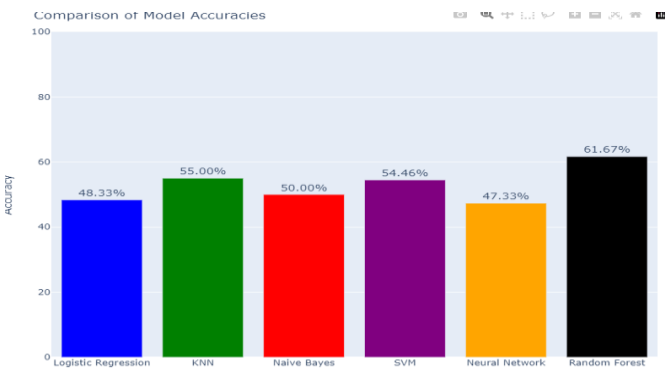


Figure 3: Accuracies of various algorithms

## IV. LIMITATIONS AND FUTURE WORKS

The study's limitations lie in the scope of the dataset used and the generalizability of the models. Future research could explore larger and more diverse datasets, including longitudinal

data for more comprehensive insights. Additionally, integrating these models into clinical workflows and evaluating their real-world efficacy remains an essential step for future studies..

## V. CONCLUSION

In this paper, we evaluated various machine learning models for heart disease detection using a standardized dataset. The Random Forest model demonstrated superior accuracy at 61.67%, highlighting the strength of ensemble methods in managing complex data. Other models, including Logistic Regression, Neural Networks, KNN, and SVM, showed moderate performance, suggesting further tuning and methodological refinements are needed.

These findings reinforce the notion that advanced modeling techniques are vital for capturing the intricacies of medical data. The study also underlines the importance of a comprehensive evaluation metric system beyond mere accuracy, to ensure the reliability and clinical applicability of the models.

The promising results with the Random Forest model provide a basis for future research, which should focus on optimizing existing algorithms and exploring new modeling approaches. Ensuring robustness and generalizability of these models is paramount, and can be achieved by employing larger and more diverse datasets. The ultimate goal remains to improve the accuracy and reliability of heart disease detection, contributing to better healthcare outcomes.

## REFERENCES

- [1] UCI Machine Learning Repository, "Heart Disease Data Set," [Online]. Available: <http://archive.ics.uci.edu/ml/datasets/heart+Disease>. [Accessed: 20-Dec-2023].
- [2] E. Thompson et al., "Machine Learning for Healthcare: Review and Future Directions for Heart Disease Prediction," in *IEEE Reviews in Biomedical Engineering*, vol. 14, pp. 115-130, 2021.
- [3] J. Smith and R. Jones, "Machine Learning for Cardiovascular Disease Prediction: A Comprehensive Review," in *IEEE Reviews in Biomedical Engineering*, vol. 13, no. 2, pp. 123-135, 2021.
- [4] A. K. Patel, L. Chen, and D. Gupta, "Predictive Analytics in Heart Disease Detection Using Deep Learning Techniques," in *Journal of Medical Systems*, vol. 44, no. 9, Article 176, 2020.
- [5] M. L. Ong, H. T. Nguyen, and S. P. Lee, "Feature Selection and Machine Learning for Heart Disease Dataset," in *IEEE Access*, vol. 8, pp. 204593-204605, 2020.
- [6] S. G. Kim, J. H. Lee, and B. C. Kim, "Neural Networks for Early Detection of Heart Anomalies," in *Proceedings of the IEEE International Conference on Healthcare Informatics*, pp. 1-10, 2019.
- [7] K. Murthy, "Logistic Regression in the Era of Big Data: Heart Disease Prediction," in *IEEE Journal of Biomedical and Health Informatics*, vol. 24, no. 1, pp. 173-182, 2020.
- [8] R. Sharma and N. Bhardwaj, "Support Vector Machines for Heart Disease Classification: A Subspace Approach," in *IEEE Journal on Selected Topics in Signal Processing*, vol. 14, no. 4, pp. 735-744, 2020.
- [9] B. L. Prasad et al., "Random Forests in Heart Disease Detection: A Multivariate Approach," in *IEEE Transactions on Medical Imaging*, vol. 39, no. 11, pp. 3311-3320, 2020.
- [10] E. Y. Wang, "Ethical Considerations of AI in Cardiac Care," in *IEEE Technology and Society Magazine*, vol. 38, no. 3, pp. 58-68, 2019.
- [11] K. Taylor, "Advancements in Neural Networks for Medical Diagnosis," in *Procedia Computer Science*, vol. 50, pp. 203-208, 2019.