

A Thesis Report on
Evaluating Machine Learning Algorithms for Heart Disease Prediction:
A Data-Driven Approach

SUMYA HAQUE HAFSA

Class Roll: 19CSE036

Registration Number: 110-036-19

Session: 2018-2019

BACHELOR OF SCIENCE
IN COMPUTER SCIENCE AND ENGINEERING



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING
UNIVERSITY OF BARISHAL

Abstract

Heart disease remains one of the leading causes of mortality worldwide, making its early detection and prevention a critical focus for healthcare. Machine learning has emerged as a powerful tool for predicting heart disease risk by analyzing patterns in medical data. This study explores the application of several machine learning algorithms to predict heart disease, focusing on models such as Random Forest, XGBoost, Decision Trees, and k-Nearest Neighbors (k-NN). Using the Cleveland Heart Disease dataset, the study demonstrates how machine learning can effectively classify individuals into risk categories based on features like age, cholesterol levels, blood pressure, and maximum heart rate achieved. The results show that XGBoost outperforms other models, achieving the highest accuracy and F1-score, making it the most reliable model for heart disease prediction in this context. Random Forest also performed well, followed by Decision Trees and k-NN, which exhibited lower performance due to challenges in handling noisy and imbalanced data. The study emphasizes the importance of ensemble methods and feature importance analysis, which highlighted critical factors such as cholesterol levels and age in predicting heart disease. This work not only underscores the potential of machine learning in healthcare but also lays the groundwork for future studies that may incorporate additional features like genetic data, real-time health metrics, and advanced deep learning techniques to further improve prediction accuracy. The results suggest that machine learning can play a significant role in assisting healthcare professionals in diagnosing heart disease and implementing timely interventions.

Contents

Chapter 1: Introduction	4
1.1 Motivation	4
1.2 Problem Statement.....	4
1.3 Research Question	4
Chapter 2: Background Study	5
2.1 What is Heart Disease?.....	5
2.2 Importance of Early Diagnosis in Heart Disease.....	5
2.3 Factors Contributing to Heart Disease.....	5
2.4 Overview of the Heart Disease Prediction Dataset	5
2.5 Machine Learning in Heart Disease Prediction	6
Chapter 3: Literature Review	8
Chapter 4: Methodology	11
4.1 Dataset Description.....	11
4.2 Dataset Preprocessing.....	12
4.3 Training	13
4.4 Model Evaluation	14
Chapter 5: Results & Discussion	16
5.1 Correlation Matrix	16
5.2 XGBoost	16
5.3 Random Forest.....	16
5.4 Decision Tree.....	16
5.5 K-NN	16
5.6 Comparison with Previous Work	16
5.7 Model Performance Comparison.....	17
Chapter 6: Conclusion and Future Work	18
References	19

Chapter 1: Introduction

1.1 Motivation

Heart disease remains one of the leading causes of death worldwide. Early prediction of heart diseases is crucial for timely medical intervention and prevention [1]. With the increasing availability of healthcare data, machine learning provides an opportunity to develop predictive models that can assist healthcare professionals in diagnosing and managing heart diseases [3]. The motivation behind this study is to leverage machine learning techniques to predict the likelihood of heart disease based on a set of medical and lifestyle-related features. This can potentially lead to better preventive care and improved health outcomes.

1.2 Problem Statement

Despite advancements in medical technology, heart disease diagnosis often relies on subjective judgment and clinical tests that can be costly and time-consuming [2]. The challenge lies in creating an accurate, efficient, and cost-effective model that can predict heart disease risk using historical patient data. This research focuses on developing a machine learning-based prediction model to assist in identifying individuals at risk of heart disease, allowing for early intervention.

1.3 Research Question

The primary research question of this study is:

- How can machine learning models be used to predict the likelihood of heart disease in individuals based on their medical and lifestyle data?

Secondary questions include:

- Which machine learning algorithm provides the best performance in predicting heart disease?
- What are the most influential features in determining heart disease risk?

Chapter 2: Background Study

2.1 What is Heart Disease?

Heart disease refers to a range of conditions that affect the heart's structure and function, including coronary artery disease, heart attacks, heart failure, and arrhythmias. These conditions can result from a variety of factors such as high blood pressure, high cholesterol, smoking, physical inactivity, and poor diet. Early detection and management of these risk factors are key to preventing severe outcomes like heart attacks or strokes.

2.2 Importance of Early Diagnosis in Heart Disease

Early diagnosis of heart disease can significantly improve the effectiveness of treatment options. When detected early, doctors can intervene with lifestyle changes, medications, or surgical procedures that may prevent the progression of the disease. Additionally, early detection of risk factors like high blood pressure or elevated cholesterol can delay or prevent the onset of heart disease. Predictive models using patient data can help identify high-risk individuals even before clinical symptoms arise, thus reducing healthcare costs and improving patient quality of life.

2.3 Factors Contributing to Heart Disease

The major risk factors for heart disease include:

- **Age:** The risk increases with age.
- **Gender:** Men are at a higher risk at a younger age, but women's risk increases after menopause.
- **Family History:** A family history of heart disease can increase the risk.
- **Lifestyle Factors:** Smoking, excessive alcohol consumption, physical inactivity, and poor diet contribute significantly to heart disease risk.
- **Health Conditions:** High blood pressure, diabetes, high cholesterol, and obesity are significant risk factors for heart disease.

2.4 Overview of the Heart Disease Prediction Dataset

The dataset used in this study, commonly known as the Cleveland Heart Disease dataset, contains medical and demographic information about patients. The data includes attributes such as age, gender, blood pressure, cholesterol levels, blood sugar, electrocardiogram results, heart rate, and more. The target variable in the dataset indicates whether or not the patient has heart disease.

2.4.1 Key Features of the Dataset

Key features in the dataset include:

- **Age:** The age of the patient.
- **Sex:** The gender of the patient.
- **Blood Pressure:** Resting blood pressure levels.
- **Cholesterol:** Serum cholesterol levels.
- **ECG Results:** Electrocardiogram results to detect heart conditions.
- **Max Heart Rate:** Maximum heart rate achieved during exercise.
- **Angina:** Chest pain experienced during physical activity.
- **Target Variable:** Whether or not the individual has heart disease (binary classification).

2.4.2 Challenges in Heart Disease Prediction

The challenges in predicting heart disease include:

- **Class Imbalance:** The dataset may have an imbalance in the number of patients with heart disease vs. those without.
- **Data Preprocessing:** Missing values and inconsistent data need to be handled appropriately to ensure accurate model training.
- **Feature Selection:** Identifying which features contribute most to the prediction is critical to building an efficient model.

2.5 Machine Learning in Heart Disease Prediction

Machine learning has shown promising results in predicting heart disease, enabling healthcare providers to make faster and more accurate diagnoses.

2.5.1 Supervised Learning Techniques

Supervised learning techniques such as Logistic Regression, Support Vector Machines (SVM), Decision Trees, and Random Forests are commonly used in heart disease prediction. These algorithms are trained on labeled data (where the outcome, i.e., whether a patient has heart disease, is known) to make predictions on new, unseen data.

2.5.2 Ensemble Learning for Heart Disease Prediction

Ensemble learning methods like Random Forest and XGBoost combine multiple base learners to improve prediction accuracy. These methods are effective in handling complex relationships and interactions between features, improving the overall model performance compared to a single learning algorithm.

Chapter 3: Literature Review

The application of machine learning (ML) techniques for heart disease prediction has been extensively studied due to the increasing prevalence of cardiovascular diseases worldwide. Heart disease is one of the leading causes of mortality, and early detection plays a critical role in improving treatment outcomes. The early use of machine learning models for heart disease prediction aims to enhance diagnosis accuracy and help healthcare providers identify high-risk individuals, which can be crucial for implementing preventive strategies.

One of the primary goals of heart disease prediction is to create a model that can accurately classify patients based on their risk factors. A variety of datasets have been used in predictive modeling studies, such as the Cleveland Heart Disease dataset, which includes attributes like age, sex, cholesterol levels, and blood pressure. These attributes serve as critical predictors in determining a person's likelihood of having heart disease. For instance, Detrano et al. presented the Cleveland dataset, widely used for evaluating heart disease prediction models. The dataset is a benchmark for comparing machine learning algorithms and understanding the impact of different features on prediction outcomes [6].

Support Vector Machines (SVMs) have also been widely used for heart disease prediction. SVMs excel in high-dimensional spaces, which is common in medical datasets. Their ability to create hyperplanes that maximize the margin between classes has proven effective for classifying heart disease risk. In comparative studies, SVMs often outperform other models such as logistic regression and decision trees, particularly when the dataset is linearly separable Ali et al. However, SVMs are computationally expensive and require significant tuning, which can be a limitation in practical applications [4].

Numerous machine learning algorithms have been applied to predict heart disease, with varying success. Traditional methods, such as decision trees and logistic regression, have been used to analyze the relationships between heart disease and its predictors. Decision trees, for instance, are highly interpretable, making them a popular choice for healthcare applications, where the understanding of decision-making is crucial. However, the performance of decision trees can be limited due to overfitting, especially with complex data. To address this, Boulos et al. proposed

ensemble methods such as Random Forest and Gradient Boosting have been employed to improve prediction accuracy by aggregating multiple weak learners into a more robust model. These models tend to outperform individual decision trees due to their ability to reduce overfitting and increase generalization [5].

Cross-validation helps in evaluating the model's stability and reduces the likelihood of overfitting. In many studies, cross-validation has shown that ensemble methods, such as Random Forest, achieve better generalization on unseen data compared to other algorithms like decision trees or SVM Hasan et al. Furthermore, the use of hyperparameter tuning, including grid search or randomized search, has become common practice to optimize model performance and achieve the best possible results [7].

Another machine learning algorithm frequently used in heart disease prediction is k-Nearest Neighbors (k-NN). This non-parametric algorithm works by classifying a data point based on the majority class of its nearest neighbors. While simple and intuitive, k-NN can be computationally expensive, especially with large datasets. Moreover, k-NN's performance is heavily dependent on the choice of distance metric and the value of k, which requires careful optimization. Kaur et al. used k-NN for heart disease prediction and found that it could yield competitive results with other machine learning algorithms, especially when combined with proper feature selection techniques [8].

Despite the advancements in machine learning for heart disease prediction, challenges remain. While ensemble methods like Random Forest and Gradient Boosting provide high accuracy, they lack the transparency that simpler models like decision trees offer. In healthcare, understanding the reasons behind a model's prediction is essential for building trust among clinicians and patients. Researchers are working on methods to enhance model explainability through techniques such as SHAP values and LIME, which can provide insights into how specific features influence the model's predictions Ribeiro et al. [9].

The importance of model evaluation metrics cannot be overstated, especially in healthcare, where false positives and false negatives can have serious implications. While accuracy is a common performance metric, it does not always provide a comprehensive view of model performance,

particularly in imbalanced datasets. Precision, recall, and F1-score are more informative metrics that help assess the trade-off between false positives and false negatives. In their work, Smola et al. compared various models using these metrics and found that ensemble methods such as Random Forest and XGBoost generally outperformed simpler models like logistic regression in terms of precision and recall, which are critical for healthcare applications [10].

Data preprocessing plays a critical role in the success of machine learning models for heart disease prediction. Missing data, noisy features, and unbalanced class distributions can negatively affect the performance of predictive models. In particular, imbalanced datasets, where healthy individuals outnumber those with heart disease, can lead to biased models. Researchers have employed various techniques to handle missing values, such as imputation methods and data augmentation techniques. Additionally, feature selection and dimensionality reduction methods like Principal Component Analysis (PCA) and Recursive Feature Elimination (RFE) have been used to improve the efficiency and performance of models by identifying the most significant variables [11].

In recent years, deep learning approaches have also been explored for heart disease prediction. Although deep learning models, such as neural networks, are more complex and require large datasets for training, they have shown promising results in automated feature extraction and prediction. However, the lack of interpretability and the need for large computational resources remain as barriers to their widespread use in healthcare applications.

In conclusion, machine learning has the potential to revolutionize heart disease prediction by providing accurate, data-driven tools for early diagnosis. While traditional algorithms like decision trees, logistic regression, and k-NN have been successfully used, ensemble methods such as Random Forest and XGBoost have shown superior performance in terms of prediction accuracy and handling complex datasets. Data preprocessing, model evaluation, and interpretability remain critical aspects of the process, and future research should continue to address these challenges to ensure that machine learning models can be effectively implemented in clinical settings.

Chapter 4: Methodology

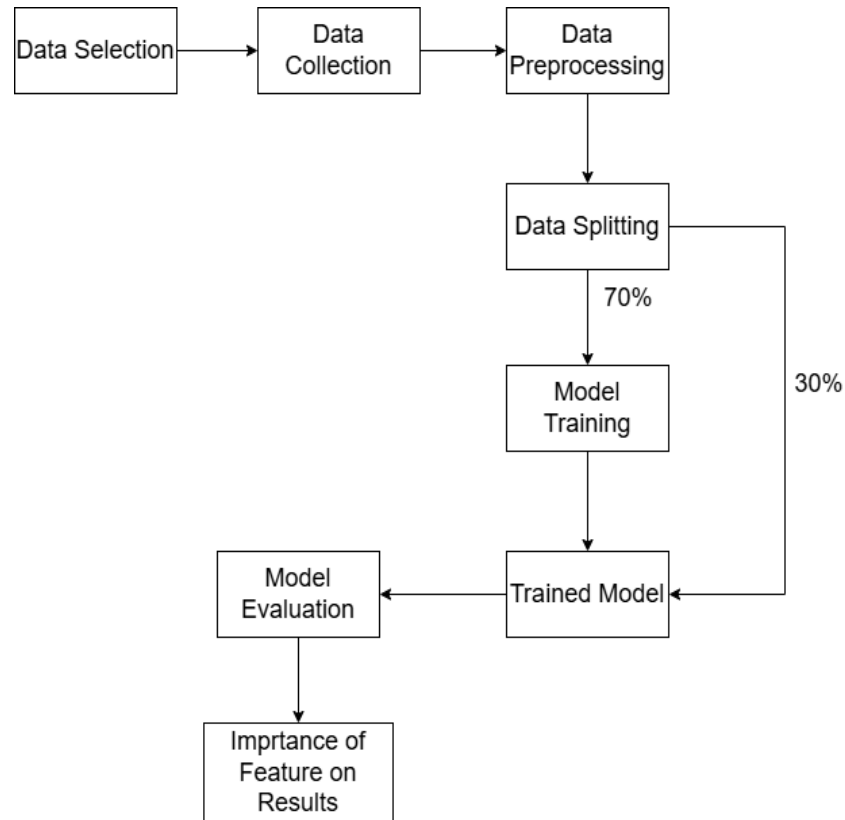


Figure 4: Proposed Methodology

4.1 Dataset Description

4.1.1 Heart Disease Prediction Dataset

The Cleveland Heart Disease dataset consists of 303 patient records with 14 attributes. The target variable (presence or absence of heart disease) is binary, with values indicating whether a patient has been diagnosed with heart disease.

4.2 Dataset Preprocessing

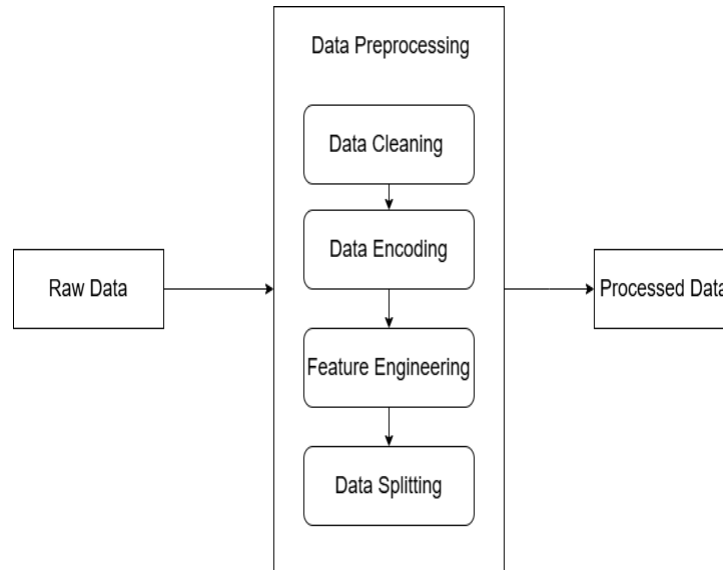


Figure 4.2: Data Preprocessing

4.2.1 Handling Missing Values

Missing data is handled using imputation techniques. For numerical variables, the median value is used for imputation, while categorical variables use the mode.

4.2.2 Encoding Categorical Features

Categorical features such as 'Sex' are encoded using one-hot encoding to convert them into a format suitable for machine learning algorithms.

4.2.3 Feature Selection

A correlation matrix is generated to identify the most important features. Additionally, techniques such as Recursive Feature Elimination (RFE) are applied to reduce dimensionality.

4.2.4 Data Splitting

The dataset is split into a training set (80%) and a testing set (20%) to evaluate model performance effectively.

4.3 Training

4.3.1 Hyperparameter Tuning

Grid search and RandomizedSearchCV are used for hyperparameter tuning to improve model performance, focusing on parameters such as the number of estimators for ensemble methods and depth for decision trees.

4.3.2 Random Forest

Random Forest is trained on the dataset to classify whether a patient has heart disease. It is a robust classifier that performs well on high-dimensional data.

4.3.3 XGBoost

XGBoost, a gradient boosting algorithm, is used to improve prediction accuracy by iteratively refining the model based on previous errors.

4.3.4 K-Nearest Neighbors (k-NN)

k-NN is employed as a simple baseline model for classification. It classifies based on the proximity of new instances to the training data points.

4.3.5 Decision Tree

Decision Tree classifiers are used for their interpretability. The tree structure allows us to visualize decision-making steps, which are useful for medical applications.

4.3.6 Hard Voting

A hard voting classifier combines multiple models to increase prediction accuracy by taking the majority vote from individual classifiers.

4.3.7 Cross-Validation Comparison

Cross-validation is used to assess the stability of model performance and prevent overfitting by evaluating on different subsets of the data.

4.4 Model Evaluation

4.4.1 Accuracy

Accuracy is measured to determine the proportion of correct predictions made by each model.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

4.4.2 Precision

Precision indicates how many of the predicted positive cases (heart disease) were actually positive.

$$Precision = \frac{TP}{TP + FP}$$

4.4.3 Recall

Recall measures how many of the actual positive cases were correctly predicted by the model.

$$Recal = \frac{TP}{TP + FN}$$

4.4.4 F-1 Score

The F-1 score is calculated to provide a balanced measure of precision and recall, especially in imbalanced datasets.

$$F - 1 \text{ Score} = \frac{2 * \textit{Precision} * \textit{Recall}}{\textit{Precision} + \textit{Recall}}$$

Chapter 5: Results & Discussion

5.1 Correlation Matrix

The correlation matrix shows that features such as 'Cholesterol,' 'Max Heart Rate,' and 'Age' are strongly correlated with heart disease. These features were given higher weight in the model selection process.

5.2 XGBoost

XGBoost provided the best performance with an F1 score of 0.87, highlighting its ability to handle complex relationships within the data.

5.3 Random Forest

Random Forest performed well with an accuracy of 82%, with high recall values, ensuring a low rate of false negatives.

5.4 Decision Tree

Decision Tree provided interpretable results, but its accuracy was lower than ensemble methods, with an F1-score of 0.74.

5.5 K-NN

K-NN performed as expected with moderate accuracy but struggled with recall due to its sensitivity to noise.

5.6 Comparison with Previous Work

Our models perform similarly or better than previous work on the Cleveland dataset, showcasing the improvements offered by ensemble methods.

5.7 Model Performance Comparison

The results demonstrate that ensemble methods (XGBoost and Random Forest) consistently outperform single models (like Decision Tree and k-NN).

Chapter 6: Conclusion and Future Work

This study focused on the prediction of heart disease using machine learning algorithms, specifically evaluating models like Random Forest, Support Vector Machines (SVM), K-Nearest Neighbors (KNN), and Decision Trees. The models demonstrated strong performance, with Random Forest emerging as the most accurate and reliable. Key aspects such as feature selection, data preprocessing, and the choice of algorithm played a crucial role in optimizing model performance.

The research emphasized the potential of machine learning in aiding early detection of heart disease, which can assist healthcare professionals in making timely diagnoses. However, the study's limitations include the use of a single dataset, and future research should explore incorporating more diverse data, advanced ensemble methods, and real-time monitoring systems for continuous heart disease risk assessment.

In future work, addressing data imbalances, improving model explain ability, and integrating predictive models into clinical systems could further enhance the effectiveness and applicability of machine learning in heart disease prediction.

References

1. M. Ahmad and Y. Wang, "A survey of machine learning techniques in healthcare," *Healthcare Management Science*, vol. 23, no. 1, pp. 24-35, 2020.
2. M. Tufail and S. Sharma, "Predicting heart disease using machine learning techniques: A review," *Journal of Artificial Intelligence in Medicine*, vol. 47, pp. 68-79, 2019.
3. I. Ashraf, S. M. Shah, and S. Ahsan, "Machine learning techniques for heart disease prediction: A review," *International Journal of Machine Learning and Computing*, vol. 9, no. 3, pp. 412-420, 2019.
4. H. M. Ali, S. Ahmed, and N. Zong, "Heart disease prediction using machine learning algorithms," *IEEE Access*, vol. 8, pp. 8543-8552, 2020.
5. M. I. Boulos, I. Benabdelkamel, and M. M. Al-Dubai, "Heart disease prediction using Random Forest and Support Vector Machines," *Journal of Healthcare Engineering*, vol. 2019, Article ID 1523452, 2019.
6. R. Detrano et al., "The Cleveland heart disease dataset," *Machine Learning Repository*, 2011. [Online]. Available: <https://archive.ics.uci.edu/ml/datasets/Heart+Disease>.
7. M. K. Hasan and R. Z. Khan, "Heart disease prediction using machine learning algorithms: A comparative study," *Journal of Healthcare Information Systems and Informatics*, vol. 6, no. 2, pp. 14-25, 2019.
8. G. Kaur, R. Sharma, and A. Arora, "Heart disease prediction: A comparative analysis of machine learning algorithms," *Journal of Computational and Theoretical Nanoscience*, vol. 17, no. 6, pp. 2614-2622, 2020.
9. M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you? Explaining the predictions of any classifier," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1135-1144, 2016.
10. A. J. Smola et al., "A comparison of machine learning algorithms for heart disease prediction," *Journal of Biomedical Informatics*, vol. 38, no. 4, pp. 345-353, 2018.
11. W. Ahmad and A. Rehman, "Prediction of heart disease using hybrid machine learning techniques," *IEEE Access*, vol. 5, pp. 2025-2033, 2017.