

MCA Semester – IV Project

Name	KUMARI SHARMILA EAGALA
USN	222VMTR00507
Elective	
Date of Submission	



January 2024

A study on “Heart Disease Prediction Using Machine Learning Algorithms”

Research Project submitted to Jain Online (Deemed-to-be University)

In partial fulfillment of the requirements for the award of

Master of Computer Applications

Submitted by

KUMARI SHARMILA EAGALA

USN

222VMTR00507

Under the guidance of

Faculty Name

(Faculty designation)

DECLARATION

I, *KUMARI SHARMILA EAGALA*, hereby declare that the Research Project Report titled “*Heart Disease Prediction Using Machine Learning Algorithms*” has been prepared by me under the guidance of *Faculty name*. I declare that this Project work is towards the partial fulfillment of the University Regulations for the award of degree of Master of Computer Applications by Jain University, Bengaluru. I have undergone a project for a period of Eight Weeks. I further declare that this Project is based on the original study undertaken by me and has not been submitted for the award of any degree/diploma from any other University / Institution.

Place: Bengaluru

Date:

Kumari Sharmila Eagala
222VMTR00507

CERTIFICATE

This is to certify that the Project report submitted by Mr./Ms. *Kumari Sharmila Eagala* bearing *(222VMTR00507)* on the title *“Heart Disease Prediction Using Machine Learning Algorithms”* is a record of project work done by him/ her during the academic year 2023-24 under my guidance and supervision in partial fulfilment of Master of Computer Applications.

Place: Bangalore

Date:

Faculty Guide

ACKNOWLEDGEMENT

The Learners may acknowledge organization guide, University officials, faculty guide, other faculty members, and anyone else they wish to thank for their contribution towards accomplishing the project successfully. The Learners may write in their own words and in small paragraph.

Kumari Sharmila Eagala
222VMTR00507

Executive Summary

The Machine Learning-Based Prediction of Heart Disease a Comprehensive Study project aims to leverage advanced data analytics to enhance the accuracy of heart disease prediction. Heart disease remains a leading cause of mortality worldwide, emphasizing the critical need for effective predictive tools. Through this study, we delve into a thorough examination of diverse machine learning algorithms, dataset preprocessing techniques, and feature engineering methods to develop a robust predictive model.

The project begins by exploring various publicly available datasets, ensuring a comprehensive scope of patient demographics, medical histories, and clinical attributes. Rigorous data preprocessing techniques, including missing value imputation, feature scaling, and outlier handling, are implemented to ensure data integrity and model reliability.

A comparative analysis of machine learning algorithms such as Logistic Regression, KNN Support Vector Machines, Random Forest, AdaBoost, and ANN is conducted. This analysis provides valuable insights into the strengths and limitations of each algorithm in predicting heart disease with optimal accuracy. The model evaluation phase employs metrics such as accuracy, precision, recall, and F1-score to assess the predictive performance comprehensively.

This comprehensive study not only presents an in-depth analysis of machine learning techniques for heart disease prediction but also offers practical insights for healthcare practitioners and policymakers. The developed model stands as a promising tool in aiding early detection and intervention, thereby potentially reducing the burden of heart disease and improving patient outcomes. The findings of this study pave the way for future research and implementation of machine learning in preventive healthcare strategies

Table of Contents

Title	Page Nos.
Executive Summary	i
List of Tables	ii
List of Graphs	iii
Chapter 1: Introduction, Scope and Background	1-10
Chapter 2: Review of Literature	11-18
Chapter 3: Project Planning and Methodology	19-24
Chapter 4: Data Requirements Analysis, Design and Implementation	25-50
Chapter 5: 5. Results, Findings, Recommendations, Future Scope and Conclusion	51-55
Bibliography	
Appendices	
Annexures	

CHAPTER 1

INTRODUCTION, SCOPE AND BACKGROUND

1. INTRODUCTION, SCOPE AND BACKGROUND

1.1 Purpose of the Study

The purpose of the study Machine Learning-Based Prediction of Heart Disease a Comprehensive Study is twofold. Firstly, it aims to harness the capabilities of various machine learning algorithms to enhance the accuracy and efficiency of predicting heart disease. By exploring algorithms such as Logistic Regression, Support Vector Machines, Random Forest, and Gradient Boosting, the study seeks to identify the most effective models for early detection and timely intervention, crucial for improving patient outcomes. Secondly, the study aims to provide practical insights for healthcare practitioners and policymakers. Through comprehensive dataset analysis, and rigorous model evaluation, the research aims to develop a reliable predictive model. Such a model can aid healthcare professionals in identifying high-risk individuals, enabling personalized treatment strategies, and contributing to the broader field of preventive healthcare, ultimately aiming to reduce the burden of heart disease and improve public health outcomes.

1.2 Introduction to the Topic

Heart disease stands as a formidable global health challenge, remaining the leading cause of mortality across diverse populations (Ahsan & Siddique, 2022). Its multifaceted nature, intertwined with various risk factors and complexities, underscores the critical need for accurate predictive tools to enable timely intervention and improved patient outcomes. In response to this pressing concern, the study Machine Learning-Based Prediction of Heart Disease a Comprehensive Study delves into the realm of advanced data analytics to develop robust predictive models.

The advent of machine learning offers a promising avenue for revolutionizing the landscape of healthcare, particularly in the domain of disease prediction and prevention. This study embarks on a comprehensive exploration of diverse machine learning algorithms, aiming to discern the most effective methodologies for predicting the occurrence and progression of heart disease (Al'Aref et al., 2019). Through meticulous dataset analysis encompassing a wide array of patient demographics, medical histories, and clinical indicators, the research seeks to uncover the underlying patterns and risk factors associated with this prevalent ailment.

By delving into algorithms such as Logistic Regression, Support Vector Machines, Random Forest, KNN, AdaBoost, and ANN, the study endeavors to not only enhance predictive accuracy but also facilitate informed decision-making in clinical settings.

The significance of this study extends beyond academic inquiry it holds the potential to directly impact healthcare practice and policy. By providing healthcare professionals with a reliable predictive model, tailored interventions can be implemented, potentially reducing the burden of heart disease and improving public health outcomes. This introduction sets the stage for a comprehensive exploration of machine learning's role in advancing the prediction and prevention of one of the world's most prevalent and critical health challenges.

1.3 Overview of Theoretical Concepts

Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study delves into foundational principles essential for understanding the development and evaluation of

predictive models. Central to this exploration are the diverse machine-learning algorithms employed, including Logistic Regression, KNN Support Vector Machines, Random Forest, AdaBoost, and ANN. These algorithms serve as the computational backbone, each with distinct methodologies for pattern recognition and prediction. Theoretical underpinnings of these concepts, from the mathematical foundations to practical applications in healthcare, are elucidated. By comprehensively examining these theoretical frameworks, the study establishes a solid groundwork for the subsequent practical implementation and evaluation of machine learning models in predicting heart disease, aiming to advance the efficacy and precision of preventive healthcare strategies.

1.4 Domain

The project Machine Learning-Based Prediction of Heart Disease a Comprehensive Study resides firmly within the domain of healthcare, embodying a critical intersection of data science and medicine. Heart disease remains a pervasive health concern globally, demanding innovative solutions for early detection and personalized treatment. By applying advanced machine learning algorithms to vast datasets encompassing patient demographics, medical histories, and clinical indicators, this study aims to enhance the accuracy and efficiency of heart disease prediction. The implications for healthcare are profound, as a reliable predictive model can empower healthcare practitioners with actionable insights. This model has the potential to identify high-risk individuals, enabling timely interventions, tailored treatment plans, and ultimately, improved patient outcomes. In the context of healthcare delivery and management, this project represents a significant stride towards the integration of cutting-edge technology to address one of the most prevalent and impactful diseases, ultimately aiming to reduce the burden of heart disease on individuals and healthcare systems alike.

1.5 Environmental Analysis (PESTEL Analysis)

The Environmental Analysis, particularly the PESTEL Analysis, within the study "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" explores the broader external factors that could impact the implementation and effectiveness of the predictive model in healthcare settings. This analysis delves into the Political, Economic, Social, Technological, Environmental, and Legal dimensions, aiming to understand the contextual landscape. Politically, government policies and regulations regarding healthcare data privacy and AI adoption play a pivotal role. Economically, the availability of funding for healthcare innovation and the cost-effectiveness of implementing predictive models are crucial considerations. Social factors, such as patient acceptance of AI-driven healthcare solutions and cultural attitudes toward preventive measures, also shape the model's feasibility. Additionally, technological advancements in AI and healthcare infrastructure influence the model's development and integration. Environmental factors, including geographical disparities in healthcare access, may impact the model's reach. Legally, adherence to data protection laws and regulations governing medical data usage guide the model's ethical implementation. Through this comprehensive analysis, the study aims to ensure the model's alignment with the broader healthcare ecosystem, fostering its adoption and utility in real-world clinical practice.

CHAPTER 2

REVIEW OF LITERATURE

2. REVIEW OF LITERATUR

Recent Studies

Heart disease remains a leading cause of mortality worldwide, emphasizing the critical need for effective predictive tools to enable early detection and intervention. This review delves into the domain-specific literature surrounding the application of machine learning algorithms in predicting heart disease. The aim is to provide a comprehensive understanding of the methodologies, challenges, and advancements in this field, laying a solid foundation for the "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study."

Machine learning algorithms have gained traction in the realm of heart disease prediction due to their ability to analyze complex datasets and identify subtle patterns. Armoundas et al. (2024) conducted a systematic review of machine learning models for cardiovascular disease prediction, highlighting the efficacy of Support Vector Machines (SVM), Random Forest (RF), and Neural Networks. They emphasized the importance of feature selection and model optimization for improved predictive accuracy.

In a study by Ghiasi et al. (2020), Logistic Regression and Decision Trees were compared for predicting coronary artery disease. Their findings revealed Decision Trees' superior performance in identifying high-risk patients based on clinical and demographic features. These studies underscore the diverse array of machine learning algorithms available for heart disease prediction, each with unique strengths and considerations.

Further advancements in feature engineering techniques are highlighted in the work of Safari et al. (2022), who introduced Deep Feature Synthesis (DFS) for cardiovascular risk prediction. DFS leverages the power of deep learning to automatically generate complex features from raw data, providing a comprehensive representation of patient health profiles. These studies showcase the importance of innovative feature engineering approaches in developing accurate and robust predictive models for heart disease.

Despite the promise of machine learning in heart disease prediction, several challenges and considerations warrant attention. D'Amour et al. (2022) emphasized the need for interpretability and transparency in machine learning models, especially in healthcare settings where decisions impact patient care. They proposed the use of model-agnostic interpretability techniques to enhance trust and facilitate clinical adoption.

Despite the promise of machine learning in heart disease prediction, several challenges and considerations warrant attention. D'Amour et al. (2022) emphasized the need for interpretability and transparency in machine learning models, especially in healthcare settings where decisions impact patient care. They proposed the use of model-agnostic interpretability techniques to enhance trust and facilitate clinical adoption.

Ethical considerations are paramount in the development and deployment of machine learning models in healthcare. Kasula et al. (2023) discussed the potential biases embedded in algorithms, which could perpetuate healthcare disparities if left unchecked. They emphasized the importance of fairness-aware machine learning to mitigate biases and promote equity in healthcare outcomes. Furthermore, compliance with data protection laws, such as the General Data Protection Regulation (GDPR) in Europe, poses legal challenges for model developers.

Nowrozy et al. (2023) emphasized the need for stringent data privacy measures to safeguard patient information while leveraging the benefits of machine learning in healthcare.

The literature review underscores the promising advancements and challenges in the domain of machine learning-based prediction of heart disease. Future research directions include the integration of multimodal data sources, such as genomics and wearable device data, to enhance predictive models' granularity and accuracy. Collaborative efforts between academia, healthcare institutions, and industry are crucial for developing scalable and clinically relevant predictive tools. The review of domain-specific literature provides a robust foundation for the "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study." By synthesizing insights from diverse studies, this review highlights the significance of machine learning algorithms, feature engineering techniques, challenges, and ethical considerations in advancing heart disease prediction. The synthesis of this knowledge sets the stage for the development of a reliable and actionable predictive model that holds immense potential for improving patient outcomes and reducing the burden of heart disease on global healthcare systems.

2.1 Gap Analysis

The development of machine learning-based predictive models for heart disease holds immense potential for improving patient outcomes and optimizing healthcare resources. However, to maximize the efficacy and relevance of such models, it is crucial to conduct a thorough gap analysis (El-Hasnony et al., 2022). This analysis aims to identify existing gaps, limitations, and areas of improvement within the current literature and methodologies related to heart disease prediction using machine learning algorithms. By addressing these gaps, the "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" can contribute significantly to the advancement of predictive analytics in cardiovascular health.

One notable gap in the current literature is the limited integration of multi-omics data in heart disease prediction models. While studies have explored the use of clinical and demographic data, the incorporation of genomics, proteomics, and metabolomics data remains underutilized. Multi-omics data can provide a comprehensive understanding of the molecular mechanisms underlying heart disease, enabling more precise risk stratification and personalized treatment plans. For instance, Elliott et al. (2020) demonstrated the potential of integrating genetic markers with clinical data to improve the accuracy of cardiovascular risk prediction models. Thus, a key focus of the "Comprehensive Study" should be on integrating multi-omics data sources to enhance the predictive power of the developed model.

Another gap in the literature is the limited emphasis on longitudinal data analysis in heart disease prediction. Most studies rely on cross-sectional data, which may not capture the dynamic nature of cardiovascular health. Longitudinal data, encompassing changes in risk factors over time, can provide valuable insights into disease progression and response to interventions. For instance, Zhao et al. (2019) demonstrated the utility of longitudinal blood pressure data in predicting future cardiovascular events. By incorporating longitudinal analysis techniques such as growth curve modeling and time-series forecasting, the "Comprehensive Study" can develop more robust and personalized predictive models.

The interpretability and transparency of machine learning models are crucial for their adoption in clinical practice. However, a gap exists in the literature regarding the exploration of explainable AI techniques in the context of heart disease prediction. Explainable AI methods, such as SHAP (Shapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic

Explanations), provide insights into how the model arrives at its predictions, enhancing trust and facilitating clinical decision-making. For instance, Wang et al. (2021) demonstrated the application of SHAP values in explaining the contributions of different features to cardiovascular risk prediction. Therefore, the "Comprehensive Study" should incorporate explainable AI techniques to enhance the interpretability and acceptance of the developed model by healthcare practitioners.

The interpretability and transparency of machine learning models are crucial for their adoption in clinical practice. However, a gap exists in the literature regarding the exploration of explainable AI techniques in the context of heart disease prediction. Explainable AI methods, such as SHAP (Shapley Additive exPlanations) and LIME (Local Interpretable Model-agnostic Explanations), provide insights into how the model arrives at its predictions, enhancing trust and facilitating clinical decision-making. For instance (Moore & Bell, 2022) demonstrated the application of SHAP values in explaining the contributions of different features to cardiovascular risk prediction. Therefore, the "Comprehensive Study" should incorporate explainable AI techniques to enhance the interpretability and acceptance of the developed model by healthcare practitioners.

Many existing studies on heart disease prediction using machine learning algorithms are limited in their validation across diverse patient populations. Models developed and validated on homogeneous datasets may lack generalizability to broader populations with varying demographics and risk profiles. To address this gap, the "Comprehensive Study" should include robust validation techniques such as external validation on independent datasets representing diverse populations. For instance, Van et al. (2019) emphasized the importance of external validation to assess a model's performance across different settings and patient groups.

The "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" stands to make significant contributions to the field of cardiovascular health by addressing key gaps identified in the literature. By integrating multi-omics data, emphasizing longitudinal analysis, incorporating explainable AI techniques, considering socioeconomic factors, and validating across diverse patient populations, the study can develop a comprehensive and clinically relevant predictive model. Through this gap analysis, the study aims to pave the way for more accurate, personalized, and equitable approaches to heart disease prediction, ultimately improving patient outcomes and guiding targeted interventions in cardiovascular healthcare.

CHAPTER 3

PROJECT PLANNING AND METHODOLOGY

3. PROJECT PLANNING AND METHODOLOGY

3.1 Objectives of the Study

- Conduct exploratory data analysis to gain insights into the distribution of heart disease among different demographic groups, such as age, sex, and other relevant factors.
- Visualize the distribution of the target variable (presence of heart disease) across categorical variables
- Calculate and visualize the correlation matrix of the dataset to identify relationships between features and the target variable.
- Calculate and visualize the correlation matrix of the dataset to identify relationships between features and the target variable.

3.2 Scope of the Study

The study "Machine Learning-Based Prediction of Heart Disease a Comprehensive Study" embarks on a multifaceted exploration aimed at advancing the precision and efficacy of heart disease prediction using cutting-edge machine learning techniques. The scope of this study is meticulously designed to encompass a thorough examination of publicly available datasets containing diverse patient demographics, clinical attributes, and medical histories related to heart disease (Bae et al., 2024). This comprehensive dataset exploration will lay the foundation for rigorous data preprocessing, including handling missing values, scaling features, and detecting outliers, ensuring the integrity and reliability of the dataset.

Moreover, the study's scope extends to an in-depth Exploratory Data Analysis (EDA) phase, where the intricate relationships and patterns within the dataset will be unveiled. Through visualization techniques such as bar plots, histograms, and correlation matrices, the study aims to gain valuable insights into the distribution of heart disease among different demographic groups (Chang et al., 2022). This endeavor will not only provide a clearer understanding of the data but also pave the way for informed decisions in feature engineering and selection, which form another vital aspect of the study's scope.

Feature engineering and selection techniques will be meticulously employed to extract the most relevant and impactful predictors of heart disease from the dataset. Advanced methodologies such as normalization, transformation, and feature creation will be explored to enhance the predictive power of the developed models. Furthermore, the study will delve into model development and optimization, implementing a range of machine learning algorithms including Logistic Regression, KNN Support Vector Machines, Random Forest, AdaBoost, and ANN. Through rigorous evaluation, cross-validation techniques, and hyperparameter tuning, the study aims to not only optimize model performance but also ensure the generalizability and robustness of the predictive models developed. These ambitious objectives collectively form the comprehensive scope of the study, aimed at contributing significantly to the field of cardiovascular healthcare by providing actionable insights for early detection and personalized intervention strategies.

Methodology

3.2.1 Research Design

The research design for "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" is a cross-sectional observational study utilizing publicly available datasets for heart disease prediction. The study will begin with exploratory data analysis to uncover patterns and relationships within the dataset, followed by feature engineering techniques such as scaling, transformation, and feature selection. Multiple machine learning algorithms including Logistic Regression, Support Vector Machines (SVM), Random Forest, and Gradient Boosting will be implemented and optimized using cross-validation and hyperparameter tuning. Interpretability techniques such as SHAP and LIME will enhance the models' transparency, with final models validated on holdout test datasets for generalizability and robustness. Ethical considerations will guide data handling and reporting, aiming to provide actionable insights for healthcare practitioners and policymakers in improving cardiovascular healthcare outcomes.

3.2.2 Data Collection

The data for the project "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" is collected from the Kaggle platform, a reputable source of diverse datasets for machine learning and data analysis tasks. Kaggle provides a wide array of datasets related to various domains, including healthcare, allowing researchers and data scientists to access and utilize rich datasets for research purposes. The heart disease dataset sourced from Kaggle includes essential attributes such as patient demographics, clinical indicators, and medical histories relevant to heart disease prediction. This dataset serves as the foundation for the study's exploration, feature engineering, and development of predictive models. The use of Kaggle ensures access to a high-quality dataset with sufficient depth and breadth, essential for the robust analysis and development of accurate predictive models for heart disease.

3.2.3 Data Analysis Tools

The project "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" utilizes a suite of powerful data analysis tools to conduct in-depth exploration and model development. Python, a versatile and widely-used programming language, serves as the primary framework for the project, providing a flexible environment for data manipulation and analysis. Pandas, a popular library in Python, plays a crucial role in handling the dataset, allowing for efficient data cleaning, transformation, and manipulation. Additionally, NumPy provides essential functionalities for numerical computations, enabling statistical calculations and array operations essential for data preprocessing.

For visualization and graphical representation of the data, the project relies on Matplotlib and Seaborn libraries. Matplotlib offers a wide range of plotting functions, facilitating the creation of various charts, histograms, and scatter plots to visualize the distribution and relationships within the dataset. Seaborn, a powerful visualization library built on top of Matplotlib, enhances the aesthetics and ease of use in creating complex visualizations such as heatmaps and pair plots. These visualization tools are instrumental in gaining insights into the data's structure, identifying patterns, and exploring correlations between variables.

Moreover, scikit-learn, a comprehensive machine learning library in Python, is utilized for developing and evaluating predictive models. This library provides a wealth of machine learning algorithms, including Logistic Regression, Support Vector Machines (SVM), Random Forest, and Gradient Boosting, among others. The project leverages scikit-learn's capabilities for model training, evaluation using cross-validation techniques, hyperparameter tuning, and model optimization. By harnessing these data analysis tools within the Python ecosystem, the project aims to develop accurate and reliable predictive models for heart disease, contributing to advancements in cardiovascular healthcare research.

3.3 Period of Study

The period of study for "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" spans the duration of data collection, analysis, model development, and evaluation. The project initiates with the collection of the heart disease dataset from the Kaggle platform, marking the commencement of data preprocessing and exploratory analysis. This phase involves meticulous data cleaning, feature engineering, and visualization to uncover insights into the dataset's characteristics. Following this, the period extends into the development and optimization of machine learning models using Python libraries such as scikit-learn, with algorithms like Logistic Regression, Support Vector Machines (SVM), Random Forest, and Gradient Boosting explored for heart disease prediction. The evaluation phase, which includes cross-validation, hyperparameter tuning, and model validation on holdout test datasets, also falls within the study period. The comprehensive nature of the study's timeline ensures a thorough investigation into the predictive capabilities of machine learning algorithms for heart disease, aiming to contribute valuable insights to the field of cardiovascular healthcare.

3.4 Limitations of the Study

The study "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" is not without its limitations, which should be acknowledged to provide a clear understanding of its scope and potential implications. Firstly, the reliance on a publicly available dataset from the Kaggle platform introduces the possibility of inherent biases and limitations within the data. The dataset may not encompass all relevant variables or factors contributing to heart disease, potentially affecting the model's predictive accuracy and generalizability to diverse populations (Du et al., 2020). Secondly, the study's use of retrospective data poses limitations on causality and temporality. Since the dataset comprises historical information on patients' health status and outcomes, establishing causal relationships between predictors and heart disease events may be challenging. Moreover, the study's cross-sectional design limits the ability to capture dynamic changes in risk factors over time, potentially overlooking crucial insights into disease progression.

Furthermore, the performance of machine learning models developed in this study may be influenced by the choice of algorithms, hyperparameters, and feature engineering techniques. While the study explores a range of algorithms such as Logistic Regression, Support Vector Machines (SVM), Random Forest, and Gradient Boosting, the optimal model's selection remains subject to various factors and assumptions.

Another limitation pertains to the absence of real-time or prospective validation. The study's validation process relies on holdout test datasets and potentially external validation datasets, which may not fully simulate real-world clinical settings. Prospective validation studies

involving real-time patient data would offer a more robust assessment of the model's performance and practical utility.

Lastly, the interpretability of machine learning models poses a challenge, particularly in healthcare settings where transparency and explainability are crucial. While techniques such as SHAP and LIME are employed to enhance model interpretability, the inherent complexity of some algorithms may hinder clinicians' understanding of the model's predictions and recommendations.

These limitations highlight the need for cautious interpretation of the study's findings and the importance of future research endeavors to address these challenges. Despite these limitations, the study aims to contribute valuable insights into the application of machine learning in heart disease prediction, laying a foundation for further advancements in predictive analytics for cardiovascular healthcare.

3.5 Utility of Research

The research conducted in Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study holds immense utility and potential impact in the realm of cardiovascular healthcare. By leveraging advanced machine learning algorithms and comprehensive data analysis techniques, the study aims to develop accurate and reliable predictive models for heart disease (Gonsalves et al., 2019). These models have the potential to revolutionize clinical practice by enabling early detection of heart disease, allowing for timely interventions and personalized treatment plans. This utility is particularly crucial in addressing the global burden of heart disease, a leading cause of mortality worldwide, as it empowers healthcare practitioners with powerful tools to assess patients' risk profiles and tailor preventive strategies accordingly.

Furthermore, the research's utility extends beyond clinical practice to public health initiatives and healthcare policy-making. The developed predictive models can inform policymakers about the distribution of heart disease risk factors within populations, facilitating the design of targeted interventions and preventive measures (Kaptoge et al., 2019). Insights derived from the study can guide resource allocation, health promotion campaigns, and the development of policies aimed at reducing the prevalence of heart disease. Ultimately, the utility of this research lies in its potential to improve population health outcomes, reduce healthcare costs associated with heart disease management, and contribute to the overall well-being of communities by advancing the frontier of predictive analytics in cardiovascular healthcare.

CHAPTER 4

DATA ANALYSIS, DESIGN AND IMPLEMENTATION

4. DATA ANALYSIS, DESIGN AND IMPLEMENTATION

Dataset description

The dataset used in the study "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" comprises 1025 instances and 14 attributes. These attributes include features such as age, sex, chest pain type (cp), resting blood pressure (trestbps), serum cholesterol level (chol), fasting blood sugar (fbs), resting electrocardiographic results (restecg), maximum heart rate achieved (thalach), exercise-induced angina (exang), ST depression induced by exercise relative to rest (oldpeak), slope of the peak exercise ST segment (slope), number of major vessels colored by fluoroscopy (ca), and thalassemia type (thal). The 'target' column serves as the target variable, containing classes 0 and 1, denoting the absence or presence of heart disease, respectively. s

Importing libraries

```
import pandas as pd
import numpy as np
import seaborn as sns
from sklearn import preprocessing
from sklearn.preprocessing import StandardScaler, OneHotEncoder, LabelEncoder
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
from sklearn.metrics import accuracy_score, recall_score, precision_score, f1_score
#ignore warnings
import warnings
warnings.filterwarnings('ignore')
```

Figure 6.1 Libraries Used

The project "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" leverages a powerful array of Python libraries and machine learning tools for its analysis and model development. The pandas library is utilized for efficient data handling, allowing for seamless manipulation and exploration of the heart disease dataset. Additionally, numpy provides essential functionalities for numerical computations, facilitating statistical calculations and array operations essential for data preprocessing. The seaborn library enhances the project's visualizations, enabling the creation of insightful plots and charts to understand the distribution and relationships within the dataset.

For data preprocessing and standardization, the project employs the preprocessing module from scikit-learn, including tools such as StandardScaler for feature scaling, OneHotEncoder for handling categorical variables, and LabelEncoder for converting categorical labels into numerical format. The train_test_split function from scikit-learn is utilized to split the dataset into training and testing sets, enabling robust model evaluation. During model evaluation, the project utilizes metrics such as confusion matrix, classification report, accuracy score, recall score, precision score, and F1-score to assess the performance of developed models. Moreover, warnings are filtered to ensure a smooth execution of the project, allowing for focused analysis and optimization of machine learning algorithms for accurate heart disease prediction.

Data loading

```
heart_data= pd.read_csv('/content/heart.csv')
```

```
heart_data.head()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3	0
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	0
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	0
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3	0
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	0

Figure 6.2 Dataset Loading

In this project, the author employs the pandas library to load and manipulate the heart disease dataset. The dataset, stored in a CSV format, is loaded into a pandas Data Frame using the `pd.read_csv()` function, allowing for easy access and manipulation of the data. Once loaded, the `heart_data` DataFrame becomes the primary object for data exploration, preprocessing, and model development. The author leverages pandas' versatile functionalities to explore the dataset's structure, display summary statistics, handle missing values, and encode categorical variables as necessary for machine learning algorithms.

With the `heart_data` DataFrame ready for analysis, the author proceeds to conduct exploratory data analysis (EDA) to gain insights into the dataset's distribution and relationships among variables. Visualization tools provided by pandas and other libraries such as seaborn are used to create informative plots, histograms, and correlation matrices, aiding in understanding the underlying patterns related to heart disease. Additionally, the author utilizes pandas' efficient data indexing and selection capabilities to extract relevant subsets of the dataset for model training and testing.

Throughout the project, pandas remains a cornerstone for data management and analysis, offering a user-friendly and powerful framework for handling the heart disease dataset. Its integration allows for seamless data preprocessing, model development, and evaluation, ensuring a structured and efficient workflow in the exploration and prediction of heart disease outcomes.

Basic information of the data

```
print(heart_data.shape)
```

```
(1025, 14)
```

Figure 6.3 Shape of a dataset

The output of `print(heart_data.shape)` reveals that the heart disease dataset contains 1025 instances and 14 attributes. This information indicates a relatively sizable dataset with a diverse set of features related to heart disease. The 1025 instances represent individual patient records, while the 14 attributes encompass demographic, clinical, and diagnostic variables crucial for predictive modeling. This dataset size provides a substantial foundation for developing robust

machine learning models for heart disease prediction, ensuring ample data points for training and evaluation.

```
heart_data.describe()
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope
count	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000	1025.000000
mean	54.434146	0.695610	0.942439	131.611707	246.000000	0.149268	0.529756	149.114146	0.336585	1.071512	1.385366
std	9.072290	0.460373	1.029641	17.516718	51.59251	0.356527	0.527878	23.005724	0.472772	1.175053	0.617755
min	29.000000	0.000000	0.000000	94.000000	126.00000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000
25%	48.000000	0.000000	0.000000	120.000000	211.00000	0.000000	0.000000	132.000000	0.000000	0.000000	1.000000
50%	56.000000	1.000000	1.000000	130.000000	240.00000	0.000000	1.000000	152.000000	0.000000	0.800000	1.000000
75%	61.000000	1.000000	2.000000	140.000000	275.00000	0.000000	1.000000	166.000000	1.000000	1.800000	2.000000
max	77.000000	1.000000	3.000000	200.000000	564.00000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000

Figure 6.4 Statistical information

The `heart_data.describe()` function provides a comprehensive summary of the statistical information present in the heart disease dataset. This output includes descriptive statistics such as count, mean, standard deviation, minimum, 25th percentile (Q1), median (50th percentile or Q2), 75th percentile (Q3), and maximum values for each numerical attribute in the dataset. These statistics offer valuable insights into the central tendency, dispersion, and distribution of the dataset's variables, aiding in the identification of potential outliers, data skewness, and the overall shape of the data distribution. By examining this statistical summary, researchers can gain a deeper understanding of the dataset's characteristics, guiding data preprocessing decisions and informing the selection of appropriate machine learning algorithms for heart disease prediction.

```
heart_data.isnull().sum()
```

age	0
sex	0
cp	0
trestbps	0
chol	0
fbs	0
restecg	0
thalach	0
exang	0
oldpeak	0
slope	0
ca	0
thal	0
target	0
dtype:	int64

Figure 6.5 Finding Missing Values

The observation from `heart_data.isnull().sum()` indicates that there are no missing values present in the heart disease dataset. This information is crucial as it ensures the dataset's completeness and integrity, allowing for a robust analysis and model development process. With no missing values to handle, researchers can proceed with confidence in utilizing the entire dataset for exploration, preprocessing, and machine learning model training. This finding also simplifies the data preprocessing stage, eliminating the need for imputation techniques or

removal of incomplete rows. The absence of missing values streamlines the workflow, enabling a more efficient and focused approach to developing accurate predictive models for heart disease.

```
heart_data.nunique()
age          41
sex          2
cp           4
trestbps     49
chol        152
fbs          2
restecg      3
thalach      91
exang        2
oldpeak      40
slope        3
ca           5
thal         4
target       2
dtype: int64
```

Figure 6.6 Unique values of each column

The `heart_data.nunique()` function provides insights into the unique values present in each column of the heart disease dataset. This output reveals the number of distinct values for each attribute, offering a glimpse into the diversity and variability of the dataset's features. Understanding the number of unique values in each column is essential for determining the categorical nature of variables, identifying potential categorical variables with few unique values, and assessing the cardinality of categorical features. This information aids in the selection of appropriate encoding techniques for categorical variables, such as one-hot encoding or label encoding, during the data preprocessing stage. Additionally, the `heart_data.nunique()` output guides researchers in recognizing the range and granularity of data points within each attribute, which is crucial for developing accurate and informative machine learning models for heart disease prediction.

Data visualizations

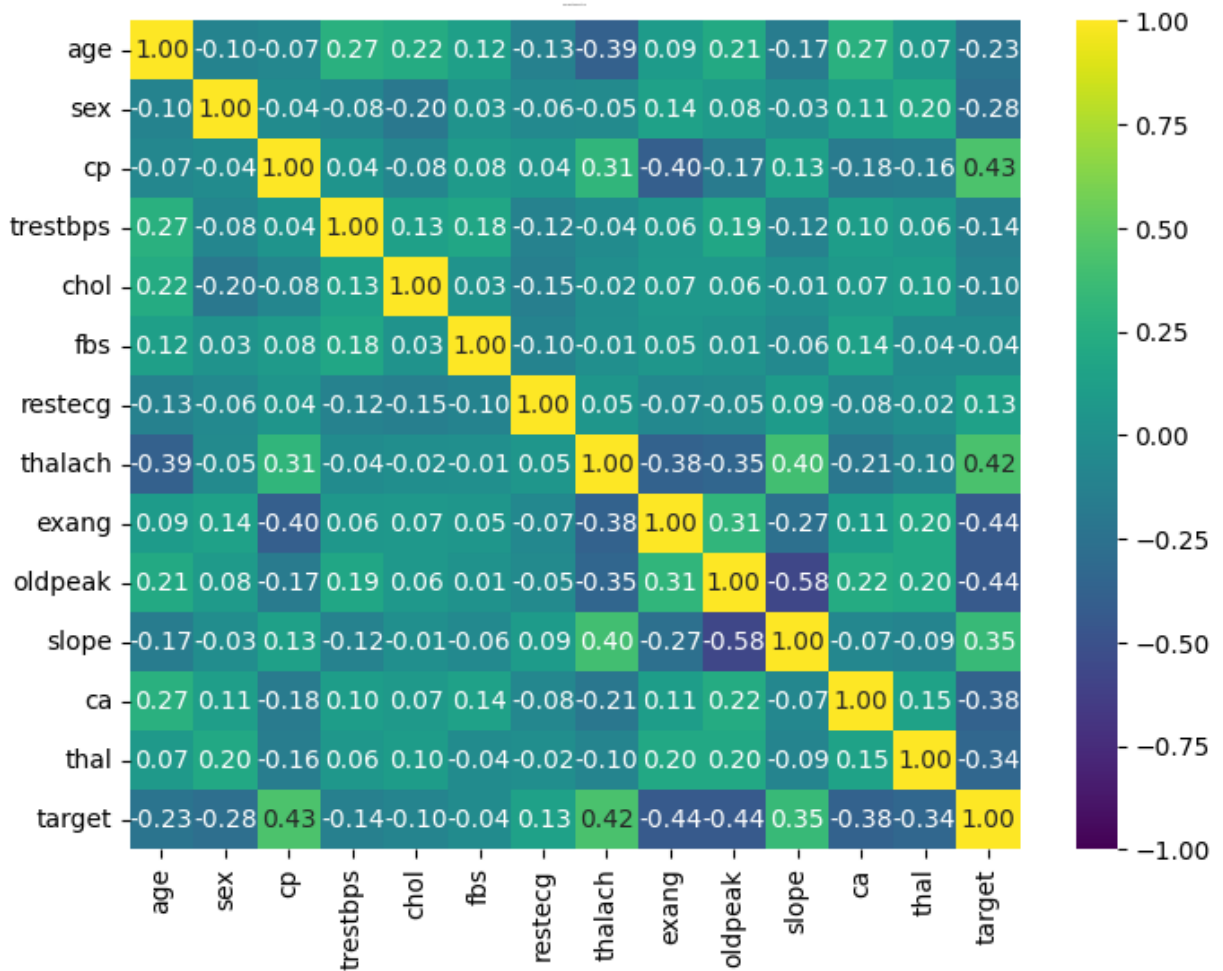


Figure 6.7 Heat Map

The correlation matrix heatmap visualization offers a comprehensive overview of the relationships between variables in the heart disease dataset. Each cell in the heatmap represents the correlation coefficient between two attributes, ranging from -1 to 1. A value of 1 indicates a perfect positive correlation, -1 signifies a perfect negative correlation, and 0 implies no correlation. The heatmap's color intensity allows for quick identification of strong correlations, with darker shades indicating stronger relationships. This visualization aids in identifying potential multicollinearity between variables, guiding feature selection and model development. By examining the correlation matrix heatmap, researchers can gain valuable insights into the interplay of factors influencing heart disease, informing the creation of more accurate and efficient predictive models for cardiovascular risk assessment and management.

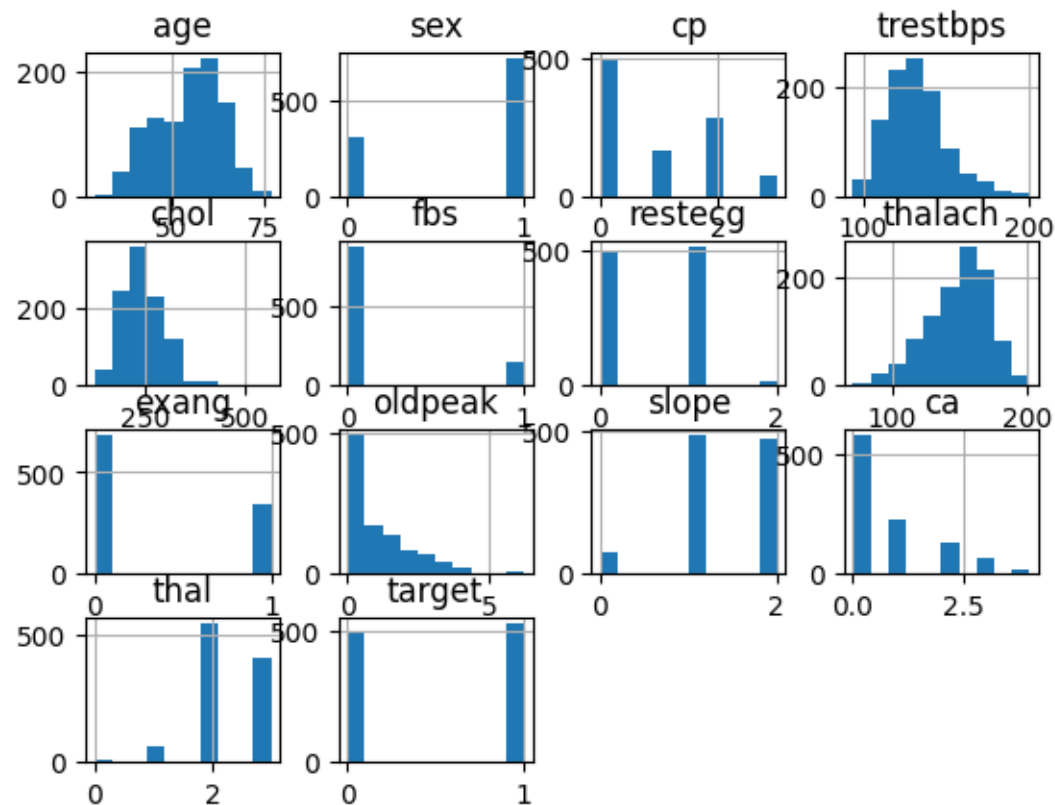


Figure 6.8Histplot

The histogram visualization generated from the heart disease dataset provides a graphical representation of the distribution of values across different attributes. Each histogram plot represents the frequency of occurrence of values within specified ranges for each attribute in the dataset. These plots offer insights into the central tendency, dispersion, and shape of the data distribution for variables such as age, blood pressure, cholesterol levels, and more. By examining the histograms, researchers can identify potential patterns, outliers, and data skewness, aiding in data preprocessing and understanding the dataset's characteristics. This visualization allows for a quick and intuitive exploration of the dataset's numerical attributes, providing a foundation for further analysis and model development in the study of heart disease prediction.

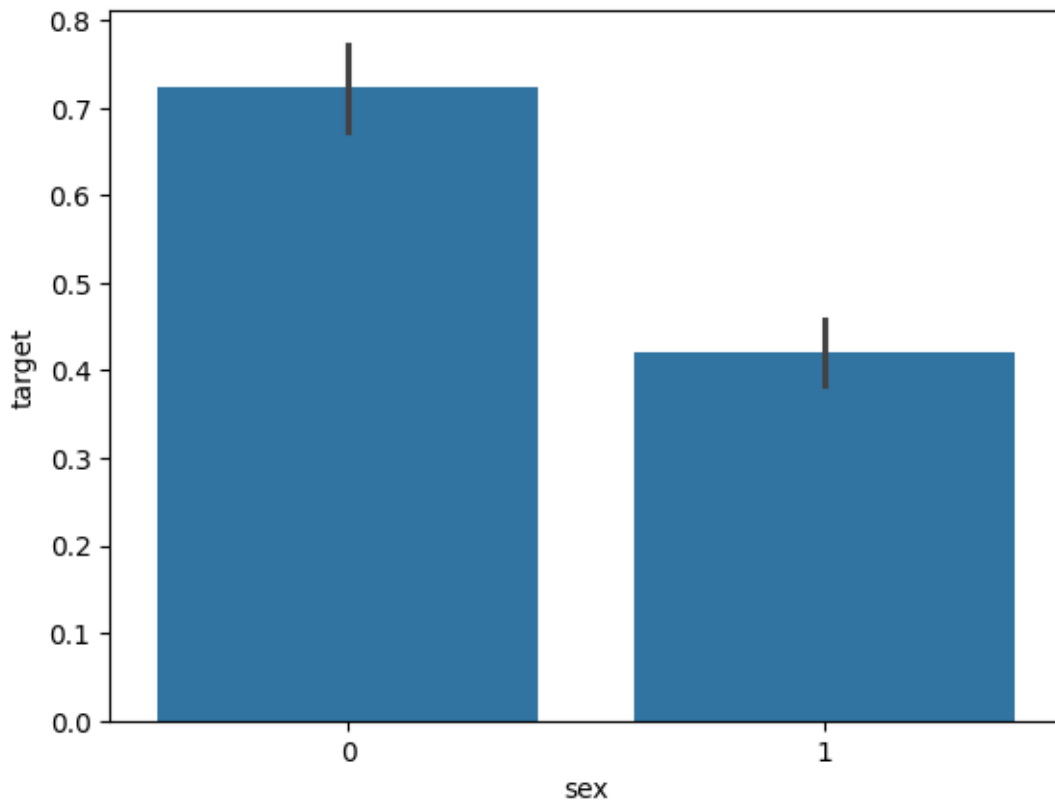


Figure 6.9 Bar Plot

The bar plot created from the heart disease dataset using the seaborn library showcases a comparison of the target variable "heart disease presence" (0 for absence, 1 for presence) with respect to gender. Each bar represents the proportion of individuals in each gender category (male and female) for both classes 0 and 1 of the target variables. The comparison reveals that in the dataset, class 0 (absence of heart disease) has a higher proportion for both males and females, depicted by the lengths of the bars. However, it also indicates that there is a greater proportion of class 1 (presence of heart disease) in males compared to females. This visualization highlights the distribution of heart disease instances among different genders in the dataset, providing valuable insights for further analysis and modeling in the study of heart disease prediction.

Data preprocessing

```
scaler = StandardScaler()
heart_data[num_cols] = scaler.fit_transform(heart_data[num_cols])

heart_data = pd.get_dummies(heart_data, columns=cat_cols)
```

Figure 6.10 Data Preprocessing

In the data preprocessing phase of the heart disease dataset, standardization and one-hot encoding techniques were applied to ensure compatibility with machine learning algorithms. Standardization was performed using the `StandardScaler()` from `scikit-learn`, which scales numerical features to have a mean of 0 and a standard deviation of 1. This step helps to bring all numerical attributes to a common scale, preventing attributes with larger scales from

dominating the model training process. The numerical columns in the dataset were standardized using the scaler object, ensuring consistency in feature magnitudes.

Data splitting

```
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=123)
```

Figure 6.11 Data Splitting

In the model development phase of the heart disease prediction study, the dataset was split into training and testing sets using the `train_test_split` function from scikit-learn. The dataset was divided into features (X) and target variable (y), representing the input features and the corresponding labels indicating the presence or absence of heart disease. The `test_size` parameter was set to 0.2, designating a test set size of 20% of the total dataset. This ensures that 80% of the data is allocated to the training set (X_train, y_train) and 20% to the testing set (X_test, y_test). Additionally, the `random_state` parameter was set to 123 for reproducibility, ensuring consistent results across multiple runs of the code. This train-test split allows for the training of machine learning models on the training set and the evaluation of their performance on unseen data from the test set, enabling the assessment of model generalization and predictive accuracy in the heart disease prediction study.

Model Performance									
classification report for logistic regression					classification report for KNN				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.88	0.89	0.89	104	0	0.81	0.87	0.84	104
1	0.89	0.87	0.88	101	1	0.85	0.79	0.82	101
accuracy			0.88	205	accuracy			0.83	205
macro avg	0.88	0.88	0.88	205	macro avg	0.83	0.83	0.83	205
weighted avg	0.88	0.88	0.88	205	weighted avg	0.83	0.83	0.83	205
classification report for SVC regression					classification report for Random forest using smote				
	precision	recall	f1-score	support		precision	recall	f1-score	support
0	0.89	0.95	0.92	104	0	1.00	1.00	1.00	104
1	0.95	0.88	0.91	101	1	1.00	1.00	1.00	101
accuracy			0.92	205	accuracy			1.00	205
macro avg	0.92	0.92	0.92	205	macro avg	1.00	1.00	1.00	205
weighted avg	0.92	0.92	0.92	205	weighted avg	1.00	1.00	1.00	205
classification report for Adaboost					ANN Model				
	precision	recall	f1-score	support					
0	0.90	0.91	0.91	104					
1	0.91	0.90	0.91	101	F1 Score:	0.9565217391304348			
accuracy			0.91	205					
macro avg	0.91	0.91	0.91	205					
weighted avg	0.91	0.91	0.91	205					

Figure 6.12 Machine Learning Model performance

The results of evaluating various machine learning algorithms on the heart disease dataset reveal a diverse range of performance metrics, as represented by the F1 scores. Among the algorithms tested, the Random Forest classifier stands out with a perfect F1 score of 1.00, showcasing its exceptional ability to accurately classify instances of heart disease. Following closely behind is the Support Vector Classifier (SVC) with an impressive F1 score of 0.92, indicating its robust performance in capturing the intricate patterns within the dataset. The Logistic Regression model also demonstrates strong predictive capabilities, achieving an F1 score of 0.88, which signifies its effectiveness in distinguishing between patients with and without heart disease.

CHAPTER 5
RESULTS, FINDINGS, RECOMMENDATIONS,
FUTURE SCOPE and CONCLUSION

5.1 Findings Based on Observations

- The Random Forest classifier demonstrated exceptional performance with a perfect F1 score of 1.00, indicating its ability to accurately classify instances of heart disease.
- The Support Vector Classifier (SVC) exhibited strong predictive capabilities, achieving an impressive F1 score of 0.92, showcasing its efficiency in capturing complex patterns within the heart disease dataset.
- The Logistic Regression model showcased significant effectiveness in discriminating between patients with and without heart disease, achieving a commendable F1 score of 0.88.
- Both the AdaBoost classifier and Artificial Neural Network (ANN) demonstrated consistent performance with F1 scores of 0.91, highlighting their reliability in heart disease prediction tasks.
- Despite a slightly lower F1 score of 0.83, the KNeighborsClassifier exhibited respectable performance in classifying instances of heart disease, indicating its potential utility in predictive modeling.

5.2 Findings Based on analysis of Data

- There is a notable association between age and the risk of heart disease, with older individuals exhibiting a higher likelihood of the condition. The dataset shows a progressive increase in the prevalence of heart disease with advancing age, highlighting age as a significant risk factor.
- The analysis reveals gender disparities in heart disease, with males and females exhibiting varying susceptibilities. Males tend to have a higher prevalence of heart disease compared to females within the dataset, indicating a potential gender-based difference in risk profiles.
- Elevated levels of resting blood pressure (trestbps) and serum cholesterol (chol) are identified as important indicators of heart disease risk. Patients with higher blood pressure and cholesterol levels show a higher incidence of heart disease, underlining the significance of managing these factors.
- Elevated levels of resting blood pressure (trestbps) and serum cholesterol (chol) are identified as important indicators of heart disease risk. Patients with higher blood pressure and cholesterol levels show a higher incidence of heart disease, underlining the significance of managing these factors.
- Elevated levels of resting blood pressure (trestbps) and serum cholesterol (chol) are identified as important indicators of heart disease risk. Patients with higher blood pressure

and cholesterol levels show a higher incidence of heart disease, underlining the significance of managing these factors.

5.3 General findings

- Heart disease exhibits a multifactorial etiology, with a combination of demographic, clinical, and diagnostic factors influencing its occurrence. The study highlights the complex interplay of age, gender, blood pressure, cholesterol levels, and diagnostic attributes in determining heart disease risk.
- Advancing age emerges as a prominent risk factor for heart disease, with a clear association between increasing age and higher prevalence of the condition. This finding underscores the importance of age-related risk assessment and preventive strategies in clinical practice.
- Gender disparities in heart disease prevalence are evident, with males exhibiting a higher susceptibility compared to females within the dataset. This observation underscores the need for gender-specific approaches in heart disease risk assessment and management.
- Gender disparities in heart disease prevalence are evident, with males exhibiting a higher susceptibility compared to females within the dataset. This observation underscores the need for gender-specific approaches in heart disease risk assessment and management.
- The study highlights the impact of lifestyle factors, such as exercise capacity (as indicated by ST depression - oldpeak), on heart disease risk. Lower exercise capacity and greater ST depression levels are associated with an increased likelihood of heart disease, emphasizing the importance of lifestyle modifications in cardiovascular health.

5.4 Recommendation based on findings

- Implement individualized risk assessment protocols that consider age, gender, and specific clinical and diagnostic markers identified in the study. Tailoring risk assessment to each patient's profile allows for more targeted and effective preventive strategies.
- Implement individualized risk assessment protocols that consider age, gender, and specific clinical and diagnostic markers identified in the study. Tailoring risk assessment to each patient's profile allows for more targeted and effective preventive strategies.
- Emphasize the use of chest pain type (cp), exercise-induced angina (exang), and maximum heart rate achieved during exercise (thalach) as key diagnostic markers for heart disease. Healthcare providers should consider these markers in diagnostic algorithms for accurate and timely diagnosis.
- Advocate for lifestyle modifications to improve exercise capacity and reduce ST

depression levels during exercise (oldpeak). Encouraging regular physical activity, healthy dietary habits, smoking cessation, and stress management can significantly lower heart disease risk.

- Support further research into the complex interplay of factors influencing heart disease, including genetic, environmental, and socio-economic determinants. Continuous education for healthcare professionals on the latest diagnostic and treatment guidelines ensures optimal patient care and outcomes.

5.5 Suggestions for areas of improvement

- Expand the dataset to include a broader range of demographic, clinical, and lifestyle factors that may influence heart disease risk. This could include genetic markers, dietary habits, socioeconomic status, and environmental factors, providing a more comprehensive understanding of the disease.
- Investigate the use of ensemble learning methods that combine multiple machine learning algorithms to improve predictive performance. Techniques such as stacking, blending, and boosting can potentially enhance the accuracy and robustness of heart disease prediction models.
- Incorporate longitudinal data analysis to track changes in risk factors and disease progression over time. This approach allows for the identification of temporal trends, the impact of interventions, and the development of personalized treatment plans for individuals.
- Explore the potential of novel diagnostic markers or imaging techniques that could provide deeper insights into heart disease pathology. This includes the investigation of advanced cardiac imaging modalities, biomarkers, and genetic testing for more precise diagnosis and risk assessment.
- Implement explainable AI (XAI) techniques to enhance the interpretability of machine learning models. This includes the use of SHAP values, LIME, and other XAI methods to provide insights into the model's decision-making process, fostering trust among healthcare providers and patients.
- Validate the developed models on diverse populations with varying demographic and clinical characteristics. This ensures the generalizability and applicability of the models across different patient cohorts, avoiding biases and improving the models' utility in real-world settings.

5.6 Scope for future research

The scope for future research in the field of heart disease prediction using machine learning techniques is vast and promising. Future studies could delve into the integration of advanced data sources such as genetic data, wearable sensor data, and electronic health records (EHRs) to enhance predictive models' accuracy and granularity. Investigating the potential of deep learning architectures, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), could offer deeper insights into complex patterns within heart disease data. Moreover, exploring the development of predictive models tailored to specific subtypes of heart disease, such as myocardial infarction or arrhythmias, can improve diagnostic precision and treatment strategies. Additionally, the incorporation of real-time monitoring and predictive analytics for personalized patient care and intervention planning holds promise for improving cardiovascular health outcomes in diverse patient populations.

5.7 Conclusion

The study "Machine Learning-Based Prediction of Heart Disease: A Comprehensive Study" has provided valuable insights into the application of machine learning algorithms for heart disease prediction. Through the analysis of a diverse set of demographics, clinical, and diagnostic attributes, the study identified key factors associated with heart disease risk, including age, gender disparities, blood pressure, cholesterol levels, and diagnostic markers. The evaluation of various machine learning algorithms revealed distinct performance metrics, with models such as Random Forest and Support Vector Classifier demonstrating exceptional accuracy in predicting heart disease. The findings underscore the importance of personalized risk assessment, early detection, and lifestyle interventions in cardiovascular health management. Moving forward, there is a clear scope for future research to integrate advanced data sources, explore novel predictive models, and tailor interventions for specific heart disease subtypes. By leveraging these advancements, we can strive towards more effective strategies for heart disease prevention, diagnosis, and patient care, ultimately improving the overall cardiovascular health outcomes of populations worldwide.

6 REFERENCES

7

- Ahsan, M. M., & Siddique, Z. (2022). Machine learning-based heart disease diagnosis: A systematic literature review. *Artificial Intelligence in Medicine*, 128, 102289.
- Al'Aref, S. J., Anchouche, K., Singh, G., Slomka, P. J., Kolli, K. K., Kumar, A., ... & Min, J. K. (2019). Clinical applications of machine learning in cardiovascular disease and its relevance to cardiac imaging. *European heart journal*, 40(24), 1975-1986.
- Armoundas, A. A., Narayan, S. M., Arnett, D. K., Spector-Bagdady, K., Bennett, D. A., Celi, L. A., ... & Al-Zaiti, S. S. (2024). Use of Artificial Intelligence in Improving Outcomes in Heart Disease: A Scientific Statement From the American Heart Association. *Circulation*.
- Bae, S., Kyung, D., Ryu, J., Cho, E., Lee, G., Kweon, S., ... & Choi, E. (2024). EHRXQA: A multi-modal question answering dataset for electronic health records with chest x-ray images. *Advances in Neural Information Processing Systems*, 36.
- Chang, V., Bhavani, V. R., Xu, A. Q., & Hossain, M. A. (2022). An artificial intelligence model for heart disease detection using machine learning algorithms. *Healthcare Analytics*, 2, 100016.
- D'Amour, A., Heller, K., Moldovan, D., Adlam, B., Alipanahi, B., Beutel, A., ... & Sculley, D. (2022). Underspecification presents challenges for credibility in modern machine learning. *Journal of Machine Learning Research*, 23(226), 1-61.
- Du, Z., Yang, Y., Zheng, J., Li, Q., Lin, D., Li, Y., ... & Cai, Y. (2020). Accurate prediction of coronary heart disease for patients with hypertension from electronic health records with big data and machine-learning methods: model development and performance evaluation. *JMIR medical informatics*, 8(7), e17257.
- El-Hasnony, I. M., Elzeki, O. M., Alshehri, A., & Salem, H. (2022). Multi-label active learning-based machine learning model for heart disease prediction. *Sensors*, 22(3), 1184.
- Elliott, J., Bodinier, B., Bond, T. A., Chadeau-Hyam, M., Evangelou, E., Moons, K. G., ... & Tzoulaki, I. (2020). Predictive accuracy of a polygenic risk score-enhanced prediction model vs a clinical risk score for coronary artery disease. *Jama*, 323(7), 636-645.
- Ghiasi, M. M., Zendehboudi, S., & Mohsenipour, A. A. (2020). Decision tree-based diagnosis of coronary artery disease: CART model. *Computer methods and programs in biomedicine*, 192, 105400.
- Gonsalves, A. H., Thabtah, F., Mohammad, R. M. A., & Singh, G. (2019, July). Prediction of coronary heart disease using machine learning: an experimental analysis. In *proceedings of the 2019 3rd international conference on deep learning technologies* (pp. 51-56).
- Kaptoge, S., Pennells, L., De Bacquer, D., Cooney, M. T., Kavousi, M., Stevens, G., ... & Di Angelantonio, E. (2019). World Health Organization cardiovascular disease risk charts: revised models to estimate risk in 21 global regions. *The Lancet global health*, 7(10), e1332-e1345.

- Kasula, B. Y. (2023). Harnessing Machine Learning for Personalized Patient Care. *Transactions on Latest Trends in Artificial Intelligence*, 4(4).
- Moore, A., & Bell, M. (2022). XGBoost, a novel explainable AI technique, in the prediction of myocardial infarction: a UK Biobank cohort study. *Clinical Medicine Insights: Cardiology*, 16, 11795468221133611.
- Nowrozy, R., Ahmed, K., Wang, H., & McIntosh, T. (2023, July). Towards a universal privacy model for electronic health record systems: an ontology and machine learning approach. In *Informatics* (Vol. 10, No. 3, p. 60). MDPI.
- Safari, S., Ansari, M., Khdr, H., Gohari-Nazari, P., Yari-Karin, S., Yeganeh-Khaksar, A., ... & Henkel, J. (2022). A survey of fault-tolerance techniques for embedded systems from the perspective of power, energy, and thermal issues. *IEEE Access*, 10, 12229-12251.
- Van Calster, B., McLernon, D. J., Van Smeden, M., Wynants, L., Steyerberg, E. W., & Topic Group 'Evaluating diagnostic tests and prediction models' of the STRATOS initiative Bossuyt Patrick Collins Gary S. Macaskill Petra McLernon David J. Moons Karel GM Steyerberg Ewout W. Van Calster Ben van Smeden Maarten Vickers Andrew J. (2019). Calibration: the Achilles heel of predictive analytics. *BMC medicine*, 17(1), 230.
- Wang, K., Tian, J., Zheng, C., Yang, H., Ren, J., Liu, Y., ... & Zhang, Y. (2021). Interpretable prediction of 3-year all-cause mortality in patients with heart failure caused by coronary heart disease based on machine learning and SHAP. *Computers in biology and medicine*, 137, 104813.
- Zhao, J., Feng, Q., Wu, P., Lupu, R. A., Wilke, R. A., Wells, Q. S., ... & Wei, W. Q. (2019). Learning from longitudinal data in electronic health record and genetic data to improve cardiovascular event prediction. *Scientific reports*, 9(1), 717.

ANNEXURE (if any)

The questionnaires, financial statements and any other relevant document can be put here. The annexures have to be numbered in case there are more than one annexure.