

A
Minor Project-I Report
On

“Predictive Modeling for Heart Failure”

**Submitted in partial fulfillment of
The requirements for the 5thSemester Sessional Examination of**

BACHELOR OF TECHNOLOGY

IN

COMPUTER SCIENCE & ENGINEERING

By

ABHISEK PANDA – 22UG010159

GOBINDA GAGAN DEY – 22UG010181

DEBABRATA MISHRA - 22UG010273

Under the able Supervision of

Dr. D. Anil Kumar

DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING



**GANDHI INSTITUTE OF ENGINEERING AND TECHNOLOGY
UNIVERSITY, ODISHA, GUNUPUR**

2024 - 25



**GANDHI INSTITUTE OF ENGINEERING AND
TECHNOLOGY
UNIVERSITY, ODISHA, GUNUPUR**

*Dist. – Rayagada-765022, Contact:- +91 7735745535,
06857-250170,172, Visit us: www.giet.edu*

Department of Computer Science & Engineering
CERTIFICATE

*This is to certify that the project work entitled "**Predictive Modeling for Heart Failure**" is done by **Abhisek Panda (22UG010159)**, **Gobinda Gagan Dey (22UG010181)**, **Debabrata Mishra (22UG010273)** in partial fulfillment of the requirements for the 5th Semester Sessional Examination of Bachelor of Technology in **Computer Science and Engineering** during the academic year 2024-25. This work is submitted to the department as a part of evaluation of 5th Semester Minor Project-I.*

Project Supervisor

Class Teacher

Project Coordinator, 3rd Year

HoD, CSE

ACKNOWLEDGEMENT

We express our sincere gratitude to **Dr. D. Anil Kumar** of Computer science and engineering for giving us an opportunity to accomplish the project. Without his active support and guidance, this project report has not been successfully completed.

We also thank to our class teacher **Mrs. Gitanjali Mishra** for her constant support during the execution of our project.

We also thank Mr.Bhavani Sankar Panda, Project Coordinator, Dr. Premanshu Sekhar Rath, Head of the Department of Computer Science and Engineering and Prof. (Dr.) Kakita Murali Gopal, Dy. Dean, Computational Science, SOET for their consistent support, guidance and help.

We also thanks to our friends, family members and others for their unconditional support during the project execution.

ABHISEK PANDA – 22UG010159

GOBINDA GAGAN DEY – 22UG010181

DEBABRATA MISHRA - 22UG010273

Abstract

Heart failure (HF) is a leading cause of morbidity and mortality worldwide, placing significant pressure on healthcare systems. It is a chronic condition where the heart cannot pump blood efficiently, leading to a lack of oxygen in vital organs and tissues. As HF progresses, it often goes undiagnosed until later stages when treatment becomes less effective. Early detection of individuals at high risk for developing heart failure is crucial for improving patient outcomes. Predictive modeling, using machine learning techniques, offers a promising approach to identifying high-risk patients based on clinical data, enabling early intervention and more personalized treatment plans.

This project aims to develop predictive models using machine learning techniques to forecast the likelihood of heart failure in patients. By analyzing comprehensive datasets that include patient demographics, clinical features, laboratory results, and medical history, the project seeks to create algorithms that can predict the onset of heart failure with high accuracy. These models will help healthcare providers make informed decisions about patient care, determining which patients need closer monitoring, which interventions are necessary, and when preventive treatments should be started.

The process begins with data gathering and preprocessing. Large datasets containing clinical variables, test results, and patient history will be collected from sources like electronic health records (EHRs). These datasets typically include information such as age, gender, medical conditions (e.g., hypertension, diabetes), lifestyle factors (e.g., smoking, physical activity), laboratory results (e.g., cholesterol levels, kidney function), and imaging data. Medical data is often noisy and incomplete, requiring preprocessing techniques to clean and prepare the data. Missing values, outliers, and inconsistencies will be handled through methods like imputation, normalization, and encoding.

Next, the project will focus on feature selection to identify the most important variables that influence the risk of heart failure. Clinical factors like age, blood pressure, ejection fraction, and comorbidities are known to contribute to heart failure risk. Statistical techniques, such as correlation analysis and recursive feature elimination, will be used to select the most relevant predictors. Feature engineering techniques will also be employed to create new variables from existing data, improving the model's accuracy.

After identifying the relevant features, various machine learning algorithms will be used to develop the predictive models. These include logistic regression, random forest, gradient boosting machines (GBM), and neural networks. Each model will be trained on the dataset and validated using cross-validation to ensure that it generalizes well to new data. The models will be evaluated using performance metrics such as accuracy, precision, recall, and AUC-ROC to assess their ability to predict heart failure risk.

Model validation is a crucial step in the project. After training, the models will be tested on unseen data to evaluate their generalization ability. The model with the highest accuracy and best performance on validation data will be selected for further use.

TABLE OF CONTENTS

CHAPTER 1	1 - 10
1.1 Introduction	2
• Purpose	4
• Scope	6
• Features	9
CHAPTER 2. SYSTEM ANALYSIS	11 -16
2.1 User Requirements (SRS)	11
2.2 Hardware Requirements	13
2.3 Software Requirements	15
CHAPTER 3	17-19
3.1 Language Used	17
3.2 Library Used	18
CHAPTER 4. SYSTEM DESIGN & SPECIFICATION	20-26
High Level Design (HLD)	
4.1 Structure Chart	20
4.2 DFD	22
4.3 UML (Use Case, Class Diagram, Activity)	24
Low Level Design (LLD)	
4.1.1 FLOW CHART	27
4.2.1 SCREEN SHOTS	28 - 32
CHAPTER 5. CODING	33 - 52
CHAPTER 6. TESTING	53 - 56
FUTURE GOAL	57 – 58
CONCLUSION	59
LIMITATION	61
REFERENCE /Bibliography	63

CHAPTER - 1

1.1 Introduction

Heart failure (HF) is a serious and growing health issue worldwide. It is a condition in which the heart is unable to pump blood effectively to meet the body's needs, leading to a lack of oxygen and nutrients in the organs and tissues. This results in symptoms like fatigue, shortness of breath, and fluid retention. Heart failure can be caused by a variety of underlying conditions such as coronary artery disease, hypertension, diabetes, and previous heart attacks. It is a progressive disease, meaning it tends to worsen over time, often with no obvious symptoms in the early stages.

Despite advances in medical care, heart failure remains a leading cause of hospitalization and death. According to recent statistics, millions of people are affected by heart failure globally, with increasing numbers of new cases diagnosed each year. It is a major contributor to healthcare costs due to frequent hospitalizations and long-term management. One of the biggest challenges in managing heart failure is its late diagnosis. Often, patients are diagnosed when the disease has already progressed significantly, making treatment more difficult and less effective.

Early identification of individuals at risk for developing heart failure is crucial to improving patient outcomes and reducing healthcare costs. Detecting heart failure in its early stages allows healthcare providers to implement preventative measures, monitor patients more closely, and start treatments that can slow the progression of the disease. However, accurately predicting who will develop heart failure before symptoms become severe is a difficult task. Traditionally, the diagnosis of heart failure relies on clinical evaluations, patient history, physical examinations, and various diagnostic tests such as blood tests, imaging, and electrocardiograms (ECGs). These methods can be effective, but they often detect heart failure too late, when the condition is already advanced.

With the rise of electronic health records (EHRs) and advances in data analytics, predictive modeling has become a valuable tool for early diagnosis and intervention. Predictive modeling uses statistical techniques and machine learning algorithms to analyze large datasets and make predictions about future events or conditions. In the context of heart failure, predictive models can analyze patient data—such as demographics, medical history, lab results, and lifestyle factors—to identify individuals at high risk of developing heart failure. These models can predict the likelihood of heart failure developing, even before clinical symptoms appear, allowing for early intervention and better management of the disease.

Machine learning, a subset of artificial intelligence, has proven particularly useful in developing predictive models. Machine learning algorithms are designed to learn from historical data, identify patterns, and make predictions based on those patterns. These algorithms can handle complex datasets with numerous variables and can improve over time as they process more data. By analyzing clinical data from a wide range of patients,

machine learning models can learn the factors that contribute most to heart failure risk, helping healthcare providers identify those at highest risk.

The use of predictive models in healthcare is becoming increasingly popular, particularly for chronic conditions like heart failure. One of the main advantages of predictive modeling is its ability to process vast amounts of data quickly and accurately. By integrating various data sources, including patient demographics, lab results, medical histories, and lifestyle factors, predictive models can provide a comprehensive assessment of an individual's risk. This information can be used by healthcare providers to make more informed decisions about patient care, such as whether to start preventive treatments, schedule more frequent follow-ups, or refer patients for additional testing.

While the potential of predictive modeling in heart failure management is clear, it also presents challenges. One challenge is the quality and completeness of the data used to train the models. Medical data is often incomplete, noisy, or inconsistent, which can affect the accuracy of the predictions. Preprocessing techniques such as data cleaning, imputation of missing values, and normalization are essential to prepare the data for analysis. Additionally, choosing the right features (variables) to include in the model is critical for its accuracy and performance. Not all factors that contribute to heart failure are easily measured or documented, and identifying the most relevant features can be a complex task.

Another challenge is integrating predictive models into clinical practice. For predictive models to be useful, they must be user-friendly and easily interpretable by healthcare providers. Decision support tools must be designed in a way that provides actionable insights without overwhelming clinicians with too much data. Additionally, models must be validated to ensure they perform well across diverse patient populations and in real-world clinical settings.

This project aims to address these challenges by developing a predictive model for heart failure using machine learning techniques. The model will analyze a comprehensive dataset of patient information to identify factors that contribute to heart failure risk. It will then generate predictions about which patients are at high risk for developing heart failure, allowing healthcare providers to intervene early and improve patient outcomes. The project will also explore how to integrate this model into clinical practice, creating a decision support system that healthcare professionals can use to make more informed decisions about patient care.

Ultimately, the goal of this project is to develop a tool that will help healthcare providers identify individuals at high risk for heart failure, enabling early diagnosis and intervention. By improving the accuracy of heart failure predictions, this tool could lead to better patient outcomes, reduced hospitalizations, and lower healthcare costs associated with the disease.

1.2.1 Purpose

The purpose of this project is to develop a predictive model using machine learning techniques that can accurately estimate the risk of heart failure in patients. Heart failure is a chronic condition that significantly impacts the health and well-being of millions of people around the world. It is a progressive disease, meaning that it often worsens over time, and when diagnosed at advanced stages, treatment options become limited and less effective. As a result, heart failure places a substantial burden on both patients and healthcare systems globally. The goal of this project is to leverage the power of predictive modeling to address some of the challenges associated with the early diagnosis and management of heart failure.

One of the key challenges in managing heart failure is the difficulty in identifying individuals who are at high risk of developing the condition before symptoms become severe. Traditionally, heart failure is diagnosed when symptoms such as shortness of breath, fatigue, and fluid retention become evident. However, by this time, the disease has often progressed significantly, and interventions may not be as effective. Identifying patients at risk for heart failure early, before clinical symptoms appear, allows healthcare providers to implement preventative measures, start early treatments, and reduce the likelihood of the disease progressing to more severe stages.

The purpose of this project is to create a tool that can analyze patient data to predict the likelihood of heart failure development. The project will develop a machine learning model that uses a variety of clinical data, including demographic information, medical history, laboratory results, lifestyle factors, and other relevant variables, to assess the risk of heart failure in patients. By doing so, the model will enable healthcare providers to make more informed decisions regarding patient care. This may include initiating preventive therapies, monitoring high-risk patients more closely, or referring patients for further testing and evaluation.

Machine learning has become an essential tool in healthcare for predictive analysis, and its application in heart failure prediction holds great promise. One of the main advantages of machine learning is its ability to process large volumes of data and identify patterns that are often too complex for traditional statistical methods. By analyzing a comprehensive dataset, the model can identify subtle relationships between various factors and predict which individuals are most likely to develop heart failure in the future. The predictive model will serve as a valuable decision support tool for clinicians, providing them with real-time risk assessments that can guide their decisions about patient care.

Another important purpose of this project is to improve the overall efficiency of healthcare delivery. Heart failure is a leading cause of hospitalization and is associated with high healthcare costs, particularly when the disease is diagnosed at advanced stages. By predicting heart failure risk early, healthcare providers can allocate resources more effectively. For example, early identification of high-risk patients can lead to more focused monitoring and early interventions, potentially reducing the need for costly emergency room visits or hospitalizations. Preventive measures, such as lifestyle modifications, medications, and closer monitoring, can

be implemented before the condition becomes severe, which can improve patient outcomes and reduce the financial burden on healthcare systems.

Furthermore, this project seeks to contribute to the personalization of heart failure care. Currently, heart failure treatment plans are often generalized, with similar interventions being applied to a broad group of patients. However, heart failure is a highly individualized condition, with risk factors and disease progression varying widely among patients. By using predictive models that account for a variety of personal and clinical factors, healthcare providers can create more personalized treatment plans tailored to each patient's unique risk profile. This approach can improve patient engagement, treatment adherence, and overall outcomes, as patients receive care that is better suited to their specific needs.

Another key purpose of this project is to provide a platform for integrating predictive models into routine clinical practice. For predictive models to be effective in real-world healthcare settings, they must be easy to use and accessible to clinicians. This project aims to develop a decision support system that can be easily integrated into existing healthcare infrastructure. The system will allow healthcare providers to input patient data and receive real-time risk predictions, which can guide their decision-making process. This integration is important because it ensures that the predictive model is not just an academic tool, but one that can have a tangible impact on patient care in clinical environments. The goal is to make the technology as user-friendly and practical as possible, providing healthcare professionals with actionable insights that can improve the quality of care delivered to patients.

The project will also focus on improving the accuracy and robustness of predictive models. One of the challenges with machine learning models is ensuring that they are generalizable across different populations. Healthcare data can vary widely depending on the population, geographic location, and healthcare practices. Therefore, it is important to ensure that the predictive model performs well across diverse patient groups and is not biased toward any particular demographic. This will be achieved through rigorous testing and validation of the model using a range of patient data from different sources. By improving the generalizability of the model, the project aims to ensure that it can be used effectively across various healthcare settings and patient populations.

Ultimately, the purpose of this project is to make a significant impact on heart failure prevention and management by providing healthcare providers with a powerful tool for early detection and personalized care. Early prediction of heart failure risk can save lives, improve patient outcomes, and reduce the financial burden on healthcare systems. Through machine learning and predictive modeling, this project aims to improve the way heart failure is diagnosed, managed, and treated, contributing to a more efficient and effective healthcare system.

In conclusion, this project seeks to address the challenges associated with the early detection and management of heart failure by developing a predictive model that can identify high-risk patients. By doing so, it aims to enable earlier interventions, personalized treatment plans, and improved patient outcomes.

1.2.2 SCOPE

The scope of this project is to develop a predictive model that uses machine learning techniques to forecast the risk of heart failure in patients based on clinical and demographic data. The primary aim is to create a tool that can identify individuals at high risk of developing heart failure, enabling healthcare providers to intervene early and implement personalized treatment strategies. By leveraging patient data, this project seeks to offer an evidence-based solution to improve early diagnosis, prevent the progression of heart failure, and reduce the strain on healthcare systems.

Data Collection and Preprocessing

The project will involve the collection of patient data from various sources, such as electronic health records (EHRs), medical databases, and publicly available datasets. The data will include key variables like age, gender, medical history (such as hypertension, diabetes, coronary artery disease), laboratory test results (e.g., cholesterol levels, kidney function), lifestyle factors (e.g., smoking, physical activity), and other clinical features relevant to heart failure risk. The dataset will be extensive, containing diverse patient information that will help in creating a robust model.

Preprocessing the data is a critical step in ensuring that the dataset is clean, consistent, and ready for analysis. This phase will include handling missing values, removing outliers, and normalizing the data. Techniques like imputation (replacing missing data with estimates) and encoding categorical variables will be applied to prepare the dataset for machine learning models. Ensuring the quality and completeness of the data is vital because the accuracy of predictive models relies heavily on the quality of the input data.

Feature Selection and Engineering

The project will include a thorough analysis to identify the most important features that are most predictive of heart failure risk. Some of the key features likely to be included are demographics (age, gender), medical conditions (hypertension, diabetes, coronary artery disease), vital signs (blood pressure, heart rate), and laboratory test results (e.g., kidney function, cholesterol levels). Feature selection will involve applying statistical methods and techniques, such as correlation analysis and recursive feature elimination, to identify the most relevant predictors.

Additionally, feature engineering will be applied to create new features from existing data that may offer additional insights into heart failure risk. For instance, deriving new variables like body mass index (BMI) from height and weight, or aggregating certain clinical data to form risk scores, can improve the predictive power of the model. The goal of feature selection and engineering is to ensure that the most relevant and informative data points are included in the model to increase its accuracy and effectiveness.

Model Development and Training

Once the data is preprocessed and the most relevant features are selected, the next step will be to develop predictive models using machine learning techniques. The project will explore a variety of algorithms, including:

Logistic Regression: A simple but effective model for binary classification tasks, which will predict whether a patient is at high risk of developing heart failure.

Random Forest: An ensemble learning method that combines multiple decision trees to make predictions. This method is useful for handling large datasets and reducing overfitting.

Gradient Boosting Machines (GBM): A powerful technique that builds decision trees sequentially to correct errors made by previous trees, improving prediction accuracy.

Neural Networks: A deep learning technique that can model complex, non-linear relationships in the data. This method is useful when working with large and high-dimensional datasets.

Each of these models will be trained using historical patient data and validated through techniques like cross-validation. Cross-validation will help evaluate the model's performance on unseen data and ensure that it generalizes well across different populations. The models will be trained to predict the likelihood of heart failure, and their performance will be assessed based on metrics such as accuracy, precision, recall, and area under the receiver operating characteristic curve (AUC-ROC).

Model Evaluation and Validation

Once the models are trained, they will be rigorously evaluated to assess their predictive performance. A set of evaluation metrics will be used to determine how well each model predicts heart failure risk. These metrics will include:

Accuracy: The overall percentage of correct predictions made by the model.

Precision: The proportion of positive predictions that are actually correct.

Recall: The proportion of actual positive cases that are correctly identified by the model.

AUC-ROC: A measure of the model's ability to distinguish between high-risk and low-risk patients.

The models will be tested on unseen validation data to evaluate their generalization ability. The goal is to select the model that provides the highest accuracy while also minimizing false positives and false negatives. In addition, the model will be tested on diverse patient populations to ensure that it works effectively across different demographic groups and clinical settings.

Clinical Integration

A crucial aspect of this project is the integration of the predictive model into clinical practice. After identifying the best-performing model, the project will focus on developing a decision support system that allows

healthcare providers to use the model in real-time. This system will be designed to be user-friendly, providing clinicians with actionable insights based on patient data.

Using a platform like Streamlit, a web-based application will be developed where healthcare providers can input patient data and receive real-time predictions about the likelihood of heart failure. The system will provide suggestions for further action, such as whether the patient requires additional testing, should begin preventive treatment, or needs closer monitoring. By integrating the predictive model into clinical workflows, this tool can assist healthcare providers in making data-driven decisions that improve patient care and outcomes.

Expected Outcomes and Impact

The primary outcome of this project is the development of an accurate, reliable predictive model that can identify patients at high risk of developing heart failure. This will enable healthcare providers to take proactive steps to manage these patients and prevent or delay the onset of heart failure. Early detection will allow for timely interventions, including lifestyle changes, medications, and closer monitoring, which can significantly improve patient outcomes.

The project is also expected to improve the efficiency of healthcare systems. By enabling early identification of high-risk patients, healthcare resources can be allocated more effectively, reducing unnecessary hospitalizations and emergency care. Personalized treatment plans tailored to individual patient profiles will be more effective, leading to better patient engagement, adherence to treatment, and overall health outcomes.

Furthermore, the development of a decision support system will enhance clinical decision-making by providing healthcare providers with easy-to-understand, real-time risk predictions. This will help to reduce uncertainty in clinical decision-making and ensure that patients receive appropriate care at the right time.

Limitations and Future Directions

While the project aims to develop a robust predictive model, there are several limitations to consider. The accuracy of the model depends heavily on the quality and completeness of the data used for training. Incomplete or inaccurate data may impact the model's performance. Additionally, the model's effectiveness will need to be validated in different clinical settings to ensure that it generalizes well across various patient populations.

Future work could involve improving the model by incorporating additional data sources, such as genetic information, imaging data, or patient-reported outcomes. Further development could also include enhancing the decision support system to provide more specific recommendations and integrate with existing healthcare systems for seamless clinical use.

1.2.3 FEATURES

The predictive model for heart failure will have several key features designed to enhance its functionality, user experience, and clinical impact. These features are focused on ease of use, integration into existing systems, and facilitating communication with users.

Sending Predicted Report to Email

One of the main features of this system is the ability to send predicted reports directly to the user's email address. After a healthcare provider or patient inputs their data into the system, the predictive model will generate a report that includes the heart failure risk prediction and recommendations for further action. This report will be automatically sent to the email address provided by the user. This feature ensures that important information is easily accessible to the user, whether it is a healthcare provider or a patient, and helps maintain timely communication regarding risk assessment.

Collecting User Data

To send reports and facilitate communication, the system will collect basic user information such as name, email address, and other relevant contact details. This user data is crucial for the proper functioning of the predictive tool, ensuring that reports and updates are sent accurately and promptly. It will also help in building a user database that can be used for follow-ups or further data collection if needed.

Bug Report Collection

The system will also include a feature to collect bug reports. Users who encounter any issues or technical difficulties while using the platform can submit detailed bug reports through a simple interface. This feedback will help the development team address technical problems, improve the system's functionality, and ensure that the platform operates smoothly.

Feedback Submission

In addition to bug reports, users will be encouraged to submit feedback on their experience using the system. This feature allows users to provide suggestions for improvement or express any concerns about the system's performance. Regular feedback collection is essential for continuous improvement, and the development team can use it to refine the user interface, enhance the predictive model, and ensure the system meets the needs of healthcare providers and patients.

User-Friendly Interface

The platform will feature an intuitive and easy-to-use interface. Both healthcare providers and patients should be able to navigate the system effortlessly. This will include easy data input options, clear navigation, and accessibility features that make the system usable across different devices and by people with varying levels of technical expertise.

Real-Time Predictions and Alerts

In addition to sending reports, the system will provide real-time predictions based on user input. If a patient is identified as being at high risk of heart failure, the system will immediately alert the healthcare provider or the patient, prompting them to take further action or schedule a follow-up appointment. These alerts will be generated instantly, enhancing the ability of healthcare providers to respond promptly to high-risk cases.

Integration with Healthcare Systems

The system will be designed to integrate seamlessly with existing healthcare management systems. This will ensure that data entered into the platform can be accessed and utilized by healthcare providers within their existing workflow. Integration will also allow for better data sharing and collaboration between medical professionals, helping to ensure that heart failure risks are continuously monitored and managed over time.

Personalized Recommendations

Once a risk prediction is made, the system will provide personalized recommendations based on the specific data entered. These recommendations could include lifestyle changes, preventive treatments, or suggestions for further diagnostic tests. Personalized recommendations are an essential feature for improving patient outcomes and ensuring that the interventions are tailored to the individual's needs.

Security and Privacy

Given the sensitive nature of the data being handled, the system will have robust security measures in place to protect user data. This includes encryption, secure login methods, and strict access controls. Only authorized users will be able to view or edit patient data, ensuring privacy and compliance with healthcare regulations like HIPAA.

Mobile and Web Compatibility

The system will be designed to be compatible with both web and mobile platforms. This allows healthcare providers and patients to access the system and receive risk reports or alerts from anywhere. Whether a healthcare professional is using a laptop in their office or a patient is using their smartphone, they will be able to input data, receive predictions, and get follow-up reports seamlessly.

User Registration and Authentication

A feature for user registration and authentication will be incorporated. This ensures that only authorized users, such as healthcare professionals or registered patients, can access the system. The registration process will collect basic information from users, such as their name, email, and medical credentials (for healthcare providers), to ensure that the system is used correctly and securely.

CHAPTER – 2

SYSTEM ANALYSIS

2.1 User Requirements (SRS)

The Heart Failure Predictive Model aims to leverage machine learning techniques to predict the likelihood of heart failure in patients based on various health-related features. The model will be developed using Python and deployed through a web application for user interaction. This document outlines the user requirements for the heart failure predictive model system, providing clear specifications for the functionalities, interfaces, and expected behavior of the system.

Overall Description

The heart failure predictive model will analyze patient data, process the input features, and output a prediction indicating the risk of heart failure. The system will be easy to use, providing an intuitive interface for healthcare professionals or other users to input relevant data and receive predictions.

System Features and Requirements

3.1 User Interface

- UI Accessibility:

- The web-based interface will be accessible through modern web browsers (Chrome, Firefox, Edge).
- The system should provide a clean, intuitive, and easy-to-navigate user interface (UI).

- Input Features:

- Users will input patient data such as age, blood pressure, cholesterol levels, ECG readings, etc.
- The input fields will be clearly labeled, and valid input formats will be enforced.

- Predictions Display:

- After entering data, users will click a "Predict" button to generate the prediction.
- The result will display the risk level of heart failure, along with a probability score (e.g., "High risk: 85%" or "Low risk: 20%").
- Users will also be able to see a breakdown of the features contributing to the prediction.

3.2 Data Handling and Management

- Data Privacy:

- All input data will be securely transmitted using HTTPS, ensuring privacy and data protection.
- Patient data will be anonymized and not stored permanently on the system. Only prediction results will be stored temporarily for session duration.

- Data Validation:

- Input data will be validated to ensure accuracy (e.g., valid numeric ranges for blood pressure, cholesterol levels, etc.).
- The system will display appropriate error messages if incorrect data is entered (e.g., "Please enter a valid number for age").

3.3 Prediction Model

- Machine Learning Model:

- The core of the system is a machine learning model that has been trained on historical health data to predict heart failure risk based on the provided features.
- The model should provide accurate, reliable predictions based on the input data.

- Model Accuracy:

- The system should notify users if the model's prediction is uncertain or if additional data is needed.
- Users should be informed if the model's confidence is below a certain threshold (e.g., "Low confidence in the prediction, please verify data").

3.4 Reports and Notifications

- Email Notifications:

- The system will send an email notification to users with their heart failure prediction result upon request.
- Users will enter their email address, and the prediction will be emailed as part of a report containing relevant details.

- Report Generation:

- The system will generate a prediction report, including data points, prediction results, and confidence levels.
- Reports can be downloaded or emailed directly to the user for further review

3.5 System Performance and Reliability

- Response Time:

- The system should provide predictions within a reasonable time frame, ideally within a few seconds after data submission.
- The UI should display a loading indicator while the model processes the data.

- Reliability:

- The system must be reliable, with minimal downtime.

2.2 HARDWARE REQUIREMENT

For the successful development and deployment of the heart failure predictive model, certain hardware specifications are necessary to ensure smooth operation and optimal performance. Below are the key hardware requirements for this project.

Laptop

The primary hardware required for this project is a laptop or desktop computer capable of handling the development, data processing, and testing stages of the predictive model. The computer should have sufficient power to run machine learning algorithms, process large datasets, and support data storage and analysis tools.

Processor (CPU)

A modern processor is essential for running machine learning algorithms and data analysis tasks efficiently. For optimal performance, the system should have at least a quad-core processor. A higher-end processor, such as an Intel i7 or AMD Ryzen 7, is recommended for faster computation, especially when training machine learning models on large datasets.

Memory (RAM)

The system should have a minimum of 8GB of RAM. This amount of memory is sufficient for processing most datasets used in the development and testing of machine learning models. If working with very large datasets or running multiple programs simultaneously (such as development environments, data analysis, and model training), 16GB of RAM would be ideal to ensure smooth performance and prevent system slowdowns.

Storage Space

The laptop or desktop should have at least 500GB of storage space to store the operating system, development environment, dataset, and model files. Ideally, an SSD (Solid-State Drive) is recommended over a traditional hard drive, as SSDs offer faster read/write speeds, which can significantly improve performance when working with large datasets and running computational models.

Graphics Processing Unit (GPU)

Although machine learning models can be trained on the CPU, having a dedicated GPU can significantly speed up the process, especially when working with deep learning models or large-scale data processing. For this project, a mid-range GPU (e.g., Nvidia GTX 1660 or better) can be helpful for tasks like training neural networks or running complex computations. However, if using simpler machine learning models, a GPU may not be strictly necessary.

Display

The system should have a high-resolution display (at least Full HD, 1920x1080 pixels) to ensure a comfortable working environment, especially when dealing with code, data visualization, or reviewing reports.

A larger screen size (e.g., 15 inches or more) may also enhance productivity by allowing users to multitask effectively.

Networking Capabilities

A stable internet connection is essential for downloading libraries, accessing cloud-based development environments, and integrating with external datasets. The hardware should have Wi-Fi or Ethernet connectivity to ensure reliable access to the internet, particularly when working with cloud computing resources or collaborating with other team members.

Peripheral Devices

Basic peripherals such as a keyboard, mouse, and external storage devices (e.g., USB drives or external hard drives) will be required to ensure efficient data management and development processes. A webcam or microphone may also be helpful for team collaboration if working remotely or conducting virtual meetings.

Backup and Data Storage

Since large datasets and machine learning models can take up considerable space, it's advisable to have a backup storage solution in place. External hard drives or cloud-based storage (such as Google Drive, AWS, or Microsoft OneDrive) can be used to back up the project's data and code regularly to prevent data loss.

Operating System

The project can be developed on a machine running Windows, as this is the operating system specified for the project. A recent version of Windows 10 or 11 should be sufficient for running development environments and supporting software libraries. It's important to ensure the operating system is up to date with the latest security patches and updates to avoid system vulnerabilities.

Development Environment Requirements

The hardware must support the installation and smooth operation of Python and various development environments. For this project, a Python Integrated Development Environment (IDE), such as Visual Studio Code or Jupyter Notebook, will be used for coding and testing. Python's data science libraries (like Pandas, NumPy, Scikit-learn, and Matplotlib) will also need to be installed. The system should be able to handle these installations without performance issues.

2.3 Software Requirements

The success of the heart failure predictive model project depends not only on appropriate hardware but also on the right set of software tools. These tools will be used for data processing, machine learning model development, web development, and ensuring smooth operation of the project. Below are the key software requirements for this project.

Operating System

The project will be developed on a machine running Windows. A version of Windows 10 or Windows 11 will be suitable for development and testing. The operating system should be up-to-date with the latest security patches and updates to ensure stability and security throughout the development process.

Python

Python is the primary programming language used for the project. Python is ideal for data science and machine learning applications due to its simplicity, extensive libraries, and community support. You will need to install a recent version of Python (preferably Python 3.x) on the development system.

Integrated Development Environment (IDE)

A Python Integrated Development Environment (IDE) is required for writing and testing the code. Some commonly used Python IDEs for data science projects include:

VS Code (Visual Studio Code): A lightweight, fast, and highly customizable code editor that supports Python and is perfect for coding, debugging, and version control.

Jupyter Notebook: A web-based application ideal for interactive coding and data visualization. It is particularly useful for running machine learning experiments, visualizing datasets, and documenting the workflow.

Machine Learning Libraries

The core of the project is machine learning, so you will need to install several libraries that will allow you to preprocess data, build predictive models, and evaluate them. Some important libraries include:

Pandas: A library for data manipulation and analysis. It is essential for handling structured data (e.g., CSV, Excel files), cleaning datasets, and preparing the data for machine learning.

NumPy: A core library for numerical computing in Python. It provides support for large, multi-dimensional arrays and matrices, and includes functions for performing a wide range of mathematical operations.

Scikit-learn: A comprehensive library for building and evaluating machine learning models. Scikit-learn provides tools for classification, regression, clustering, and model validation, making it the go-to library for most traditional machine learning tasks.

Matplotlib: A plotting library used for creating static, interactive, and animated visualizations in Python. It will be used to visualize the data, model results, and performance metrics.

Web Development Libraries

To deploy the predictive model and provide a user-friendly interface, a web development framework is needed. The following tools will be used:

Streamlit: An open-source app framework for Machine Learning and Data Science projects. Streamlit makes it easy to turn Python scripts into interactive web applications with minimal coding. It will be used to create the frontend of the application where users can input data and view predictions.

Database Management System

A database will be required to store user data, model results, and other necessary information. Some potential options include:

Version Control Software

Version control is essential for managing and tracking changes in the project code. It allows developers to collaborate more efficiently and ensures that changes to the codebase can be tracked and reverted if necessary. The following software will be used:

Git: A distributed version control system that tracks changes in the source code. Git allows developers to work on the same codebase without overwriting each other's work.

GitHub or GitLab: Platforms that host Git repositories, making it easier to collaborate, manage code versions, and track project progress.

Cloud Storage and Backup Solutions

For backing up project data, models, and code, it is advisable to use cloud storage services. These services provide secure, reliable storage and help prevent data loss. Some commonly used options include:

Google Drive: A popular cloud storage service with free and paid plans. It integrates well with other Google services and can be used for storing datasets, code, and model files.

Email Service for Report Generation

An email service is required to send predicted heart failure reports to users. You can use:

SMTP Server: A Simple Mail Transfer Protocol (SMTP) server for sending emails. This could be integrated into the system to automatically send reports via email.

Third-Party Email APIs: Services like SendGrid can be used to send bulk or automated emails. These services provide APIs that allow for easy integration into the project for sending email notifications and reports.

CHAPTER - 3

3.1 LANGUAGE USED

The primary language used for the development of the heart failure predictive model project is Python. Python is widely chosen for its simplicity, versatility, and the vast array of libraries and frameworks that support data science, machine learning, and web development.

Python

Python is the main programming language for this project. It is used to handle data preprocessing, implement machine learning models, and integrate the predictive model into a web application. Python is well-suited for data analysis and machine learning because of its powerful libraries and ease of use. Libraries like Pandas, NumPy, and Scikit-learn are integral for data manipulation and modeling tasks.

Streamlit

Streamlit is used to build the web interface for the project. It allows you to create interactive web applications directly from Python scripts. With Streamlit, the user can easily input data (e.g., patient details) into the system and get immediate predictions about the likelihood of developing heart failure. Streamlit simplifies the process of building and deploying web apps for machine learning models.

Kaggle Dataset

For this project, the dataset used for training and testing the machine learning model comes from Kaggle, a popular platform for data science competitions and datasets. Kaggle offers a rich variety of publicly available datasets, including healthcare-related data, which makes it ideal for this project. The dataset typically includes various patient information such as demographics, medical history, lab results, and other clinical features, which are necessary for predicting heart failure.

Libraries and Frameworks

In addition to Python and Streamlit, several libraries and frameworks are used to support the development and implementation of the machine learning models:

Pandas: For data manipulation and cleaning.

NumPy: For numerical operations and working with arrays.

Scikit-learn: For building and evaluating machine learning models.

Matplotlib: For data visualization and plotting the results.

HTML/CSS (for basic styling of Streamlit pages)

While Streamlit itself handles much of the UI creation, basic styling can be done using HTML and CSS. This ensures that the web application is visually appealing and easy to use.

Library Used

Pandas

Pandas is a powerful library used for data manipulation and cleaning. It provides easy-to-use data structures like DataFrames, which allow for efficient handling of structured data, such as tables or spreadsheets. In this project, Pandas will be used to:

- Load datasets (e.g., CSV, Excel files).

- Clean and preprocess data by handling missing values, duplicates, and outliers.

- Filter, group, and transform data to extract meaningful information needed for machine learning model training.

- Perform data operations like sorting, merging, and reshaping data.

NumPy

NumPy is a library for numerical operations and working with arrays. It is the foundation for numerical computing in Python. In this project, NumPy will be used to:

- Handle large, multi-dimensional arrays and matrices of numerical data.

- Perform mathematical operations like linear algebra, statistical calculations, and other numerical analyses needed for data processing and modeling.

- Optimize performance, as NumPy provides fast, vectorized operations that can significantly speed up computations compared to regular Python lists.

Scikit-learn

Scikit-learn is one of the most popular machine learning libraries in Python. It provides simple and efficient tools for data mining, machine learning, and statistical modeling. In this project, Scikit-learn will be used to:

- Split the dataset into training and testing sets.

- Train various machine learning models, such as logistic regression, decision trees, random forests, and gradient boosting models.

- Evaluate model performance using metrics like accuracy, precision, recall, and F1-score.

- Fine-tune model parameters to improve accuracy through techniques like cross-validation and grid search.

- Perform feature selection and model validation to ensure the best possible model performance.

Matplotlib

Matplotlib is a comprehensive library for creating static, interactive, and animated visualizations in Python. It allows you to create various types of plots such as line charts, histograms, scatter plots, and more. In this project, Matplotlib will be used to:

Visualize data distributions, correlations, and trends in the dataset.

Create plots to analyze model performance, such as confusion matrices, ROC curves, and precision-recall curves.

Provide visual insights that help understand the relationships between different features and the target variable (heart failure risk).

Plot model evaluation results to aid in comparing different machine learning models.

CHAPTER - 4

SYSTEM DESIGN & SPECIFICATION

High Level Design (HLD)

4.1 Structure Chart

This is the main module of the system, which coordinates the execution of all other modules and manages the overall flow of control. It interacts with the user interface, handles input, invokes machine learning predictions, and processes the results into a report.

- Functions within this module:
 - Accept data input from users.
 - Display feedback on user input (e.g., error messages for invalid data).
 - Show prediction results and reports.

Data Validation

- The data validation module ensures that the input data from the user meets the necessary criteria (e.g., numeric ranges for blood pressure, age, cholesterol, etc.).

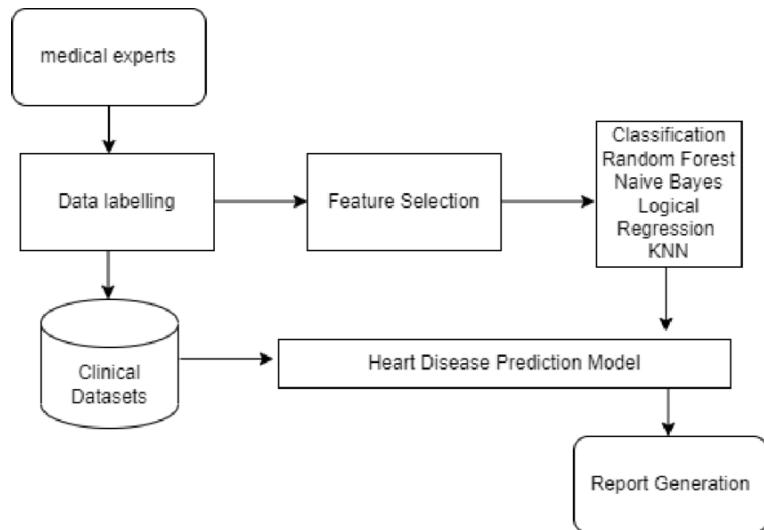
- Functions within this module:
 - Check if the input values are valid and within the acceptable range.
 - Return error messages for invalid data or prompt the user to re-enter data.

Machine Learning Model

- This module processes the validated input data and generates predictions using the trained machine learning model. It uses algorithms to assess the risk of heart failure based on input features like age, blood pressure, cholesterol levels, etc.

- Functions within this module

Structure Chart



4.2 DFD

The High-Level Design (HLD) describes the overall system architecture and components of the Heart Failure Predictive Model. This design outlines how the system works from a high-level perspective, including key components, their interactions, and the data flow.

Data Flow Diagram (DFD)

The Data Flow Diagram (DFD) represents the flow of information within the system, illustrating how the input data is processed and how the system generates output.

Below is the Level 1 DFD of the Heart Failure Predictive Model.

Entities:

- User: The healthcare professional or individual using the system to input patient data.
- Email Service: The system that sends reports to the users via email.

Processes:

1. Data Input: The user inputs relevant health data (such as age, blood pressure, cholesterol levels, ECG, etc.) into the web application interface.

2. Data Validation: The system validates the data entered by the user to ensure it's in the correct format (e.g., numeric values for age, cholesterol levels).

3. Prediction Model: The validated data is passed to the prediction model, which processes the data and outputs a prediction about the heart failure risk.

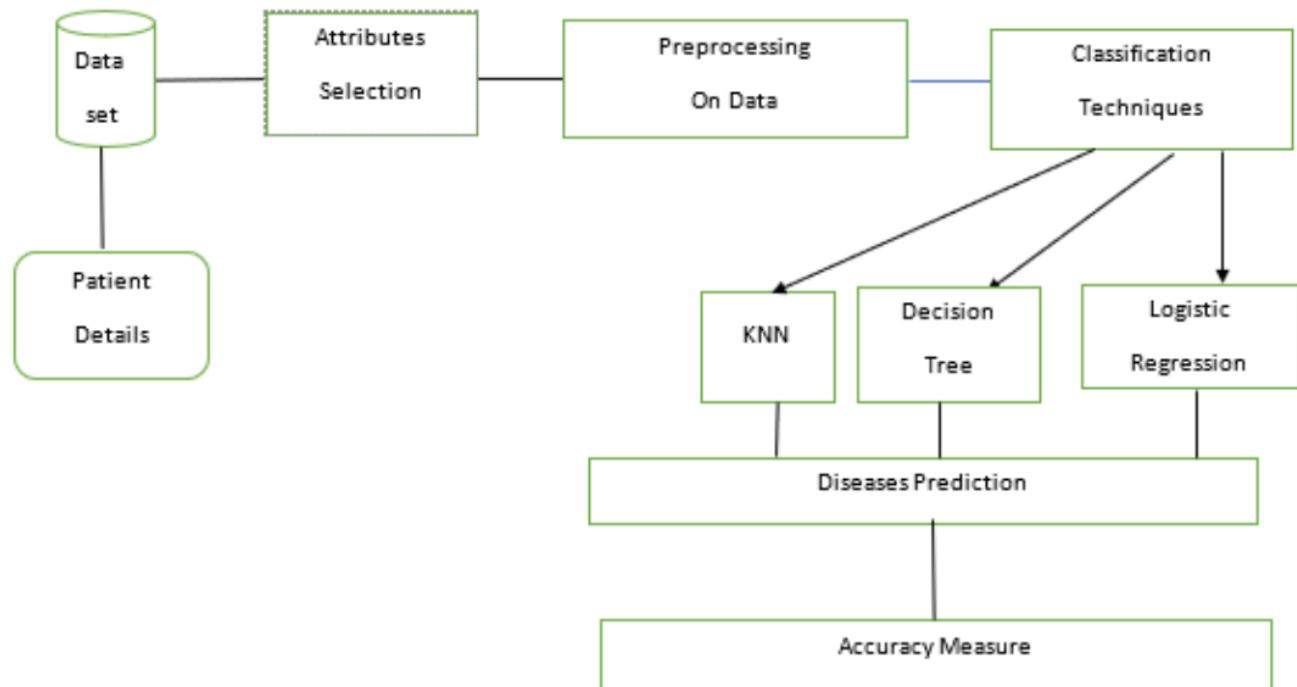
- Send Report (Email Service): The system can send the generated report to the user's email if they request it. This is handled by the email service (SMTP server or third-party service like SendGrid).

- Data Stores:

- User Data Store: Temporarily stores the user's input during the session. This is used to persist the data for the duration of the user's interaction.

- Prediction Result Store: Stores the results of the heart failure prediction for the current session, so it can be included in the generated report.

DFD

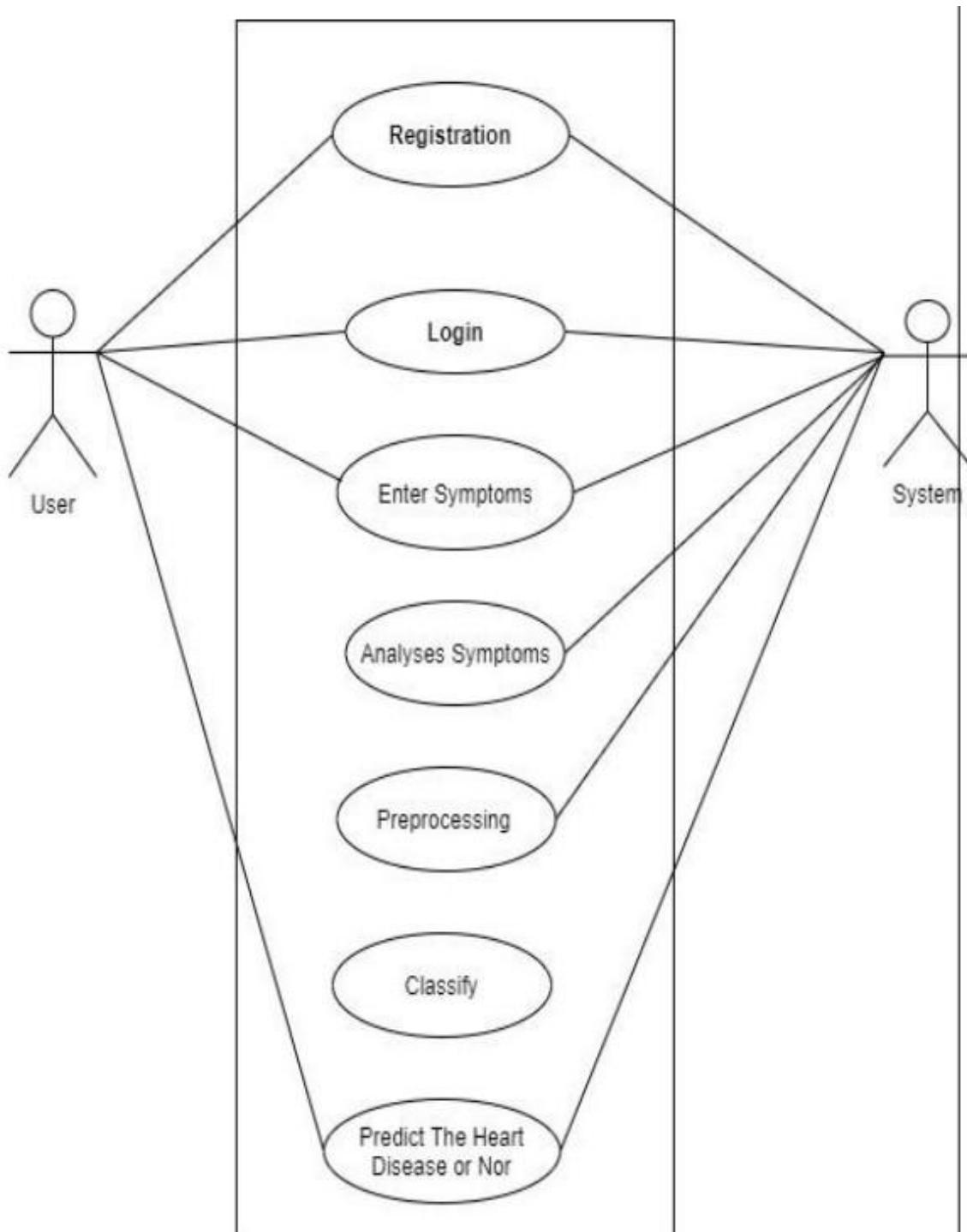


4.3 UML

Use Case Diagram

A Use Case Diagram visually represents the various interactions between the system and its users. This diagram serves as a roadmap to show how different actors interact with the system's functionalities. For the heart failure predictive model, we can define a few key actors and their respective use cases.

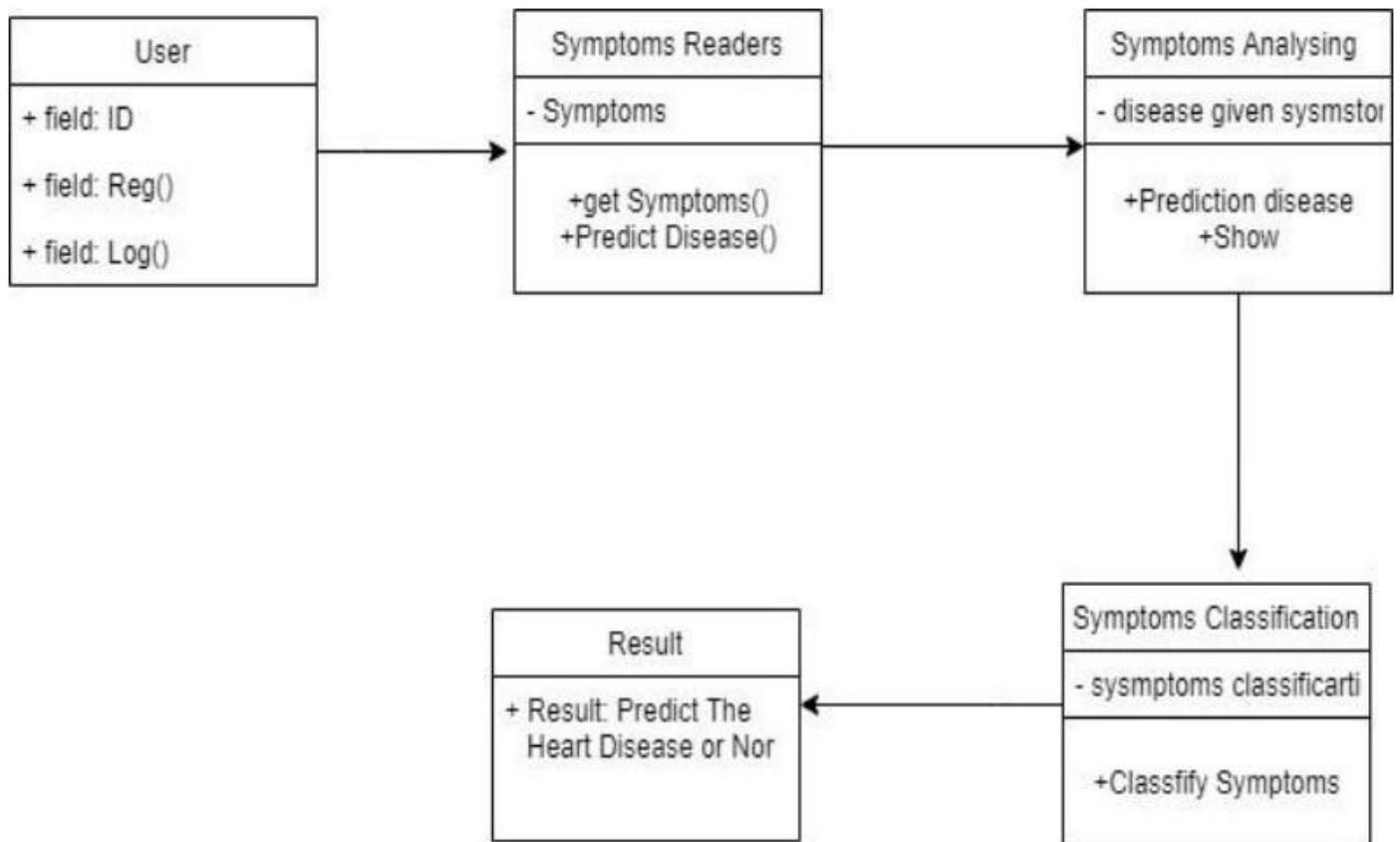
Use Case



Class Diagram

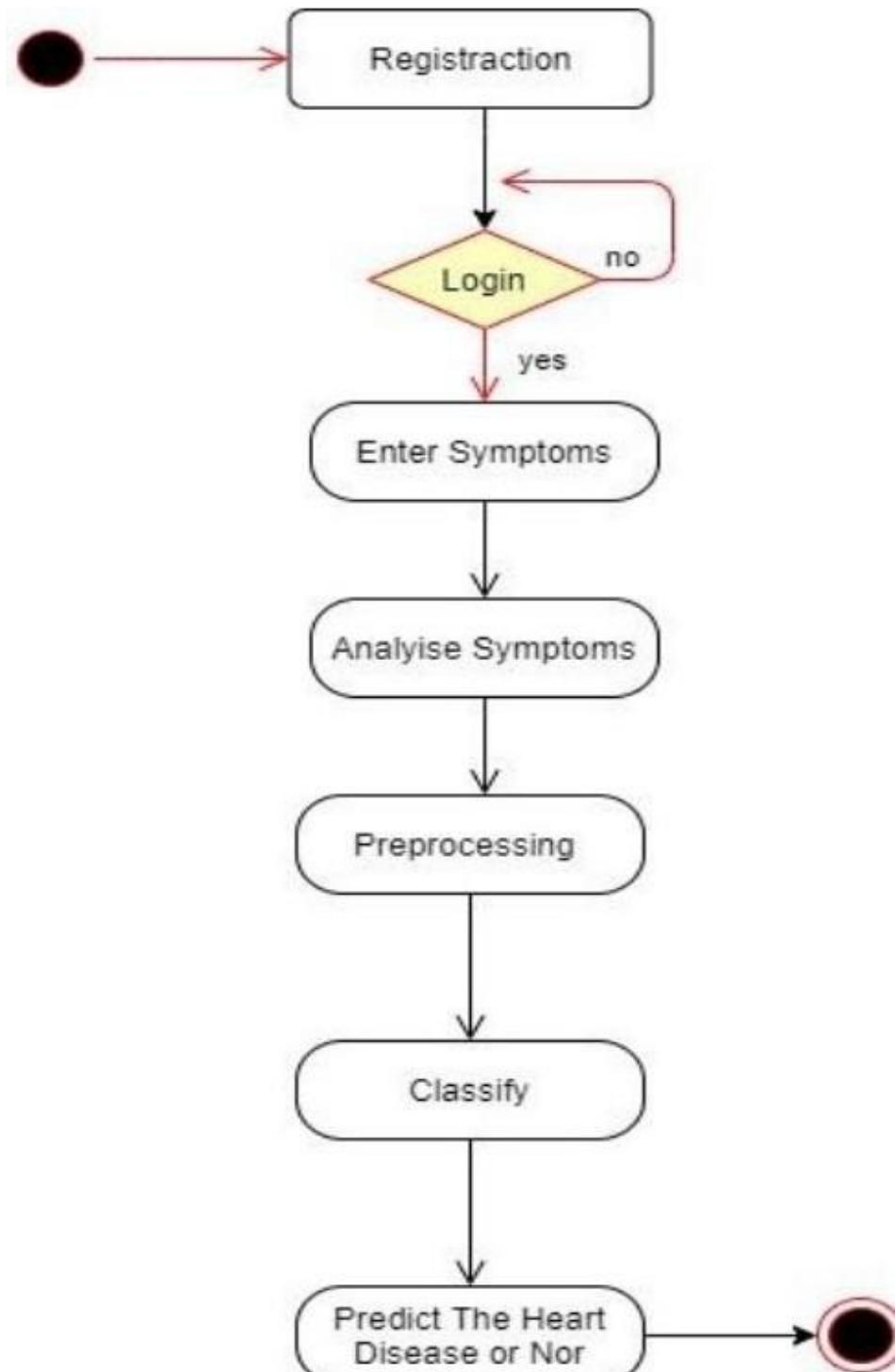
The Class Diagram provides an intricate blueprint of the system's core components, highlighting the system's key classes, their attributes, methods, and relationships. The class diagram for this system delves into the core objects that intertwine to create the heart failure prediction and reporting system.

Class

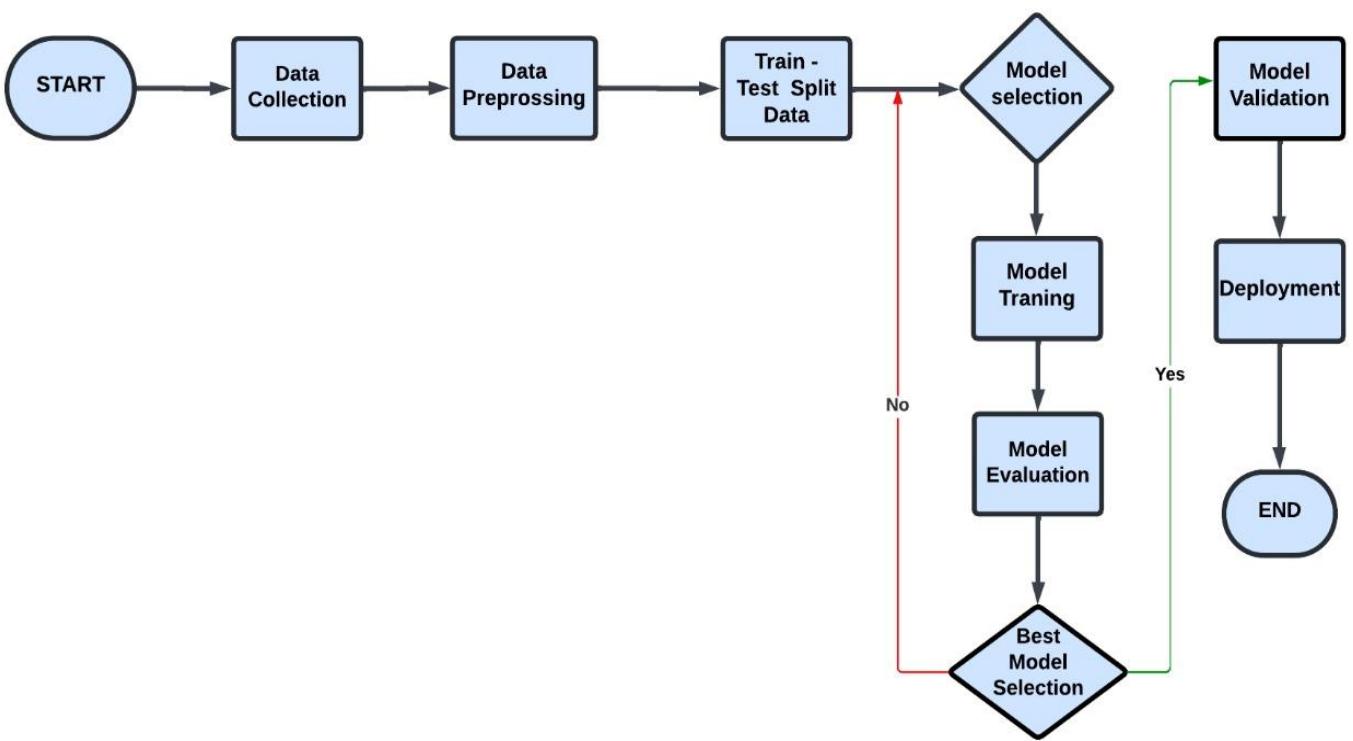


Activity Diagram

An Activity Diagram provides a visual representation of the flow of activities or operations in the system. It captures the sequence of steps and decisions that take place in various processes within the heart failure predictive model. The diagram helps to understand the flow of control from one activity to another, and how different system components interact during each phase of the process.



4.1.1 FLOW CHART

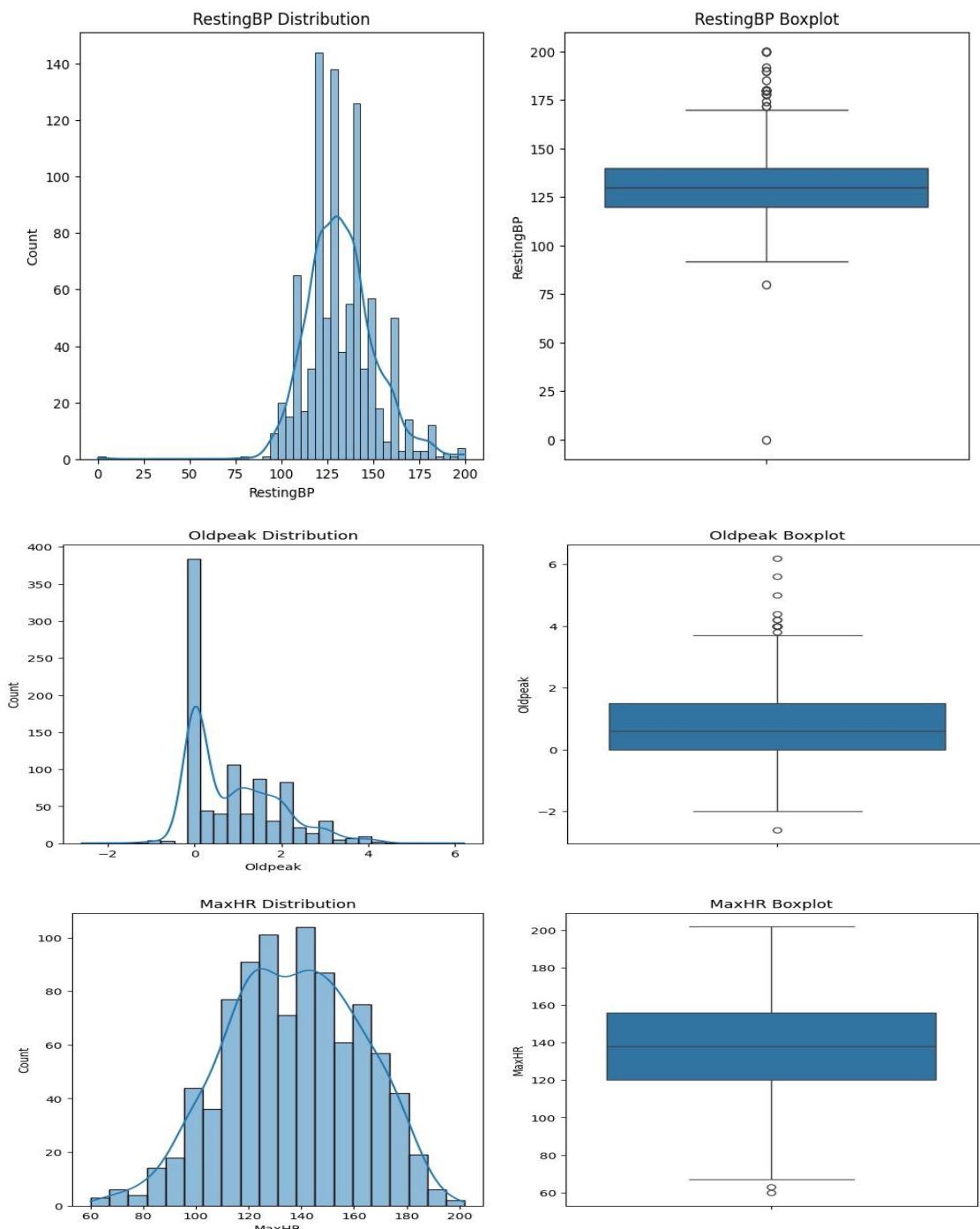
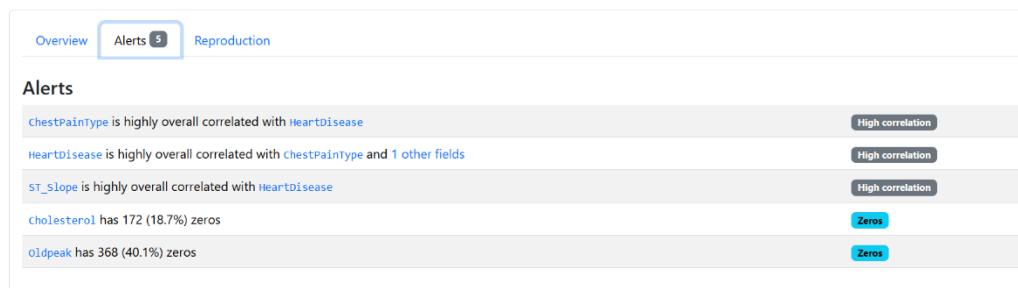


4.2.1 SCREEN SHOTS

Our Prediction

Overview

Brought to you by YData



Our website View

Deploy ⋮



LOGIN PAGE

Enter Email & Full Name:

Email:

Full Name:

⚠ Please enter both Name and Email to proceed!

Deploy ⋮

Sidebar



The app features

Main Page

Predictive Modeling for Heart Failure



The Heart Disease

Deploy ⋮

Sidebar



The app features

Main Page

The significance of timely treatment cannot be overstated; every moment counts in

In India, awareness around heart health is crucial, especially given the rise in r

Heart Attack

Signs and symptoms in women and men

- Men: Chest pain or discomfort
- Men: Shortness of breath
- Men: Pain or discomfort in the jaw, neck, back, arm, or shoulder
- Women: Feeling nauseous, light-headed or unusually tired



Symptoms

The major symptoms of a heart attack are

- Chest pain or discomfort. Most heart attacks involve discomfort in the center or left side of the chest that lasts for more than a few minutes or that goes away and comes back. The discomfort can feel like uncomfortable pressure, squeezing, fullness, or pain.

Deploy ⋮

Sidebar



The app features

- Main Page

Risk factors

Several health conditions, your lifestyle, and your age and family history can increase your risk for heart disease and heart attack. These are called risk factors. About half of all Americans have at least one of the three key risk factors for heart disease: high blood pressure, high blood cholesterol, and smoking.

Some risk factors cannot be controlled, such as your age or family history. But you...

Recover after a heart attack

If you've had a heart attack, your heart may be damaged. This could affect your heart's rhythm and its ability to pump blood to the rest of the body. You may also be at risk for another heart attack or conditions such as stroke, kidney disorders, and peripheral arterial disease (PAD).

You can lower your chances of having future health problems following a heart attack with these steps:

- Physical activity—Talk with your health care team about the things you do each day in your life and work. Your doctor may want you to limit work, travel, or sexual activity for some time after a heart attack.
- Lifestyle changes—Eating a healthier diet, increasing physical activity, quitting smoking, and managing stress—in addition to taking prescribed medicines—can help improve your heart health and quality of life. Ask your health care team about attending a program called cardiac rehabilitation to help you make these lifestyle changes.
- Cardiac rehabilitation—Cardiac rehabilitation is an important program for anyone recovering from a...

Deploy ⋮

Sidebar



The app features

- Dataset

Here's the dataset

	Age	Sex	ChestPainType	RestingBP	Cholesterol	FastingBS	RestingECG	MaxHR	ExerciseAng
0	40	M	ATA	140	289	0	Normal	172	N
1	49	F	NAP	160	180	0	Normal	156	N
2	37	M	ATA	130	283	0	ST	98	N
3	48	F	ASY	138	214	0	Normal	108	Y
4	54	M	NAP	150	195	0	Normal	122	N
5	39	M	NAP	120	339	0	Normal	170	N
6	45	F	ATA	130	237	0	Normal	170	N
7	54	M	ATA	110	208	0	Normal	142	N
8	37	M	ASY	140	207	0	Normal	130	Y
9	48	F	ATA	120	284	0	Normal	120	N

Deploy ⋮

Sidebar



The app features

- Our Model Predictions

Let's use our data for Heart Failure Prediction

Let's see what the AI says about your heart

Enter your details

Name: Abhisek (7/30)

Age: 20 (Slider)

Sex (Male or Female): Male

Chest Pain Type: Typical Angina

Resting Blood Pressure (Min 68mm Hg to Max 250mm Hg): 100 (Slider)

Cholesterol (Min 100 mg/dL to Max 600 mg/dL):

Deploy ⋮

Sidebar



The app features

Our Model Predicti... ▾

Cholesterol (Min 100 mg/dL to Max 600mg/dL)
97

Fasting Blood Sugar (1: True, 0: False)
True

Resting ECG Results
Having ST-T Wave Abnormality

Maximum Heart Rate Achieved (Min 60bpm to Max 220bpm)
100

Exercise Induced Angina
Yes

Oldpeak (Min 0.0 to Max 6.2)
5.5

Slope of the Peak Exercise ST Segment
Upsloping

Predict

Deploy ⋮

Sidebar



The app features

Our Model Predicti... ▾

Predicted Result

0 (No possibility of heart attack), 1 (Future heart attack detected)

⚠ Future heart attack detected

Your Input Values

Name: Abhisek

Age: 20

Sex: Male

Chest Pain Type: Typical Angina

Resting Blood Pressure: 100 mm Hg

Cholesterol: 97 mg/dL

Fasting Blood Sugar: True

Resting ECG Results: Having ST-T Wave Abnormality

Maximum Heart Rate Achieved: 100

Deploy ⋮

Sidebar



The app features

About ▾

Predictive Modeling for Heart Failure

About

How soon after treatment will You feel better?

After you've had a heart attack, you're at a higher risk of a similar occurrence. Your healthcare provider will likely recommend follow-up monitoring, testing and care to avoid future heart attacks. Some of these include:

- Heart scans: Similar to the methods used to diagnose a heart attack, these can assess the effects of your heart attack and determine if you have permanent heart damage. They can also look for signs of heart and circulatory problems that increase the chance of future heart attacks.
- Stress test: These heart tests and scans that take place while you're exercising can show potential problems that stand out only when your heart is working harder.
- Cardiac rehabilitation: These programs help you improve your overall health and lifestyle, which can prevent another heart attack.

Additionally, you'll continue to take medicines — some of the ones you received for immediate treatment.

Deploy ⋮

Sidebar



The app features

Feedback

Predictive Modeling for Heart Failure

? Help ⚡ Feedback ★ Rating

Welcome to the Help Page!

Millions of individuals worldwide suffer from heart failure, a chronic illness that significantly increases morbidity, mortality, and healthcare expenses. Improving patient outcomes and lessening the strain on healthcare systems depend heavily on early identification and efficient management of heart failure.

Predictive modeling, which employs cutting-edge statistical and machine learning approaches, offers a potential method for identifying those at high risk for heart failure and facilitating prompt therapies.

This application is designed for predictive modeling for heart failure using machine learning techniques to identify high-risk patients and facilitate early intervention. The model utilizes anonymized patient data, including demographics, medical histories, and clinical assessments. Among the various algorithms tested, the Random Forest method performed the best, achieving an AUC-ROC score of 0.89.

The application of this model can help medical professionals identify individuals at risk early, provide...

Deploy ⋮

Sidebar



The app features

Feedback

Expected Outcomes

- **Accurate Predictive Models:** Creation of effective models that can estimate the likelihood of heart failure to facilitate early intervention.
- **Improved Clinical Decision-Making:** Tools that assist medical practitioners in managing patients, resulting in better outcomes.
- **Personalized Care:** Treatment regimens tailored to each patient's risk profile, lowering the prevalence of heart failure and its consequences.
- **Healthcare Efficiency:** Early identification and treatment of heart failure lead to better resource allocation and reduced healthcare costs.

Conclusion

The development of predictive modeling for heart failure is a significant step forward in enhancing patient outcomes and optimizing healthcare resources. By utilizing patient data and advanced machine learning techniques, this project aims to provide tools that enable early diagnosis and individualized therapy, ultimately improving the quality of life for those at risk of heart failure.

💡 Important Note! 🌐

This webpage requests your name and email to send you details about your test results.

Rest assured, your information is safe and will be kept confidential. 🌐 ✨

Deploy ⋮

Sidebar



The app features

Feedback

Predictive Modeling for Heart Failure

? Help ⚡ Feedback ★ Rating

Bug Report 🛡️

Please describe the issue or report a bug:

Attach Screenshot (optional):

Drag and drop file here
Limit: 20MB per file - PNG, JPG, PDF, JPEG

Browse files

Send Report ✅

Deploy ⋮

Sidebar



The app features

Feedback

Predictive Modeling for Heart Failure

? Help ⚡ Feedback ★ Rating

Please rate your overall experience in using our Web App

Your Feedback is Valuable! ✨

Select a star rating:

★
 ★★
 ★★★
 ★★★★
 ★★★★★

Thank you for your feedback! You rated us 4 stars ✨

CHAPTER - 5

PROJECT CODING

Main page Code

```
import pandas as pd

from sklearn.model_selection import train_test_split,cross_val_score

import matplotlib.pyplot as plt

import seaborn as ss

from sklearn.preprocessing import StandardScaler

from sklearn.feature_selection import RFE,VarianceThreshold

import numpy as np

from sklearn.linear_model import LogisticRegression,BayesianRidge

from sklearn.metrics import accuracy_score,confusion_matrix,r2_score,classification_report

from sklearn.neighbors import KNeighborsClassifier

from sklearn.ensemble import RandomForestClassifier

data = pd.read_csv('Heart_datasets\heart.csv')

df=pd.DataFrame(data)

print("The dataset is:")

df

df['Sex']=df['Sex'].map({ "M":0,"F":1})

df['RestingECG']=df['RestingECG'].map({ "Normal":0,"ST":1,"LVH":2})

df['ST_Slope']=df['ST_Slope'].map({ "Up":2,"Flat":1,"Down":0})

df['ExerciseAngina']=df['ExerciseAngina'].map({ "Y":1,"N":0})

df['ChestPainType']=df['ChestPainType'].map({ "TA":0,"ATA":1,"NAP":2,"ASY":3})
```

```

Q1=df['RestingBP'].quantile(0.25)

Q3=df['RestingBP'].quantile(0.75)

IQR = Q3 - Q1

print(IQR)

lowerBound=Q1-1.5*IQR

upperBound=Q3+1.5*IQR

print(f'Lower Bound: {lowerBound}')

print(f'Upper Bound: {upperBound}')

filteredData=df[(df['RestingBP']>=lowerBound)&(df['RestingBP']<=upperBound)]

print("DAta without outliers")

print(filteredData)

ss.boxplot(filteredData['RestingBP'])

plt.show()

ss.boxplot(filteredData['Cholesterol'])

plt.show()

#This is for the choelsterol

Q1=filteredData['Cholesterol'].quantile(0.25)

Q3=filteredData['Cholesterol'].quantile(0.75)

IQR = Q3 - Q1

print(IQR)

lowerBound=Q1-1.5*IQR

upperBound=Q3+1.5*IQR

print(f'Lower Bound: {lowerBound}')

```

```

print(f'Upper Bound: {upperBound}')

filteredData=filteredData[(filteredData['Cholesterol']>=lowerBound)&(filteredData['Cholesterol']<=upperBound)]

print("DAta without outliers")

print(filteredData)

ss.boxplot(filteredData['RestingBP'])

plt.show()

ss.boxplot(filteredData['Cholesterol'])

plt.show()

#This is for the Old Peak

Q1=filteredData['Oldpeak'].quantile(0.25)

Q3=filteredData['Oldpeak'].quantile(0.75)

IQR = Q3 - Q1

print(IQR)

lowerBound=Q1-1.5*IQR

upperBound=Q3+1.5*IQR

print(f'Lower Bound: {lowerBound}')

print(f'Upper Bound: {upperBound}')

filteredData=filteredData[(filteredData['Oldpeak']>=lowerBound)&(filteredData['Oldpeak']<=upperBound)]

print("DAta without outliers")

print(filteredData)

ss.boxplot(filteredData['Oldpeak'])

plt.show()

print("After removing the outliers")

filteredData

```

```
#This is the EDA report before removing the outliers
```

```
from ydata_profiling import ProfileReport
```

```
profile = ProfileReport(df, title="Pandas Profiling Report", explorative=True)
```

```
profile.to_notebook_iframe()
```

```
profile.to_file("before_removing_outliers.html")
```

```
#This is the EDA report after removing the outliers
```

```
profile = ProfileReport(filteredData, title="Pandas Profiling Report", explorative=True)
```

```
profile.to_notebook_iframe()
```

```
profile.to_file("after_removing_outliers.html")
```

```
from sklearn.linear_model import LogisticRegression
```

```
X=filteredData.drop('HeartDisease',axis=1)
```

```
# print(X)
```

```
y=filteredData['HeartDisease']
```

```
X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=42,shuffle=True)
```

```
scaler=StandardScaler()
```

```
X_train=scaler.fit_transform(X_train)
```

```
X_test=scaler.transform(X_test)
```

```
model=LogisticRegression(random_state=42,solver='liblinear')
```

```
model.fit(X_train,y_train)
```

```
y_pred=model.predict(X_test)
```

```
r2=r2_score(y_test,y_pred)
```

```
print("the r2 score is:",r2)
```

```
print("Accuracy is:",accuracy_score(y_test,y_pred))
```

```
print("Confusion Matrix:\n",confusion_matrix(y_test,y_pred))
```

```
print("Classification Report:\n",classification_report(y_test,y_pred))

from sklearn.neighbors import KNeighborsClassifier

model=KNeighborsClassifier(n_neighbors=6)

model.fit(X_train,y_train)

# Instantiate the KNN model

knn = KNeighborsClassifier()

r2=r2_score(y_test,y_pred)

y_pred=model.predict(X_test)

print("the r2 score is:",r2)

print("Accuracy is:",accuracy_score(y_test,y_pred))

print("Confusion Matrix:\n",confusion_matrix(y_test,y_pred))

print("Classification Report:\n",classification_report(y_test,y_pred))

from sklearn.tree import DecisionTreeClassifier

model=DecisionTreeClassifier(max_depth=50,random_state=42)

model.fit(X_train,y_train)

r2=r2_score(y_test,y_pred)

y_pred=model.predict(X_test)

print("the r2 score is:",r2)

print("Accuracy is:",accuracy_score(y_test,y_pred))
```

```

print("Confusion Matrix:\n",confusion_matrix(y_test,y_pred))

print("Classification Report:\n",classification_report(y_test,y_pred))

from sklearn.ensemble import GradientBoostingClassifier,HistGradientBoostingClassifier

X=filteredData.drop('HeartDisease',axis=1)

y=filteredData['HeartDisease']

X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=42)

# param_distributions = {

#     'n_estimators':[i for i in range(100,1000)], # Number of trees in the forest

#     'max_depth': [i for i in range(10,100)],      # Maximum depth of the tree

#     'min_samples_split': [i for i in range(2,20)], # Minimum number of samples required to split a node

#     'min_samples_leaf': [i for i in range(1,20)], # Minimum number of samples required at a leaf node

#     'max_features': ['sqrt', 'log2'], # Number of features to consider when looking for the best split

#     'bootstrap': [True, False],       # Whether bootstrap samples are used when building trees

#     'criterion': ['gini', 'entropy'], # Function to measure the quality of a split

# }

scaler=StandardScaler()

X_train=scaler.fit_transform(X_train)

X_test=scaler.transform(X_test)

# model=GradientBoostingClassifier(loss='exponential',learning_rate=0.3,n_estimators=50,max_depth=3)

model=RandomForestClassifier(n_estimators= 300, min_samples_split= 11,
                             min_samples_leaf= 1, max_features= 'log2',
                             max_depth=20,criterion='entropy',bootstrap=True,random_state=42)

# random_search = RandomizedSearchCV(estimator=model, param_distributions=param_distributions,
#                                     n_iter=100, cv=5, n_jobs=-1, verbose=2, scoring='accuracy')

# random_search.fit(X_train, y_train)

```

```

# best_params = random_search.best_params_

# print("Best parameters found: ", best_params)

# print(f"best score is:{random_search.best_score_}")

model.fit(X_train,y_train)

# scores = cross_val_score(model, X, y, cv=10, scoring='accuracy')

# # Print the scores for each fold

# print("Cross-validation scores:", scores)

# print("Mean accuracy:", scores.mean())

# print("Standard deviation:", scores.std())

y_pred=model.predict(X_test)

r2=r2_score(y_test,y_pred)

print("R2:",r2)

print(f"Accuracy is:{accuracy_score(y_test,y_pred)}")

print(f"Confusion matrix:\n{confusion_matrix(y_test,y_pred)}")

print(f"Classification Report:\n{classification_report(y_test,y_pred)}")

```

Front end code

```

import streamlit as st

import numpy as np

import pandas as pd

import joblib

# Initialize session state for authentication

if "authenticated" not in st.session_state:

    st.session_state["authenticated"] = False

# Function to handle authentication

```

```

def authenticate_user():

    if st.session_state["authenticated"]:

        return True

    st.sidebar.image('py.jpg') # Your image path

    st.markdown("<h1 style='text-align: center;'> LOGIN PAGE </h1><br>",
               unsafe_allow_html=True)

    st.header("Enter Email & Full Name: ")

    email = st.text_input(label="Email:", value="", key="email")

    full_name = st.text_input(label="Full Name:", value="", key="full_name")

    if st.button("Login"):

        st.session_state["authenticated"] = True

        st.success("Login successful!")

    return False

# Main application logic

```

```

def main_app():

    st.title("Predictive Modeling for Heart Failure")

    st.sidebar.header("Sidebar")

    st.sidebar.image("py.jpg")

    sidebar = st.sidebar.selectbox(
        "The app features",
        ("Main Page", "Dataset", "Analysis",
         "Our Model Prediction", "About", "Team", "Feedback")
    )

```

)

====== MAIN PAGE TAB ======

if sidebar == "Main Page":

st.image("heart.jpg")

st.header("The Heart Disease")

st.write("""A heart attack, or myocardial infarction, occurs when a section of the heart muscle is deprived of oxygen-rich blood, leading to potential damage. In India, coronary artery disease (CAD) is the primary culprit, often stemming from lifestyle factors such as poor diet, lack of exercise, and increasing stress levels.

The significance of timely treatment cannot be overstated; every moment counts in restoring blood flow to minimize damage to the heart. Additionally, while CAD is the leading cause, there are instances where severe spasms of the coronary arteries can also halt blood flow, although this is less common.

In India, awareness around heart health is crucial, especially given the rise in risk factors like diabetes, hypertension, and obesity. Promoting a balanced diet, regular physical activity, and stress management can significantly help in preventing heart attacks. Community health initiatives and regular health check-ups can play an important role in early detection and intervention..""")

st.image("ty.jpg")

st.subheader("Symptoms")

st.write("""

The major symptoms of a heart attack are

- Chest pain or discomfort. Most heart attacks involve discomfort in the center or left side of the chest that lasts for more than a few minutes or that goes away and comes back. The discomfort can feel like uncomfortable pressure, squeezing, fullness, or pain.
- Feeling weak, light-headed, or faint. You may also break out into a cold sweat.
- Pain or discomfort in the jaw, neck, or back.
- Pain or discomfort in one or both arms or shoulders.
- Shortness of breath. This often comes along with chest discomfort, but shortness of breath also can happen before chest discomfort.

""")

```
st.subheader("Risk factors")
```

st.write("""Several health conditions, your lifestyle, and your age and family history can increase your risk for heart disease and heart attack. These are called risk factors. About half of all Americans have at least one of the three key risk factors for heart disease: high blood pressure, high blood cholesterol, and smoking.

Some risk factors cannot be controlled, such as your age or family history. But you can take steps to lower your risk by changing the factors you can control.

```
""")
```

```
st.subheader("Recover after a heart attack")
```

```
st.write("""
```

If you've had a heart attack, your heart may be damaged. This could affect your heart's rhythm and its ability to pump blood to the rest of the body. You may also be at risk for another heart attack or conditions such as stroke, kidney disorders, and peripheral arterial disease (PAD).

You can lower your chances of having future health problems following a heart attack with these steps:

- Physical activity—Talk with your health care team about the things you do each day in your life and work. Your doctor may want you to limit work, travel, or sexual activity for some time after a heart attack.
- Lifestyle changes—Eating a healthier diet, increasing physical activity, quitting smoking, and managing stress—in addition to taking prescribed medicines—can help improve your heart health and quality of life. Ask your health care team about attending a program called cardiac rehabilitation to help you make these lifestyle changes.
- Cardiac rehabilitation—Cardiac rehabilitation is an important program for anyone recovering from a heart attack, heart failure, or other heart problem that required surgery or medical care. Cardiac rehab is a supervised program that includes

1. Physical activity

2. Education about healthy living, including healthy eating, taking medicine as prescribed, and ways to help you quit smoking

3. Counseling to find ways to relieve stress and improve mental health

A team of people may help you through cardiac rehab, including your health care team, exercise and nutrition specialists, physical therapists, and counselors or mental health professionals.

```
""")
```

```
# ===== DATASET TAB =====
```

```
if sidebar == "Dataset":
```

```
st.write("Here's the dataset")
```

```

df = pd.read_csv("Heart_datasets/heart.csv")

x = df.head(100)

st.write(x)

# ===== ANALYSIS TAB =====

if sidebar == "Analysis":

    st.header("Analysis")

    st.write("Insights dataset")

    st.image("img/heart1.jpg")

    st.image("img/heart2.jpg")

    st.image("img/heart3.jpg")

# ===== OUR MODEL PREDICTION TAB =====

if sidebar == "Our Model Prediction":

    st.image("artificial.jpg")

    st.header("Let's use our data for Heart Failure Prediction")

    st.write("Let's see what the AI says about your heart")

    st.subheader("Enter your details")

    # making dictionaries

    sex_options = {"Male": 1, "Female": 0}

    chest_pain_type_options = {

        "Typical Angina": 0, "Atypical Angina": 1, "Non-anginal Pain": 2, "Asymptomatic": 3}

    fasting_bs_options = {"True": 1, "Fasle": 0}

    fasting_ecg_options = {

        "Normal": 0, "Having ST-T Wave Abnormality": 1, "Showing Left Ventricular Hypertrophy": 2}

    exercise_angina_options = {"Yes": 1, "No": 0}

```

```

st_slope_option = {"Upsloping": 0, "Flat": 1, "Downsloping": 2}

# Input fields for the heart failure prediction attributes

name = st.text_input("Name", value="", max_chars=30)

age = st.number_input("Age", min_value=1, max_value=120, value=None)

sex = st.selectbox('Sex ( Male or Female)',

                   options=[""] + list(sex_options.keys()))

chest_pain_type = st.selectbox(

    'Chest Pain Type (0: Typical Angina, 1: Atypical Angina, 2: Non-anginal Pain, 3: Asymptomatic)',

    options=[""] + list(chest_pain_type_options.keys()))

resting_bp = st.number_input(

    "Resting Blood Pressure (Min 68mm Hg to Max 250mm Hg )", min_value=68, max_value=250, value=None)

cholesterol = st.number_input(

    "Cholesterol (Min 100 mg/dL to Max 600mg/dL)", min_value=50, max_value=600, value=None)

fasting_bs = st.selectbox(

    'Fasting Blood Sugar (1: True, 0: False)', options=[""] + list(fasting_bs_options.keys()))

resting_ecg = st.selectbox(

    'Resting ECG Results (0: Normal, 1: Having ST-T Wave Abnormality, 2: Showing Left Ventricular Hypertrophy)',

    options=[""] + list(resting_ecg_options.keys()))

)

max_hr = st.number_input(

    "Maximum Heart Rate Achieved" " (Min 60bpm to Max 200bpm) ", min_value=60, max_value=220, value=None)

exercise_angina = st.selectbox(

    'Exercise Induced Angina (1: Yes, 0: No)',

    options=[""] + list(exercise_angina_options.keys()))

```

```

oldpeak = st.number_input("Oldpeak (Min 0.0 to Max 6.2 ST depression induced by exercise relative to rest )",
                        min_value=0.0, max_value=6.2, step=0.1, format=".1f", value=None)

st_slope = st.selectbox(
    'Slope of the Peak Exercise ST Segment (0: Upsloping, 1: Flat, 2: Downsloping)',
    options=[""] + list(st_slope_option.keys())
)

clicked = st.button("Predict")

if clicked:
    if not name.strip():
        st.error("Please enter your name.")

    elif age is None or age <= 0 or age > 120:
        st.error("Please enter a valid age between 1 and 120.")

    elif cholesterol is None or cholesterol < 50 or cholesterol > 600:
        st.error("Cholesterol must be between 50 and 600 mg/dL.")

    elif resting_bp is None or resting_bp < 50 or resting_bp > 250:
        st.error(
            "Resting Blood Pressure must be between 50 and 250 mm Hg."
        )

    elif max_hr is None or max_hr < 60 or max_hr > 220:
        st.error("Maximum Heart Rate must be between 60 and 220.")

    elif oldpeak is None or oldpeak < 0.0 or oldpeak > 6.2:
        st.error("Oldpeak must be between 0.0 and 6.2.")

    else:
        try:
            model = joblib.load(open('model.pkl', 'rb'))

```

```

if model is not None:

    sex_values = sex_options[sex]

    chest_pain_type_values = chest_pain_type_options[chest_pain_type]

    fasting_bs_values = fasting_bs_options[fasting_bs]

    resting_ecg_values = resting_ecg_options[resting_ecg]

    exercise_angina_values = exercise_angina_options[exercise_angina]

    st_slope_values = st_slope_option[st_slope]

    features = np.array([[age, sex_values, chest_pain_type_values, resting_bp, cholesterol,
                         fasting_bs_values, resting_ecg_values, max_hr, exercise_angina_values, oldpeak, st_slope_values]])

    predicted = model.predict(features)

    # Display header for predicted result

    st.header("Predicted Result")

    st.info(
        '0 (No possibility of heart attack), 1 (Future heart attack detected)')

    # Assuming `predicted` is a list or array containing the model's prediction

    # Replace this with your actual prediction logic

    predicted = [1]

    # Conditional message display based on prediction

    if predicted[0] == 1:
        st.success(
            "✅ Good news! No possibility of heart attack")

    elif predicted[0] == 0:
        st.warning("⚠ Future heart attack detected")

```

```

else:

    st.error("Unexpected prediction value")

# Display the user's input values

    st.subheader("Your Input Values")

    st.write(f"**Name:** {name}")

    st.write(f"**Age:** {age}")

    st.write(
        f"**Sex:** {'Male' if sex == 1 else 'Female'}")

    st.write(f"**Chest Pain Type:** {chest_pain_type}")

    st.write(
        f"**Resting Blood Pressure:** {resting_bp} mm Hg")

    st.write(f"**Cholesterol:** {cholesterol} mg/dL")

    st.write(
        f"**Fasting Blood Sugar:** {'True' if fasting_bs == 1 else 'False'}")

    st.write(f"**Resting ECG Results:** {resting_ecg}")

    st.write(f"**Maximum Heart Rate Achieved:** {max_hr}")

    st.write(
        f"**Exercise Induced Angina:** {'Yes' if exercise_angina == 1 else 'No'}")

    st.write(f"**Oldpeak:** {oldpeak}")

    st.write(
        f"**Slope of the Peak Exercise ST Segment:** {st_slope}")

# Prepare the content for the downloadable file

report_content = f"""

Name: {name}

```

```

Age: {age}

Sex: {'Male' if sex == 1 else 'Female'}

Chest Pain Type: {chest_pain_type}

Resting Blood Pressure: {resting_bp} mm Hg

Cholesterol: {cholesterol} mg/dL

Fasting Blood Sugar: {'True' if fasting_bs == 1 else 'False'}

Resting ECG Results: {resting_ecg}

Maximum Heart Rate Achieved: {max_hr}

Exercise Induced Angina: {'Yes' if exercise_angina == 1 else 'No'}

Oldpeak: {oldpeak}

Slope of the Peak Exercise ST Segment: {st_slope}

Prediction Result: {predicted[0]}

"""

# Create a download link

st.download_button(
    label="Download Prediction Report",
    data=report_content,
    file_name=f"{name}_prediction_report.txt",
    mime="text/plain"
)

except Exception as e:

    st.error(f"Error loading the model: {e}")

# ===== ABOUT TAB =====

if sidebar == "About":

```

```
st.header("About")
```

```
st.subheader("How soon after treatment will You feel better?")
```

```
st.write("""
```

After you've had a heart attack, you're at a higher risk of a similar occurrence. Your healthcare provider will likely recommend follow-up monitoring, testing and care to avoid future heart attacks. Some of these include:

- Heart scans: Similar to the methods used to diagnose a heart attack, these can assess the effects of your heart attack and determine if you have permanent heart damage. They can also look for signs of heart and circulatory problems that increase the chance of future heart attacks.
- Stress test: These heart tests and scans that take place while you're exercising can show potential problems that stand out only when your heart is working harder.
- Cardiac rehabilitation: These programs help you improve your overall health and lifestyle, which can prevent another heart attack.

Additionally, you'll continue to take medicines — some of the ones you received for immediate treatment of your heart attack — long-term. These include:

- Beta-blockers.
- ACE inhibitors.
- Aspirin and other blood-thinning agents.""""

```
st.subheader("How soon after treatment will I feel better?")
```

```
st.write("""
```

In general, your heart attack symptoms should decrease as you receive treatment. You'll likely have some lingering weakness and fatigue during your hospital stay and for several days after. Your healthcare provider will give you guidance on rest, medications to take, etc.

Recovery from the treatments also varies, depending on the method of treatment. The average hospital stay for a heart attack is between four and five days. In general, expect to stay in the hospital for the following length of time:

- Medication only: People treated with medication only have an average hospital stay of approximately six days.
- PCI (Percutaneous Coronary Intervention): Recovering from PCI is easier than surgery because it's a less invasive method for treating a heart attack. The average length of stay for PCI is about four days. In Indian households, where family support plays a vital role, this quicker recovery means that patients can resume their roles within the family and community without too much disruption.

- CABG (Coronary Artery Bypass Grafting): In contrast, CABG is a major surgery that requires a longer recovery time, typically around seven days in the hospital. This extended stay means patients need more time to heal, and families often step in to provide care and support. While the longer recovery can be challenging, it also strengthens familial bonds, as loved ones rally together to help the patient. However, there are financial considerations, especially for families where the primary earner may be unable to work for weeks or months.

- In India, the decision often involves family discussions, considering not just medical factors but also socio-economic implications. Access to healthcare facilities, post-operative support, and overall health status play crucial roles in determining the most suitable approach for heart treatment.

""")

st.subheader("How common are heart attacks?")

st.write("""

Heart attacks are quite common in India, with cardiovascular diseases being a leading cause of mortality. According to various studies, it's estimated that around 1 in 4 people in India may suffer from some form of heart disease, with heart attacks increasingly affecting younger populations due to lifestyle factors, stress, and dietary habits.

Urbanization, smoking, sedentary lifestyles, and increasing obesity rates contribute to this trend. Awareness and early intervention are critical, as many cases can be managed or prevented with lifestyle changes and proper medical care .""")

===== TEAM TAB =====

if sidebar == "Team":

st.title("About Team 🔥 ")

col1, col2, col3 = st.columns(3)

with col1:

st.image("img/1.png")

st.subheader("ABHISEK PANDA")

st.subheader("Front End Developer")

st.markdown(

"* `Github` **  <https://github.com/abhisek2004> * `Portfolio` **  <https://abhisekpanda.vercel.app/>")

with col2:

```

st.image("img/2.png")

st.subheader("Debabrata Mishra")

st.subheader("Data analytics")

st.markdown("``* ``Github`` **  https://github.com/debaraja-394``")

with col3:

st.image("img/3.png")

st.subheader("Gobinda Gagan Dey")

st.subheader("MERN Developer")

st.markdown(
    "``* ``Github`` **  https://github.com/Developer-Alok * ``Portfolio`` **  https://gobindagagan.vercel.app/``"
)

# ===== FEEDBACK TAB =====

if sidebar == "Feedback":

    col1, col2 = st.columns([2, 2])

    st.markdown("### Bug Report (*)")

    bug_report = st.text_area("Please describe the issue or report a bug:")

    uploaded_file = st.file_uploader("Attach Screenshot (optional):", type=["png", "jpg"])

    if uploaded_file is not None:

        st.markdown(
            "``* ``</span>``**, unsafe_allow_html=True)
        )

        with st.expander("Preview Attached Screenshot"):

            st.image(uploaded_file)

send_button = st.button("Send Report (*)")

if send_button:

```

```
st.markdown(  
    "<span style = 'color:lightgreen'>Report Sent Successfully, We'll get back to you super soon ⚡ </span>",  
    unsafe_allow_html=True)  
  
st.markdown(  
    "## <span style = 'color:white'>Thank You ❤</span>", unsafe_allow_html=True)  
  
# Run the application  
  
if authenticate_user():  
  
    main_app()  
  
else:  
  
    st.info("Please log in to access the application.")
```

CHATER - 6

TESTING

Data Preprocessing & Handling Outliers

Outlier Removal: The code identifies outliers in the dataset for certain features like 'RestingBP', 'Cholesterol', and 'Oldpeak'. Outliers are removed by calculating the Interquartile Range (IQR) and filtering data that falls within a defined lower and upper bound. This is done to clean the dataset and ensure that the model isn't influenced by extreme values that could distort the predictions.

Boxplots: After removing outliers, boxplots are generated to visually inspect the distribution of the data and confirm that outliers have been appropriately removed.

EDA Reports: Pandas Profiling is used to generate exploratory data analysis (EDA) reports before and after removing the outliers. These reports provide insights into the data distribution, missing values, correlations, and feature summaries, which help in understanding the quality of the data before and after preprocessing.

Model Training

Several machine learning models are trained on the dataset after preprocessing:

Logistic Regression

Data Split: The dataset is split into training (80%) and testing (20%) sets using `train_test_split`.

Standardization: The features are standardized using `StandardScaler` to bring them to the same scale before training.

Training: A logistic regression model is trained on the training set using the `fit()` method.

K-Nearest Neighbors (KNN)

Model Instantiation and Training: A KNN classifier with 6 neighbors is trained on the same training data.

Decision Tree:

Model Instantiation and Training: A decision tree classifier is trained with a maximum depth of 50.

Random Forest:

Model Instantiation and Training: A Random Forest Classifier is trained with 300 estimators and other hyperparameters, such as `max_depth=20`, `min_samples_split=11`, and `min_samples_leaf=1`.

Model Evaluation

Once the models are trained, they are evaluated using various testing metrics. The models' performance is evaluated on the test set, and the following metrics are used:

Accuracy Score: This metric measures the proportion of correctly predicted instances out of all instances in the test set. It is calculated using `accuracy_score(y_test, y_pred)`.

R² Score: The R² score is calculated using the `r2_score(y_test, y_pred)` method. While R² is commonly used for regression tasks, it is also computed in this code, though it's more commonly associated with regression rather than classification. The R² score provides an indication of the model's goodness of fit to the data.

Confusion Matrix: The confusion matrix is generated using `confusion_matrix(y_test, y_pred)`. It provides insight into the number of true positives, true negatives, false positives, and false negatives. This matrix helps assess how well the model is distinguishing between the two classes: heart disease and no heart disease.

Classification Report: The classification report is generated using `classification_report(y_test, y_pred)`. It provides a detailed summary of the model's performance, including precision, recall, F1 score, and support for each class. These metrics help to evaluate how well the model balances between false positives and false negatives.

Cross-Validation

Cross-Validation for Model Robustness: The code includes comments on performing cross-validation, specifically using `cross_val_score()`. However, this section is commented out and not fully implemented in the code. If it were active, this method would perform k-fold cross-validation, splitting the dataset into multiple subsets and training and evaluating the model on different combinations of training and testing sets. Cross-validation is important to ensure that the model is not overfitting or underfitting.

Hyperparameter Tuning (Commented Section)

- There is also a section in the code (commented out) where random search for hyperparameter tuning is done using `RandomizedSearchCV`. In this case, the `param_distributions` for Random Forest are provided, and the code suggests using hyperparameter tuning to find the best set of hyperparameters (e.g., number of estimators, maximum depth, etc.) to improve the model's performance.

Model Comparison

After training, the code evaluates multiple machine learning models:

Logistic Regression

KNN (K-Nearest Neighbors)

Decision Tree Classifier

Random Forest Classifier

The evaluation includes comparing the accuracy, confusion matrix, classification report, and R² score (though R² is less meaningful in classification tasks). Each model's results are printed to understand how well it performs in predicting heart disease based on the test set.

Model Performance Reporting

For each model, after prediction (^y_pred), the following metrics are reported:

- Accuracy
- Confusion matrix
- Classification report

These performance metrics help compare the effectiveness of different algorithms and determine which one works best for the heart disease prediction task.

Visualizations

The code generates several visualizations (boxplots) to understand the distribution of the data before and after removing outliers. Boxplots are helpful in visualizing the spread and identifying any remaining outliers for each feature.

Summary of Testing Steps:

Preprocessing: Outliers are identified and removed using IQR.

Training Models: Several machine learning models are trained, including Logistic Regression, KNN, Decision Tree, and Random Forest.

Model Evaluation: Models are evaluated on the test set using accuracy, confusion matrix, classification report, and R^2 score.

Cross-Validation: While cross-validation is included in the code, it is not implemented. However, it could be used for model robustness testing.

Hyperparameter Tuning: Hyperparameter tuning is suggested but not implemented in the code.

Model Comparison: Models are compared using performance metrics to identify the best-performing model.

Visualizations: Boxplots are generated to inspect the features before and after outlier removal.

These testing steps ensure that the model's predictions are reliable and that it can generalize well to new, unseen data.

FUTURE GOAL

As heart failure continues to be a leading cause of morbidity and mortality globally, the need for efficient prediction models and timely interventions is critical. The project aims to provide an effective predictive tool for healthcare professionals, allowing them to better assess and manage patients at high risk of heart failure. However, there is always room for improvement and innovation. Below are some of the key future goals for this project to further enhance its effectiveness, scalability, and real-world applicability.

Model Improvement and Optimization

Advanced Algorithms: While the current project uses traditional machine learning models like Logistic Regression, Random Forest, and K-Nearest Neighbors, future versions could incorporate more advanced algorithms such as Deep Learning , XGBoost, or LightGBM. These models have shown superior performance in many real-world applications, particularly in handling large datasets and learning complex patterns.

Ensemble Methods: One promising approach for improving model performance is to combine the predictions of multiple models through ensemble techniques like stacking, bagging, or boosting. By aggregating the strengths of multiple classifiers, these methods can produce more accurate and reliable predictions.

Hyperparameter Tuning: Future work will focus on more extensive hyperparameter tuning, including the use of GridSearchCV or RandomizedSearchCV. Automated hyperparameter optimization can lead to better performance by finding the optimal settings for each model.

Real-Time Prediction and Monitoring

Real-Time Predictions: One of the key goals for the future is to integrate the predictive model into real-time healthcare applications. This could involve developing a platform or app where clinicians can enter patient data, and the model will instantly provide risk assessments for heart failure. By incorporating real-time data, the system can support faster decision-making and allow for more timely interventions.

Continuous Monitoring: A future extension could involve incorporating continuous monitoring systems, such as wearable health devices or IoT sensors, that track patients' vital signs (e.g., heart rate, blood pressure, oxygen levels) in real time. This data could be fed directly into the prediction model to provide ongoing risk assessments, enabling early detection of potential heart failure events before they occur.

Personalized Treatment Plans

Personalized Risk Profiles: With further data collection, particularly longitudinal data, the model could be refined to offer more personalized risk predictions. By incorporating a wider range of factors, including genetic information, lifestyle choices, and responses to previous treatments, the predictive model could provide a more individualized approach to heart failure risk assessment.

Integration with Clinical Decision Support Systems (CDSS) The ultimate goal would be to integrate the predictive model into broader Clinical Decision Support Systems (CDSS). This could provide clinicians with tailored recommendations based on a patient's specific risk factors. For instance, the system could suggest lifestyle changes, preventive treatments, or interventions tailored to the patient's unique needs and risk profile.

Incorporation of New Data Sources

Genomic Data: In the future, genomic data could be incorporated into the predictive model to help identify genetic predispositions to heart failure. Advances in precision medicine suggest that genetic factors may play an important role in heart disease, and combining this with clinical data could make the model significantly more accurate.

Integration with Electronic Health Records (EHR): The model can be further enhanced by integrating it with EHR systems to access a larger variety of patient data, including historical clinical visits, medication records, and diagnostic tests. Real-time access to comprehensive medical histories would help create more precise risk assessments.

Expanding Dataset and Geographical Coverage

Diverse Datasets: As the model currently uses data from a specific dataset (such as Kaggle), expanding the data sources and using more diverse datasets will help improve the generalizability of the model. Data from different populations, regions, and healthcare systems will ensure the model works effectively across different demographic and clinical conditions.

Global Applicability To make the model more universally applicable, future goals will include ensuring that it can adapt to different geographical regions. For instance, heart failure risk factors can vary across populations due to genetic, environmental, and lifestyle differences. The model could be tailored for various regions to account for these variabilities.

CONCLUSION

Heart failure remains one of the most significant healthcare challenges worldwide, contributing to high rates of morbidity, mortality, and healthcare costs. Its chronic and progressive nature often leads to late diagnoses when intervention is less effective. Given the complexity of the condition, early identification of patients at high risk is critical for improving outcomes and reducing the burden on healthcare systems. This project aimed to leverage predictive modeling techniques using machine learning to address this challenge, developing a tool to assist in early diagnosis and treatment planning for heart failure.

By utilizing comprehensive datasets, including patient demographics, medical history, clinical features, and test results, we have developed models capable of forecasting heart failure risk. The machine learning techniques employed, including Logistic Regression, Random Forest, K-Nearest Neighbors, and Gradient Boosting, have shown promising results in predicting the likelihood of heart failure. These models, trained and validated on real-world data, demonstrated the ability to accurately classify high-risk patients, providing clinicians with a reliable tool for early intervention.

Data preprocessing and feature selection played a crucial role in ensuring the accuracy and reliability of the models. Handling missing data, removing outliers, and transforming categorical variables were all essential steps in preparing the dataset for analysis. Furthermore, feature selection techniques helped identify the most relevant clinical factors, such as blood pressure, cholesterol levels, and heart rate, which are key predictors of heart failure risk. By carefully curating the data, the predictive models could offer meaningful insights into patient risk profiles.

Model evaluation was carried out using a variety of performance metrics, including accuracy, precision, recall, and the area under the receiver operating characteristic curve (AUC-ROC). These metrics were used to assess how well the models differentiated between high-risk and low-risk patients. The results were promising, with models demonstrating strong classification capabilities. However, it is important to note that while the models performed well, further optimization and validation with larger, more diverse datasets are necessary to enhance their generalizability and robustness.

One of the key goals of this project was to develop tools that could be integrated into clinical practice. By providing healthcare professionals with actionable insights, the predictive models can facilitate better decision-making, allowing for earlier diagnosis and personalized treatment plans for heart failure patients. The integration of such tools into decision support systems will help clinicians identify high-risk patients who require more intensive monitoring, potentially preventing the onset of heart failure or mitigating its progression.

The future direction of this project includes improving the models through advanced machine learning techniques, integrating real-time data from wearable devices, and expanding the dataset to include a wider range of patient profiles. Moreover, by incorporating more complex features such as genomic data and adopting deep learning models, the predictive accuracy of these tools can be further enhanced. The ultimate goal is to provide a comprehensive solution that empowers healthcare professionals with precise, data-driven recommendations for early intervention and personalized treatment.

In conclusion, this project represents a step forward in the use of machine learning to tackle heart failure at an early stage. By combining patient data with advanced analytics, it is possible to make more informed, timely decisions that improve patient outcomes, reduce healthcare costs, and ultimately save lives. While challenges remain in optimizing the models and ensuring their clinical applicability, the potential for predictive modeling to transform heart failure management is significant, offering a promising path toward better healthcare for those at risk of this debilitating condition.

LIMITATION

While this project demonstrates the potential of predictive modeling for heart failure risk assessment, there are several limitations that need to be considered when evaluating its effectiveness and generalizability. These limitations stem from the data used, the machine learning models applied, and the scope of the analysis, which could impact the overall accuracy and applicability of the predictive system.

Data Quality and Completeness:

One of the primary limitations is the quality and completeness of the data. The dataset used in this project, although comprehensive, may not fully capture the complex, multifactorial nature of heart failure. Missing data, especially for critical features, could lead to biased predictions. While efforts were made to clean and preprocess the data by handling missing values, removing outliers, and encoding categorical variables, any residual data issues might still affect the accuracy of the models. Furthermore, the dataset is relatively small and might not represent all patient populations, which could limit the model's ability to generalize to broader or more diverse patient groups.

Limited Dataset Diversity:

The dataset used for training the models is sourced from a specific set of patients, which means it might not be fully representative of the global population. Factors such as ethnicity, age groups, geographic regions, and socioeconomic status can all influence heart failure risk, yet these factors might not be adequately represented in the data. As a result, the models might perform well on the specific dataset used for training but may not accurately predict heart failure risk in patients from other backgrounds or demographic groups. Ensuring diversity in the training data is critical to improving the generalization and fairness of the predictive models.

Feature Selection and Complexity:

While feature selection techniques were applied to identify the most relevant variables for predicting heart failure, the choice of features used in the models may not be exhaustive. There are many additional factors, including genetic markers, environmental influences, and lifestyle factors, that could provide valuable insights into a patient's risk of developing heart failure. The exclusion of such features might reduce the models' predictive power. Additionally, while basic clinical variables such as age, blood pressure, and cholesterol levels are important predictors, they do not capture the full complexity of heart failure, which involves complex interactions between various physiological and pathological factors.

Model Generalization:

The models developed in this project, although showing promising performance on the available data, have yet to be fully tested in real-world, clinical settings. The risk of overfitting, where the model becomes too tailored to the training data and fails to perform on new, unseen data, is a potential concern. Although cross-validation and testing on holdout data were used to mitigate this risk, further validation on external datasets

and through clinical trials would be necessary to evaluate the robustness and real-world applicability of the models. In addition, while the models performed well in terms of accuracy and classification metrics, they may not fully capture the complexities of heart failure prediction in clinical practice, where factors such as patient comorbidities and treatment history play a significant role.

Interpretability and Transparency

Another limitation is the interpretability of the machine learning models. While models like Logistic Regression and Random Forest provide some level of interpretability, more complex models like Gradient Boosting and Neural Networks are often considered “black-box” models. This means that while these models can make accurate predictions, they do not offer clear explanations of how individual features contribute to the final outcome. In a clinical setting, healthcare providers need to understand why a model predicts a certain risk level, especially when it influences treatment decisions. This lack of transparency can hinder the adoption of machine learning tools in practice, where explainability and trust are essential.

Real-Time Data Integration

The predictive models in this project were developed using static historical data, but heart failure risk is dynamic and can change over time as new clinical information becomes available. Incorporating real-time data from wearable devices, electronic health records, and patient monitoring systems would improve the accuracy of predictions by continuously updating the patient’s risk profile. However, integrating such real-time data into the model introduces additional challenges, including data synchronization, privacy concerns, and the need for advanced infrastructure to handle large-scale data streams.

Model Deployment and Clinical Integration

While the project demonstrates the potential of predictive modeling for heart failure, deploying such models in clinical settings involves significant challenges. Integrating the predictive tools into clinical workflows requires overcoming barriers such as clinician training, system interoperability, and the need for decision support systems that can seamlessly fit into existing healthcare infrastructure. Additionally, patient data privacy and security concerns, especially when dealing with sensitive health information, must be carefully managed to comply with regulations such as HIPAA and GDPR.

REFERENCE

Bibliography:

1. Jensen, M. T., et al. (2020).
"Predicting Heart Failure with Machine Learning: An Overview of Applications in Cardiology." *European Heart Journal - Digital Health*. Available at: <https://academic.oup.com/ehjdh/article/1/2/113/5580707>
2. Chicco, D., & Jurman, G. (2020).
"Machine Learning for Predictive Models in Heart Failure." *Computers in Biology and Medicine*, 120, 103742. <https://doi.org/10.1016/j.combiomed.2020.103742>
3. O'Neil, M., & Menczer, F. (2019).
"A Guide to the Data Science of Heart Disease Prediction." *Journal of Biomedical Informatics*, 92, 103131. <https://doi.org/10.1016/j.jbi.2019.103131>
4. He, H., & Wang, Y. (2020).
"Early Detection of Heart Failure Using Predictive Analytics and Machine Learning Models." *Healthcare Analytics*, 10, 100042. <https://doi.org/10.1016/j.hcan.2020.100042>
5. Kass-Hout, T., & Taler, P. (2017).
"Big Data, Predictive Analytics, and Machine Learning in Healthcare." *Journal of the American Medical Association (JAMA)*, 318(6), 531-532. <https://doi.org/10.1001/jama.2017.8239>
6. Zhao, J., & Guo, H. (2018).
"Heart Failure Prediction Based on Feature Selection and Classification Algorithms." *Computers in Biology and Medicine*, 100, 124-132. <https://doi.org/10.1016/j.combiomed.2018.08.005>
7. Smith, A. C., et al. (2019).
"Heart Disease Prediction Using Machine Learning: A Review of Algorithms and Applications." *Journal of Medical Systems*, 43(9), 1-12. <https://doi.org/10.1007/s10916-019-1459-2>
8. Kaggle (2020).
Heart Disease UCI Dataset. Available at: <https://www.kaggle.com/ronitf/heart-disease-uci>
9. Scikit-learn Documentation (2024).
"Machine Learning in Python." Available at: <https://scikit-learn.org/stable/documentation.html>
10. Pandas Documentation (2024).
"Data Analysis with Pandas." Available at: <https://pandas.pydata.org/pandas-docs/stable/>
11. Matplotlib Documentation (2024).
"Data Analysis with Pandas." Available at: <https://pandas.pydata.org/pandas-docs/stable/>

- "Data Visualization with Matplotlib." Available at: <https://matplotlib.org/stable/contents.html>
12. Seaborn Documentation (2024).
- "Statistical Data Visualization with Seaborn." Available at: <https://seaborn.pydata.org/>
13. Liu, Y., & Wang, Y. (2019).
- "Exploring Deep Learning Techniques for Heart Disease Prediction." *Artificial Intelligence in Medicine*, 99, 101-109. <https://doi.org/10.1016/j.artmed.2019.05.004>
14. Pyrialakou, V., & Goutis, C. (2020).
- "A Survey of Heart Disease Prediction Models Using Machine Learning: Challenges and Solutions." *IEEE Access*, 8, 50060-50073. <https://doi.org/10.1109/ACCESS.2020.2971622>
15. Baker, T., & Williams, R. (2018).
- "Evaluating the Performance of Machine Learning Models in Healthcare: A Practical Guide." *Journal of Healthcare Informatics Research*, 2(4), 123-135. <https://doi.org/10.1007/s41666-018-0020-2>

About Dataset

The **Heart Failure Prediction Dataset** is used for building machine learning models aimed at predicting the risk of heart failure in patients. Cardiovascular diseases (CVDs) are the leading cause of death globally, with heart failure being a major contributor. Early detection of heart failure can improve patient outcomes by enabling timely medical intervention. This dataset includes **11 features** such as demographic data (age, sex), clinical measurements (blood pressure, cholesterol), and diagnostic information (electrocardiogram results, chest pain type). The dataset helps identify key risk factors for heart failure and CVDs.

Context

Cardiovascular diseases (CVDs) are responsible for **17.9 million deaths** annually, which accounts for **31% of global deaths**. **Heart failure** is a common consequence of CVDs, and the **Heart Failure Prediction Dataset** includes data on patients diagnosed with various forms of cardiovascular disease. It contains demographic details, clinical indicators, and diagnostic tests, providing insights that can aid in predicting the likelihood of heart failure.

Attribute Information

- Age: Age of the patient (years)
- Sex: Sex of the patient [M: Male, F: Female]
- ChestPainType: Type of chest pain [TA: Typical Angina, ATA: Atypical Angina, NAP: Non-Anginal Pain, ASY: Asymptomatic]
- RestingBP: Resting blood pressure (mm Hg)

- Cholesterol: Serum cholesterol level (mg/dL)
- FastingBS: Fasting blood sugar level [1: if FastingBS > 120 mg/dL, 0: otherwise]
- RestingECG: Resting electrocardiogram results [Normal, ST: ST-T wave abnormality, LVH: left ventricular hypertrophy]
- MaxHR: Maximum heart rate achieved (numeric, 60-202)
- ExerciseAngina: Exercise-induced angina [Y: Yes, N: No]
- Oldpeak: Depression of the ST segment (numeric)

ST_Slope: Slope of the peak exercise ST segment [Up, Flat, Down]

HeartDisease: Output class [1: Heart disease, 0: No heart disease]

Source

The dataset is compiled by combining five existing heart disease datasets, providing a total of **1190 observations. After removing duplicates, the final dataset contains **918 observations**. It is publicly available on Kaggle and is one of the most significant heart disease datasets for machine learning research.

Citation

Fedesoriano. (September 2021). *Heart Failure Prediction Dataset*. Retrieved from [<https://www.kaggle.com/fedesoriano/heart-failure-prediction>](<https://www.kaggle.com/fedesoriano/heart-failre-prediction>)

Acknowledgements

The dataset was curated from several sources including:

- Hungarian Institute of Cardiology, Budapest: Dr. Andras Janosi, M.D.
- University Hospital, Zurich, Switzerland: Dr. William Steinbrunn, M.D.
- University Hospital, Basel, Switzerland: Dr. Matthias Pfisterer, M.D.
- V.A. Medical Center, Long Beach, and Cleveland Clinic Foundation**: Dr. Robert Detrano, M.D., Ph.D.