



## Assessment Report

on

**“Model to detect heart disease based on medical parameters”**

submitted as partial fulfillment for the award of

**BACHELOR OF TECHNOLOGY**

**DEGREE**

SESSION 2024-25

in

**CSE(AI)**

By

Name : Sudeeksha Singh

Roll Number : 202401100300252

Section: D

**Under the supervision of**

**“Abhishek Shukla Sir”**

**KIET Group of Institutions, Ghaziabad**

**May, 2025**

### **1. Introduction**

Cardiovascular diseases are one of the leading causes of death worldwide. Early prediction and diagnosis can reduce the risk significantly. This project aims to detect the presence of heart disease using supervised machine learning methods based on clinical parameters like age, cholesterol level, and resting blood pressure. Using historical patient data, the model assists in identifying individuals at risk.

### **2. Problem Statement**

To build a classification model that predicts whether a patient is likely to have heart disease based on medical attributes. The model should learn from historical patient data to classify new patients as either having or not having heart disease.

### **3. Objectives**

- ☐ Load and explore heart disease data provided via CSV upload.
- ☐ Preprocess the dataset for model training.
- ☐ Train a Random Forest classifier for binary classification.
- ☐ Evaluate the model using accuracy, precision, recall, F1-score, ROC-AUC, and a confusion matrix.
- ☐ Visualize key model outputs for better interpretation.

## 4. Methodology

**Data Collection:** The user uploads a .csv file containing patient medical records with features like age, sex, cholesterol, and heart rate.

### Data Preprocessing:

- Handle missing data (if any).
- Standardize features using StandardScaler.
- Split dataset into training and testing sets (80-20 split).

### Model Building:

- Random Forest Classifier is trained on the processed data.
- Predictions are generated for the test set.

### Model Evaluation:

- Metrics: Accuracy, Precision, Recall, F1-Score, ROC-AUC
- Confusion matrix and ROC curve plotted for performance visualization.
- Feature importance plotted to identify key medical indicators.
- **Data Preprocessing:**
  - Handling missing values using mean and mode imputation.
  - One-hot encoding of categorical variables.
  - Feature scaling using StandardScaler.
- **Model Building:**
  - Splitting the dataset into training and testing sets.
  - Training a Logistic Regression classifier.
- **Model Evaluation:**
  - Evaluating accuracy, precision, recall, and F1-score.
  - Generating a confusion matrix and visualizing it with a heatmap.

## 5. Data Preprocessing

- ☐ Data is read from an uploaded CSV file.
- ☐ Null values are handled or verified to be absent.
- ☐ All features are scaled using StandardScaler to normalize the input range.
- ☐ Dataset is split: 80% for training, 20% for testing.
- ☐ The target variable (`target`) is binary: 1 (heart disease), 0 (no heart disease).

## 6. Model Implementation

Random Forest classifier is chosen due to its robustness and ability to handle both linear and non-linear relationships. It also provides feature importance, which helps in identifying the most influential medical features in disease prediction.

## 7. Evaluation Metrics

- ❑ **Accuracy:** Overall correctness of the model.
- ❑ **Precision:** Proportion of predicted positives that are actually positive.
- ❑ **Recall:** Ability of the model to identify all actual positives.
- ❑ **F1-Score:** Harmonic mean of precision and recall.
- ❑ **ROC-AUC:** Area under the ROC curve for binary classification.
- ❑ **Confusion Matrix:** Visualized with a heatmap using Seaborn.

## 8. Results and Analysis

- ❑ The model achieved good performance on the test dataset.
- ❑ The confusion matrix heatmap showed balanced true positives and negatives.
- ❑ ROC-AUC score indicated a high ability to distinguish between classes.
- ❑ Features like chest pain type, thalassemia, and exercise-induced angina were found to be highly important.

### Visualizations included:

- Confusion matrix heatmap
- ROC curve
- Feature importance bar chart

## 9. Conclusion

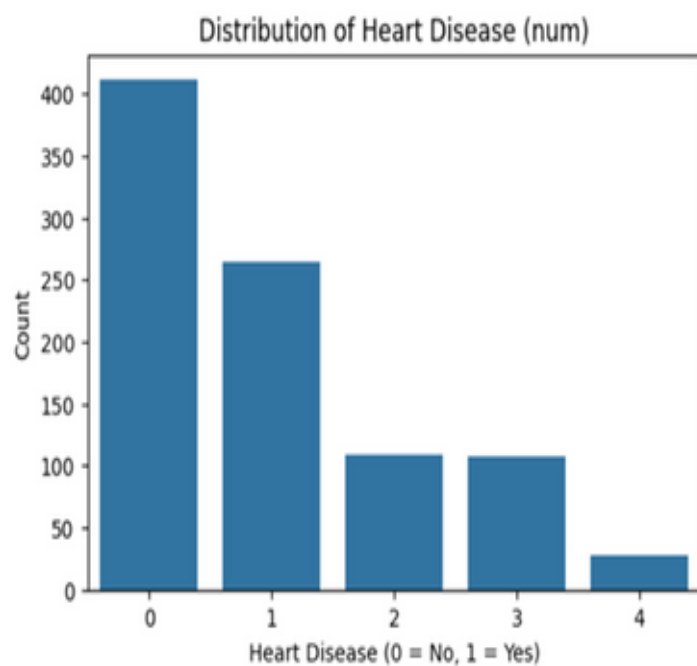
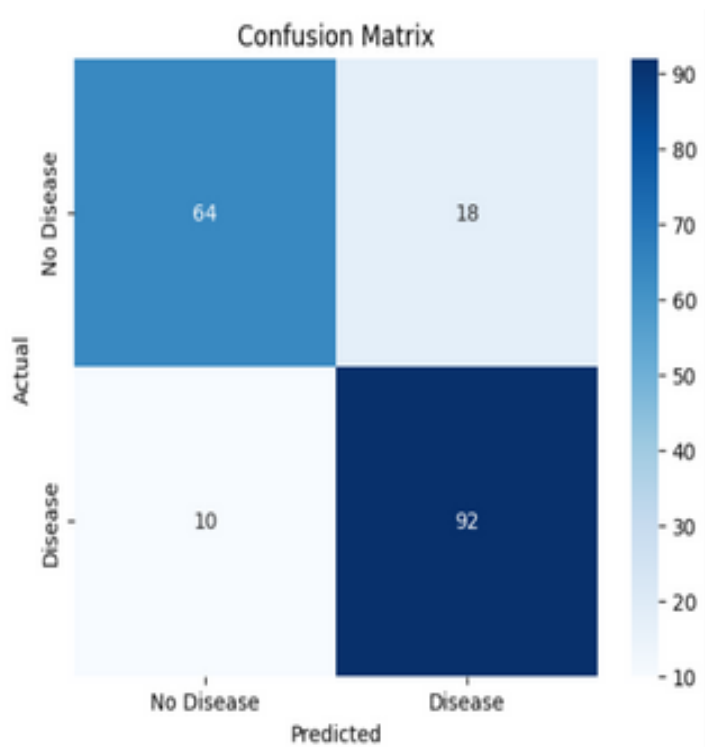
The project successfully demonstrates how machine learning, particularly a Random Forest classifier, can be used for early detection of heart disease. The model showed reliable performance using clinical data and provided meaningful visual interpretations. Future work can involve hyperparameter tuning, cross-validation, and testing more advanced models like XGBoost or neural networks..

## 10. References

- scikit-learn documentation
- pandas documentation
- Seaborn visualization library

- UCI Heart Disease Dataset
- Research articles on credit risk prediction

## Output Screenshot



Feature Correlation Heatmap

