



Project Report :- Heart Disease Prediction & Analysis Web Application

1. Executive Summary

This report details the end-to-end development of a machine learning-powered web application designed for the analysis and prediction of heart disease. The project encompasses a full data science workflow, including data preprocessing, exploratory data analysis (EDA), model training with hyperparameter optimization, and deployment via an interactive Streamlit dashboard. The primary outcome is a robust **Random Forest Classifier** that predicts patient risk with approximately **96% accuracy**. The final application serves as a practical tool for data-driven risk assessment, providing users with actionable insights and personalized health recommendations based on their clinical and lifestyle inputs.

2. Introduction

Cardiovascular diseases are a leading cause of mortality globally, making early and accurate risk identification a critical public health objective. This project addresses this challenge by leveraging machine learning to analyze complex patient data and identify high-risk individuals proactively. The goal is to translate a predictive model into a user-friendly tool that can aid in educational and preliminary screening contexts, promoting health awareness and timely medical consultation.

3. System Architecture & Dataset

3.1. Technical Architecture

The system is architected as a Python-based web application with two core components:

- **Machine Learning Backend:** A pre-trained Scikit-learn model (`best_heart_disease_model.joblib`) handles prediction logic. It is supported by a data preprocessing pipeline to ensure input data is correctly formatted.
- **Streamlit Frontend:** An interactive user interface that provides modules for data visualization and real-time risk prediction.

3.2. Dataset Description (`heart_disease_dataset.csv`)

The model was trained on a comprehensive heart disease dataset containing anonymized patient records. The dataset is characterized by a mix of numerical and categorical features, with "**Heart Disease**" (0 for No, 1 for Yes) serving as the binary target variable. Key features include:

Category	Features
Demographic	Age, Gender

Clinical Metrics	Cholesterol, Blood Pressure, Heart Rate, Blood Sugar
Behavioral Factors	Smoking, Alcohol Intake, Exercise Hours, Stress Level
Medical History	Family History, Diabetes, Obesity, Exercise Induced Angina

The dataset was found to be well-balanced, making it suitable for training a classification model without significant class imbalance issues.

4. Data Processing & Analysis

4.1. Data Preprocessing Pipeline

A standardized preprocessing pipeline was developed to prepare the data for modeling:

- Missing Value Imputation:** Null values in the Alcohol Intake feature were imputed with the string "None" to treat them as a distinct category.
- Categorical Encoding:** `sklearn.preprocessing.LabelEncoder` was applied to transform all categorical features (e.g., Gender, Smoking) into a numerical format suitable for the model.
- Feature Scaling:** All numerical features were standardized using `sklearn.preprocessing.StandardScaler`, which normalizes the data by removing the mean and scaling to unit variance. This prevents features with larger scales from dominating the model training process.

4.2. Exploratory Data Analysis (EDA)

The web application's visualization module provides several interactive EDA plots generated using Plotly and Seaborn, including:

- Univariate Analysis:** Histograms to show the distribution of key numeric features (Age, Cholesterol).
- Bivariate Analysis:** Count plots to explore the relationship between categorical predictors and the heart disease outcome.
- Multivariate Analysis:** A correlation heatmap to visualize the linear relationships between all numerical features.

5. Predictive Modeling

5.1. Model Selection & Training

A **Random Forest Classifier** was selected as the final model due to its high performance and robustness.

The model was trained on an 80/20 train-test split of the preprocessed data.

5.2. Hyperparameter Tuning

GridSearchCV was employed to systematically identify the optimal hyperparameters for the Random Forest model, maximizing its predictive power.

5.3. Performance Evaluation

The model's performance on the unseen test set was excellent, validating its reliability:

- **Overall Accuracy:** ~96%
- **Macro F1-Score:** ~0.96
- **Precision & Recall:** The model demonstrated a strong balance of precision and recall for both positive and negative classes, indicating it is effective at correctly identifying patients with and without heart disease.

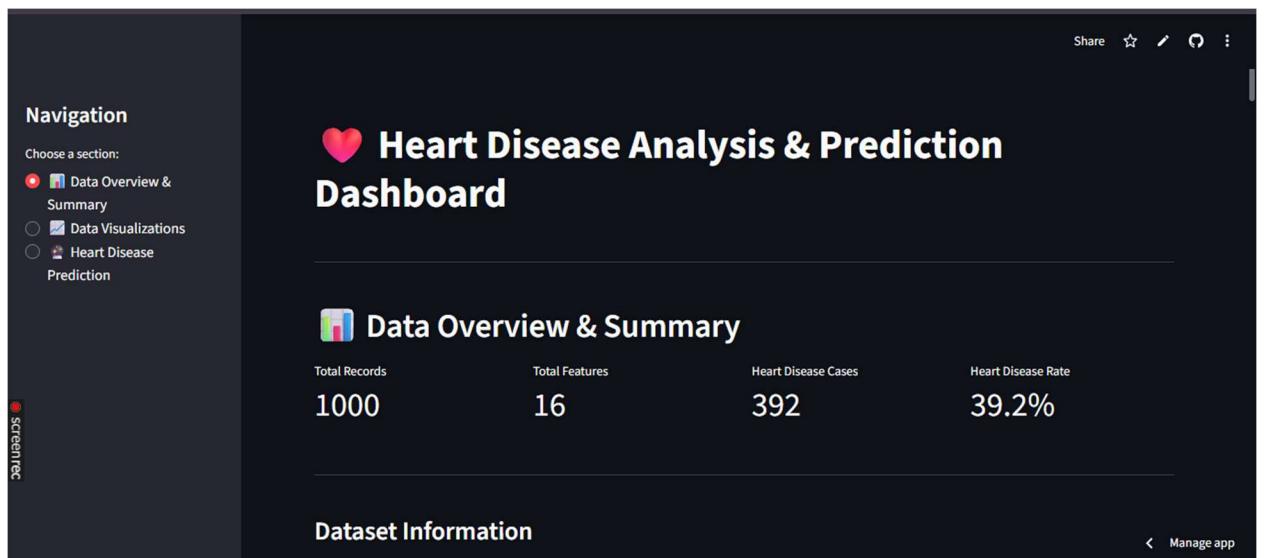
The finalized, trained model was serialized to best_heart_disease_model.joblib for deployment.

6. Deployed Application & User Interface

The application is deployed as an interactive Streamlit dashboard featuring:

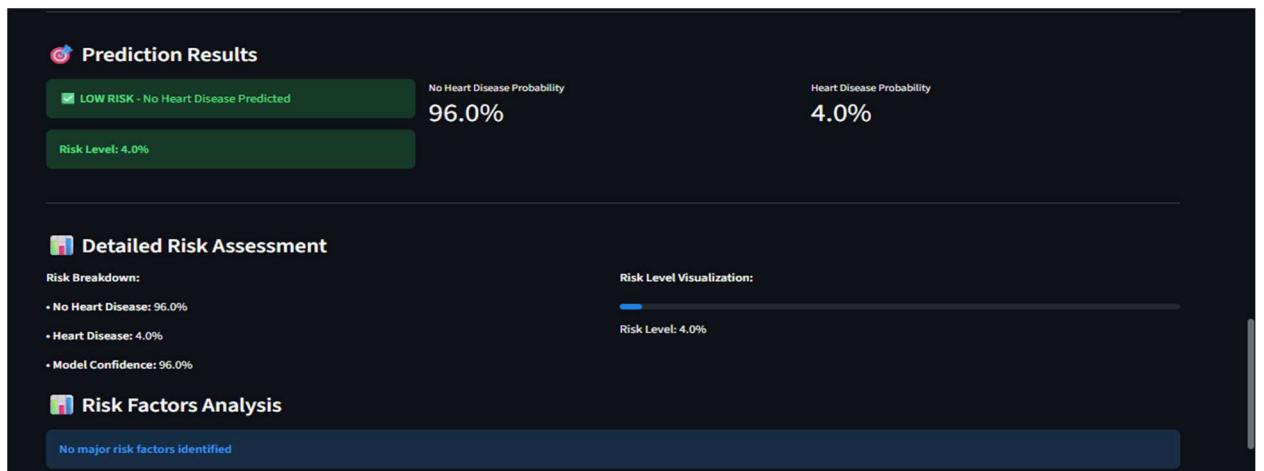
- **Modular Navigation:** A sidebar allows users to switch between the project overview, data visualizations, and the prediction engine.

Dataset Preview														
	Age	Gender	Cholesterol	Blood Pressure	Heart Rate	Smoking	Alcohol Intake	Exercise Hours	Family History	Diabetes	Obesity	S		
0	75	Female	228	119	66	Current	Heavy	1	No	No	Yes			
1	48	Male	204	165	62	Current	None	5	No	No	No	No		
2	53	Male	234	91	67	Never	Heavy	3	Yes	No	Yes	Yes		
3	69	Female	192	90	72	Current	None	4	No	Yes	No			
4	62	Female	172	163	93	Never	None	6	No	Yes	No			
5	77	Male	309	110	73	Never	None	0	No	Yes	Yes			
6	64	Female	211	105	86	Former	Heavy	8	Yes	Yes	Yes			
7	60	Female	208	148	83	Never	Moderate	4	No	Yes	Yes			
8	37	Female	317	137	66	Current	Heavy	3	No	Yes	Yes			
9	63	Male	204	141	68	Former	Heavy	8	No	Yes	No			

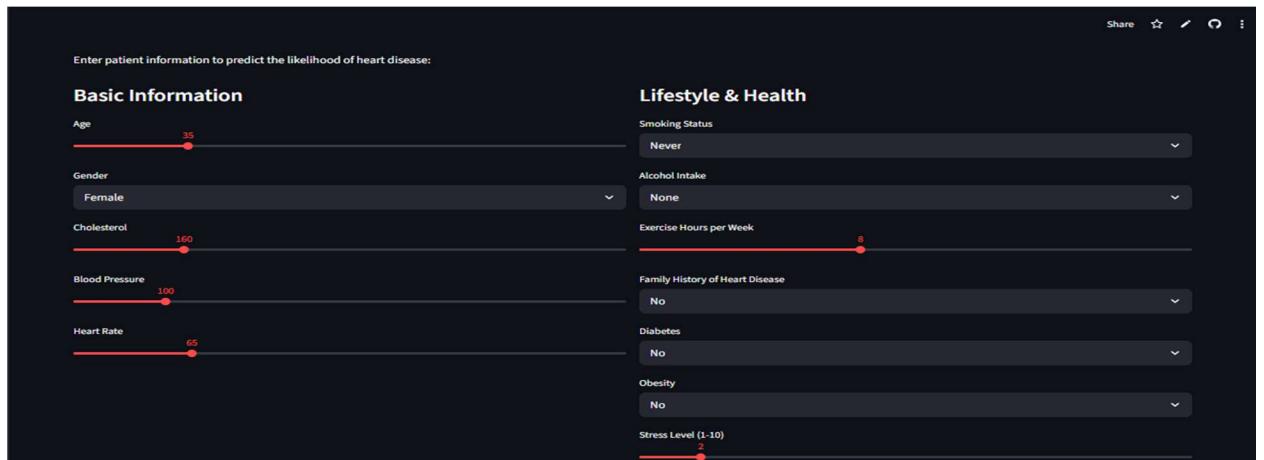


The screenshot shows a dark-themed dashboard titled "Heart Disease Analysis & Prediction Dashboard". On the left, a navigation sidebar lists "Data Overview & Summary", "Data Visualizations", and "Heart Disease Prediction". The main area features a "Data Overview & Summary" section with metrics: Total Records (1000), Total Features (16), Heart Disease Cases (392), and Heart Disease Rate (39.2%). Below this is a "Dataset Information" section. At the top right are sharing and management icons.

- **Dynamic Prediction Form:** A user-friendly form with sliders and dropdown menus for inputting patient data.
- **Results Interpretation:** The application outputs:
 - A clear risk classification ("Low Risk" or "High Risk").
 - A list of identified risk factors based on the user's inputs (e.g., "High Blood Pressure," "Obesity").
 - Actionable, context-aware recommendations (e.g., suggesting a consultation for high-risk individuals).



This section displays three main analysis components. The "Prediction Results" box shows a green bar indicating "LOW RISK - No Heart Disease Predicted" with a probability of 96.0% and a risk level of 4.0%. The "Detailed Risk Assessment" box provides a breakdown: No Heart Disease (96.0%), Heart Disease (4.0%), and Model Confidence (96.0%). It also includes a risk level visualization slider set at 4.0%. The "Risk Factors Analysis" box states "No major risk factors identified".



The screenshot shows a form for entering patient information. The "Basic Information" section includes sliders for Age (35), Cholesterol (160), Blood Pressure (100), and Heart Rate (65). The "Lifestyle & Health" section includes dropdowns and sliders for Smoking Status (Never), Alcohol Intake (None), Exercise Hours per Week (8), Family History of Heart Disease (No), Diabetes (No), Obesity (No), and Stress Level (1-10).

Heart Disease Analysis & Prediction Dashboard

Heart Disease Prediction

Example Values for Low Risk (Normal)

Low Risk Example Values:

- Age: 35 years
- Gender: Female
- Cholesterol: 160 mg/dL
- Blood Pressure: 100 mmHg
- Heart Rate: 65 bpm
- Smoking: Never
- Alcohol Intake: None
- Exercise Hours: 8 hours/week

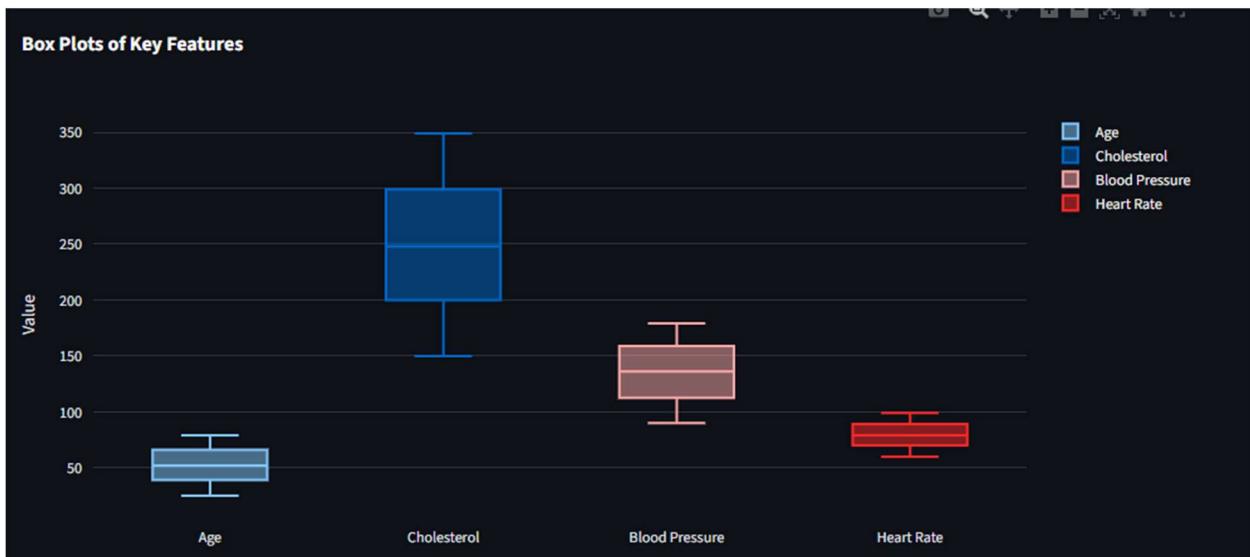
Low Risk Example Values (continued):

- Family History: No
- Diabetes: No
- Obesity: No
- Stress Level: 2 (Low)
- Blood Sugar: 80 mg/dL
- Exercise Induced Angina: No
- Chest Pain Type: Asymptomatic

- Demonstration Scenarios:** Pre-filled buttons for "Low Risk" and "High Risk" examples to showcase the tool's functionality.

7. Conclusion

This project successfully demonstrates the creation of a full-stack data science application, from raw data to a deployed, interactive tool. The Heart Disease Prediction Web Application stands as a robust proof-of-concept that effectively integrates a high-accuracy machine learning model with a user-centric interface. It serves as a valuable resource for educational purposes and highlights the potential of AI in promoting proactive health management.



Correlation Heatmap of Numerical Features

	Age	-0.0106725	0.002092796	x: Blood Pressure y: Age color: 0.002092796	-0.02136605	-0.04555547	-0.0416764	0.6468706
Age	1							
Cholesterol	-0.0106725	1	0.02184138		-0.008526932	0.01612381	0.09045762	0.3650408
Blood Pressure	0.002092796	0.02184138	1		-0.001675468	0.01192404	0.00225742	0.006900152
Heart Rate	0.02902667	-0.008526932	-0.001675468	1		-0.01354102	-0.04050449	0.0102396
Exercise Hours	-0.02136605	0.01612381	0.01192404	-0.01354102	1		-0.006957077	-0.03450328
Stress Level	-0.04555547	0.09045762	0.00225742	-0.04050449	-0.006957077	1	-0.007918463	0.007071059
Blood Sugar	-0.0416764	0.002483562	-0.05351618	0.0102396	-0.03450328	-0.007918463	1	-0.01300403
Heart Disease	0.6468706	0.3650408	0.006900152	0.01320899	-0.01422575	0.007071059	-0.01300403	1

