

# Heart Disease Risk Prediction

**Presented By :** Mostafa Gamal Abdel-Fatah Fouda

**Date :** 15/5/2025

## ➤ Problem Definition:

Heart disease remains the leading cause of death globally. Many patients remain undiagnosed until severe symptoms appear. There is a crucial need for early, accessible, and affordable prediction tools that help identify individuals at risk before complications occur.

## ➤ Proposed Solution:

This project delivers a complete machine learning pipeline to predict the presence of heart disease based on structured patient data. The system is deployed as a user-friendly web application using Streamlit, allowing users to input clinical parameters and receive instant feedback.

## ➤ Objectives:

- Automate heart disease risk prediction using machine learning.
- Provide an interactive and accessible web interface.
- Improve early detection and support health decisions.

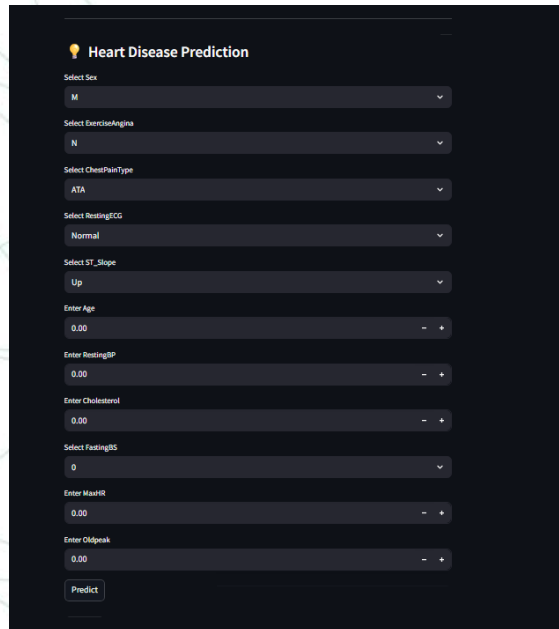
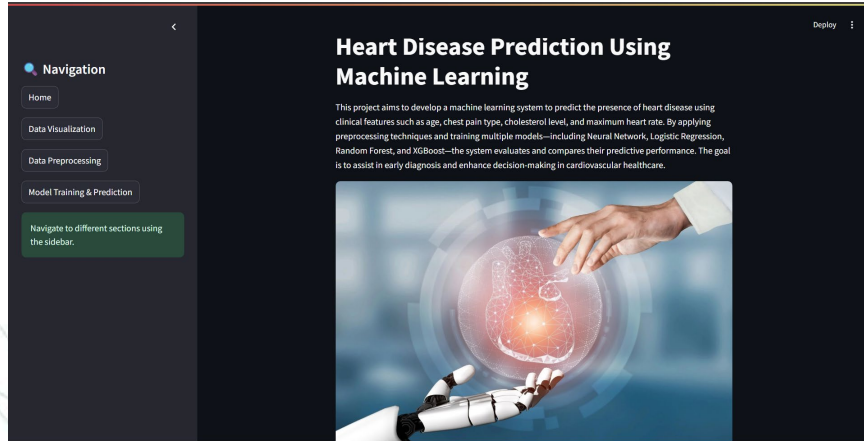
## ➤ Unique Value Proposition:

- **End-to-End Pipeline:** Covers all steps from data preprocessing to deployment.
- **Multi-Model Comparison:** Trained multiple models Random Forest, Logistic Regression, XGBoost and Neural Network.
- **Real-Time Prediction:** Immediate feedback via a live Web app.
- **Visual Insights:** Includes intuitive data visualizations for users and analysts.
- **Scalable & Extendable:** Can be enhanced with more features or applied to other diseases.

## ➤ Real-World Impact:

- Assists doctors with a quick second opinion.
- Empowers individuals to assess their risk anytime.
- Can be integrated into health screening tools or mobile apps in the future.

# Project Wireframe



## User Journey Overview:

### 1. Home Page – Project Overview

1. Brief introduction to the purpose of the app.
2. Explains the problem and the goal of the prediction system.
3. Navigation menu for moving between pages.

### 2. Preprocessing Page

1. Shows how raw data is cleaned and prepared.
2. Displays steps like: missing value handling, outlier removal, feature encoding & scaling.
3. Visual feedback for each preprocessing stage.

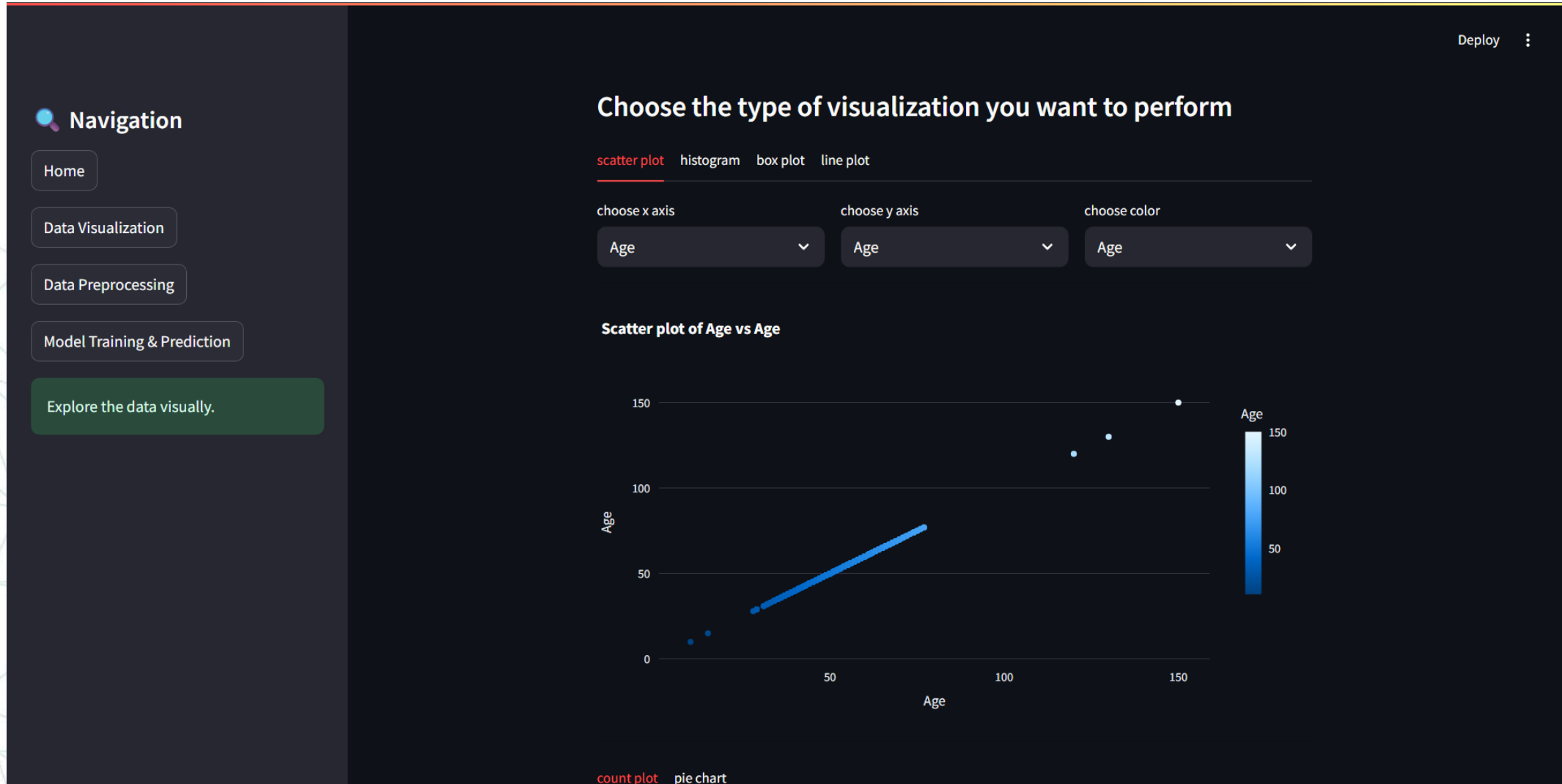
### 3. Visualizations Page

1. Exploratory Data Analysis (EDA) to understand feature distributions.
2. Charts like bar plots, pie charts, correlation heatmaps.
3. Helps users explore important risk indicators.

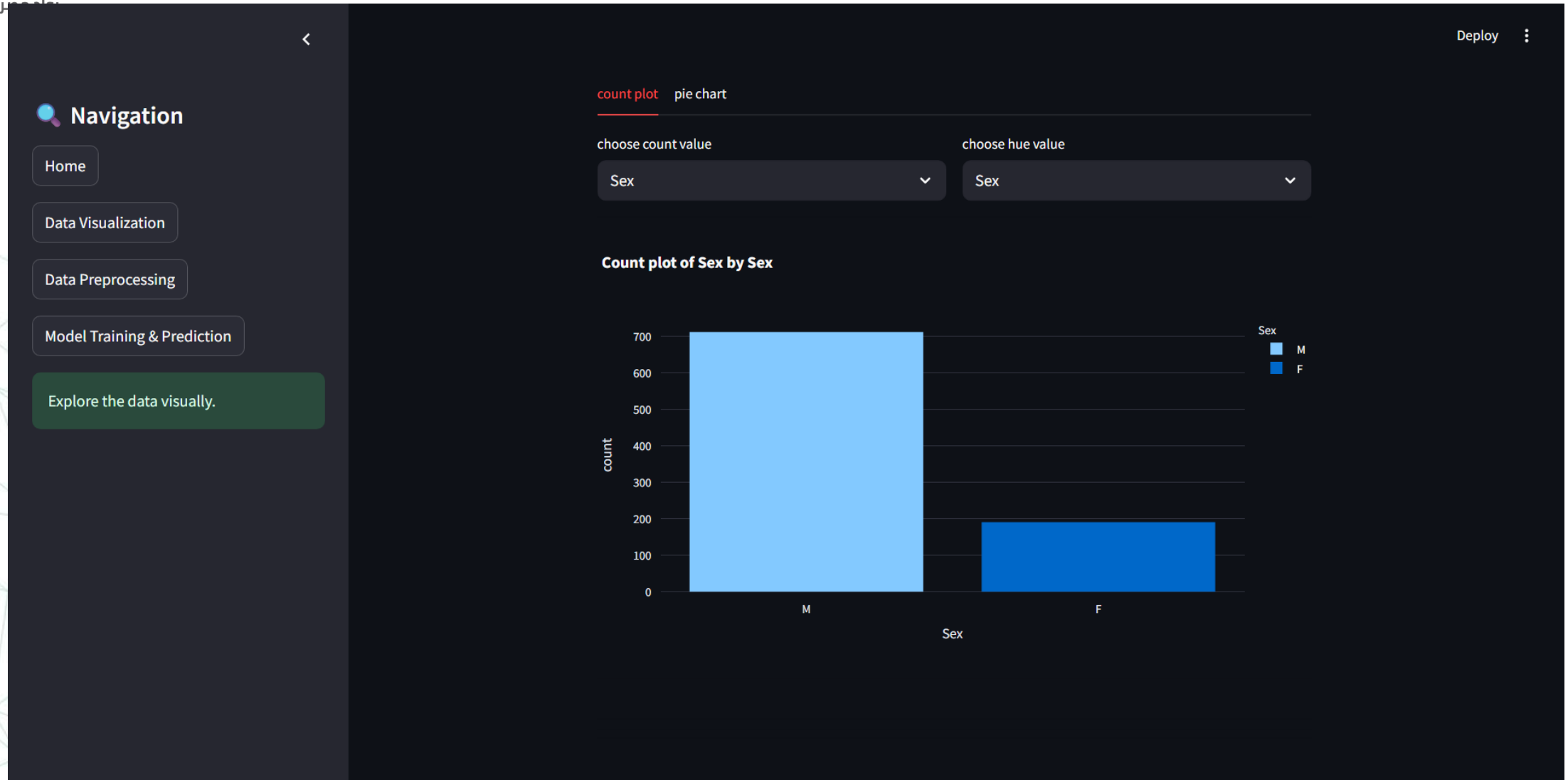
### 4. Modeling Page

1. Allows users to input medical data manually.
2. Choose from multiple models: Random Forest, Logistic Regression, XGBoost, Neural Network.
3. Click "Predict" to get real-time risk prediction and probability.
4. Shows evaluation metrics like accuracy, precision, recall.

# Project Wireframe



# Project Wireframe



# Project Wireframe

Navigation

Home

Data Visualization

Data Preprocessing

Model Training & Prediction

Data preprocessing completed.

Deploy

## Description of Numeric Columns

	count	mean	std	min	25%	50%	75%	max
Age	869	53.6686	10.7711	10	47	54	60	150
RestingBP	918	132.3965	18.5142	0	120	130	140	200
Cholesterol	839	201.6996	109.5172	0	177	223	268	700
FastingBS	918	0.2331	0.423	0	0	0	0	1
MaxHR	918	136.8094	25.4603	60	120	138	156	202
Oldpeak	918	0.8874	1.0666	-2.6	0	0.6	1.5	6.2
HeartDisease	918	0.5534	0.4974	0	0	1	1	1

## Description of Object Columns

	count	unique	top	freq
Sex	918	4	M	712
ChestPainType	918	4	ASY	496
RestingECG	918	3	Normal	552
ExerciseAngina	918	2	N	547
ST_Slope	918	3	Flat	460

5/15/25

3Youth<sup>®</sup>

6

# Project Wireframe

<

Navigation

Home

Data Visualization

Data Preprocessing

Model Training & Prediction

Data preprocessing completed.

Deploy

⋮

## Percentage of Missing Values per Column

	Missing Percentage %
Age	5.3377
Sex	0
ChestPainType	0
RestingBP	0
Cholesterol	8.6057
FastingBS	0
RestingECG	0
MaxHR	0
ExerciseAngina	0
Oldpeak	0

Handling Missing Values (Imputation) ▾

Handling Outliers ▾

Encoding Categorical Features (Label Encoding) ▾

Encoding Categorical Features (One-Hot Encoding) ▾

Correlation Heatmap ▾


Feature Scaling ▾

# Project Wireframe

Navigation

- Home
- Data Visualization
- Data Preprocessing
- Model Training & Prediction

Deploy



## Model Training

Select Model

Random Forest

Train Model

### Training Random Forest

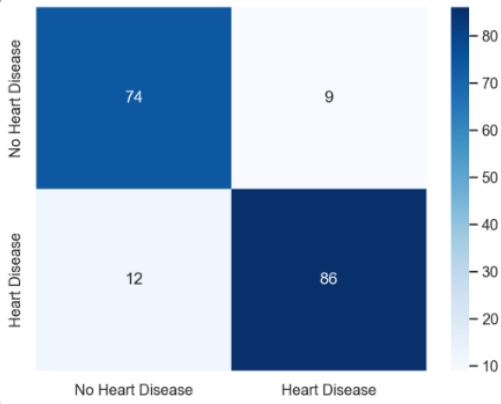
### Random Forest - Evaluation

Accuracy: 8.88

Classification Report:

	precision	recall	f1-score	support
0	0.86	0.89	0.88	83.00
1	0.91	0.88	0.89	98.00
accuracy	0.88	0.88	0.88	0.88
macro avg	0.88	0.88	0.88	181.00
weighted avg	0.88	0.88	0.88	181.00

Confusion Matrix:





## End Users + Features

### ➤ Target Users:

#### - Doctors & Medical Professionals

Use as a fast decision-support tool to identify patients at risk and prioritize further testing.

#### - Health-Conscious Individuals

Allows non-experts to check their risk level based on simple medical inputs, promoting early awareness.

#### - Healthcare Programs & Institutions

Useful for public health screening and can be integrated into digital awareness platforms.

### ➤ Key Features:

- Real-time heart disease prediction using clinical inputs.
- Multiple models supported: Random Forest, Logistic Regression, XGBoost, Neural Network.
- Clean and user-friendly interface built with Streamlit.
- Visual feedback (charts, probabilities) to enhance understanding.
- Fully deployed and accessible online.

### ➤ How It Helps:

- Speeds up early diagnosis and supports clinical decisions.
- Encourages individuals to monitor their health proactively.
- Scales well for public use and educational demonstrations

## ➤ Data Source:

A structured CSV file containing labeled medical records of patients. The dataset includes various clinical measurements and diagnosis labels.

## ➤ Data Type:

- Mixture of **numerical features**  
(e.g., Age, RestingBP, Cholesterol, MaxHR)
- And **categorical features**  
(e.g., Sex, ChestPainType, ExerciseAngina)

## ➤ Target Variable:

- **Heart Disease** (Binary Classification):
  - 0 : No heart disease
  - 1 : heart disease

## ➤ Preprocessing Steps:

### - Handling Categorical Variables:

- Applied **Label Encoding** for binary categories (e.g., Sex, ExerciseAngina)
- Used **One-Hot Encoding** for multi-class features (e.g., ChestPainType, RestingECG)

### - Data Cleaning:

- Removed or treated unrealistic values (e.g., 0 in blood pressure or cholesterol)
- Checked for duplicates and missing values

### - Feature Scaling:

- Standardized numerical features to improve model performance and convergence

# Programming Languages + Frameworks

## ➤ Language:

**Python**

## ➤ Frameworks & Libraries:

**Streamlit**

**Scikit-learn**

**TensorFlow/Keras**

**XGBoost**

## ➤ Tools & Hosting:

**Pandas**

**NumPy**

**Matplotlib**

**Seaborn**

**Streamlit Cloud**

**plotly**

# Live Application + Test

Section	Details
<b>Application Status</b>	<ul style="list-style-type: none"> <li>- Live and fully functional.</li> <li>- Publicly accessible via Streamlit Cloud</li> <li>- Four main pages: Home, Preprocessing, Visualizations, Modeling.</li> </ul>
<b>Deployment</b>	<ul style="list-style-type: none"> <li>- Hosted on Streamlit Cloud.</li> <li>- No installation required (browser-based).</li> <li>- GitHub integrated for version control and updates.</li> </ul>
<b>Testing – Unit</b>	<ul style="list-style-type: none"> <li>- Core functions tested individually (e.g., encoding, scaling, prediction logic).</li> <li>- Ensured each part behaves as expected.</li> </ul>
<b>Testing – Validation</b>	<ul style="list-style-type: none"> <li>- Train/Test split used to evaluate model.</li> <li>- Calculated accuracy, precision, recall on unseen data.</li> </ul>
<b>Testing – User (UX/UI)</b>	<ul style="list-style-type: none"> <li>- Interface tested with non-technical users.</li> <li>- Feedback applied to improve usability, layout, and navigation.</li> </ul>
<b>Outcome</b>	<ul style="list-style-type: none"> <li>- Reliable predictions with clear feedback.</li> <li>- Ready for educational use or real-time health risk checks.</li> </ul>

# Deliverables (Reports, etc.)

Deliverable	Description
Streamlit Web Application	Fully functional, interactive web app for real-time heart disease prediction.
Cleaned Dataset	Preprocessed and cleaned version of the original dataset, ready for analysis and modeling.
Source Code & Notebooks	All code files and Jupyter notebooks used in data processing, EDA, model training, and testing.
Documentation	Includes: <ul style="list-style-type: none"> <li>• Data Dictionary</li> <li>• Preprocessing steps</li> <li>• Model explanation</li> <li>• User Guide</li> </ul>
Deployment Link	Publicly accessible link to the deployed web app (hosted on Streamlit Cloud). <a href="https://heart-disease-predictor-lyz98qemywla9hy8fxz76y.streamlit.app/">https://heart-disease-predictor-lyz98qemywla9hy8fxz76y.streamlit.app/</a>

# Project Team + Roles

Name	Role	Key Responsibilities	Tools Used
Mostafa Gamal Abdel-Fatah Fouda	Data Scientist	Modeling, Deployment , Presentation	GitHub, Jupyter, Streamlit, Power point , Scikit-learn (sklearn), XGBoost , Python, TensorFlow/Keras, Streamlit Cloud
Mohamed Emad Ibrahim Mostafa	Data Scientist		
Lucas Alkomos Philopater Zakher Hanna	Data Scientist	EDA, Visualization, Documentations	Python, Seaborn, Matplotlib, Plotly, Jupyter, Microsoft Word, PowerBI
Youssef Alaa El Din Metwally Musallam	Data Scientist		
Maryam Gamal Ahmed Kamal Askar	Data Scientist	Data Cleaning, Data preprocessing, User testing	Python, Pandas, NumPy, Jupyter, Scikit-learn (sklearn), Seaborn
Marco Wael Issa Ibrahim	Data Scientist		

# Thank You

Thank you for your attention and interest in this project!  
Feel free to ask any questions or share your feedback.

## Contact:

Email: [mostafagamalf9@gmail.com](mailto:mostafagamalf9@gmail.com)

GitHub: <https://github.com/MostafaGmalFouda>

LinkedIn: [linkedin.com/in/mostafa-gamal-fouda-2645212a4](https://www.linkedin.com/in/mostafa-gamal-fouda-2645212a4)