# Heart Disease Prediction

**Name:** Shrusti Chintawar
**Roll No:** 09
**Course:** DMBI CA-2
**Date:** April 16, 2025

---

## 1. Abstract

This project applies a machine learning approach to predict the likelihood of heart disease using the Random Forest Classifier. With the UCI Heart Disease dataset, various features such as chest pain type, cholesterol levels, and maximum heart rate are analyzed. Exploratory Data Analysis (EDA) is performed, followed by model training and evaluation. Results are visualized with relevant graphs including a confusion matrix and feature importance chart.

---

## 2. Introduction

Heart disease remains a top health concern globally. Early prediction can reduce risk through early interventions. Machine learning offers automated and accurate prediction models. This project utilizes a Random Forest Classifier due to its robustness and interpretability.

---

## 3. Objective

- Predict heart disease presence using a machine learning model.

- Visualize feature influence on disease.

- Evaluate classifier performance with suitable metrics.

---

## 4. Dataset Description

**Source:** UCI Repository / Kaggle
**Features:**

- `age`, `sex`, `cp`, `trestbps`, `chol`, `fbs`, `restecg`, `thalach`, `exang`, `oldpeak`, `slope`, `ca`, `thal` **Target:** `target` (1 = heart disease, 0 = no disease)

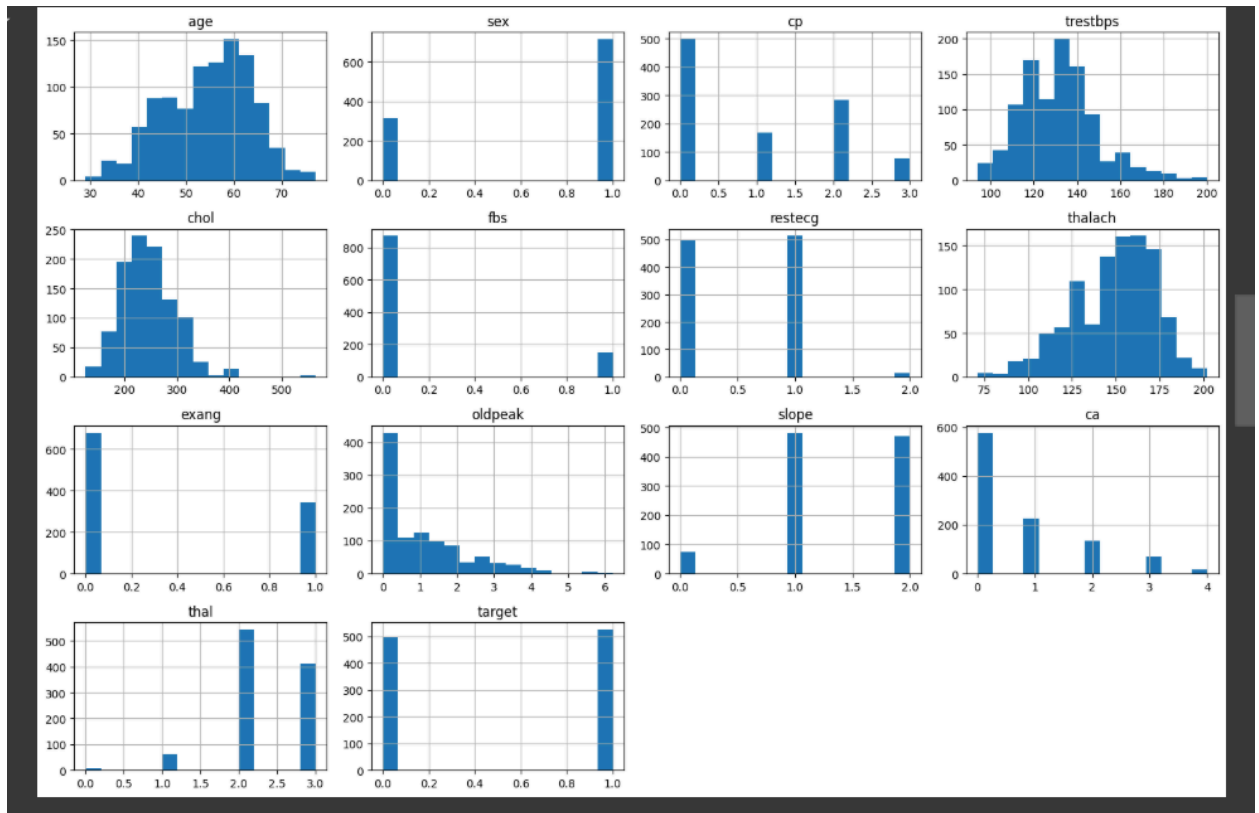---

## 5. Tools and Libraries Used

- Python (Colab)

- Libraries: `pandas`, `numpy`, `matplotlib`, `seaborn`, `sklearn`

---

## 6. Exploratory Data Analysis (EDA)
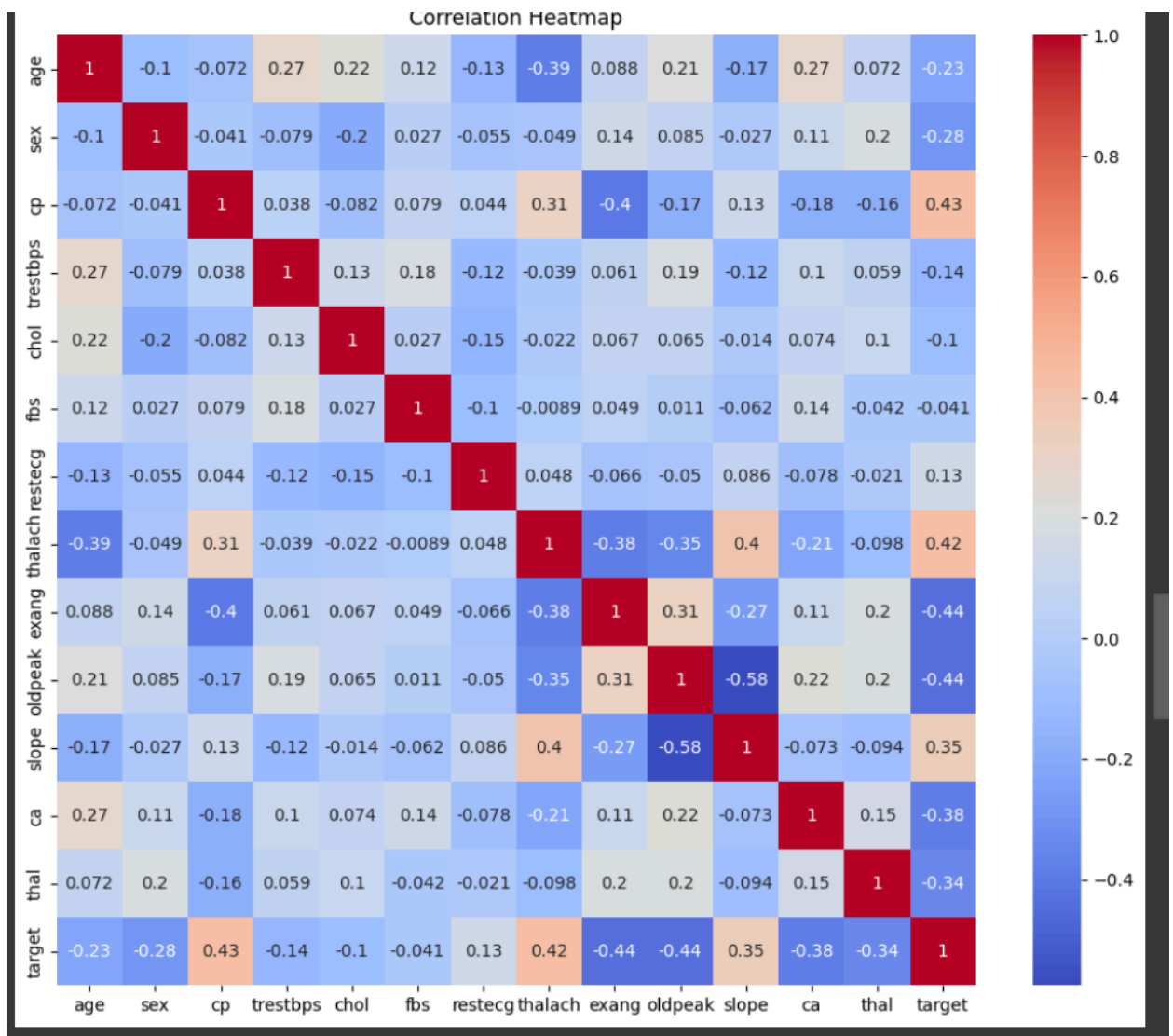
- **Missing Values:** None present.

Exploratory Data Analysis is a critical step in understanding the structure and patterns within the dataset. It helps identify trends, spot anomalies, and test assumptions using visual and quantitative methods. In this project, EDA is used to assess the distribution of variables, check for class balance, and understand feature correlations.

- **Histograms:**

df.hist(bins=15, figsize=(15, 10))

- **Correlation Heatmap:**

Correlation Heatmap

```
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')
```

---

## 7. Data Preprocessing

Data preprocessing involves preparing the raw data for machine learning algorithms. This includes splitting the data into features and target variables, scaling the features to bring them to a common range, and dividing the dataset into training and testing sets. These steps ensure that the model learns effectively and performs accurately on unseen data.

.

## 8. Model Implementation

Model implementation refers to the process of training a machine learning algorithm on the preprocessed dataset. In this case, the Random Forest Classifier is chosen for its ability to handle high-dimensional data and avoid overfitting. The model is trained on the training dataset and then used to predict outcomes on the test dataset

```
               precision    recall  f1-score   support

           0       0.97      1.00      0.99       102
           1       1.00      0.97      0.99       103

    accuracy                           0.99       205
   macro avg       0.99      0.99      0.99       205
weighted avg       0.99      0.99      0.99       205
```

---

## 9. Evaluation Metrics

- **Classification Report:**

print(classification_report(y_test, y_pred))

```
               precision    recall  f1-score   support

           0       0.97      1.00      0.99       102
           1       1.00      0.97      0.99       103

    accuracy                           0.99       205
   macro avg       0.99      0.99      0.99       205
weighted avg       0.99      0.99      0.99       205
```
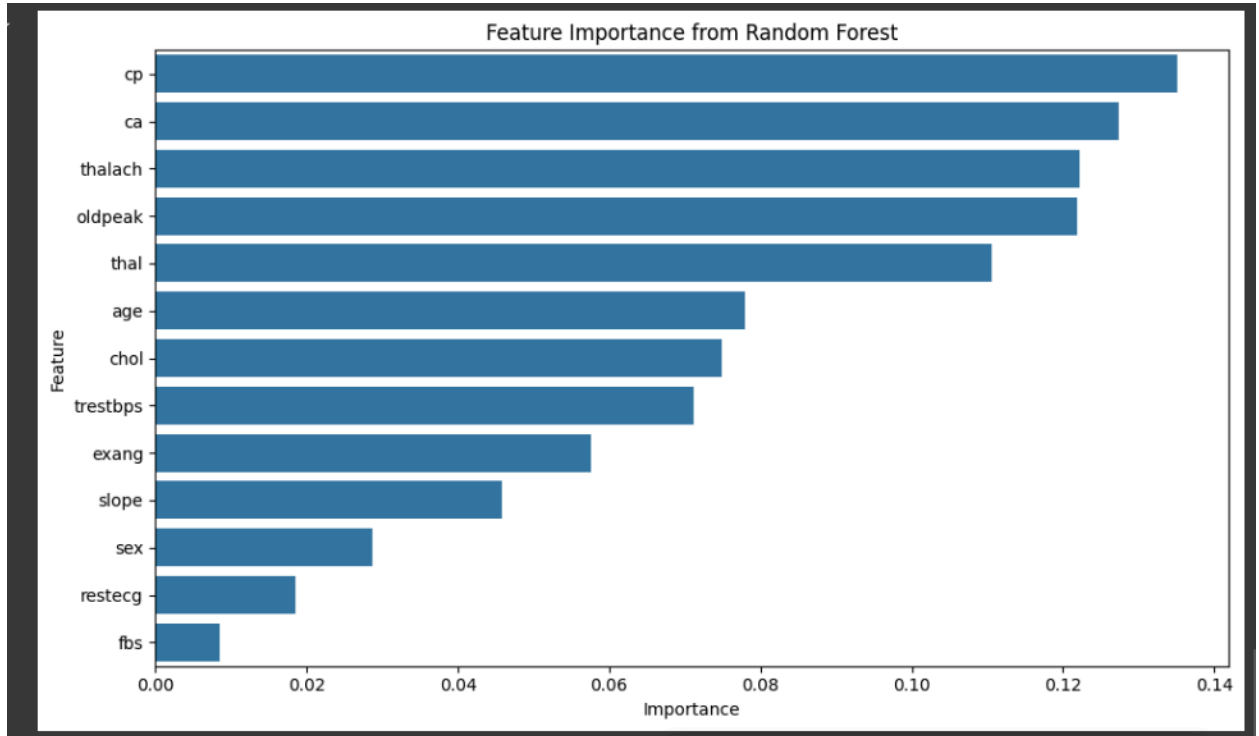
- **Feature Importance Plot:**


Feature Importance from Random Forest

---

## 10. Results and Insights

- Most significant features: `cp`, `thal`, `ca`

- High recall and precision indicate robust performance.

- Visualizations helped in understanding both data and model behavior.

---

## 11. Conclusion

The Random Forest classifier achieved effective and interpretable predictions for heart disease. With good accuracy and insightful feature importance, this model can support healthcare decisions.

---

## 12. Future Scope

- Use of grid search for hyperparameter tuning

- Ensemble comparisons (e.g., Gradient Boosting)

- Real-time integration with wearable data

---

## 13. References

- UCI Machine Learning Repository

- Kaggle: Heart Disease Dataset

- Scikit-learn Documentation