# Heart Disease Prediction Using Machine Learning - Report

## Project Overview:

This project involves building a machine learning model to predict heart disease based on clinical parameters using Python. The process includes data preprocessing, exploratory data analysis (EDA), feature scaling, model building, evaluation, and interpreting results.

## Step 1: Import Required Libraries

**pandas, numpy:** Data handling and numerical operations.
**matplotlib.pyplot, seaborn:** Visualization libraries for EDA.
**scikit-learn:** Machine learning models and evaluation tools.

Purpose: These libraries provide all the necessary tools for data manipulation, visualization, model training, and evaluation.

## Step 2: Load the Dataset

Command:

> *df = pd.read_csv(r"D:\Internship-DEN\HDP\heart-disease.csv")*

Result: Displays the first 5 rows of the dataset, confirming successful load.

Example Output:

```
  age sex cp trestbps chol fbs restecg thalach exang oldpeak slope ca thal target
0 63  1   3  145      233  1   0       150     0     2.3     0     0  1    1
...
```

Purpose: Ensures dataset is correctly loaded and structured.

## Step 3: Data Preprocessing

Check for missing values and duplicates:
df.isnull().sum()
df.duplicated().sum()

Result:
* No missing values.
* One duplicate row found and removed.

Purpose: Clean data is essential for accurate model performance.

## Step 4: Exploratory Data Analysis (EDA)

Generate a heatmap:
sns.heatmap(df.corr(), annot=True, cmap='coolwarm')

Purpose:
* Understand correlations between features.
* Identify which features strongly relate to the target (heart disease).

## Step 5: Feature Scaling

Standardization applied using:
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

Purpose: Normalizes feature values to improve model performance.

## Step 6: Split Data into Training and Testing Sets

Command:
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.2, random_state=42)

Purpose: Provides a way to train the model and then test its accuracy on unseen data.

## Step 7 & 8: Build Models

* Logistic Regression Model
* Random Forest Classifier Model

Purpose: Compare two models to select the one with better performance.

## Step 9 & 10: Model Evaluation Results

Logistic Regression:

* Accuracy: 82%
* Precision: 84%
* Recall: 81%
* F1 Score: 83%
* Confusion Matrix:
[[24  5]
 [ 6 26]]

Random Forest Classifier:

* Accuracy: 87%
* Precision: 90%
* Recall: 84%
* F1 Score: 87%

* Confusion Matrix:
[[26  3]
 [ 5 27]]

Purpose: Random Forest outperforms Logistic Regression in this case with higher overall metrics.

## Conclusion:

- The heart disease prediction model works effectively.
- Random Forest provides the best balance of accuracy, precision, recall, and F1 score.
- The project demonstrates the full machine learning pipeline, including data preprocessing, EDA, model training, evaluation, and interpretation of results.

## Recommendations:

- Further testing with cross-validation or additional hyperparameter tuning could enhance performance.
- Deployment options include using Streamlit or Flask for building a user interface.

## Result:

```
First 5 rows of the dataset:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
0   63    1   3       145   233    1        0      150      0      2.3      0   0     1       1
1   37    1   2       130   250    0        1      187      0      3.5      0   0     2       1
2   41    0   1       130   204    0        0      172      0      1.4      2   0     2       1
3   56    1   1       120   236    0        1      178      0      0.8      2   0     2       1
4   57    0   0       120   354    0        1      163      1      0.6      2   0     2       1

Checking for missing values:
age        0
sex        0
cp         0
trestbps   0
chol       0
fbs        0
restecg    0
thalach    0
exang      0
oldpeak    0
slope      0
ca         0
thal       0
target     0
dtype: int64

Checking for duplicates:
Number of duplicate rows: 1
```

Correlation Heatmap