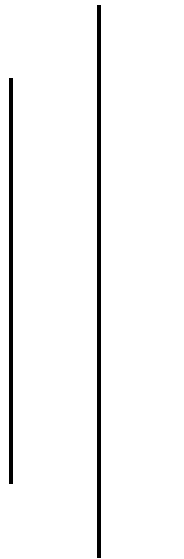




**A  
Lab Report  
On**

**Predicting Heart Disease Using Machine  
Learning**



**Submitted By:**

**Name:** Bivek Chaudhary

**Faculty:** BSc. CSIT

**Roll No:** 23871/076

**Batch:** 2076

**Year:** IV

**Semester:** VIII

**Submitted To:**

Department of CSIT

Orchid International College

---

This report tells discusses various Python-based ML and Data Science libraries in an attempt to build a Machine Learning model capable of predicting whether or not someone has heart disease based on their medical attributes. We're going to take the following approach:

1. Problem Definition
2. Understanding Features
3. Data Preparation and its tools
4. Exploratory Data Analysis
5. Modelling
6. Model Evaluation

## **Problem Definition**

Given clinical parameters about a patient, can we predict whether or not they have heart disease? Our aim is to reach the model accuracy of more than 85% . If the model scores better than 85%, we will select the model.

## **Understanding Features**

- a. Heart\_disease: It is the target variable indicating the presence of heart disease (object type) into the value of:  
0 = No  
1 = Yes
- b. BMI: It Stands for Body Mass Index, a measure of body fat based on height and weight (float) which is in numeric values.
- c. Smoking: It indicates if the individual smokes (object) or not categorized into:  
0=No  
1= Yes
- d. Alcohol\_drinking: It indicates if the individual consumes alcohol (object) or not categorized into:  
0=No  
1= Yes
- e. Stroke: It indicates if the individual has had a stroke (object) ) or not categorized into:  
0=No  
1= Yes
- f. Physical\_health: Number of days the individual experienced poor physical health in the past 30 days (float).

- g. Mental\_health: Number of days the individual experienced poor mental health in the past 30 days (float).
- h. Diff\_walking: It indicates if the individual has difficulty walking (object) whose format:
  - 0=No
  - 1= Yes
- i. Sex: Gender of the individual (object)
  - 0=Female
  - 1= Male
- j.
- k. Age\_category: It displays the age category of the individual (object) whose format:
  - 18-24': 21,
  - '25-29': 27,
  - '30-34': 32,
  - '35-39': 37,
  - '40-44': 42,
  - '45-49': 47,
  - '50-54': 52,
  - '55-59': 57,
  - '60-64': 62,
  - '65-69': 67,
  - '70-74': 72,
  - '75-79': 77,
  - '80 or older': 85
- l. Race: It displays the race of the individual (object), categorized into:
  - 'White': 1,
  - 'Black': 2,
  - 'Asian': 3,
  - 'American Indian/Alaskan Native': 4,
  - 'Other': 5,
  - 'Hispanic': 6
- m. Diabetic: It indicates if the individual has diabetes (object) whose format:
  - 0=No
  - 1= Yes

- n. Physical\_activity: Indicates if the individual engages in physical activity (object).
- o. General\_health: Self-assessed general health status (object).
- p. Sleep\_time: Average hours of sleep per night (float).
- q. Asthma: Indicates if the individual has asthma (object) whose format:
  - 0=No
  - 1= Yes
- r. Kidney\_disease: Indicates if the individual has kidney disease (object) whose format:
  - 0=No
  - 1= Yes
- s. Skin\_cancer: Indicates if the individual has skin cancer (object) whose format:
  - 0=No
  - 1= Yes

This dataset includes a mix of categorical and numerical features, capturing various health-related attributes and behaviors that can influence heart disease risk.

## **Data preparation and its tools**

- a. Pandas & NumPy: Pandas offers powerful data structures like Series and DataFrame for structured data manipulation, while NumPy provides support for numerical computing with multi-dimensional arrays and mathematical functions.
- b. Matplotlib & Seaborn: Matplotlib is a versatile library for creating static, interactive, and animated visualizations, while Seaborn simplifies the process of generating complex statistical graphics on top of Matplotlib.
- c. Scikit-Learn: Scikit-Learn is a comprehensive machine learning library with various algorithms for supervised and unsupervised learning tasks, along with utilities for model evaluation and selection.

```
In [1]: 1 # Importing necessary libraries
2 import pandas as pd
3 import numpy as np
4 import seaborn as sns
5 import matplotlib.pyplot as plt
6 from sklearn.model_selection import train_test_split
```

```
In [2]: 1 # Load the dataset
2 df = pd.read_csv('heart_disease.csv')
```

```
In [3]: 1 # Display the first few rows of the dataframe
2 df.head()
```

```
Out[3]:
```

	HeartDisease	BMI	Smoking	AlcoholDrinking	Stroke	PhysicalHealth	MentalHealth	DiffWalking	Sex	AgeCategory	Race	Diabetic	PhysicalActivity	Ge
0	No	16.60	Yes	No	No	3.0	30.0	No	Female	55-59	White	Yes	Yes	V
1	No	20.34	No	No	Yes	0.0	0.0	No	Female	80 or older	White	No	Yes	V
2	No	26.58	Yes	No	No	20.0	30.0	No	Male	65-69	White	Yes	Yes	V
3	No	24.21	No	No	No	0.0	0.0	No	Female	75-79	White	No	No	V
4	No	23.71	No	No	No	28.0	0.0	Yes	Female	40-44	White	No	Yes	V

### Remaining the columns

```
In [4]: 1 df.rename(columns = {'HeartDisease':'Heart_disease','AlcoholDrinking':'Alcohol_drinking','PhysicalHealth':'Physical_health',
2
```

```
In [5]: 1 # View of the Renamed Dataframe
2 df.head()
```

```
Out[5]:
```

	Heart_disease	BMI	Smoking	Alcohol_drinking	Stroke	Physical_health	Mental_health	Diff_walking	Sex	Age_category	Race	Diabetic	Physical_activi
0	No	16.60	Yes	No	No	3.0	30.0	No	Female	55-59	White	Yes	Yi
1	No	20.34	No	No	Yes	0.0	0.0	No	Female	80 or older	White	No	Yi
2	No	26.58	Yes	No	No	20.0	30.0	No	Male	65-69	White	Yes	Yi
3	No	24.21	No	No	No	0.0	0.0	No	Female	75-79	White	No	Yi
4	No	23.71	No	No	No	28.0	0.0	Yes	Female	40-44	White	No	Yi

```
In [6]: 1 # Replace "Very good" with "Very_good" in General_health
2 df['General_health'] = df['General_health'].replace('Very good', 'Very_good')
3
4 # Replace "No, borderline diabetes" with "No" in Diabetic
5 df['Diabetic'] = df['Diabetic'].replace('No, borderline diabetes', 'No')
6
7 # Replace "Yes (during pregnancy)" with "Yes" in Diabetic
8 df['Diabetic'] = df['Diabetic'].replace('Yes (during pregnancy)', 'Yes')
```

## Exploratory Data Analysis

```
In [7]: 1 # Exploratory Data Analysis (EDA)
        2 df.info()
        3 df.describe()
        4 df.isnull().sum()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 319795 entries, 0 to 319794
Data columns (total 18 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Heart_disease          319795 non-null object
1   BMI                    319789 non-null float64
2   Smoking                319721 non-null object
3   Alcohol_drinking       319680 non-null object
4   Stroke                  319774 non-null object
5   Physical_health        319778 non-null float64
6   Mental_health          319781 non-null float64
7   Diff_walking           319762 non-null object
8   Sex                    319482 non-null object
9   Age_category           319751 non-null object
10  Race                    319756 non-null object
11  Diabetic                319762 non-null object
12  Physical_activity       319779 non-null object
13  General_health          319302 non-null object
14  Sleep_time              319576 non-null float64
15  Asthma                  319795 non-null object
16  Kidney_disease          319778 non-null object
17  Skin_cancer             319795 non-null object
dtypes: float64(4), object(14)
memory usage: 43.9+ MB
```

```
Out[7]: Heart_disease          0
        BMI                    6
        Smoking                74
        Alcohol_drinking       115
        Stroke                  21
        Physical_health        17
        Mental_health          14
        Diff_walking           33
        Sex                    313
        Age_category           44
        Race                    39
        Diabetic                33
        Physical_activity       16
        General_health          493
        Sleep_time              219
        Asthma                  0
        Kidney_disease          17
        Skin_cancer             0
dtype: int64
```

## Handling the missing data using SimpleImputer

```
In [8]: 1 # Numerical columns to impute
2 from sklearn.impute import SimpleImputer
3 numerical_cols = ['BMI', 'Physical_health', 'Mental_health', 'Sleep_time']
4 imputer = SimpleImputer(strategy='median')
5 df[numerical_cols] = imputer.fit_transform(df[numerical_cols])
```

```
In [9]: 1 # Categorical columns to impute
2 categorical_cols = ['Smoking', 'Alcohol_drinking', 'Stroke', 'Diff_walking',
3                   'Sex', 'Age_category', 'Race', 'Diabetic', 'Physical_activity',
4                   'General_health', 'Kidney_disease']
5 imputer = SimpleImputer(strategy='most_frequent')
6 df[categorical_cols] = imputer.fit_transform(df[categorical_cols])
```

```
In [10]: 1 # Confirm all missing values are handled
2 print(df.isnull().sum())
```

```
Heart_disease      0
BMI                0
Smoking            0
Alcohol_drinking   0
Stroke             0
Physical_health     0
Mental_health      0
Diff_walking       0
Sex               0
Age_category       0
Race              0
Diabetic           0
Physical_activity   0
General_health     0
Sleep_time         0
Asthma             0
Kidney_disease     0
Skin_cancer        0
dtype: int64
```

```
In [11]: 1 df.head()
```

```
Out[11]:
```

	Heart_disease	BMI	Smoking	Alcohol_drinking	Stroke	Physical_health	Mental_health	Diff_walking	Sex	Age_category	Race	Diabetic	Physical_acti
0	No	16.60	Yes	No	No	3.0	30.0	No	Female	55-59	White	Yes	Y
1	No	20.34	No	No	Yes	0.0	0.0	No	Female	80 or older	White	No	Y
2	No	26.58	Yes	No	No	20.0	30.0	No	Male	65-69	White	Yes	Y
3	No	24.21	No	No	No	0.0	0.0	No	Female	75-79	White	No	Y
4	No	23.71	No	No	No	28.0	0.0	Yes	Female	40-44	White	No	Y

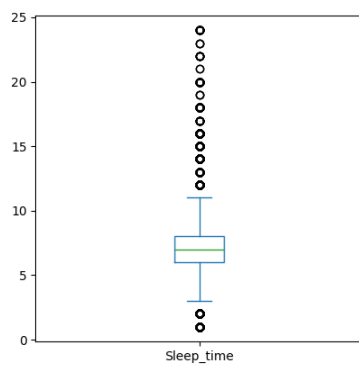
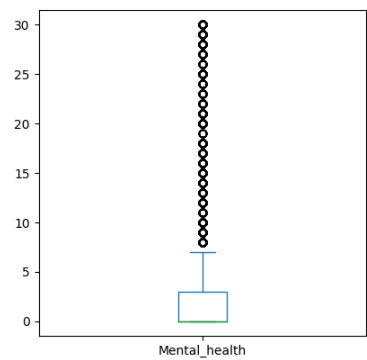
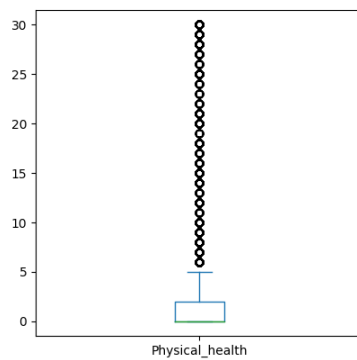
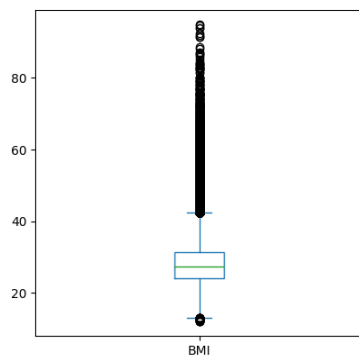
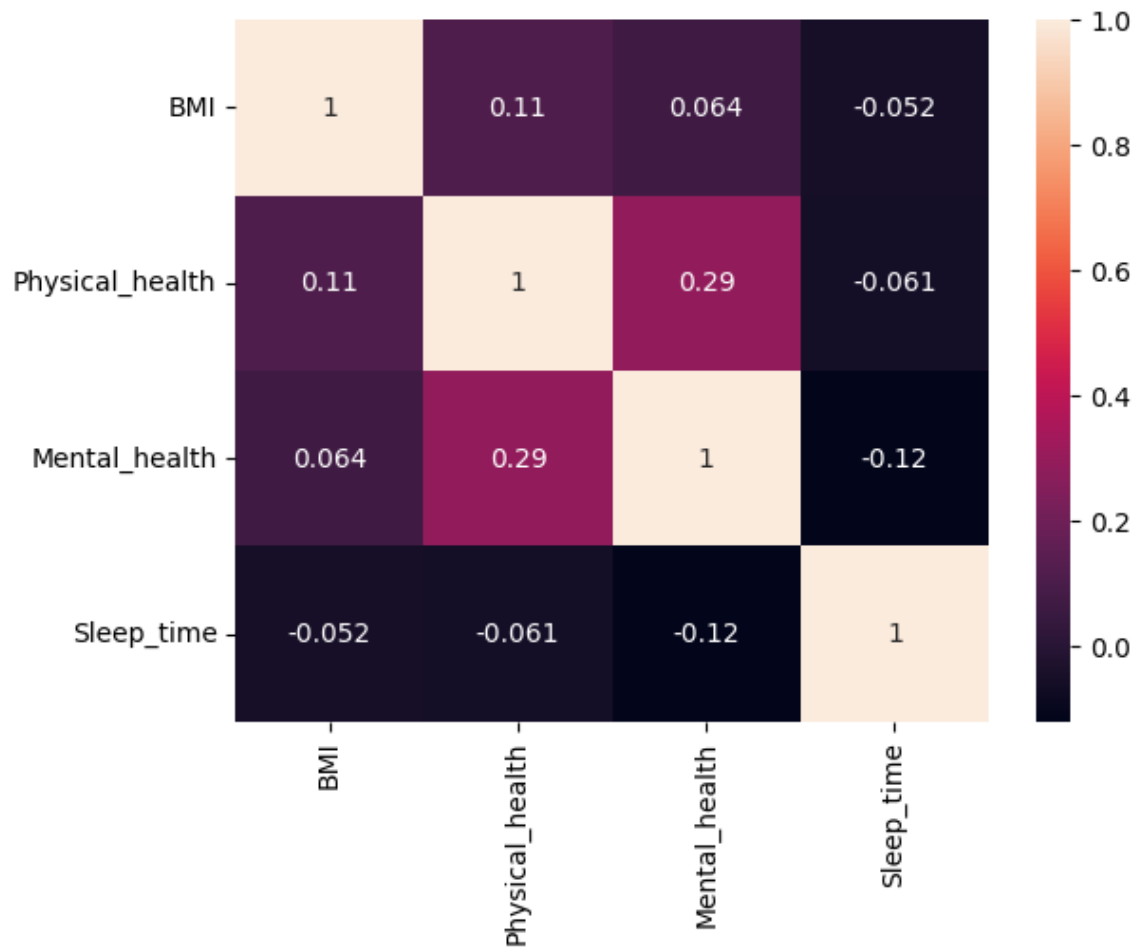
## Correlation matrix & Matrix Visualisation

```
In [12]: 1 df.corr(numeric_only=True)
```

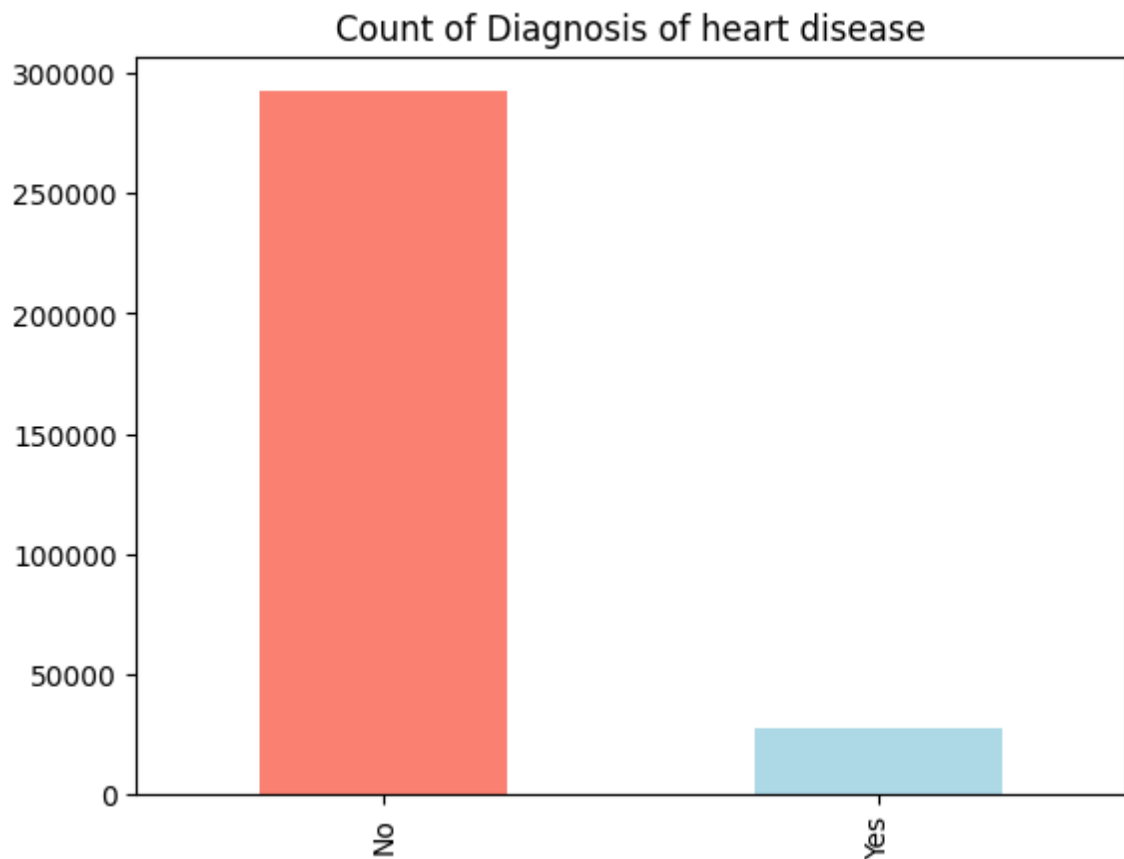
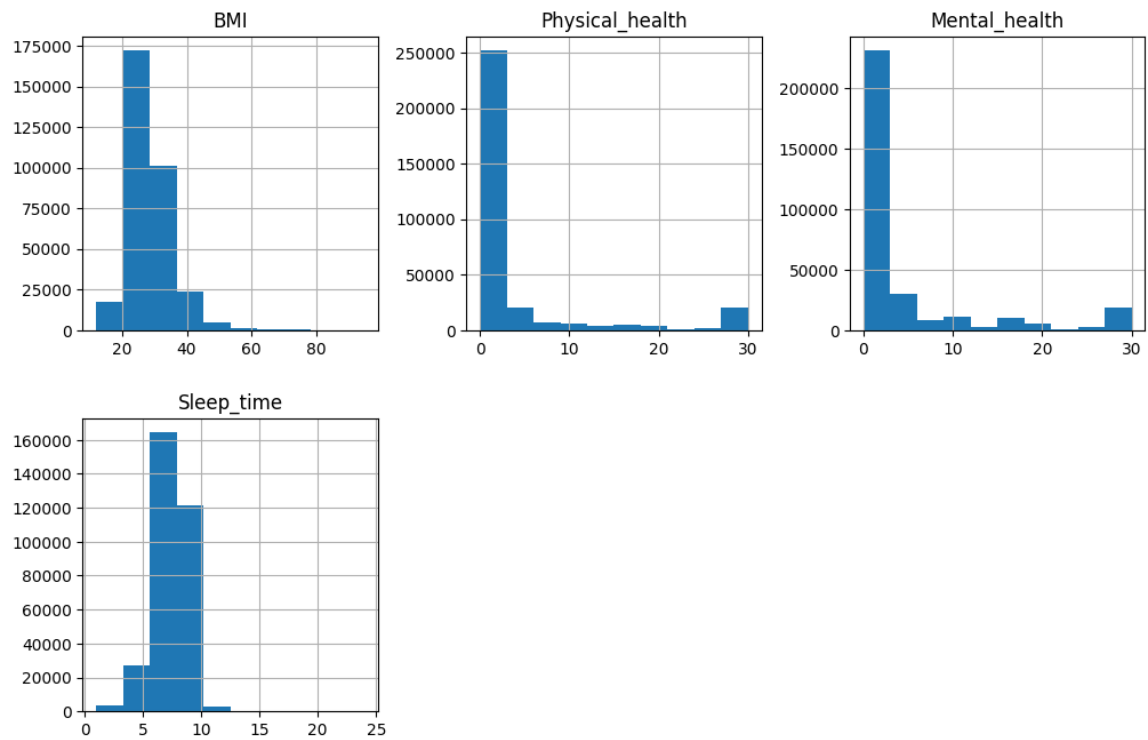
```
Out[12]:
```

	BMI	Physical_health	Mental_health	Sleep_time
BMI	1.000000	0.109787	0.064127	-0.051801
Physical_health	0.109787	1.000000	0.287948	-0.061384
Mental_health	0.064127	0.287948	1.000000	-0.119692
Sleep_time	-0.051801	-0.061384	-0.119692	1.000000

```
In [13]: 1 # Visualizing the correlation matrix
2 # Create the correlation heatmap
3 sns.heatmap(df.corr(numeric_only=True), annot = True);
```







**Interpretation:**

*The above plot shows the count of population having disease and not having disease where the count of having disease is low and the count of not having disease is high*

## Heart Disease vs Sex

```
In [17]: 1 # Male and Female
        2 df['Sex'].value_counts()
```

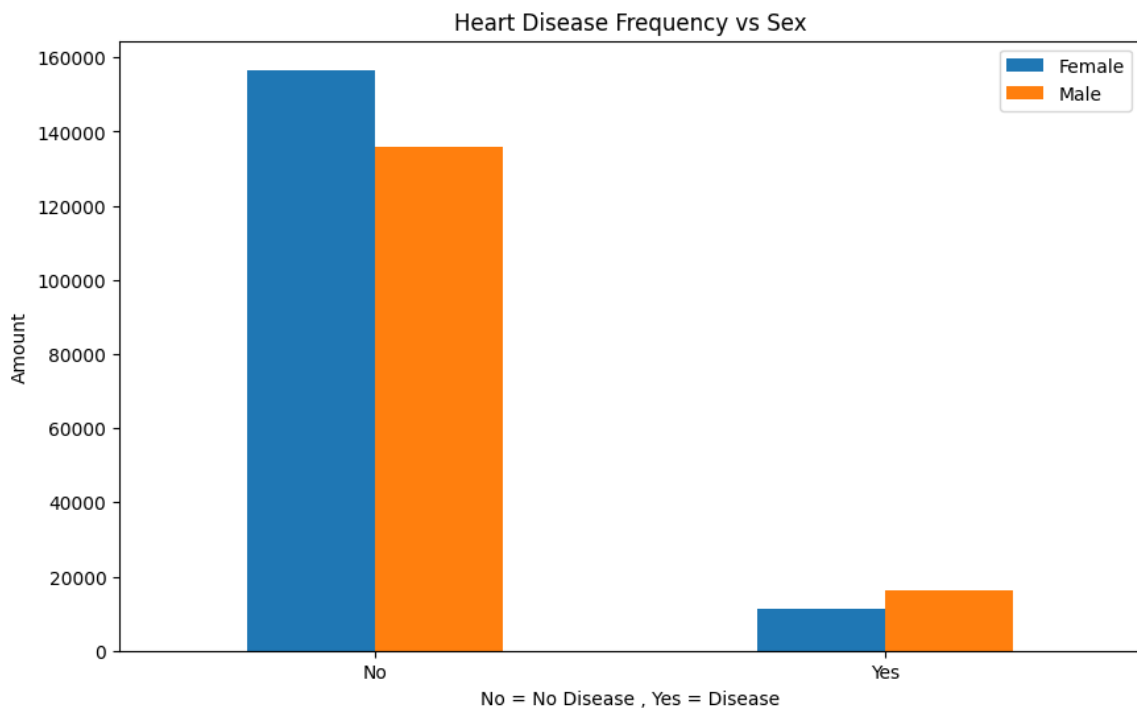
```
Out[17]: Female    167904
        Male      151891
        Name: Sex, dtype: int64
```

```
In [18]: 1 #Compare Target column i.e. Heart Disease with Sex Column
        2 pd.crosstab(df['Heart_disease'],df['Sex'])
```

```
Out[18]:
```

	Sex	Female	Male
Heart_disease			
	No	156648	135774
	Yes	11256	16117

```
In [19]: 1 #Create a plot for crosstab
        2 pd.crosstab(df['Heart_disease'],df['Sex']).plot(kind="bar",figsize=(10,6));
        3 plt.title("Heart Disease Frequency vs Sex")
        4 plt.xlabel("No = No Disease , Yes = Disease")
        5 plt.ylabel("Amount")
        6 plt.legend(["Female","Male"])
        7 plt.xticks(rotation=0)
```

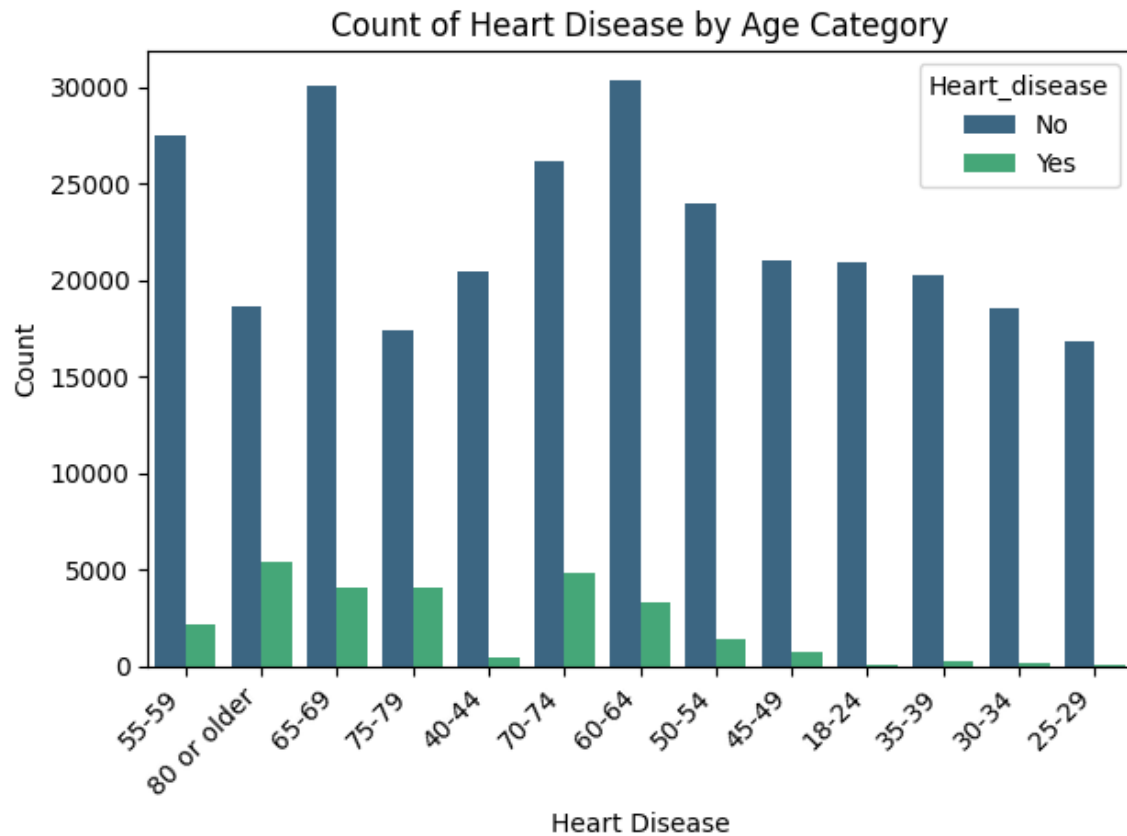


### Interpretation:

*The above plot shows the ratio of population having disease and not having disease wrt sex*

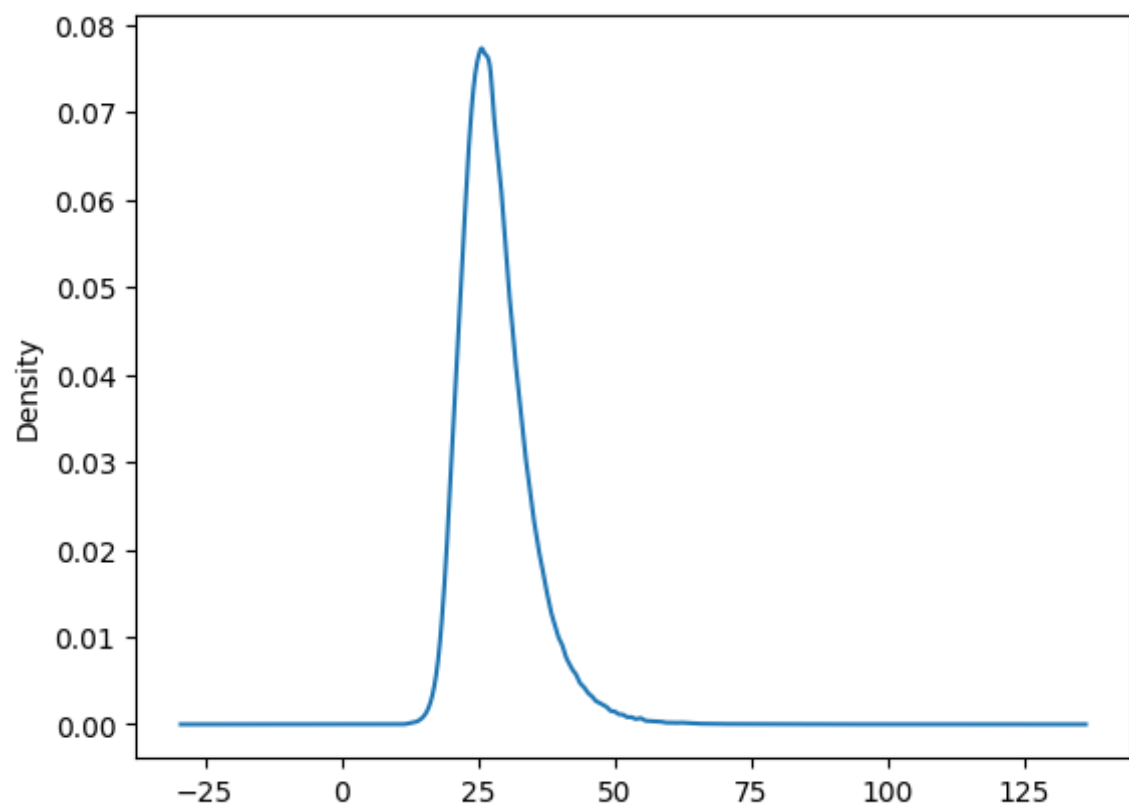
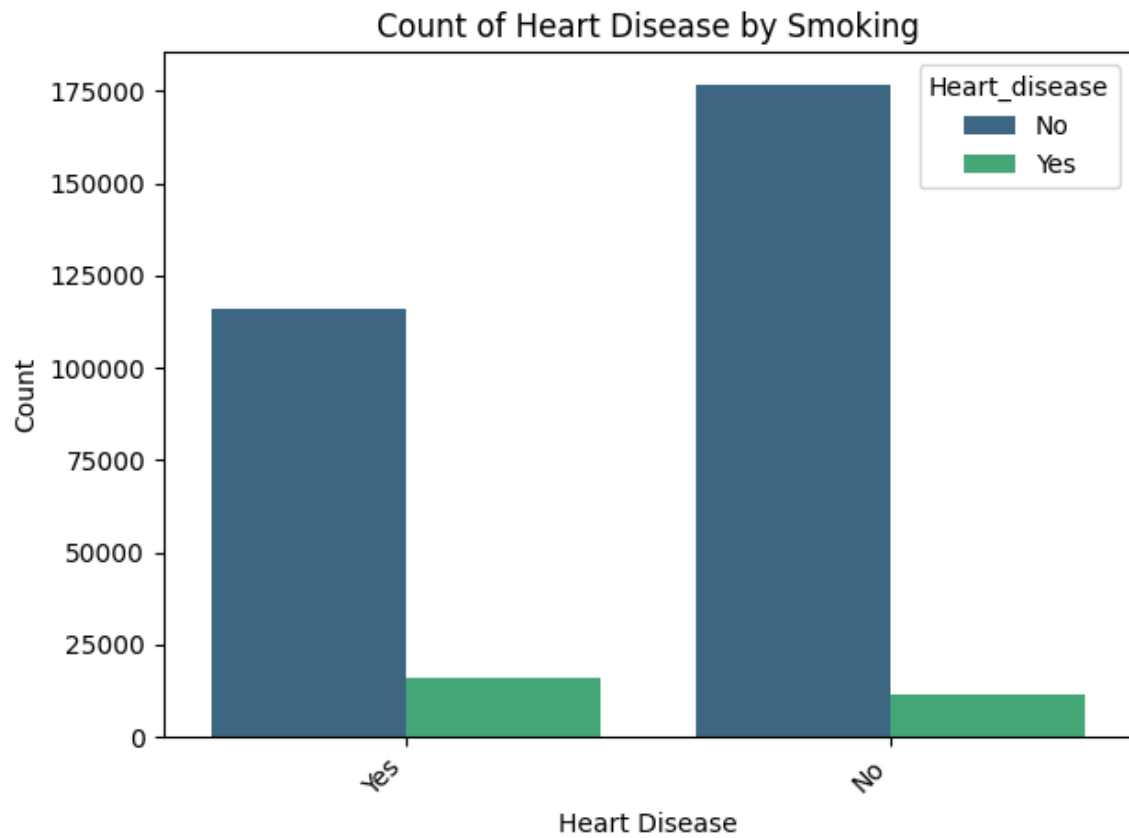
## Heart Disease vs Age

```
In [20]: 1 sns.countplot(x='Age_category', hue='Heart_disease', palette='viridis', data=df)
2
3 # Display the plot
4 plt.xlabel('Heart Disease')
5 plt.ylabel('Count')
6 plt.title('Count of Heart Disease by Age Category')
7 plt.xticks(rotation=45, ha='right')
8 plt.tight_layout()
9 plt.show()
```



## Count plot of Heart Disease vs Smoking

```
In [21]: 1 # Create a plot for countplot of Heart_disease vs Smoking
2 sns.countplot(x='Smoking', hue='Heart_disease', palette='viridis', data=df)
3
4 # Display the plot
5 plt.xlabel('Heart Disease')
6 plt.ylabel('Count')
7 plt.title('Count of Heart Disease by Smoking')
8 plt.xticks(rotation=45, ha='right')
9 plt.tight_layout()
10 plt.show()
```



**Interpretation:**  
The above plot shows the density of BMI and the maximum density is around 23-26

```
In [23]: 1 # Define the mapping for age categories to continuous values
2 age_mapping = {
3     '18-24': 21,
4     '25-29': 27,
5     '30-34': 32,
6     '35-39': 37,
7     '40-44': 42,
8     '45-49': 47,
9     '50-54': 52,
10    '55-59': 57,
11    '60-64': 62,
12    '65-69': 67,
13    '70-74': 72,
14    '75-79': 77,
15    '80 or older': 85
16 }
17
18 # Apply the mapping to the Age_category column
19 df['Age_category'] = df['Age_category'].map(age_mapping)|
```

```
In [24]: 1 print('Heart_disease:', df['Heart_disease'].nunique(),df['Heart_disease'].unique())
2 print('Smoking:', df['Smoking'].nunique(),df['Smoking'].unique())
3 print('Alcohol_drinking:',df['Alcohol_drinking'].nunique(),df['Alcohol_drinking'].unique())
4 print('Stroke:',df['Stroke'].nunique(),df['Stroke'].unique())
5 print('Diff_waliking:',df['Diff_waliking'].nunique(),df['Diff_waliking'].unique())
6 print('Race:',df['Race'].nunique(),df['Race'].unique())
7 print('Diabetic:',df['Diabetic'].nunique(),df['Diabetic'].unique())
8 print('Physical_health:',df['Physical_health'].nunique(),df['Physical_health'].unique())
9 print('General_health:',df['General_health'].nunique(),df['General_health'].unique())
10 print('Asthma:',df['Asthma'].nunique(),df['Asthma'].unique())
11 print('Kidney_disease:',df['Kidney_disease'].nunique(),df['Kidney_disease'].unique())
12 print('Skin_cancer:',df['Skin_cancer'].nunique(),df['Skin_cancer'].unique())
```

```
Heart_disease: 2 ['No' 'Yes']
Smoking: 2 ['Yes' 'No']
Alcohol_drinking: 2 ['No' 'Yes']
Stroke: 2 ['No' 'Yes']
Diff_waliking: 2 ['No' 'Yes']
Race: 6 ['White' 'Black' 'Asian' 'American Indian/Alaskan Native' 'Other'
'Hispanic']
Diabetic: 2 ['Yes' 'No']
Physical_health: 31 [ 3.  0. 20. 28.  6. 15.  5. 30.  7.  1.  2. 21.  4. 10. 14. 18.  8. 25.
16. 29. 27. 17. 24. 12. 23. 26. 22. 19.  9. 13. 11.]
General_health: 5 ['Very_good' 'Fair' 'Good' 'Poor' 'Excellent']
Asthma: 2 ['Yes' 'No']
Kidney_disease: 2 ['No' 'Yes']
Skin_cancer: 2 ['Yes' 'No']
```

```

In [25]: 1 # Convert categorical columns to numerical values
2 df['Heart_disease'] = df['Heart_disease'].map({'No': 0, 'Yes': 1})
3 df['Smoking'] = df['Smoking'].map({'No': 0, 'Yes': 1})
4 df['Alcohol_drinking'] = df['Alcohol_drinking'].map({'No': 0, 'Yes': 1})
5 df['Stroke'] = df['Stroke'].map({'No': 0, 'Yes': 1})
6 df['Diff_walking'] = df['Diff_walking'].map({'No': 0, 'Yes': 1})
7 df['Sex'] = df['Sex'].map({'Male': 0, 'Female': 1})
8 race_map = {
9     'White': 1,
10    'Black': 2,
11    'Asian': 3,
12    'American Indian/Alaskan Native': 4,
13    'Other': 5,
14    'Hispanic': 6
15 }
16
17 df['Race'] = df['Race'].map(race_map)
18 df['Asthma'] = df['Asthma'].map({'No': 0, 'Yes': 1})
19 df['Kidney_disease'] = df['Kidney_disease'].map({'No': 0, 'Yes': 1})
20 df['Skin_cancer'] = df['Skin_cancer'].map({'No': 0, 'Yes': 1})
21
22 # Map General_health to numerical values
23 health_map = {
24     'Excellent': 1,
25     'Very_good': 2,
26     'Good': 3,
27     'Fair': 4,
28     'Poor': 5
29 }
30 df['General_health'] = df['General_health'].map(health_map)
31
32 # Map Diabetic to numerical values
33 df['Diabetic'] = df['Diabetic'].map({'No': 0, 'Yes': 1})
34 df['Physical_activity'] = df['Physical_activity'].map({'No': 0, 'Yes': 1})
35
36 # Display the unique values after mapping
37 print('Heart_disease:', df['Heart_disease'].nunique(), df['Heart_disease'].unique())
38 print('Smoking:', df['Smoking'].nunique(), df['Smoking'].unique())
39 print('Alcohol_drinking:', df['Alcohol_drinking'].nunique(), df['Alcohol_drinking'].unique())
40 print('Stroke:', df['Stroke'].nunique(), df['Stroke'].unique())
41 print('Diff_walking:', df['Diff_walking'].nunique(), df['Diff_walking'].unique())
42 print('Race:', df['Race'].nunique(), df['Race'].unique())
43 print('Diabetic:', df['Diabetic'].nunique(), df['Diabetic'].unique())
44 print('Physical_health:', df['Physical_health'].nunique(), df['Physical_health'].unique())
45 print('General_health:', df['General_health'].nunique(), df['General_health'].unique())
46 print('Asthma:', df['Asthma'].nunique(), df['Asthma'].unique())
47 print('Kidney_disease:', df['Kidney_disease'].nunique(), df['Kidney_disease'].unique())
48 print('Skin_cancer:', df['Skin_cancer'].nunique(), df['Skin_cancer'].unique())
49

```

```

Heart_disease: 2 [0 1]
Smoking: 2 [1 0]
Alcohol_drinking: 2 [0 1]
Stroke: 2 [0 1]
Diff_walking: 2 [0 1]
Race: 6 [1 2 3 4 5 6]
Diabetic: 2 [1 0]
Physical_health: 31 [ 3.  0. 20. 28.  6. 15.  5. 30.  7.  1.  2. 21.  4. 10. 14. 18.  8. 25.
 16. 29. 27. 17. 24. 12. 23. 26. 22. 19.  9. 13. 11.]
General_health: 5 [2 4 3 5 1]
Asthma: 2 [1 0]
Kidney_disease: 2 [0 1]
Skin_cancer: 2 [1 0]

```

```
In [26]: 1 df['Heart_disease'].value_counts()
```

```

Out[26]: 0    292422
         1     27373
         Name: Heart_disease, dtype: int64

```

```

In [27]: 1 #Representation of Target variable in Percentage
2
3 countNoDisease = len(df[df.Heart_disease == 0])
4 countHaveDisease = len(df[df.Heart_disease == 1])
5 print("Percentage of Patients Haven't Heart Disease: {:.2f}%".format((countNoDisease / (len(df.Heart_disease))*100)))
6 print("Percentage of Patients Have Heart Disease: {:.2f}%".format((countHaveDisease / (len(df.Heart_disease))*100)))

```

```

Percentage of Patients Haven't Heart Disease: 91.44%
Percentage of Patients Have Heart Disease: 8.56%

```

## Modelling

We will experiment with the models, trying the model and getting the results from them.

Since in the target variable, Percentage of Patients Haven't Heart Disease: 91.44%  
Percentage of Patients Have Heart Disease: 8.56% so to imbalance the ratio we use the SMOTETomek which helps in balancing the data by importing Synthetic Minority Oversampling Technique by doing both downsampling and upsampling of class.

```
In [28]: 1 # Splitting the X and y values we called X as independent variables and y as target variable
2 X = df.drop('Heart_disease',axis=1)
3 y = df['Heart_disease']

In [30]: 1 #balancing the data by importing Synthetic Minority Oversampling Technique(SMOTE)
2 from imblearn.combine import SMOTETomek
3 #smotetomek does both downsampling and upsampling of class

In [31]: 1 # Instantiate SMOTETomek with the desired sampling strategy (float between 0 and 1)
2 # As specified as 0.9 means minprity has 90% of majority
3 smote_tomek = SMOTETomek(sampling_strategy=0.9)

In [32]: 1 # Fit and transform data
2 X_resampled, y_resampled = smote_tomek.fit_resample(X, y)

In [33]: 1 df.head()

Out[33]:
```

	Heart_disease	BMI	Smoking	Alcohol_drinking	Stroke	Physical_health	Mental_health	Diff_walking	Sex	Age_category	Race	Diabetic	Physical_activity
0	0	16.60	1	0	0	3.0	30.0	0	1	57	1	1	1
1	0	20.34	0	0	1	0.0	0.0	0	1	85	1	0	1
2	0	26.58	1	0	0	20.0	30.0	0	0	67	1	1	1
3	0	24.21	0	0	0	0.0	0.0	0	1	77	1	0	0
4	0	23.71	0	0	0	28.0	0.0	1	1	42	1	0	1

```


In [34]: 1 # Checking the number of majority and minority values
2 from collections import Counter
3 Counter(y_resampled)

Out[34]: Counter({0: 289909, 1: 260666})

In [35]: 1 # Splitting the dataset into train and test data
2 from sklearn.model_selection import train_test_split
3 X_train, X_test, y_train, y_test = train_test_split(X_resampled, y_resampled, test_size = 0.2, random_state = 0)
4
5 # Display the shapes of the training and testing sets
6 print(f"\nX_train shape: {X_train.shape}, y_train shape: {y_train.shape}")
7 print(f"\nX_test shape: {X_test.shape}, y_test shape: {y_test.shape}")
8
9 # Save training and testing sets to CSV files
10 X_train.to_csv('train.csv', index=False)
11 X_test.to_csv('test.csv', index=False)
12 y_train.to_csv('train_data_labels.csv', index=False, header=True)
13 y_test.to_csv('test_data_labels.csv', index=False, header=True)
```

Now we have got our data split into training and test sets, it is time to build a Machine Learning model. We will train it (find the patterns) on the training set. And we will test it (use the patterns) on the test set.

### Building Model

```
In [37]: 1 from sklearn.ensemble import RandomForestClassifier
2 from xgboost import XGBClassifier
3 from sklearn.linear_model import LogisticRegression
4 from sklearn.neighbors import KNeighborsClassifier
5 from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, recall_score, f1_score, roc_auc_score, classif

In [38]: 1 # Random Forest
2 print("Random Forest")
3 rf_model = RandomForestClassifier(random_state=42)
4 rf_model.fit(X_train, y_train)

Random Forest

Out[38]: RandomForestClassifier(random_state=42)
```

## Model Evaluation

```
In [52]: 1 # Function to evaluate and print the performance of a model
2 def evaluate_model(model,X_train, y_train, X_test, y_test):
3     y_pred = model.predict(X_test)
4     y_pred_proba = model.predict_proba(X_test)[: , 1]
5     # Predictions on training set for ROC curve
6     y_pred_train_proba = model.predict_proba(X_train)[: , 1]
7
8     accuracy = accuracy_score(y_test, y_pred)
9     precision = precision_score(y_test, y_pred)
10    recall = recall_score(y_test, y_pred)
11    f1 = f1_score(y_test, y_pred)
12    roc_auc = roc_auc_score(y_test, y_pred)
13    cm = confusion_matrix(y_test, y_pred)
14
15    print(f'Accuracy: {accuracy:.4f}')
16    print(f'Precision: {precision:.4f}')
17    print(f'Recall: {recall:.4f}')
18    print(f'F1 Score: {f1:.4f}')
19    print(f'ROC AUC Score: {roc_auc:.4f}')
20
21
22    print("\nClassification Report:")
23    print(classification_report(y_test, y_pred))
24
25    # Plot the confusion matrix using seaborn heatmap
26    sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", cbar=False,
27               xticklabels=["Predicted Positive", "Predicted Negative"],
28               yticklabels=["Actual Positive", "Actual Negative"])
29    plt.xlabel("Predicted")
30    plt.ylabel("Actual")
31    plt.title("Confusion Matrix")
32    plt.show()
33
34    # ROC Curve for Training Dataset
35    fpr_train, tpr_train, _ = roc_curve(y_train, y_pred_train_proba)
36    roc_auc_train = auc(fpr_train, tpr_train)
37
38    plt.figure()
39    lw = 2
40    plt.plot(fpr_train, tpr_train, color='darkorange',
41            lw=lw, label='ROC curve (area = %0.2f)' % roc_auc_train)
42    plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
43    plt.xlim([0.0, 1.0])
44    plt.ylim([0.0, 1.05])
45    plt.xlabel('False Positive Rate')
46    plt.ylabel('True Positive Rate')
47    plt.title('Receiver Operating Characteristic Curve (Training Dataset)')
48    plt.legend(loc="lower right")
49    plt.show()
50
51    # ROC Curve for Testing Dataset
52    fpr, tpr, _ = roc_curve(y_test, y_pred_proba)
53    roc_auc = auc(fpr, tpr)
54
55    plt.figure()
56    lw = 2
57    plt.plot(fpr, tpr, color='darkorange',
58            lw=lw, label='ROC curve (area = %0.2f)' % roc_auc)
59    plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
60    plt.xlim([0.0, 1.0])
61    plt.ylim([0.0, 1.05])
62    plt.xlabel('False Positive Rate')
63    plt.ylabel('True Positive Rate')
64    plt.title('Receiver Operating Characteristic Curve')
65    plt.legend(loc="lower right")
66    plt.show()
```

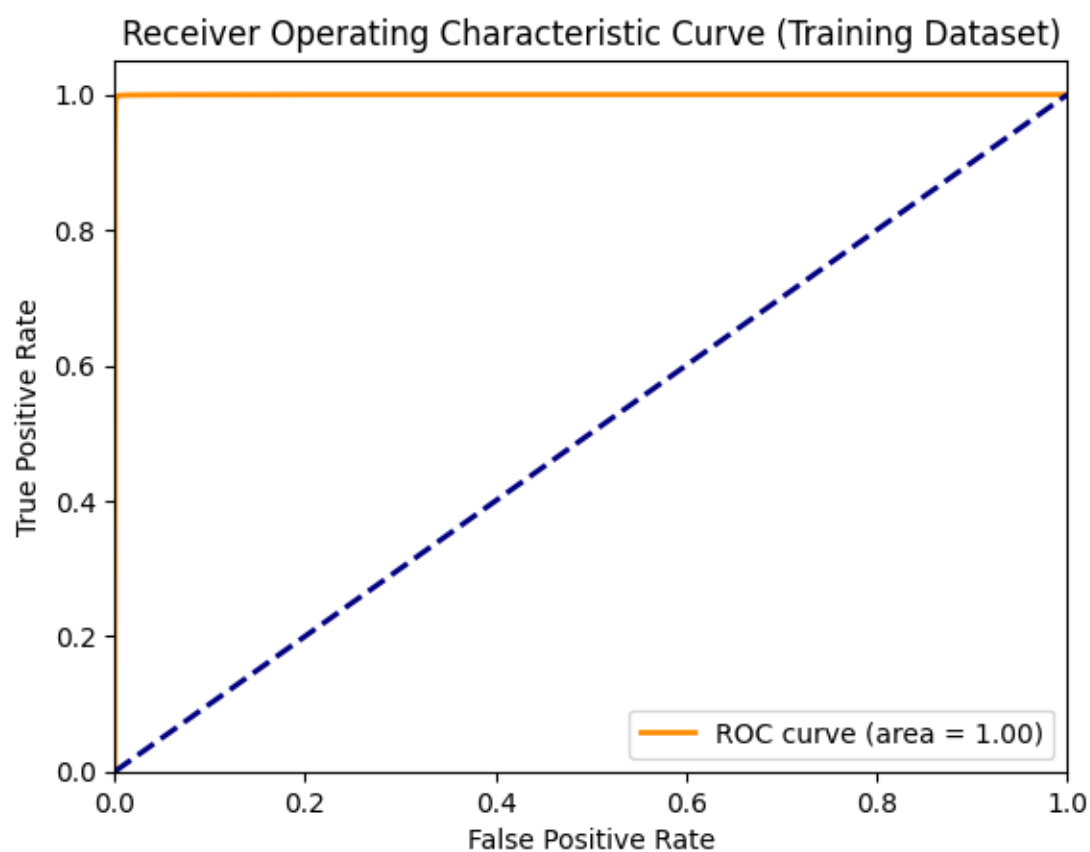
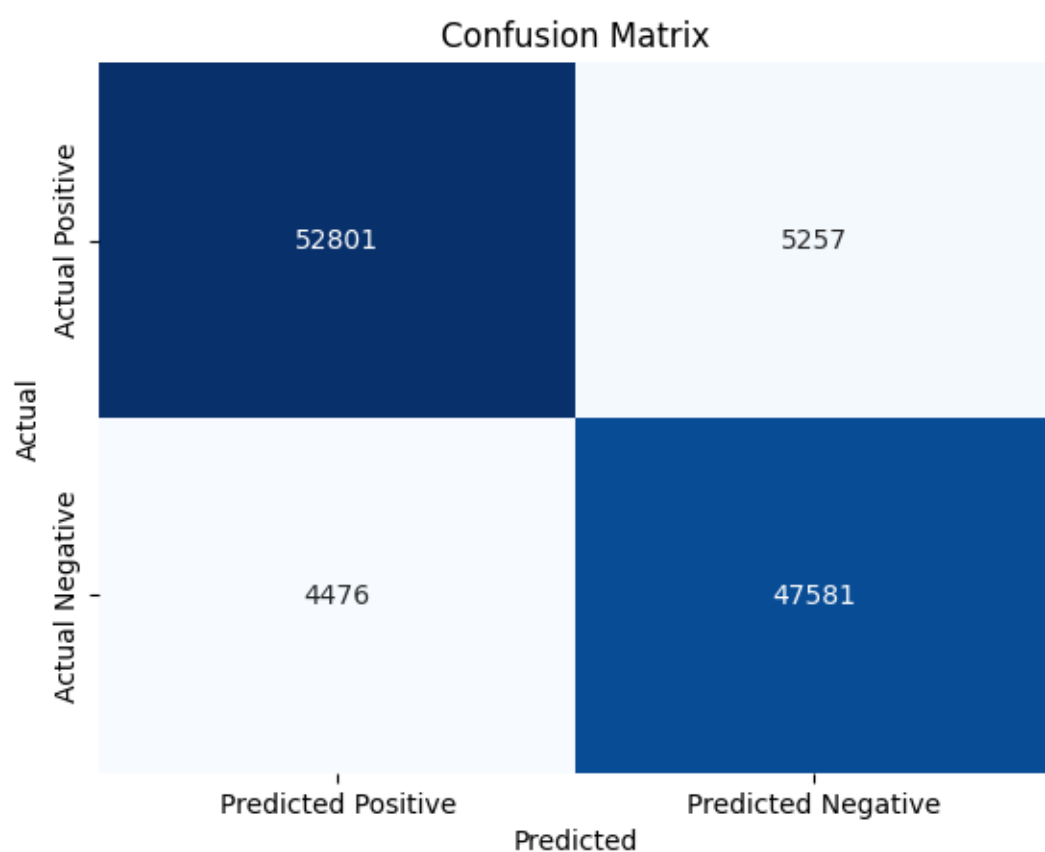
```
In [53]: 1 evaluate_model(rf_model, X_train, y_train, X_test, y_test)
```

```
Accuracy: 0.9116
Precision: 0.9005
Recall: 0.9140
F1 Score: 0.9072
ROC AUC Score: 0.9117
```

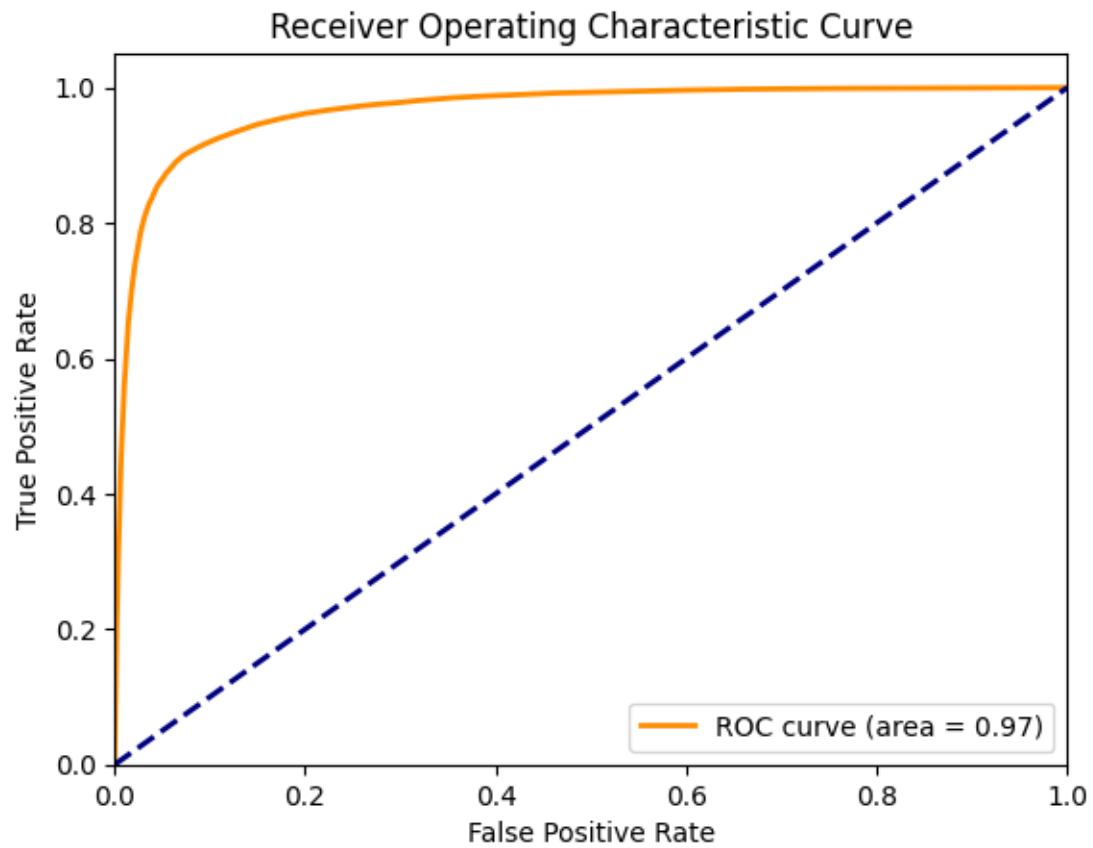
```
Classification Report:
```

	precision	recall	f1-score	support
0	0.92	0.91	0.92	58058
1	0.90	0.91	0.91	52057
accuracy			0.91	110115
macro avg	0.91	0.91	0.91	110115
weighted avg	0.91	0.91	0.91	110115





In training dataset. ROC Curve shows that area = 1 which is equals to the AUC of 1 that indicates perfect performance.



In testing dataset, ROC Curve shows that area = **0.97** which means nearer to the AUC of 1 that indicates better performance.

# Form For the Heart Disease Prediction

← → ↺ ⌚ 127.0.0.1:5000

DAA OIC ChatGPT DFA C# Office Coursera Alex Java SmartDraw SQL ORACLE Mathematical Oper... Glassmorphism Logi...

### Heart Disease Prediction

BMI:

Smoking:

--Select Option--

Alcohol Drinking:

--Select Option--

Stroke:

--Select Option--

Physical Health:

15

Mental Health:

15

Difficulty Walking:

--Select Option--

Sex:

--Select Option--

Age Category:

--Select Option--

Race:

--Select Option--

Diabetic:

--Select Option--

Physical Activity:

--Select Option--

General Health:

--Select Option--

Sleep Time (hours):

Asthma:

--Select Option--

Kidney Disease:

--Select Option--

Skin Cancer:

--Select Option--

Predict

**Prediction Result:**

Result:

### Heart Disease Prediction

BMI:

30

Smoking:

Yes

Alcohol Drinking:

No

Stroke:

No

Physical Health:

6

Mental Health:

2

Difficulty Walking:

Yes

Sex:

Male

Age Category:

75-79

Race:

White

Diabetic:

Yes

Physical Activity:

Yes

General Health:

fair

Sleep Time (hours):

8

Asthma:

No

Kidney Disease:

Yes

Skin Cancer:

No

Predict

**Prediction Result:**  
Heart disease probability: 0.77

## Heart Disease Prediction

BMI:	<input type="text" value="30"/>
Smoking:	<input type="text" value="Yes"/>
Alcohol Drinking:	<input type="text" value="No"/>
Stroke:	<input type="text" value="No"/>
Physical Health:	<input type="range" value="30"/>
Mental Health:	<input type="range" value="10"/>
Difficulty Walking:	<input type="text" value="Yes"/>
Sex:	<input type="text" value="Female"/>
Age Category:	<input type="text" value="55-59"/>
Race:	<input type="text" value="White"/>
Diabetic:	<input type="text" value="No"/>
Physical Activity:	<input type="text" value="No"/>
General Health:	<input type="text" value="fair"/>
Sleep Time (hours):	<input type="text" value="6"/>
Asthma:	<input type="text" value="No"/>
Kidney Disease:	<input type="text" value="No"/>
Skin Cancer:	<input type="text" value="No"/>
<input type="button" value="Predict"/>	

### Prediction Result:

Heart disease probability: 0.5

## Conclusion

From the above prediction analysis, we were able to predict the probability of the heart disease or not as in the result there is the heart disease probability is **0.77** which is suffering from heart disease where as there is also the heart disease probability is **0.5** which is not suffering from heart disease using Random Forest Classifier Model.