

```
In [4]: import pandas as pd
#
In [6]: data = pd.read_csv('Heart Disease data.csv')
In [8]: data
Out[8]:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
0    52   1   0      125   212   0         1    168     0     1.0    2   2   3     0
1    53   1   0      140   203   1         0    155     1     3.1   0   0   3     0
2    70   1   0      145   174   0         1    125     1     2.6   0   0   3     0
3    61   1   0      148   203   0         1    161     0     0.0   2   1   3     0
4    62   0   0      138   294   1         1    106     0     1.9   1   3   2     0
...    ...  ...  ...    ...    ...    ...    ...    ...    ...    ...    ...  ...  ...
1020  59   1   1      140   221   0         1    164     1     0.0   2   0   2     1
1021  60   1   0      125   258   0         0    141     1     2.8   1   1   3     0
1022  47   1   0      110   275   0         0    118     1     1.0   1   1   2     0
1023  50   0   0      110   254   0         0    159     0     0.0   2   0   2     1
1024  54   1   0      120   188   0         1    113     0     1.4   1   1   3     0
1025 rows x 14 columns

In [18]: data.shape
Out[18]: (1025, 14)
In [12]: data.head()
Out[12]:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  target
0    52   1   0      125   212   0         1    168     0     1.0    2   2   3     0
1    53   1   0      140   203   1         0    155     1     3.1   0   0   3     0
2    70   1   0      145   174   0         1    125     1     2.6   0   0   3     0
3    61   1   0      148   203   0         1    161     0     0.0   2   1   3     0
4    62   0   0      138   294   1         1    106     0     1.9   1   3   2     0

In [14]: data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1025 entries, 0 to 1024
Data columns (total 14 columns):
 #   Column  Non-Null Count  Dtype
---  ---
 0   age     1025 non-null     int64
 1   sex     1025 non-null     int64
 2   cp      1025 non-null     int64
 3   trestbps 1025 non-null     int64
 4   chol    1025 non-null     int64
 5   fbs     1025 non-null     int64
 6   restecg 1025 non-null     int64
 7   thalach 1025 non-null     int64
 8   exang    1025 non-null     int64
 9   oldpeak 1025 non-null     float64
10  slope    1025 non-null     int64
11  ca       1025 non-null     int64
12  thal     1025 non-null     int64
13  target   1025 non-null     int64
dtypes: float64(1), int64(13)
memory usage: 112.1 KB

In [16]: data.describe()
Out[16]:
   age      sex      cp  trestbps      chol      fbs  restecg  thalach  exang  oldpeak  slope      ca      thal  target
count  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000  1025.000000
mean    54.434146    0.695610    0.942439   131.611707   246.000000    0.149268    0.529756   149.114146    0.336585    1.071512    1.385366    0.754146    2.323902    0.513171
std     9.072290    0.460373    1.029641   17.516718   51.592511    0.356527    0.527878   23.005724    0.472772    1.175053    0.617755    1.030798    0.620660    0.500070
min     29.000000    0.000000    0.000000    94.000000   126.000000    0.000000    0.000000    71.000000    0.000000    0.000000    0.000000    0.000000    0.000000    0.000000
25%    48.000000    0.000000    0.000000   120.000000   211.000000    0.000000    0.000000   132.000000    0.000000    0.000000    1.000000    0.000000    2.000000    0.000000
50%    56.000000   1.000000    1.000000   130.000000   240.000000    0.000000    1.000000   152.000000    0.000000    0.800000    1.000000    0.000000    2.000000    1.000000
75%    61.000000   1.000000    1.000000   140.000000   275.000000    0.000000    1.000000   166.000000    1.000000    1.800000    2.000000    1.000000    3.000000    1.000000
max     77.000000   1.000000    3.000000   200.000000   564.000000    1.000000    2.000000   202.000000    1.000000    6.200000    2.000000    4.000000    3.000000    1.000000

In [18]: missing_values_count = data.isnull().sum()
# Print the count of missing values in each column
print("Missing Values in Each Column:")
print(missing_values_count)
Missing Values in Each Column:
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

In [28]: data.columns
Out[28]: Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach', 'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
              dtype='object')
In [22]: data.columns = [str(col) for col in data.columns]
In [24]: # Scale numerical features if necessary
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
data[['age', 'trestbps', 'chol', 'thalach', 'oldpeak']] = scaler.fit_transform(data[['age', 'trestbps', 'chol', 'thalach', 'oldpeak']])
In [26]: import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
# Correlation matrix
plt.figure(figsize=(12, 8))
sns.heatmap(data.corr(), annot=True, cmap='coolwarm')
plt.show()
age      sex      cp  trestbps      chol      fbs  restecg  thalach  exang  oldpeak  slope      ca      thal  target
sex      1.000000  -0.072  -0.041  0.038  -0.082  0.079  0.044  0.31  -0.4  -0.17  0.13  -0.18  -0.16  0.43
cp      -0.072  -0.041      1.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
trestbps 0.27  0.079  0.038      1.000000  0.13  0.18  -0.12  -0.039  0.061  0.19  -0.12  0.1  0.059  -0.14
chol     0.22  -0.2  -0.082  0.13      1.000000  0.027  -0.15  -0.022  0.067  0.065  -0.014  0.074  0.1  -0.1
fbs      0.12  0.027  0.079  0.18  0.027      1.000000  -0.1  -0.0089  0.049  0.011  -0.062  0.14  -0.042  -0.041
restecg  -0.13  -0.055  0.044  -0.12  -0.15  -0.1      1.000000  0.048  -0.066  -0.05  0.086  -0.078  -0.021  0.13
thalach  -0.39  -0.049  0.31  -0.039  -0.022  -0.0089  0.048      1.000000  -0.38  -0.35  0.4  -0.21  -0.098  0.42
exang    -0.088  0.14  -0.4  0.061  0.067  0.049  -0.066  -0.38      1.000000  0.31  -0.27  0.11  0.2  -0.44
oldpeak  -0.21  0.085  -0.17  0.19  0.065  0.011  -0.05  -0.35  0.31      1.000000  -0.58  0.22  0.2  -0.44
slope    -0.17  -0.027  0.13  -0.12  -0.014  -0.062  0.086  0.4  -0.27  -0.58      1.000000  -0.073  -0.094  0.35
ca       0.27  0.11  -0.18  0.1  0.074  0.14  -0.078  -0.21  0.11  0.22  -0.073      1.000000  0.15  -0.38
thal     -0.072  0.2  -0.16  0.059  0.1  -0.042  -0.021  -0.098  0.2  0.2  -0.094  0.15      1.000000  -0.34
target   -0.23  -0.28  0.43  -0.14  -0.1  -0.041  0.13  0.42  -0.44  -0.44  0.35  -0.38  -0.34      1.000000

In [27]: # Distribution plots
age
sns.histplot(data['age'], kde=True)
plt.title('Age Distribution')
plt.show()
Age Distribution
Count
160
140
120
100
80
60
40
20
0
-3 -2 -1 0 1 2
age

In [68]: import plotly.express as px
import matplotlib.pyplot as plt
import seaborn as sns
# Create a bar plot using Plotly Express
fig = px.bar(data, x='sex', color='target', barmode='group',
             category_orders={'sex': ['M', 'F']},
             labels={'sex': 'Gender', 'target': 'Heart Disease'},
             title='Heart Disease by Gender')
# Update the layout for better visualization
fig.update_layout(
    xaxis=dict(
        tickmode='array',
        tickvals=[0, 1],
        ticktext=['Female', 'Male']
    ),
    xaxis_title='Gender',
    yaxis_title='Count'
)
# Show the plot
fig.show()
# Pairplot
sns.pairplot(data, hue='target')
plt.show()

In [38]: # Analyzing the relationship between age and maximum heart rate achieved
sns.scatterplot(x='age', y='thalach', hue='target', data=data)
plt.title('Age vs. Maximum Heart Rate by Heart Disease Status')
plt.show()
# Analyzing the impact of chest pain type on heart disease
sns.catplot(x='cp', hue='target', data=data)
plt.title('Heart Disease by Chest Pain Type')
plt.show()
Age vs. Maximum Heart Rate by Heart Disease Status
thalach
2
1
0
-1
-2
-3
-3 -2 -1 0 1 2
age
target
0
1

Heart Disease by Chest Pain Type
count
350
300
250
200
150
100
50
0
0 1 2 3
cp
target
0
1

In [38]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import classification_report
# Encode categorical variables if necessary
# For this example, we assume that all categorical variables are already numerical
# Define features and target
X = data.drop('target', axis=1)
y = data['target']
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Standardize the numerical features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
In [31]: # Train a Random Forest classifier
clf = RandomForestClassifier(n_estimators=100, random_state=42)
clf.fit(X_train, y_train)
# Evaluate the model
y_pred = clf.predict(X_test)
print(classification_report(y_test, y_pred))
precision    recall  f1-score   support
0       0.97       1.00       0.99         182
1       1.00       0.97       0.99         103
accuracy          0.99         0.99         285
macro avg          0.99         0.99         285
weighted avg          0.99         0.99         285
In [32]: def predict_heart_disease(patient_data):
    """
    Predict whether a patient has heart disease based on their data.
    Parameters:
        patient_data (list or array): A list or array of patient features in the same order as the training data.
    Returns:
        str: 'Yes' if the patient is predicted to have heart disease, 'No' otherwise.
    """
    # Ensure the input data is a numpy array
    patient_data = np.array(patient_data).reshape(1, -1)
    # Standardize the patient data
    patient_data = scaler.transform(patient_data)
    # Make a prediction
    prediction = clf.predict(patient_data)
    return 'Yes' if prediction[0] == 1 else 'No'
# Example usage
new_patient = [63, 1, 3, 145, 233, 1, 0, 150, 0, 2.3, 0, 0, 1] # patient data
print(predict_heart_disease(new_patient))
Yes
C:\Users\usman\anaconda3\lib\site-packages\sklearn\base.py:493: UserWarning:
X does not have valid feature names, but StandardScaler was fitted with feature names
In [32]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
import pickle
# Check for missing values and handle them
data = data.dropna()
# Define features and target
X = data.drop('target', axis=1)
y = data['target']
# Split the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Standardize the numerical features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
# Train a Random Forest classifier
clf = RandomForestClassifier(n_estimators=100, random_state=42)
clf.fit(X_train, y_train)
# Save the trained model and scaler as a pickle file
with open('heart_disease_model.pkl', 'wb') as model_file:
    pickle.dump((clf, scaler), model_file)
print('Model and scaler saved')
Model and scaler saved
In [ ]:
```