

```

1 Objective:
2     1) To analysis which factor is most affects for heart diseases
3     2) To analyse which age group has more chances to get heart attack

```

```

In [3]: 1 import pandas as pd
        2 import numpy as np
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5 from sklearn.linear_model import LogisticRegression
        6 from sklearn.neighbors import KNeighborsClassifier
        7 from sklearn.ensemble import RandomForestClassifier
        8 from sklearn.model_selection import train_test_split, cross_val_score
        9 from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
       10 from sklearn.metrics import confusion_matrix, classification_report
       11 from sklearn.metrics import precision_score, recall_score, f1_score
       12 from sklearn.metrics import roc_curve, auc

```

```

In [4]: 1 df= pd.read_csv("C:\\Users\\mahes\\Downloads\\heart.csv")

```

```

In [5]: 1 df

```

```

Out[5]:
   age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  ca  thal  ta
0    52    1   0     125   212    0         1     168     0        1.0      2   2    3
1    53    1   0     140   203    1         0     155     1        3.1      0   0    3
2    70    1   0     145   174    0         1     125     1        2.6      0   0    3
3    61    1   0     148   203    0         1     161     0        0.0      2   1    3
4    62    0   0     138   294    1         1     106     0        1.9      1   3    2
...   ...  ...  ...   ...   ...   ...   ...   ...   ...   ...   ...  ...  ...
1020  59    1   1     140   221    0         1     164     1        0.0      2   0    2
1021  60    1   0     125   258    0         0     141     1        2.8      1   1    3
1022  47    1   0     110   275    0         0     118     1        1.0      1   1    2
1023  50    0   0     110   254    0         0     159     0        0.0      2   0    2
1024  54    1   0     120   188    0         1     113     0        1.4      1   1    3

```

1025 rows × 14 columns



```

In [6]: 1 #data first 5 row
        2 df.head()

```

```

Out[6]:
   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
1    53    1   0     140   203    1         0     155     1        3.1      0   0    3
2    70    1   0     145   174    0         1     125     1        2.6      0   0    3
3    61    1   0     148   203    0         1     161     0        0.0      2   1    3
4    62    0   0     138   294    1         1     106     0        1.9      1   3    2

```



In [7]:

```
1 #data last 5 rows
2 df.tail()
```

Out[7]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	ta
1020	59	1	1	140	221	0	1	164	1	0.0	2	0	2	
1021	60	1	0	125	258	0	0	141	1	2.8	1	1	3	
1022	47	1	0	110	275	0	0	118	1	1.0	1	1	2	
1023	50	0	0	110	254	0	0	159	0	0.0	2	0	2	
1024	54	1	0	120	188	0	1	113	0	1.4	1	1	3	

In [8]:

```
1 df.describe()
```

Out[8]:

	age	sex	cp	trestbps	chol	fbs	restecg
count	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	0.529268
std	9.072290	0.460373	1.029641	17.516718	51.59251	0.356527	0.527268
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000
25%	48.000000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000
50%	56.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000
75%	61.000000	1.000000	2.000000	140.000000	275.000000	0.000000	1.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000

In [9]:

```
1 #Let, we get some basic information about the dataset.
2 df.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: 112.2 KB
```

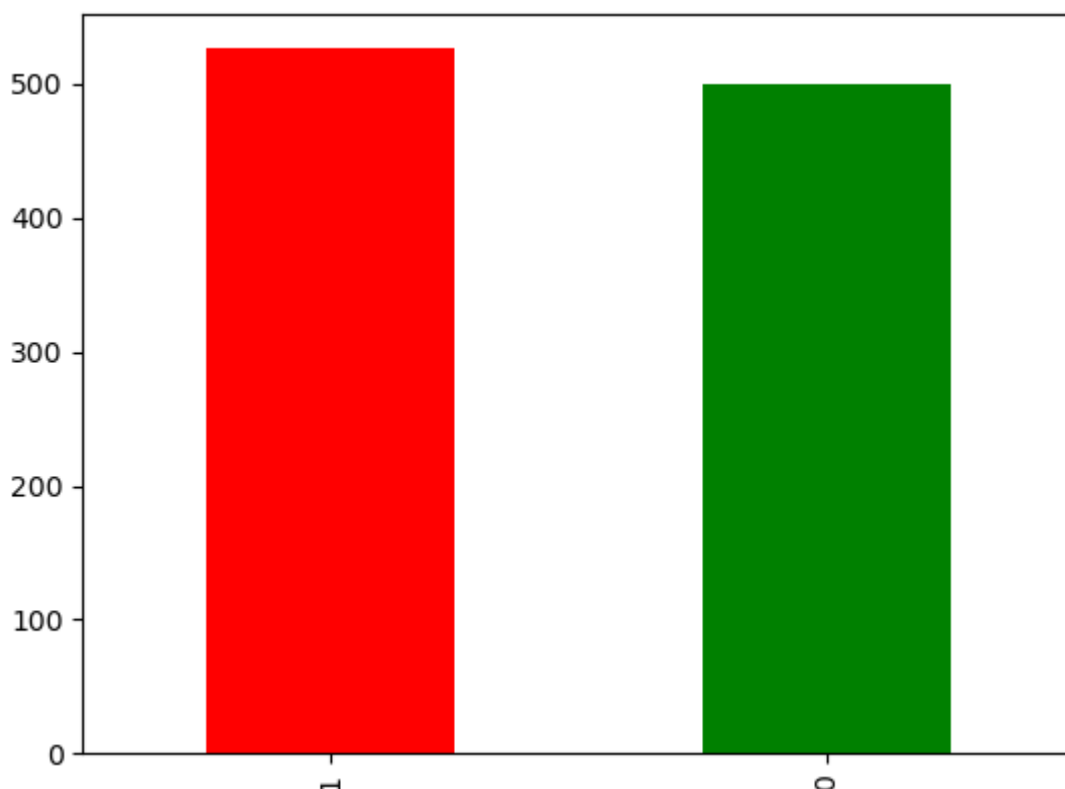
```
In [10]: 1 #no of row and no of column  
2 df.shape
```

```
Out[10]: (1025, 14)
```

```
In [11]: 1 df['target'].value_counts()
```

```
Out[11]: 1    526  
0    499  
Name: target, dtype: int64
```

```
In [12]: 1 df['target'].value_counts().plot(kind='bar', color=['red','green']);
```



```
In [13]: 1 # Let us look at whether the dataset has null values or not.  
2 df.isna().sum()
```

```
Out[13]: 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
```

From this output, our data does not contain null values and duplicates. So, the data is good which will be further analyzed.

```
In [14]: 1 df_dup = df.duplicated().any()
```

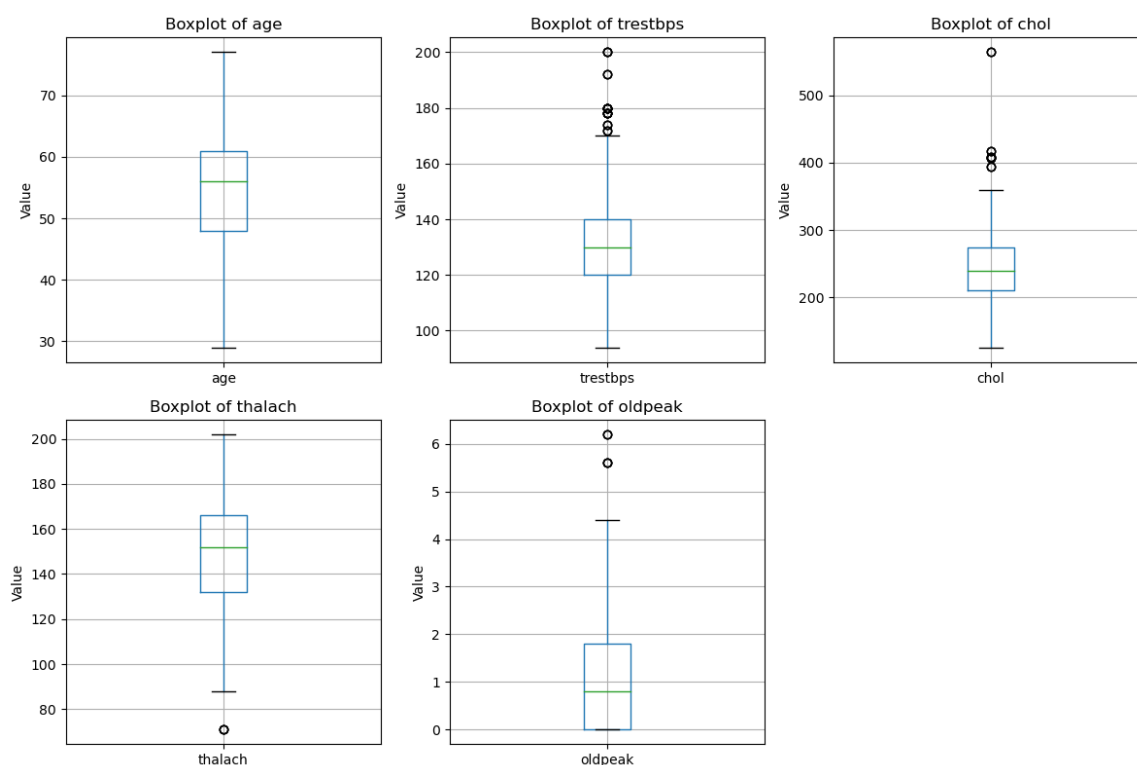
```
In [15]: 1 df_dup
```

```
Out[15]: True
```

```
In [16]: 1 df.sex.value_counts()
```

```
Out[16]: 1    713
          0    312
          Name: sex, dtype: int64
```

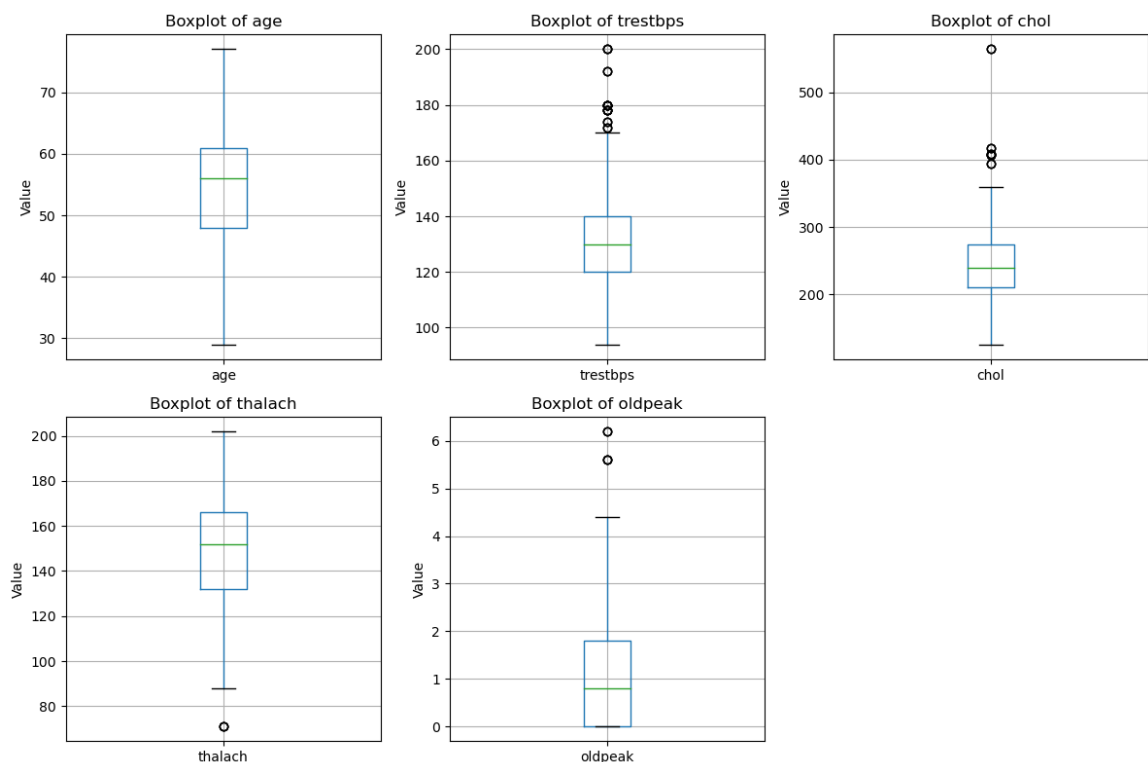
```
In [17]: 1 import pandas as pd
          2 import matplotlib.pyplot as plt
          3
          4 # Assuming 'data' is your heart disease prediction dataset loaded into a variable
          5
          6 # Select relevant features for visualization
          7 features = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
          8
          9 # Create boxplots for each feature to visualize outliers
          10 plt.figure(figsize=(12, 8))
          11 for i, feature in enumerate(features, start=1):
          12     plt.subplot(2, 3, i)
          13     df.boxplot(column=feature)
          14     plt.title(f'Boxplot of {feature}')
          15     plt.ylabel('Value')
          16
          17 plt.tight_layout()
          18 plt.show()
```



```

In [18]: 1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 # Load your heart disease prediction dataset into a pandas DataFrame
5 # Assuming 'data' is your DataFrame containing the dataset
6
7 # Select relevant features for visualization
8 features = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
9
10 # Create boxplots for each feature to visualize outliers
11 plt.figure(figsize=(12, 8))
12 for i, feature in enumerate(features, start=1):
13     plt.subplot(2, 3, i)
14     df.boxplot(column=feature)
15     plt.title(f'Boxplot of {feature}')
16     plt.ylabel('Value')
17
18 plt.tight_layout()
19 plt.show()
20
21 # After visual inspection, determine the range for each feature to filter
22 # For simplicity, let's assume a range for each feature based on boxplot
23
24 # Define the range for each feature based on your observations from boxplot
25 # Example ranges (you need to adjust these based on your data and insights)
26 ranges = {
27     'age': (20, 80),
28     'trestbps': (80, 180),
29     'chol': (100, 400),
30     'thalach': (60, 220),
31     'oldpeak': (0, 6)
32 }
33
34 # Apply the range filter to remove outliers from the dataset
35 for feature, (lower, upper) in ranges.items():
36     data = df[(df[feature] >= lower) & (df[feature] <= upper)]

```



```

In [19]: 1 features = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
          2
          3 # Calculate the IQR for each feature
          4 Q1 = data[features].quantile(0.25)
          5 Q3 = data[features].quantile(0.75)
          6 IQR = Q3 - Q1
          7
          8 # Define the lower and upper bounds for outliers
          9 lower_bound = Q1 - 1.5 * IQR
         10 upper_bound = Q3 + 1.5 * IQR
         11
         12 # Create a mask to identify outliers
         13 outlier_mask = ~((data[features] < lower_bound) | (data[features] > upper_bound))
         14
         15 # Filter the dataset to remove outliers
         16 data_no_outliers = data[outlier_mask]
         17
         18 # Print the number of outliers removed
         19 num_outliers_removed = len(data) - len(data_no_outliers)
         20 print(f"Number of outliers removed: {num_outliers_removed}")
         21 plt.show()

```

Number of outliers removed: 54

In [20]: 1 df

Out[20]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	ta
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3	
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3	
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	
...
1020	59	1	1	140	221	0	1	164	1	0.0	2	0	2	
1021	60	1	0	125	258	0	0	141	1	2.8	1	1	3	
1022	47	1	0	110	275	0	0	118	1	1.0	1	1	2	
1023	50	0	0	110	254	0	0	159	0	0.0	2	0	2	
1024	54	1	0	120	188	0	1	113	0	1.4	1	1	3	

1025 rows × 14 columns



In [21]:

1 data_no_outliers

Out[21]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	ta
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3	
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3	
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3	
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3	
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2	
...
1020	59	1	1	140	221	0	1	164	1	0.0	2	0	2	
1021	60	1	0	125	258	0	0	141	1	2.8	1	1	3	
1022	47	1	0	110	275	0	0	118	1	1.0	1	1	2	
1023	50	0	0	110	254	0	0	159	0	0.0	2	0	2	
1024	54	1	0	120	188	0	1	113	0	1.4	1	1	3	

968 rows × 14 columns

In [22]:

1 len(df)

Out[22]:

1025

In [23]:

1 713/1025

Out[23]:

0.6956097560975609

In [24]:

1 312/1025

Out[24]:

0.304390243902439

In [25]:

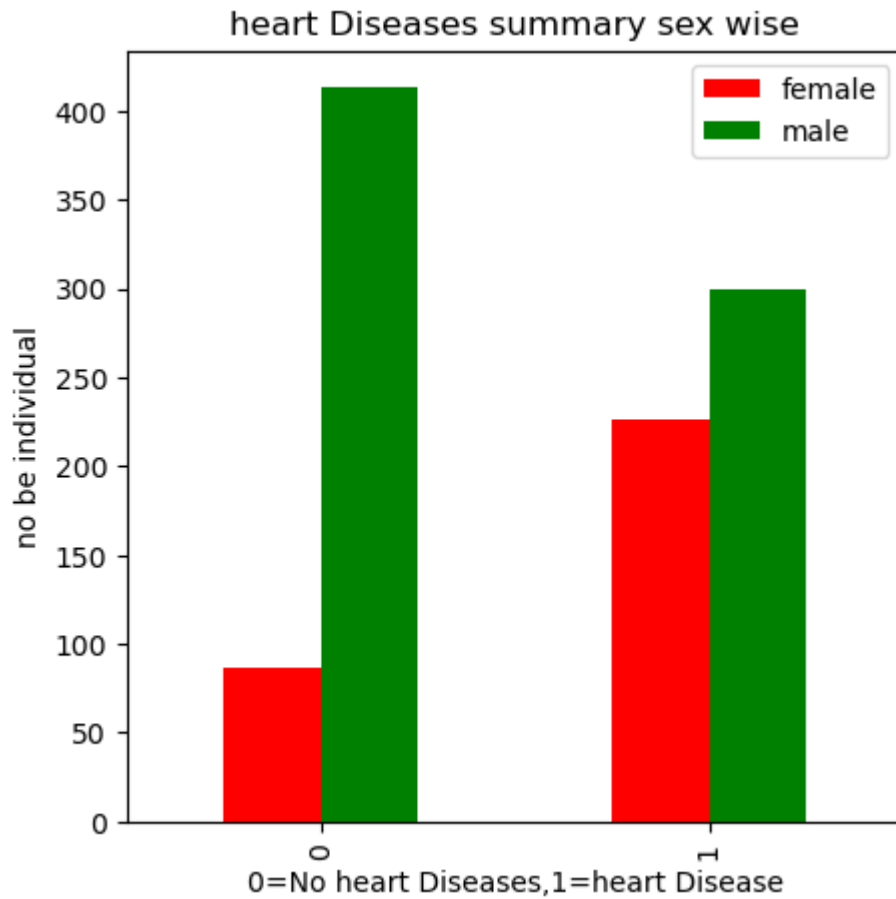
1 pd.crosstab(df.target,df.sex)

Out[25]:

sex	0	1
target		
0	86	413
1	226	300

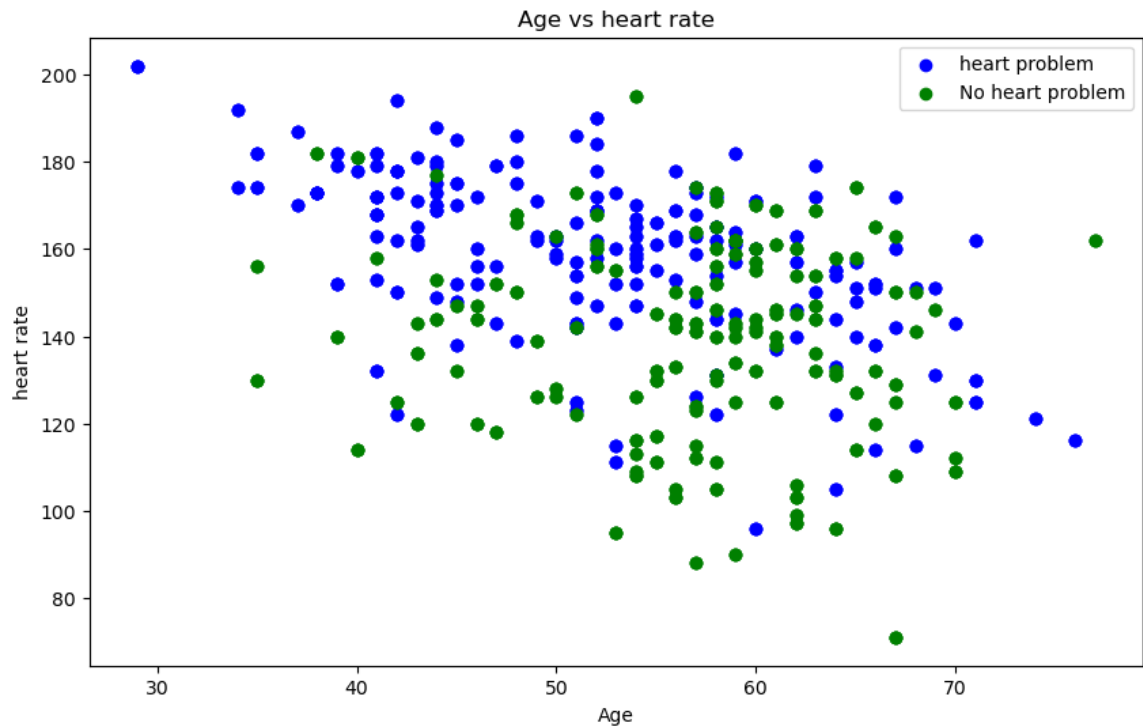
```
In [26]: 1 pd.crosstab(df.target,df.sex).plot(kind='bar',
2         figsize=(5,5),
3         color=['red','green'])
4 plt.title('heart Diseases summary sex wise')
5 plt.xlabel('0=No heart Diseases,1=heart Disease')
6 plt.ylabel('no be individual')
7 plt.legend(['female','male'])
8
```

Out[26]: <matplotlib.legend.Legend at 0x17f05817550>

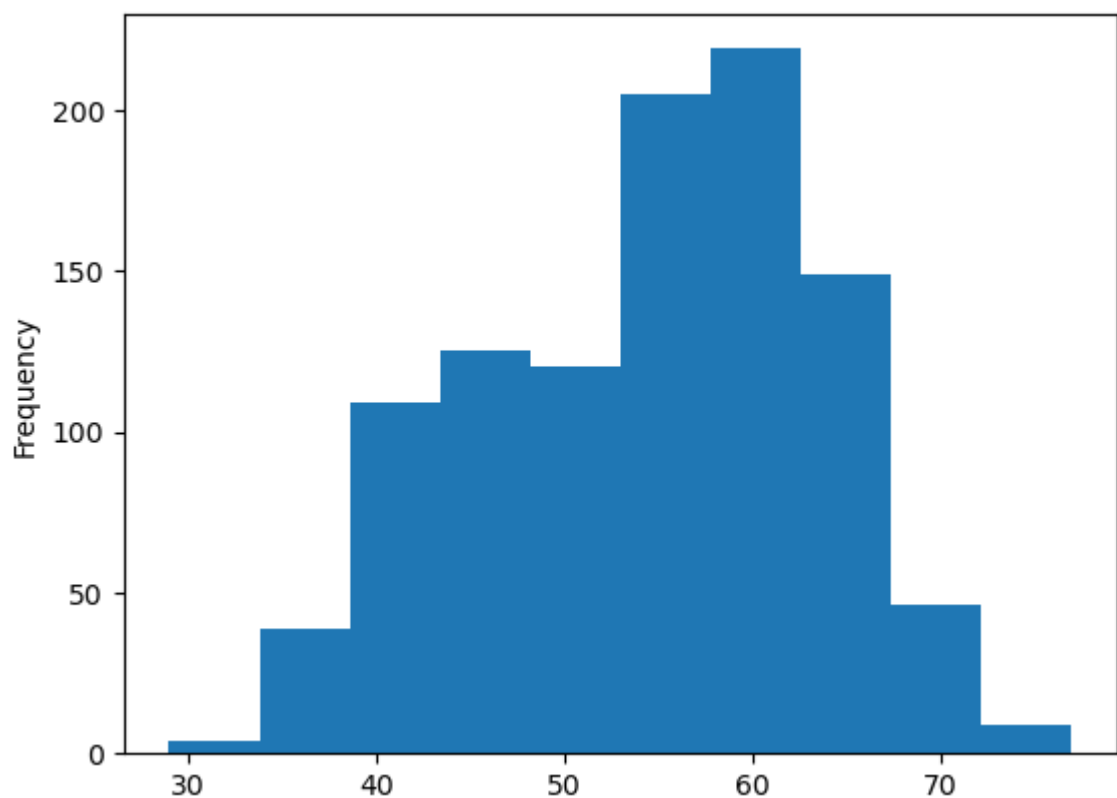



```
In [27]: 1 plt.figure(figsize=(10,6))
2 plt.scatter(df.age[df.target==1],df.thalach[df.target==1],c='blue')
3 plt.scatter(df.age[df.target==0],df.thalach[df.target==0],c='green')
4 plt.title('Age vs heart rate')
5 plt.xlabel('Age')
6 plt.ylabel('heart rate')
7 plt.legend(['heart problem', 'No heart problem'])
```

Out[27]: <matplotlib.legend.Legend at 0x17f05356f70>

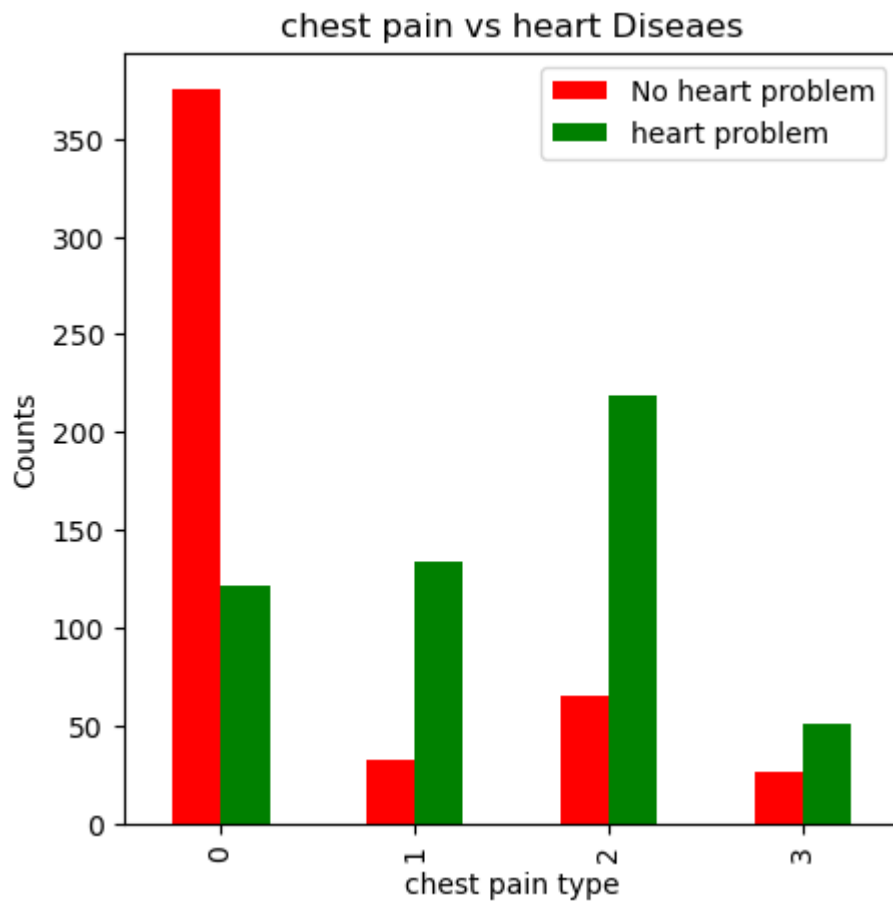


```
In [28]: 1 df.age.plot.hist();
```



```
In [29]: 1 pd.crosstab(df.cp,df.target).plot(kind='bar',  
2           figsize=(5,5),  
3           color=['red','green'])  
4 plt.title('chest pain vs heart Diseases')  
5 plt.xlabel('chest pain type')  
6 plt.ylabel('Counts')  
7 plt.legend(['No heart problem','heart problem'])
```

Out[29]: <matplotlib.legend.Legend at 0x17f05437970>

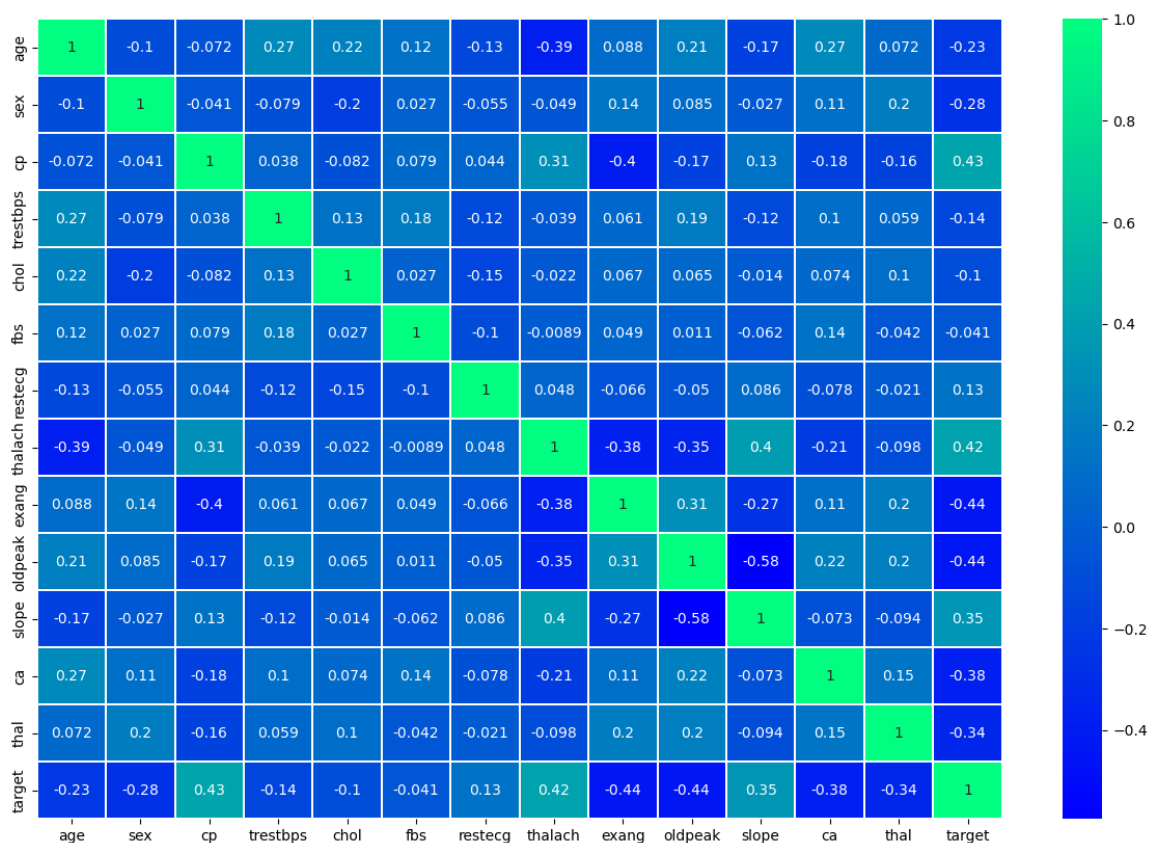


In [30]: 1 df.corr()

Out[30]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach
age	1.000000	-0.103240	-0.071966	0.271121	0.219823	0.121243	-0.132696	-0.390227
sex	-0.103240	1.000000	-0.041119	-0.078974	-0.198258	0.027200	-0.055117	-0.049365
cp	-0.071966	-0.041119	1.000000	0.038177	-0.081641	0.079294	0.043581	0.306839
trestbps	0.271121	-0.078974	0.038177	1.000000	0.127977	0.181767	-0.123794	-0.039264
chol	0.219823	-0.198258	-0.081641	0.127977	1.000000	0.026917	-0.147410	-0.021772
fbs	0.121243	0.027200	0.079294	0.181767	0.026917	1.000000	-0.104051	-0.008866
restecg	-0.132696	-0.055117	0.043581	-0.123794	-0.147410	-0.104051	1.000000	0.048411
thalach	-0.390227	-0.049365	0.306839	-0.039264	-0.021772	-0.008866	0.048411	1.000000
exang	0.088163	0.139157	-0.401513	0.061197	0.067382	0.049261	-0.065606	-0.380281
oldpeak	0.208137	0.084687	-0.174733	0.187434	0.064880	0.010859	-0.050114	-0.349796
slope	-0.169105	-0.026666	0.131633	-0.120445	-0.014248	-0.061902	0.086086	0.395308
ca	0.271551	0.111729	-0.176206	0.104554	0.074259	0.137156	-0.078072	-0.207888
thal	0.072297	0.198424	-0.163341	0.059276	0.100244	-0.042177	-0.020504	-0.098068
target	-0.229324	-0.279501	0.434854	-0.138772	-0.099966	-0.041164	0.134468	0.422895

In [31]: 1 plt.figure(figsize=(15,10))
 2 sns.heatmap(df.corr(),linewidth=.01,annot=True,cmap="winter")
 3 plt.show()
 4 plt.savefig('correlationfigure')



<Figure size 640x480 with 0 Axes>

we can understand that Chest pain(cp) and target have a positive correlation. It means that whose has a large risk of chest pain results in a greater chance to have heart disease. In addition to chest pain, thalach, slope, and resting have a positive correlation with the target. exang and the target have a negative correlation which means when we exercise, the heart requires more blood, but narrowed arteries slow down the blood flow. In addition to ca, oldpeak, thal have a negative correlation with the target.

Splitting the dataset into the Training Set and Test Set

```
In [32]: 1 x=df.drop('target', axis=1)
          2 y=df['target']
```

```
In [33]: 1 x
```

```
Out[33]:
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2
...
1020	59	1	1	140	221	0	1	164	1	0.0	2	0	2
1021	60	1	0	125	258	0	0	141	1	2.8	1	1	3
1022	47	1	0	110	275	0	0	118	1	1.0	1	1	2
1023	50	0	0	110	254	0	0	159	0	0.0	2	0	2
1024	54	1	0	120	188	0	1	113	0	1.4	1	1	3

1025 rows × 13 columns

```
In [34]: 1 y
```

```
Out[34]: 0      0
          1      0
          2      0
          3      0
          4      0
          ..
         1020    1
         1021    0
         1022    0
         1023    1
         1024    0
          Name: target, Length: 1025, dtype: int64
```

In [35]: 1 df.corr()

Out[35]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach
age	1.000000	-0.103240	-0.071966	0.271121	0.219823	0.121243	-0.132696	-0.390227
sex	-0.103240	1.000000	-0.041119	-0.078974	-0.198258	0.027200	-0.055117	-0.049365
cp	-0.071966	-0.041119	1.000000	0.038177	-0.081641	0.079294	0.043581	0.306839
trestbps	0.271121	-0.078974	0.038177	1.000000	0.127977	0.181767	-0.123794	-0.039264
chol	0.219823	-0.198258	-0.081641	0.127977	1.000000	0.026917	-0.147410	-0.021772
fbs	0.121243	0.027200	0.079294	0.181767	0.026917	1.000000	-0.104051	-0.008866
restecg	-0.132696	-0.055117	0.043581	-0.123794	-0.147410	-0.104051	1.000000	0.048411
thalach	-0.390227	-0.049365	0.306839	-0.039264	-0.021772	-0.008866	0.048411	1.000000
exang	0.088163	0.139157	-0.401513	0.061197	0.067382	0.049261	-0.065606	-0.380281
oldpeak	0.208137	0.084687	-0.174733	0.187434	0.064880	0.010859	-0.050114	-0.349796
slope	-0.169105	-0.026666	0.131633	-0.120445	-0.014248	-0.061902	0.086086	0.395308
ca	0.271551	0.111729	-0.176206	0.104554	0.074259	0.137156	-0.078072	-0.207888
thal	0.072297	0.198424	-0.163341	0.059276	0.100244	-0.042177	-0.020504	-0.098068
target	-0.229324	-0.279501	0.434854	-0.138772	-0.099966	-0.041164	0.134468	0.422895

In [36]: 1 dataset=pd.get_dummies(df,columns=["sex","cp","fbs","restecg","exang","s

In [37]: 1 from sklearn .model_selection import train_test_split
 2 from sklearn .preprocessing import StandardScaler
 3 standardscalar = StandardScaler()
 4 Columns_to_Scale=["age","trestbps","chol","thalach","oldpeak"]
 5

In [38]: 1 dataset.head()

Out[38]:

	age	trestbps	chol	thalach	oldpeak	target	sex_0	sex_1	cp_0	cp_1	...	slope_2	ca_0
0	52	125	212	168	1.0	0	0	1	1	0	...	1	0
1	53	140	203	155	3.1	0	0	1	1	0	...	0	1
2	70	145	174	125	2.6	0	0	1	1	0	...	0	1
3	61	148	203	161	0.0	0	0	1	1	0	...	1	0
4	62	138	294	106	1.9	0	1	0	1	0	...	0	0

5 rows × 31 columns

In [39]: 1 x=df.drop('target', axis=1)
 2 y=df['target']

In [40]:

1 x

Out[40]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	52	1	0	125	212	0	1	168	0	1.0	2	2	3
1	53	1	0	140	203	1	0	155	1	3.1	0	0	3
2	70	1	0	145	174	0	1	125	1	2.6	0	0	3
3	61	1	0	148	203	0	1	161	0	0.0	2	1	3
4	62	0	0	138	294	1	1	106	0	1.9	1	3	2
...
1020	59	1	1	140	221	0	1	164	1	0.0	2	0	2
1021	60	1	0	125	258	0	0	141	1	2.8	1	1	3
1022	47	1	0	110	275	0	0	118	1	1.0	1	1	2
1023	50	0	0	110	254	0	0	159	0	0.0	2	0	2
1024	54	1	0	120	188	0	1	113	0	1.4	1	1	3

1025 rows × 13 columns

In [41]:

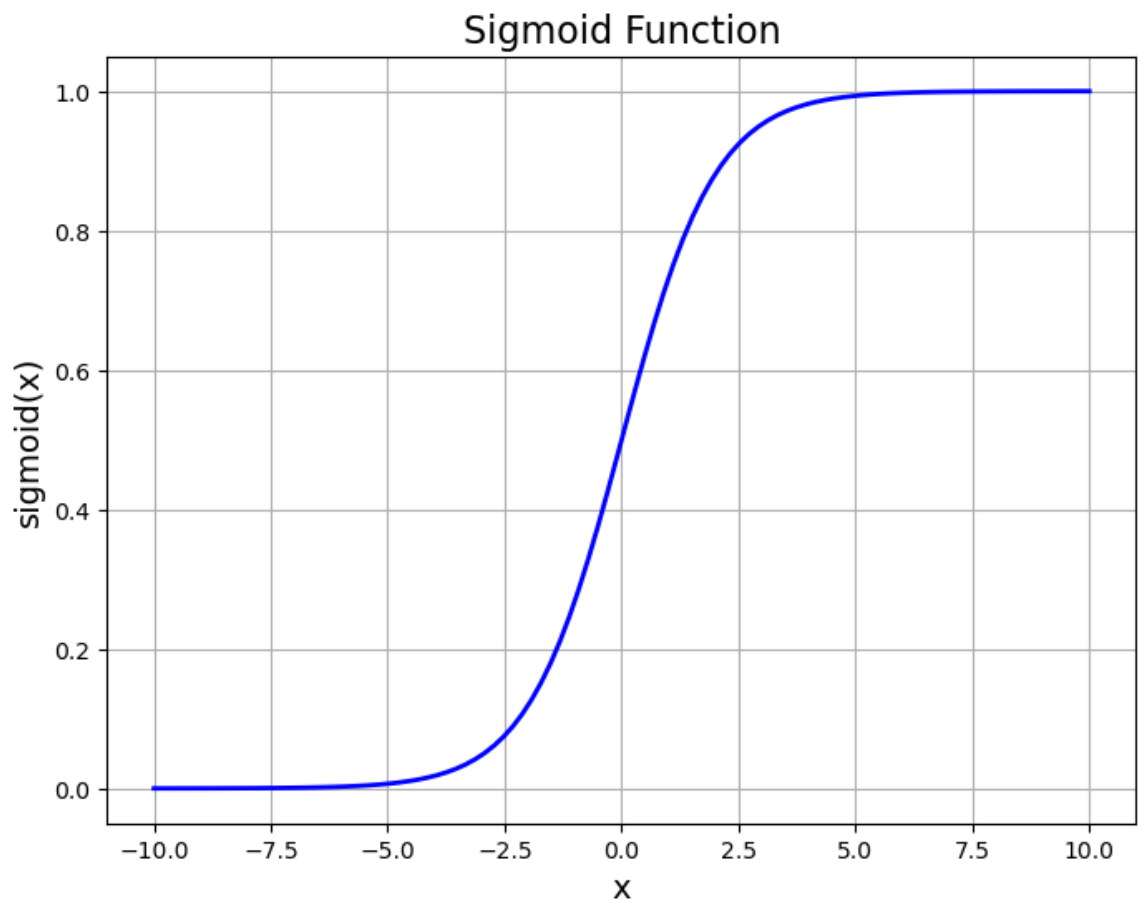
```

1 np.random.seed(7)
2 x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.2)

```

logistic Regression

```
In [42]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 # Define the sigmoid function
5 def sigmoid(x):
6     return 1 / (1 + np.exp(-x))
7
8 # Generate values for x
9 x = np.linspace(-10, 10, 100)
10
11 # Compute the corresponding y values using the sigmoid function
12 y = sigmoid(x)
13
14 # Plot the sigmoid function
15 plt.figure(figsize=(8, 6))
16 plt.plot(x, y, color='blue', linewidth=2)
17
18 # Add labels and title
19 plt.title('Sigmoid Function', fontsize=16)
20 plt.xlabel('x', fontsize=14)
21 plt.ylabel('sigmoid(x)', fontsize=14)
22 plt.grid(True)
23
24 # Show the plot
25 plt.show()
```



```
In [43]: 1 from sklearn.preprocessing import StandardScaler
2 from sklearn.linear_model import LogisticRegression
3 from sklearn.metrics import accuracy_score, classification_report, confu
4 scaler = StandardScaler()
5 x_train = scaler.fit_transform(x_train)
6 x_test = scaler.transform(x_test)
7
8 # Initialize and train the logistic regression model
9 log_reg_model = LogisticRegression()
10 log_reg_model.fit(x_train, y_train)
11
12 # Predict on the test set
13 y_pred = log_reg_model.predict(x_test)
14
15 # Evaluate the model
16 accuracy = accuracy_score(y_test, y_pred)
17 print('Accuracy:', accuracy)
18
19 # Classification report
20 print(classification_report(y_test, y_pred))
21
22 # Confusion matrix
23 print('Confusion Matrix:')
24 print(confusion_matrix(y_test, y_pred))
25
```

Accuracy: 0.8682926829268293

	precision	recall	f1-score	support
0	0.84	0.87	0.85	89
1	0.89	0.87	0.88	116
accuracy			0.87	205
macro avg	0.87	0.87	0.87	205
weighted avg	0.87	0.87	0.87	205

Confusion Matrix:

```
[[ 77  12]
 [ 15 101]]
```

Decision tree

```
In [44]: 1 from sklearn import tree
```

```
In [45]: 1 learner = tree.DecisionTreeClassifier(random_state=0)
2 learner.fit(x_train,y_train)
```

Out[45]: DecisionTreeClassifier(random_state=0)

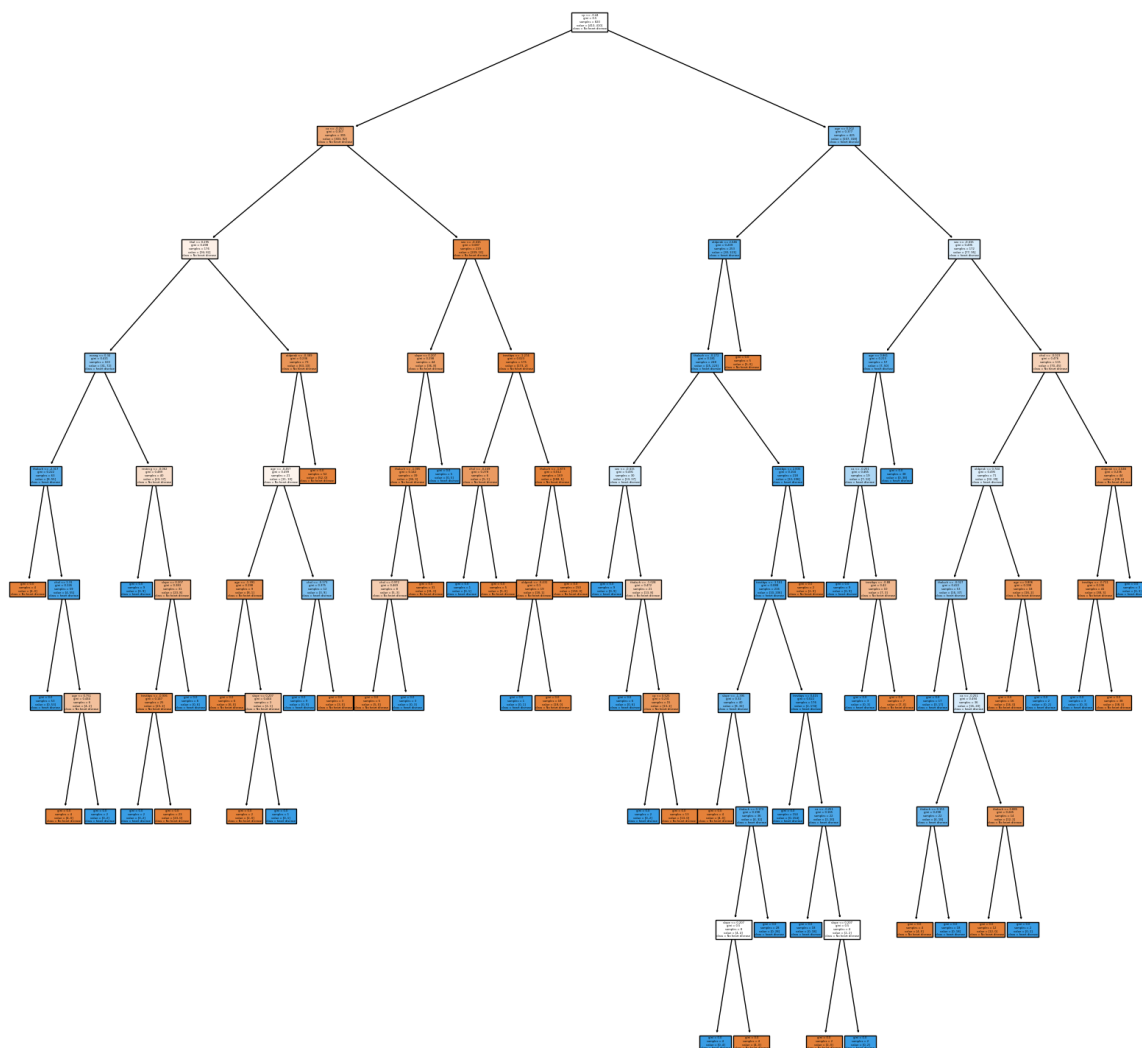
```
In [46]: 1 y_train_pred = learner.predict(x_train)
```

```
In [47]: 1 #prediction for New data
2 y_test_pred = learner.predict(x_test)
```



```
In [48]: 1 import seaborn as sns
2 import matplotlib.pyplot as plt
```

```
In [49]: 1 #visualization the decision tree
2 plt.figure(figsize=(20,20))
3 features=df.columns
4 classes = ['No heart disease', 'heart disease']
5 tree.plot_tree(learner,feature_names=features,class_names=classes,filled=True)
6 plt.show()
```



```
In [50]: 1 from sklearn.metrics import accuracy_score, confusion_matrix
```

```
In [51]: 1 print(f'training accuracy{accuracy_score(y_train_pred,y_train)}')
```

training accuracy1.0

```
In [52]: 1 print(f'Test accuracy{accuracy_score(y_test_pred,y_test)}')
```

Test accuracy1.0

In [53]:

```

1  #post pruning - Cost complexity pruning approach
2  path= learner.cost_complexity_pruning_path(x_train,y_train)
3  #path variable gives two values:ccp_alphas and impurities
4  ccp_alphas,impurities=path.ccp_alphas,path.impurities
5  print("ccp alpha provides list of values:",ccp_alphas)
6  print("Impurities in Decision Tree:",impurities)

```

```

ccp alpha provides list of values: [0.          0.00108401 0.00120557 0.0016
0754 0.00325203 0.00337711
0.00418118 0.00422764 0.00433604 0.00433604 0.0044878  0.00472083
0.00499216 0.00530391 0.00584263 0.00585001 0.00631929 0.00686992
0.00768742 0.00793979 0.00827327 0.0092042  0.00966549 0.00967683
0.01129238 0.01129761 0.0119775  0.02194491 0.02704736 0.03291565
0.04203248 0.13260164]
Impurities in Decision Tree: [0.          0.00216802 0.00699032 0.01181294
0.01506497 0.02181919
0.02600038 0.03022802 0.03890011 0.04323615 0.04772395 0.06188644
0.07686291 0.08747073 0.09331337 0.10501338 0.11765196 0.1313918
0.15445405 0.16239384 0.17066712 0.17987132 0.1895368  0.20889046
0.22018284 0.23148045 0.24345796 0.26540287 0.29245023 0.32536588
0.36739836 0.5          ]

```

In [54]:

```

1  All_learners=[]
2  for ccp_alpha in ccp_alphas:
3      learner=tree.DecisionTreeClassifier(random_state=0,ccp_alpha=ccp_alpha)
4      learner.fit(x_train,y_train)
5      All_learners.append(learner)
6  print("Last node in Decision tree is {} and ccp_alpha for last node is {}".format(

```

```

Last node in Decision tree is <sklearn.tree._tree.Tree object at 0x0000017F
05C559D0> and ccp_alpha for last node is 0.13260163812360393

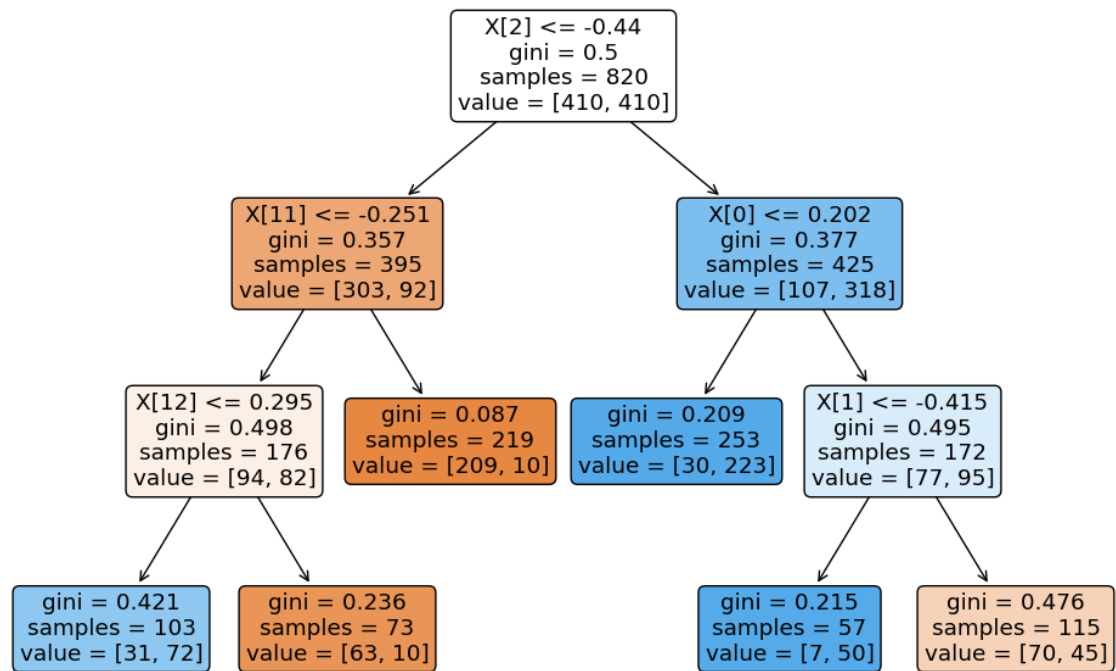
```

```
In [55]: 1 #visulizing the accuracy for train and test set
2 train_scores = [learner.score(x_train,y_train) for learner in All_learner]
3 test_scores = [learner.score(x_test,y_test) for learner in All_learners]
4 fig, ax= plt.subplots()
5 ax.set_xlabel("alpha")
6 ax.set_ylabel("accuracy")
7 ax.set_title("Accuracy vs alpha for training and testing sets")
8 ax.plot(ccp_alphas,train_scores,marker='o',label='train',drawstyle='step')
9 ax.plot(ccp_alphas,test_scores,marker='o',label='train',drawstyle='step')
10 ax.legend()
11 plt.show()
```



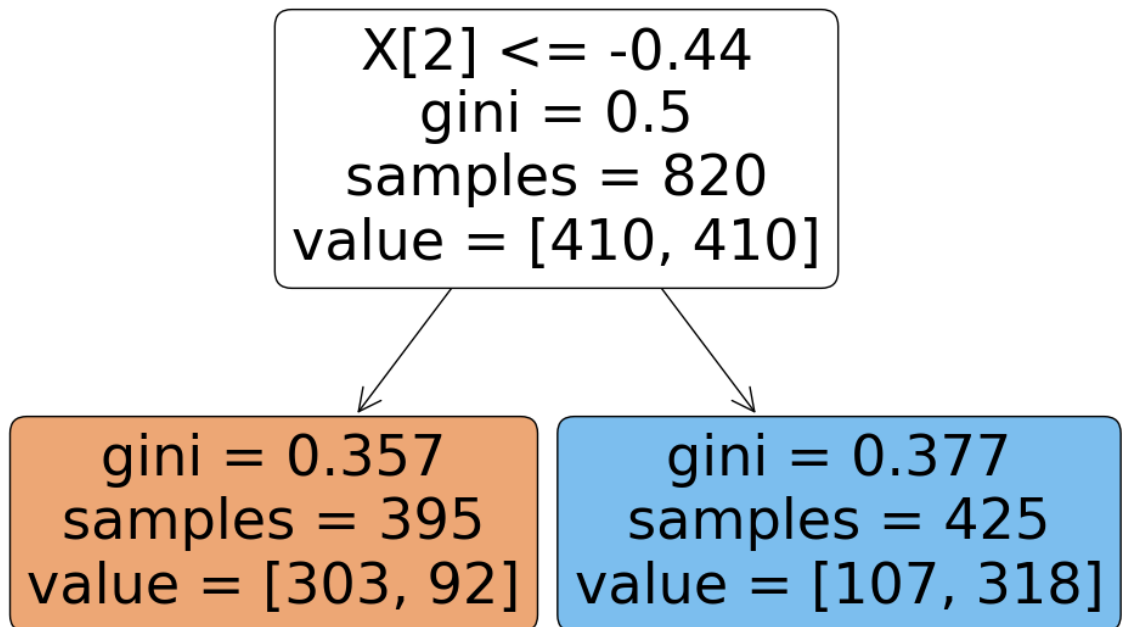
The plot displays the accuracy of the model on the training set and the testing set as a function of the complexity parameter (ccp_alpha). As the complexity parameter increases (meaning more aggressive pruning), the accuracy of the model on the training set tends to decrease. This is because overly aggressive pruning may lead to underfitting the training data. On the other hand, the accuracy on the testing set may initially increase with increasing ccp_alpha as the model generalizes better, but after a certain point, it may start to decrease due to excessive pruning, leading to overfitting. The goal is to find the optimal value of ccp_alpha that maximizes the accuracy on the testing set without sacrificing too much accuracy on the training set, thus achieving good generalization performance.

```
In [56]: 1 #choose the point with Low bias(Low training error) and Low variance (Low  
2 learner=tree.DecisionTreeClassifier (random_state=0,ccp_alpha=0.02)  
3 learner.fit(x_train,y_train)  
4 plt.figure(figsize=(12,8))  
5 tree.plot_tree(learner,rounded=True,filled=True)  
6 plt.show()  
7  
8 print(accuracy_score(y_test,learner.predict(x_test)))
```



0.848780487804878

```
In [57]: 1 learner=tree.DecisionTreeClassifier (random_state=0,ccp_alpha=0.12)
2 learner.fit(x_train,y_train)
3 plt.figure(figsize=(12,8))
4 tree.plot_tree(learner,rounded=True,filled=True)
5 plt.show()
6
7 print(accuracy_score(y_test,learner.predict(x_test)))
```



0.7707317073170732

```
In [58]: 1 learner=tree.DecisionTreeClassifier (random_state=0)
2 learner.fit(x_train,y_train)
3
```

Out[58]: DecisionTreeClassifier(random_state=0)

```
In [59]: 1 from sklearn.model_selection import GridSearchCV
```

```
In [60]: 1 #PRE-pruning
2 grid_param ={"criterion":["gini","entropy"],
3             "splitter":["best","random"],
4             "max_depth":range(2,5,1),
5             "min_samples_leaf":range(1,15,1),
6             "min_samples_split":range(2,20,1)}
7
8 grid_search=GridSearchCV(estimator=learner,param_grid=grid_param,cv=5,n
9 grid_search.fit(x_train,y_train)
```

Out[60]: GridSearchCV(cv=5, estimator=DecisionTreeClassifier(random_state=0), n_jobs=-1,

```

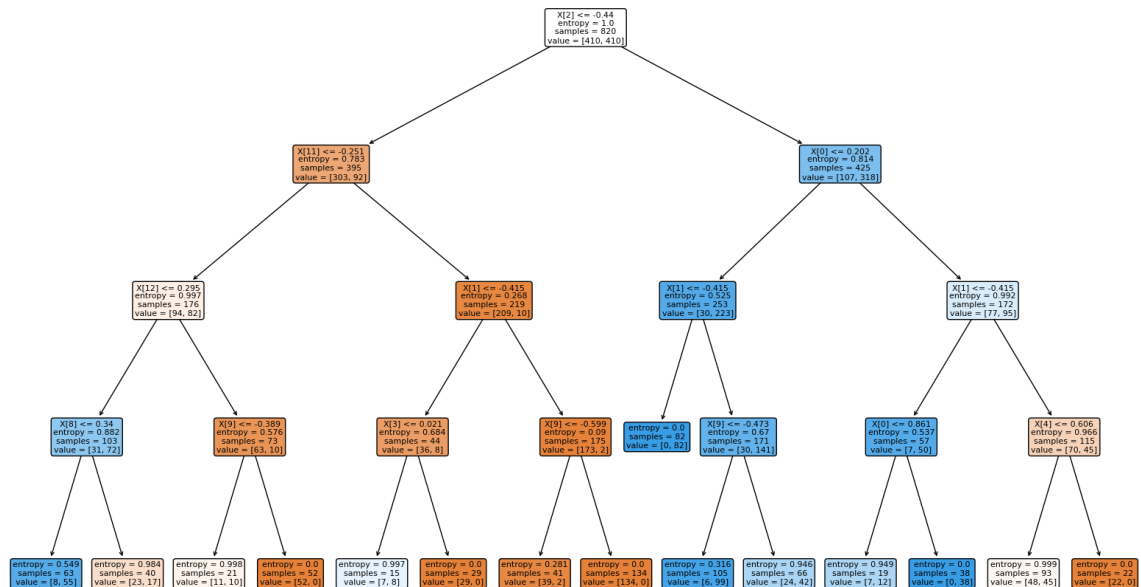
    param_grid={'criterion': ['gini', 'entropy'],
    'max_depth': range(2, 5),
    'min_samples_leaf': range(1, 15),
    'min_samples_split': range(2, 20),
    'splitter': ['best', 'random']})

```

```
In [61]: 1 #Extract best parameters
        2 print(grid_search.best_params_)
```

```
{'criterion': 'gini', 'max_depth': 4, 'min_samples_leaf': 1, 'min_samples_split': 2, 'splitter': 'best'}
```

```
In [62]: 1 learner=tree.DecisionTreeClassifier(criterion= 'entropy', max_depth= 4,
        2 learner.fit(x_train,y_train)
        3 plt.figure(figsize=(20,12))
        4 tree.plot_tree(learner, rounded=True, filled=True)
        5 plt.show()
```



```
In [65]: 1 x = df.drop('target', axis=1) # Features
2 y = df['target'] # Target variable
3
4 # Split the dataset into training and testing sets
5 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2,
6
7 # Standardize features by removing the mean and scaling to unit variance
8 scaler = StandardScaler()
9 x_train = scaler.fit_transform(x_train)
10 x_test = scaler.transform(x_test)
11
12 # Initialize and train the decision tree classifier
13 dt_classifier = DecisionTreeClassifier(random_state=42)
14 dt_classifier.fit(x_train, y_train)
15
16 # Predict on the test set
17 y_pred = dt_classifier.predict(x_test)
18
19 # Evaluate the model
20 accuracy = accuracy_score(y_test, y_pred)
21 print('Accuracy:', accuracy)
22
23 # Classification report
24 print(classification_report(y_test, y_pred))
25
26 # Confusion matrix
27 print('Confusion Matrix:')
28 print(confusion_matrix(y_test, y_pred))
29
```

Accuracy: 0.9853658536585366

	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

Confusion Matrix:

```
[[102  0]
 [ 3 100]]
```

Random forest

```

In [70]: 1 from sklearn.naive_bayes import GaussianNB
2 x = df.drop('target', axis=1) # Features
3 y = df['target'] # Labels
4
5 # Split the data into training and testing sets
6 x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2)
7
8 # Initialize and train the Naive Bayes classifier
9 nb_classifier = GaussianNB()
10 nb_classifier.fit(x_train, y_train)
11 y_pred = nb_classifier.predict(x_test)
12
13 # Evaluate the model
14 print("Accuracy:", accuracy_score(y_test, y_pred))
15 print("Classification Report:")
16 print(classification_report(y_test, y_pred))
17
18

```

Accuracy: 0.8

Classification Report:

	precision	recall	f1-score	support
0	0.87	0.71	0.78	102
1	0.75	0.89	0.82	103
accuracy			0.80	205
macro avg	0.81	0.80	0.80	205
weighted avg	0.81	0.80	0.80	205

In []:

1