

Heart Disease Data EDA and Prediction with Machine Learning Models

Import libraries

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
import os
import yellowbrick
import pickle

from matplotlib.collections import PathCollection
from statsmodels.graphics.gofplots import qqplot
from sklearn.preprocessing import MinMaxScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC
from sklearn.naive_bayes import GaussianNB
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier, GradientBoostingClassifier, AdaBoostClassifier, ExtraTreesClassifier
from sklearn.metrics import classification_report, accuracy_score
from xgboost import XGBClassifier
from yellowbrick.classifier import PrecisionRecallCurve, ROCAUC, ConfusionMatrix
from yellowbrick.style import set_palette
from yellowbrick.model_selection import LearningCurve, FeatureImportances
from yellowbrick.contrib.wrapper import wrap

# --- Libraries Settings ---
warnings.filterwarnings('ignore')
sns.set_style('whitegrid')
plt.rcParams['figure.dpi']=100
set_palette('dark')
```

Color Palettes

This section will create some color palettes that will be used in this notebook.

```
In [3]: # --- Create List of Color Palettes ---
red_grad = ['#FF0000', '#BF0000', '#800000', '#400000', '#000000']
pink_grad = ['#A00030', '#BA1141', '#FF5C8A', '#FF99B9', '#FFDEEB']
purple_grad = ['#4C0028', '#7F0043', '#8E004C', '#A80059', '#C10067']
color_mix = ['#F38BB2', '#FFB9CF', '#FFD7D7', '#F17881', '#E7525B']
black_grad = ['#100C07', '#3E3B39', '#6D6A6A', '#9B9A9C', '#CAC9CD']

# --- Plot Color Palettes --
sns.palplot(red_grad)
sns.palplot(pink_grad)
sns.palplot(purple_grad)
sns.palplot(color_mix)
sns.palplot(black_grad)
```



Read data

```
In [4]: df = pd.read_csv("/input/heart-disease-dataset/heart.csv")
df.head().style.background_gradient(cmap='Reds').set_properties(**{'font-family': 'Segoe UI'}).hide_index()
```

Out[4]:

age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
52	1	0	125	212	0	1	168	0	1.000000	2	2	3	0
53	1	0	140	203	1	0	155	1	3.100000	0	0	3	0
70	1	0	145	174	0	1	125	1	2.600000	0	0	3	0
61	1	0	148	203	0	1	161	0	0.000000	2	1	3	0
62	0	0	138	294	1	1	106	0	1.900000	1	3	2	0

```
In [5]: # --- Print Dataset Info ---
print('\033[1m'+'.: Dataset Info :.'+'\033[0m')
print('*' * 30)
print('Total Rows:'+'\033[1m', df.shape[0])
print('\033[0m'+ 'Total Columns:'+'\033[1m', df.shape[1])
print('\033[0m'*'*' * 30)
print('\n')

# --- Print Dataset Detail ---
print('\033[1m'+'.: Dataset Details :.'+'\033[0m')
print('*' * 30)
df.info(memory_usage = False)
```

```
.: Dataset Info :.
*****
Total Rows: 1025
Total Columns: 14
*****
```



```
.: Dataset Details :.
*****
<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)
```

It is evident that the dataset has been successfully imported. The dataset comprises 14 columns with 303 observations, and there are no null values present. Above, the details of each variable can also be observed.

However, it appears that the data types for some columns do not match. The following steps will address fixing the data types for those columns before performing any analysis.

```
In [6]: # --- Fix Data Types ---
lst=['sex', 'cp', 'fbs', 'restecg', 'exang', 'slope', 'ca', 'thal']
df[lst] = df[lst].astype(object)
```

Data exploration

Sex (Gender)

```
In [7]: # --- Setting Colors, Labels, Order ---
colors=color_mix[2:4]
labels=['Female', 'Male']
order=df['sex'].value_counts().index

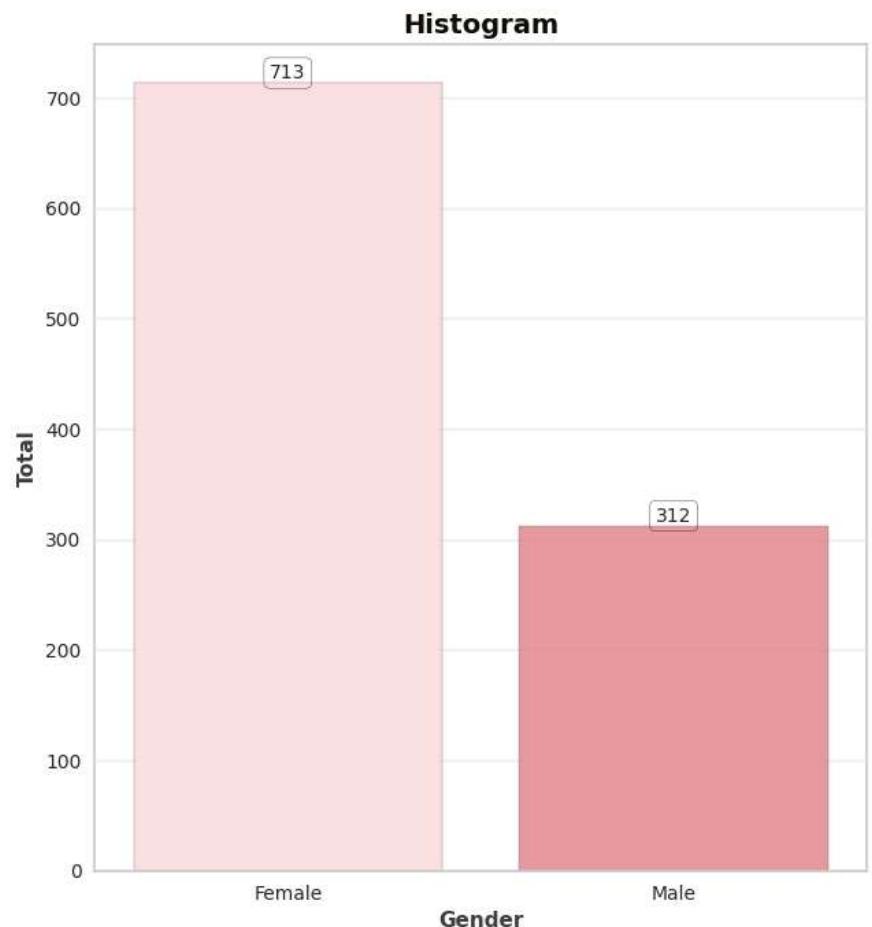
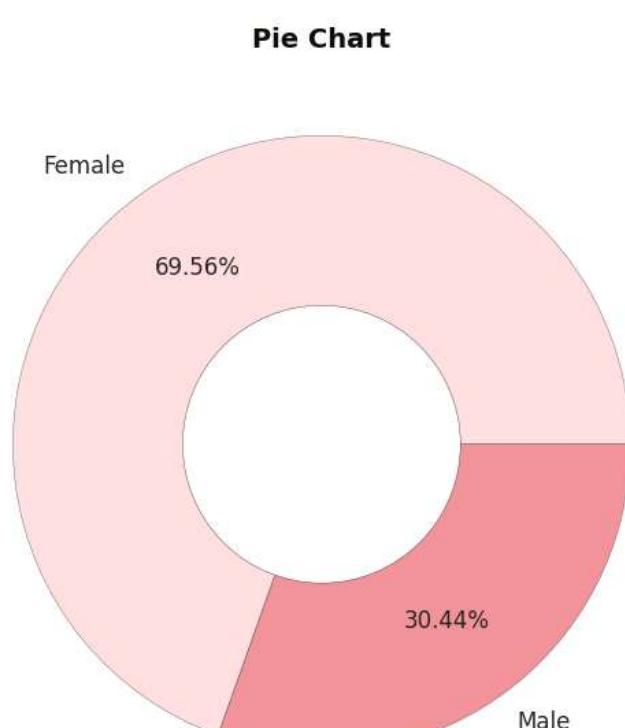
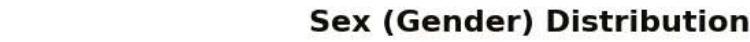
# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Sex (Gender) Distribution', fontweight='heavy',
            fontsize='16', fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1, 2, 1)
plt.title('Pie Chart', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[0])
plt.pie(df['sex'].value_counts(), labels=labels, colors=colors, pctdistance=0.7,
        autopct='%.2f%%', wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]),
        textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)
```

```
# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[0])
ax = sns.countplot(x='sex', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width() / 2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0],
                      linewidth=0.25, boxstyle='round'))
plt.xlabel('Gender', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.xticks([0, 1], labels)
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 25)
print('\u033[1m' + ': Sex (Gender) Total :.' + '\u033[0m')
print('*' * 25)
df.sex.value_counts(dropna=False)

*****
.: Sex (Gender) Total :.
*****
```



cp (Chest Pain Type)

```
In [8]: # --- Setting Colors, Labels, Order ---
colors=pink_grad[0:4]
labels=['Type 0', 'Type 2', 'Type 1', 'Type 3']
order=df['cp'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Chest Pain Type Distribution', fontweight='heavy', fontsize=16,
             fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1, 2, 1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
plt.pie(df['cp'].value_counts(), labels=labels, colors=colors, pctdistance=0.7,
        autopct='%.2f%%', textprops={'fontsize':12},
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]))
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
ax = sns.countplot(x='cp', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width()/2,
```

```

        rect.get_height()+4.25,rect.get_height(),
        horizontalalignment='center', fontsize=10,
        bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                  boxstyle='round')))

plt.xlabel('Pain Type', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.xticks([0, 1, 2, 3], labels)
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 30)
print('\u033[1m'+'.: Chest Pain Type Total :.'+'\u033[0m')
print('*' * 30)
df.cp.value_counts(dropna=False)

*****
.*: Chest Pain Type Total :.*:
*****

```

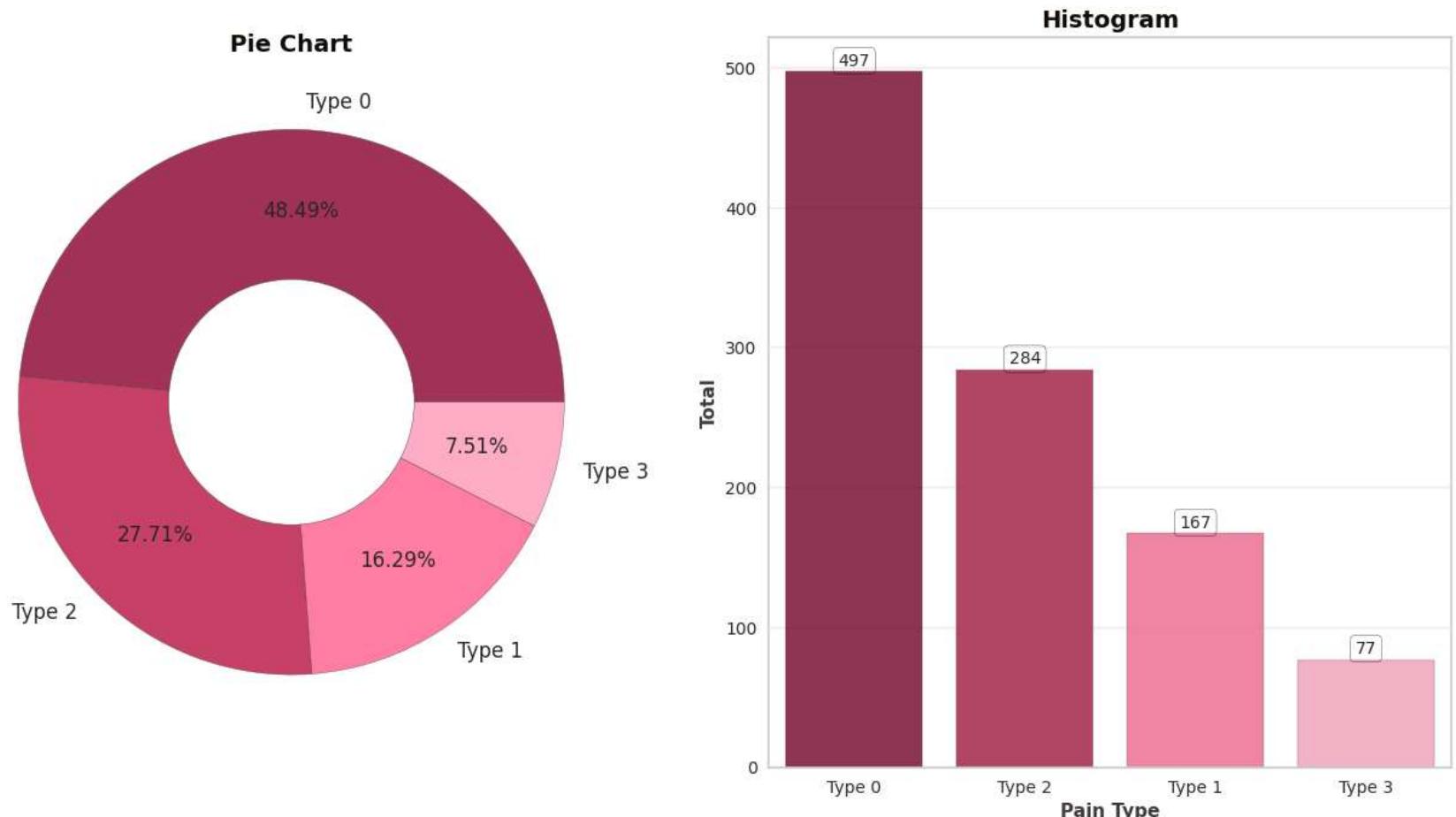
Out[8]:

```

0    497
2    284
1    167
3     77
Name: cp, dtype: int64

```

Chest Pain Type Distribution



fbs (Fasting Blood Sugar)

```

In [9]: # --- Setting Colors, Labels, Order ---
colors=color_mix[0:2]
labels=['< 120 mg/dl', '> 120 mg/dl']
order=df['fbs'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Fasting Blood Sugar Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1, 2, 1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
plt.pie(df['fbs'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]), autopct='%.2f%%',
        pctdistance=0.7, textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
ax = sns.countplot(x='fbs', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width() / 2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

plt.xlabel('Fasting Blood Sugar', fontweight='bold', fontsize=11,

```

```

        fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.xticks([0, 1], labels)
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 32)
print('\033[1m'+'.: Fasting Blood Sugar Total :: '+'\033[0m')
print('*' * 32)
df.fbs.value_counts(dropna=False)

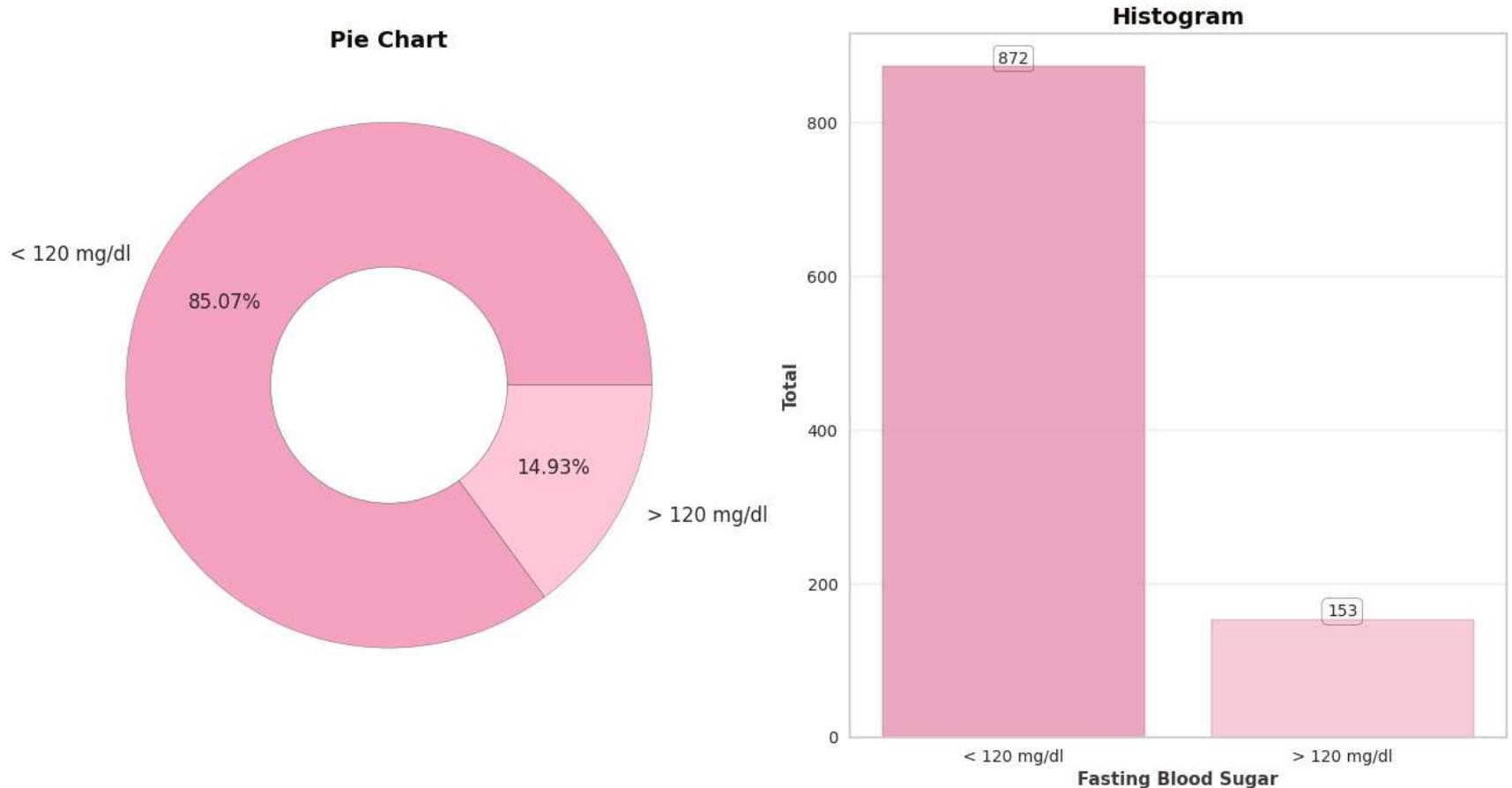
*****
.:: Fasting Blood Sugar Total ::.
*****
```

Out[9]:

	Total
0	872
1	153

Name: fbs, dtype: int64

Fasting Blood Sugar Distribution



restecg (Resting Electrocardiographic Results)

```

In [10]: # --- Setting Colors, Labels, Order ---
colors=pink_grad[1:4]
labels=['1', '0', '2']
order=df['restecg'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Resting Electrocardiographic Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1,2,1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
plt.pie(df['restecg'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]), autopct='%.2f%%',
        pctdistance=0.7, textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
ax = sns.countplot(x='restecg', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width() / 2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

plt.xlabel('Resting Electrocardiographic', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 50)
```

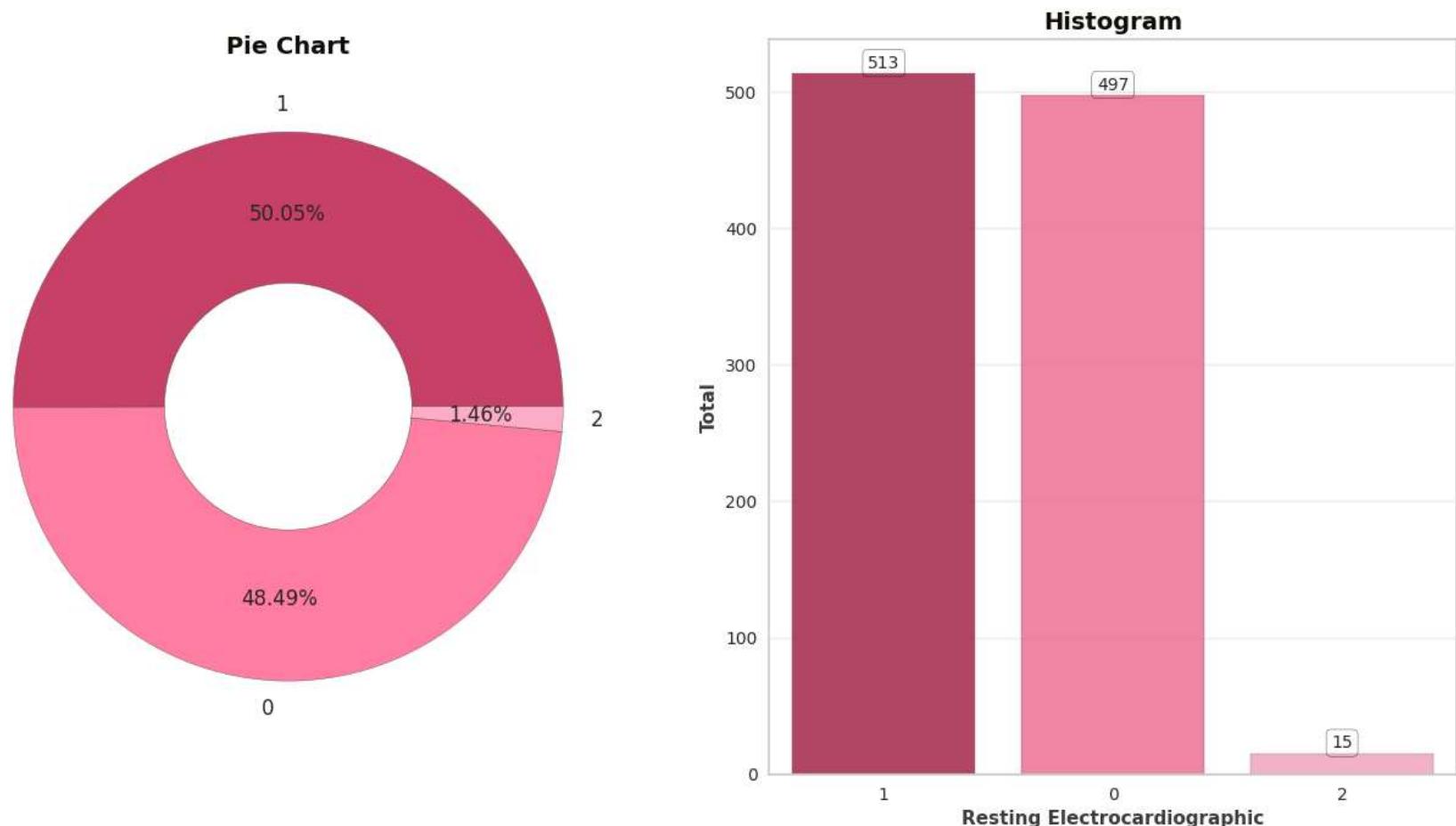
```

print('\033[1m'+'.: Resting Electrocardiographic Results Total :.'+'\033[0m')
print('*' * 50)
df.restecg.value_counts(dropna=False)

*****
.*: Resting Electrocardiographic Results Total :.
*****
Out[10]: 1    513
0    497
2     15
Name: restecg, dtype: int64

```

Resting Electrocardiographic Distribution



exang (Exercise Induced Angina)

```

In [11]: # --- Setting Colors, Labels, Order ---
colors=red_grad[1:3]
labels=['False', 'True']
order=df['exang'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Exercise Induced Angina Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1,2,1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
plt.pie(df['exang'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]), autopct='%.2f%%',
        pctdistance=0.7, textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
ax = sns.countplot(x='exang', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width()/2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

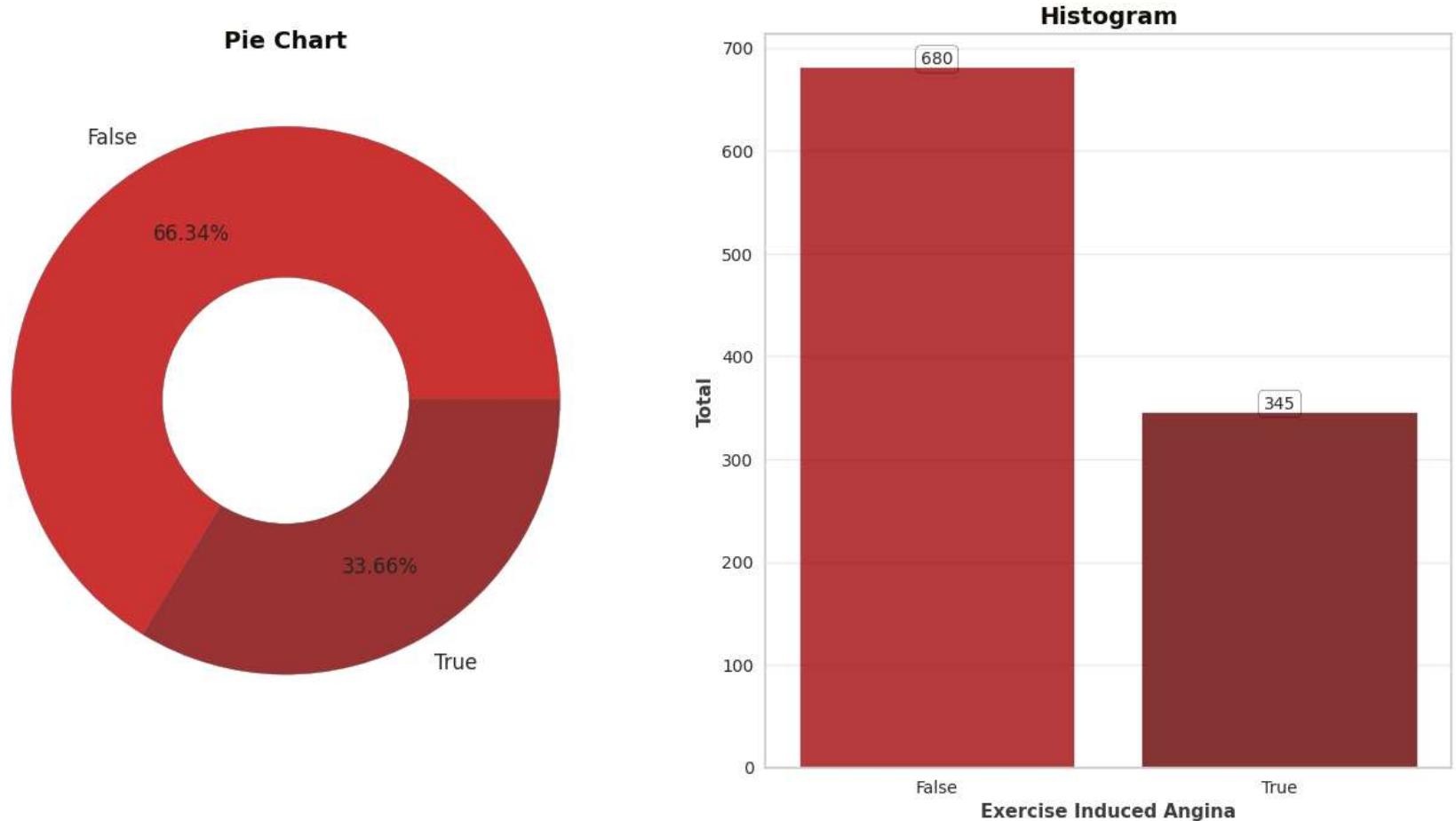
plt.xlabel('Exercise Induced Angina', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.xticks([0, 1], labels)
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 35)
print('\033[1m'+'.: Exercise Induced Angina Total :.'+'\033[0m')
print('*' * 35)
df.exang.value_counts(dropna=False)

*****
.*: Exercise Induced Angina Total :.
*****
```

```
Out[11]: 0    680
         1    345
Name: exang, dtype: int64
```

Exercise Induced Angina Distribution



slope (Slope of the Peak Exercise)

```
In [12]: # --- Setting Colors, Labels, Order ---
colors=purple_grad[2:5]
labels=['2', '1', '0']
order=df['slope'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Slope of the Peak Exercise Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1, 2, 1)
plt.title('Pie Chart', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[0])
plt.pie(df['slope'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]), autopct='%.2f%%',
        pctdistance=0.7, textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
ax = sns.countplot(x='slope', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width() / 2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

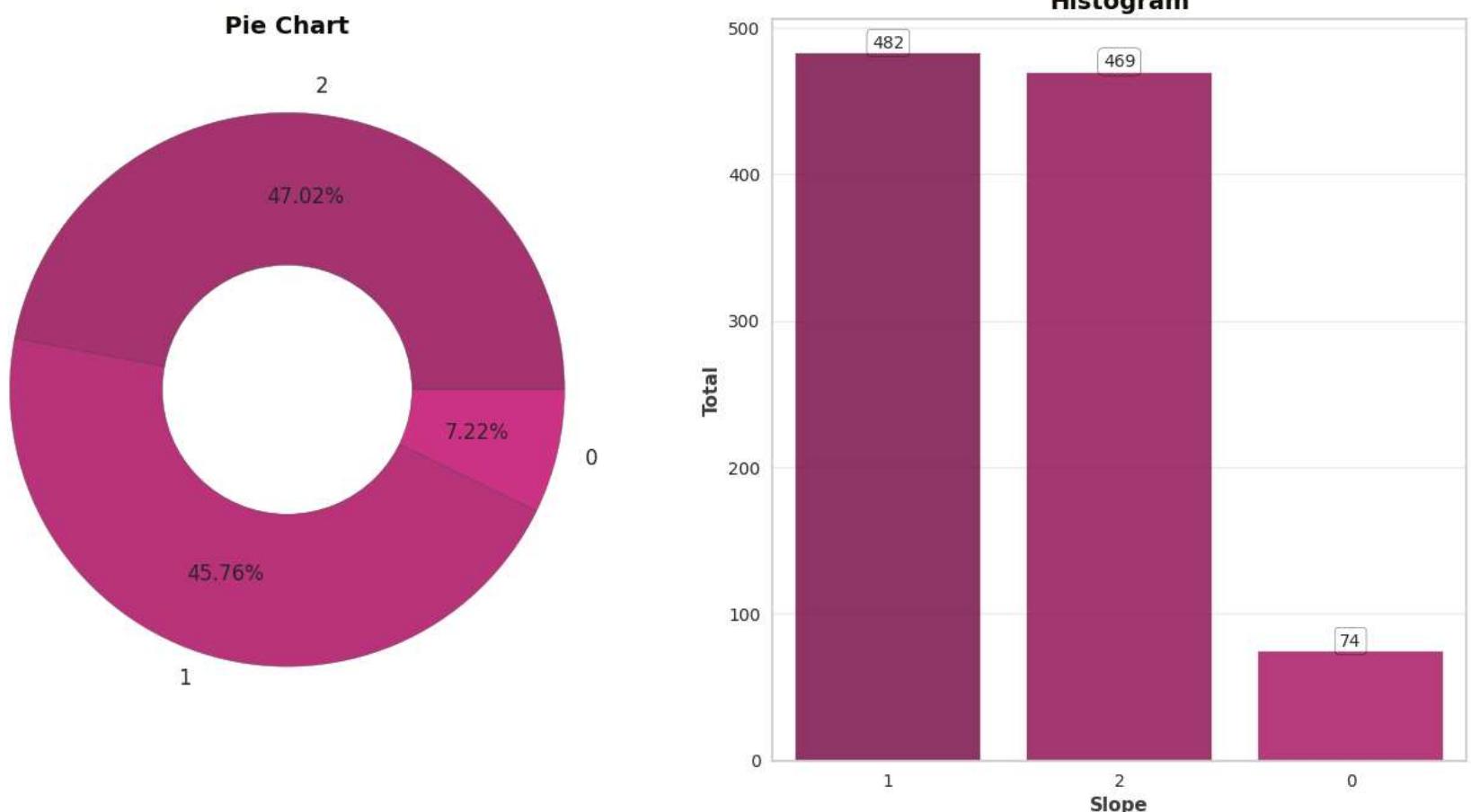
plt.xlabel('Slope', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 20)
print('\u033[1m'+' Slope Total :.'+'\u033[0m')
print('*' * 20)
df.slope.value_counts(dropna=False)
```

.: Slope Total :.

```
Out[12]: 1    482
         2    469
         0     74
Name: slope, dtype: int64
```

Slope of the Peak Exercise Distribution



ca (Number of Major Vessels)

```
In [13]: # --- Setting Colors, Labels, Order ---
colors=purple_grad
labels=['0', '1', '2', '3', '4']
order=df['ca'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16, 8))
plt.suptitle('Number of Major Vessels Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1,2,1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
plt.pie(df['ca'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]),
        autopct='%.2f%%', pctdistance=0.7, textprops={'fontsize':12})

centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
ax = sns.countplot(x='ca', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)

for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width()/2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

plt.xlabel('Number of Major Vessels', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.grid(axis='y', alpha=0.4)
countplt

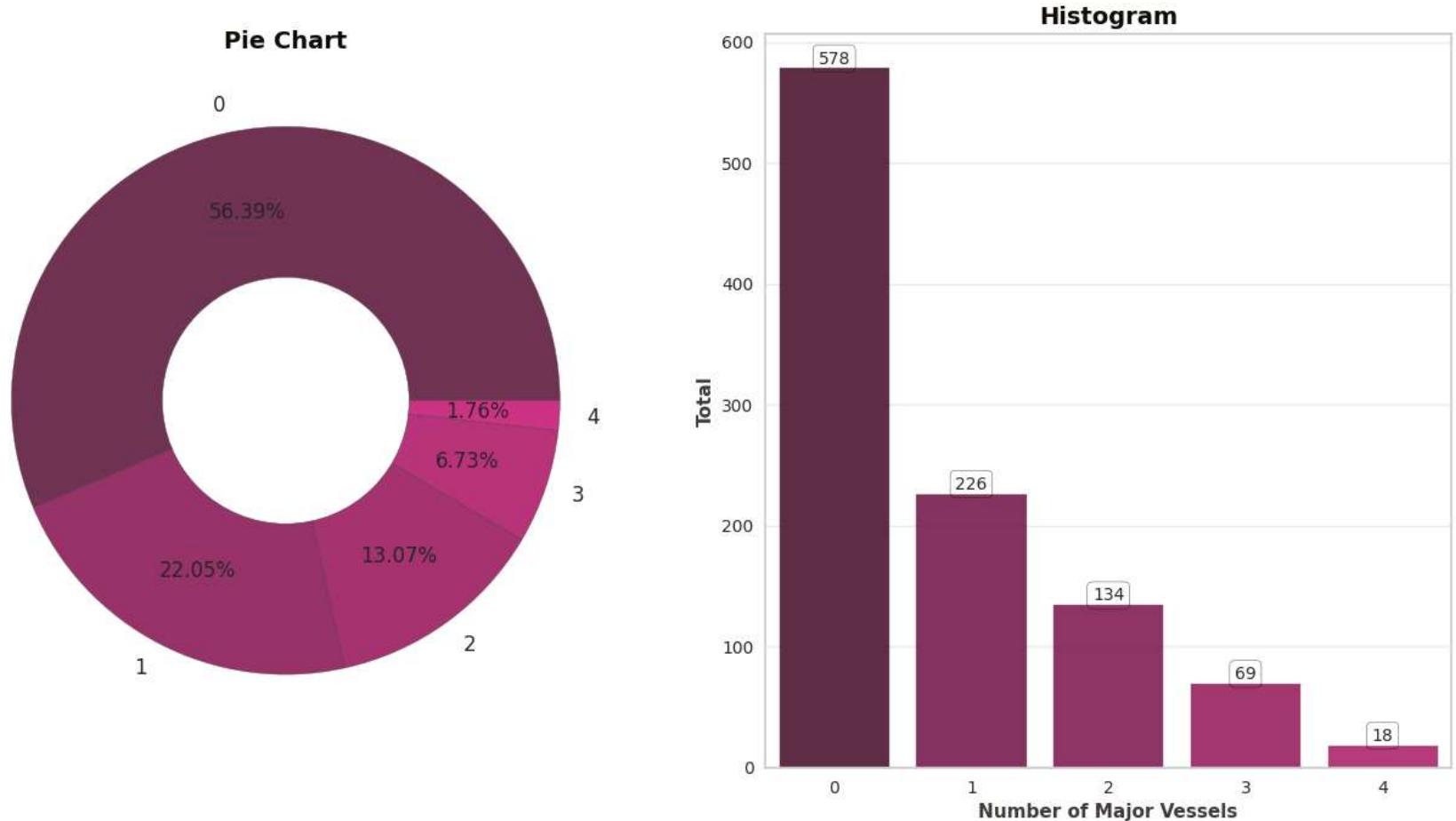
# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 40)
print('\u033[1m'+'.: Number of Major Vessels Total ::.'+'\u033[0m')
print('*' * 40)
df.ca.value_counts(dropna=False)

*****
```

.*: Number of Major Vessels Total :.*

```
Out[13]: 0    578
1    226
2    134
3     69
4     18
Name: ca, dtype: int64
```

Number of Major Vessels Distribution



thal

```
In [14]: # --- Setting Colors, Labels, Order ---
colors=red_grad[0:4]
labels=['2', '3', '1', '0']
order=df['thal'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16,8))
plt.suptitle('"thal" Distribution', fontweight='heavy', fontsize=16,
            fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1,2,1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
plt.pie(df['thal'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]),
        autopct='%.2f%%', pctdistance=0.7, textprops={'fontsize':12})

centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[0])
ax = sns.countplot(x='thal', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)

for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width()/2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

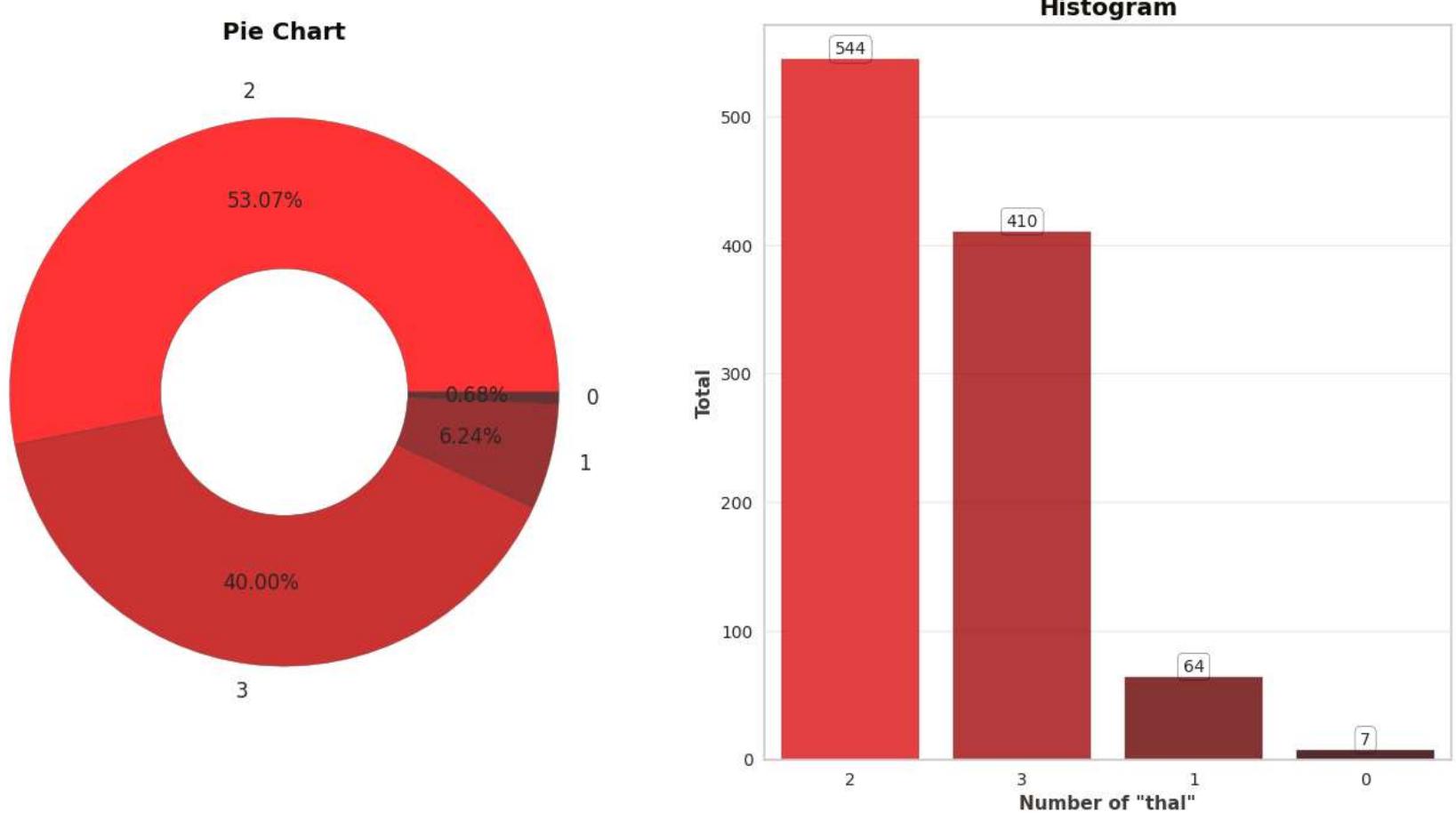
plt.xlabel('Number of "thal"', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 20)
print('\033[1m'+'.' : "thal" Total :.'+'\033[0m')
print('*' * 20)
df.thal.value_counts(dropna=False)
```

```
*****
.: "thal" Total :.
*****
```

```
Out[14]: 2      544
3      410
1      64
0       7
Name: thal, dtype: int64
```

"thal" Distribution



target (Heart Diseases Status)

```
In [15]: # --- Setting Colors, Labels, Order ---
colors=color_mix[3:5]
labels=['True', 'False']
order=df['target'].value_counts().index

# --- Size for Both Figures ---
plt.figure(figsize=(16,8))
plt.suptitle('Heart Diseases Distribution', fontweight='heavy',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])

# --- Pie Chart ---
plt.subplot(1, 2, 1)
plt.title('Pie Chart', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
plt.pie(df['target'].value_counts(), labels=labels, colors=colors,
        wedgeprops=dict(alpha=0.8, edgecolor=black_grad[1]), autopct='%.2f%%',
        pctdistance=0.7, textprops={'fontsize':12})
centre=plt.Circle((0, 0), 0.45, fc='white', edgecolor=black_grad[1])
plt.gcf().gca().add_artist(centre)

# --- Histogram ---
countplt = plt.subplot(1, 2, 2)
plt.title('Histogram', fontweight='bold', fontsize=14, fontfamily='sans-serif',
           color=black_grad[0])
ax = sns.countplot(x='target', data=df, palette=colors, order=order,
                    edgecolor=black_grad[2], alpha=0.85)
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width()/2,
            rect.get_height() + 4.25, rect.get_height(),
            horizontalalignment='center', fontsize=10,
            bbox=dict(facecolor='none', edgecolor=black_grad[0], linewidth=0.25,
                      boxstyle='round'))

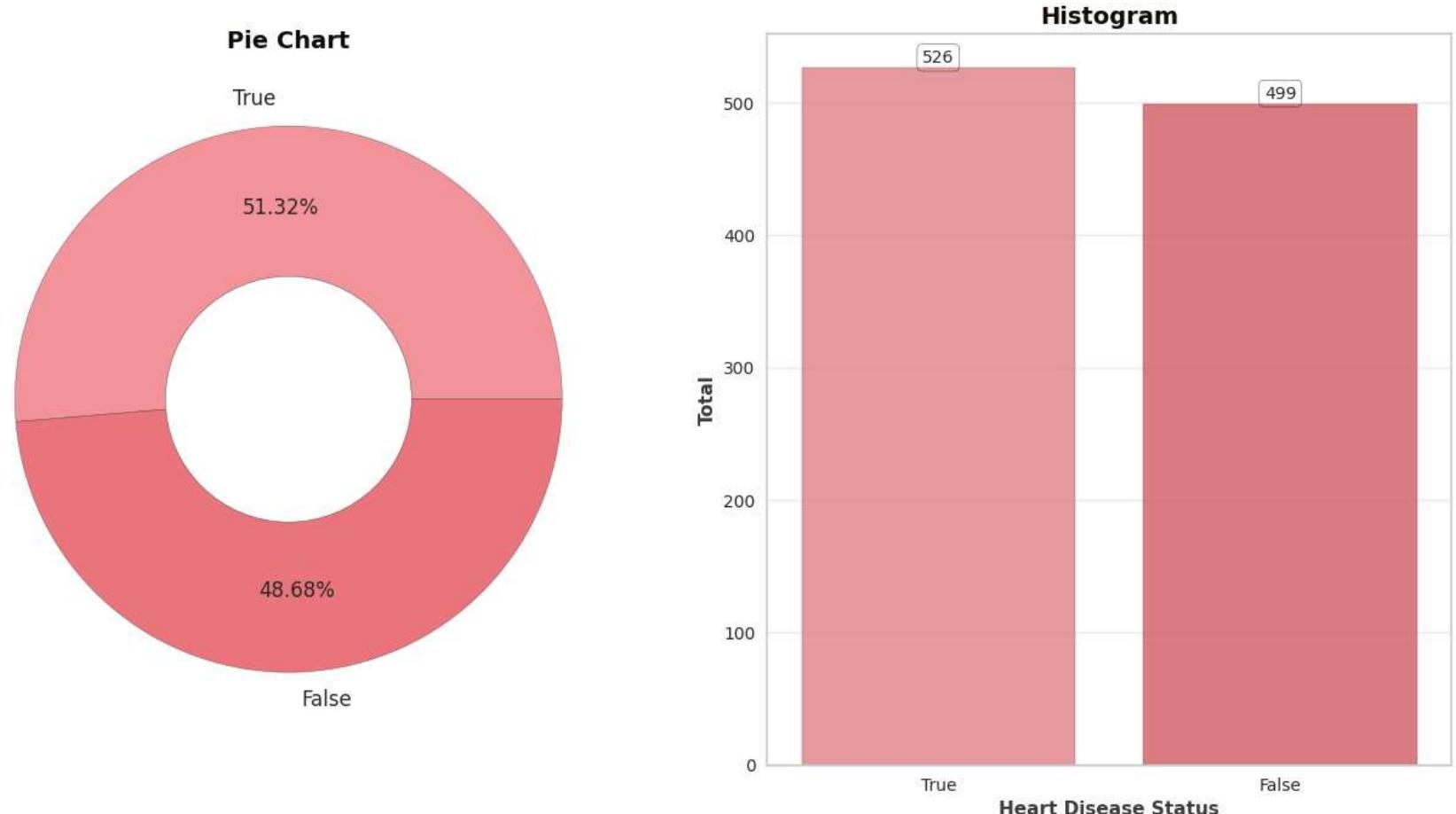
plt.xlabel('Heart Disease Status', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Total', fontweight='bold', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])
plt.xticks([0, 1], labels)
plt.grid(axis='y', alpha=0.4)
countplt

# --- Count Categorical Labels w/out Dropping Null Values ---
print('*' * 45)
print('\033[1m'+'.: Heart Diseases Status (target) Total :.'+'\033[0m')
print('*' * 45)
df.target.value_counts(dropna=False)

*****
.*: Heart Diseases Status (target) Total :.
*****
```

```
Out[15]: 1    526
          0    499
Name: target, dtype: int64
```

Heart Diseases Distribution



Numerical Variable

The second variable that will be explored is numerical variable

Descriptive Statistics

	count	mean	std	min	25%	50%	75%	max
age	1025.000000	54.434146	9.072290	29.000000	48.000000	56.000000	61.000000	77.000000
trestbps	1025.000000	131.611707	17.516718	94.000000	120.000000	130.000000	140.000000	200.000000
chol	1025.000000	246.000000	51.592510	126.000000	211.000000	240.000000	275.000000	564.000000
thalach	1025.000000	149.114146	23.005724	71.000000	132.000000	152.000000	166.000000	202.000000
oldpeak	1025.000000	1.071512	1.175053	0.000000	0.000000	0.800000	1.800000	6.200000
target	1025.000000	0.513171	0.500070	0.000000	0.000000	1.000000	1.000000	1.000000

Continuous Column Distribution

This section will show the distribution of numerical variables in histograms, boxplots, Q-Q Plots, skewness and kurtosis values.

```
In [17]: # --- Variable, Color & Plot Size ---
var = 'age'
color = color_mix[0]
fig=plt.figure(figsize=(12, 12))

# --- Skewness & Kurtosis ---
print('\u033[1m'+'.'*40)
print('Age Column Skewness & Kurtosis :.'+'\033[0m')
print('*' * 40)
print('Skewness:'+'\033[1m {:.3f}'.format(df[var].skew(axis = 0, skipna = True)))
print('\033[0m'+ 'Kurtosis:'+'\033[1m {:.3f}'.format(df[var].kurt(axis = 0, skipna = True)))
print('\n')

# --- General Title ---
fig.suptitle('Age Column Distribution', fontweight='bold', fontsize=16,
            fontfamily='sans-serif', color=black_grad[0])
fig.subplots_adjust(top=0.9)

# --- Histogram ---
ax_1=fig.add_subplot(2, 2, 2)
plt.title('Histogram Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.histplot(data=df, x=var, kde=True, color=color)
plt.xlabel('Total', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Age', fontweight='regular', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])

# --- Q-Q Plot ---
ax_2=fig.add_subplot(2, 2, 4)
plt.title('Q-Q Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
qqplot(df[var], fit=True, line='45', ax=ax_2, markerfacecolor=color,
       markeredgecolor=color, alpha=0.6)
plt.xlabel('Theoretical Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Sample Quantiles', fontweight='regular', fontsize=11,
```

```

        fontfamily='sans-serif', color=black_grad[1])

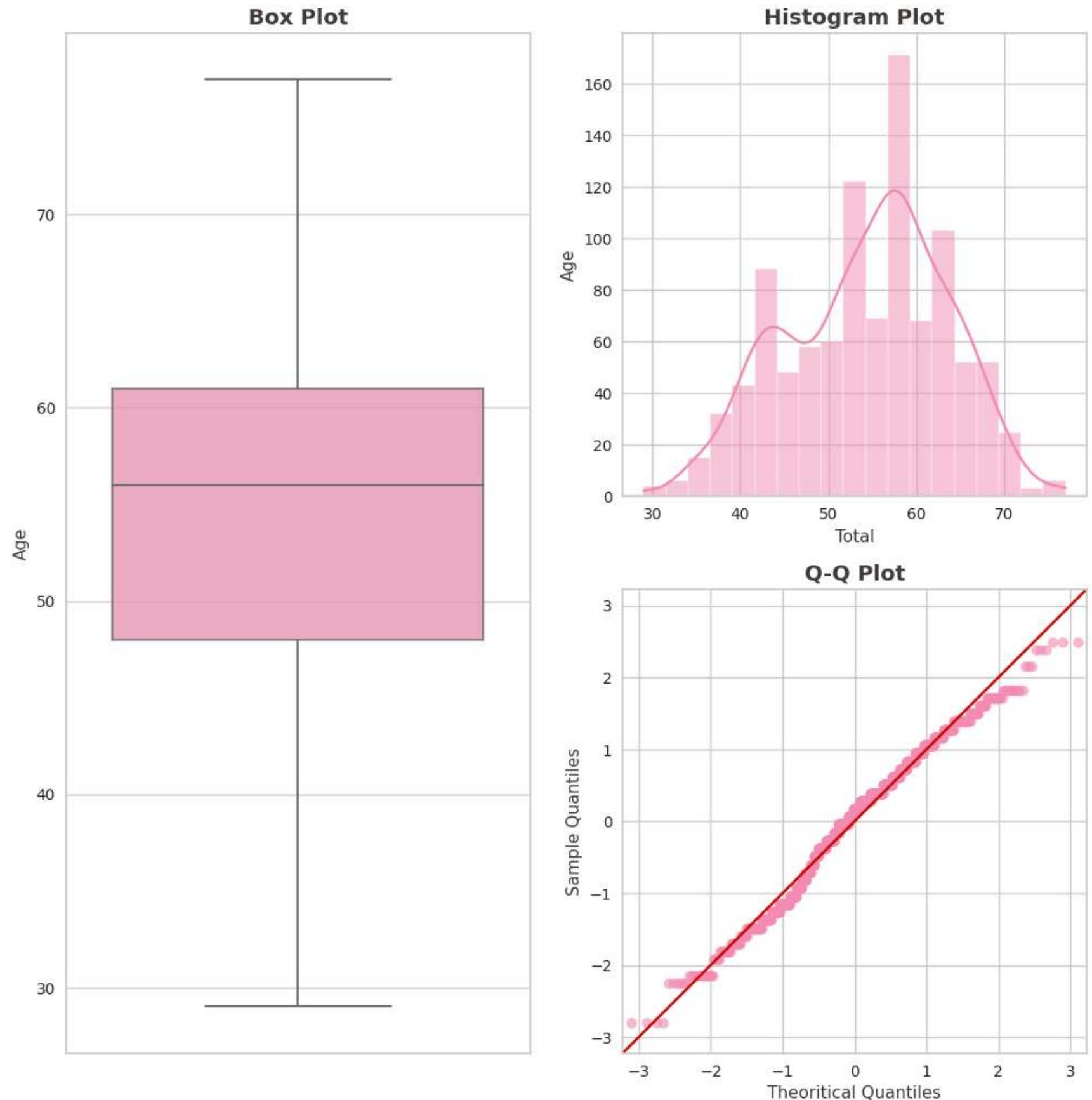
# --- Box Plot ---
ax_3=plt.subplot(1, 2, 1)
plt.title('Box Plot', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[1])
sns.boxplot(data=df, y=var, color=color, boxprops=dict(alpha=0.8), linewidth=1.5)
plt.ylabel('Age', fontweight='regular', fontsize=11, fontfamily='sans-serif',
           color=black_grad[1])

plt.show()

.: Age Column Skewness & Kurtosis :.
*****
Skewness: -0.249
Kurtosis: -0.526

```

Age Column Distribution



trestbps (Resting Blood Pressure in mm Hg)

```

In [18]: # --- Variable, Color & Plot Size ---
var = 'trestbps'
color = color_mix[2]
fig=plt.figure(figsize=(12, 12))

# --- Skewness & Kurtosis ---
print('\u033[1m '+'.: Resting Blood Pressure Column Skewness & Kurtosis :.'+'\u033[0m')
print('*' * 55)
print('Skewness:'+'\u033[1m {:.3f}'.format(df[var].skew(axis = 0, skipna = True)))
print('Kurtosis:'+'\u033[1m {:.3f}'.format(df[var].kurt(axis = 0, skipna = True)))
print('\n')

# --- General Title ---
fig.suptitle('Resting Blood Pressure Column Distribution', fontweight='bold',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])
fig.subplots_adjust(top=0.9)

```

```

# --- Histogram ---
ax_1=fig.add_subplot(2, 2, 2)
plt.title('Histogram Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.histplot(data=df, x=var, kde=True, color=color)
plt.xlabel('Total', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Resting Blood Pressure', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Q-Q Plot ---
ax_2=fig.add_subplot(2, 2, 4)
plt.title('Q-Q Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
qqplot(df[var], fit=True, line='45', ax=ax_2, markerfacecolor=color,
       markeredgecolor=color, alpha=0.6)
plt.xlabel('Theoretical Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Sample Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Box Plot ---
ax_3=fig.add_subplot(1, 2, 1)
plt.title('Box Plot', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[1])
sns.boxplot(data=df, y=var, color=color, boxprops=dict(alpha=0.8), linewidth=1.5)
plt.ylabel('Resting Blood Pressure', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

plt.show()

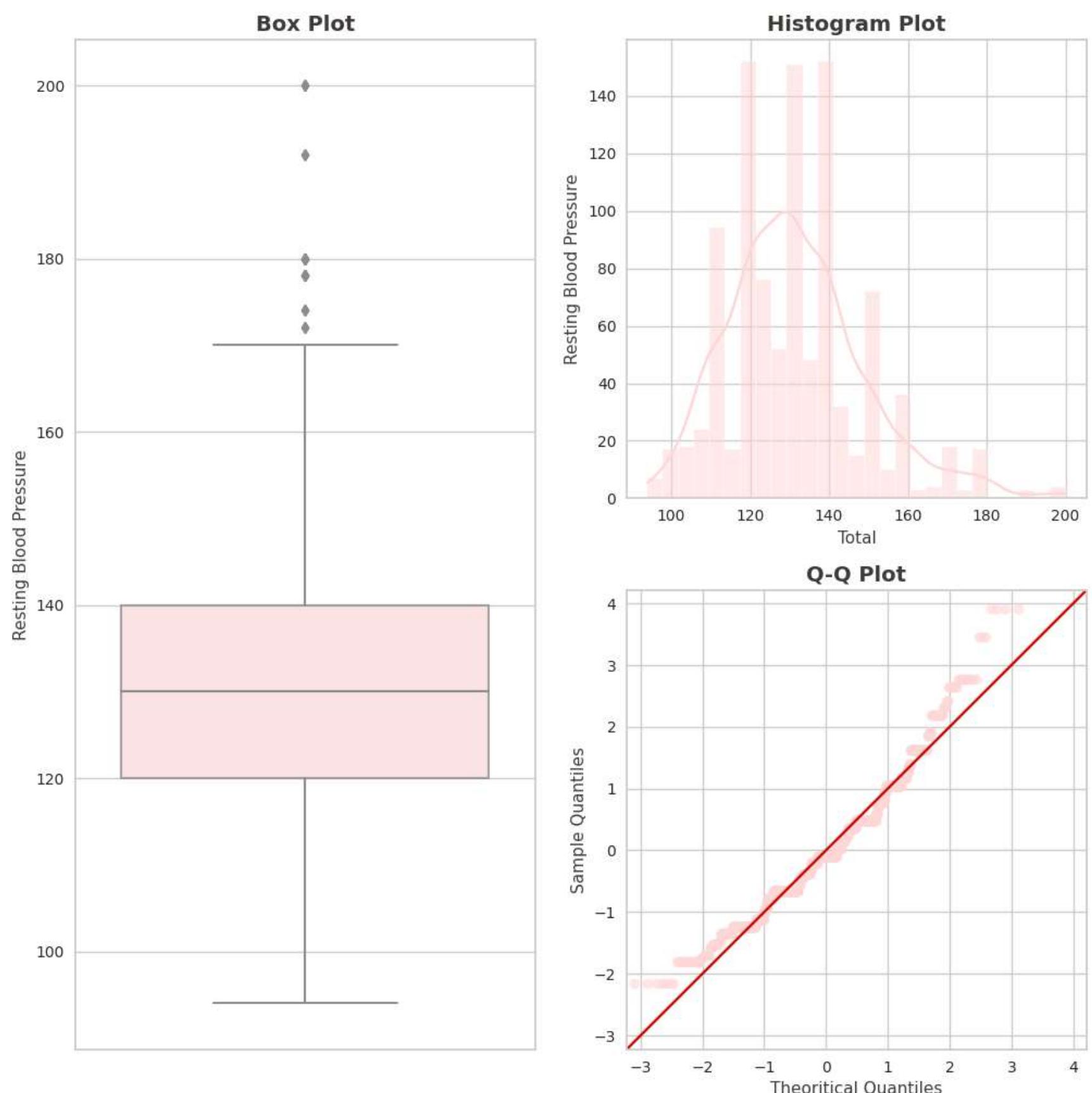
```

.: Resting Blood Pressure Column Skewness & Kurtosis :.

Skewness: 0.740

Kurtosis: 0.991

Resting Blood Pressure Column Distribution



chol (Serum Cholestorol in mg/dl)

```
In [19]: # --- Variable, Color & Plot Size ---
var = 'chol'
color = color_mix[4]
fig=plt.figure(figsize=(12, 12))

# --- Skewness & Kurtosis ---
print(':: Serum Cholestorol Column Skewness & Kurtosis ::')
print('*' * 45)
print('Skewness:' + '{:.3f}'.format(df[var].skew(axis = 0, skipna = True)))
print('Kurtosis:' + '{:.3f}'.format(df[var].kurt(axis = 0, skipna = True)))
print('\n')

# --- General Title ---
fig.suptitle('Serum Cholestorol Column Distribution', fontweight='bold',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])
fig.subplots_adjust(top=0.9)

# --- Histogram ---
ax_1=fig.add_subplot(2, 2, 2)
plt.title('Histogram Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.histplot(data=df, x=var, kde=True, color=color)
plt.xlabel('Total', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Serum Cholestorol', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
# --- Q-Q Plot ---
ax_2=fig.add_subplot(2, 2, 4)
plt.title('Q-Q Plot', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[1])
qqplot(df[var], fit=True, line='45', ax=ax_2, markerfacecolor=color,
       markeredgecolor=color, alpha=0.6)
plt.xlabel('Theoretical Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Sample Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

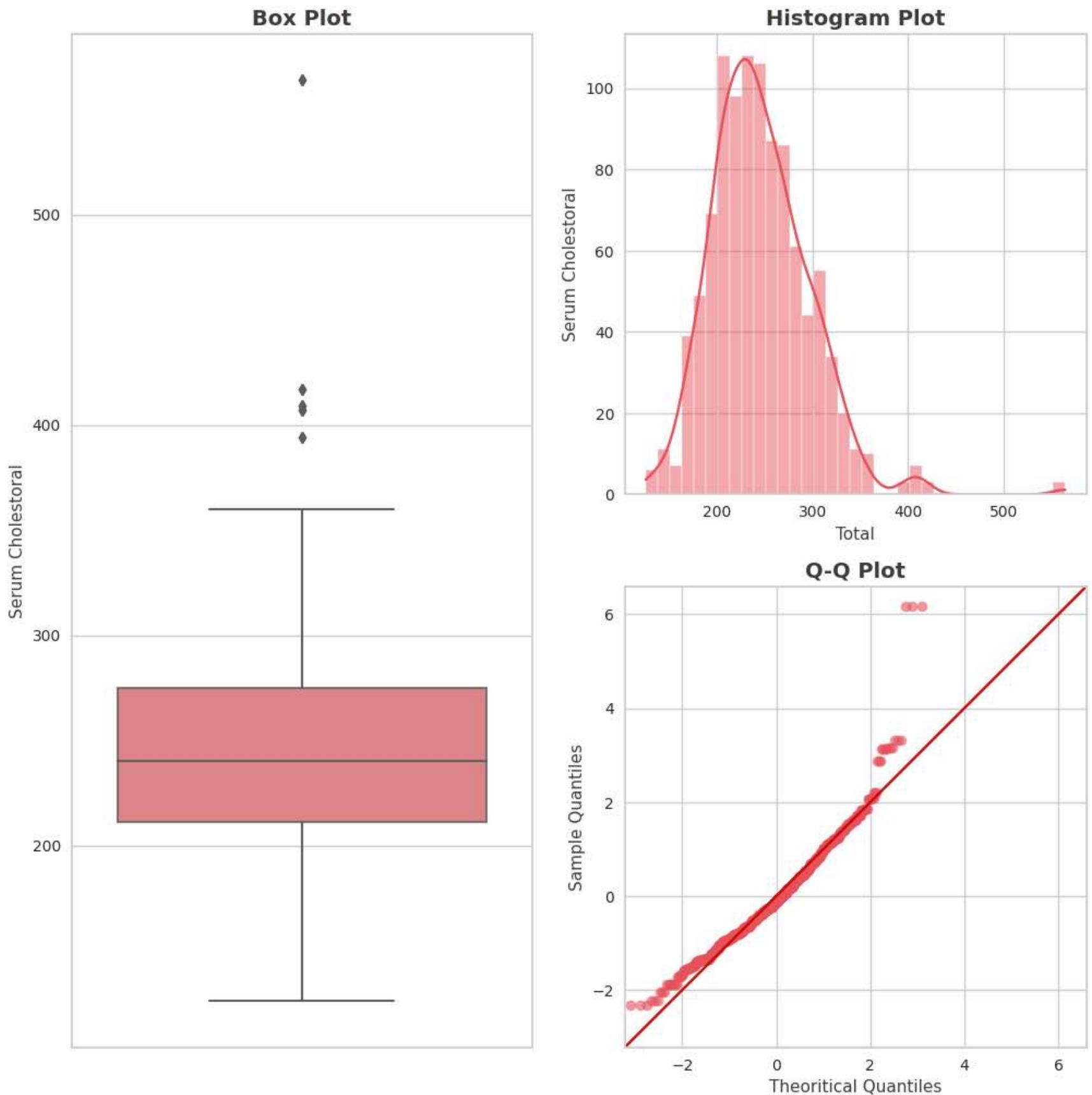
# --- Box Plot ---
ax_3=fig.add_subplot(1, 2, 1)
plt.title('Box Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.boxplot(data=df, y=var, color=color, boxprops=dict(alpha=0.8), linewidth=1.5)
plt.ylabel('Serum Cholestorol', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

plt.show()

.: Serum Cholestorol Column Skewness & Kurtosis ::
```

Skewness: 1.074
Kurtosis: 3.997

Serum Cholestral Column Distribution



thalach (Maximum Heart Rate)

```
In [20]: # --- Variable, Color & Plot Size ---
var = 'thalach'
color = purple_grad[1]
fig=plt.figure(figsize=(12, 12))

# --- Skewness & Kurtosis ---
print('\033[1m'+' Maximum Heart Rate Column Skewness & Kurtosis :.'+'\033[0m')
print('*' * 50)
print('Skewness:'+'\033[1m {:.3f}'.format(df[var].skew(axis = 0, skipna = True)))
print('\033[0m'+Kurtosis:'+'\033[1m {:.3f}'.format(df[var].kurt(axis = 0, skipna = True)))
print('\n')

# --- General Title ---
fig.suptitle('Maximum Heart Rate Column Distribution', fontweight='bold',
             fontsize=16, fontfamily='sans-serif', color=black_grad[0])
fig.subplots_adjust(top=0.9)

# --- Histogram ---
ax_1=fig.add_subplot(2, 2, 2)
plt.title('Histogram Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.histplot(data=df, x=var, kde=True, color=color)
plt.xlabel('Total', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Maximum Heart Rate', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Q-Q Plot ---
ax_2=fig.add_subplot(2, 2, 4)
plt.title('Q-Q Plot', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[1])
qqplot(df[var], fit=True, line='45', ax=ax_2, markerfacecolor=color,
       markeredgecolor=color, alpha=0.6)
plt.xlabel('Theoretical Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
```

```

plt.ylabel('Sample Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Box Plot ---
ax_3=plt.add_subplot(1, 2, 1)
plt.title('Box Plot', fontweight='bold', fontsize=14,
           fontfamily='sans-serif', color=black_grad[1])
sns.boxplot(data=df, y=var, color=color, boxprops=dict(alpha=0.8), linewidth=1.5)
plt.ylabel('Maximum Heart Rate', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

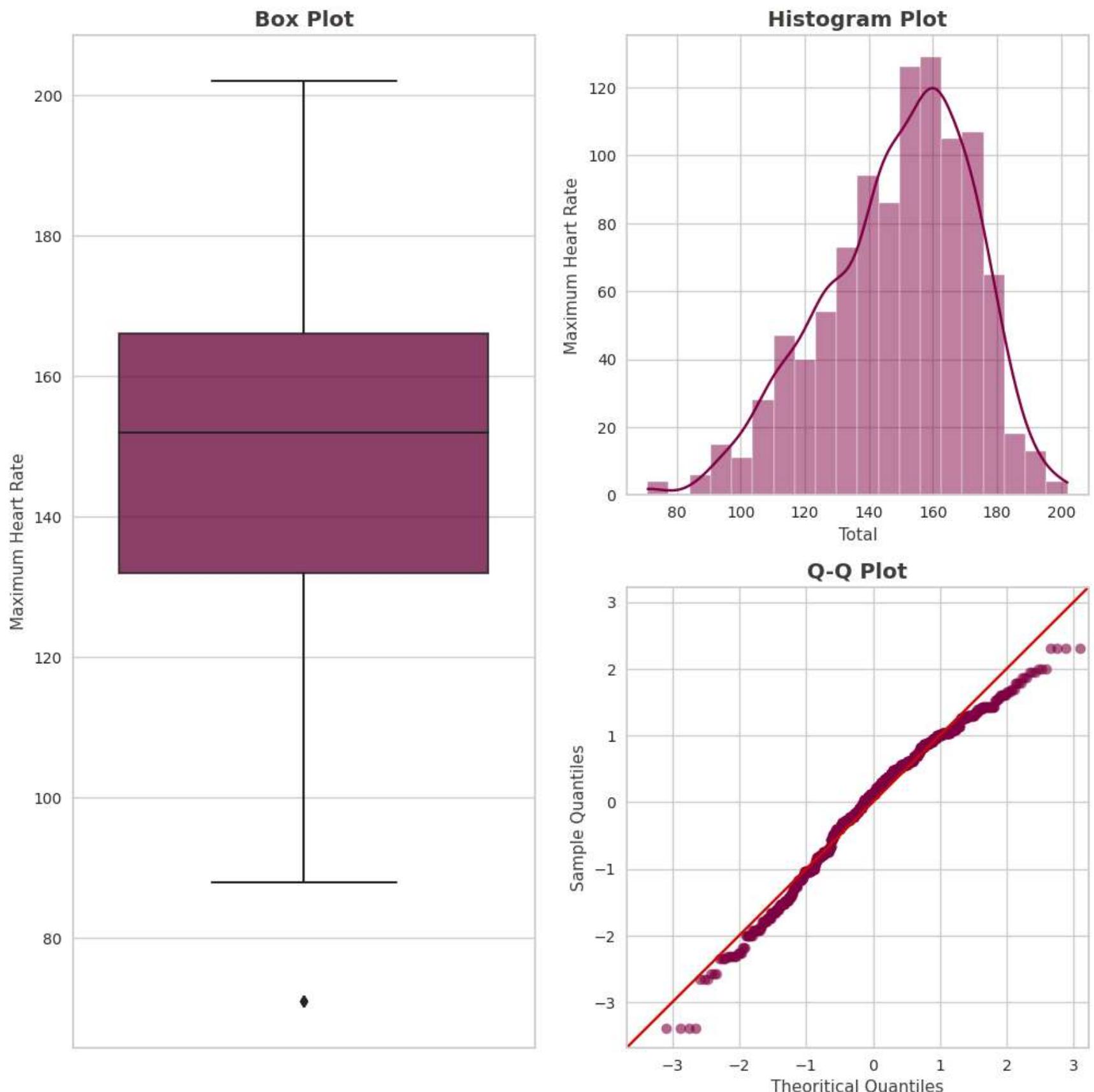
plt.show()

```

.: Maximum Heart Rate Column Skewness & Kurtosis :.

Skewness: -0.514
Kurtosis: -0.089

Maximum Heart Rate Column Distribution



oldpeak

```

In [21]: # --- Variable, Color & Plot Size ---
var = 'oldpeak'
color = red_grad[1]
fig=plt.figure(figsize=(12, 12))

# --- Skewness & Kurtosis ---
print('\u033[1m'+'.'*40 + "oldpeak" Column Skewness & Kurtosis :.'+'\033[0m')
print('*' * 40)
print('Skewness:'+'.'*40.format(df[var].skew(axis = 0, skipna = True)))
print('Kurtosis:'+'.'*40.format(df[var].kurt(axis = 0, skipna = True)))
print('\n')

# --- General Title ---
fig.suptitle('"oldpeak" Column Distribution', fontweight='bold',
            fontsize=16, fontfamily='sans-serif', color=black_grad[0])
fig.subplots_adjust(top=0.9)

```

```

# --- Histogram ---
ax_1=fig.add_subplot(2, 2, 2)
plt.title('Histogram Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.histplot(data=df, x=var, kde=True, color=color)
plt.xlabel('Total', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('oldpeak', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Q-Q Plot ---
ax_2=fig.add_subplot(2, 2, 4)
plt.title('Q-Q Plot', fontweight='bold', fontsize=14, fontfamily='sans-serif',
          color=black_grad[1])
qqplot(df[var], fit=True, line='45', ax=ax_2, markerfacecolor=color,
       markeredgecolor=color, alpha=0.6)
plt.xlabel('Theoretical Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Sample Quantiles', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

# --- Box Plot ---
ax_3=fig.add_subplot(1, 2, 1)
plt.title('Box Plot', fontweight='bold', fontsize=14,
          fontfamily='sans-serif', color=black_grad[1])
sns.boxplot(data=df, y=var, color=color, boxprops=dict(alpha=0.8), linewidth=1.5)
plt.ylabel('oldpeak', fontweight='regular', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])

plt.show()

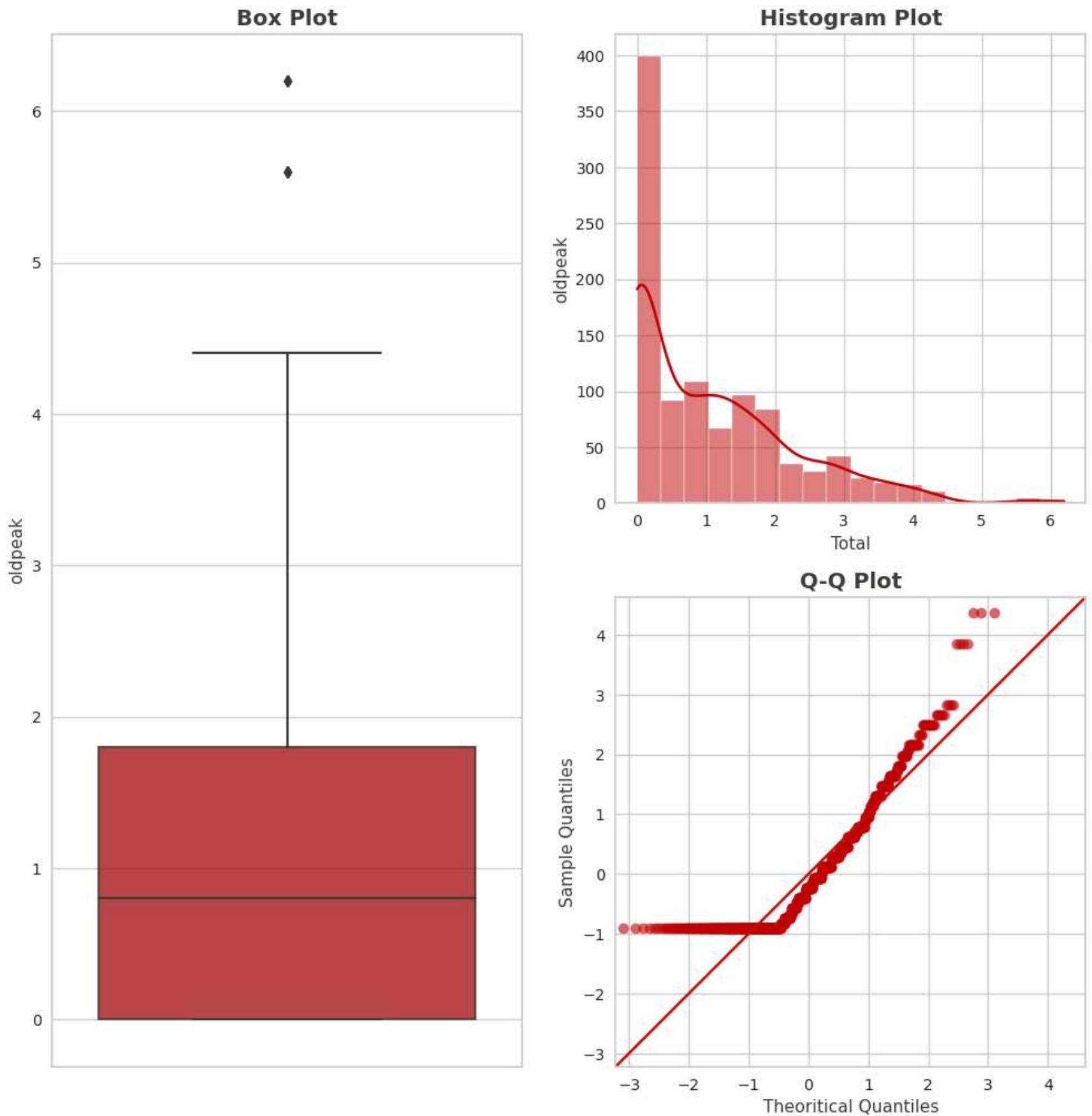
```

.: "oldpeak" Column Skewness & Kurtosis :.

Skewness: 1.211

Kurtosis: 1.314

"oldpeak" Column Distribution



EDA

This section will perform some EDA to get more insights about dataset.

Heart Disease Distribution based on Gender

```
In [22]: # --- Labels Settings ---
labels = ['False', 'True']
label_gender = np.array([0, 1])
label_gender2 = ['Male', 'Female']

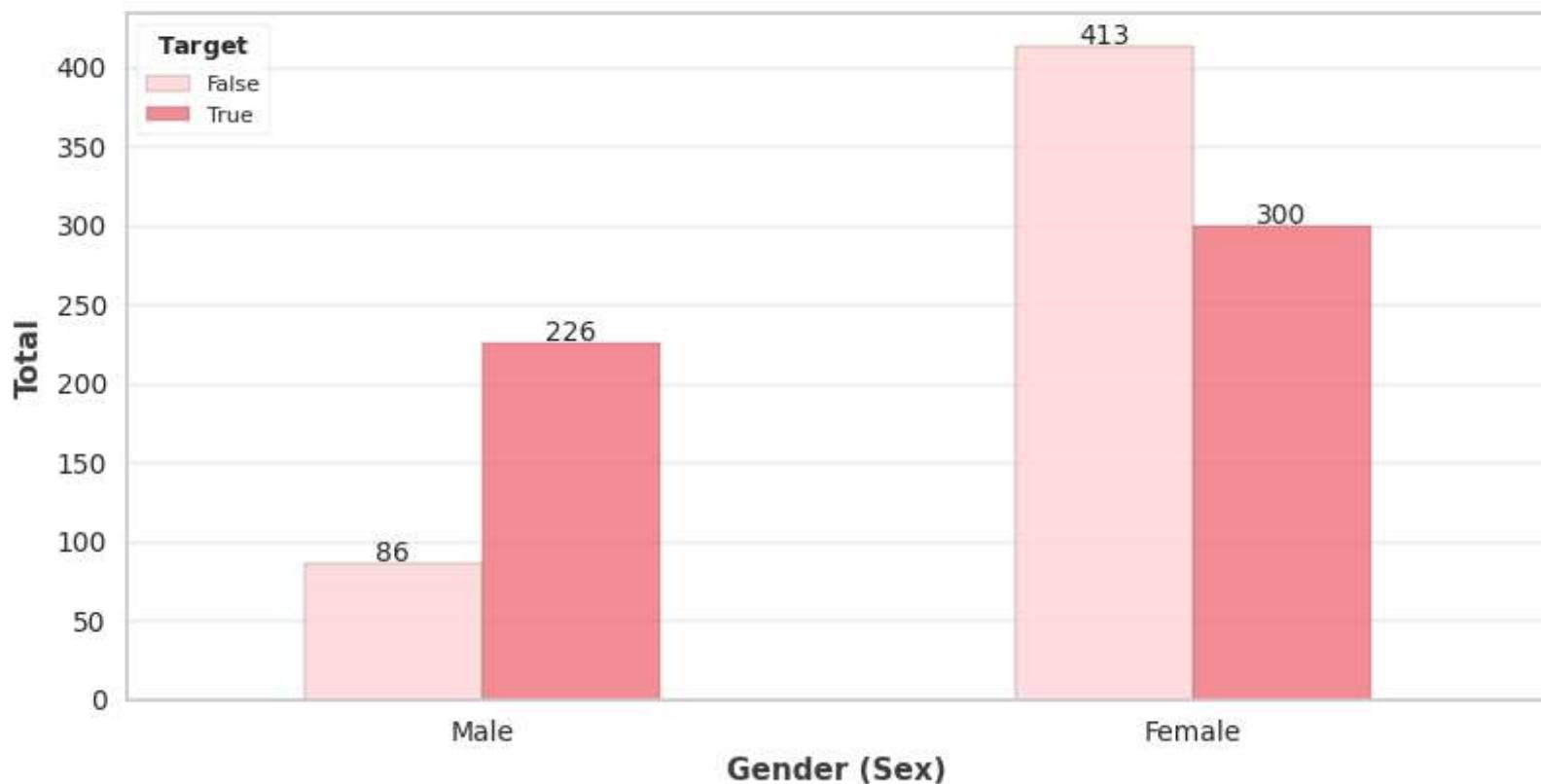
# --- Creating Bar Chart ---
ax = pd.crosstab(df.sex, df.target).plot(kind='bar', figsize=(8, 5),
                                             color=color_mix[2:4],
                                             edgecolor=black_grad[2], alpha=0.85)

# --- Bar Chart Settings ---
for rect in ax.patches:
    ax.text(rect.get_x() + rect.get_width() / 2,
            rect.get_height() + 1.25, rect.get_height(),
            horizontalalignment='center', fontsize=10)

plt.suptitle('Heart Disease Distribution based on Gender', fontweight='heavy',
             x=0.065, y=0.98, ha='left', fontsize=16, fontfamily='sans-serif',
             color=black_grad[0])
plt.title('Female tend to have heart diseases compared to Male. In male, the distribution is not imbalanced compared to',
          fontsize=8, fontfamily='sans-serif', loc='left', color=black_grad[1])
plt.tight_layout(rect=[0, 0.04, 1, 1.025])
plt.xlabel('Gender (Sex)', fontfamily='sans-serif', fontweight='bold',
           color=black_grad[1])
plt.ylabel('Total', fontfamily='sans-serif', fontweight='bold',
           color=black_grad[1])
plt.xticks(label_gender, label_gender2, rotation=0)
plt.grid(axis='y', alpha=0.4)
plt.grid(axis='x', alpha=0)
plt.legend(labels=labels, title='$\bf{Target}$', fontsize=8,
           title_fontsize=9, loc='upper left', frameon=True);
```

Heart Disease Distribution based on Gender

Female tend to have heart diseases compared to Male. In male, the distribution is not imbalanced compared to female that have almost the same distribution



Heart Disease Distribution based on Major Vessels Total

```
In [23]: # --- Labels Settings ---
labels = ['False', 'True']

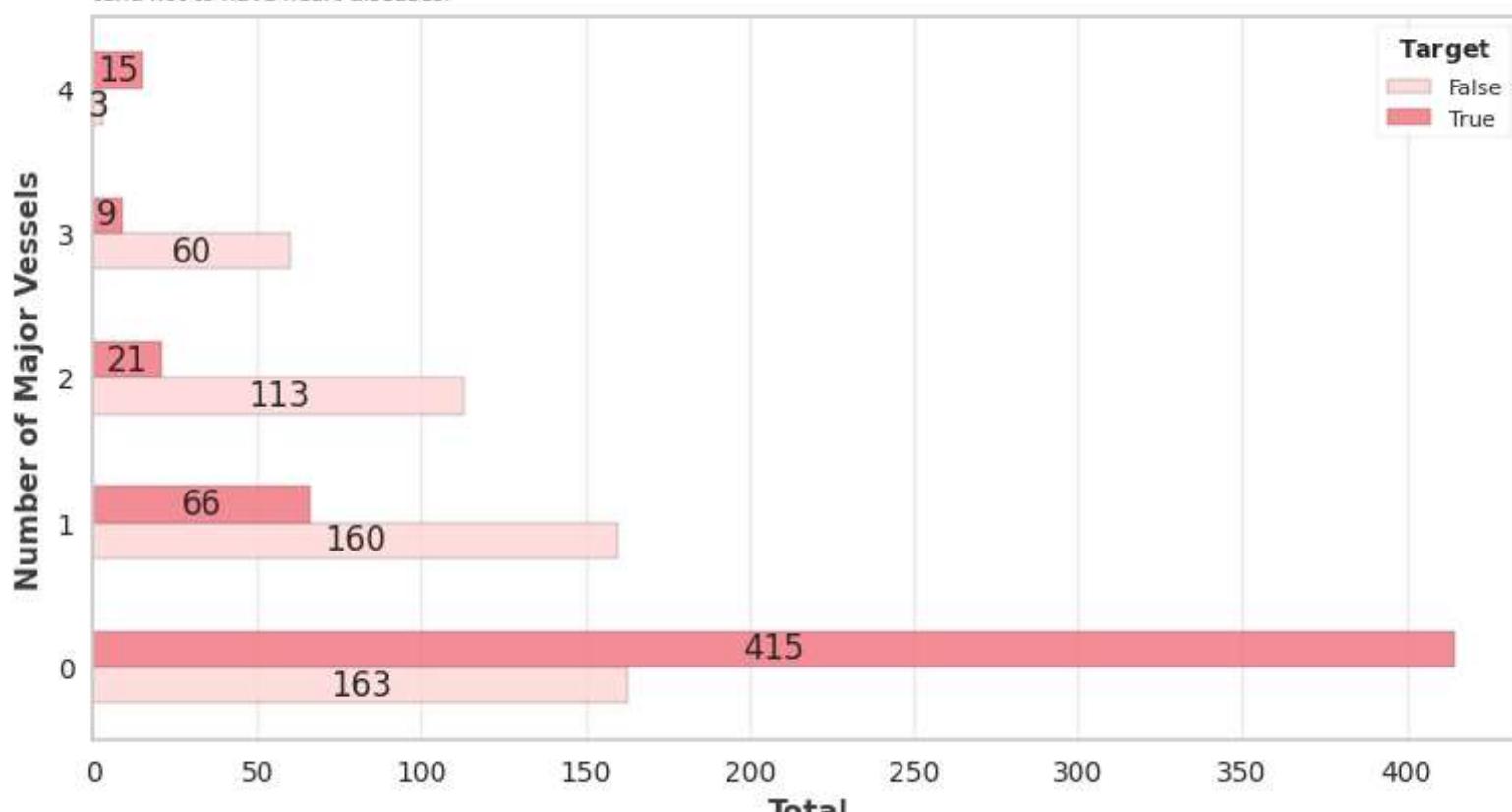
# --- Creating Horizontal Bar Chart ---
ax = pd.crosstab(df.ca, df.target).plot(kind='barh', figsize=(8, 5),
                                         color=color_mix[2:4],
                                         edgecolor=black_grad[2], alpha=0.85)

# --- Horizontal Bar Chart Settings ---
for rect in ax.patches:
    width, height = rect.get_width(), rect.get_height()
    x, y = rect.get_xy()
    ax.text(x+width/2, y+height/2, '{:.0f}'.format(width),
            horizontalalignment='center', verticalalignment='center')

plt.suptitle('Heart Disease Distribution based on Major Vessels Total',
             fontweight='heavy', x=0.069, y=0.98, ha='left', fontsize='16',
             fontfamily='sans-serif', color=black_grad[0])
plt.title('Patients with 0 and 4 major vessels tend to have heart diseases. However, patients who have a number of vessels 1 to 3 tend not to have heart diseases.', fontsize='8', fontfamily='sans-serif', loc='left', color=black_grad[1])
plt.tight_layout(rect=[0, 0.04, 1, 1.025])
plt.xlabel('Total', fontfamily='sans-serif', fontweight='bold', color=black_grad[1])
plt.ylabel('Number of Major Vessels', fontfamily='sans-serif', fontweight='bold',
           color=black_grad[1])
plt.yticks(rotation=0)
plt.grid(axis='x', alpha=0.4)
plt.grid(axis='y', alpha=0)
plt.legend(labels=labels, title='$\backslash$bf{Target}$', fontsize='8', frameon=True,
           title_fontsize='9', loc='upper right');
```

Heart Disease Distribution based on Major Vessels Total

Patients with 0 and 4 major vessels tend to have heart diseases. However, patients who have a number of vessels 1 to 3 tend not to have heart diseases.



Heart Disease Scatter Plot based on Age

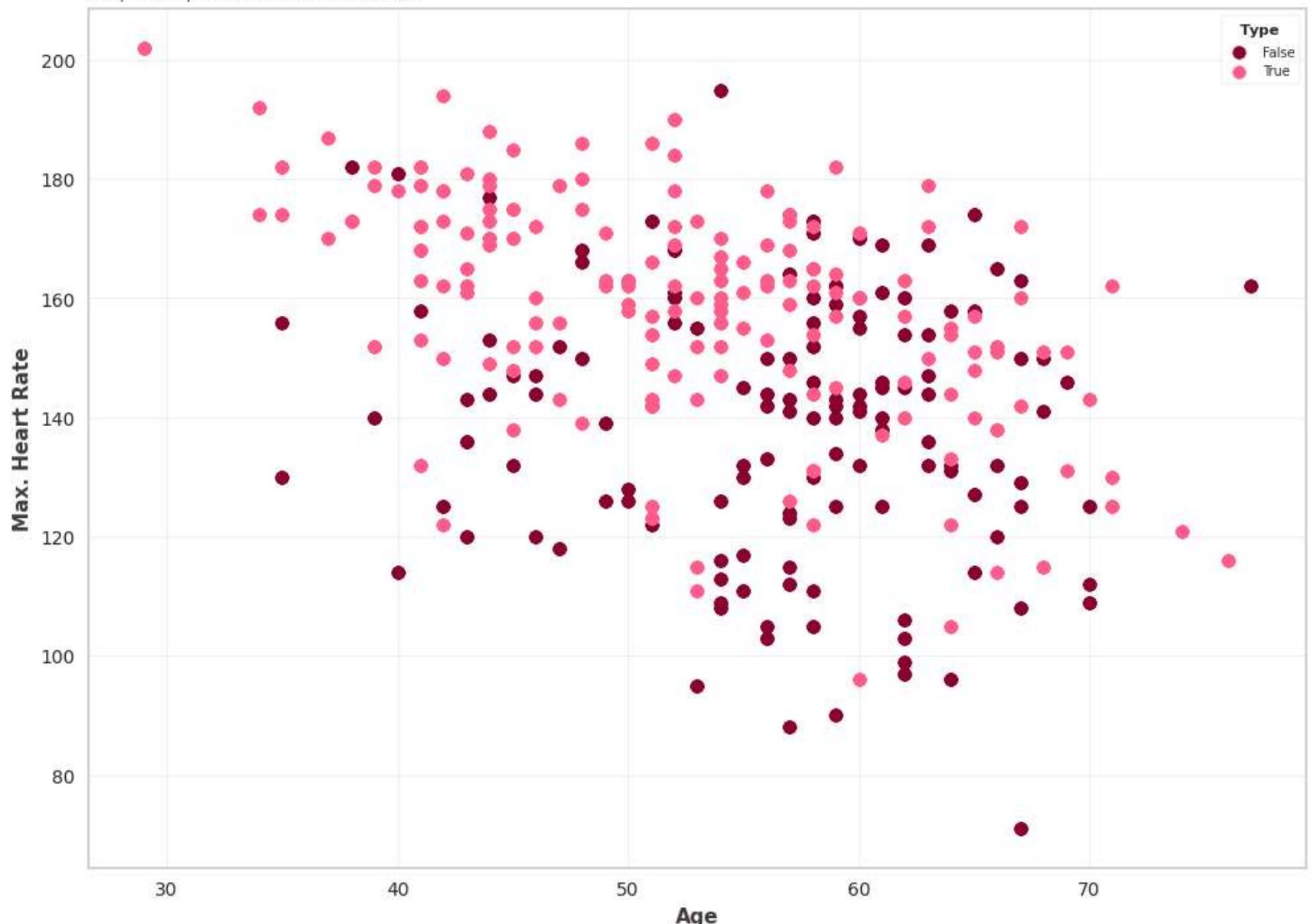
```
In [24]: # -- Scatter Plot Size & Titles Settings ---
plt.figure(figsize=(10, 8))
plt.suptitle('Heart Disease Scatter Plot based on Age', fontweight='heavy',
             x=0.048, y=0.98, ha='left', fontsize=16, fontfamily='sans-serif',
             color=black_grad[0])
plt.title('Based on age, patients with and without heart diseases mostly between 50-70 years old. Patients with heart d:',
             fontsize=8, fontfamily='sans-serif', loc='left', color=black_grad[1])
plt.tight_layout(rect=[0, 0.04, 1, 1.01])

# --- Creating Scatter Plot ---
plt.scatter(x=df.age[df.target==0], y=df.thalach[(df.target==0)], c=pink_grad[0])
plt.scatter(x=df.age[df.target==1], y=df.thalach[(df.target==1)], c=pink_grad[2])

# --- Scatter Plot Legend & Labels Settings ---
plt.legend(['False', 'True'], title='$\bf{Type}$', fontsize=7,
           title_fontsize=8, loc='upper right', frameon=True)
plt.xlabel('Age', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ylabel('Max. Heart Rate', fontweight='bold', fontsize=11,
           fontfamily='sans-serif', color=black_grad[1])
plt.ticklabel_format(style='plain', axis='both')
plt.grid(axis='both', alpha=0.4, lw=0.5)
plt.show();
```

Heart Disease Scatter Plot based on Age

Based on age, patients with and without heart diseases mostly between 50-70 years old. Patients with heart diseases tend to have high heart rate compared to patients with no heart diseases.



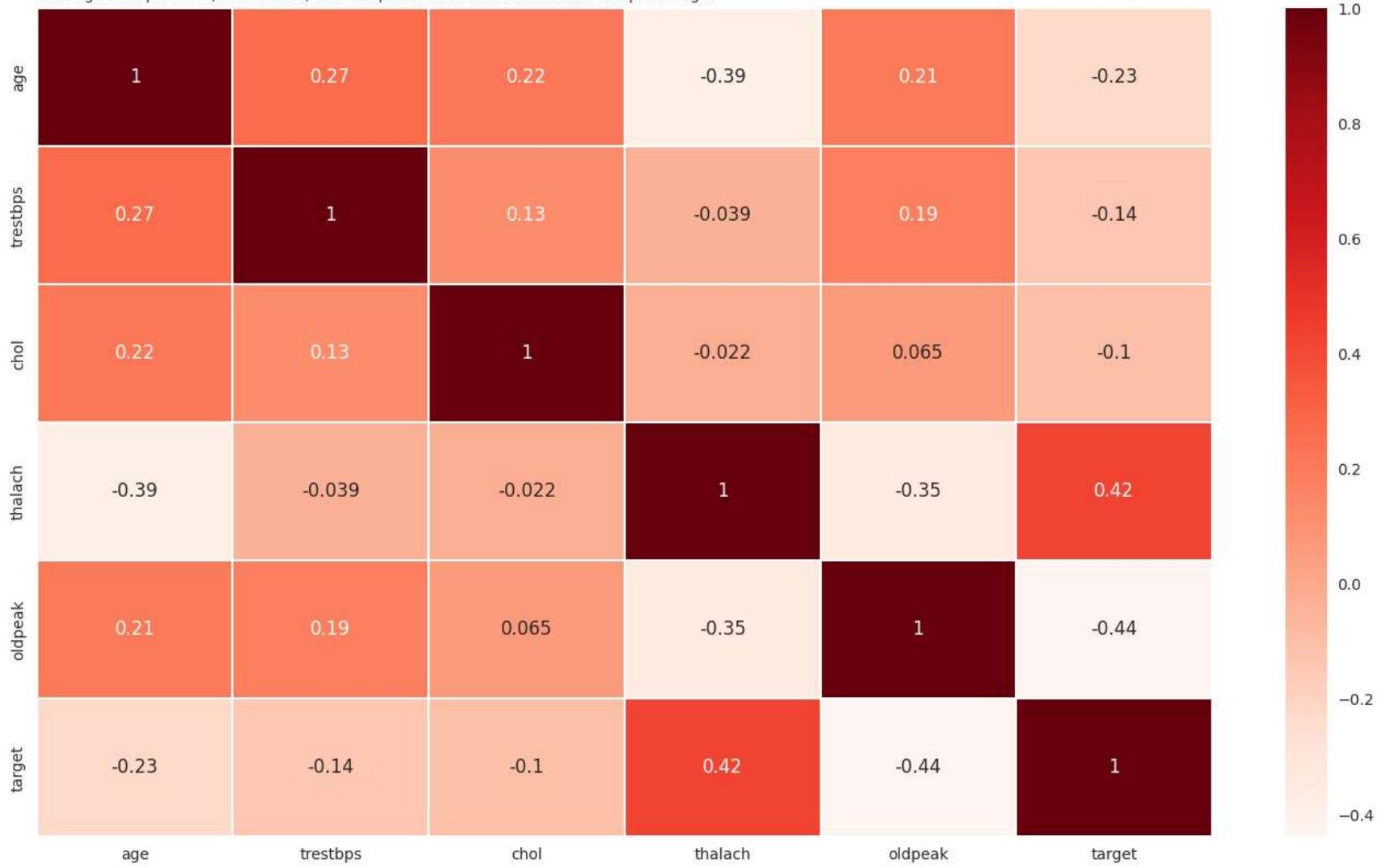
Heatmap

Below is correlation map/heatmap of numerical variables to show correlation level/values for each variables with others.

```
In [25]: # --- Correlation Map (Heatmap) ---
plt.figure(figsize=(14, 9))
sns.heatmap(df.corr(), annot=True, cmap='Reds', linewidths=0.1)
plt.suptitle('Correlation Map of Numerical Variables', fontweight='heavy',
             x=0.03, y=0.98, ha='left', fontsize=16, fontfamily='sans-serif',
             color=black_grad[0])
plt.title('Resting blood pressure, cholestorol, and "oldpeak" have moderate relationship with age.',
             fontsize=10, fontfamily='sans-serif', loc='left', color=black_grad[1])
plt.tight_layout(rect=[0, 0.04, 1, 1.01])
```

Correlation Map of Numerical Variables

Resting blood pressure, cholestral, and "oldpeak" have moderate relationship with age.



Dataset Pre-processing

One-Hot Encoding

The data pre-processing will be transforming categorical variables using one-hot encoding technique.

```
In [26]: # --- Creating Dummy Variables for cp, thal and slope ---
cp = pd.get_dummies(df['cp'], prefix='cp')
thal = pd.get_dummies(df['thal'], prefix='thal')
slope = pd.get_dummies(df['slope'], prefix='slope')

# --- Merge Dummy Variables to Main Data Frame ---
frames = [df, cp, thal, slope]
df = pd.concat(frames, axis = 1)
```

```
In [27]: # --- Display New Data Frame ---
df.head().style.background_gradient(cmap='PuRd').hide_index().set_properties(**{'font-family': 'Segoe UI'})
```

```
Out[27]: age sex cp trestbps chol fbs restecg thalach exang oldpeak slope ca thal target cp_0 cp_1 cp_2 cp_3 thal_0 thal_1 thal_2
52 1 0 125 212 0 1 168 0 1.000000 2 2 3 0 1 0 0 0 0 0 0 0 0
53 1 0 140 203 1 0 155 1 3.100000 0 0 3 0 1 0 0 0 0 0 0 0 0
70 1 0 145 174 0 1 125 1 2.600000 0 0 3 0 1 0 0 0 0 0 0 0 0
61 1 0 148 203 0 1 161 0 0.000000 2 1 3 0 1 0 0 0 0 0 0 0 0
62 0 0 138 294 1 1 106 0 1.900000 1 3 2 0 1 0 0 0 0 0 0 0 1
```

Dropping Unnecessary Variables

```
In [28]: # --- Drop Unnecessary Variables ---
df = df.drop(columns = ['cp', 'thal', 'slope'])
```

```
In [29]: # --- Display New Data Frame ---
df.head().style.background_gradient(cmap='Reds').hide_index().set_properties(**{'font-family': 'Segoe UI'})
```

```
Out[29]: age sex trestbps chol fbs restecg thalach exang oldpeak ca target cp_0 cp_1 cp_2 cp_3 thal_0 thal_1 thal_2 thal_3 slope_0
52 1 125 212 0 1 168 0 1.000000 2 0 1 0 0 0 0 0 0 0 1 0
53 1 140 203 1 0 155 1 3.100000 0 0 1 0 0 0 0 0 0 0 1 1
70 1 145 174 0 1 125 1 2.600000 0 0 1 0 0 0 0 0 0 0 1 1
61 1 148 203 0 1 161 0 0.000000 1 0 1 0 0 0 0 0 0 0 1 0
62 0 138 294 1 1 106 0 1.900000 3 0 1 0 0 0 0 0 0 0 1 0
```

Features Separating

In this section, the 'target' (dependent) column will be separated from independent columns.

```
In [30]: # --- Separating Dependent Features ---
x = df.drop(['target'], axis=1)
y = df['target']
```

Data Normalization

In this section, data normalization will be performed to normalize the range of independent variables or features of data.

Data normalization will use **min-max normalization**.

Min-max normalization is often known as feature scaling where the values of a numeric range of a feature of data, are reduced to a scale between 0 and 1.

```
In [31]: # --- Data Normalization using Min-Max Method ---
x = MinMaxScaler().fit_transform(x)
```

Splitting the Dataset

The dataset will be splitted into 80:20 ratio (80% training and 20% testing).

```
In [32]: # --- Splitting Dataset into 80:20 ---
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=4)
```

Model Implementation

Logistic Regression

```
In [33]: # --- Applying Logistic Regression ---
LRclassifier = LogisticRegression(max_iter=1000, random_state=1, solver='liblinear', penalty='l1')
LRclassifier.fit(x_train, y_train)

y_pred_LR = LRclassifier.predict(x_test)
```

Logistic Regression -84% acc

```
In [34]: # --- LR Accuracy ---
LRAcc = accuracy_score(y_pred_LR, y_test)
print('... Logistic Regression Accuracy:' + '\033[1m {:.2f}%'.format(LRAcc*100) + ' ...')

... Logistic Regression Accuracy: 83.90% ...
```

K-Nearest Neighbour (KNN)

```
In [35]: # --- Applying KNN ---
KNNclassifier = KNeighborsClassifier(n_neighbors=3)
KNNclassifier.fit(x_train, y_train)

y_pred_KNN = KNNclassifier.predict(x_test)
```

KNearest Neighbors - 96% acc

```
In [36]: # --- KNN Accuracy ---
KNNAcc = accuracy_score(y_pred_KNN, y_test)
print('... K-Nearest Neighbour Accuracy:' + '\033[1m {:.2f}%'.format(KNNAcc*100) + ' ...')

... K-Nearest Neighbour Accuracy: 95.61% ...
```

Support Vector Machine (SVM)

```
In [37]: # --- Applying SVM ---
SVMclassifier = SVC(kernel='linear', max_iter=1000, C=10, probability=True)
SVMclassifier.fit(x_train, y_train)

y_pred_SVM = SVMclassifier.predict(x_test)
```

SVM - 84% acc

```
In [38]: # --- SVM Accuracy ---
SVMAcc = accuracy_score(y_pred_SVM, y_test)
print('... Support Vector Machine Accuracy:' + '\033[1m {:.2f}%'.format(SVMAcc*100) + ' ...')

... Support Vector Machine Accuracy: 83.90% ...
```

Gaussian Naive Bayes

```
In [39]: # --- Applying Gaussian NB ---
GNBclassifier = GaussianNB(var_smoothing=0.1)
GNBclassifier.fit(x_train, y_train)

y_pred_GNB = GNBclassifier.predict(x_test)
```

Gaussian Naive Bayes - 83% acc

```
In [40]: # --- GNB Accuracy ---
GNBAcc = accuracy_score(y_pred_GNB, y_test)
print('... Gaussian Naive Bayes Accuracy:' + '\033[1m {:.2f}%'.format(GNBAcc*100) + ' ...')

... Gaussian Naive Bayes Accuracy: 82.44% ...
```

Decision Tree

```
In [41]: # --- Applying Decision Tree ---
DTCclassifier = DecisionTreeClassifier(max_depth=3, min_samples_leaf=5, criterion='entropy', min_samples_split=5,
```

```

        splitter='random', random_state=1)

DTClassifier.fit(x_train, y_train)
y_pred_DTC = DTClassifier.predict(x_test)

```

Decision Tree- 84% acc

```
In [42]: # --- Decision Tree Accuracy ---
DTAcc = accuracy_score(y_pred_DTC, y_test)
print('... Decision Tree Accuracy:'+'{:3.2f}%'.format(DTAcc*100)+' ...')

... Decision Tree Accuracy: 83.90% ...

```

Random Forest

```
In [43]: # --- Applying Random Forest ---
RFclassifier = RandomForestClassifier(n_estimators=1000, random_state=1, max_leaf_nodes=20, min_samples_split=15)

RFclassifier.fit(x_train, y_train)
y_pred_RF = RFclassifier.predict(x_test)
```

Random Forest - 89% acc

```
In [44]: # --- Random Forest Accuracy ---
RFAcc = accuracy_score(y_pred_RF, y_test)
print('... Random Forest Accuracy:'+'{:3.2f}%'.format(RFAcc*100)+' ...')

... Random Forest Accuracy: 88.78% ...

```

Gradient Boosting

```
In [45]: # --- Applying Gradient Boosting ---
GBclassifier = GradientBoostingClassifier(random_state=1, n_estimators=100, max_leaf_nodes=3, loss='exponential',
                                            min_samples_leaf=20)

GBclassifier.fit(x_train, y_train)
y_pred_GB = GBclassifier.predict(x_test)
```

Gradient Boosting - 88% acc

```
In [46]: # --- Gradient Boosting Accuracy ---
GBAcc = accuracy_score(y_pred_GB, y_test)
print('... Gradient Boosting Accuracy:'+'{:3.2f}%'.format(GBAcc*100)+' ...')

... Gradient Boosting Accuracy: 86.83% ...

```

Machine Learning Models Comparison

```
In [47]: # --- Create Accuracy Comparison Table ---
compare = pd.DataFrame({'Model': ['Logistic Regression', 'K-Nearest Neighbour', 'Support Vector Machine',
                                    'Gaussian Naive Bayes', 'Decision Tree', 'Random Forest', 'Gradient Boosting'],
                         'Accuracy': [LRAcc*100, KNNAcc*100, SVMAcc*100, GNBAcc*100, DTCAcc*100, RFAcc*100, GBAcc*100,
                                      ]})

# --- Create Accuracy Comparison Table ---
compare.sort_values(by='Accuracy', ascending=False).style.background_gradient(cmap='PuRd').hide_index().set_properties(
    {'text-align': 'center'})
```

Model	Accuracy
K-Nearest Neighbour	95.609756
Random Forest	88.780488
Gradient Boosting	86.829268
Logistic Regression	83.902439
Support Vector Machine	83.902439
Decision Tree	83.902439
Gaussian Naive Bayes	82.439024