

Heart Disease Analysis

Heart disease or Cardiovascular disease (CVD) is a class of diseases that involve the heart or blood vessels. Cardiovascular diseases are the leading cause of death globally. This is true in all areas of the world except Africa. Together CVD resulted in 17.9 million deaths (32.1%) in 2015. Deaths, at a given age, from CVD are more common and have been increasing in much of the developing world, while rates have declined in most of the developed world since the 1970s.

So, in this kernel, I have conducted Exploratory Data Analysis or EDA of the heart disease dataset. Exploratory Data Analysis or EDA is a critical first step in analyzing a new dataset. The primary objective of EDA is to analyze the data for distribution, outliers and anomalies in the dataset. It enables us to direct specific testing of the hypothesis. It includes analysing the data to find the distribution of data, its main characteristics, identifying patterns and visualizations. It also provides tools for hypothesis generation by visualizing and understanding the data through graphical representation.

Import libraries

```
In [12]: import numpy as np # Linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the input directory

import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# Any results you write to the current directory are saved as output.
```

```
In [13]: import seaborn as sns
import matplotlib.pyplot as plt
import scipy.stats as st
%matplotlib inline
```

```
sns.set(style="whitegrid")
```

```
In [14]: # ignore warnings
```

```
import warnings  
warnings.filterwarnings('ignore')
```

I have imported the libraries. The next step is to import the datasets.

Import dataset

```
In [15]: df = pd.read_csv('C:/Users/Rachana Jena/Downloads/9th- Seaborn, Eda practice/9th- Seaborn, Eda practice/EDA/heart.csv')
```

```
In [16]: df
```

Out[16]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63	1	3	145	233	1	0	150	0	2.3	0	0	1	1
1	37	1	2	130	250	0	1	187	0	3.5	0	0	2	1
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	1
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	1
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	1
...
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3	0
299	45	1	3	110	264	0	1	132	0	1.2	1	0	3	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2	0

303 rows × 14 columns

Exploratory Data Analysis

In [17]:

```
# print the shape
print('The shape of the dataset : ', df.shape)
```

The shape of the dataset : (303, 14)

Now, we can see that the dataset contains 303 instances and 14 variables.

In [18]:

```
# Preview dataset
df.head()
```

Out[18]:

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

Summary of Dataset

In [19]:

```
# summary of dataset
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype  
--- 
 0   age         303 non-null    int64  
 1   sex         303 non-null    int64  
 2   cp          303 non-null    int64  
 3   trestbps   303 non-null    int64  
 4   chol        303 non-null    int64  
 5   fbs         303 non-null    int64  
 6   restecg    303 non-null    int64  
 7   thalach    303 non-null    int64  
 8   exang       303 non-null    int64  
 9   oldpeak    303 non-null    float64 
 10  slope       303 non-null    int64  
 11  ca          303 non-null    int64  
 12  thal        303 non-null    int64  
 13  target      303 non-null    int64  
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

In [21]:

```
df.dtypes
```

```
Out[21]: age          int64  
sex           int64  
cp            int64  
trestbps     int64  
chol          int64  
fbs           int64  
restecg      int64  
thalach       int64  
exang          int64  
oldpeak      float64  
slope         int64  
ca             int64  
thal           int64  
target         int64  
dtype: object
```

```
In [22]: # statistical properties of dataset  
df.describe()
```

```
Out[22]:   age      sex      cp      trestbps      chol      fbs      restecg      thalach      exang      oldpeak  
count  303.000000  303.000000  303.000000  303.000000  303.000000  303.000000  303.000000  303.000000  303.000000  303.000000  
mean   54.366337  0.683168  0.966997  131.623762  246.264026  0.148515  0.528053  149.646865  0.326733  1.039604  
std    9.082101  0.466011  1.032052  17.538143  51.830751  0.356198  0.525860  22.905161  0.469794  1.161075  
min   29.000000  0.000000  0.000000  94.000000  126.000000  0.000000  0.000000  71.000000  0.000000  0.000000  
25%   47.500000  0.000000  0.000000  120.000000  211.000000  0.000000  0.000000  133.500000  0.000000  0.000000  
50%   55.000000  1.000000  1.000000  130.000000  240.000000  0.000000  1.000000  153.000000  0.000000  0.800000  
75%   61.000000  1.000000  2.000000  140.000000  274.500000  0.000000  1.000000  166.000000  1.000000  1.600000  
max   77.000000  1.000000  3.000000  200.000000  564.000000  1.000000  2.000000  202.000000  1.000000  6.200000
```



```
In [23]: df.columns
```

```
Out[23]: Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
   'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
  dtype='object')
```

Univariate analysis

```
In [24]: df['target'].nunique()
```

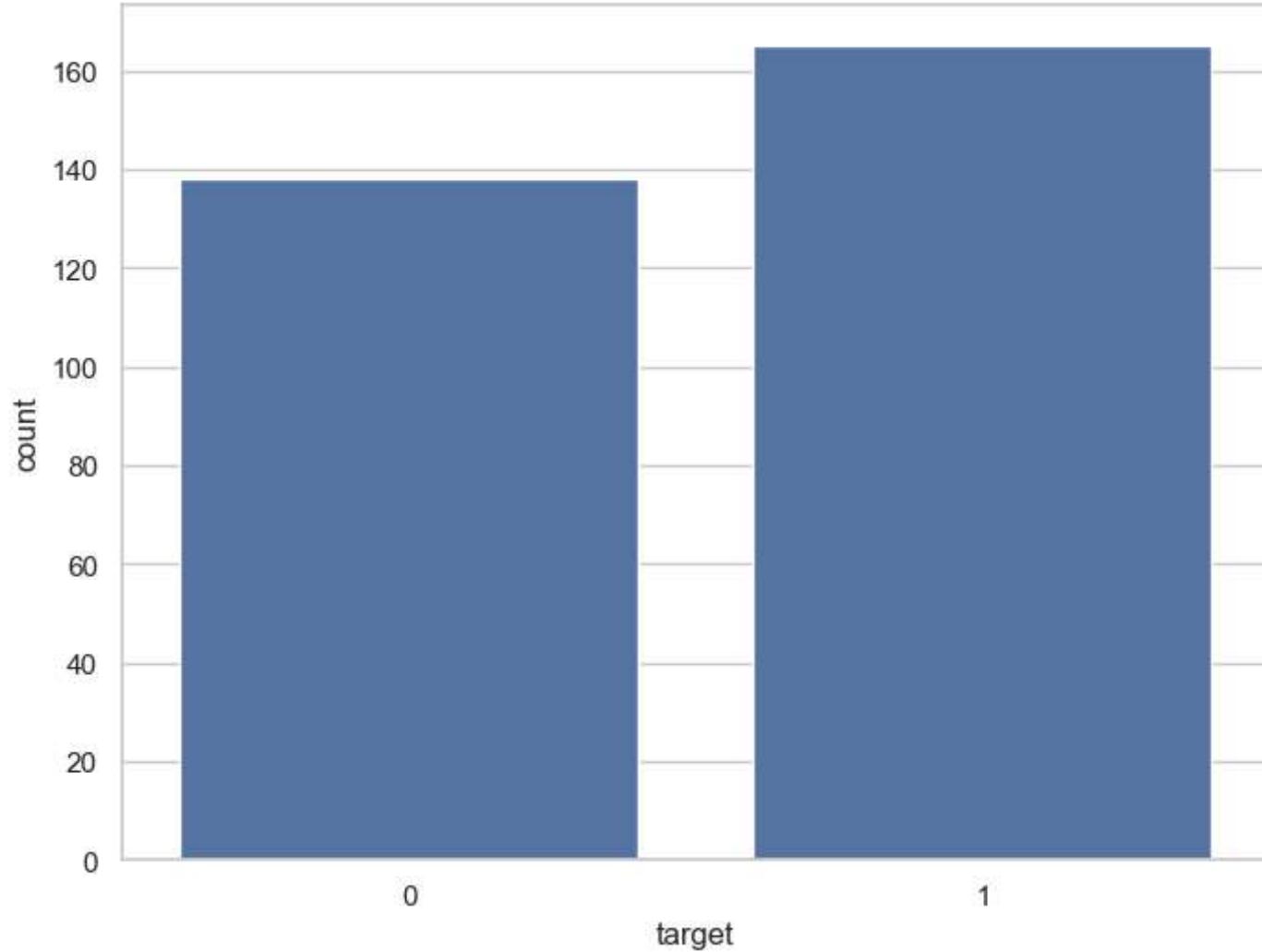
```
Out[24]: 2
```

```
In [25]: df['target'].value_counts()
```

```
Out[25]: target
1    165
0    138
Name: count, dtype: int64
```

Visualize frequency distribution of target variable

```
In [29]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="target", data=df)
plt.show()
```



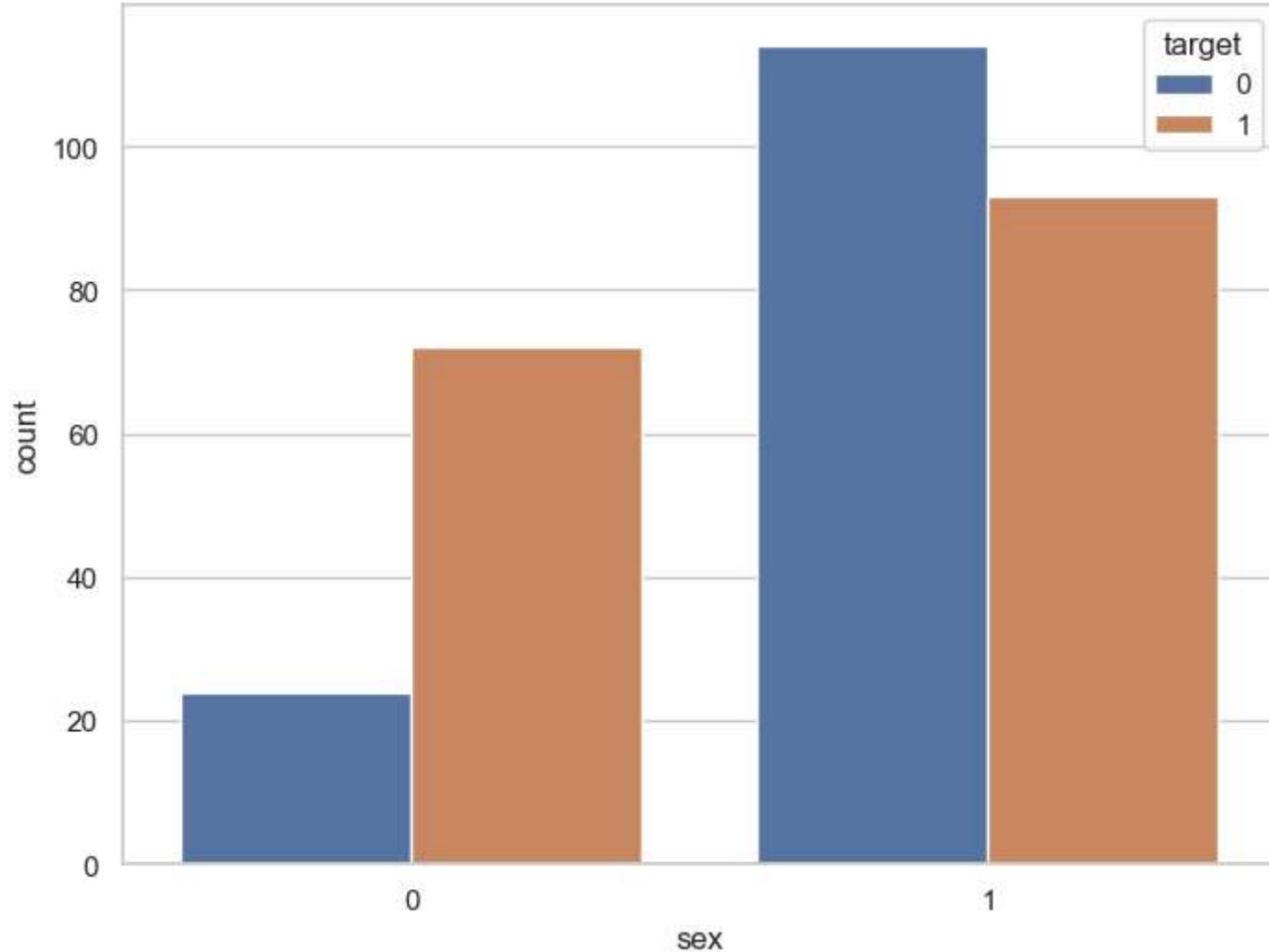
Frequency distribution of target variable wrt sex

```
In [27]: df.groupby('sex')['target'].value_counts()
```

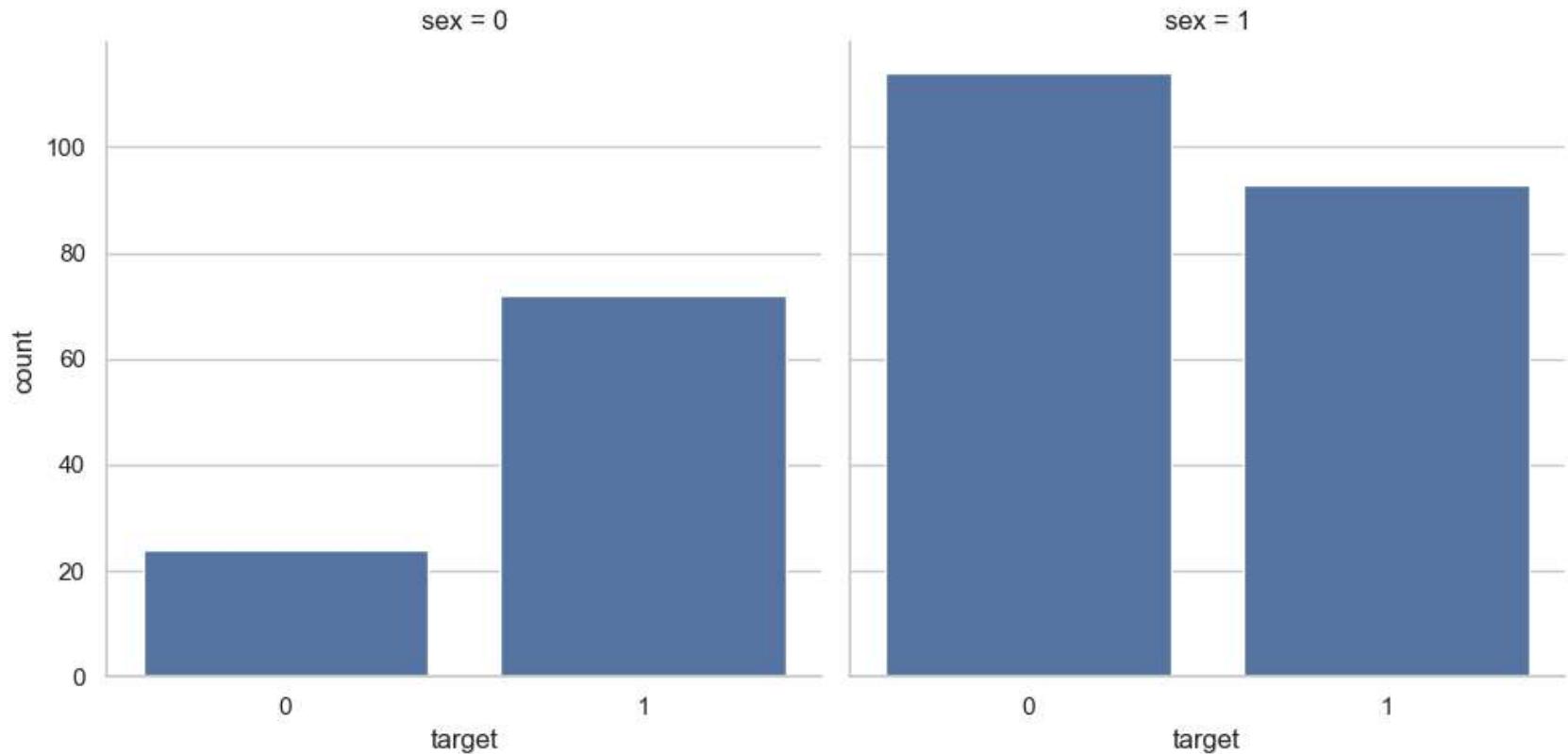
```
Out[27]: sex  target
      0      1          72
      0      0          24
      1      0         114
      1      1          93
Name: count, dtype: int64
```

We can visualize the value counts of the `sex` variable wrt `target` as follows -

```
In [28]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="sex", hue="target", data=df)
plt.show()
```

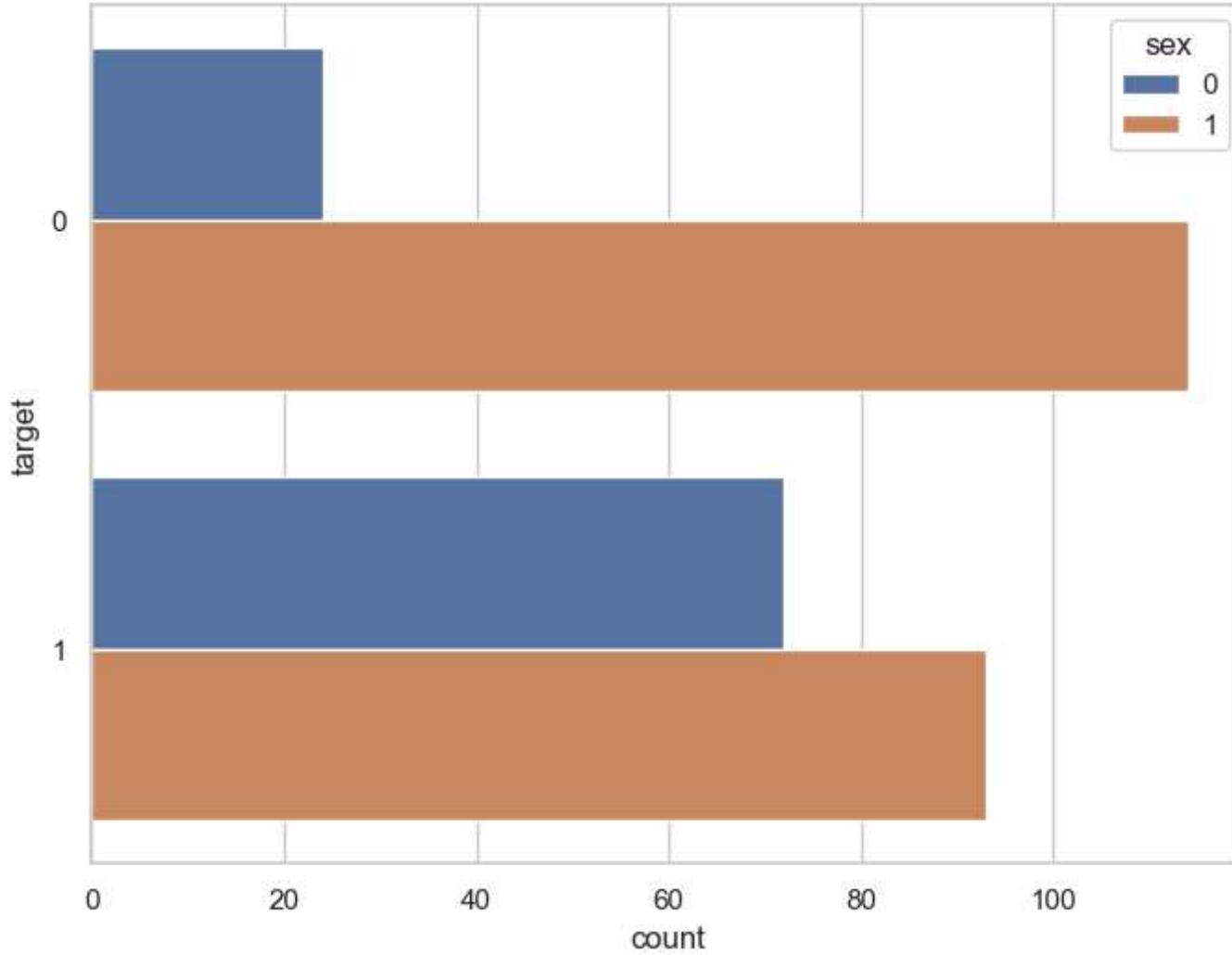


```
In [30]: ax = sns.catplot(x="target", col="sex", data=df, kind="count", height=5, aspect=1)
```



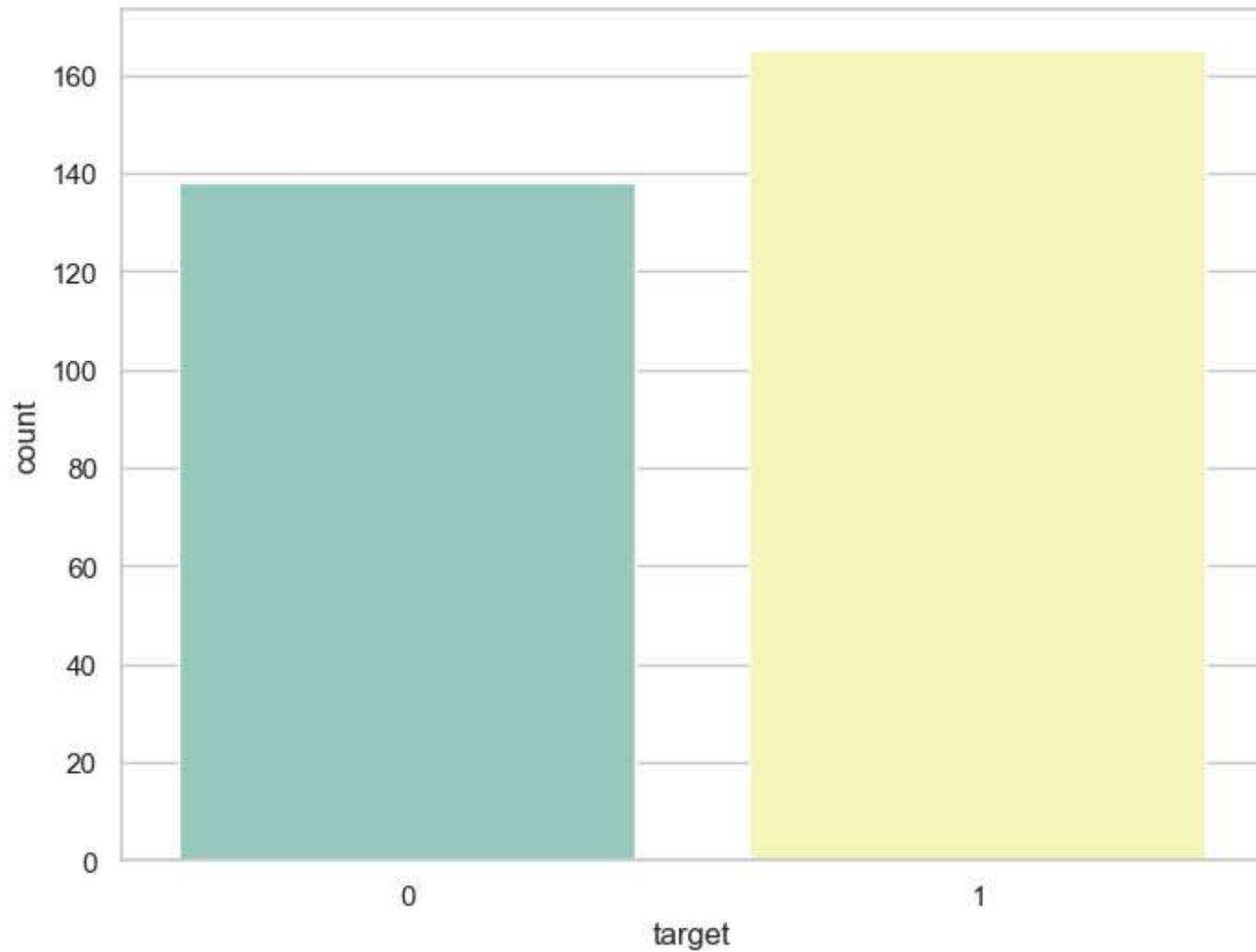
We can plot the bars horizontally as follows :

```
In [31]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(y="target", hue="sex", data=df)
plt.show()
```

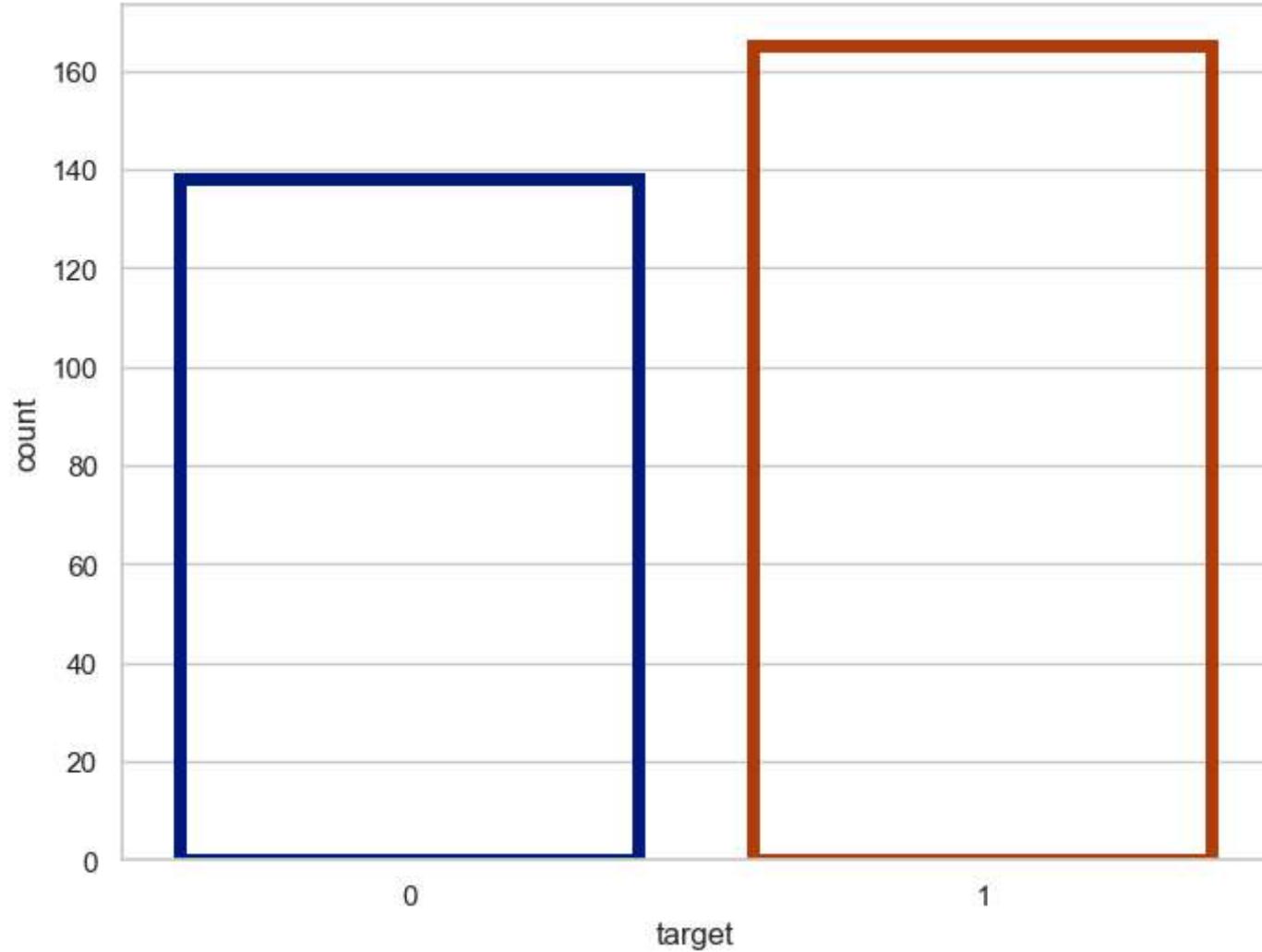


We can use a different color palette as follows :

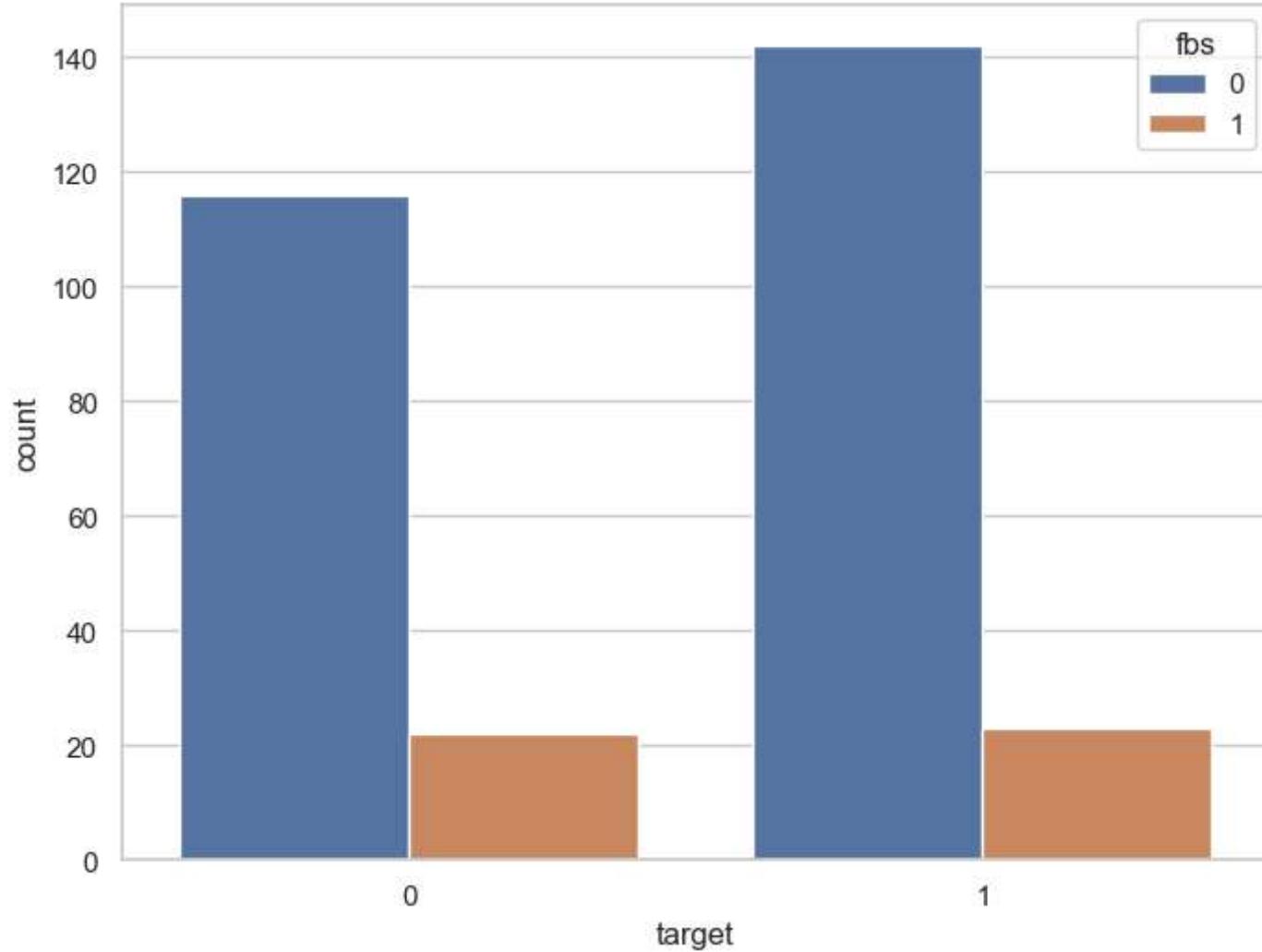
```
In [32]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="target", data=df, palette="Set3")
plt.show()
```



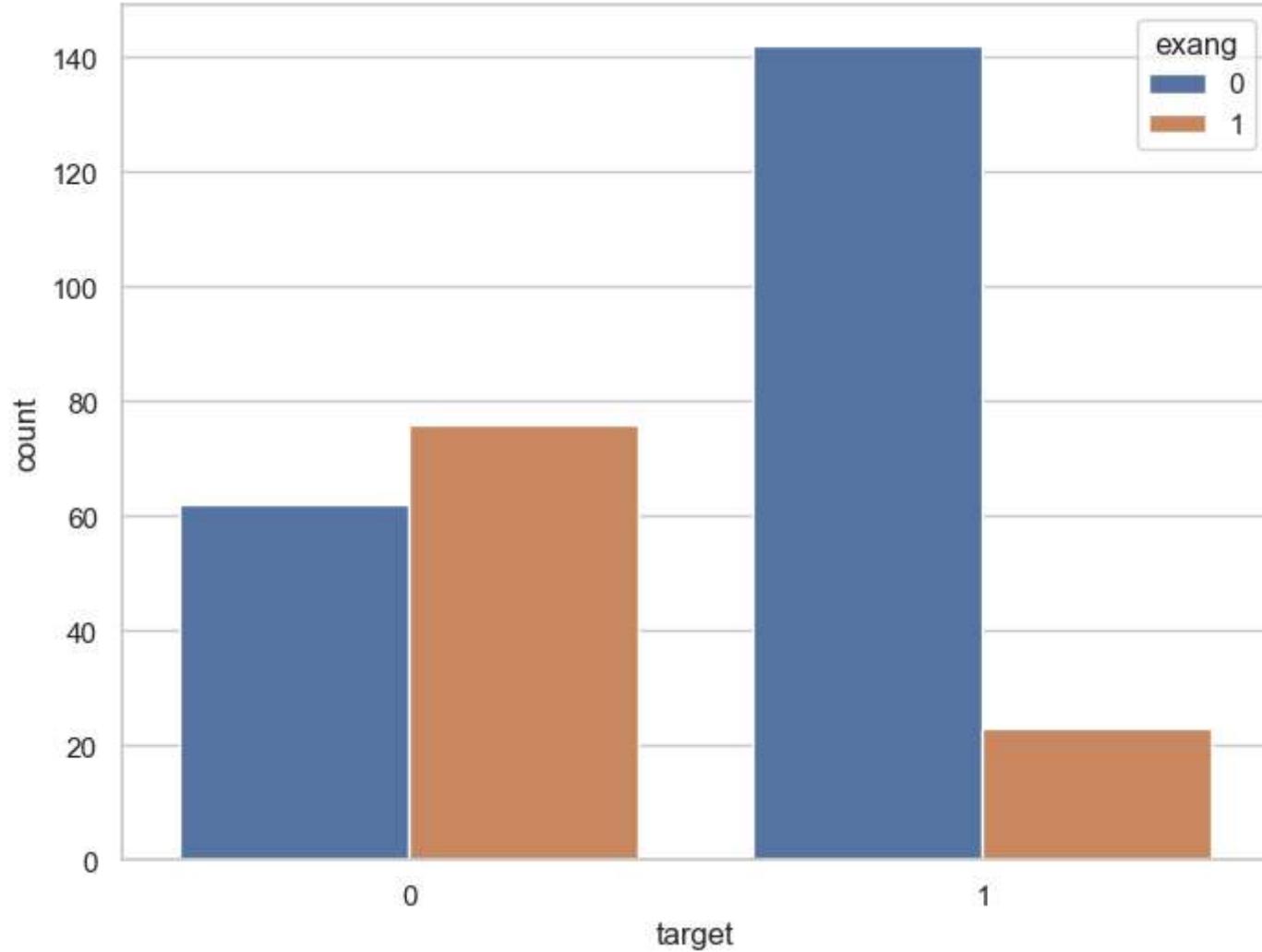
```
In [33]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="target", data=df, facecolor=(0, 0, 0, 0), linewidth=5, edgecolor=sns.color_palette("dark", 3))
plt.show()
```



```
In [34]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="target", hue="fbs", data=df)
plt.show()
```



```
In [35]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="target", hue="exang", data=df)
plt.show()
```



Bivariate Analysis

```
In [36]: correlation = df.corr()
```

```
In [37]: correlation['target'].sort_values(ascending=False)
```

```
Out[37]: target      1.000000
          cp         0.433798
          thalach    0.421741
          slope      0.345877
          restecg   0.137230
          fbs        -0.028046
          chol       -0.085239
          trestbps   -0.144931
          age        -0.225439
          sex        -0.280937
          thal       -0.344029
          ca         -0.391724
          oldpeak    -0.430696
          exang      -0.436757
          Name: target, dtype: float64
```

Analysis of target and cp variable

```
In [38]: df['cp'].nunique()
```

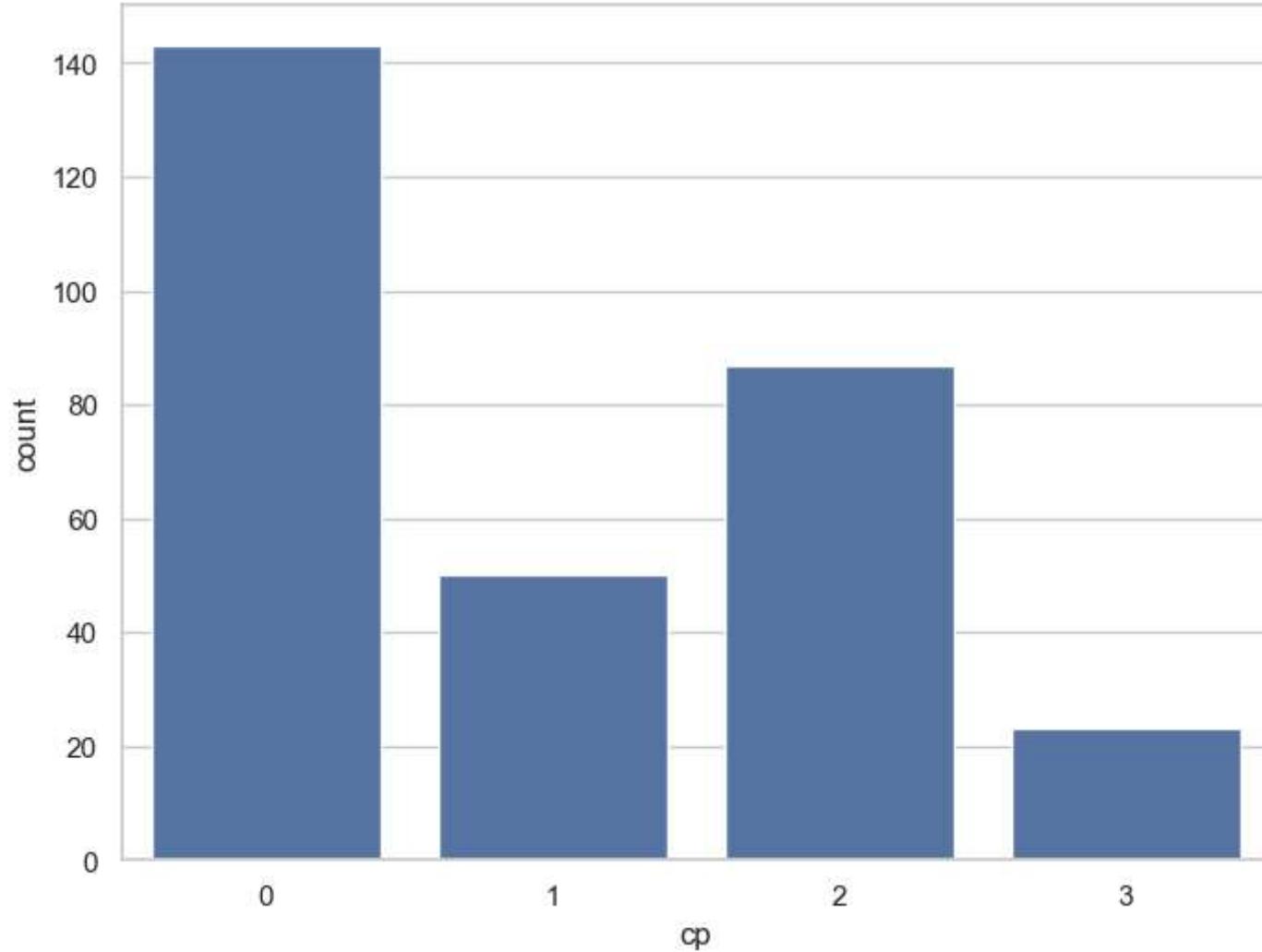
```
Out[38]: 4
```

```
In [39]: df['cp'].value_counts()
```

```
Out[39]: cp
          0    143
          2     87
          1     50
          3     23
          Name: count, dtype: int64
```

visualize frequency distribution of cp variable

```
In [40]: f, ax = plt.subplots(figsize=(8, 6))
          ax = sns.countplot(x="cp", data=df)
          plt.show()
```

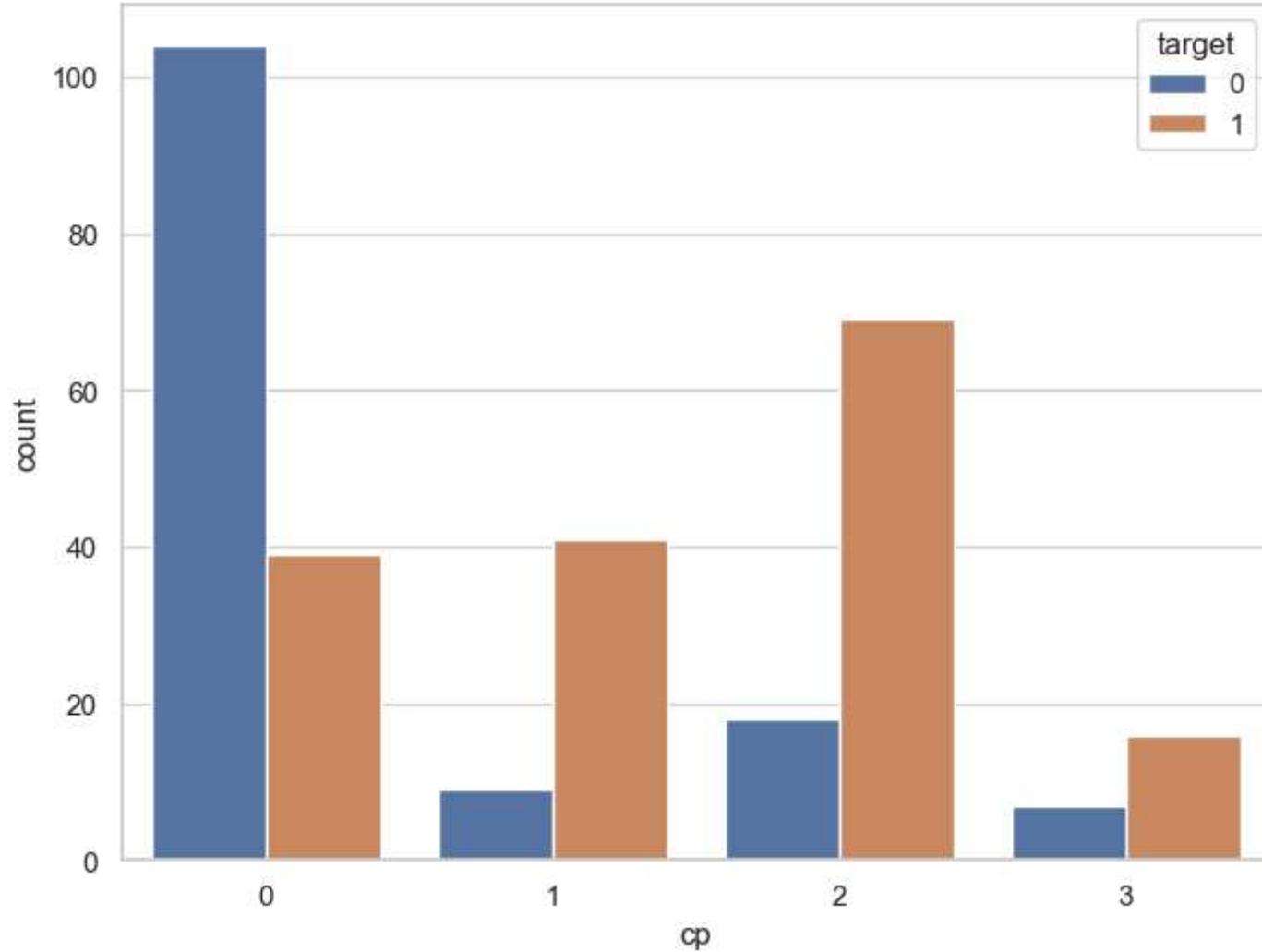


Frequency distribution of target variable wrt cp

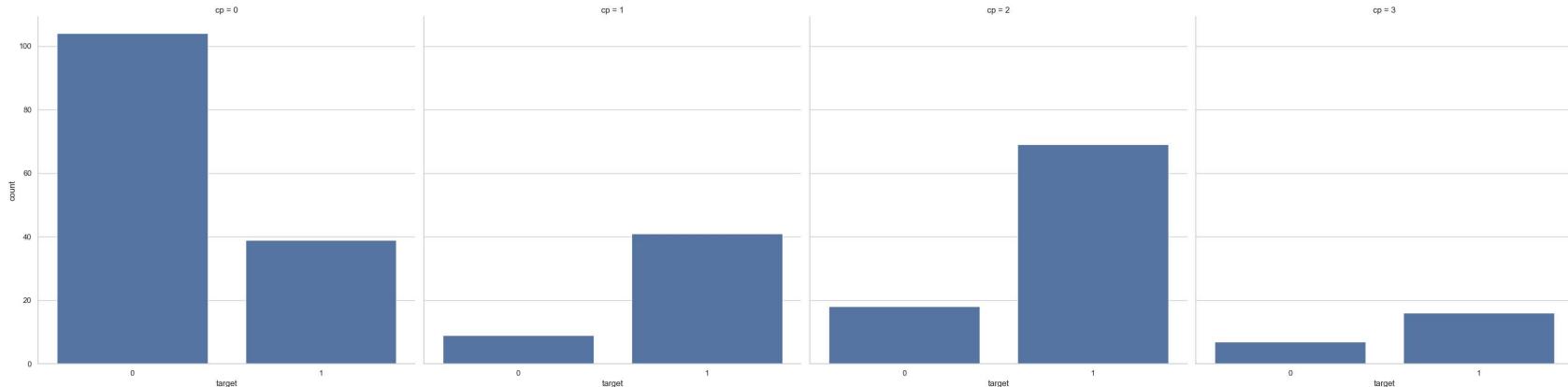
```
In [41]: df.groupby('cp')['target'].value_counts()
```

```
Out[41]: cp  target
      0      104
      1       39
      1      41
      0       9
      1      69
      0      18
      1      16
      0       7
Name: count, dtype: int64
```

```
In [42]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="cp", hue="target", data=df)
plt.show()
```



```
In [43]: ax = sns.catplot(x="target", col="cp", data=df, kind="count", height=8, aspect=1)
```



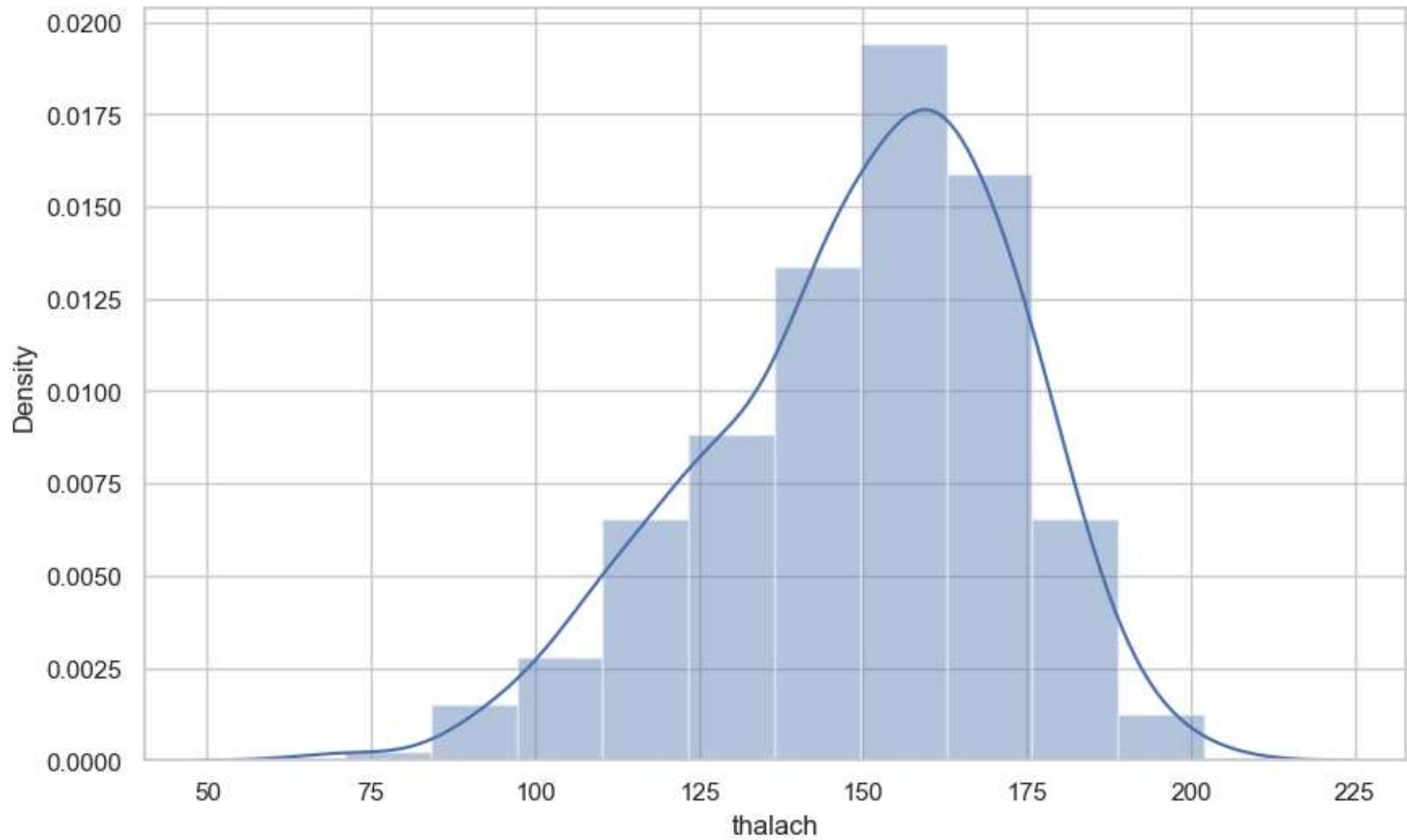
Analysis of `target` and `thalach` variable

```
In [44]: df['thalach'].nunique()
```

```
Out[44]: 91
```

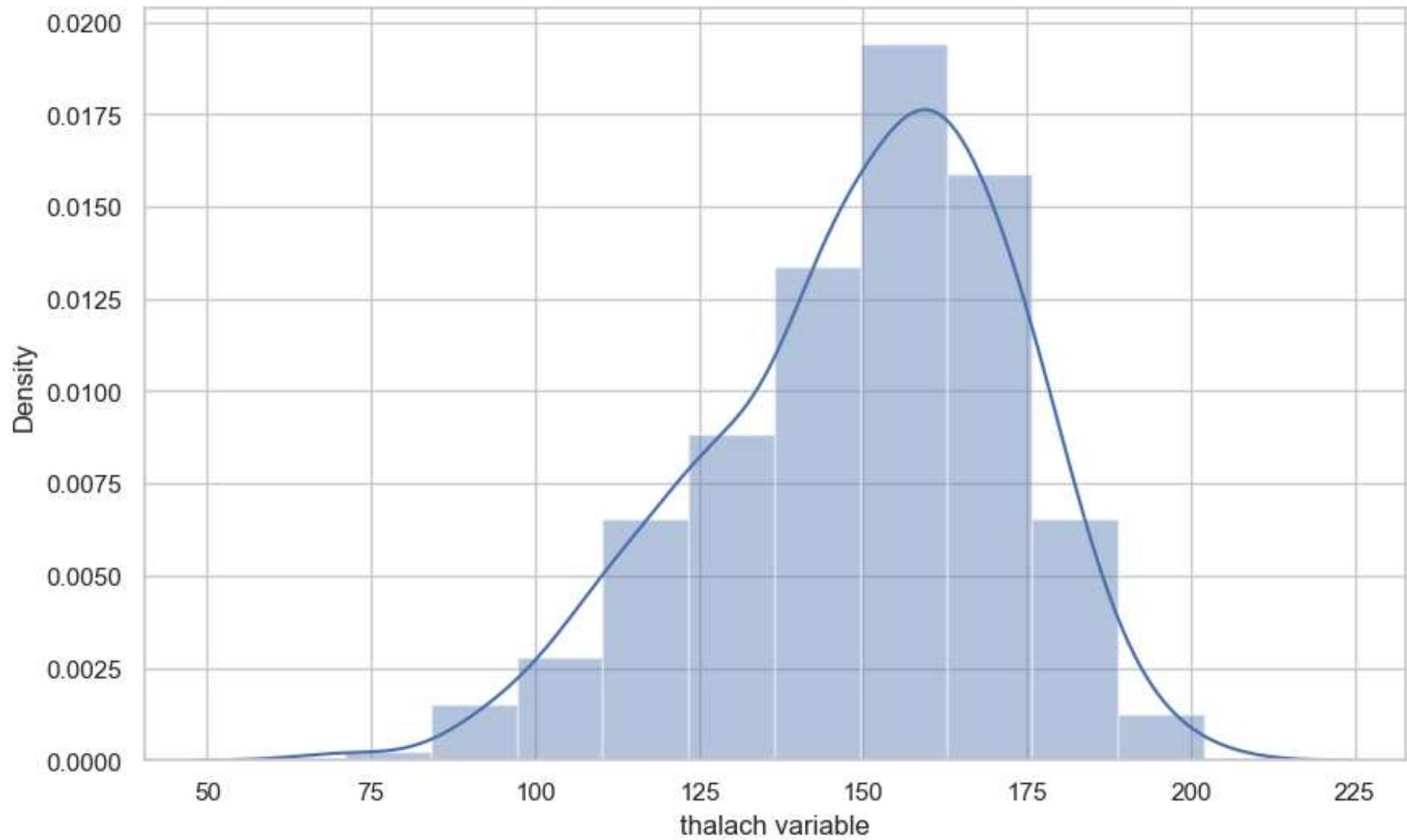
Visualize the frequency distribution of `thalach` variable

```
In [45]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
ax = sns.distplot(x, bins=10)
plt.show()
```



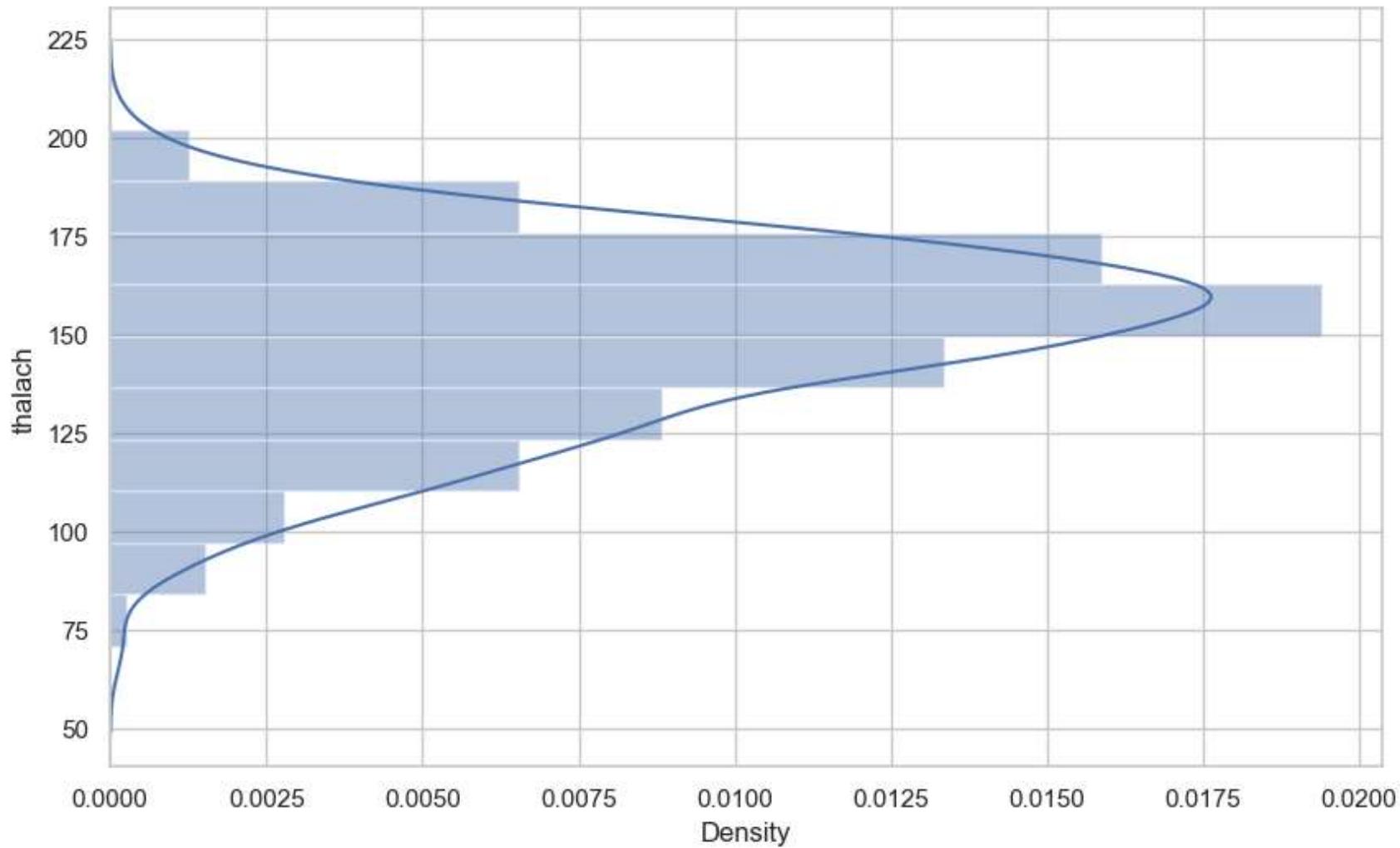
We can use Pandas series object to get an informative axis label as follows :

```
In [46]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.distplot(x, bins=10)
plt.show()
```



We can plot the distribution on the vertical axis as follows:-

```
In [47]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
ax = sns.distplot(x, bins=10, vertical=True)
plt.show()
```

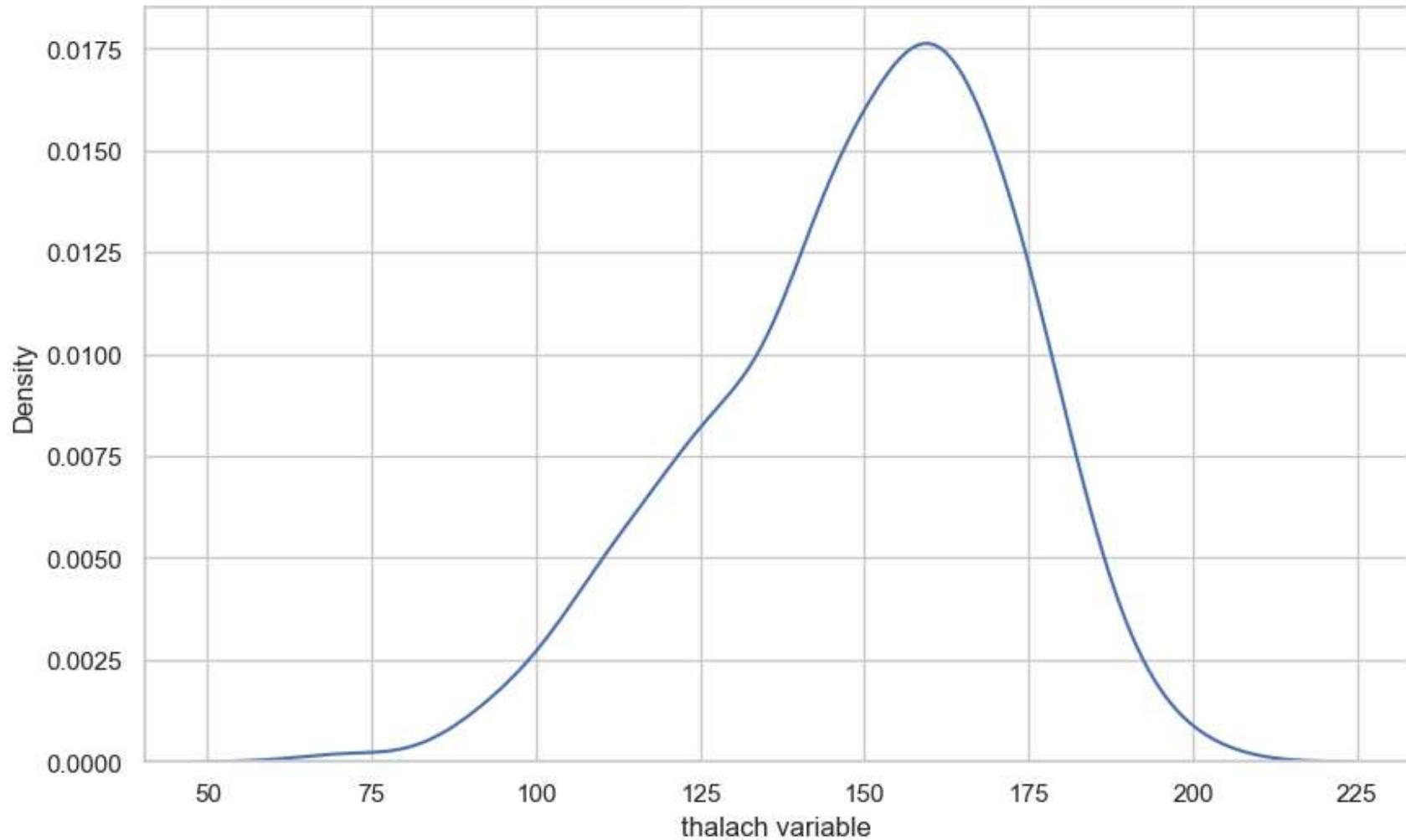


Seaborn Kernel Density Estimation (KDE) plot

- We can plot a KDE plot as follows :

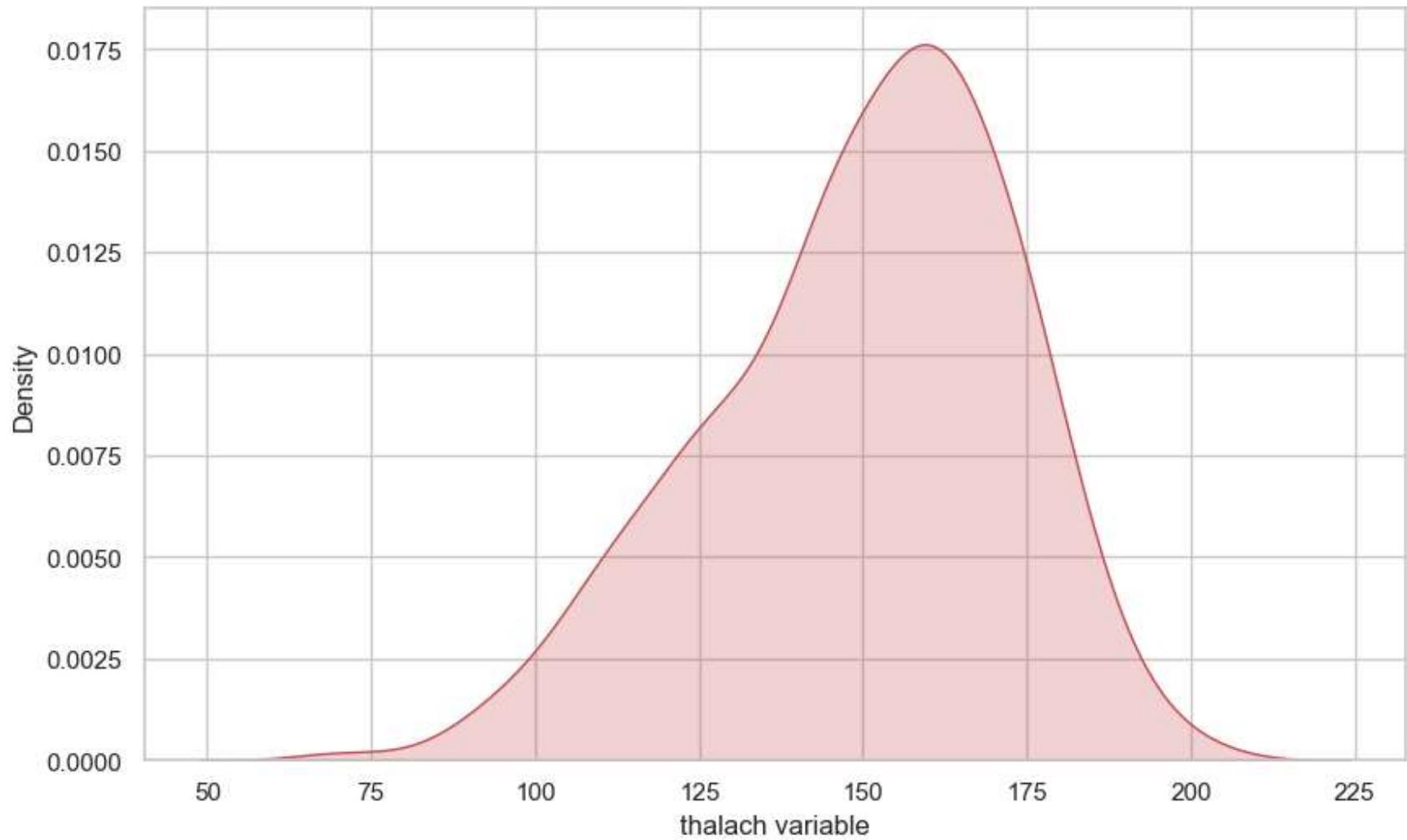
```
In [48]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
```

```
ax = sns.kdeplot(x)
plt.show()
```



We can shade under the density curve and use a different color as follows:

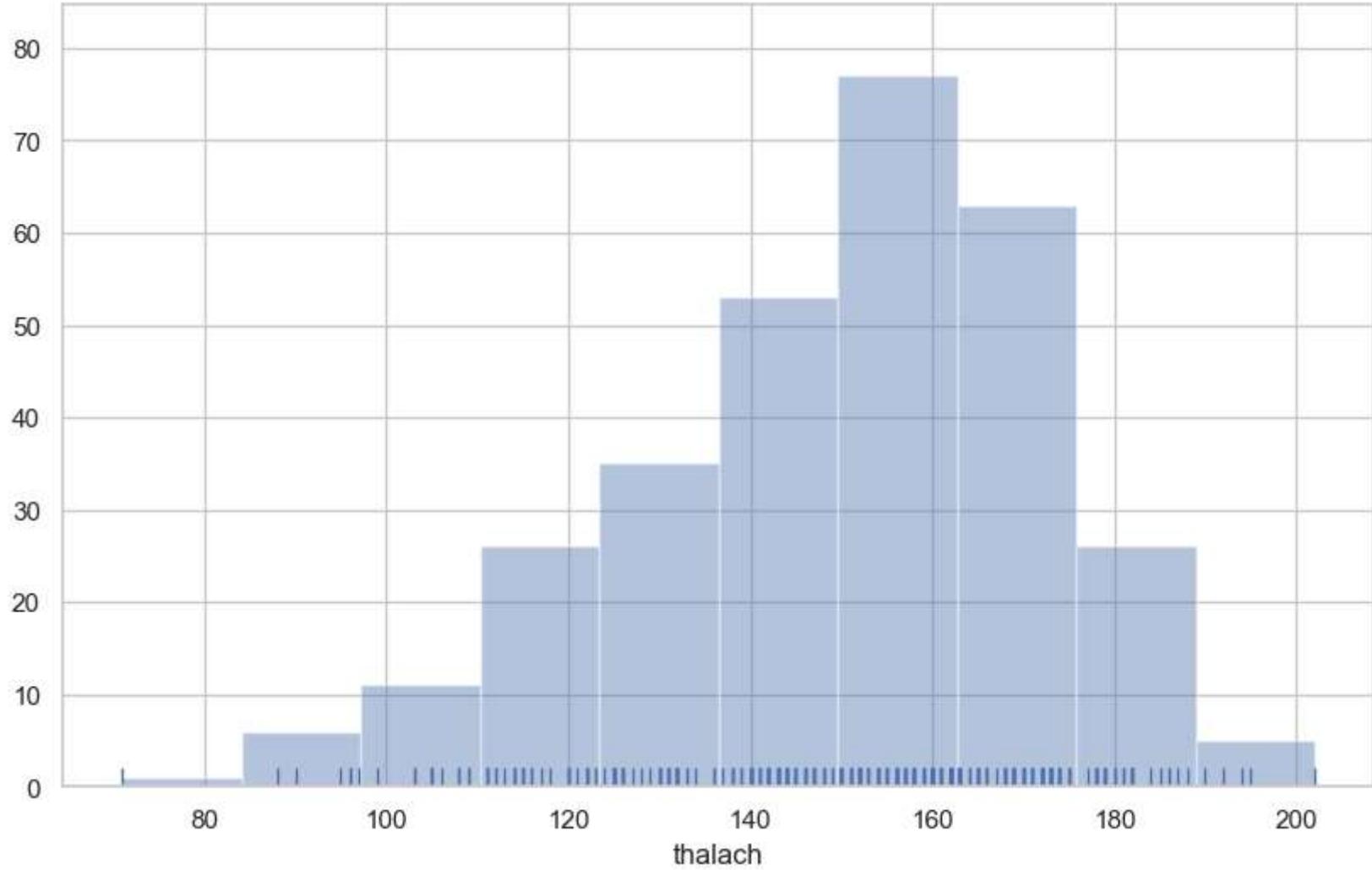
```
In [49]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.kdeplot(x, shade=True, color='r')
plt.show()
```



Histogram

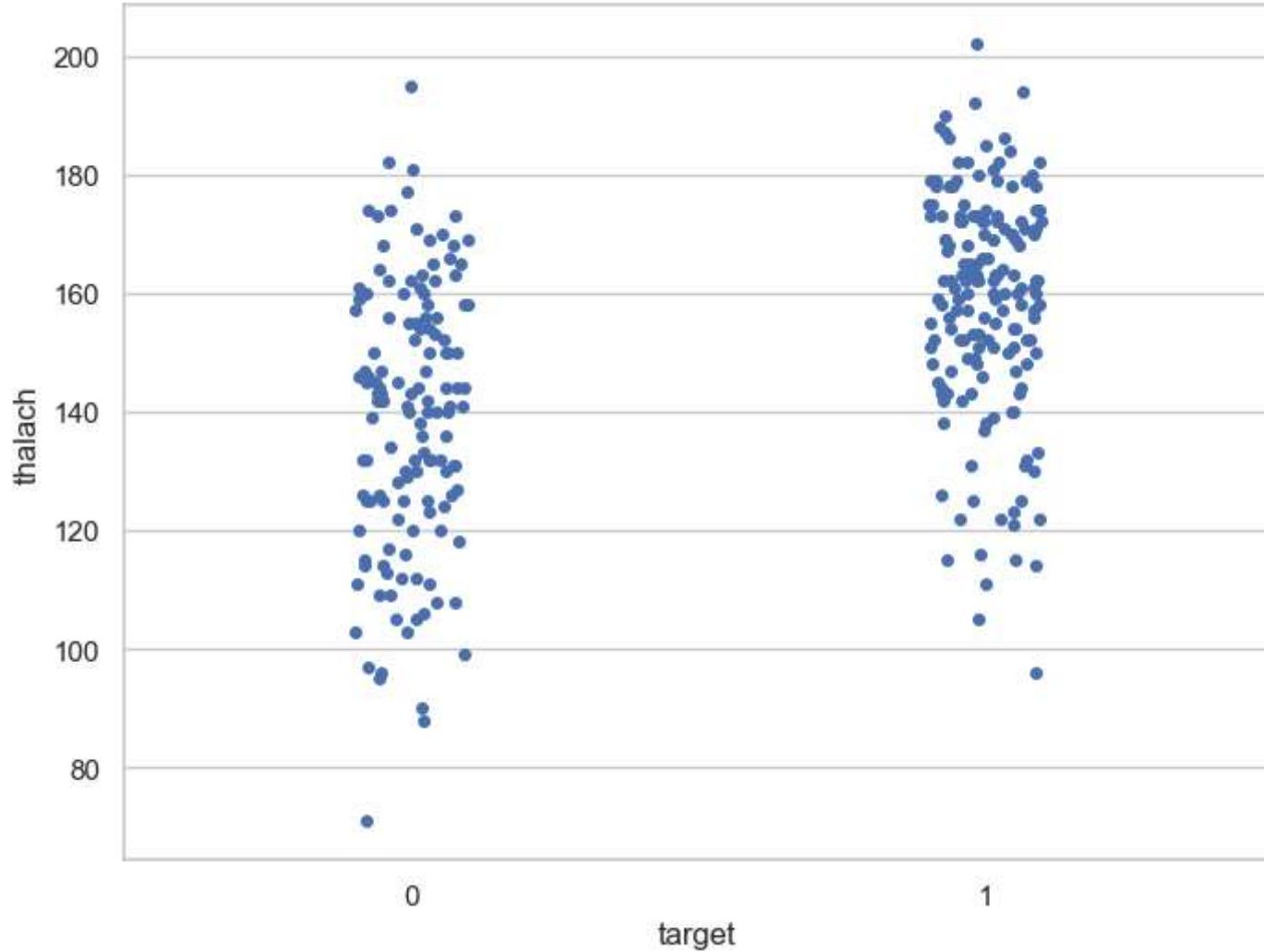
- We can plot a histogram as follows :

```
In [50]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
ax = sns.distplot(x, kde=False, rug=True, bins=10)
plt.show()
```



Visualize frequency distribution of `thalach` variable wrt `target`

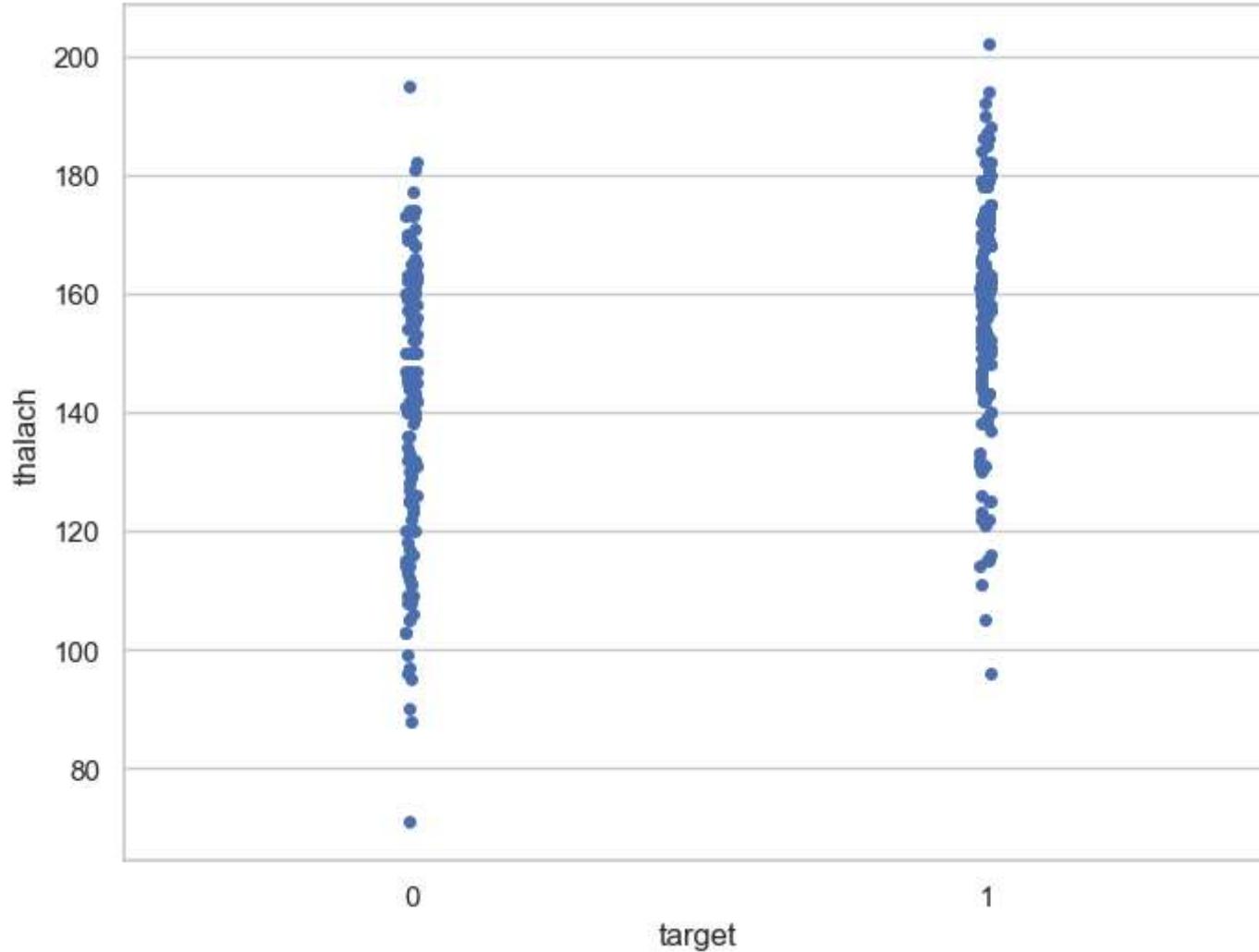
```
In [52]: f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", y="thalach", data=df)
plt.show()
```



Interpretation

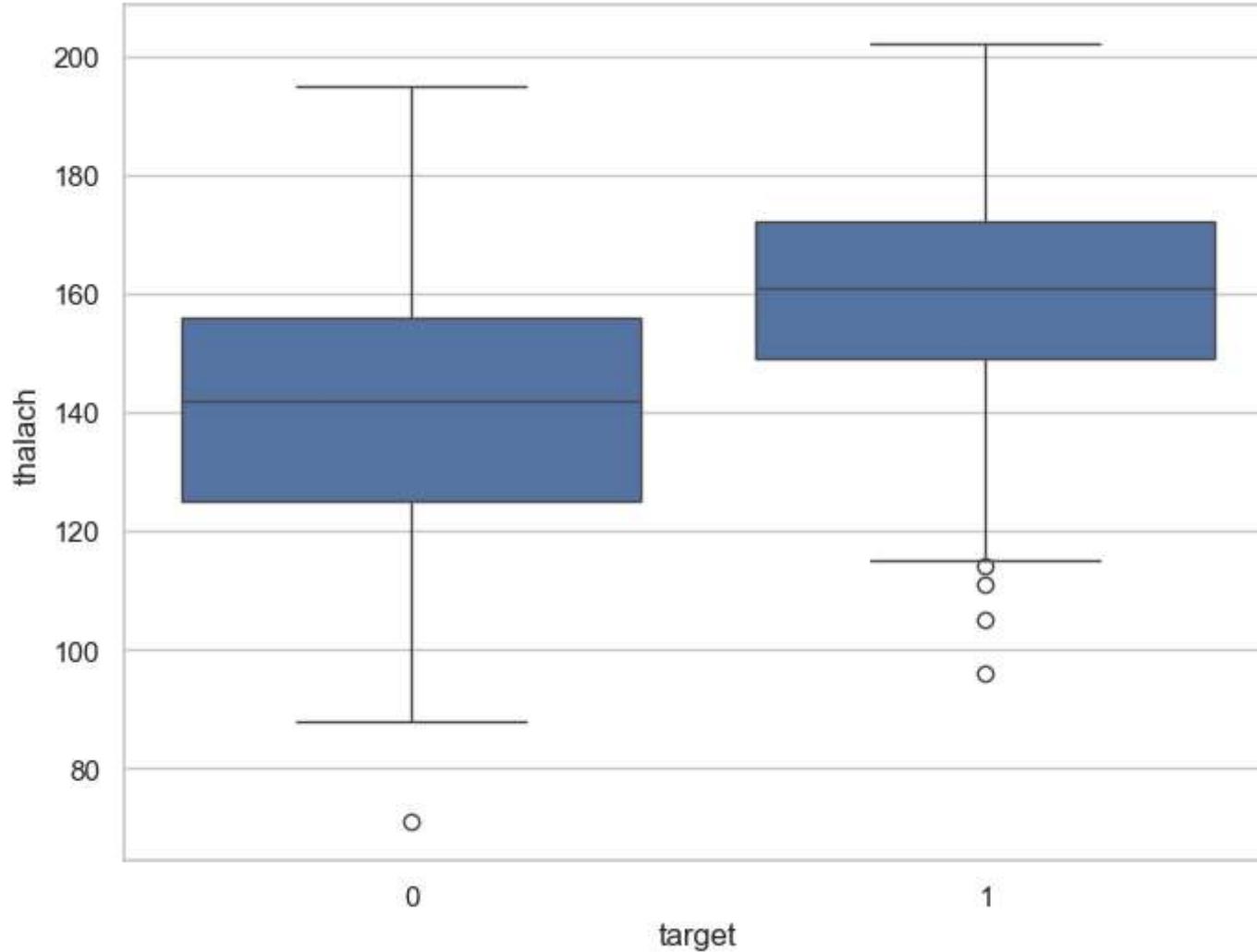
We can add jitter to bring out the distribution of values as follows :

```
In [53]: f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", y="thalach", data=df, jitter = 0.01)
plt.show()
```



Visualize distribution of `thalach` variable wrt `target` with boxplot

```
In [54]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x="target", y="thalach", data=df)
plt.show()
```



Multivariate analysis

Heat Map

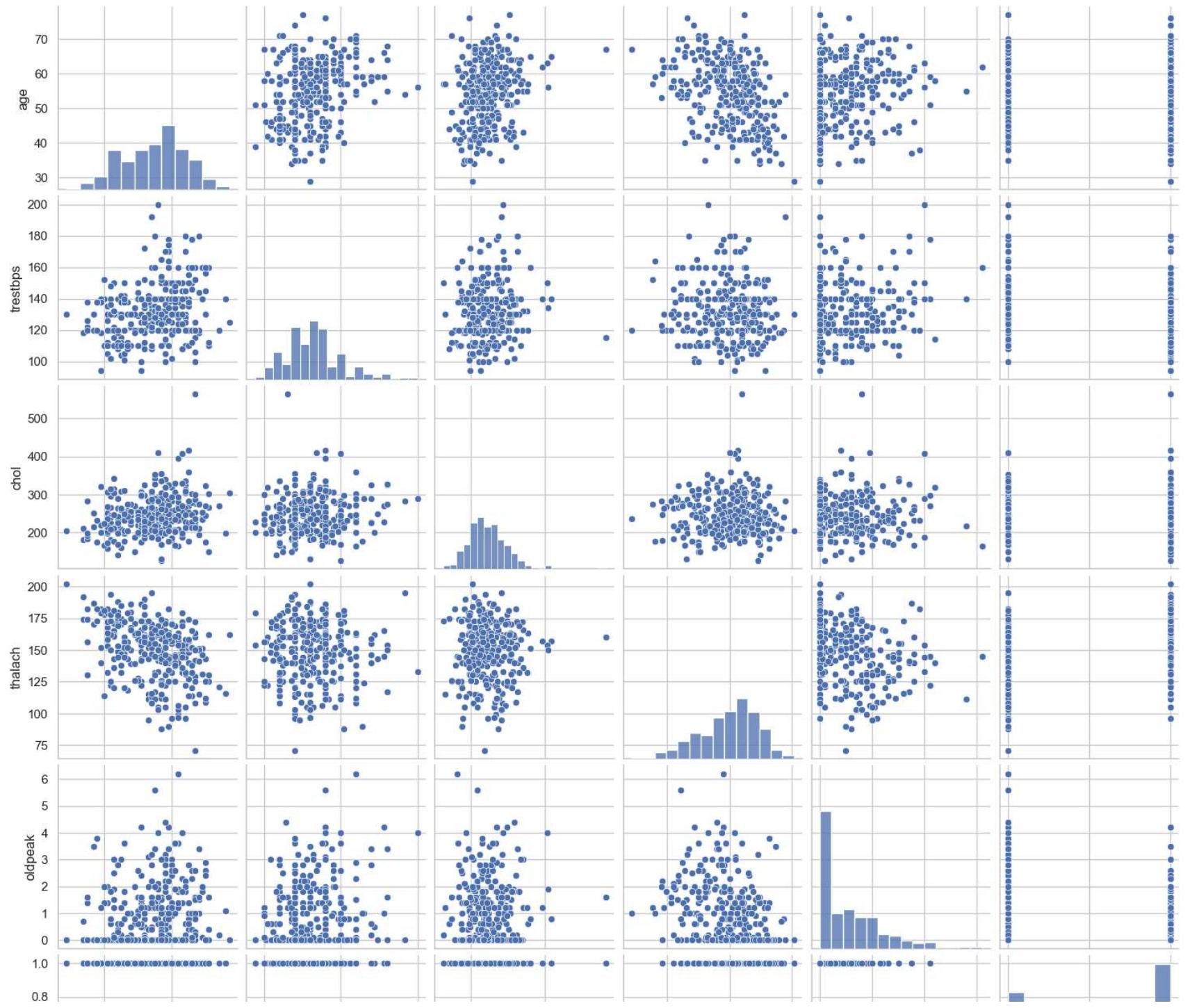
```
In [55]: plt.figure(figsize=(16,12))
plt.title('Correlation Heatmap of Heart Disease Dataset')
a = sns.heatmap(correlation, square=True, annot=True, fmt='.2f', linecolor='white')
```

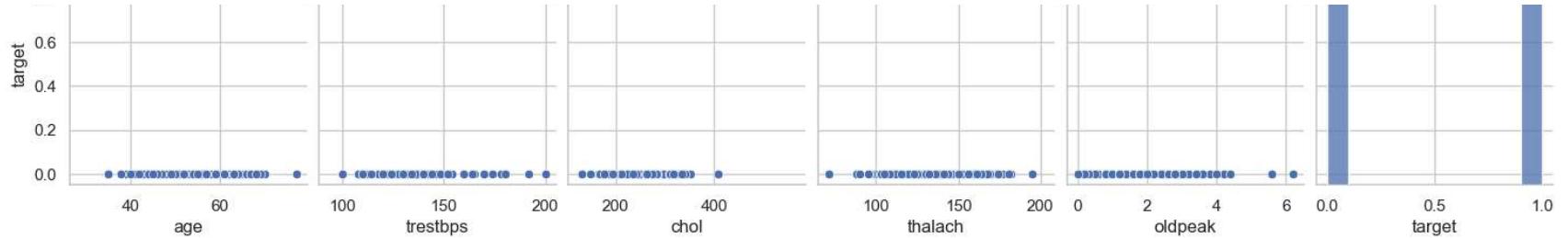
```
a.set_xticklabels(a.get_xticklabels(), rotation=90)
a.set_yticklabels(a.get_yticklabels(), rotation=30)
plt.show()
```



Pair Plot

```
In [56]: num_var = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak', 'target' ]  
sns.pairplot(df[num_var], kind='scatter', diag_kind='hist')  
plt.show()
```





Analysis of age and other variables

Check the number of unique values in age variable

```
In [57]: df['age'].nunique()
```

```
Out[57]: 41
```

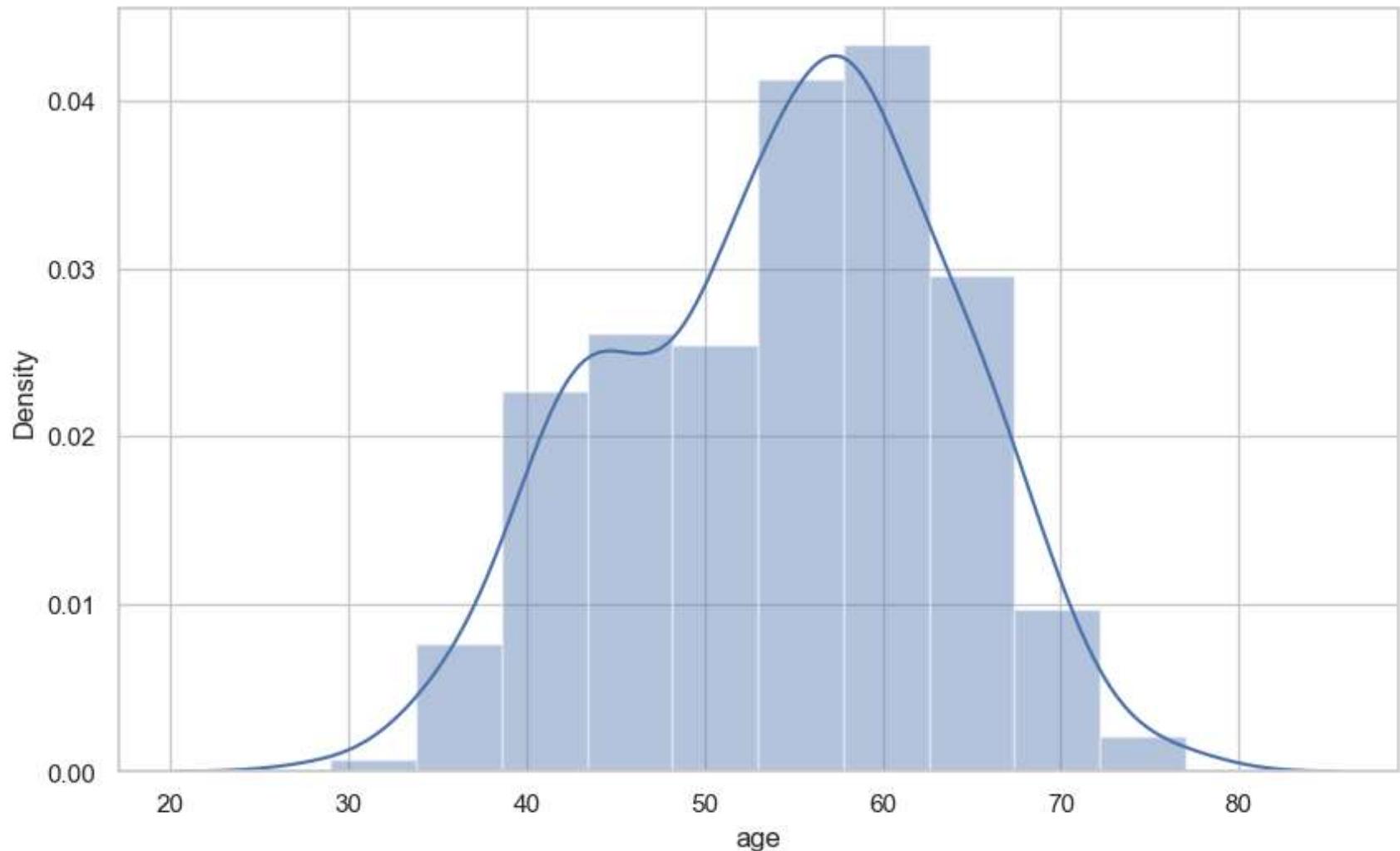
View statistical summary of age variable

```
In [58]: df['age'].describe()
```

```
Out[58]: count    303.000000
mean      54.366337
std       9.082101
min      29.000000
25%     47.500000
50%     55.000000
75%     61.000000
max     77.000000
Name: age, dtype: float64
```

Plot the distribution of age variable

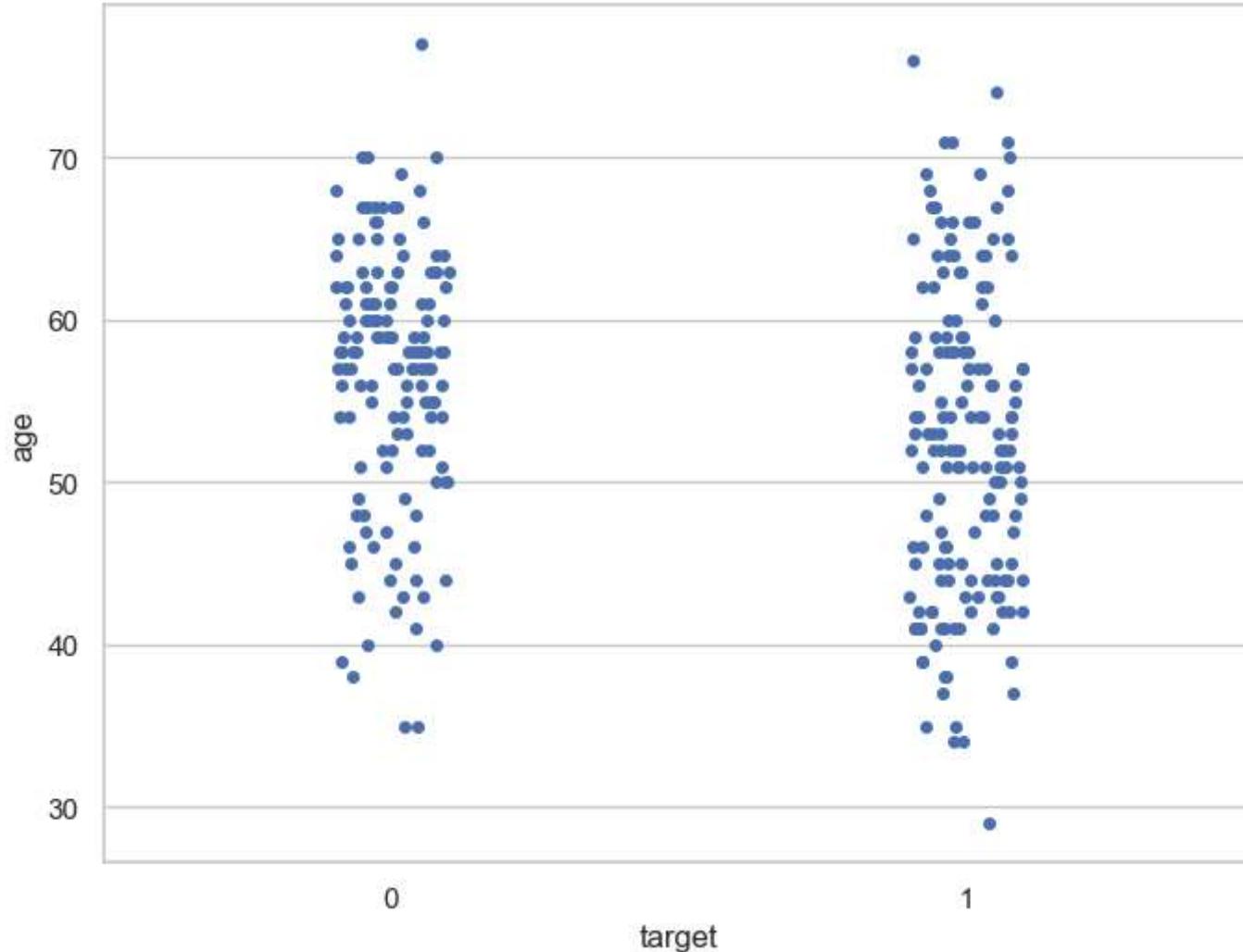
```
In [59]: f, ax = plt.subplots(figsize=(10,6))
x = df['age']
ax = sns.distplot(x, bins=10)
plt.show()
```



Analyze age and target variable

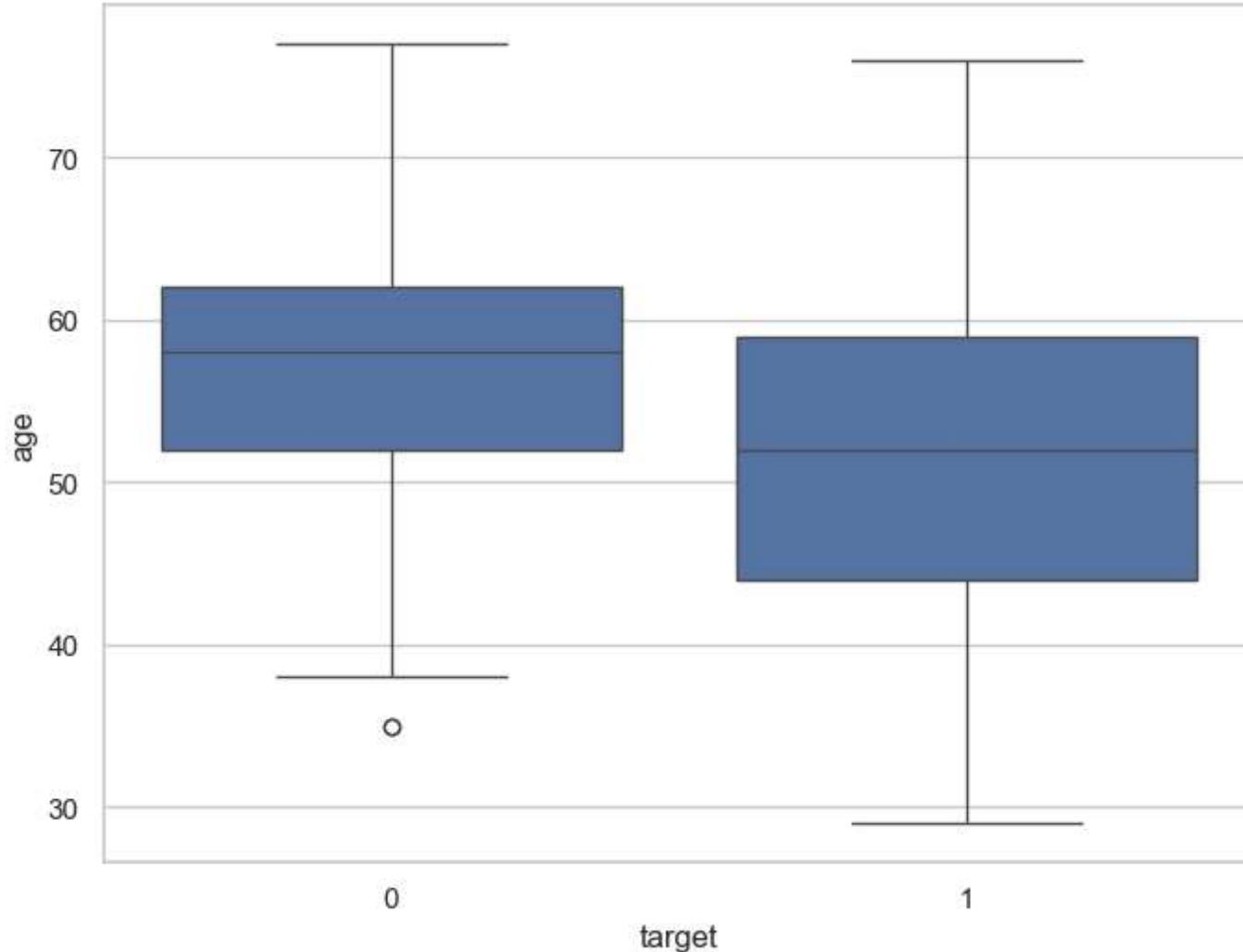
Visualize frequency distribution of age variable wrt target

```
In [60]: f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", y="age", data=df)
plt.show()
```



Visualize distribution of `age` variable wrt `target` with boxplot

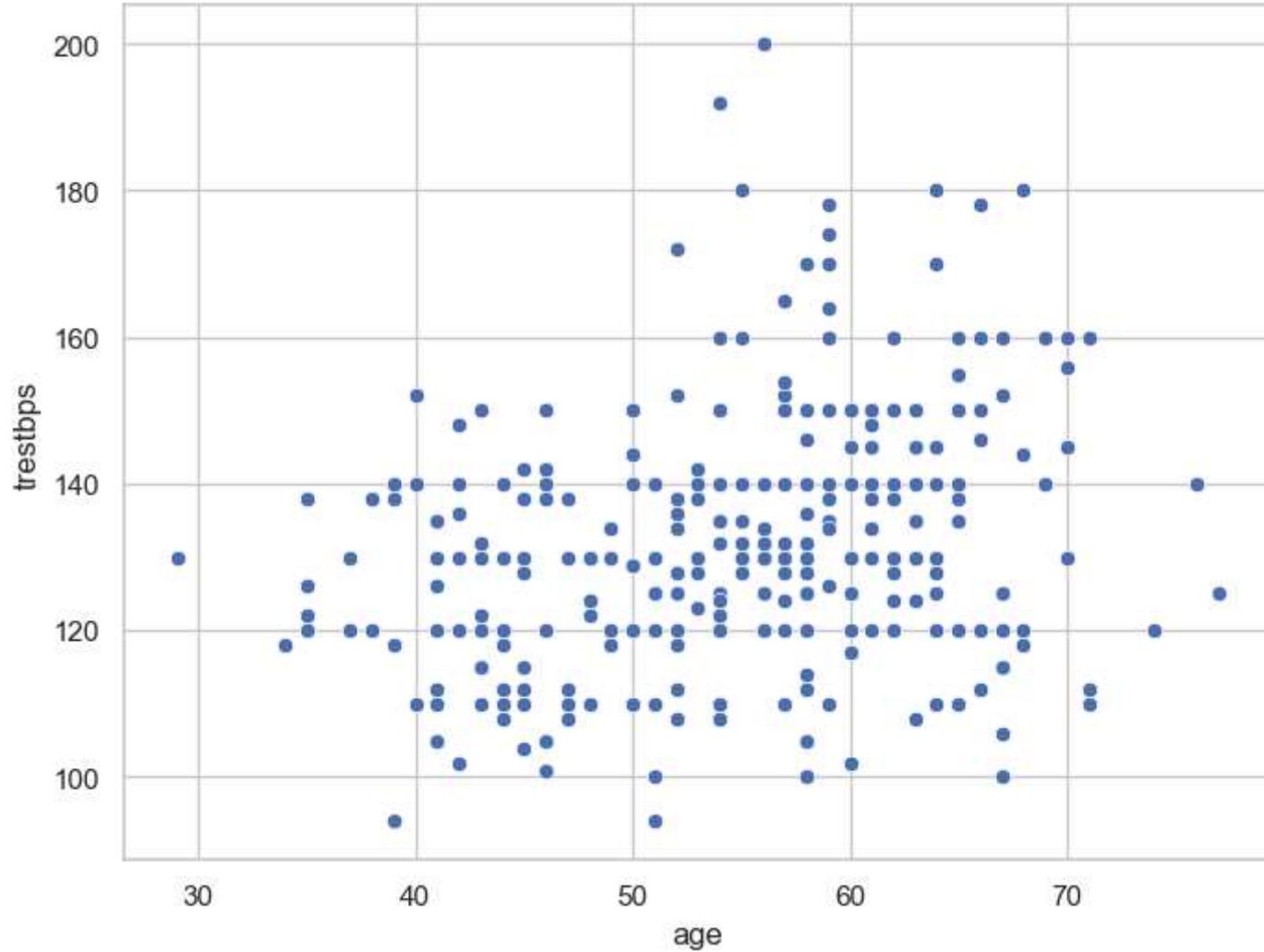
```
In [61]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x="target", y="age", data=df)
plt.show()
```



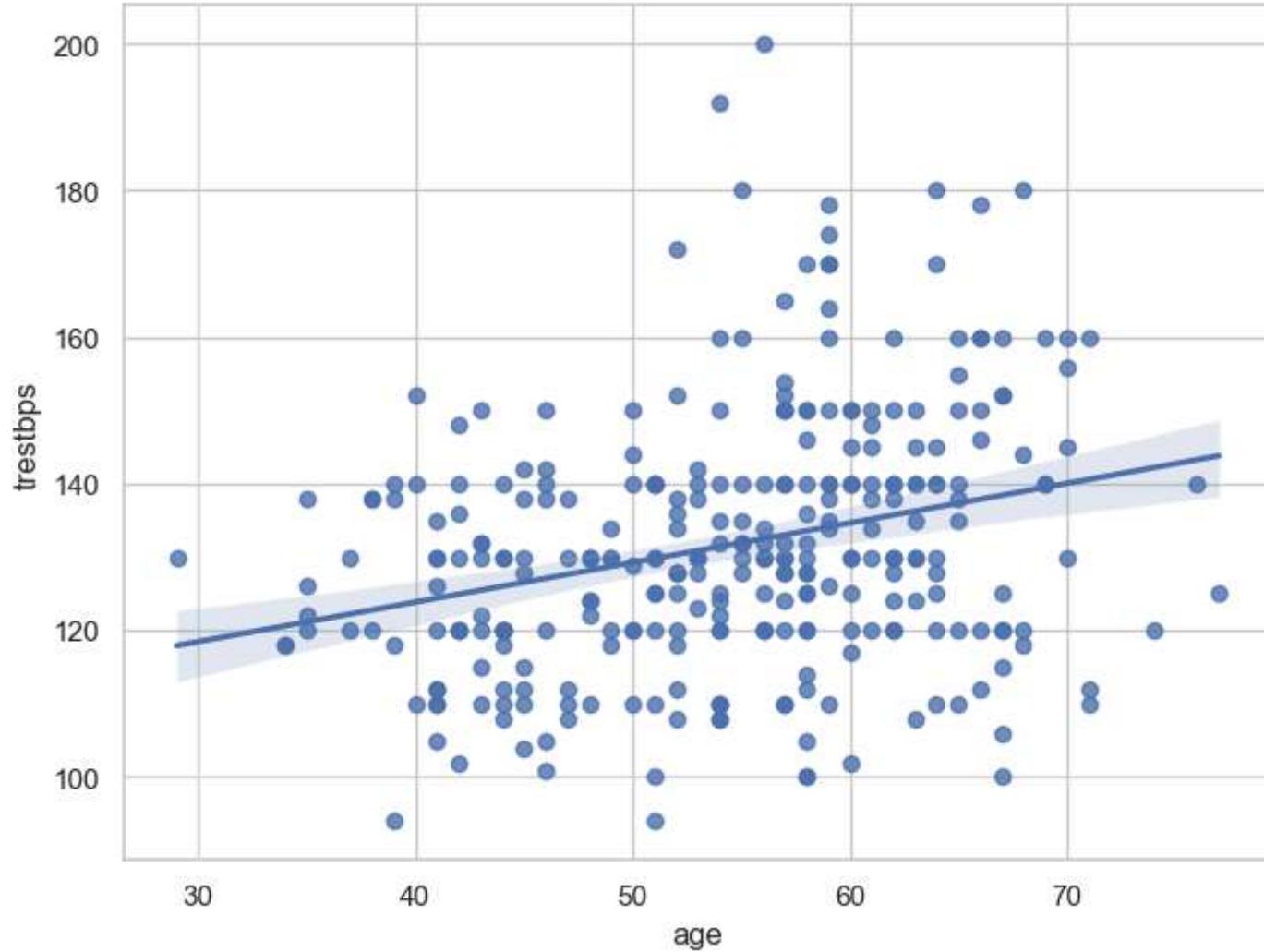
Analyze `age` and `trestbps` variable

I will plot a scatterplot to visualize the relationship between `age` and `trestbps` variable.

```
In [62]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.scatterplot(x="age", y="trestbps", data=df)
plt.show()
```

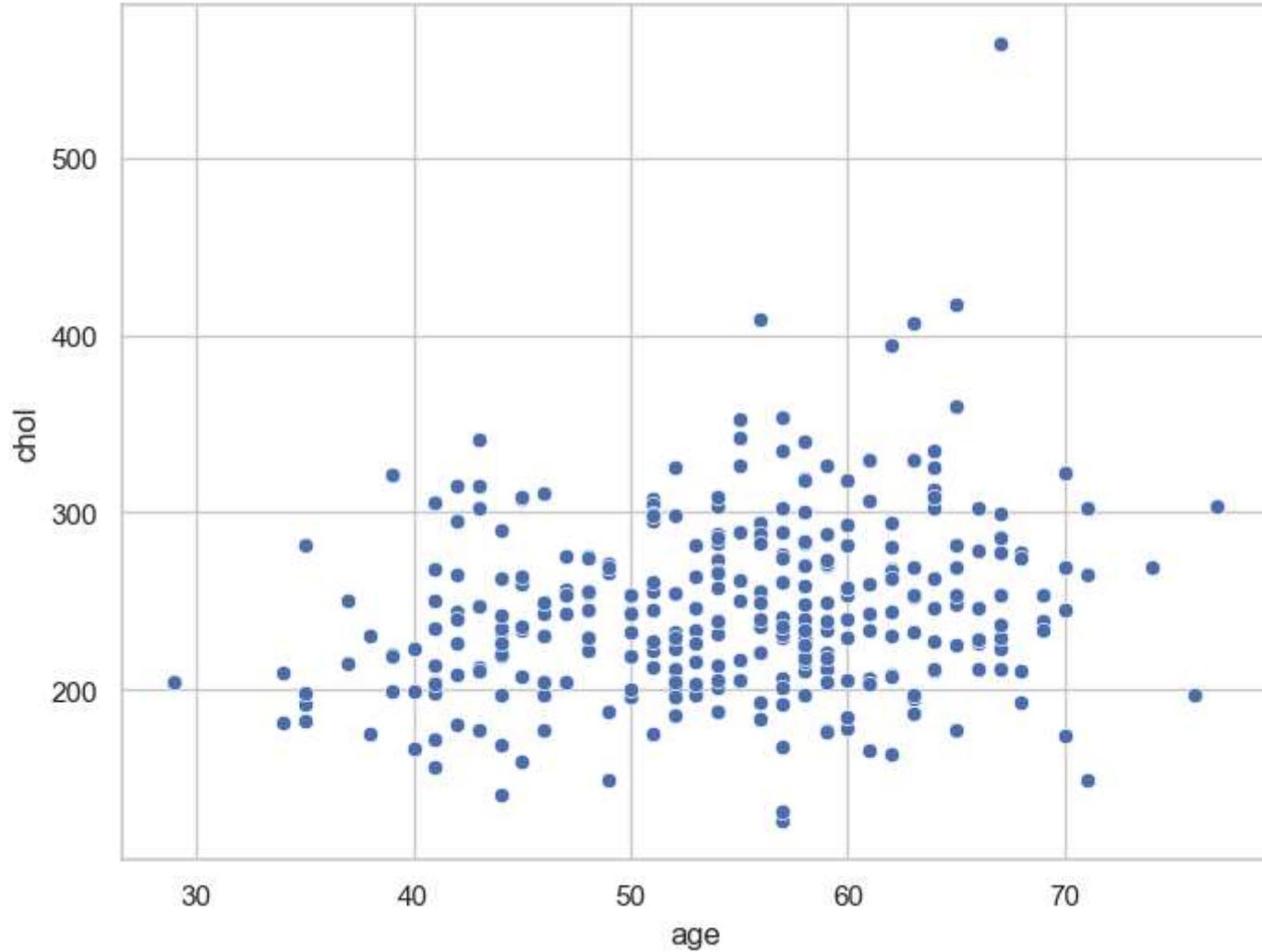


```
In [63]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.regplot(x="age", y="trestbps", data=df)
plt.show()
```

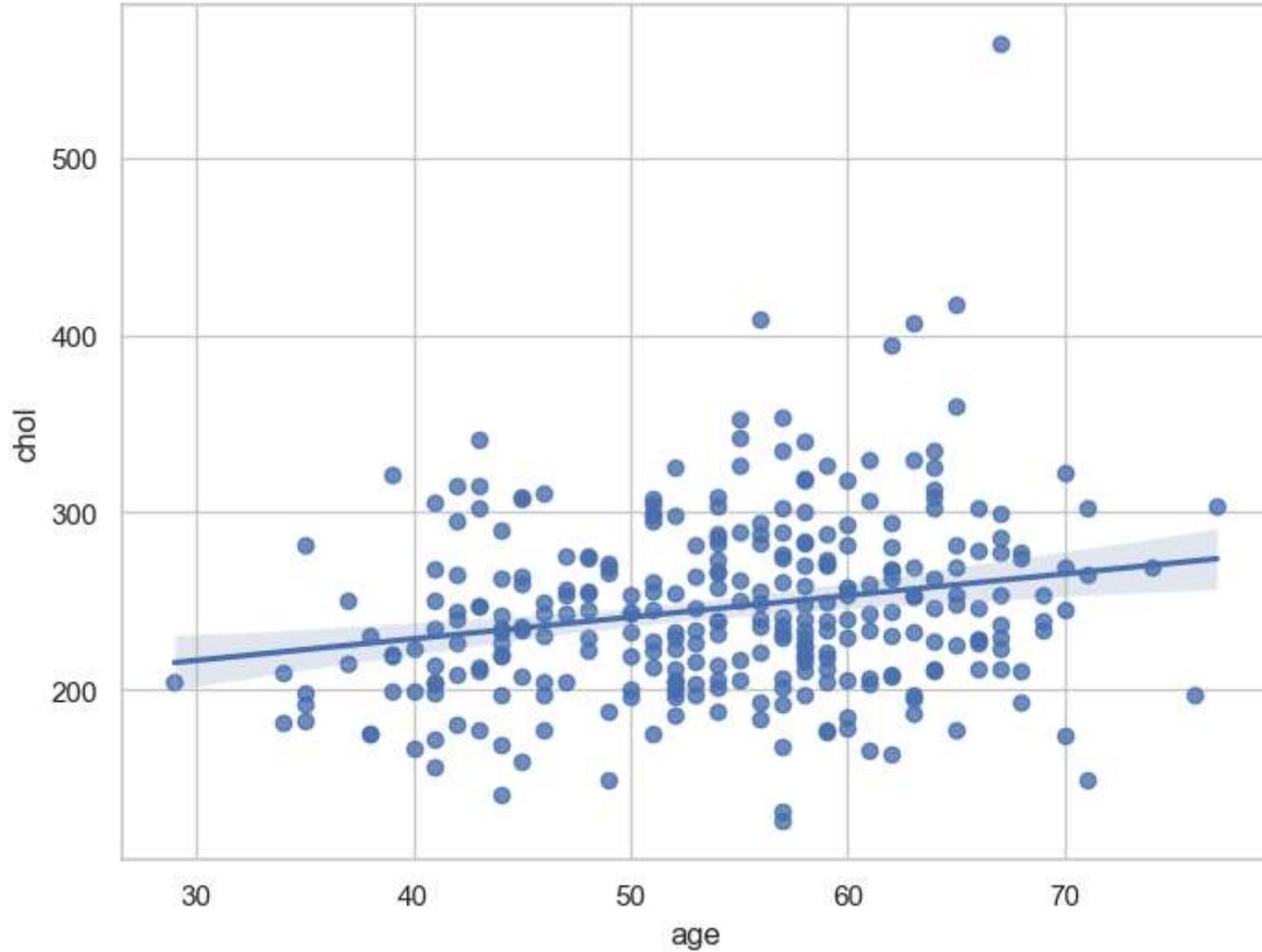


Analyze age and chol variable

```
In [64]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.scatterplot(x="age", y="chol", data=df)
plt.show()
```

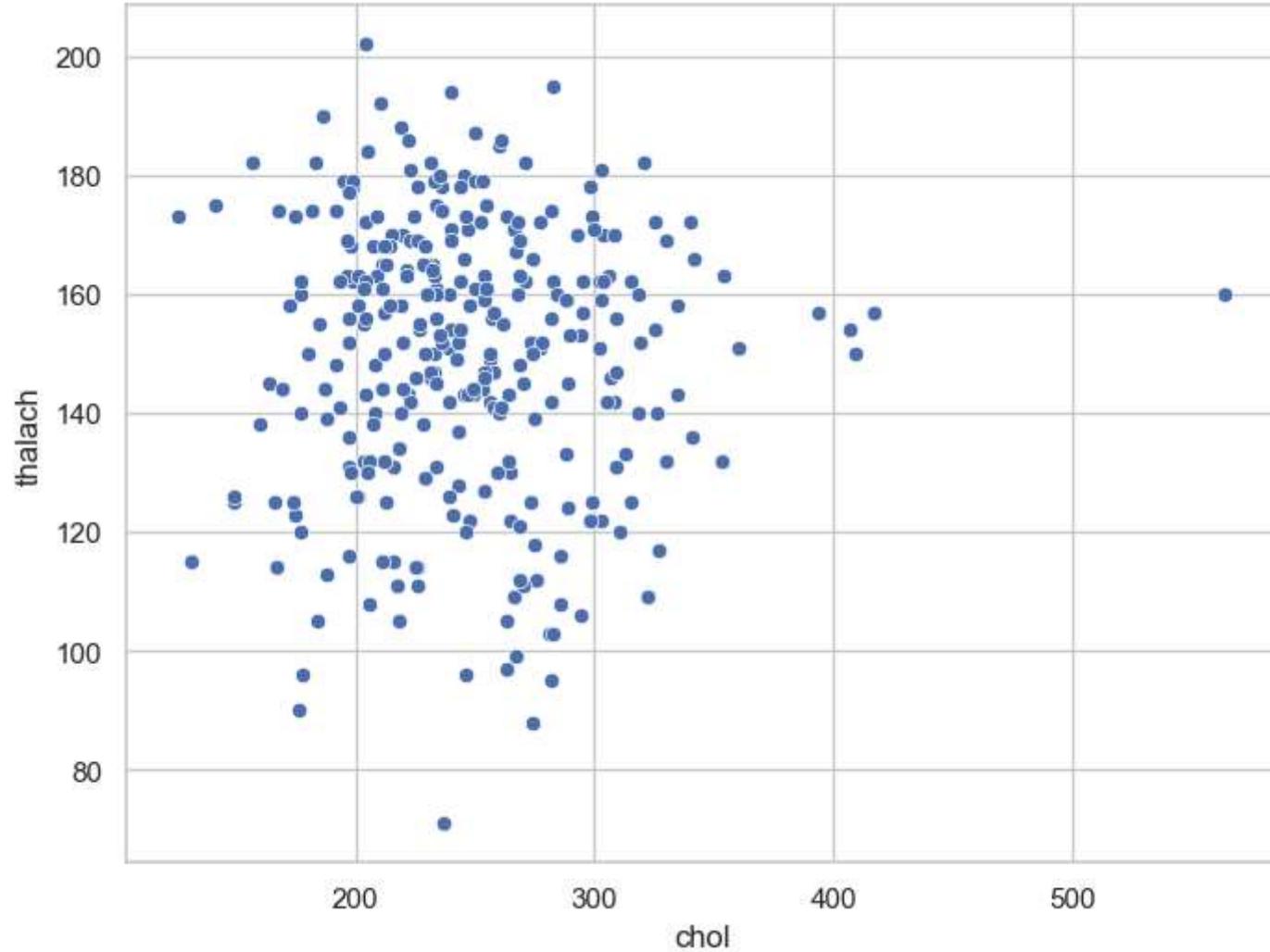


```
In [65]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.regplot(x="age", y="chol", data=df)
plt.show()
```

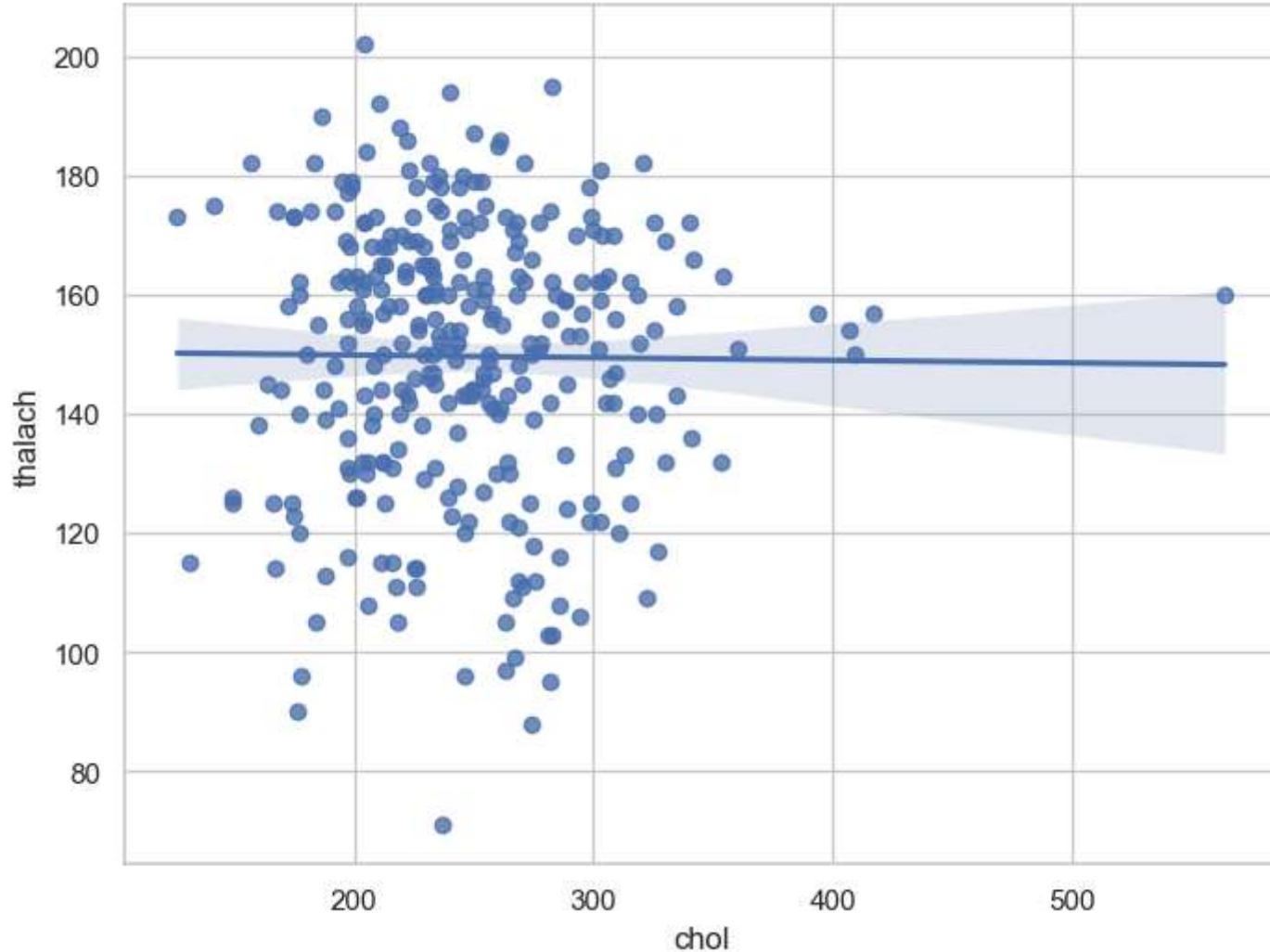


Analyze chol and thalach variable

```
In [66]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.scatterplot(x="chol", y = "thalach", data=df)
plt.show()
```



```
In [67]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.regplot(x="chol", y="thalach", data=df)
plt.show()
```



Dealing with missing values

Pandas isnull() and notnull() functions

```
In [68]: # check for missing values
```

```
df.isnull().sum()
```

```
Out[68]: 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
```

Check with ASSERT statement

```
In [69]: #assert that there are no missing values in the dataframe  
  
        assert pd.notnull(df).all().all()
```

```
In [70]: #assert all values are greater than or equal to 0  
  
        assert (df >= 0).all().all()
```

12. Outlier detection

I will make boxplots to visualise outliers in the continuous numerical variables :-

age , trestbps , chol , thalach and oldpeak variables.

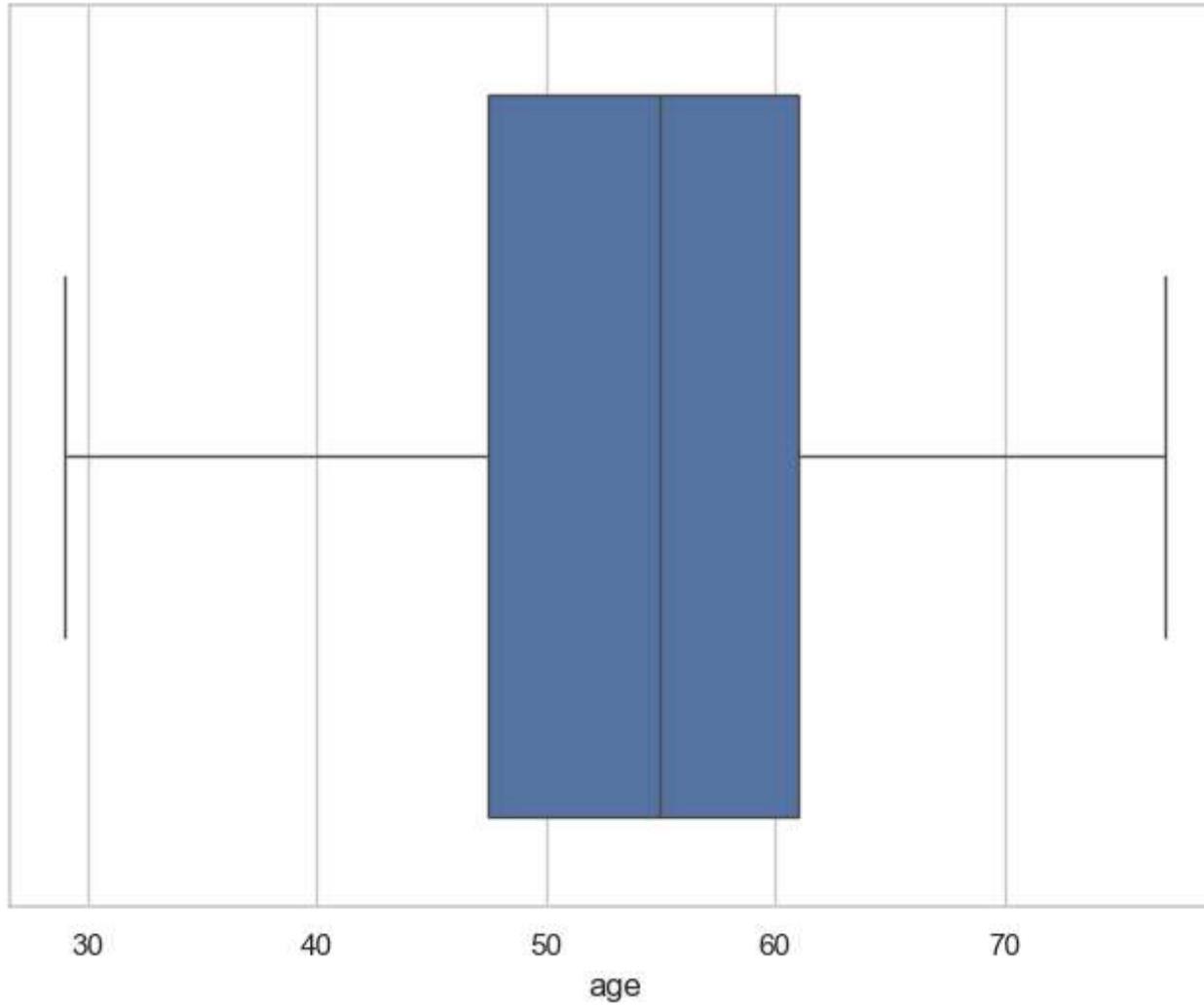
age variable

```
In [71]: df['age'].describe()
```

```
Out[71]: count    303.000000
          mean     54.366337
          std      9.082101
          min     29.000000
          25%    47.500000
          50%    55.000000
          75%    61.000000
          max     77.000000
Name: age, dtype: float64
```

Box-plot of age variable

```
In [72]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["age"])
plt.show()
```



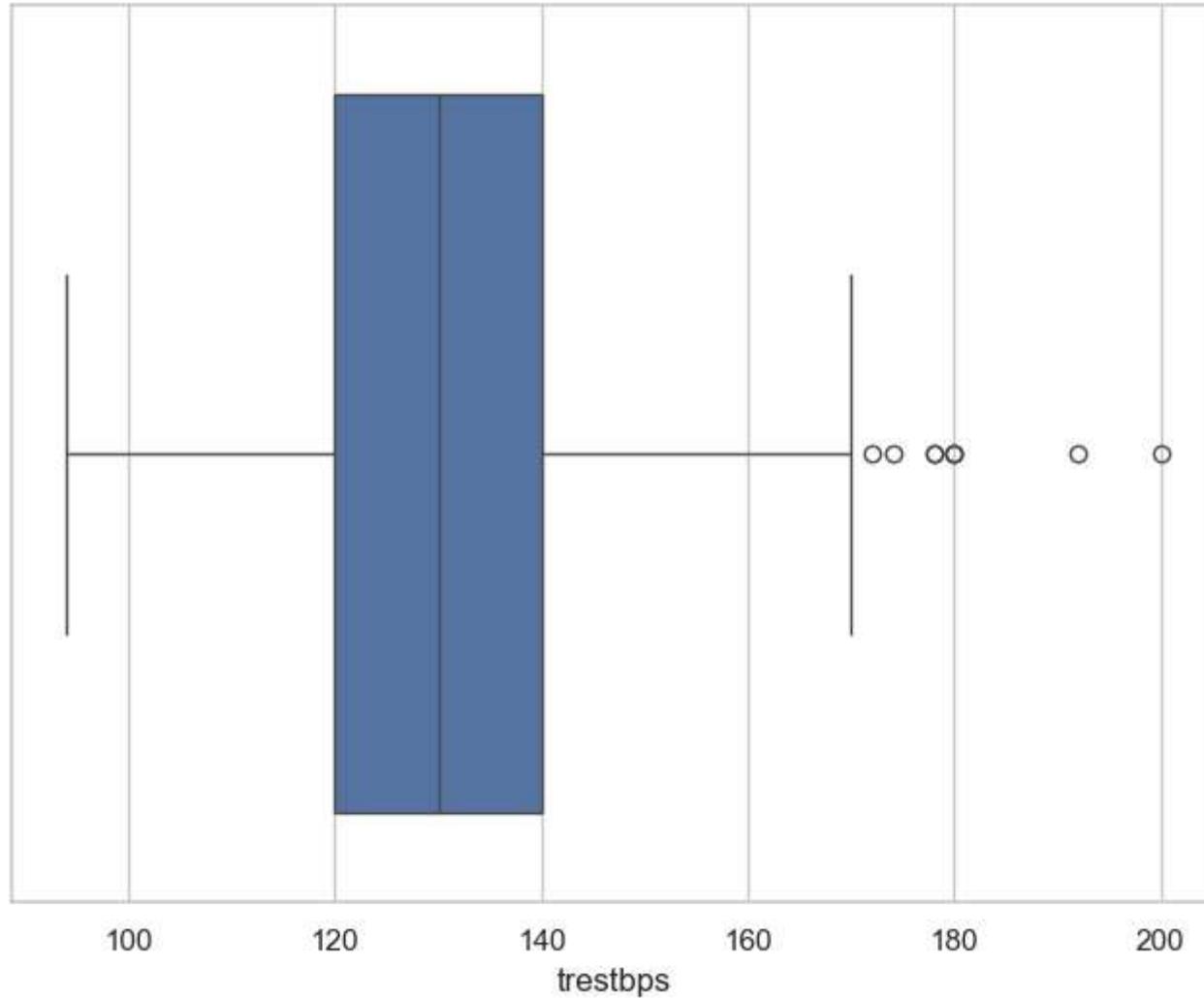
trestbps variable

```
In [73]: df['trestbps'].describe()
```

```
Out[73]: count    303.000000
          mean     131.623762
          std      17.538143
          min      94.000000
          25%     120.000000
          50%     130.000000
          75%     140.000000
          max     200.000000
Name: trestbps, dtype: float64
```

Box-plot of trestbps variable

```
In [74]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["trestbps"])
plt.show()
```



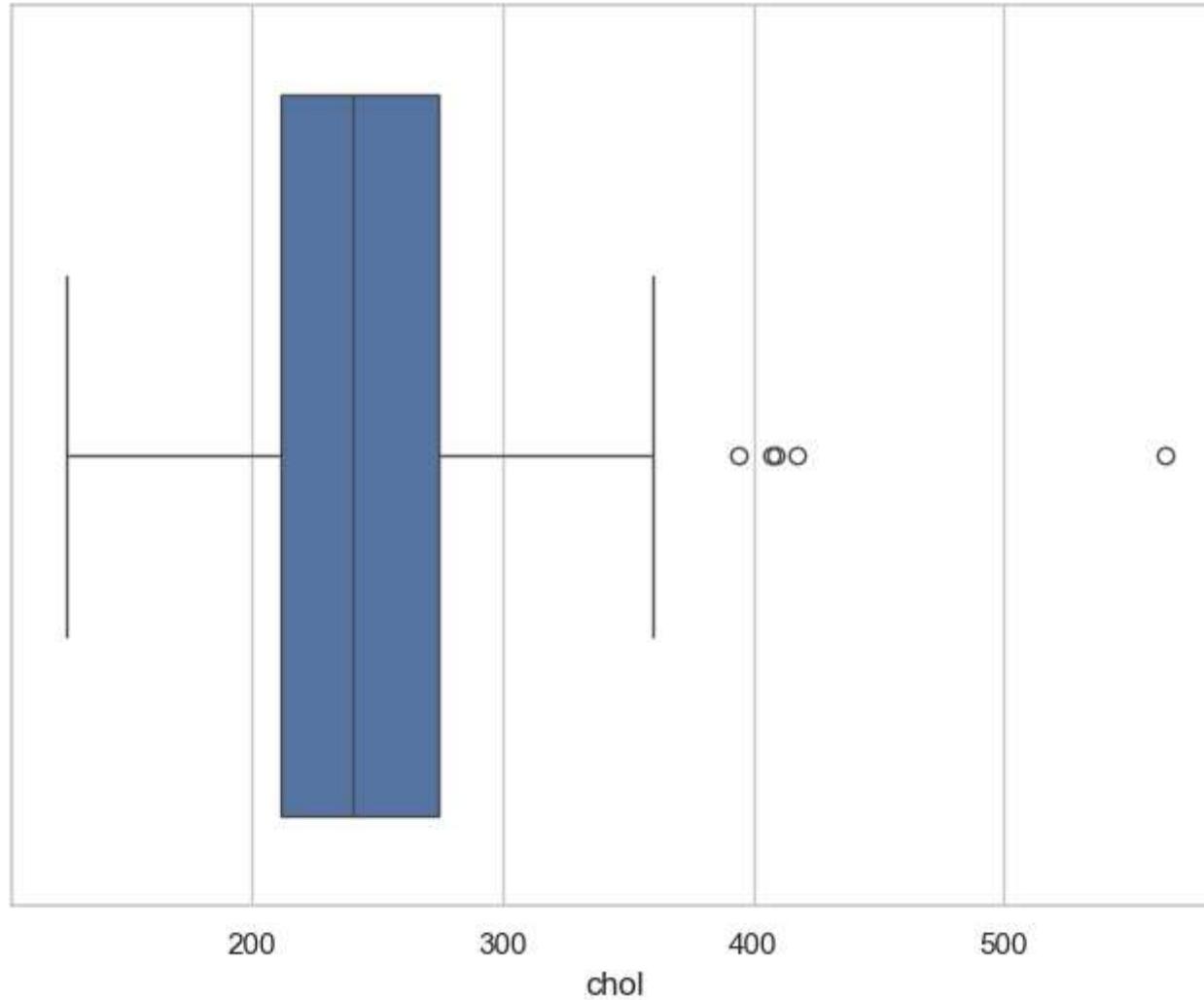
chol variable

```
In [75]: df['chol'].describe()
```

```
Out[75]: count    303.000000
          mean     246.264026
          std      51.830751
          min     126.000000
          25%    211.000000
          50%    240.000000
          75%    274.500000
          max     564.000000
Name: chol, dtype: float64
```

Box-plot of chol variable

```
In [76]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["chol"])
plt.show()
```



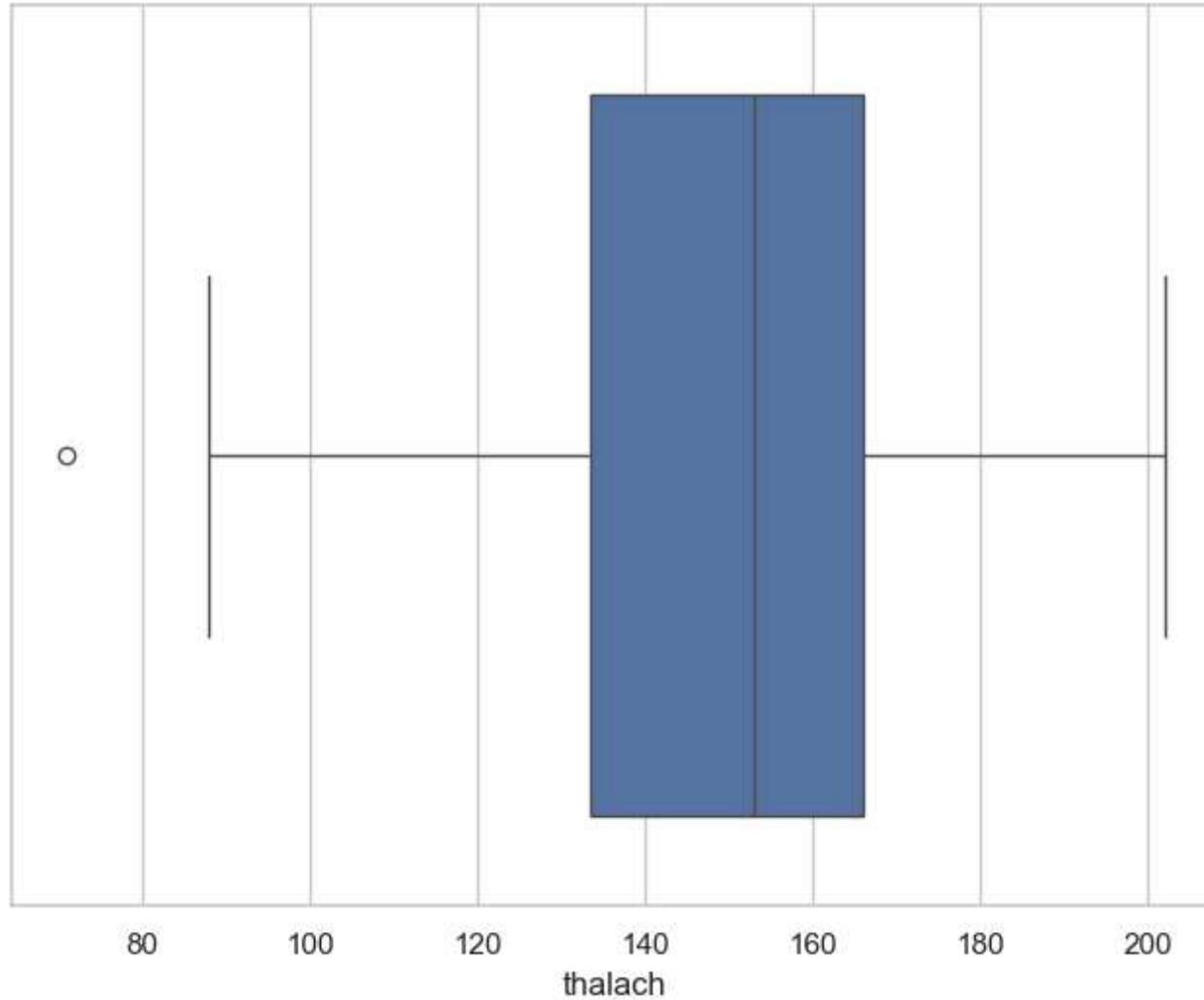
thalach variable

```
In [77]: df['thalach'].describe()
```

```
Out[77]: count    303.000000
          mean     149.646865
          std      22.905161
          min      71.000000
          25%     133.500000
          50%     153.000000
          75%     166.000000
          max     202.000000
Name: thalach, dtype: float64
```

Box-plot of thalach variable

```
In [78]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["thalach"])
plt.show()
```



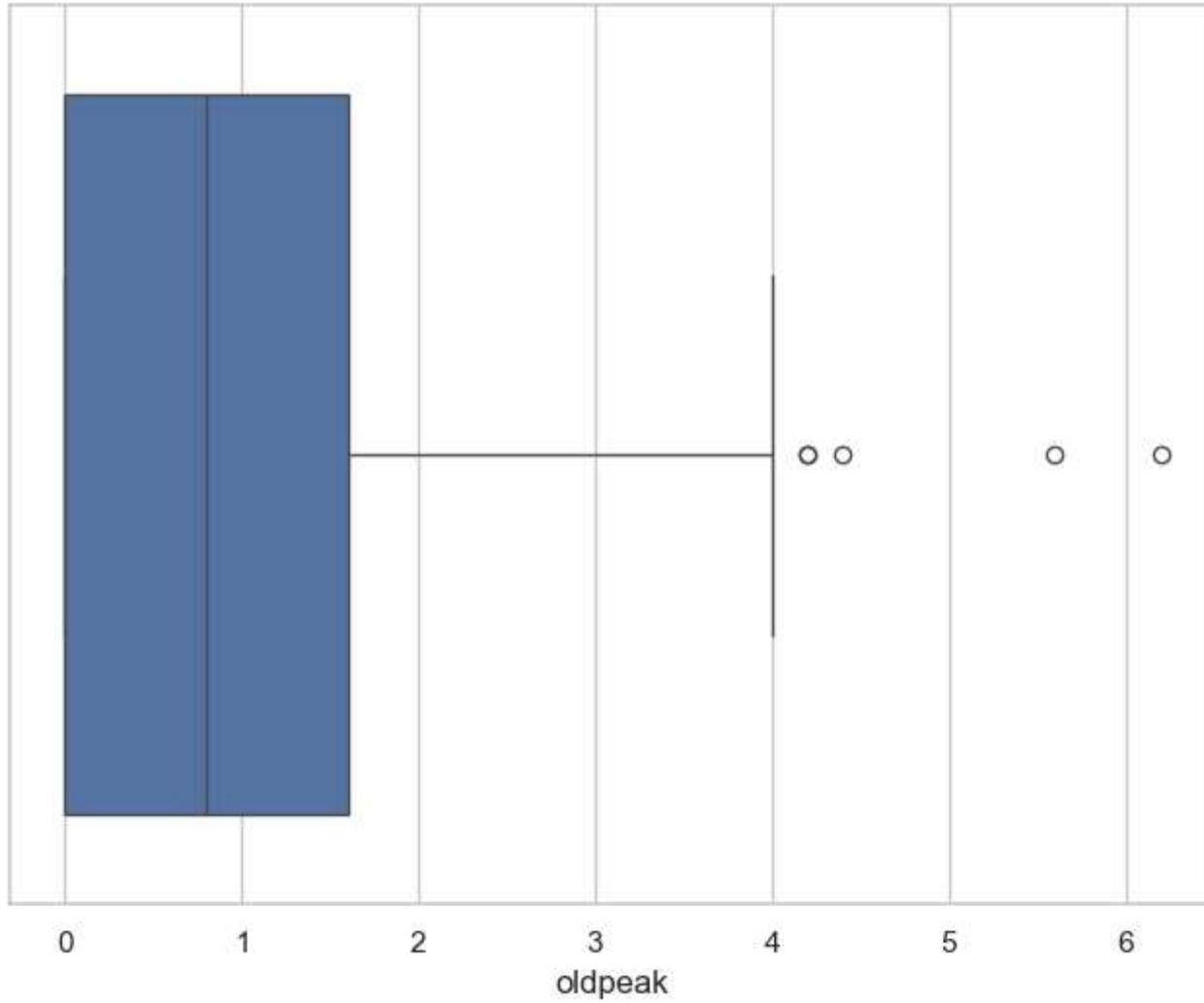
oldpeak variable

```
In [79]: df['oldpeak'].describe()
```

```
Out[79]: count    303.000000
          mean     1.039604
          std      1.161075
          min     0.000000
          25%    0.000000
          50%    0.800000
          75%    1.600000
          max     6.200000
Name: oldpeak, dtype: float64
```

Box-plot of oldpeak variable

```
In [80]: f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["oldpeak"])
plt.show()
```



In []: