

```
In [1]: import numpy as np
import pandas as pd
import os
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
In [3]: import seaborn as sns
import matplotlib.pyplot as plt
import scipy.stats as st
%matplotlib inline

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

```
In [5]: import warnings
warnings.filterwarnings('ignore')
```

```
In [7]: df=pd.read_csv(r"D:\NIT Resume Project\heart.csv")
```

```
In [9]: df
```

```
Out[9]:
```

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

303 rows × 14 columns



```
In [11]: df.shape
```

```
Out[11]: (303, 14)
```

```
In [15]: df.head()
```

Out[15]:

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

In [17]: `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 [19]: `df.dtypes`

Out[19]:

```
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 [21]: `df.describe()`

Out[21]:

	age	sex	cp	trestbps	chol	fbs	restecg
<b>count</b>	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000	303.000000
<b>mean</b>	54.366337	0.683168	0.966997	131.623762	246.264026	0.148515	0.528000
<b>std</b>	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525000
<b>min</b>	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000
<b>25%</b>	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000
<b>50%</b>	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000
<b>75%</b>	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000
<b>max</b>	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000

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 [26]: df['target'].nunique()

Out[26]: 2

In [28]: df['target'].unique()

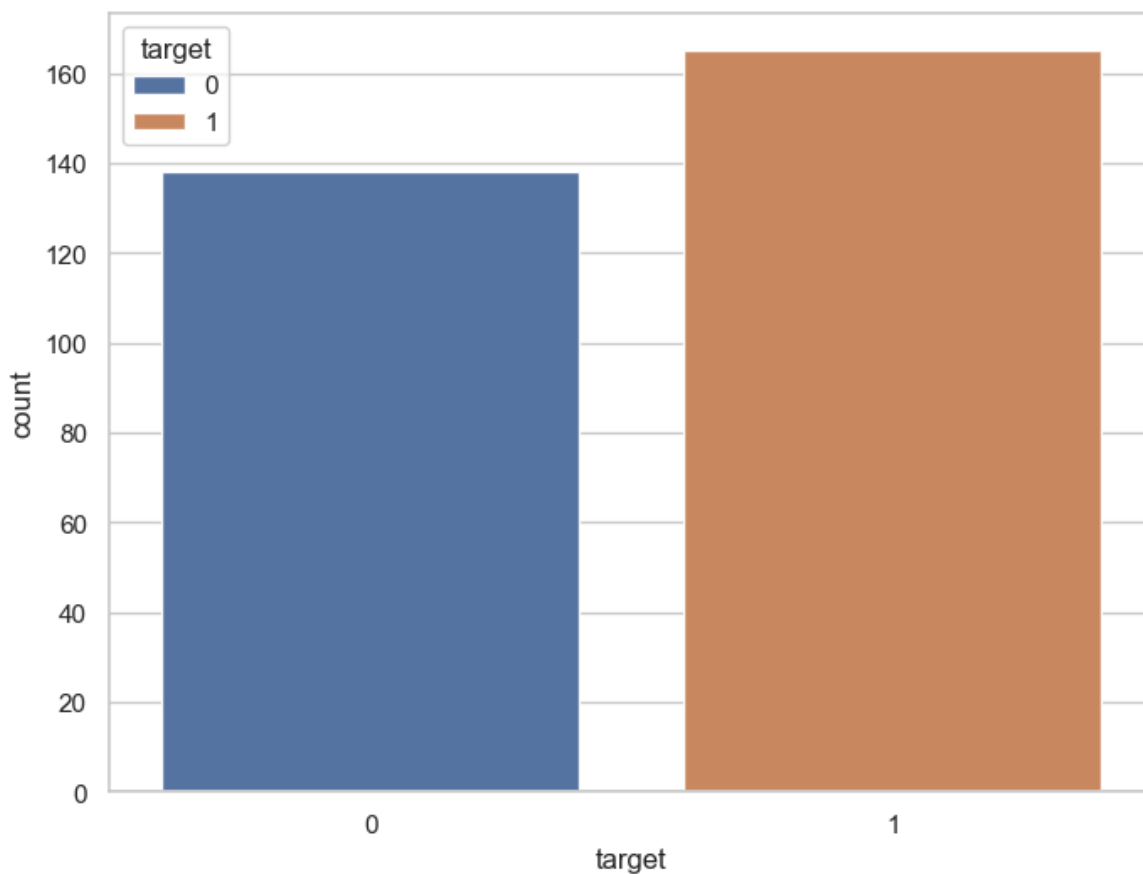
Out[28]: array([1, 0], dtype=int64)

In [30]: df['target'].value\_counts()

Out[30]: target  
 1 165  
 0 138  
 Name: count, dtype: int64

### Visualize frequency distribution of target variable

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



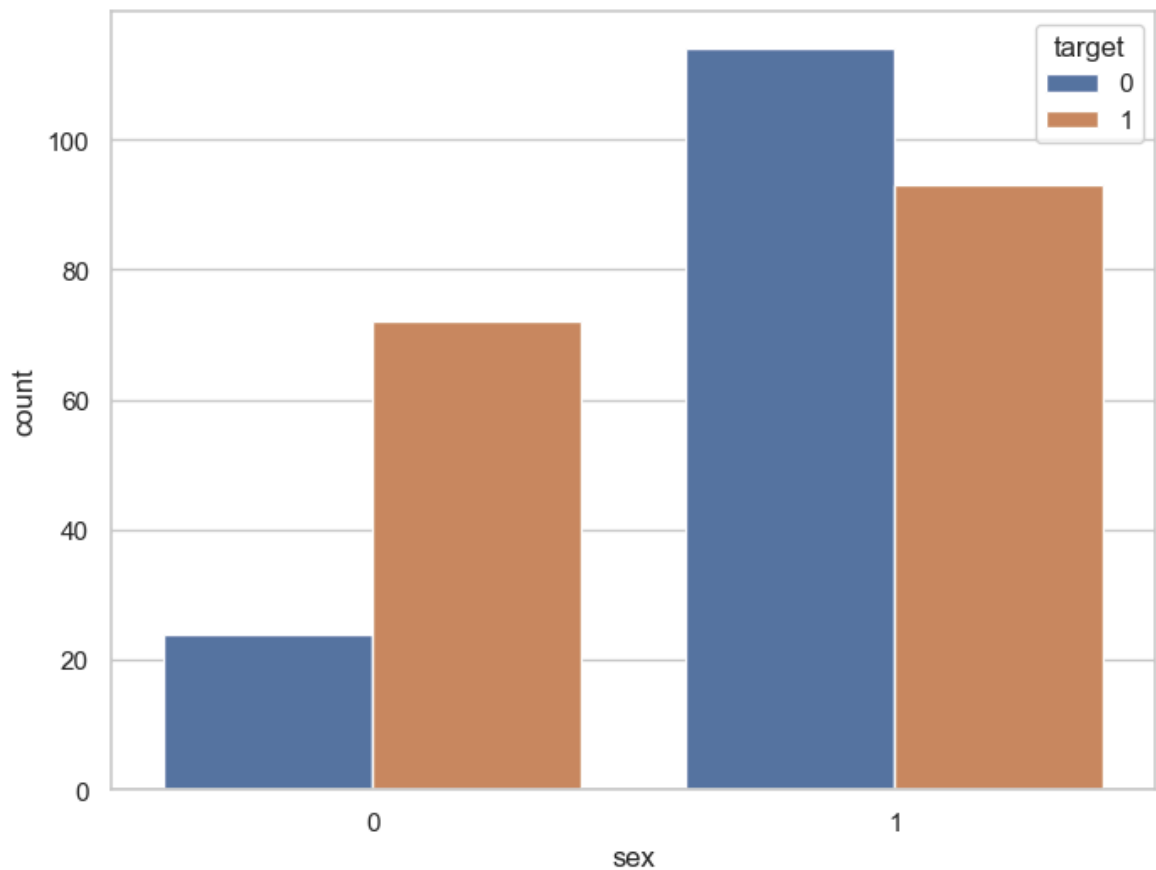
Frequency distribution of **target** variable wrt **sex**

```
In [38]: df.groupby('sex')['target'].value_counts()
```

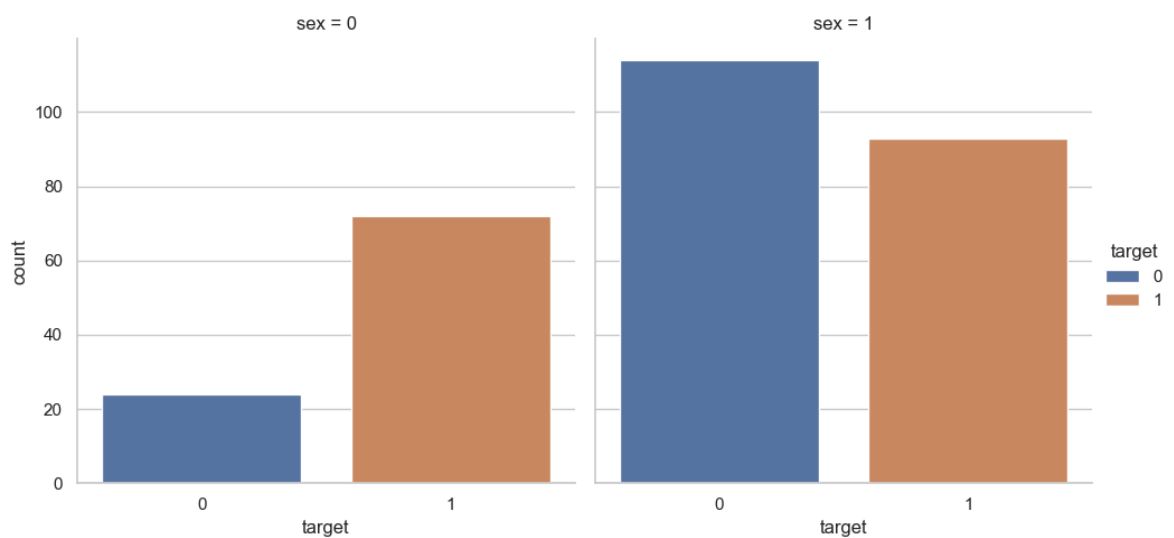
```
Out[38]: sex  target
0      1      72
      0      24
1      0     114
      1      93
Name: count, dtype: int64
```

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

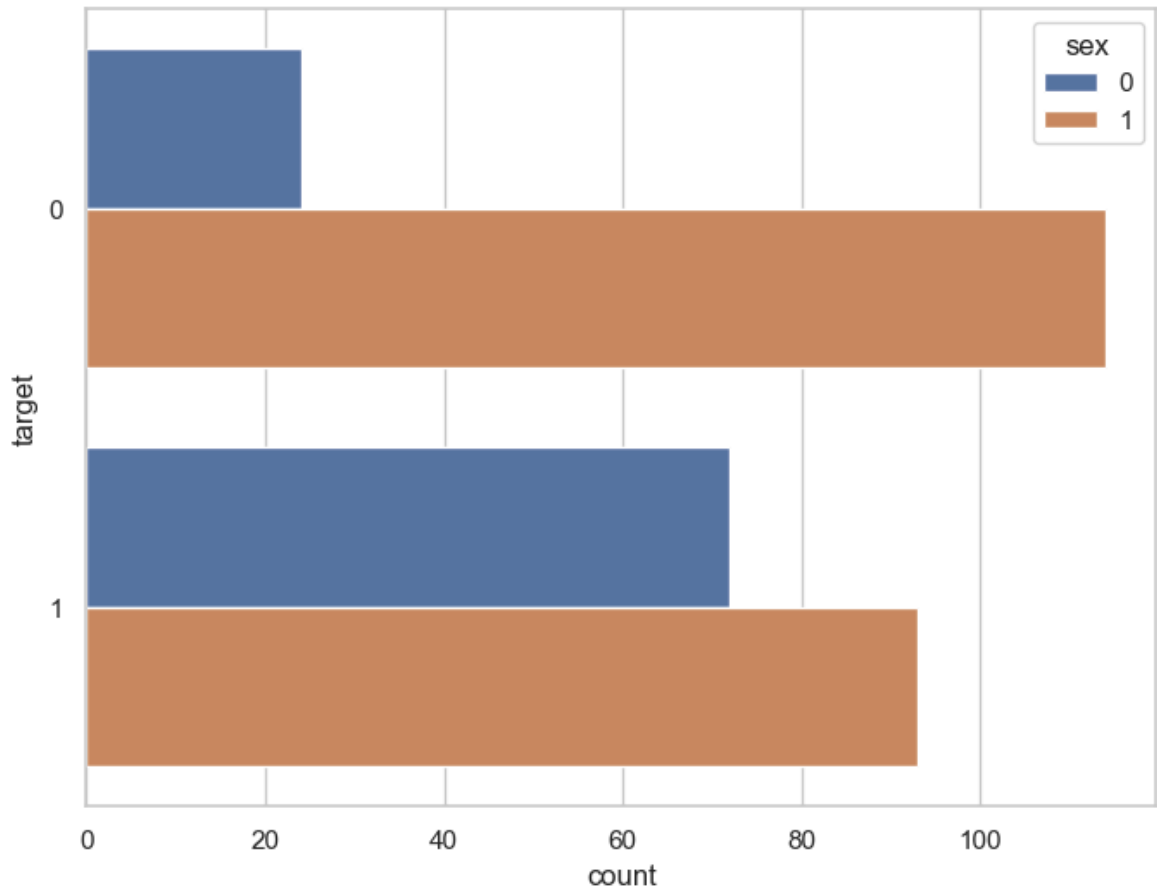
```
In [43]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="sex", hue="target", data=df)
plt.show()
```



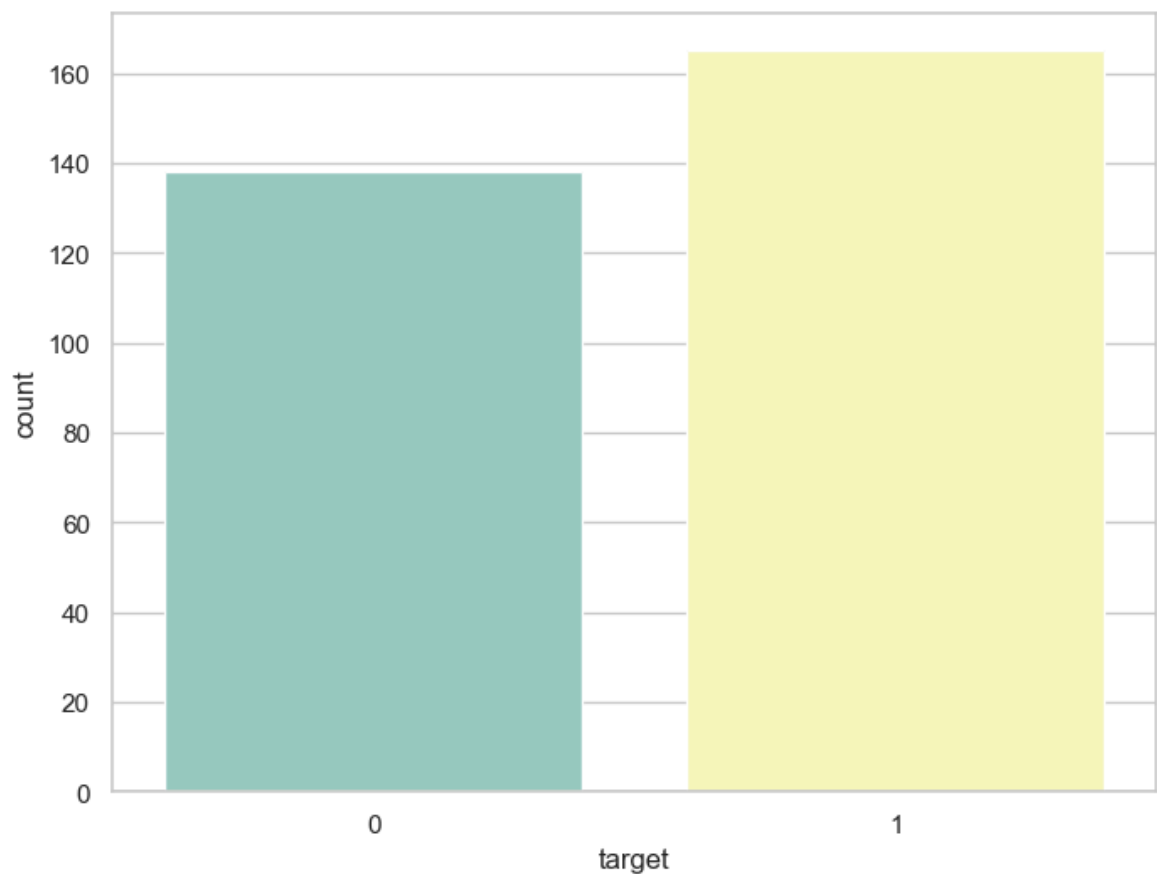
```
In [47]: ax = sns.catplot(x="target", hue='target', col="sex", data=df, kind="count", heig
```



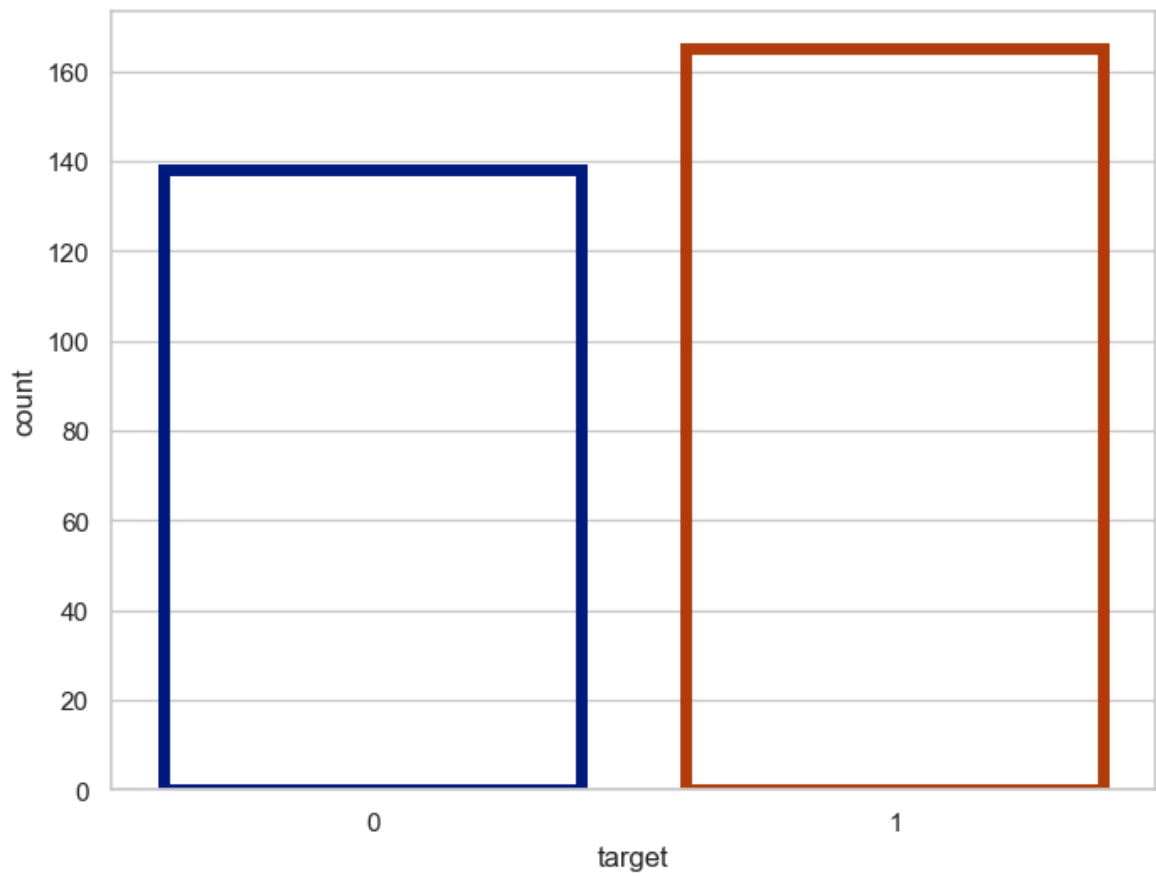
```
In [49]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(y="target", hue="sex", data=df)
plt.show()
```



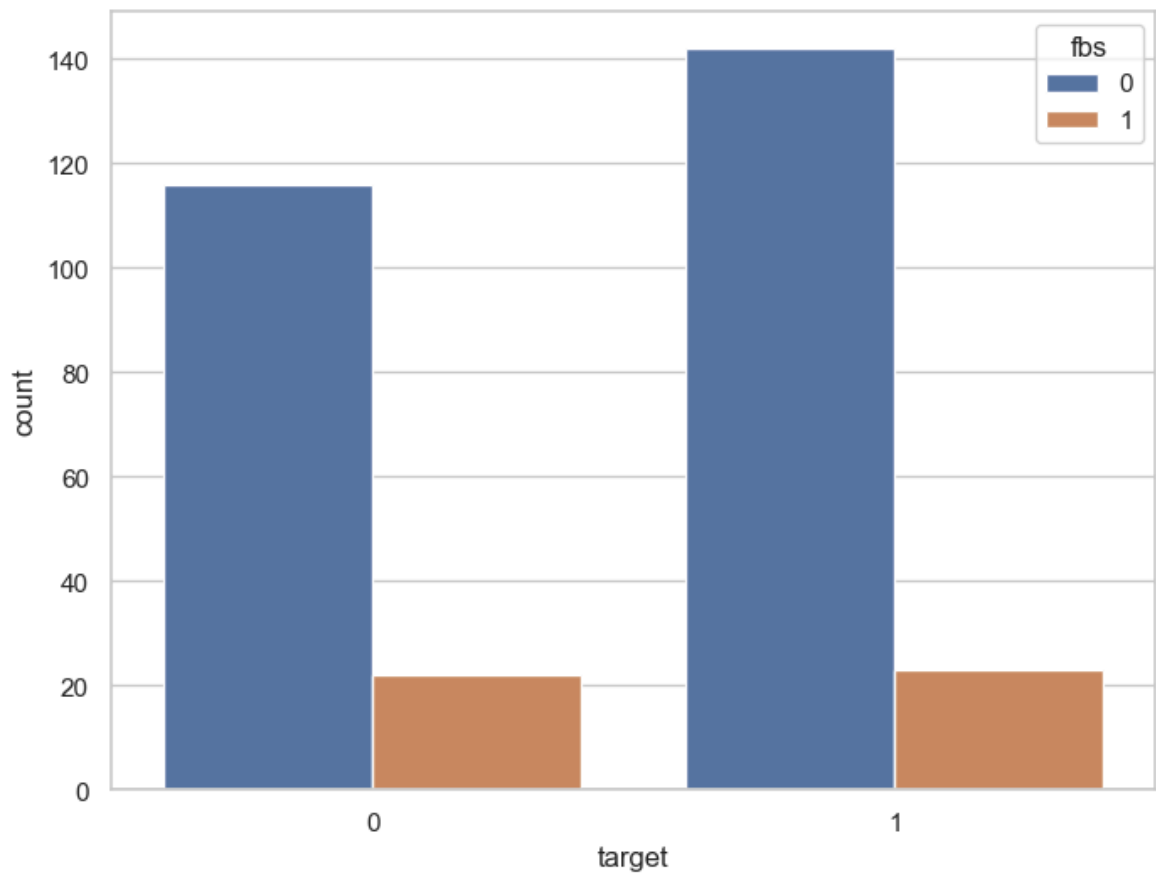
```
In [51]: f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.countplot(x="target", data=df, palette="Set3")  
plt.show()
```



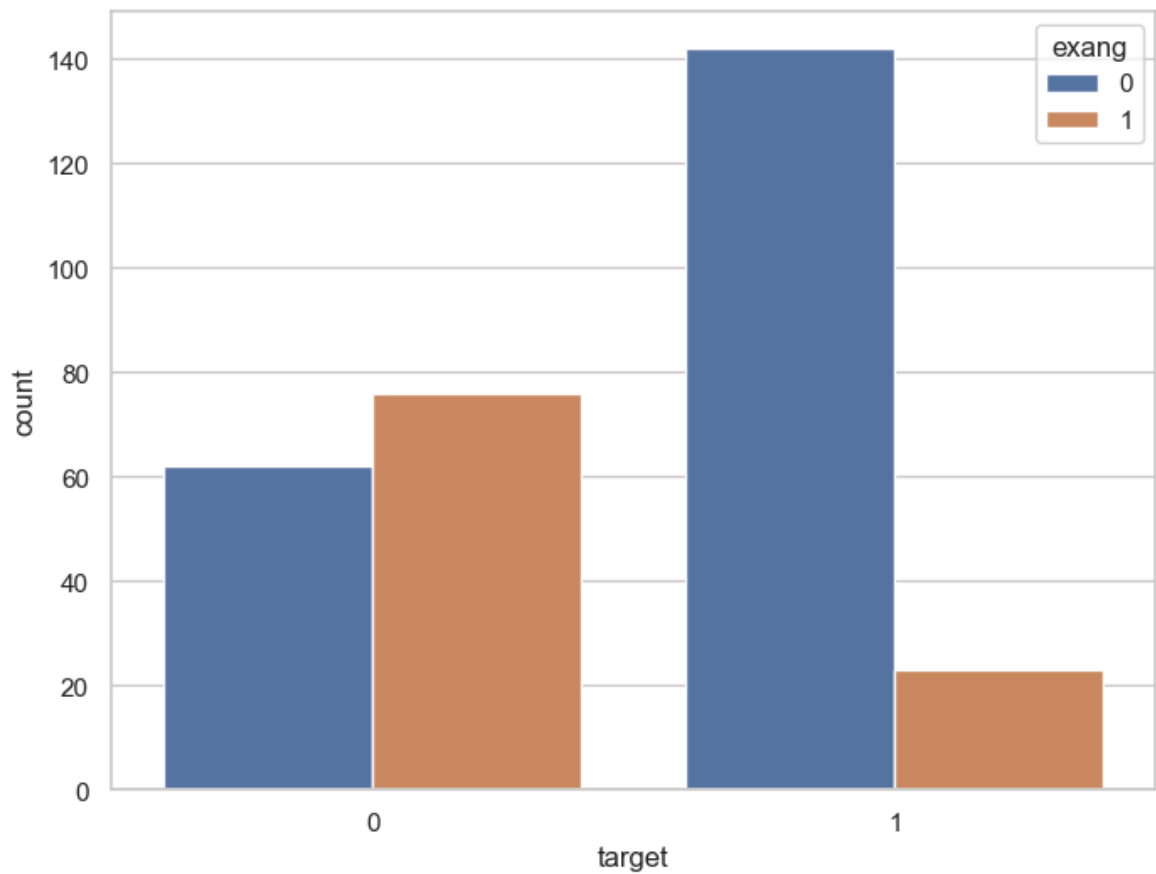
```
In [53]: f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.countplot(x="target", data=df, facecolor=(0, 0, 0, 0), linewidth=5, edge  
plt.show()
```



```
In [55]: f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.countplot(x="target", hue="fbs", data=df)  
plt.show()
```



```
In [57]: f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.countplot(x="target", hue="exang", data=df)  
plt.show()
```



## Bivariate Analysis



```
In [61]: correlation = df.corr()
```

```
In [63]: correlation['target'].sort_values(ascending=False)
```

```
Out[63]: 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 [66]: df['cp'].nunique()
```

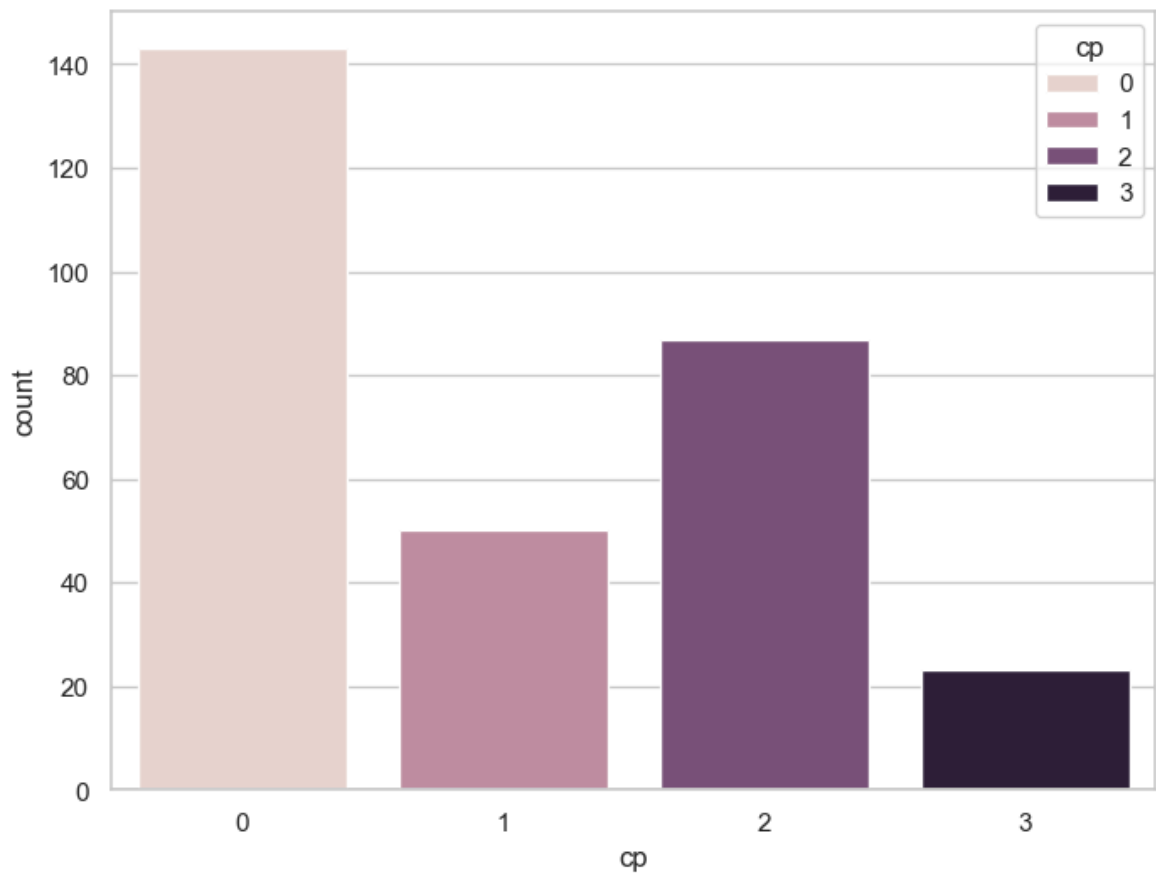
```
Out[66]: 4
```

```
In [68]: df['cp'].value_counts()
```

```
Out[68]: cp  
0      143  
2       87  
1       50  
3       23  
Name: count, dtype: int64
```

## Visualize the frequency distribution of cp variable

```
In [73]: f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.countplot(x="cp", hue='cp', data=df)  
plt.show()
```

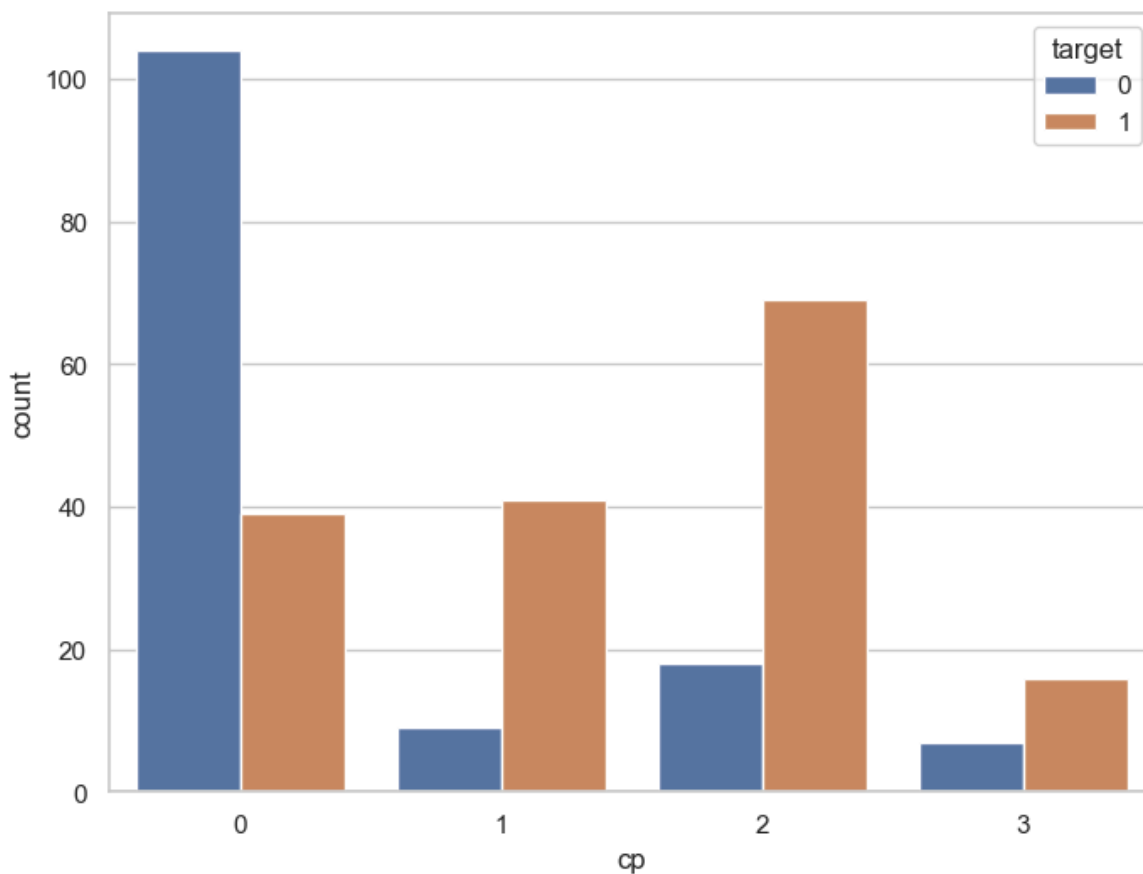


Frequency distribution of target variable wrt cp

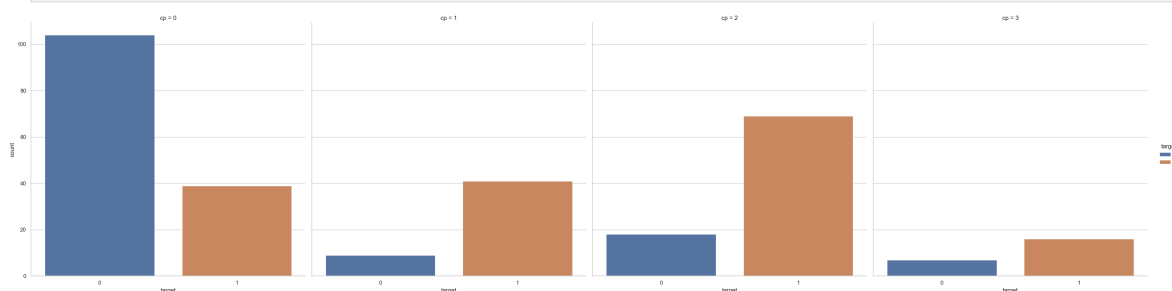
```
In [76]: df.groupby('cp')['target'].value_counts()
```

```
Out[76]: cp target
0      0      104
      1       39
1      1       41
      0        9
2      1       69
      0       18
3      1       16
      0        7
Name: count, dtype: int64
```

```
In [78]: f, ax = plt.subplots(figsize=(8, 6))
ax = sns.countplot(x="cp", hue="target", data=df)
plt.show()
```



In [84]: `ax = sns.catplot(x="target", hue='target', col="cp", data=df, kind="count", height`

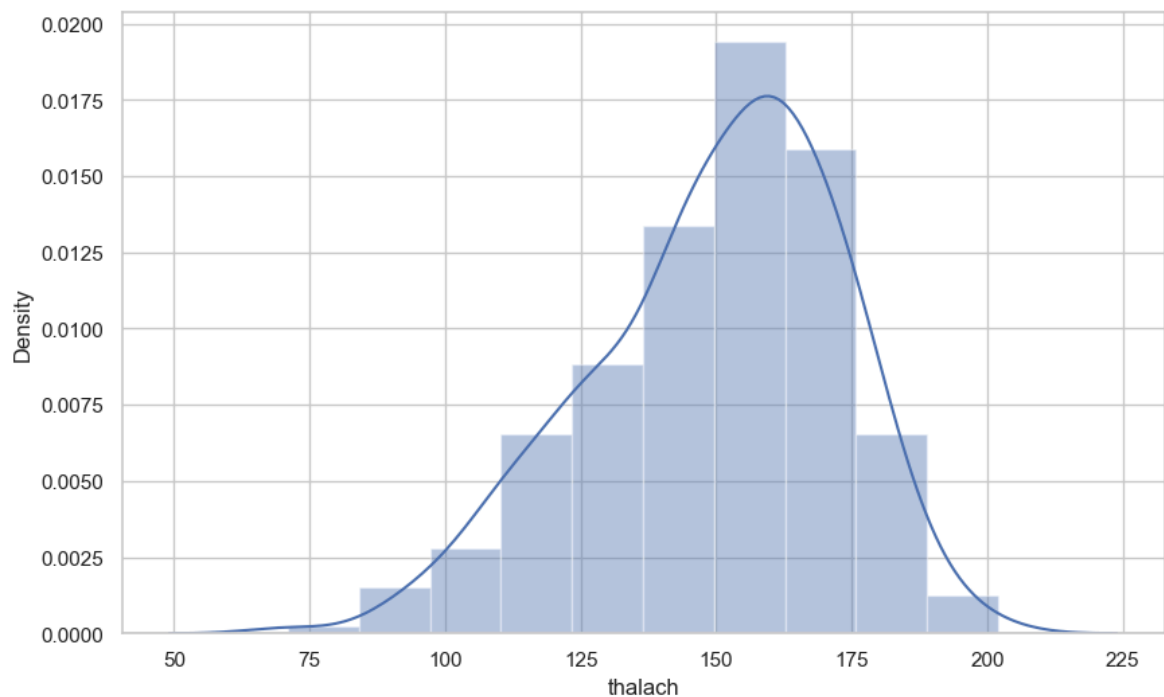


## Analysis of target and thalach variable

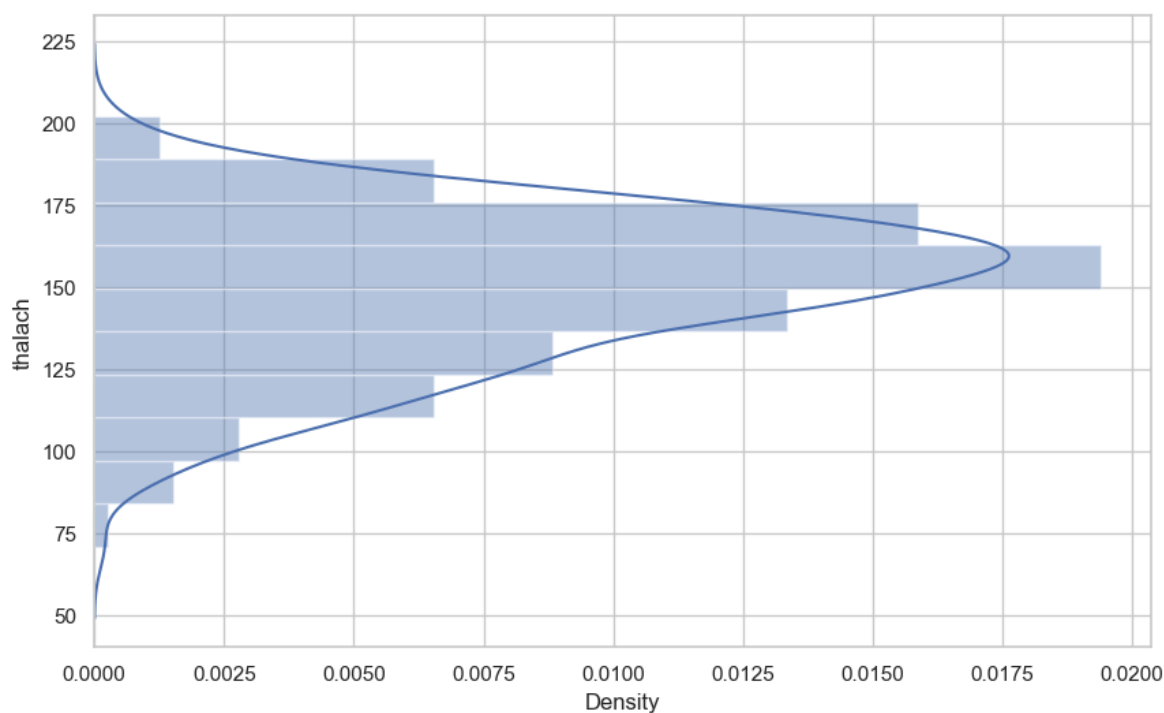
In [87]: `df['thalach'].nunique()`

Out[87]: 91

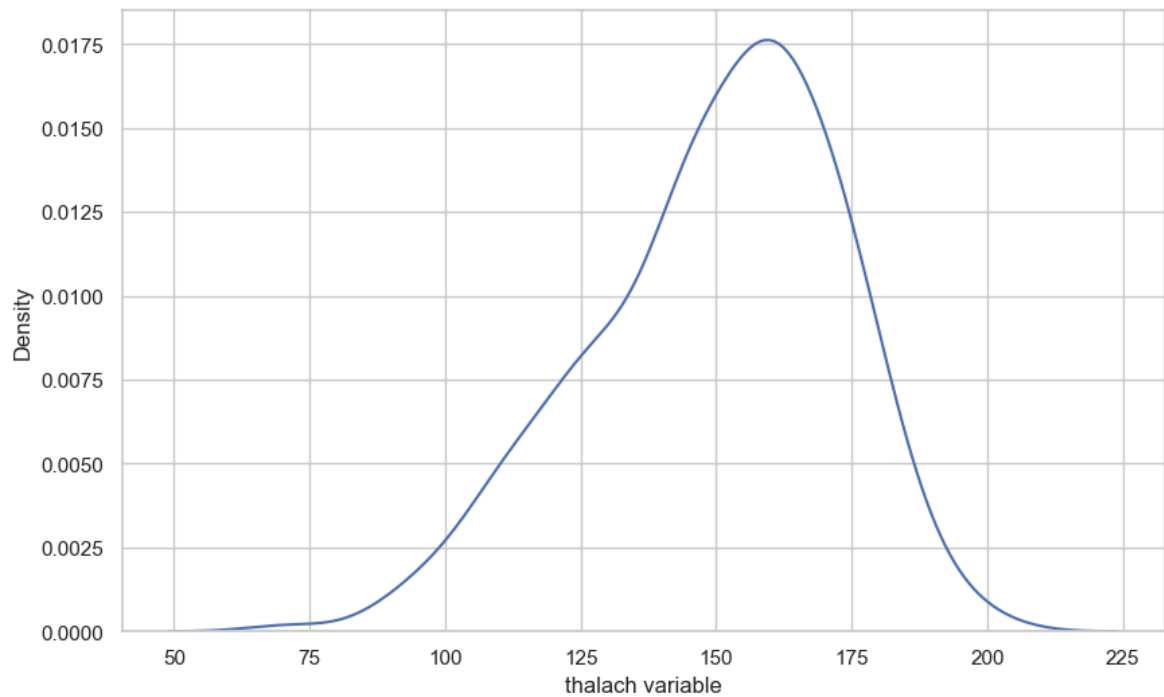
In [89]: `f, ax = plt.subplots(figsize=(10,6))  
x = df['thalach']  
ax = sns.distplot(x, bins=10)  
plt.show()`



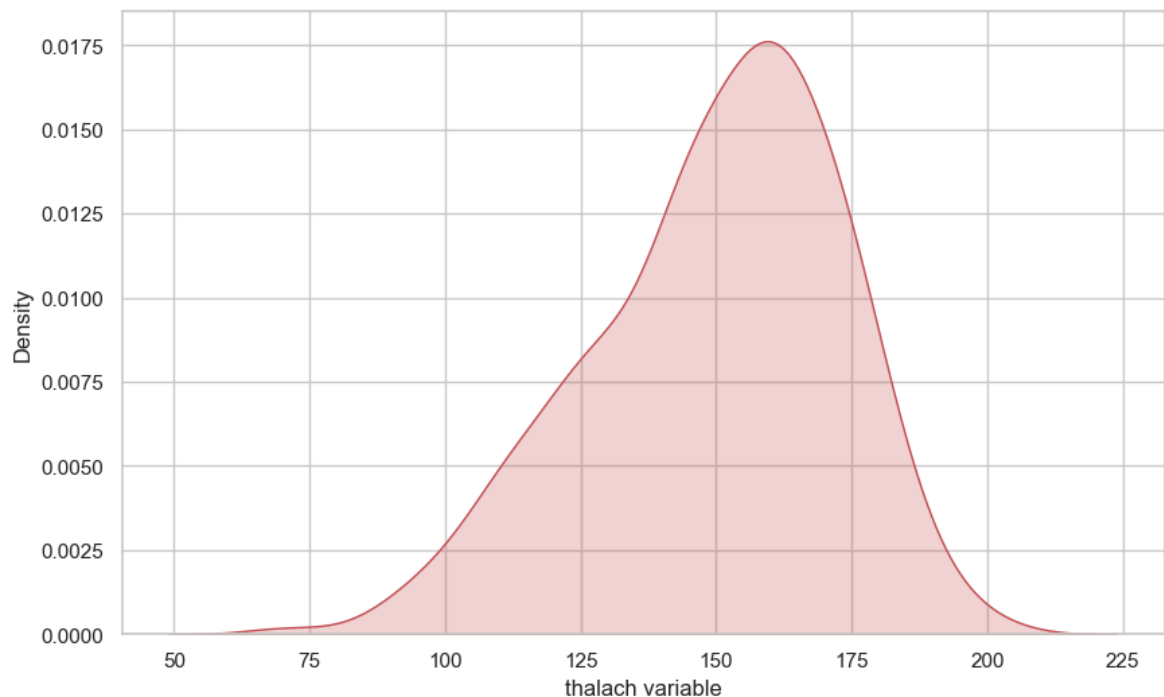
```
In [91]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
ax = sns.distplot(x, bins=10, vertical=False)
plt.show()
```



```
In [93]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.kdeplot(x)
plt.show()
```

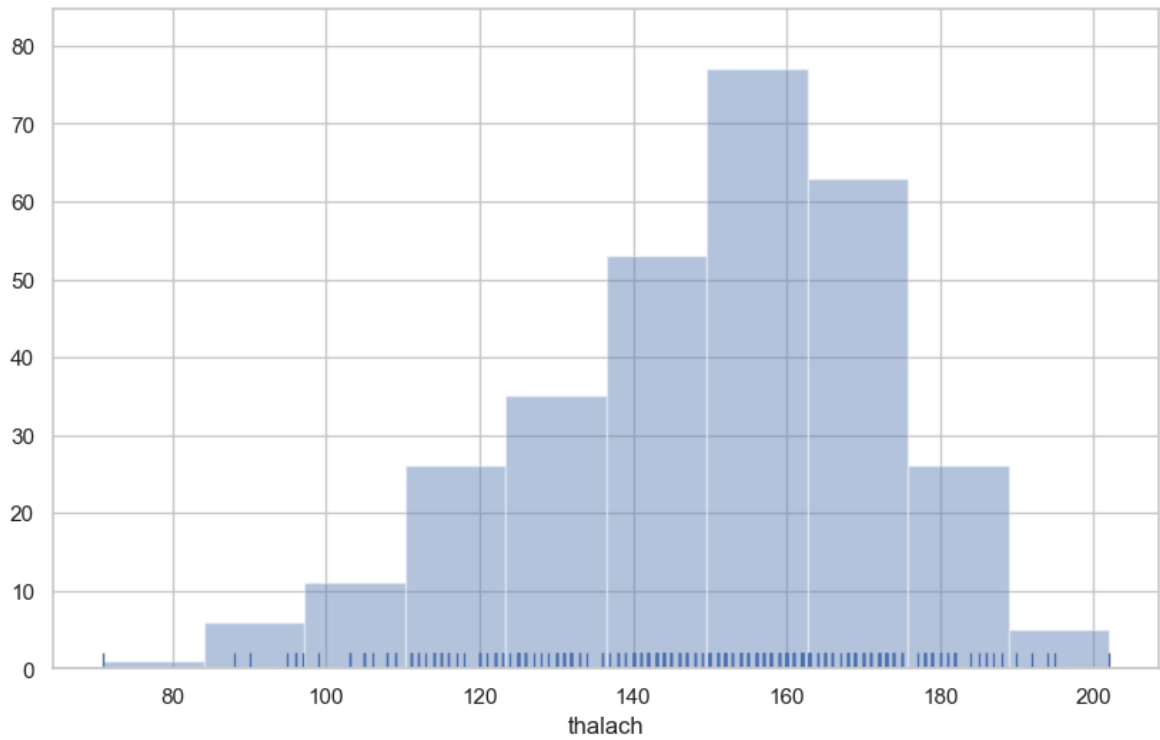


```
In [95]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.kdeplot(x, shade=True, color='r')
plt.show()
```

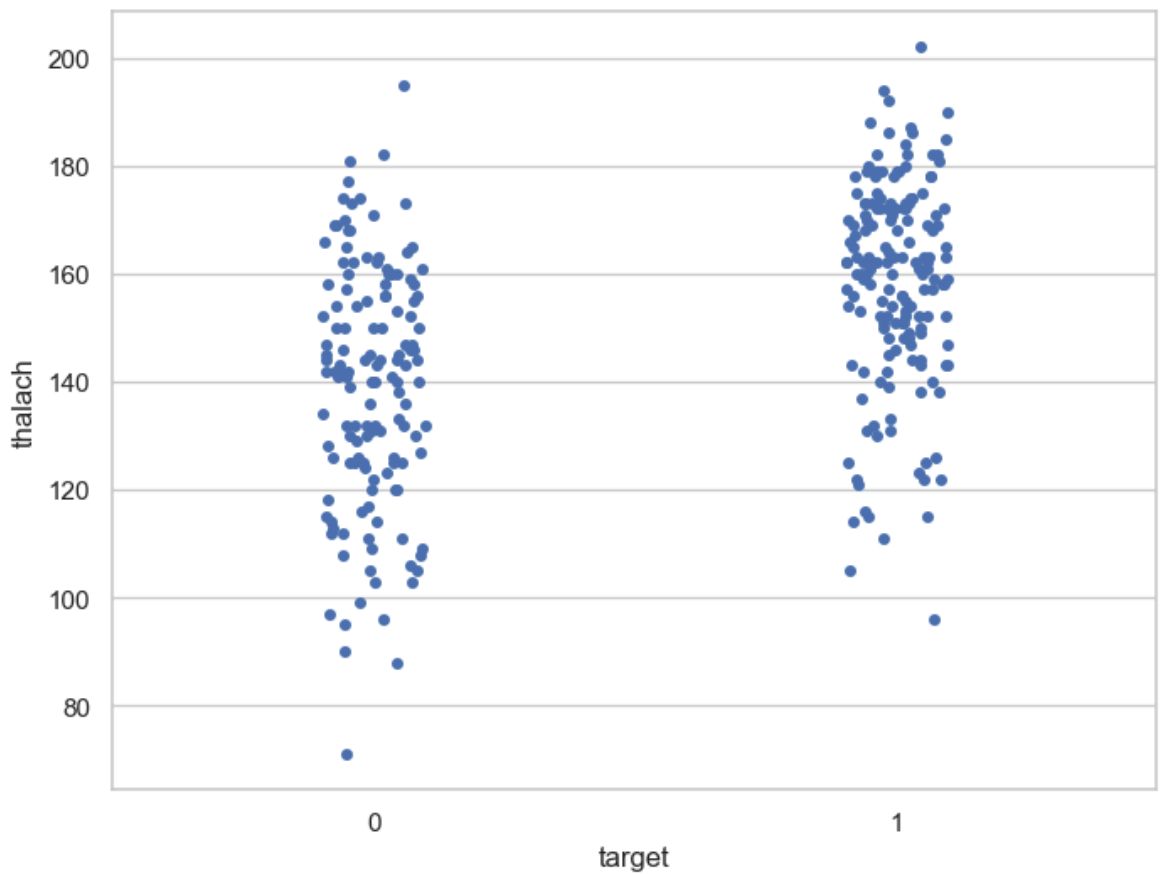


## Histogram

```
In [98]: f, ax = plt.subplots(figsize=(10,6))
x = df['thalach']
ax = sns.distplot(x, kde=False, rug=True, bins=10)
plt.show()
```

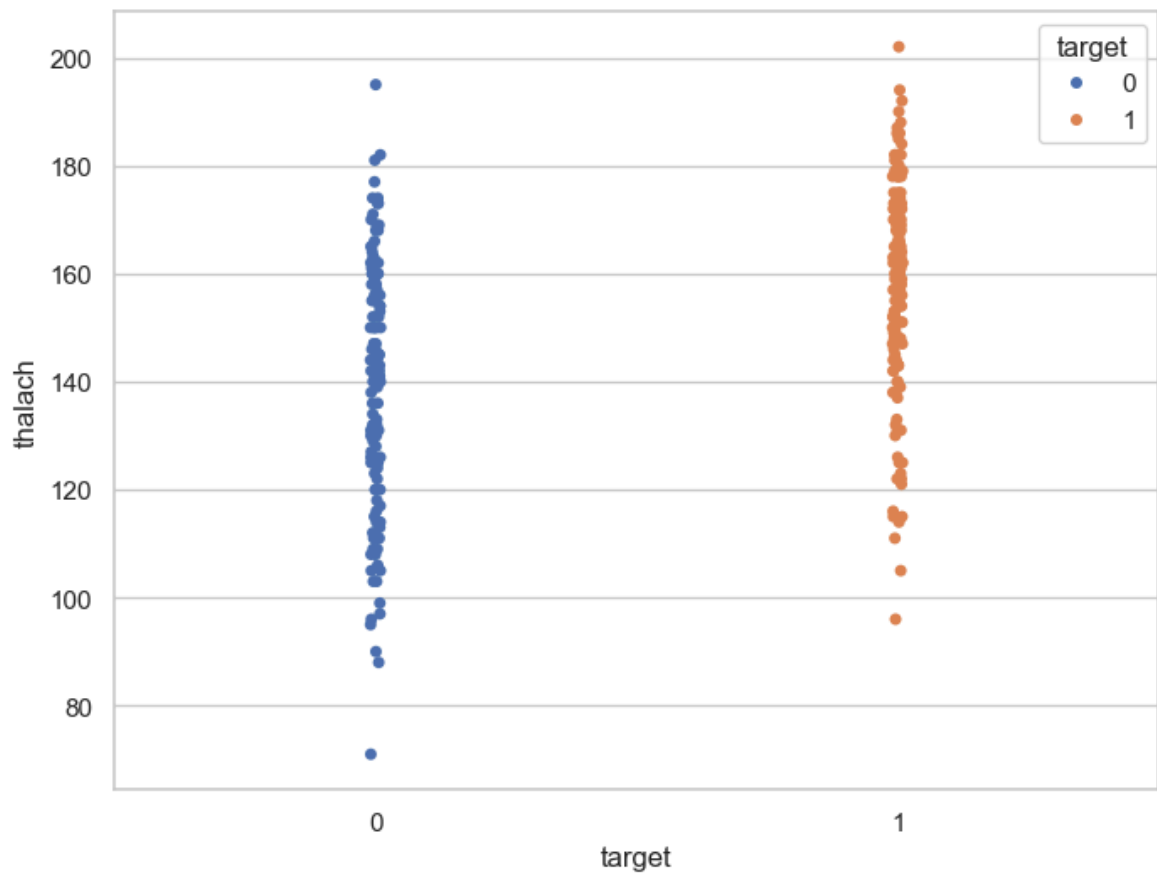


```
In [100... f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", y="thalach", data=df)
plt.show()
```



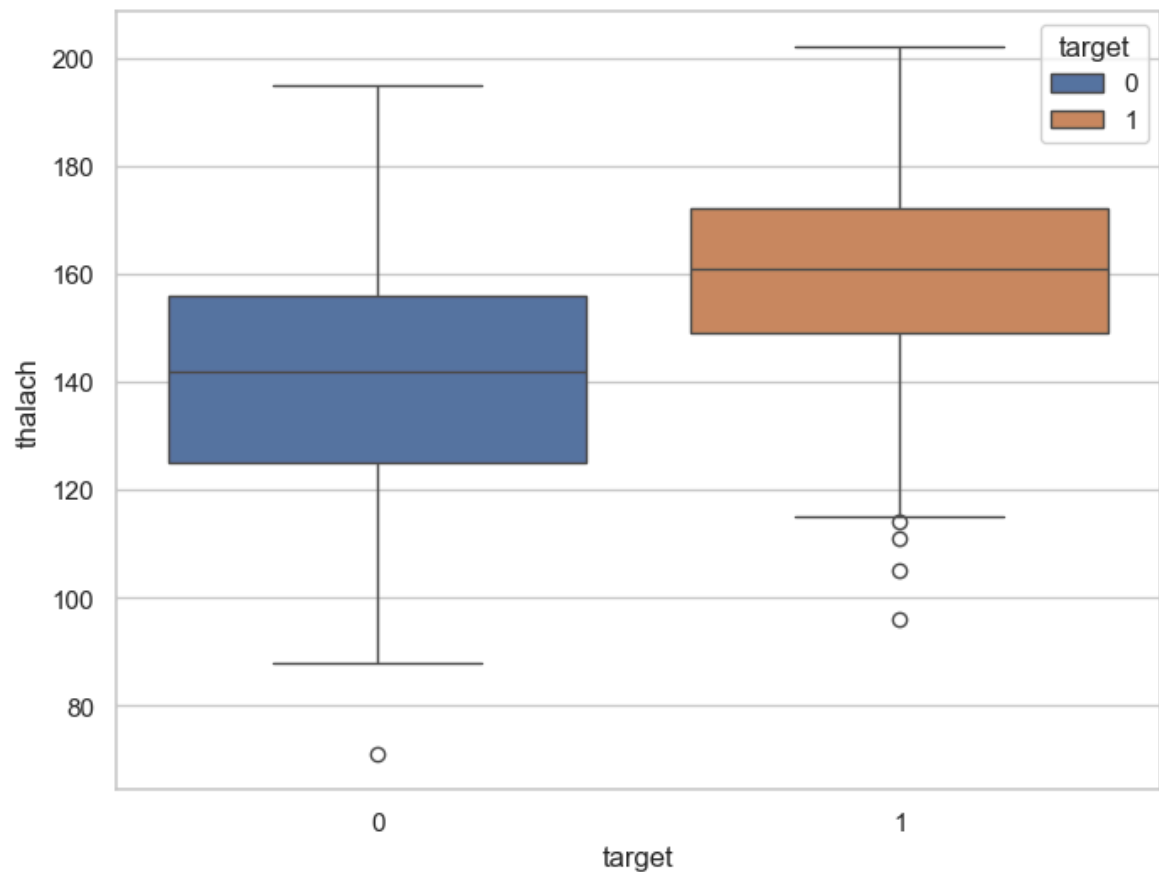
## Interpretation

```
In [105... f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", hue='target', y="thalach", data=df, jitter = 0.01)
plt.show()
```



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

```
In [110... f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x="target", hue='target', y="thalach", data=df)
plt.show()
```

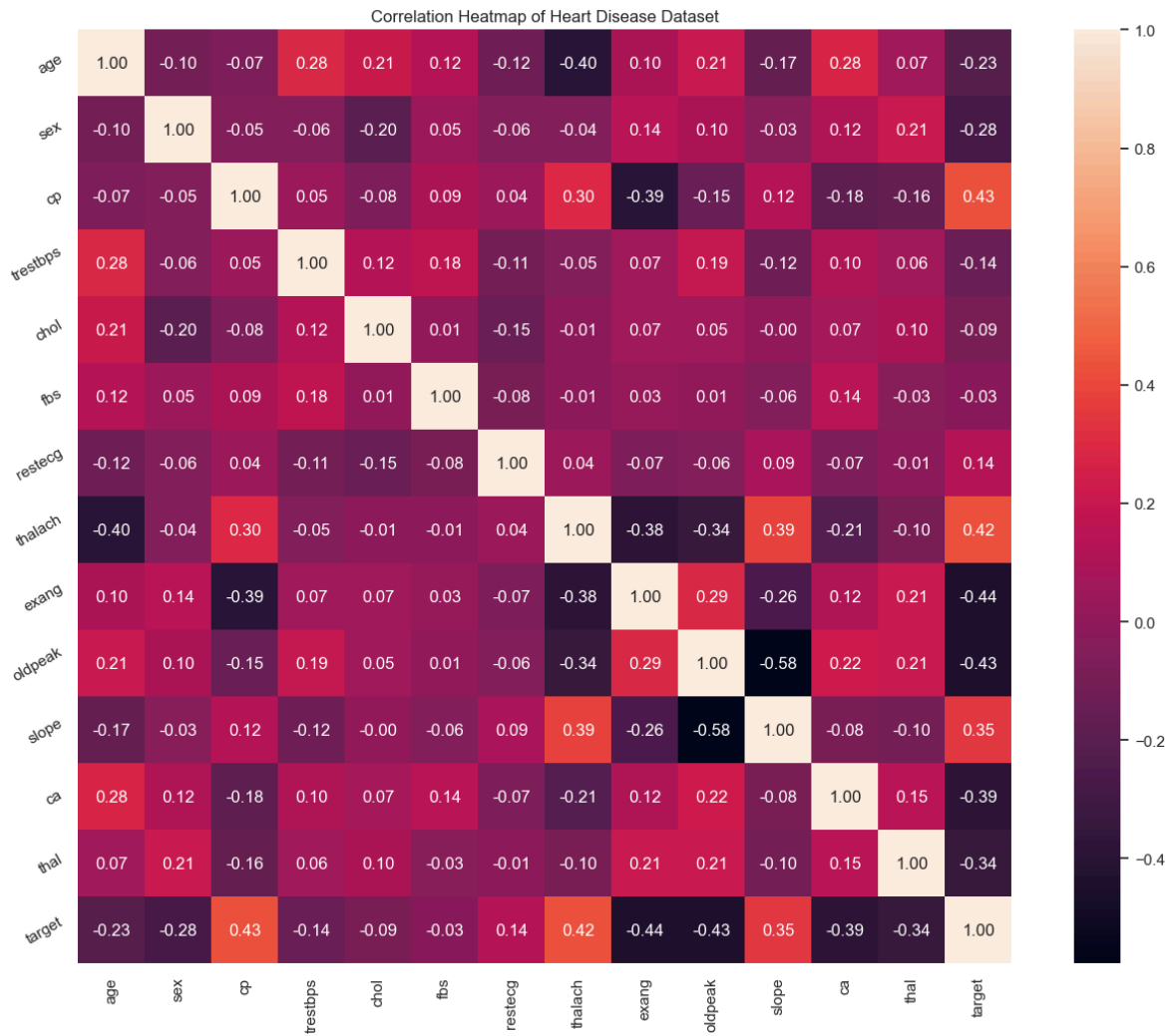


## Multivariate analysis

### Heat Map

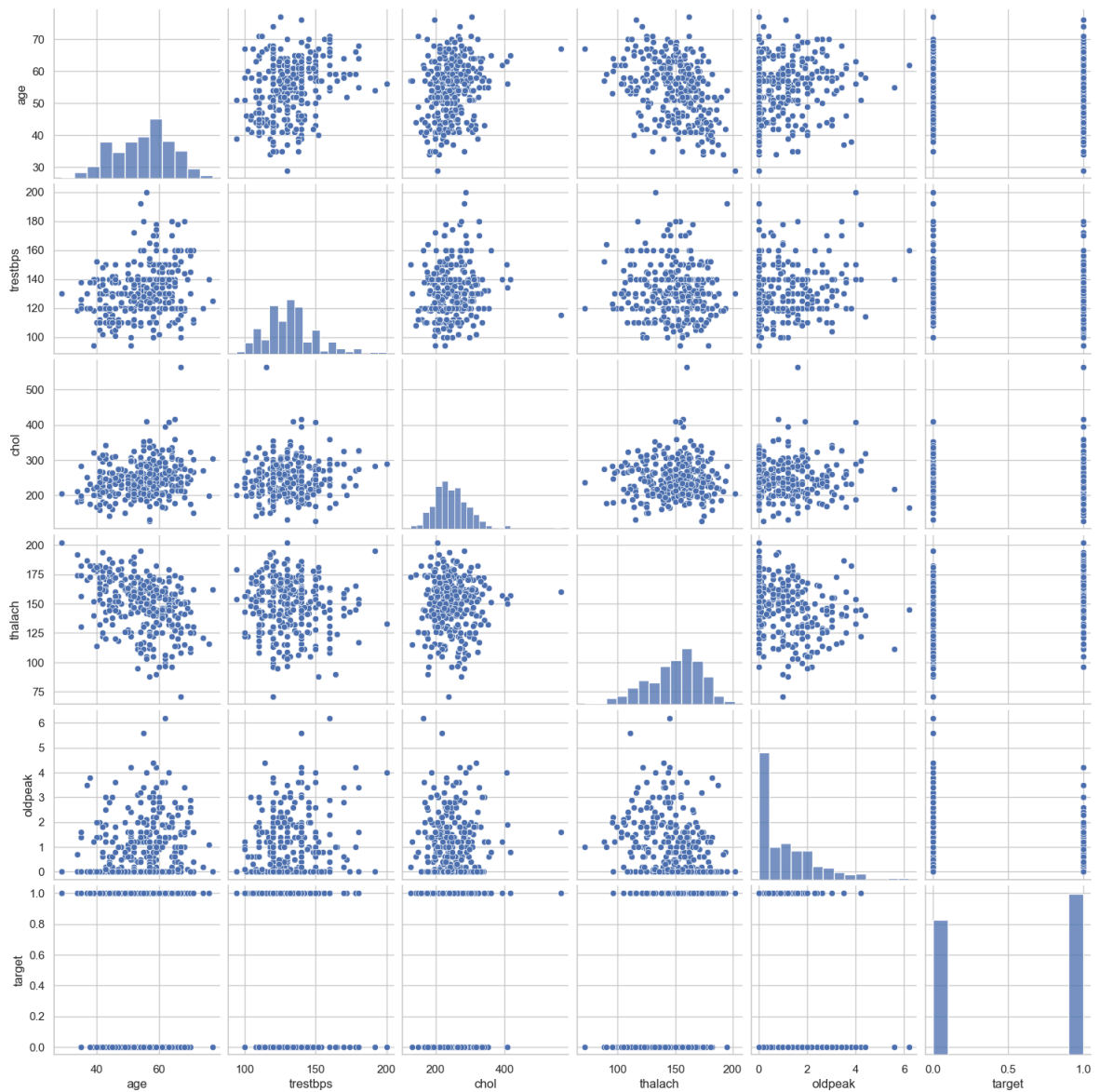
```
In [115... plt.figure(figsize=(16,12))
plt.title('Correlation Heatmap of Heart Disease Dataset')
a = sns.heatmap(correlation, square=True, annot=True, fmt='.2f', linecolor='white')
a.set_xticklabels(a.get_xticklabels(), rotation=90)
a.set_yticklabels(a.get_yticklabels(), rotation=30)
plt.show()
```





## Pair Plot

```
In [118... 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 [122...] df['age'].nunique()
```

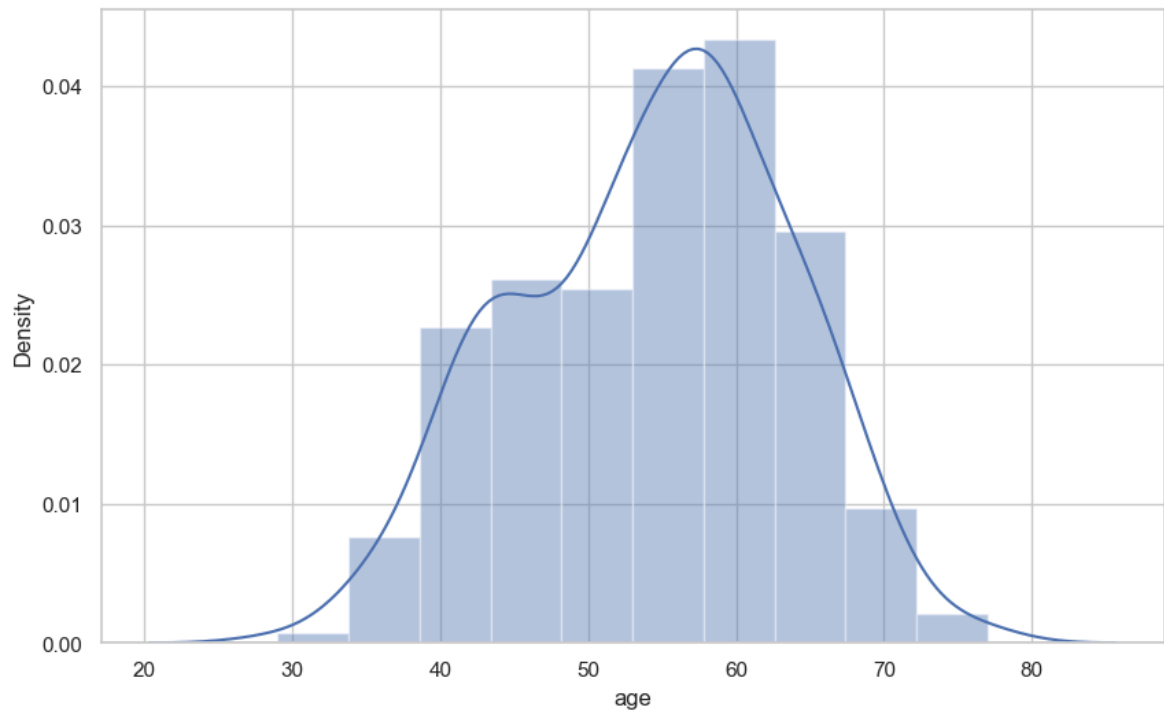
```
Out[122...] 41
```

```
In [124...] df['age'].describe()
```

```
Out[124...] 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

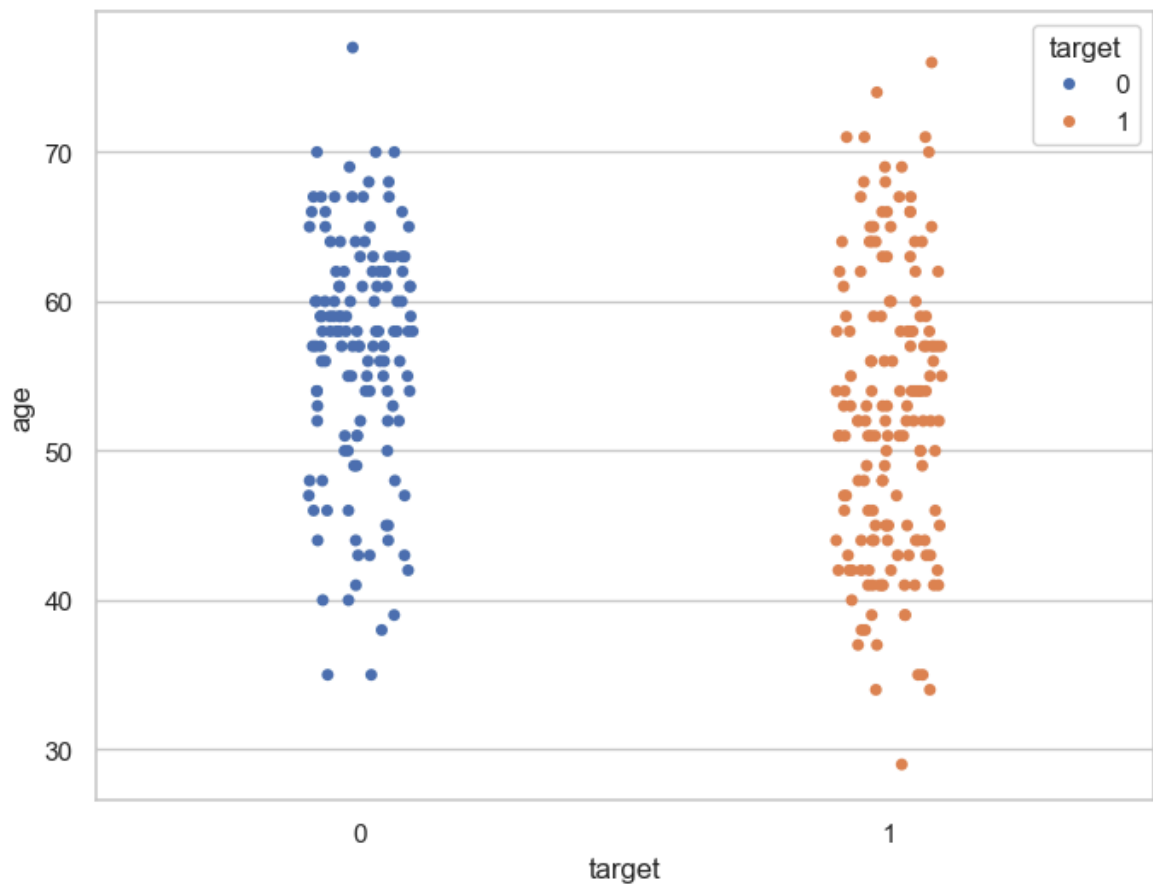
```
In [127... f, ax = plt.subplots(figsize=(10,6))
x = df['age']
ax = sns.distplot(x, bins=10)
plt.show()
```



## Analyze age and target variable

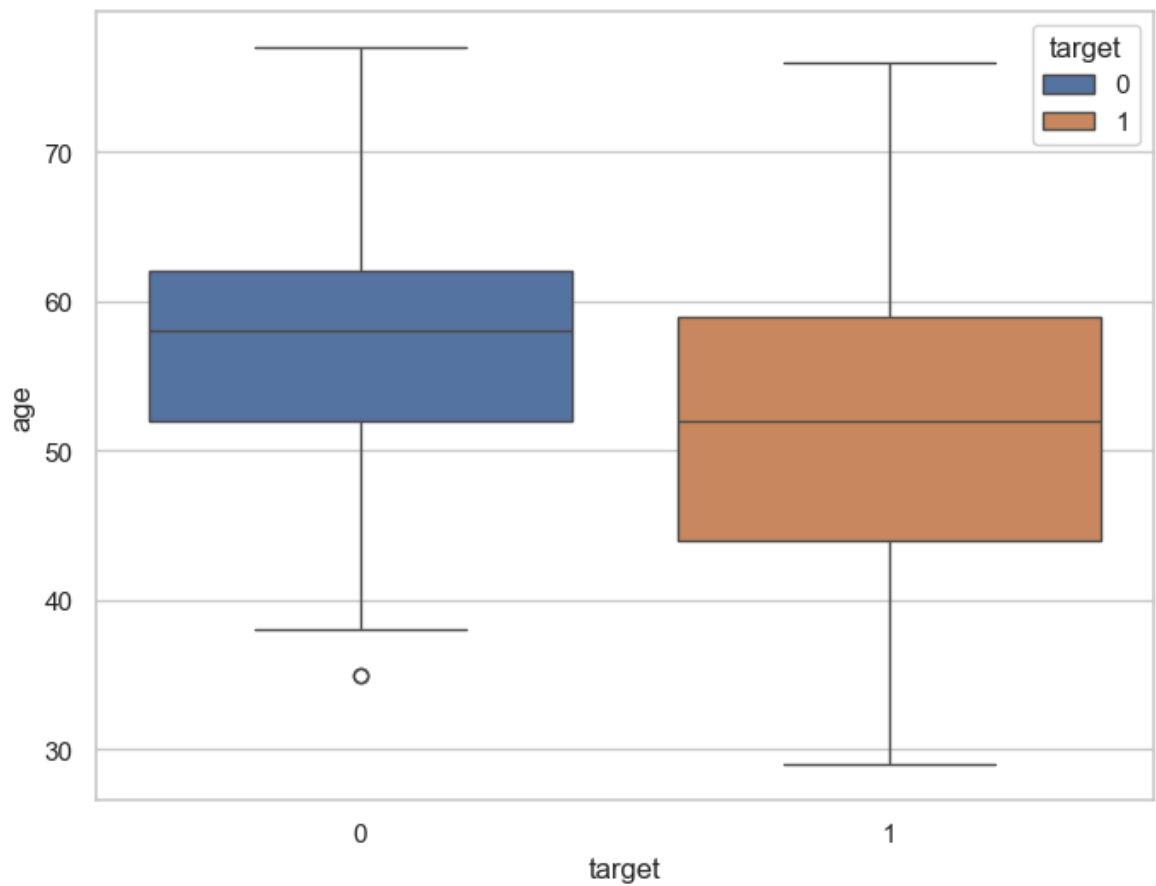
### Visualize frequency distribution of age variable wrt target

```
In [133... f, ax = plt.subplots(figsize=(8, 6))
sns.stripplot(x="target", hue='target', y="age", data=df)
plt.show()
```



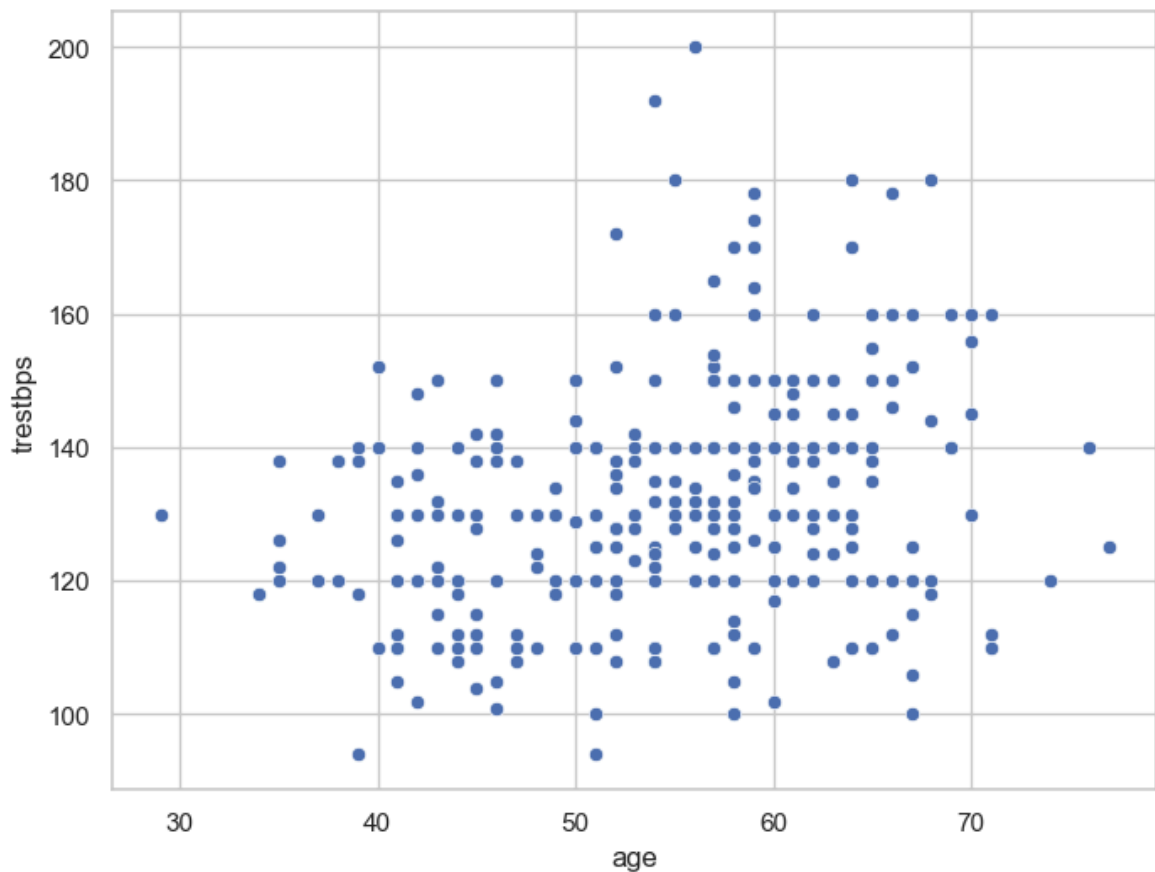
Visualize distribution of `age` variable wrt `target` with boxplot

```
In [138... f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x="target", hue='target', y="age", data=df)
plt.show()
```



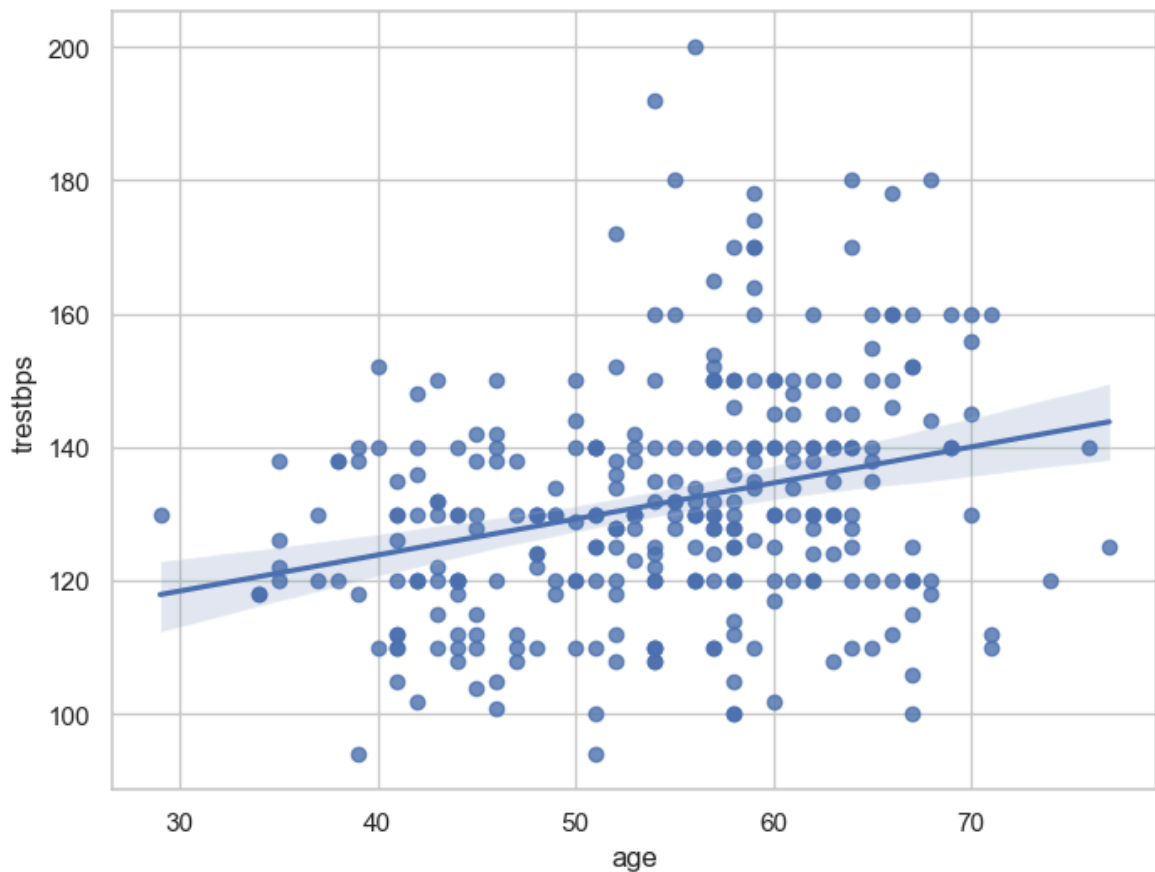
## Analyze age and trestbps variable

```
In [141... f, ax = plt.subplots(figsize=(8, 6))
ax = sns.scatterplot(x="age", y="trestbps", data=df)
plt.show()
```



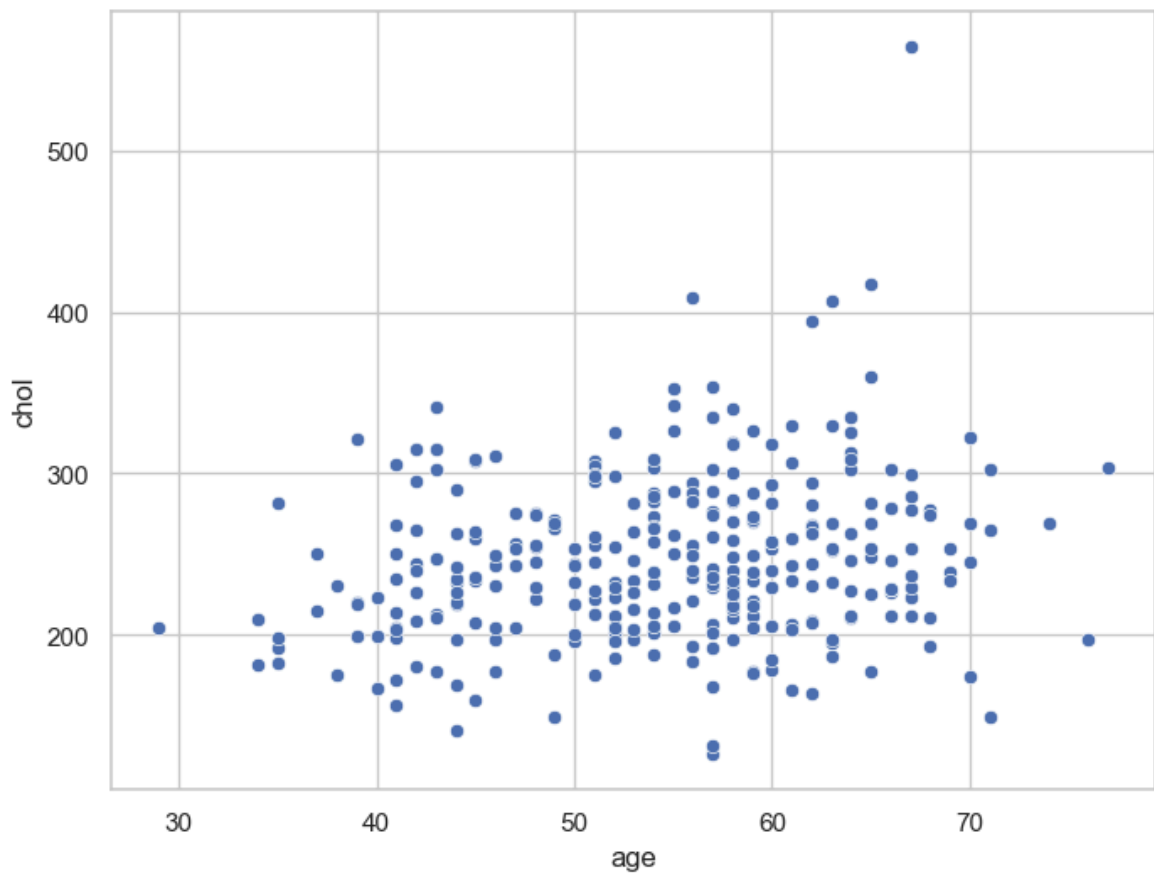
In [143...

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

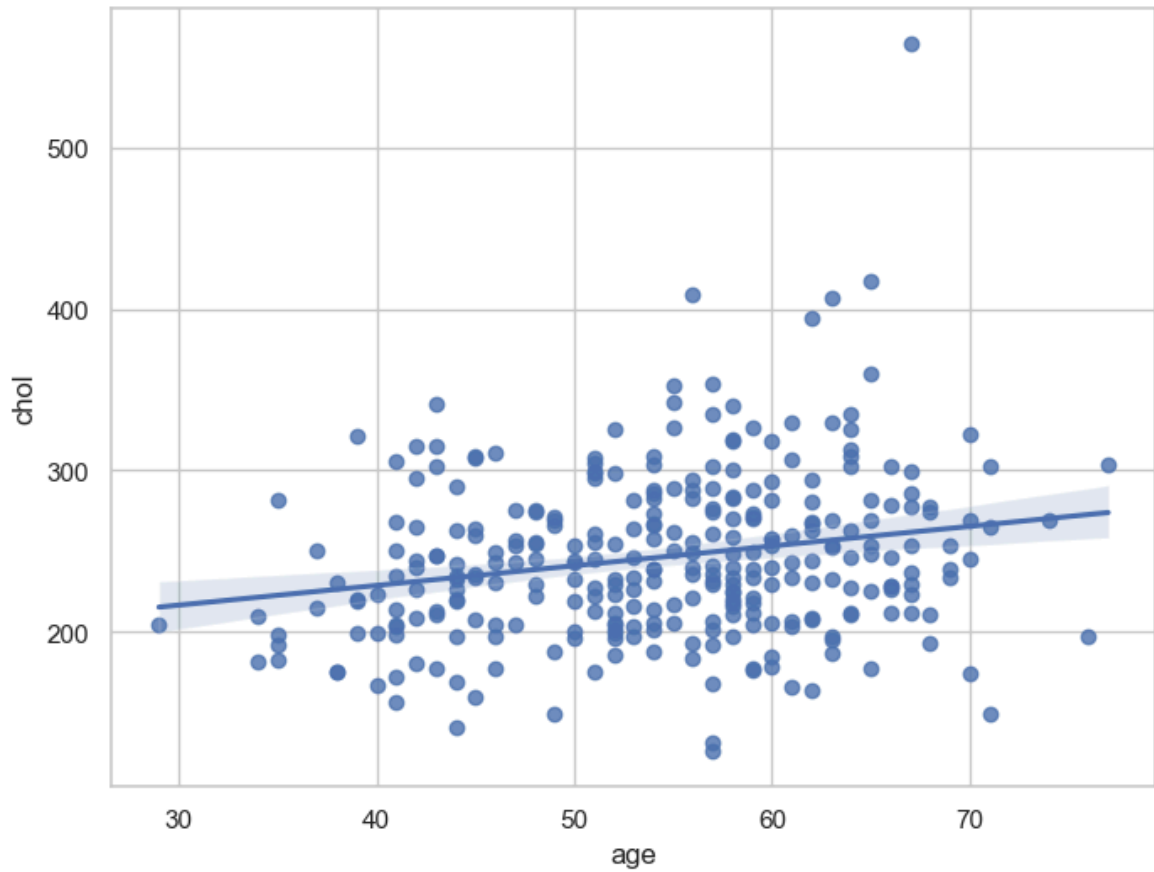


Analyze `age` and `chol` variable

```
In [146... f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.scatterplot(x="age", y="chol", data=df)  
plt.show()
```



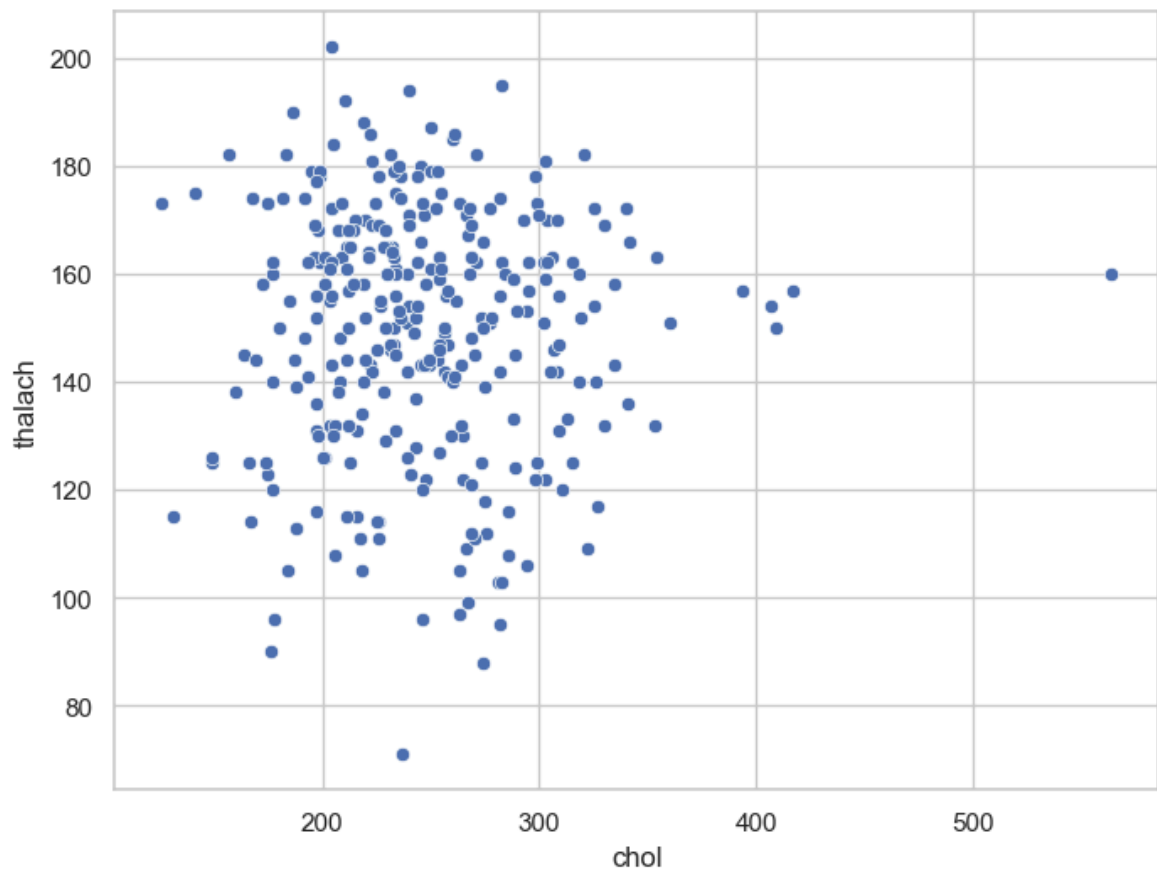
```
In [148... f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.regplot(x="age", y="chol", data=df)  
plt.show()
```



## Analyze chol and thalach variable

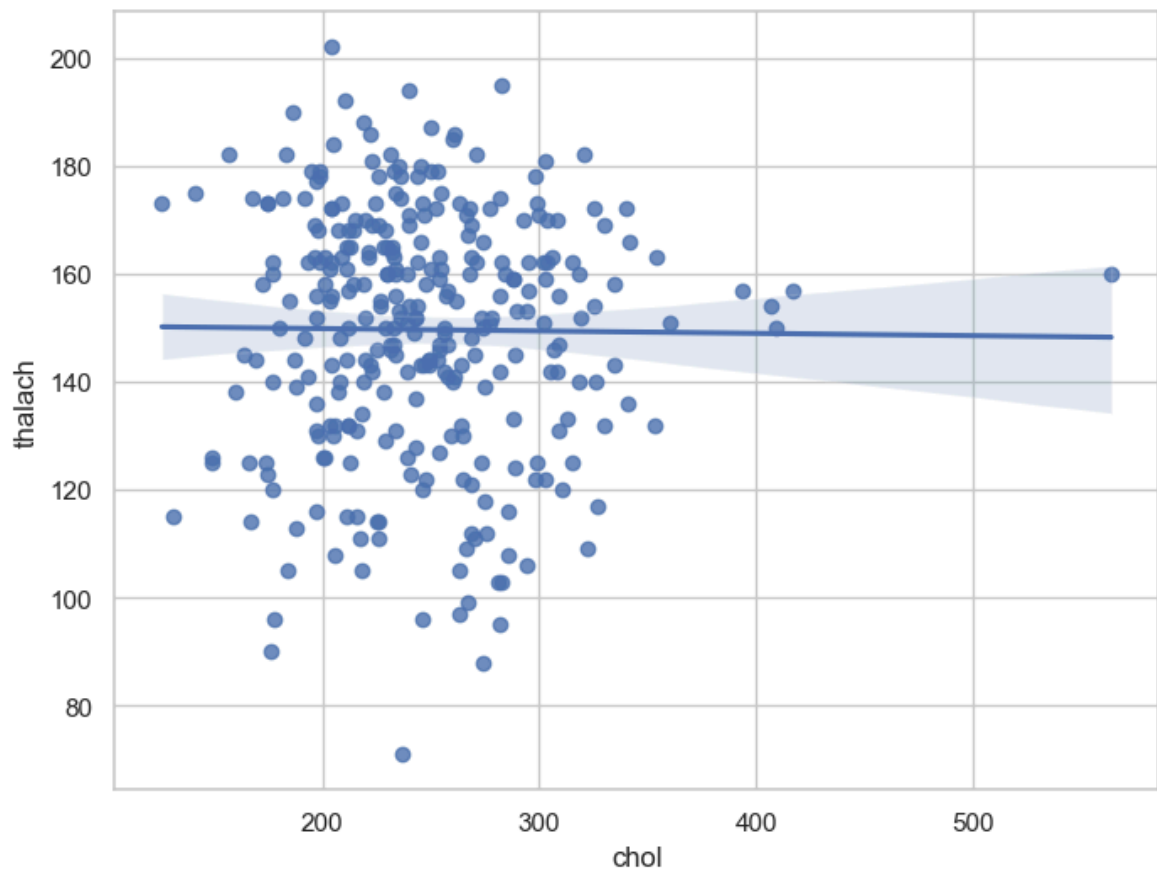
```
In [151... f, ax = plt.subplots(figsize=(8, 6))  
ax = sns.scatterplot(x="chol", y = "thalach", data=df)  
plt.show()
```





In [153...

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



## Dealing with missing values

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

## Check with ASSERT statement

```
In [161... assert pd.notnull(df).all().all()
```

```
In [163... assert (df >= 0).all().all()
```

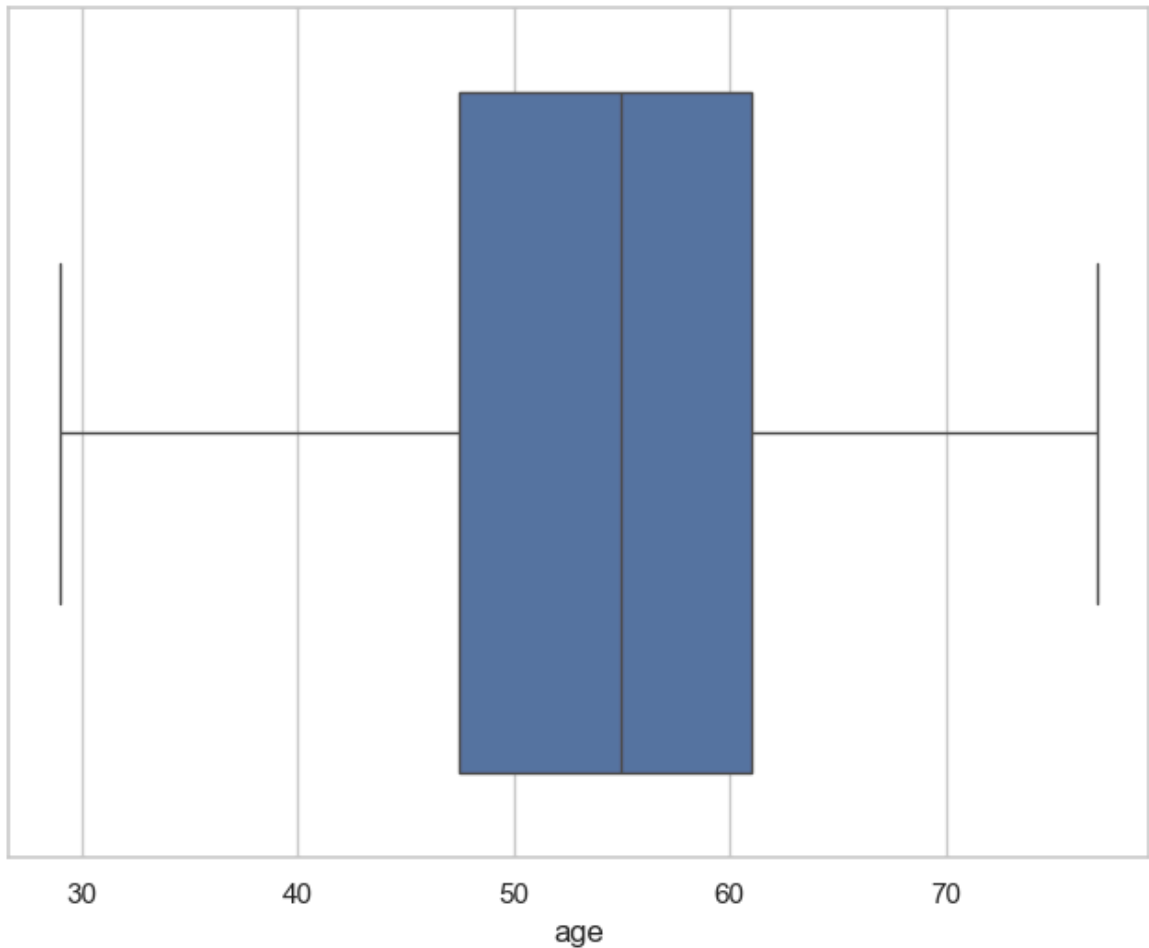
## Outlier detection

### Age Variable

```
In [167... df['age'].describe()
```

```
Out[167... 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
```

```
In [169... f, ax = plt.subplots(figsize=(8, 6))  
sns.boxplot(x=df["age"])  
plt.show()
```

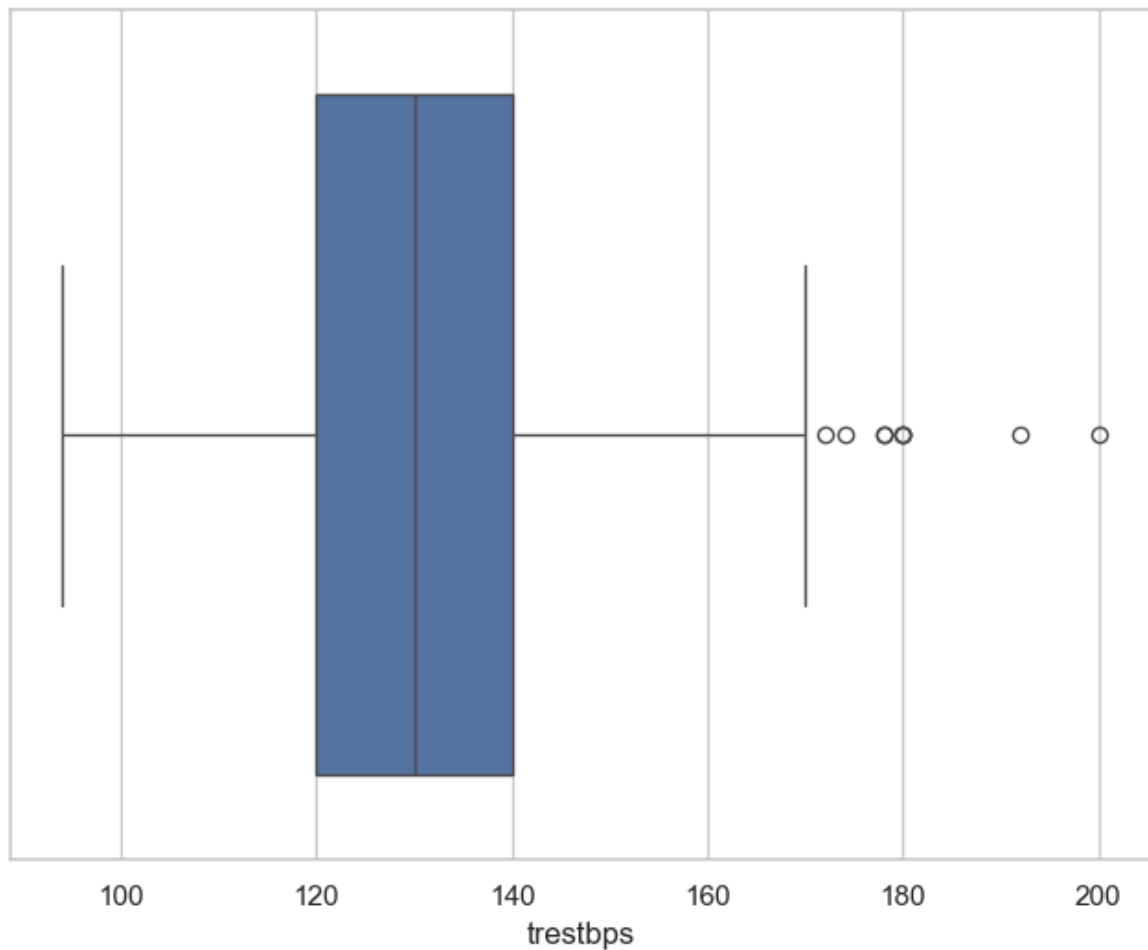


## trestbps variable

```
In [172...] df['trestbps'].describe()
```

```
Out[172...] 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
```

```
In [174...] f, ax = plt.subplots(figsize=(8, 6))  
sns.boxplot(x=df["trestbps"])  
plt.show()
```

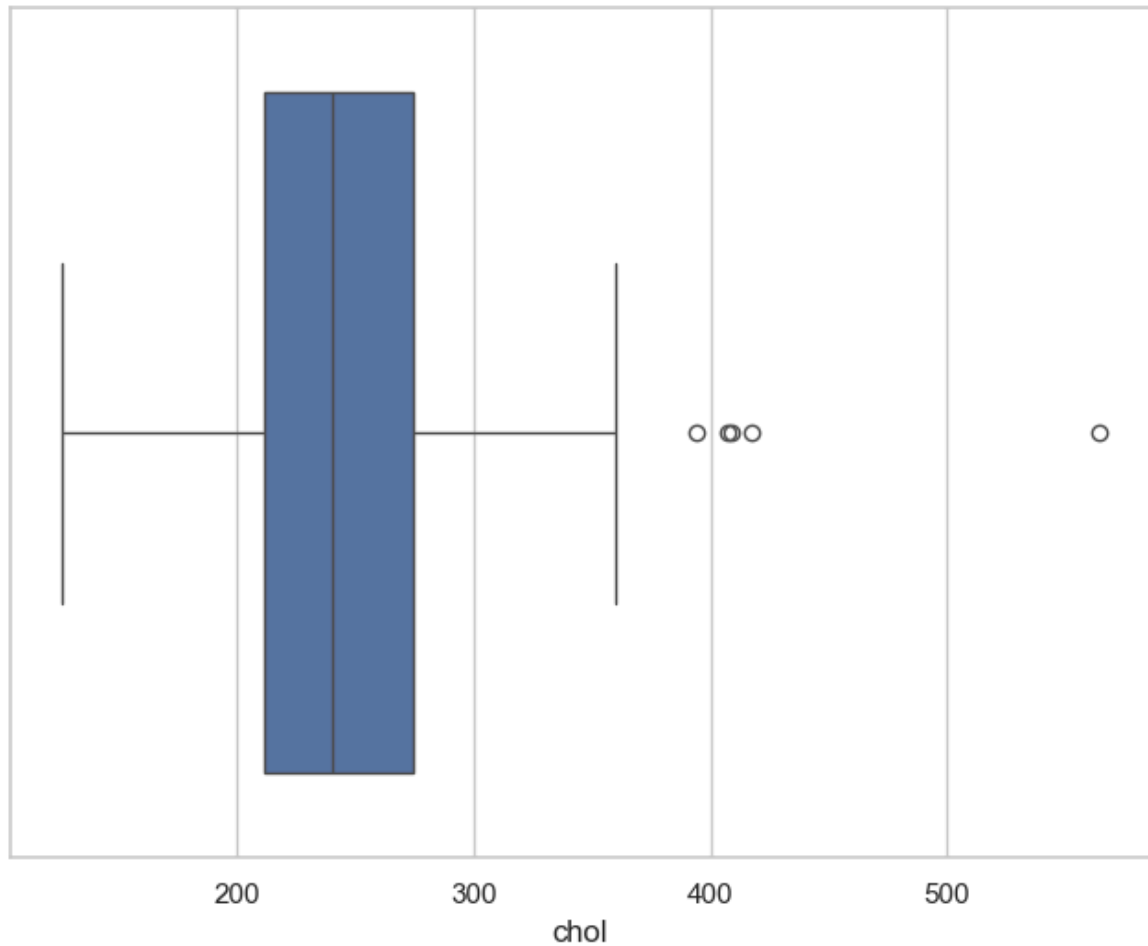


## chol variable

```
In [177...] df['chol'].describe()
```

```
Out[177...] 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
```

```
In [179...] f, ax = plt.subplots(figsize=(8, 6))  
sns.boxplot(x=df["chol"])  
plt.show()
```



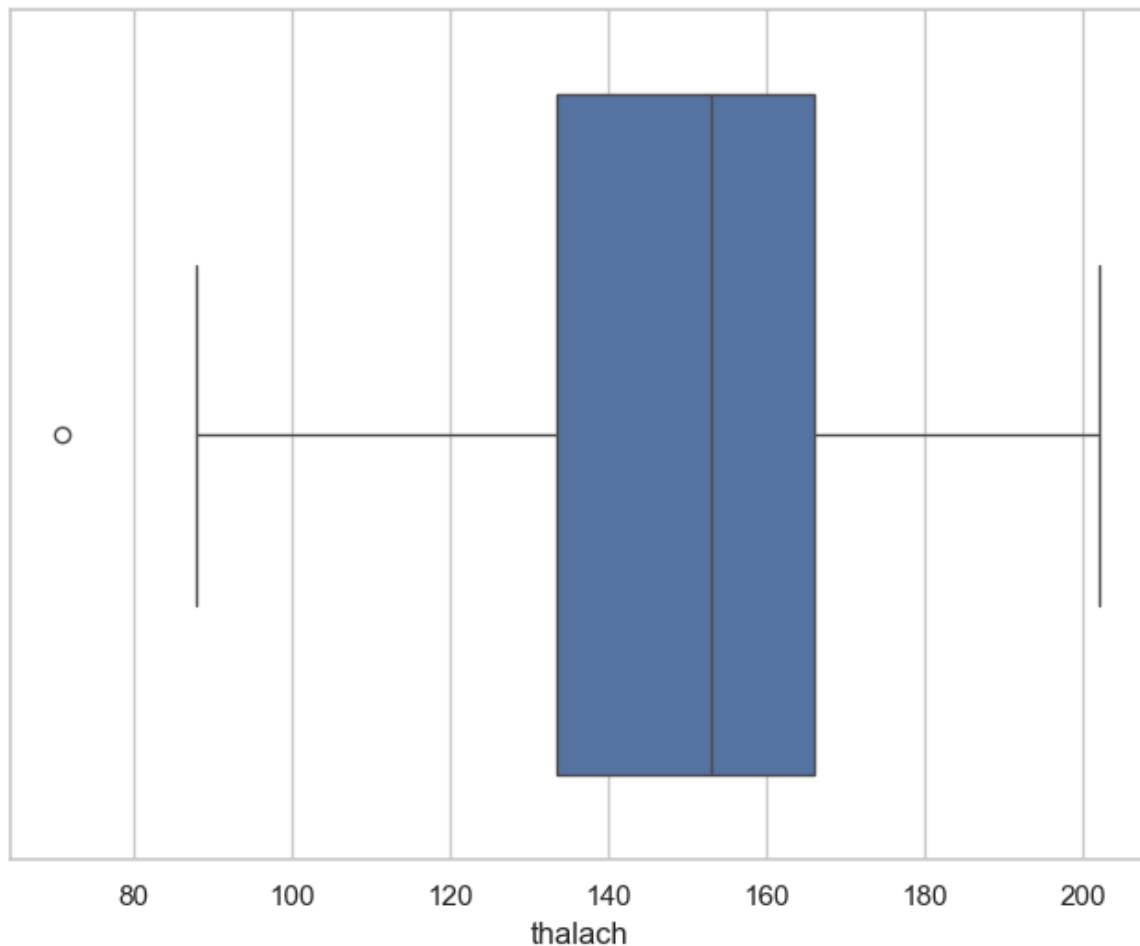
## thalach variable

In [182... `df['thalach'].describe()`

```
Out[182... 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 [185... f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["thalach"])
plt.show()
```



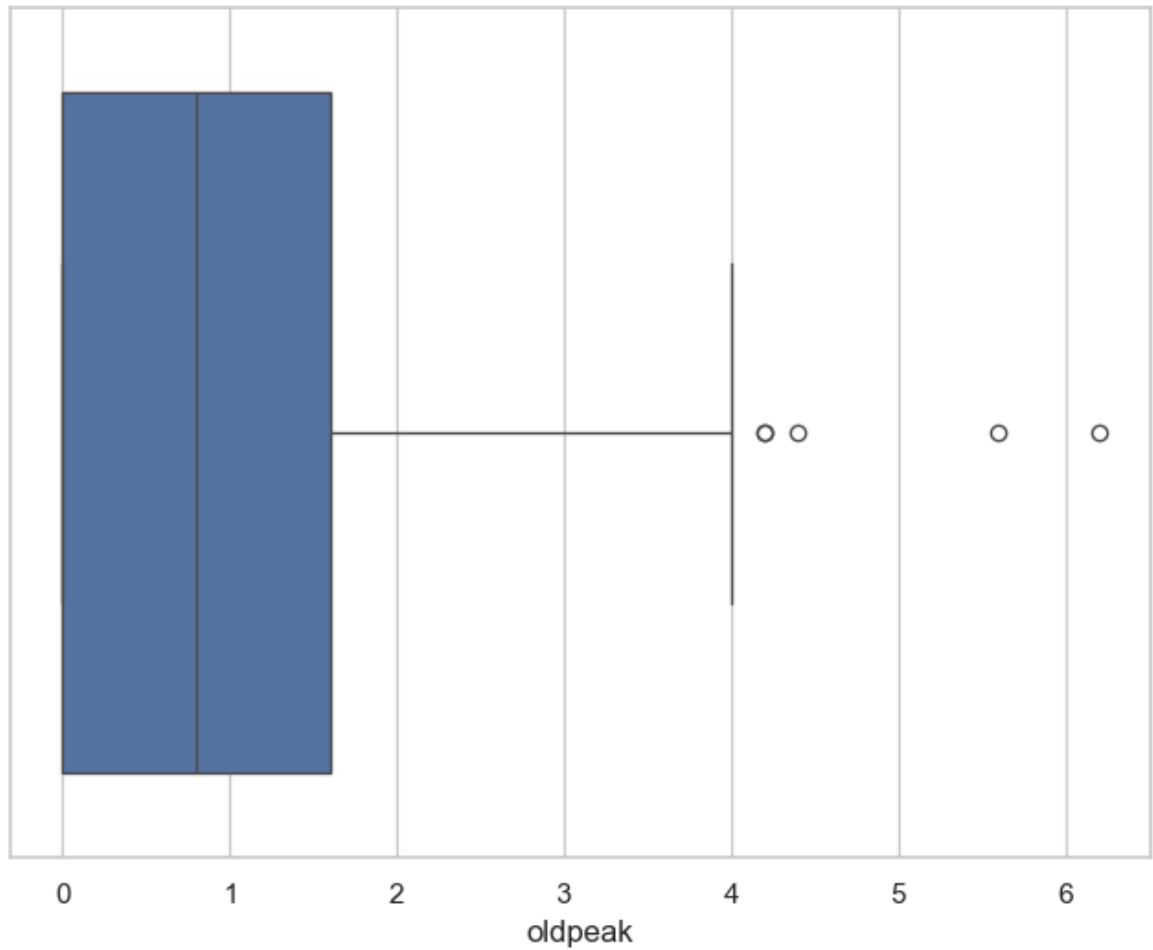
### oldpeak variable

```
In [188...] df['oldpeak'].describe()
```

```
Out[188...] 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 [191...] f, ax = plt.subplots(figsize=(8, 6))
sns.boxplot(x=df["oldpeak"])
plt.show()
```



Completed

In [ ]: