

HOSPITAL MANAGEMENT SYSTEM PROJECT USING EDA

Heart disease analysis (Domain - Healthcare)

MAJOR CAPSTONE PROJECT

Import libraries

```
In [2]: 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 [4]: import warnings
warnings.filterwarnings('ignore')
```

Import dataset

```
In [9]: df = pd.read_csv(r"D:\DA ALL NOTES\DAY24\heart.csv")
```

```
In [11]: df
```

Out[11]:

	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	0
1	37	1	2	130	250	0	1	187	0	3.5	0	0	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	0
...
298	57	0	0	140	241	0	1	123	1	0.2	1	0	0
299	45	1	3	110	264	0	1	132	0	1.2	1	0	0
300	68	1	0	144	193	1	1	141	0	3.4	1	2	0
301	57	1	0	130	131	0	1	115	1	1.2	1	1	0
302	57	0	1	130	236	0	0	174	0	0.0	1	1	1

303 rows × 14 columns



Exploratory Data Analysis

In [14]: `print('The shape of the dataset : ', df.shape)`

The shape of the dataset : (303, 14)

In [16]: `df.head()`

Out[16]:

	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 [18]: `df.tail()`

Out[18]:

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	target
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	

In [20]: 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 [22]: df.dtypes

```
Out[22]: 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 [24]: df.describe()

Out[24]:

	age	sex	cp	trestbps	chol	fbs	restecg
count	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	0.528000
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525800
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000
25%	47.500000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000



In [26]: `df.columns`

Out[26]: `Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach', 'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'], dtype='object')`

Univariate analysis

In [29]: `df['target'].nunique()`

Out[29]: 2

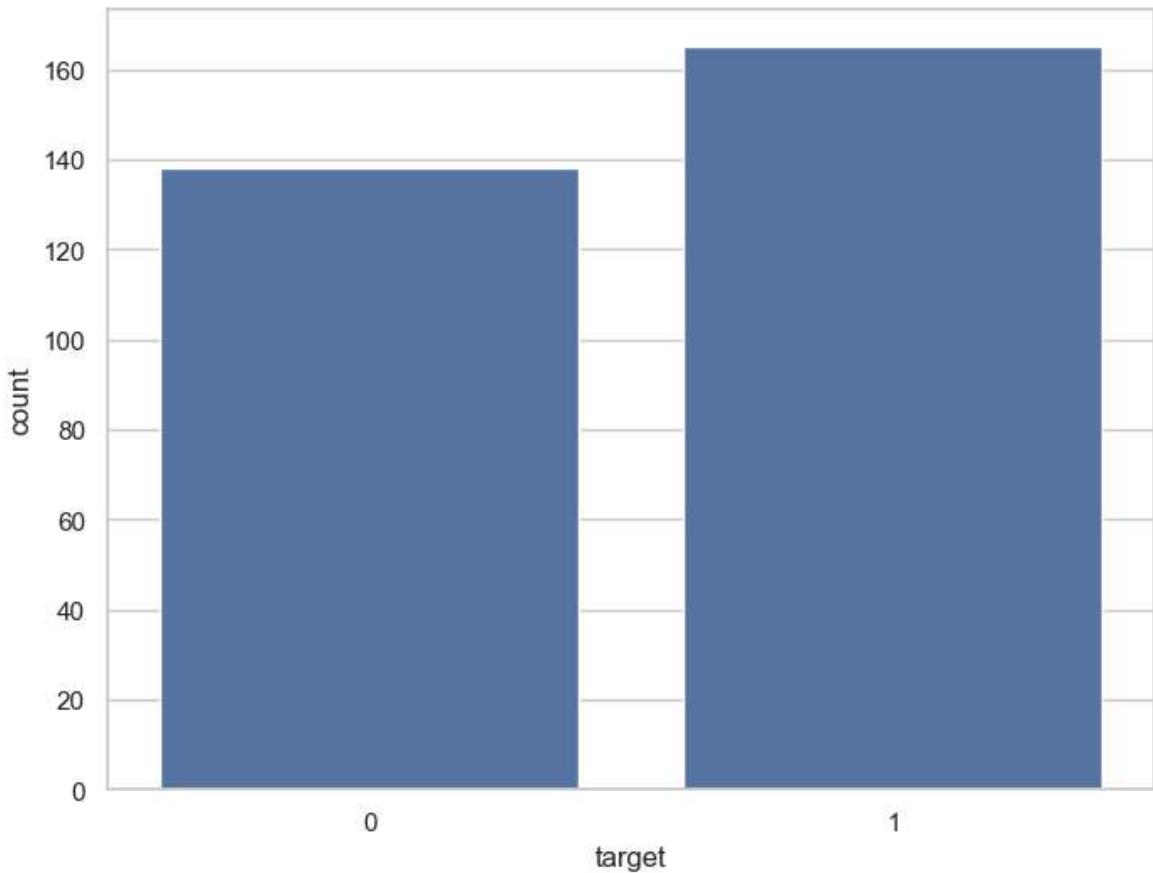
In [31]: `df['target'].unique()`

Out[31]: `array([1, 0], dtype=int64)`

In [33]: `df['target'].value_counts()`

Out[33]: `target`
1 165
0 138
Name: count, dtype: int64

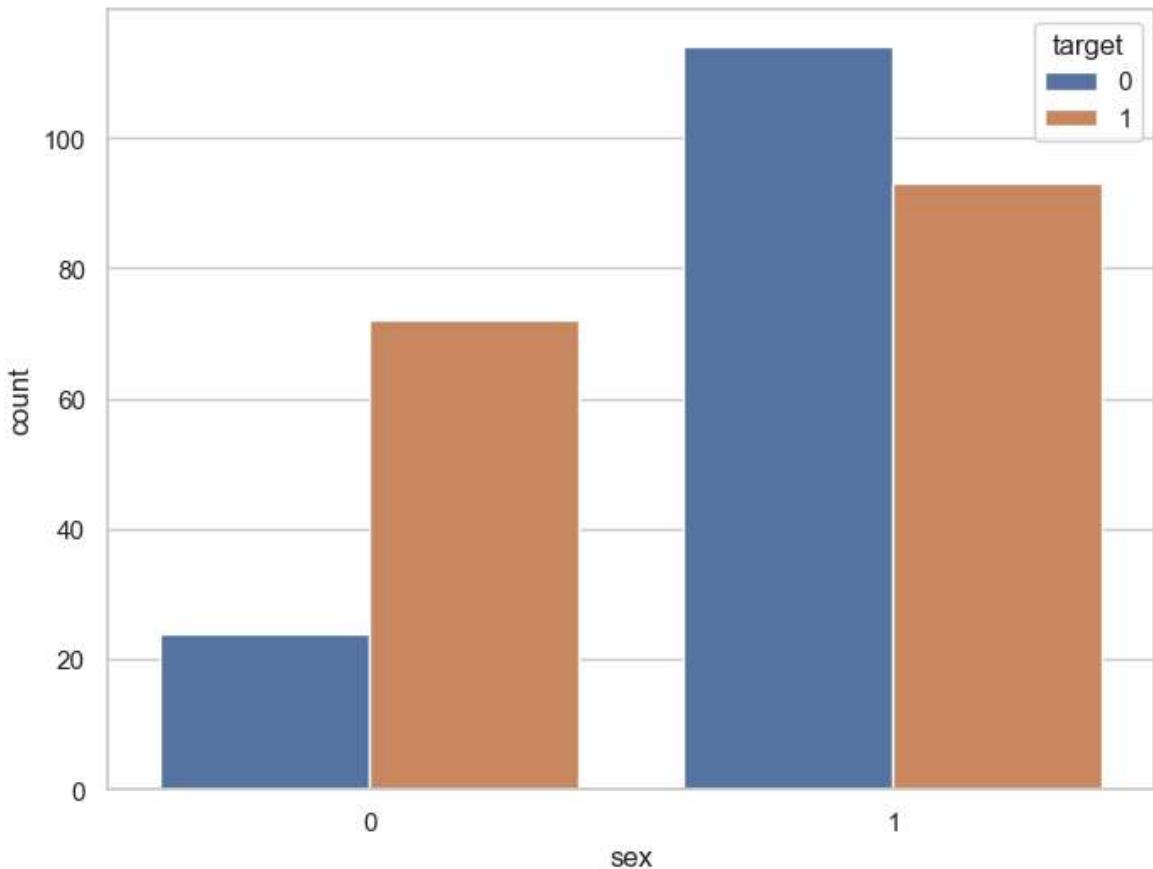
In [35]: `f,ax = plt.subplots(figsize=(8,6))
ax = sns.countplot(x='target',data=df)
plt.show()`



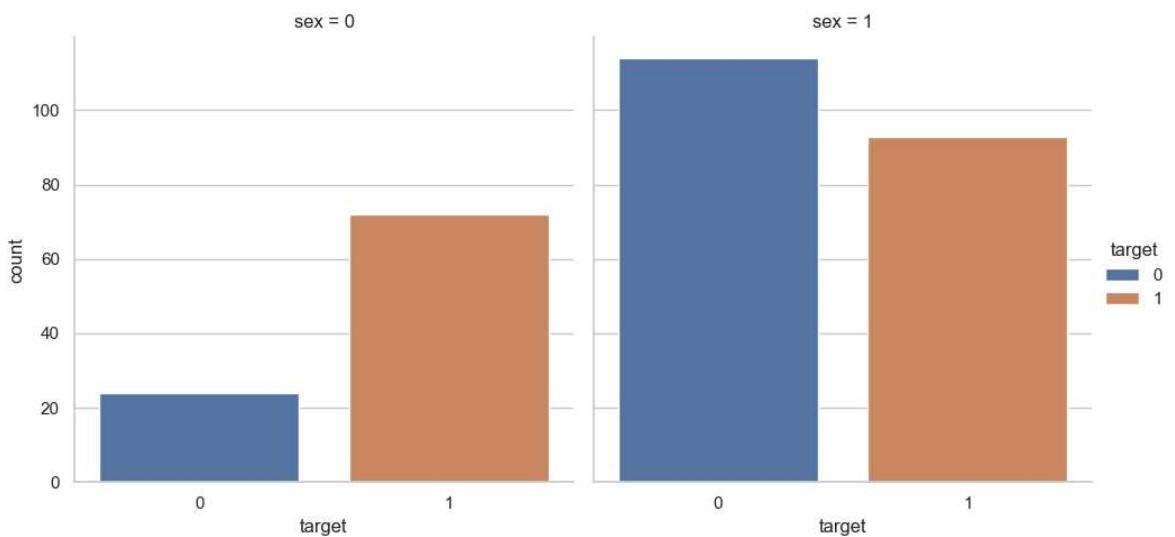
```
In [37]: df.groupby('sex')['target'].value_counts()
```

```
Out[37]:   sex  target
            0      1        72
                  0        24
            1      0       114
                  1        93
Name: count, dtype: int64
```

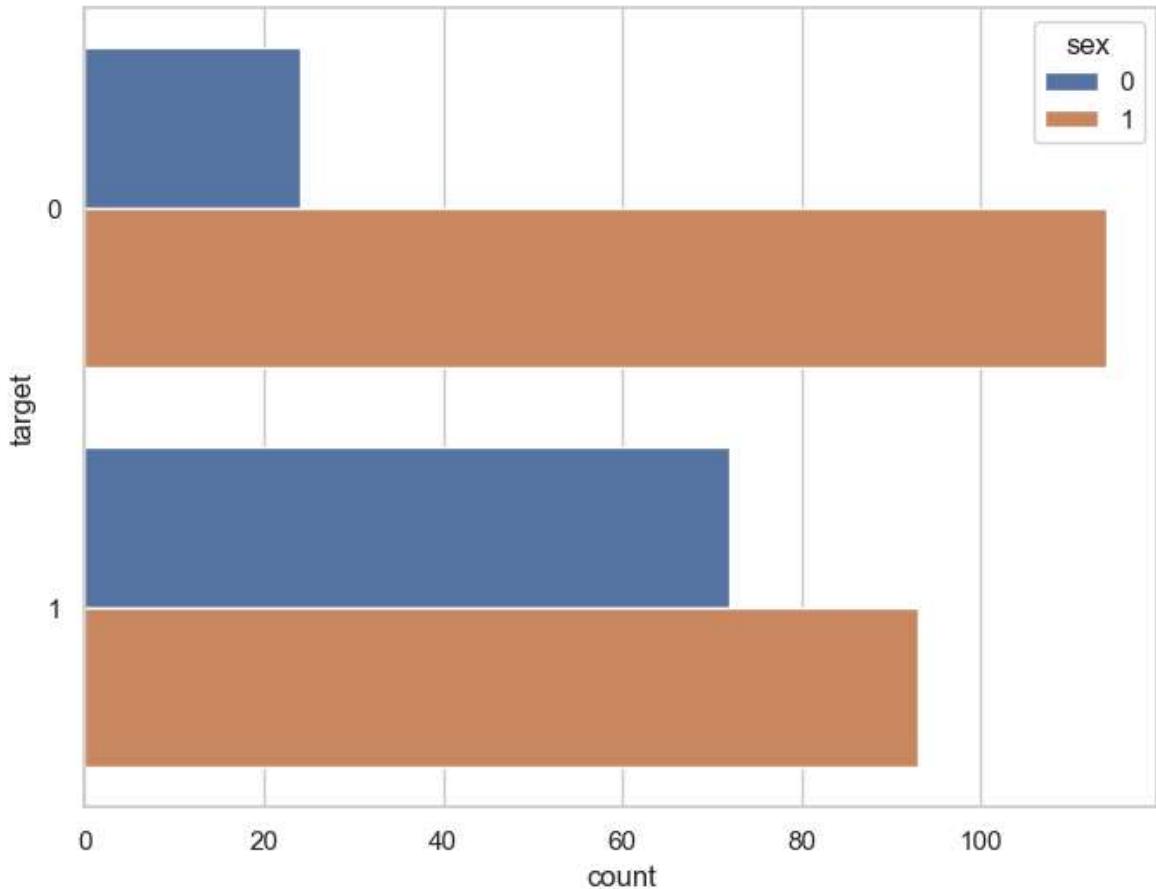
```
In [39]: f,ax = plt.subplots(figsize=(8,6))
ax=sns.countplot(x='sex',hue='target',data=df)
plt.show()
```



```
In [41]: ax = sns.catplot(x='target', col='sex', data=df, kind='count', height=5, aspect=1, hue
```

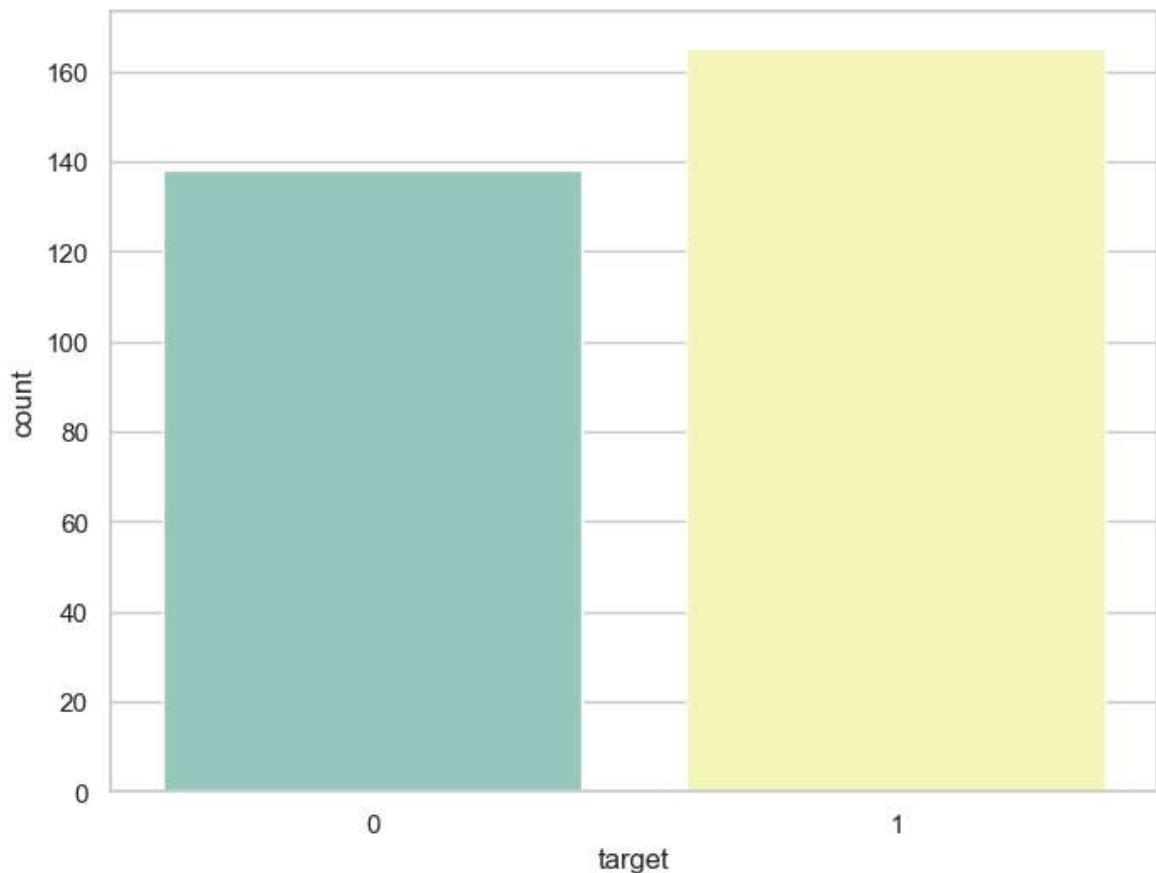


```
In [43]: f,ax = plt.subplots(figsize=(8,6))
ax = sns.countplot(y='target',hue='sex',data=df,)
plt.show()
```

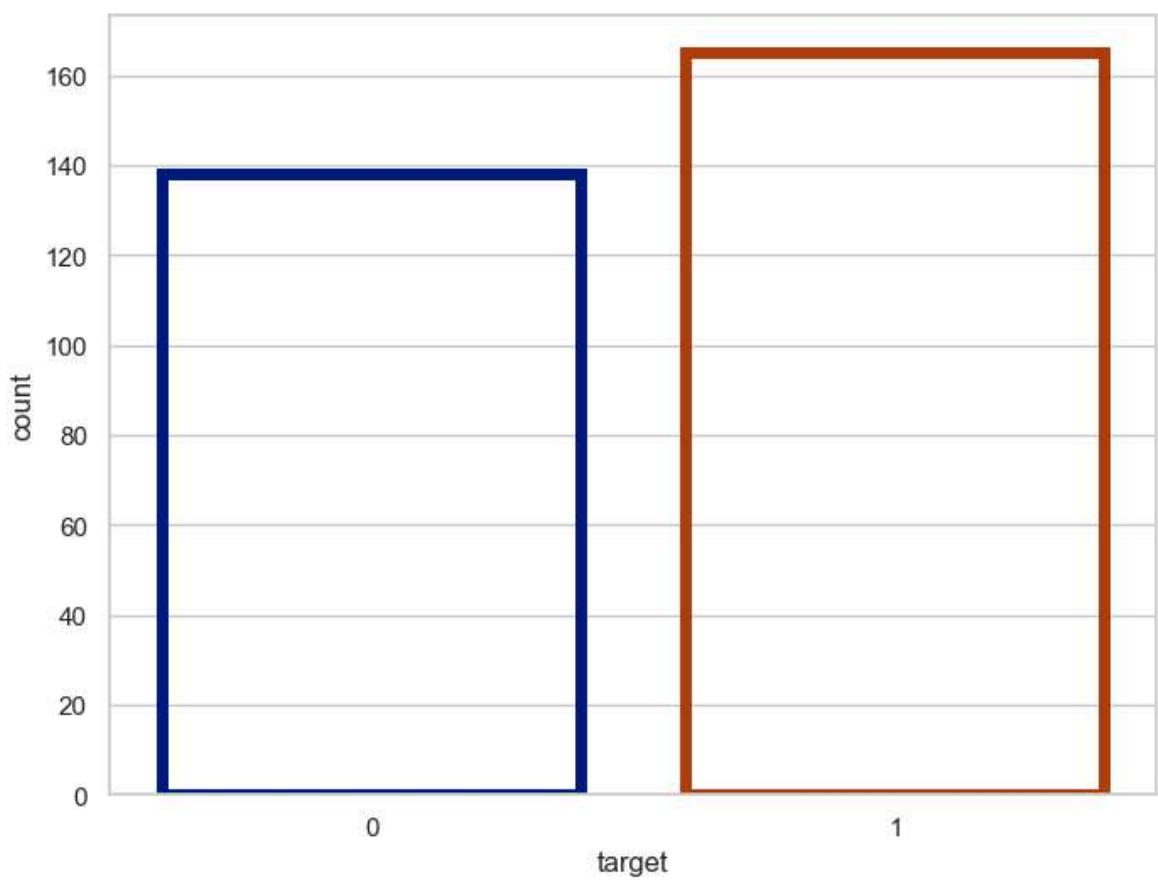


```
In [45]: f,ax = plt.subplots(figsize=(8,6))
ax = sns.countplot(x='target',data=df,palette='Set3')
plt.show
```

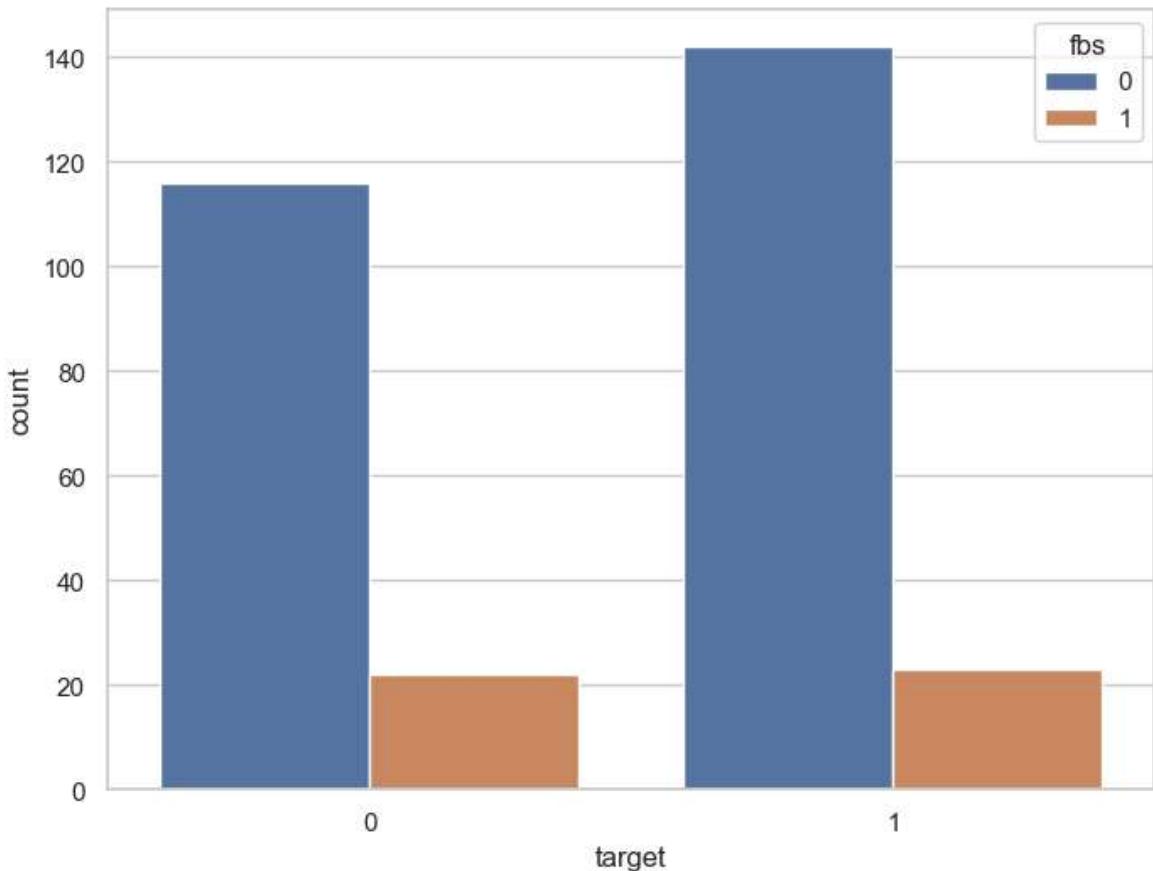
```
Out[45]: <function matplotlib.pyplot.show(close=None, block=None)>
```



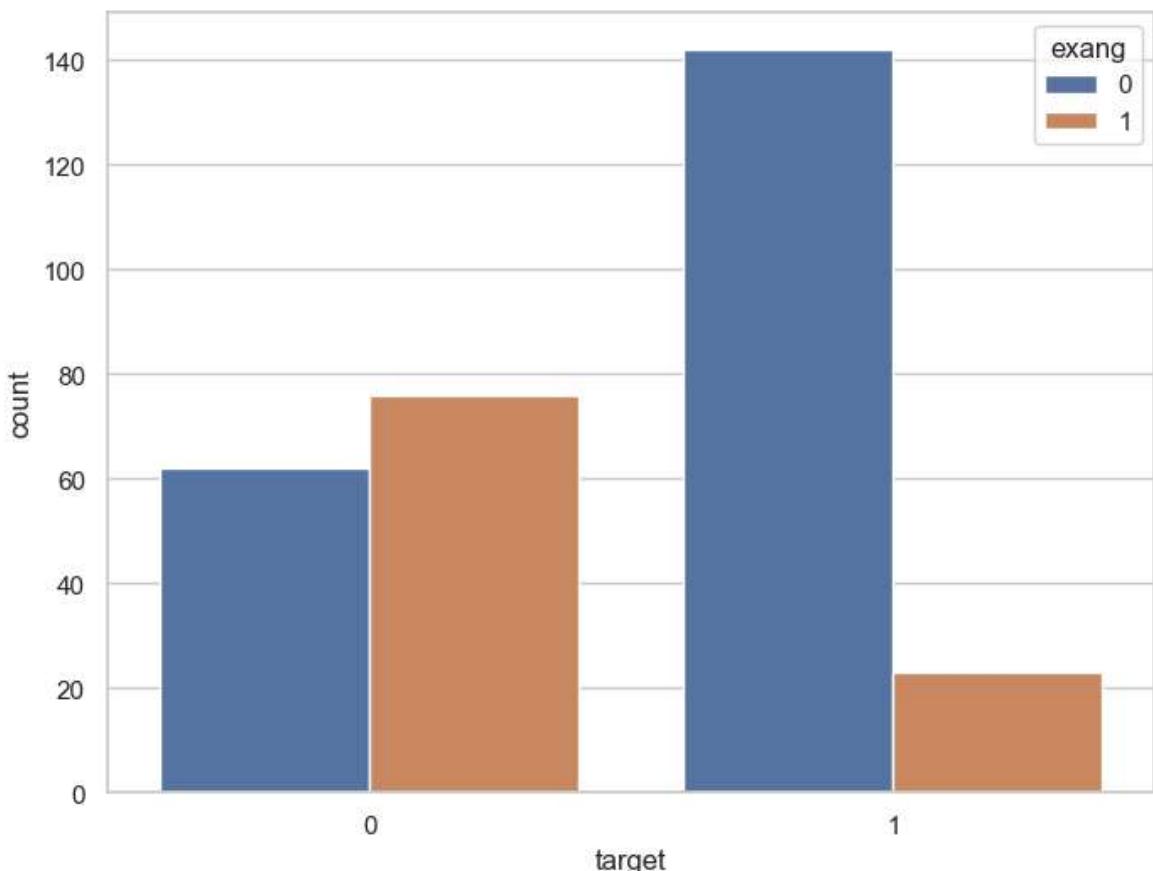
```
In [47]: f,ax = plt.subplots(figsize = (8,6))
ax=sns.countplot(x='target',data=df,facecolor=(0,0,0,0),linewidth=5,edgecolor=sns
plt.show()
```



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



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



```
In [53]: correlation =df.corr()
```

```
In [55]: correlation['target'].sort_values(ascending=False)
```

```
Out[55]: 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
```

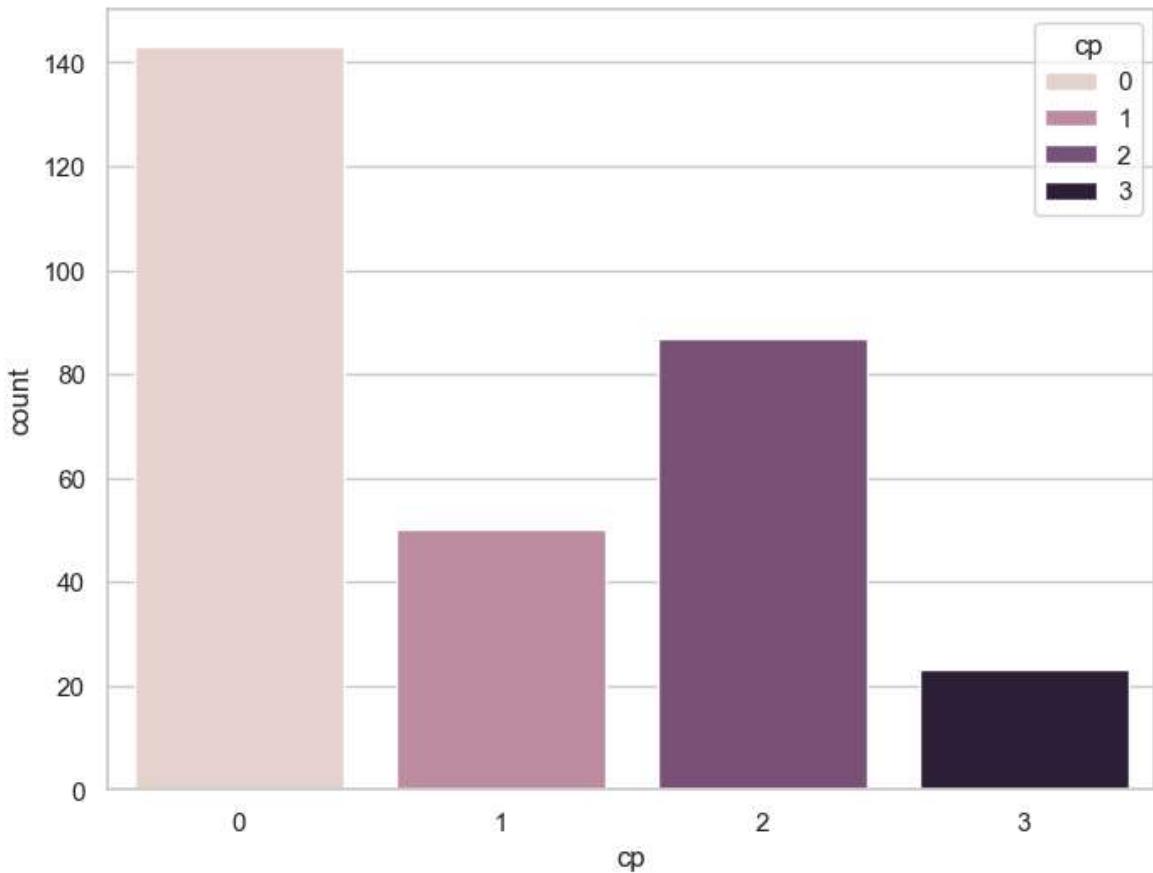
```
In [57]: df['cp'].nunique()
```

```
Out[57]: 4
```

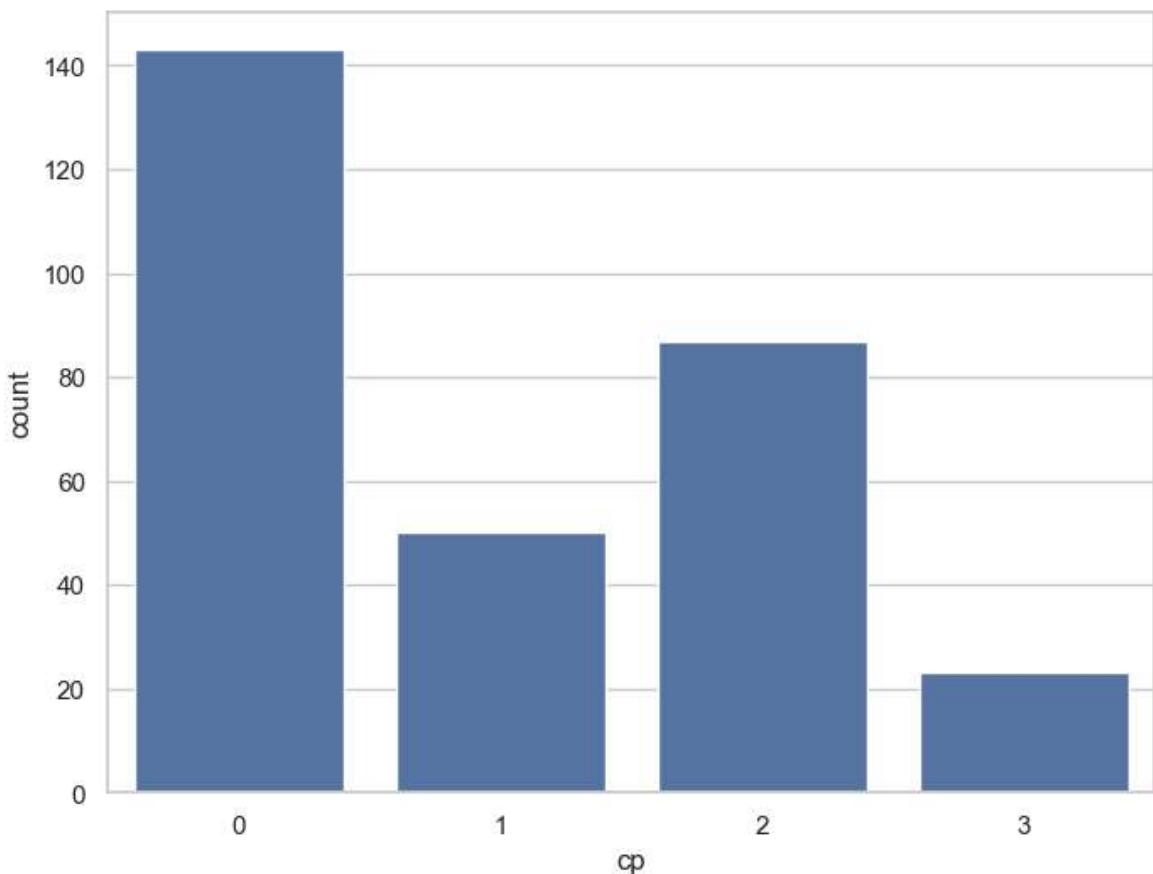
```
In [59]: df['cp'].value_counts()
```

```
Out[59]: cp
          0    143
          2     87
          1     50
          3     23
          Name: count, dtype: int64
```

```
In [61]: f,ax = plt.subplots(figsize =(8,6))
          ax = sns.countplot(x = 'cp',data=df,hue='cp')
          plt.show()
```



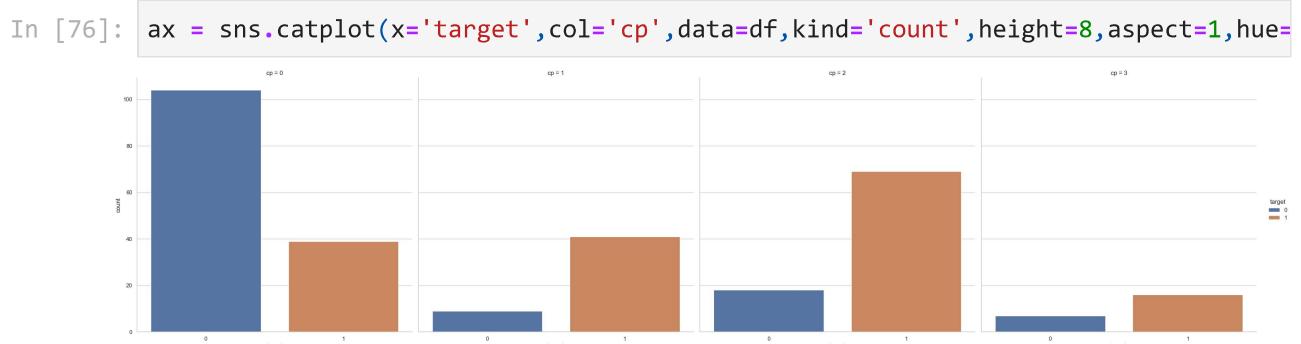
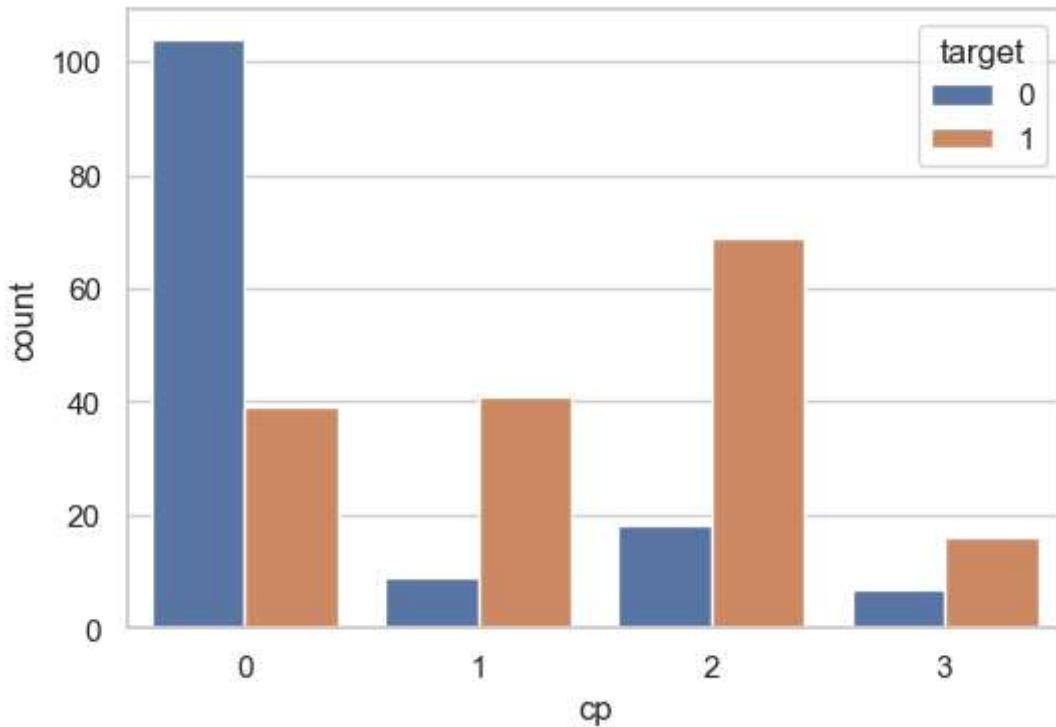
```
In [63]: f,ax = plt.subplots(figsize =(8,6))
ax = sns.countplot(x = 'cp',data=df)
plt.show()
```



```
In [65]: df.groupby('cp')[ 'target'].value_counts()
```

```
Out[65]: 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
```

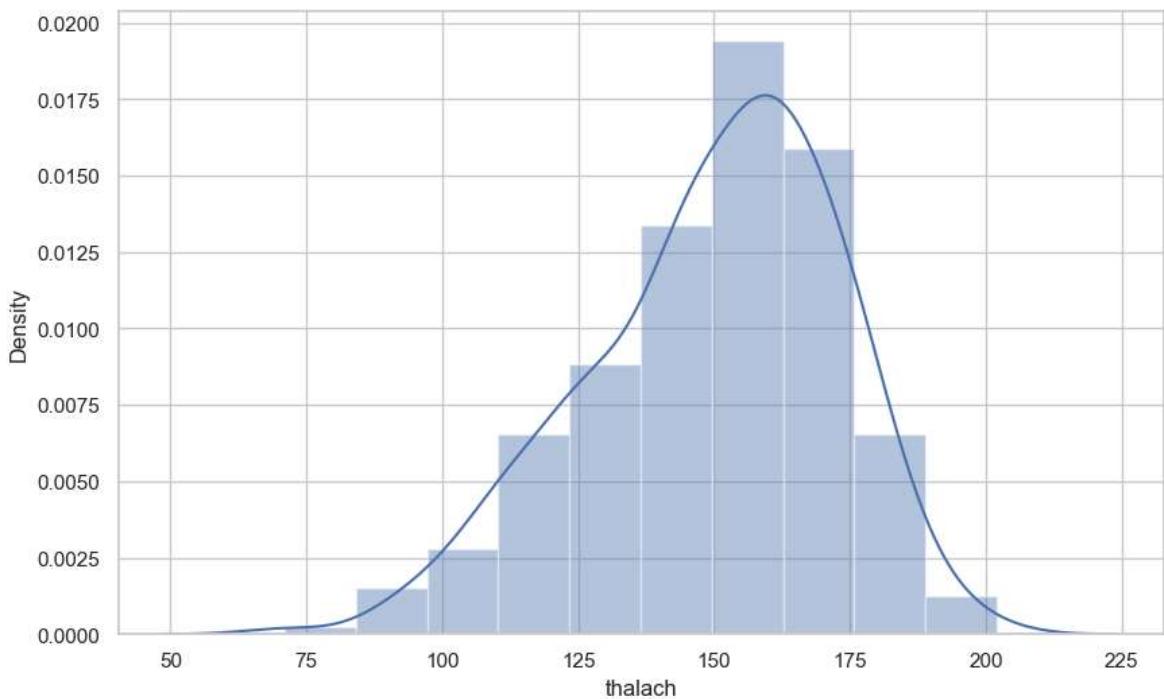
```
In [72]: f,ax =plt.subplots(figsize =(6,4))
ax = sns.countplot(x='cp',hue='target',data=df)
plt.show()
```



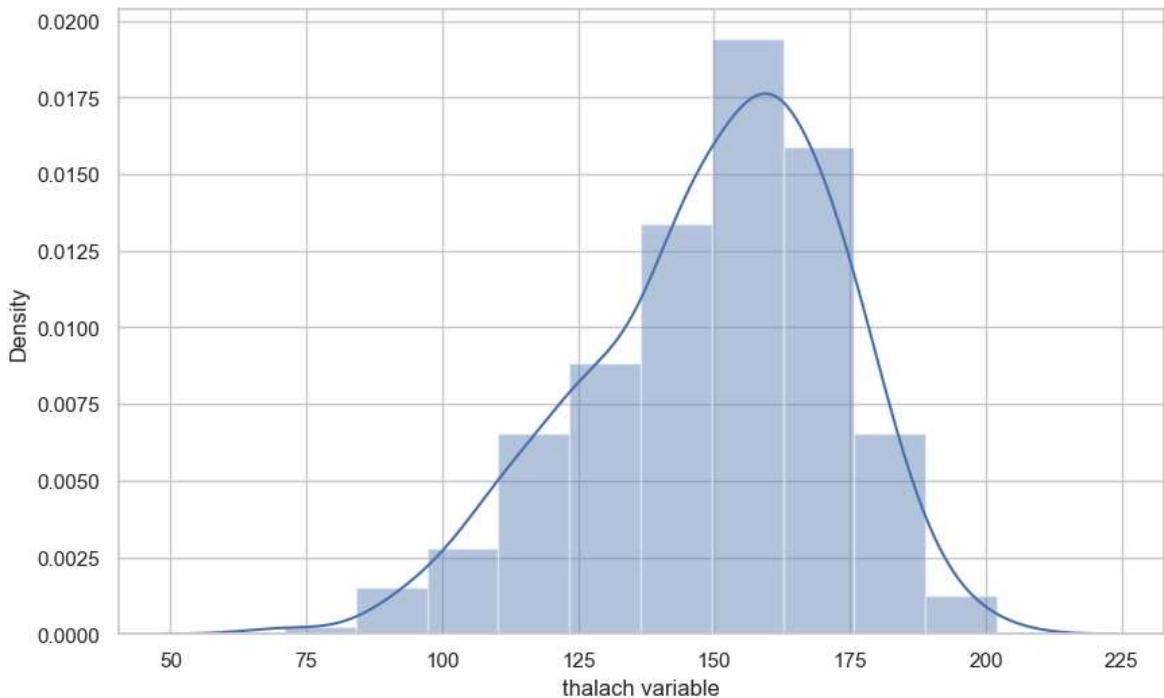
```
In [80]: df['thalach'].nunique()
```

```
Out[80]: 91
```

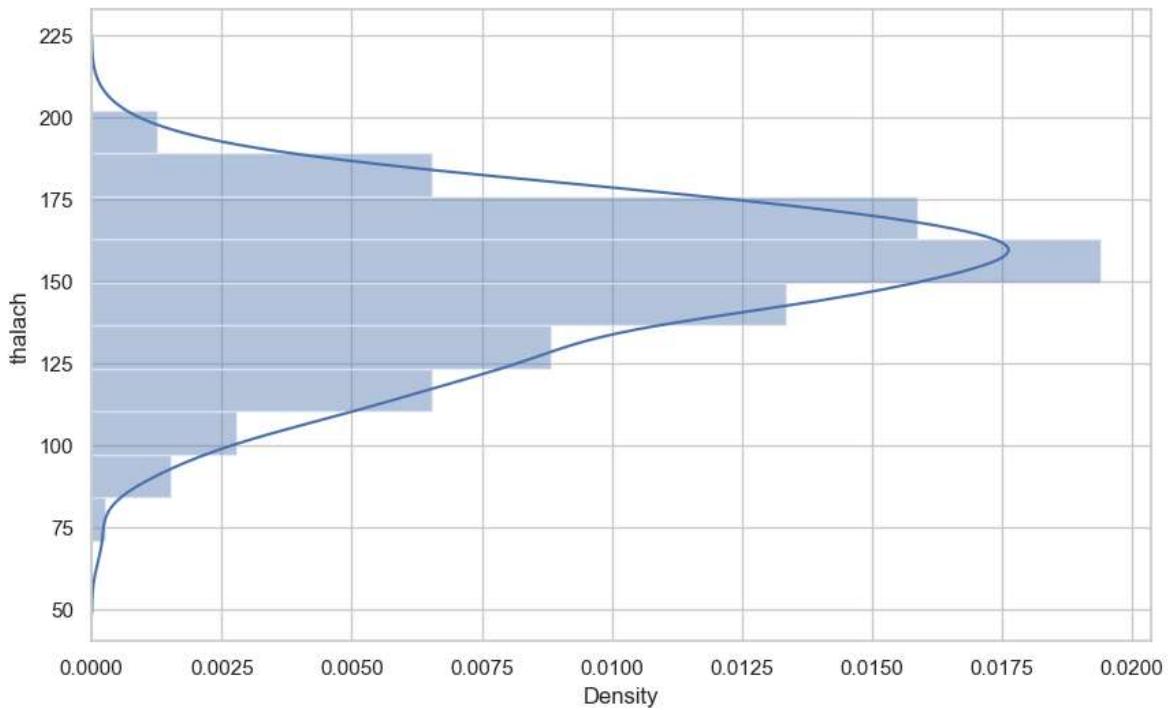
```
In [86]: f,ax=plt.subplots(figsize=(10,6))
x=df['thalach']
ax = sns.distplot(x,bins=10)
plt.show()
```



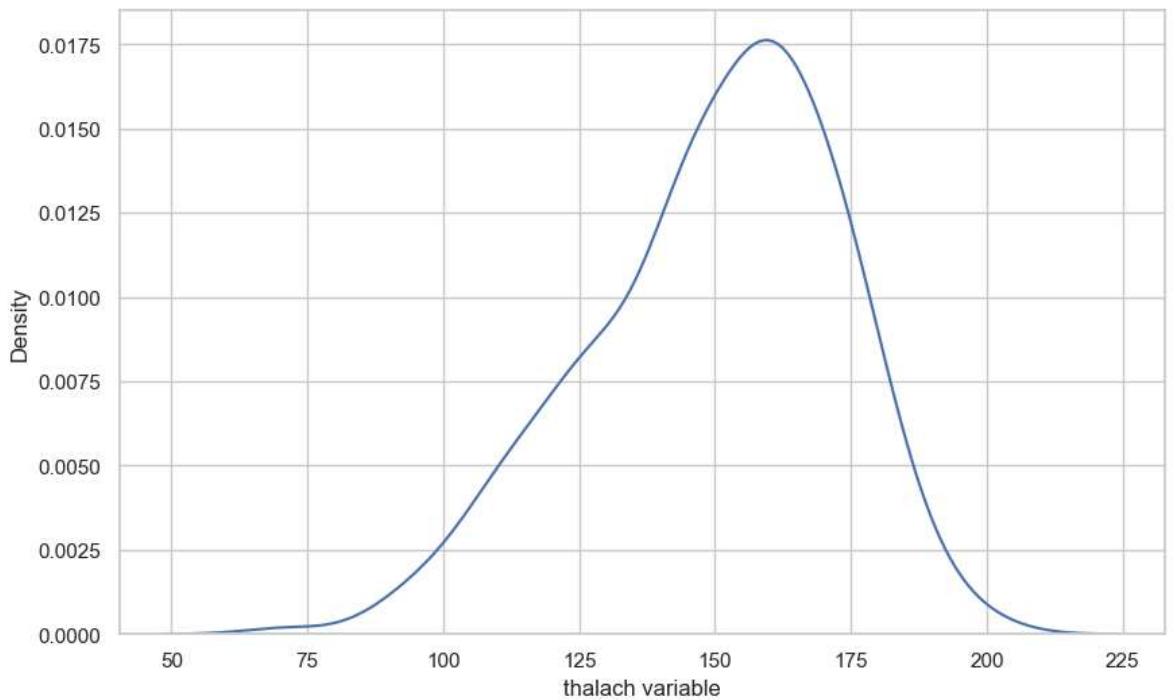
```
In [88]: f,ax=plt.subplots(figsize=(10,6))
x=df['thalach']
x=pd.Series(x,name='thalach variable')
ax=sns.distplot(x,bins=10)
plt.show()
```



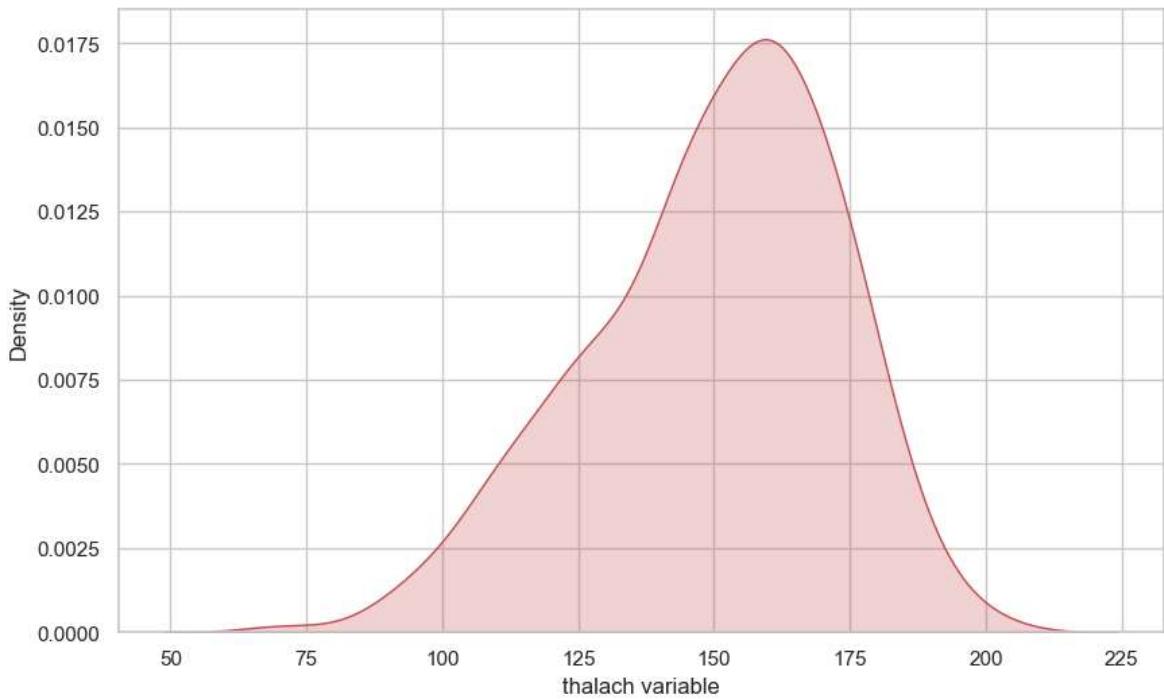
```
In [92]: f,ax= plt.subplots(figsize=(10,6))
x=df['thalach']
ax=sns.distplot(x,bins=10,vertical=True)
plt.show()
```



```
In [96]: f,ax=plt.subplots(figsize=(10,6))
x=df['thalach']
x=pd.Series(x,name="thalach variable")
ax=sns.kdeplot(x)
plt.show()
```

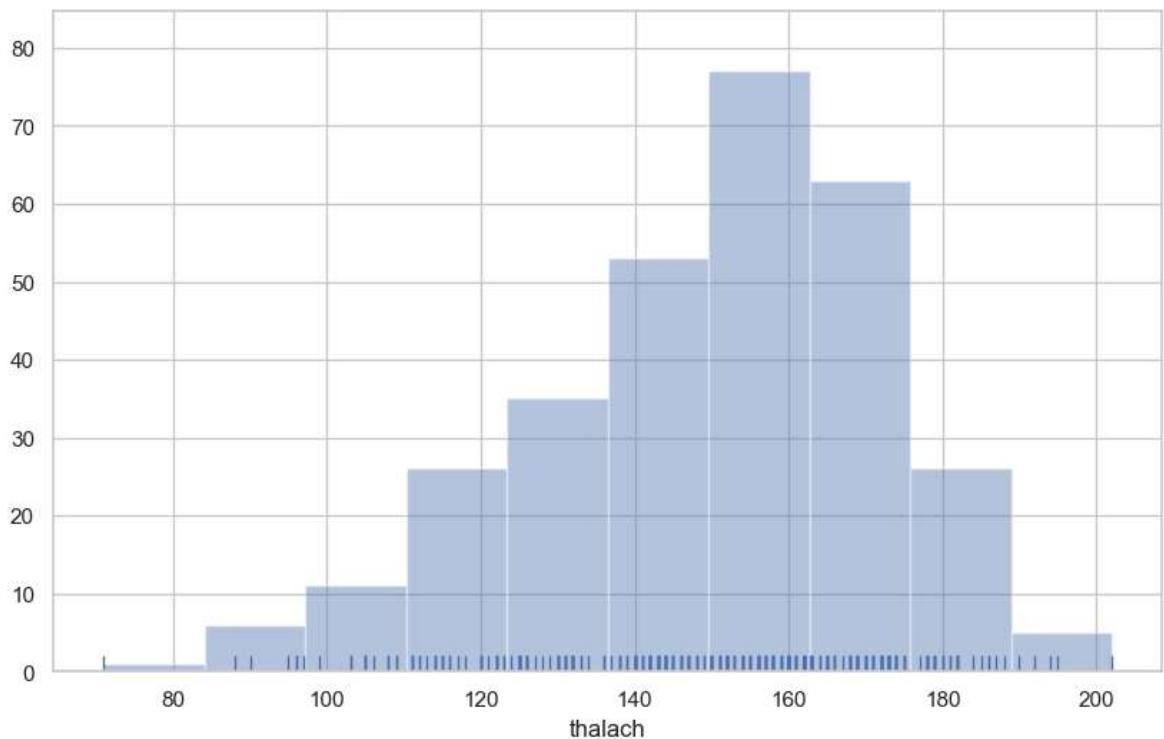


```
In [98]: f,ax=plt.subplots(figsize=(10,6))
x=df['thalach']
x=pd.Series(x,name="thalach variable")
ax=sns.kdeplot(x,shade =True,color='r')
plt.show()
```

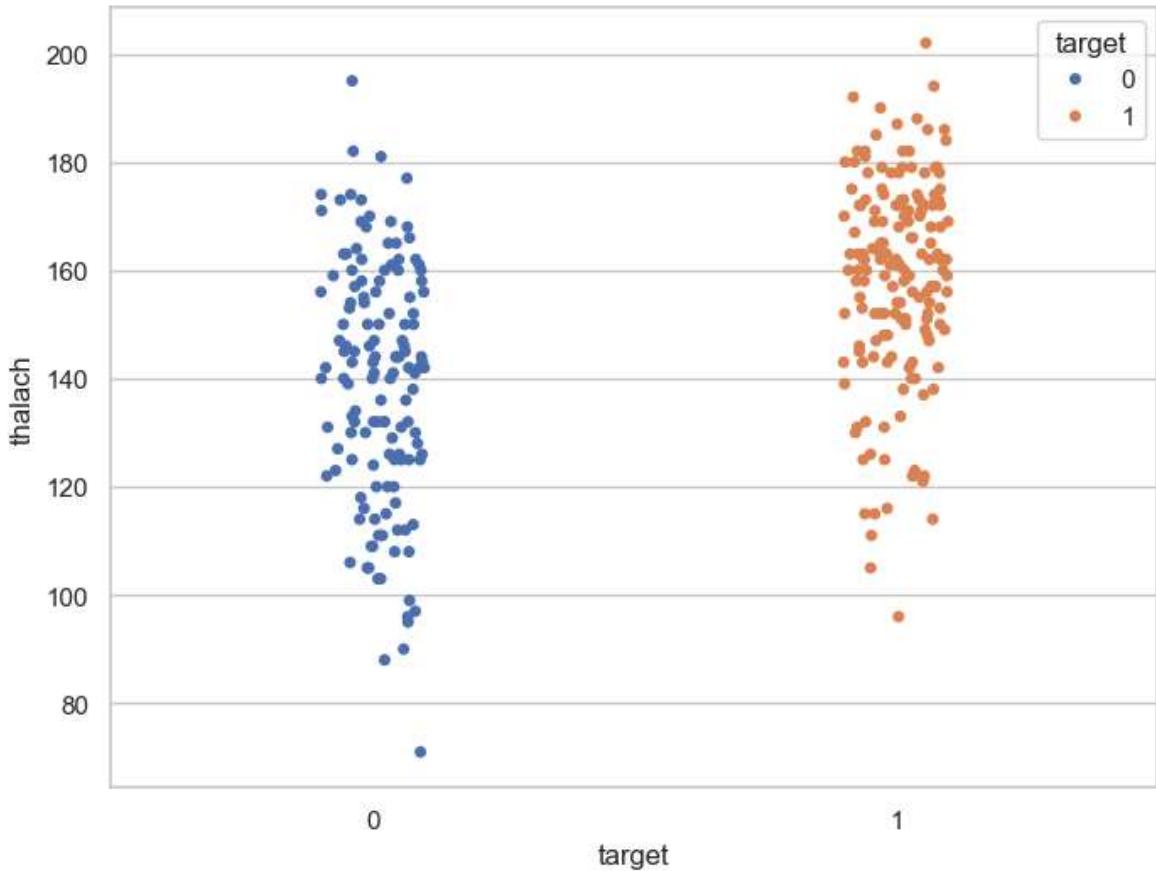


Histogram

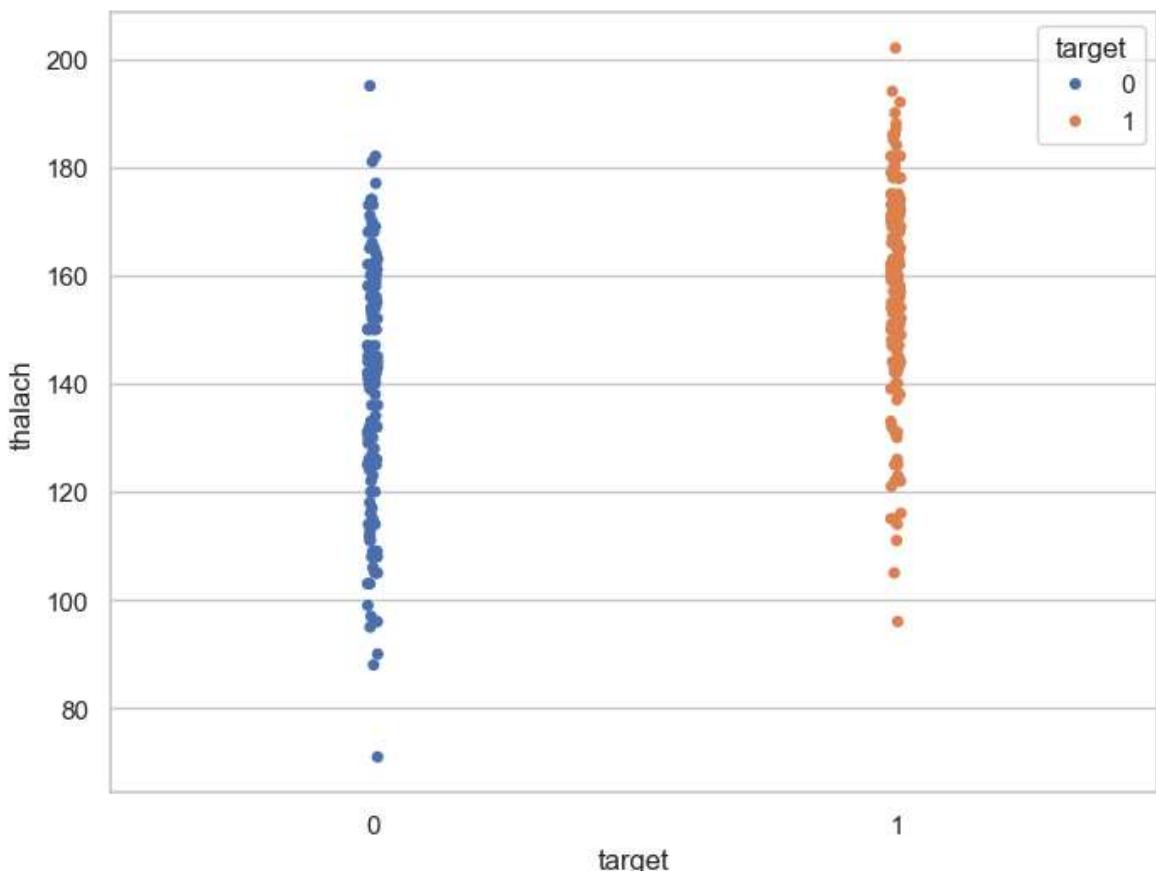
```
In [103...]:  
f,ax=plt.subplots(figsize=(10,6))  
x=df['thalach']  
ax=sns.distplot(x,kde=False,rug=True,bins=10)  
plt.show()
```



```
In [109...]:  
f,ax=plt.subplots(figsize=(8,6))  
sns.stripplot(x='target',y='thalach',data=df,hue='target')  
plt.show()
```

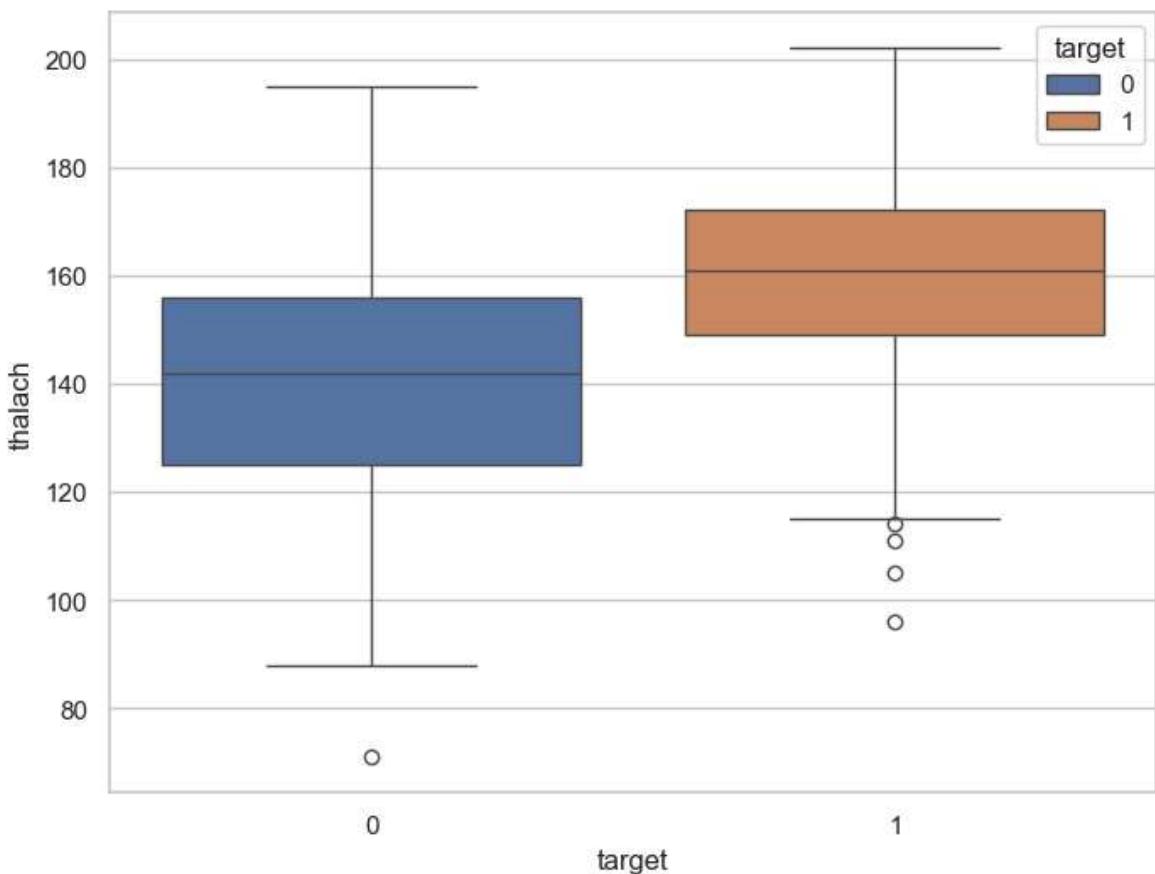


```
In [111]: f,ax=plt.subplots(figsize=(8,6))
sns.stripplot(x='target',y='thalach',data=df,jitter=0.01,hue='target')
plt.show()
```

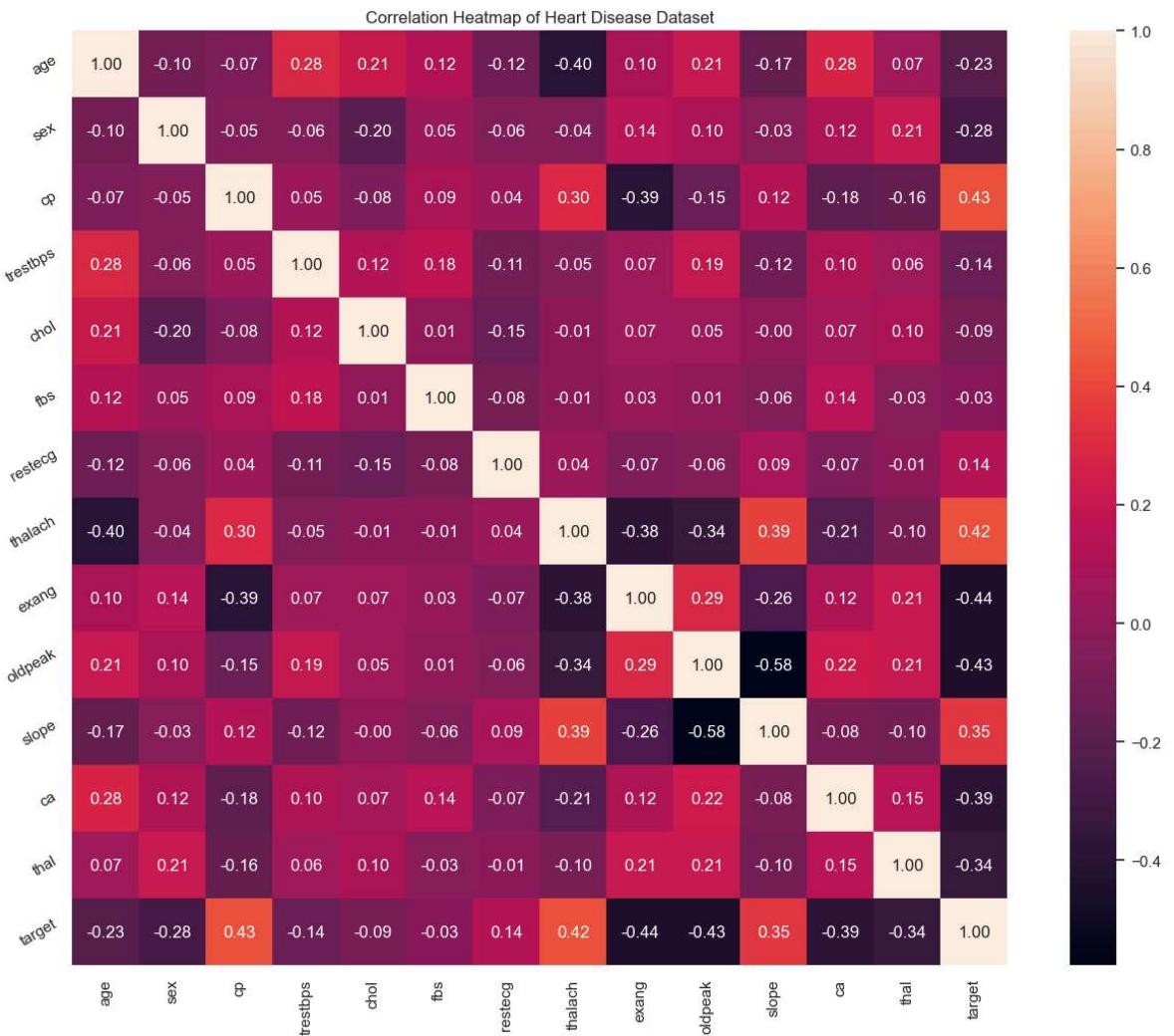


```
In [113]: f,ax=plt.subplots(figsize=(8,6))
sns.boxplot(x='target',y='thalach',data=df,hue='target')
```

```
plt.show()
```



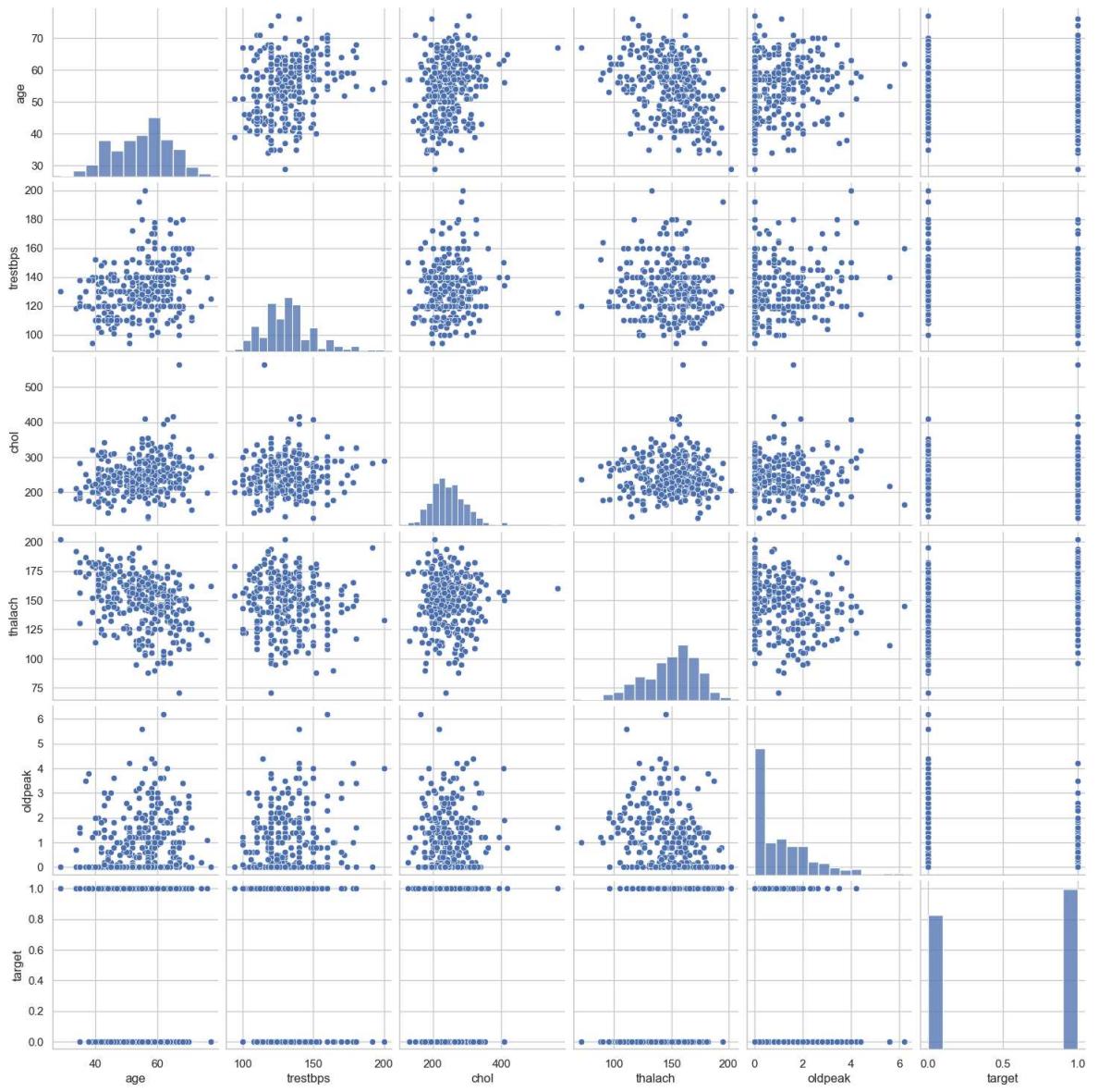
```
In [121]: plt.figure(figsize=(16,12))
plt.title('Correlation Heatmap of Heart Disease Dataset')
a=sns.heatmap(correlation,square=True,annot=True,fmt=' .2f',linecolor='white')
a.set_xticklabels(a.get_xticklabels(),rotation=90)
a.set_yticklabels(a.get_yticklabels(),rotation=30)
plt.show()
```



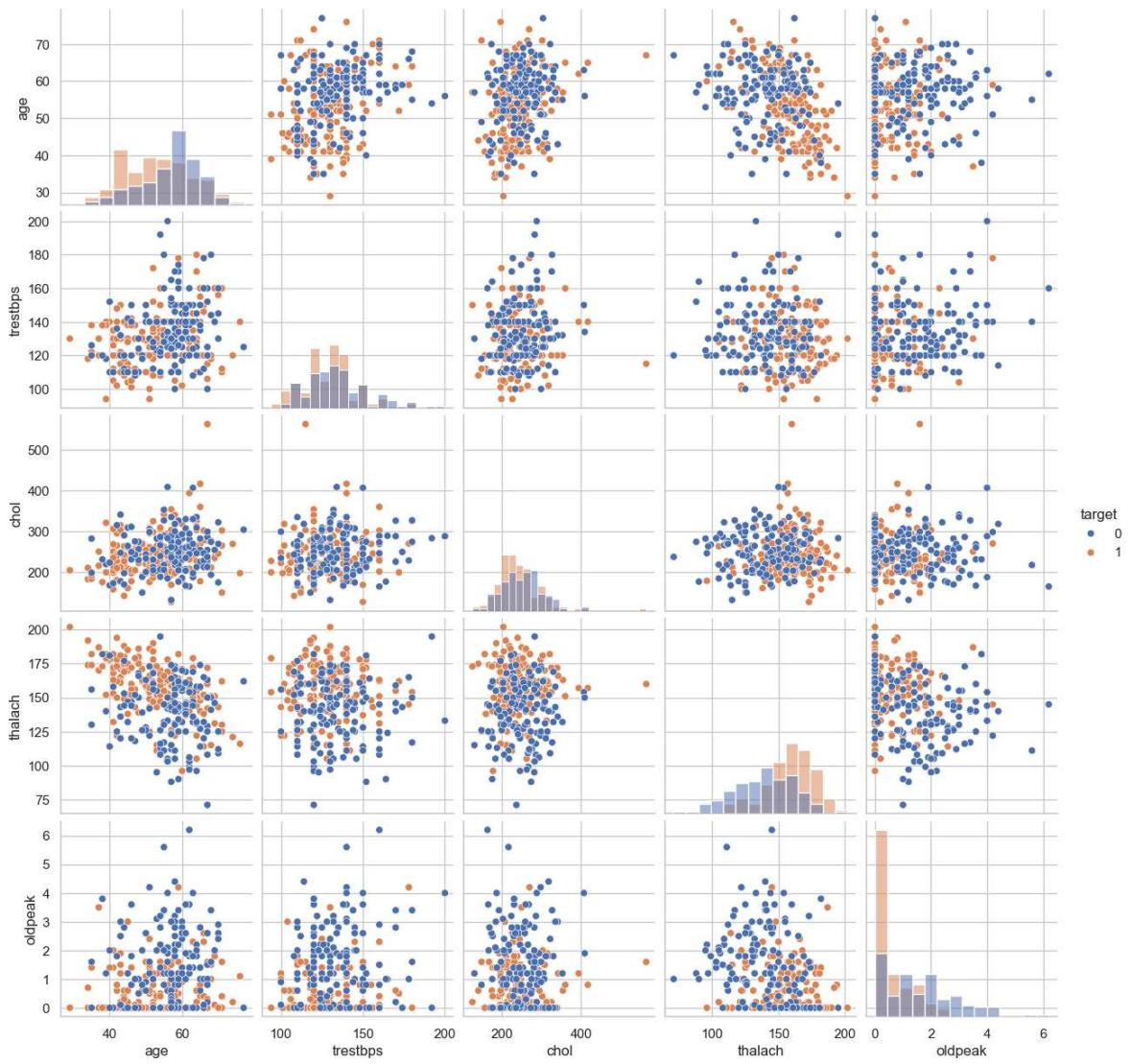
```
In [123... df.columns
```

```
Out[123... Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
       'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
      dtype='object')
```

```
In [131... num_var=['age','trestbps','chol','thalach','oldpeak','target']
sns.pairplot(df[num_var],kind='scatter',diag_kind='hist')
plt.show()
```



```
In [127]: num_var=['age','trestbps','chol','thalach','oldpeak','target']
sns.pairplot(df[num_var],kind='scatter',diag_kind='hist',hue='target')
plt.show()
```



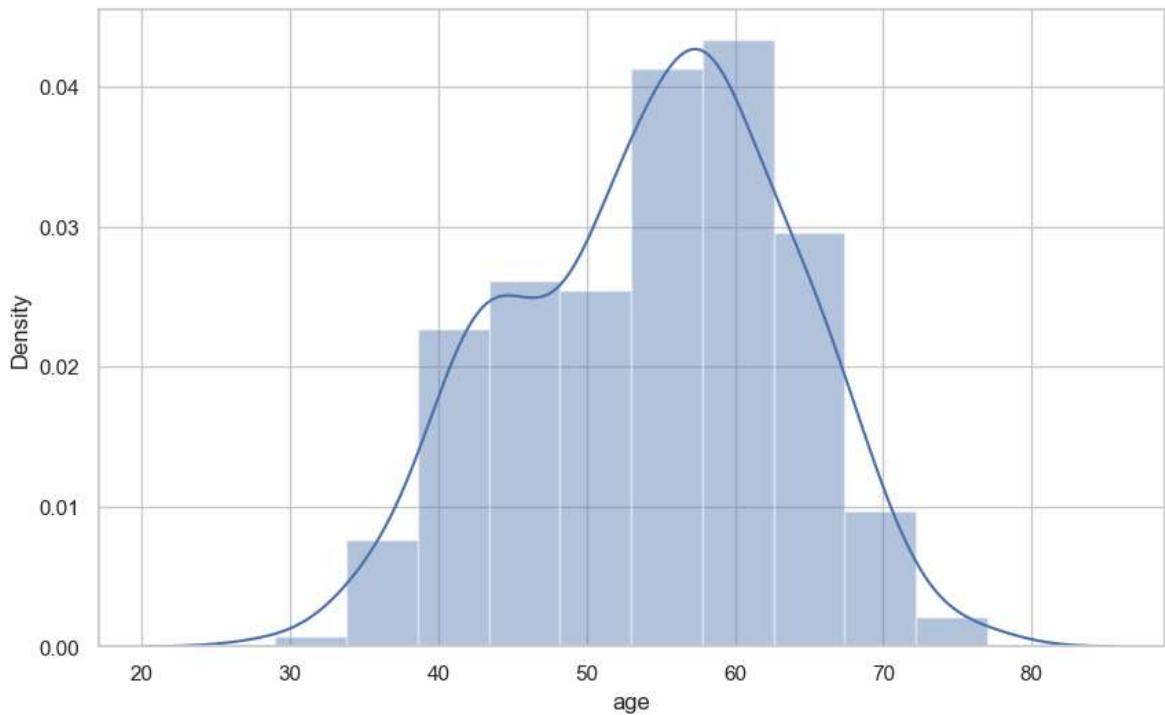
```
In [135]: df['age'].nunique()
```

```
Out[135]: 41
```

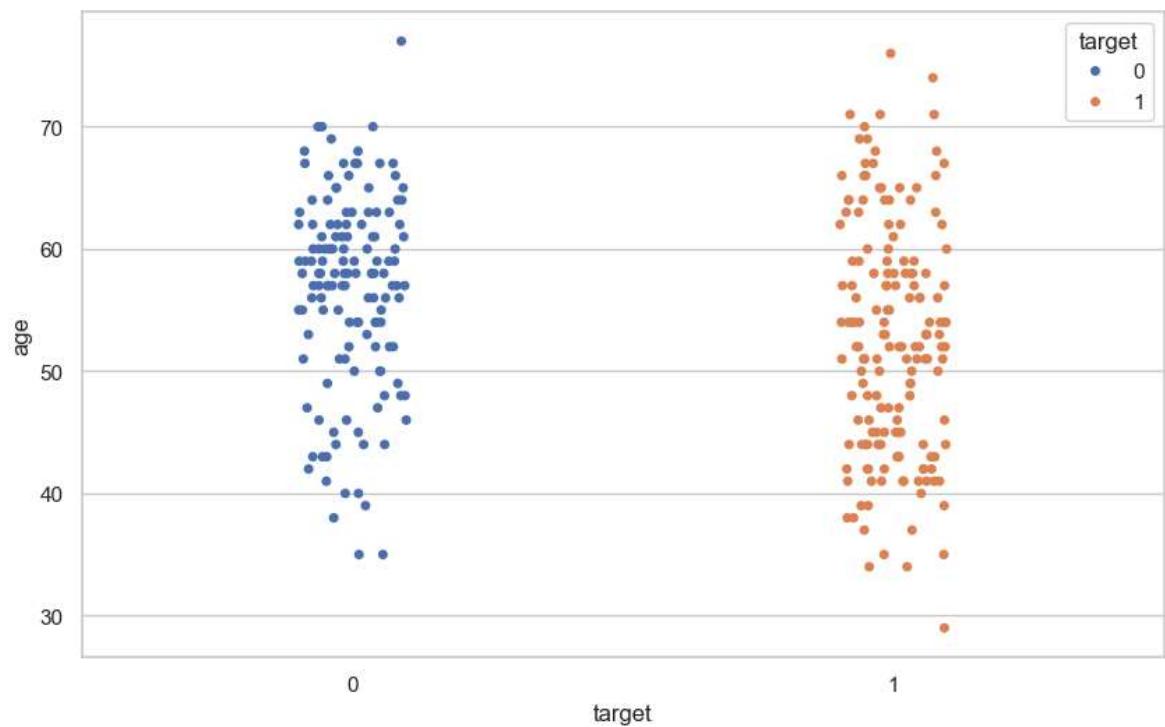
```
In [137]: df['age'].describe()
```

```
Out[137]: 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
```

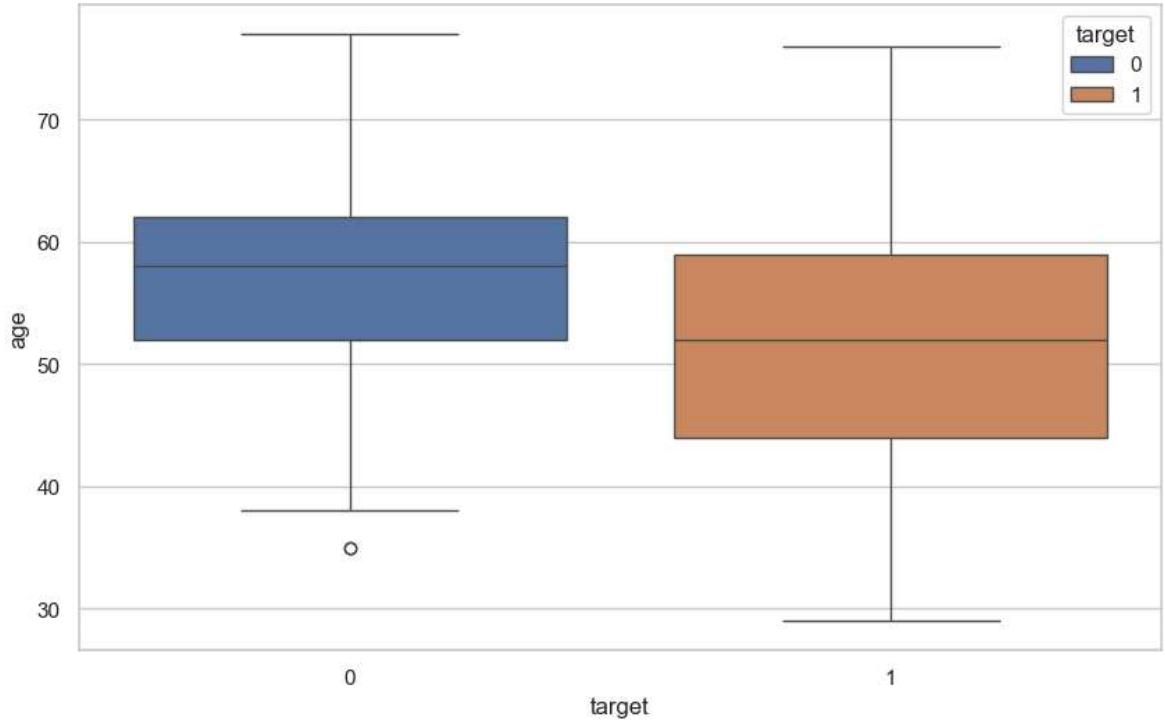
```
In [141]: f,ax=plt.subplots(figsize=(10,6))
x=df['age']
ax=sns.distplot(x,bins=10)
plt.show()
```



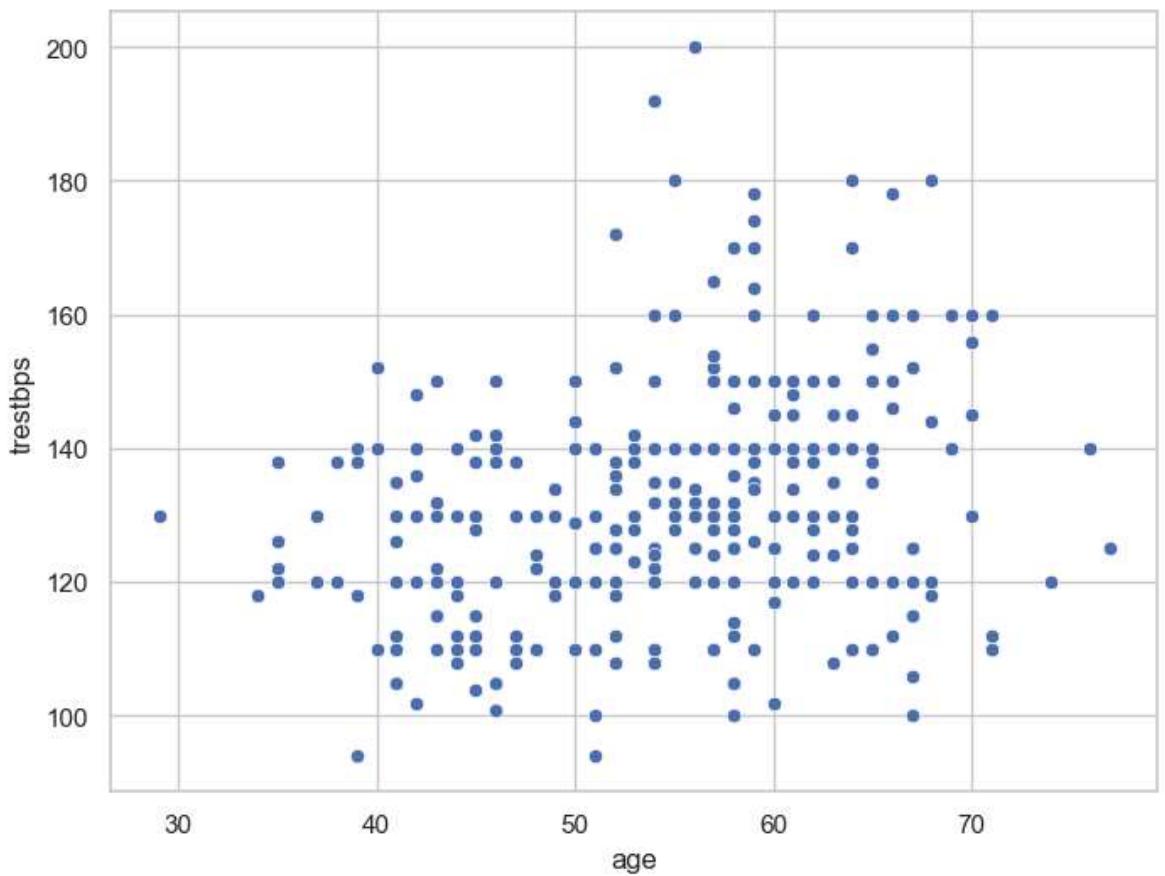
```
In [145...]:  
f,ax=plt.subplots(figsize=(10,6))  
sns.stripplot(x='target',y='age',data=df,hue='target')  
plt.show()
```



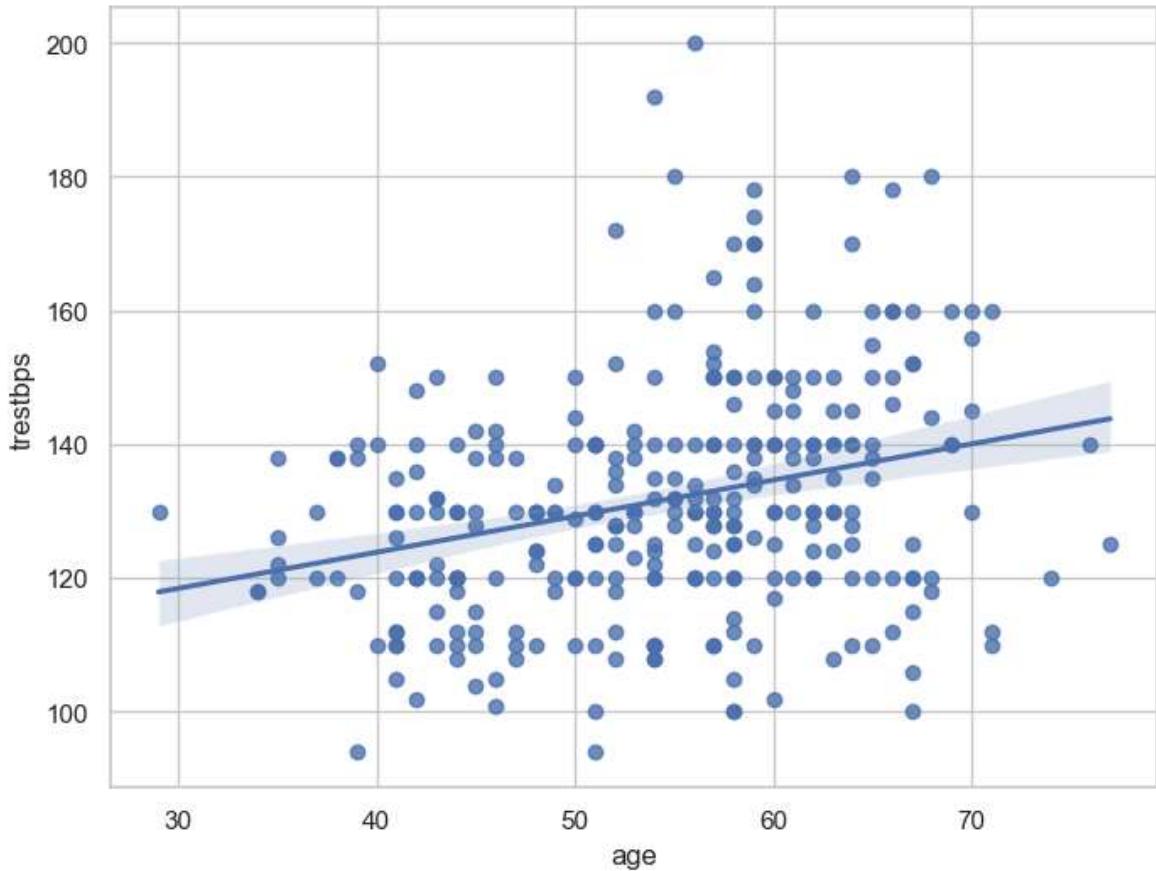
```
In [147...]:  
f,ax=plt.subplots(figsize=(10,6))  
sns.boxplot(x='target',y='age',data=df,hue='target')  
plt.show()
```



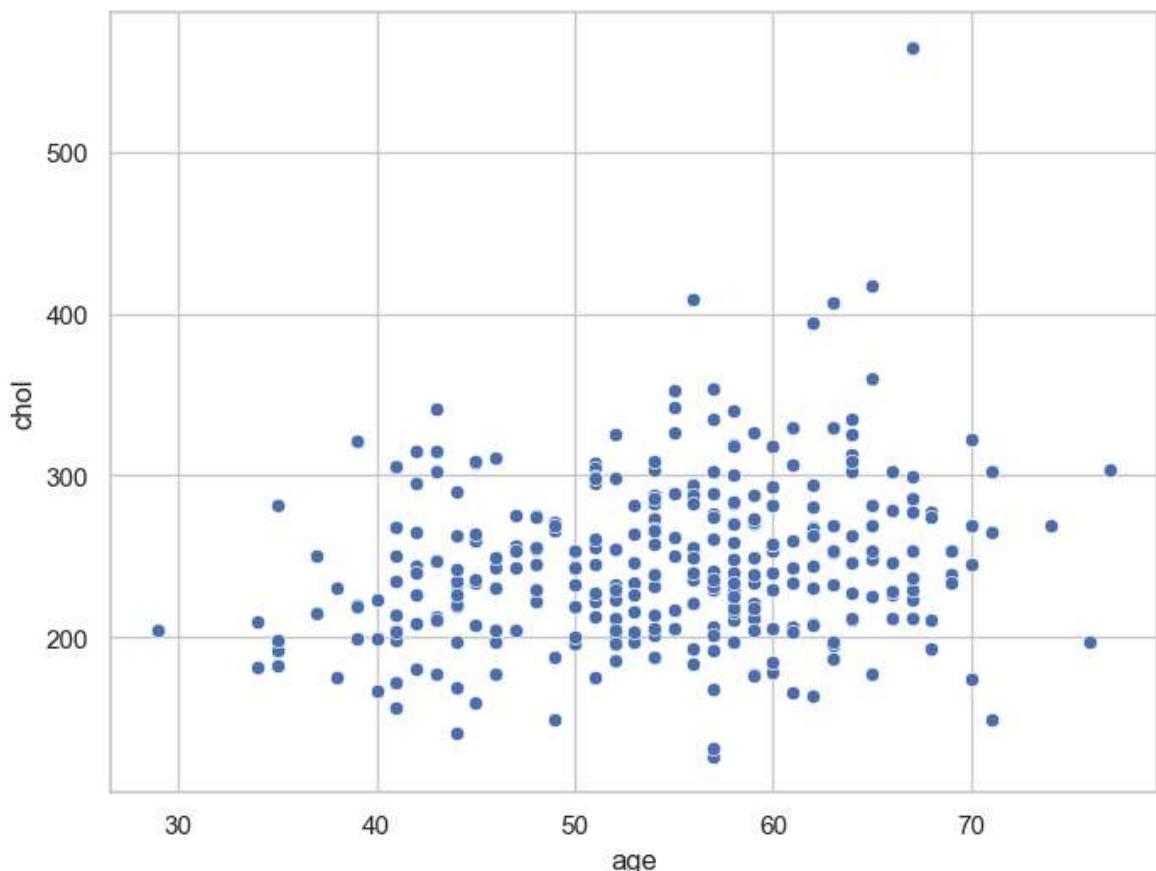
```
In [159]:  
f,ax=plt.subplots(figsize=(8,6))  
ax=sns.scatterplot(x='age',y='trestbps',data=df)  
plt.show()
```



```
In [161]:  
f,ax=plt.subplots(figsize=(8,6))  
ax=sns.regplot(x='age',y='trestbps',data=df)  
plt.show()
```

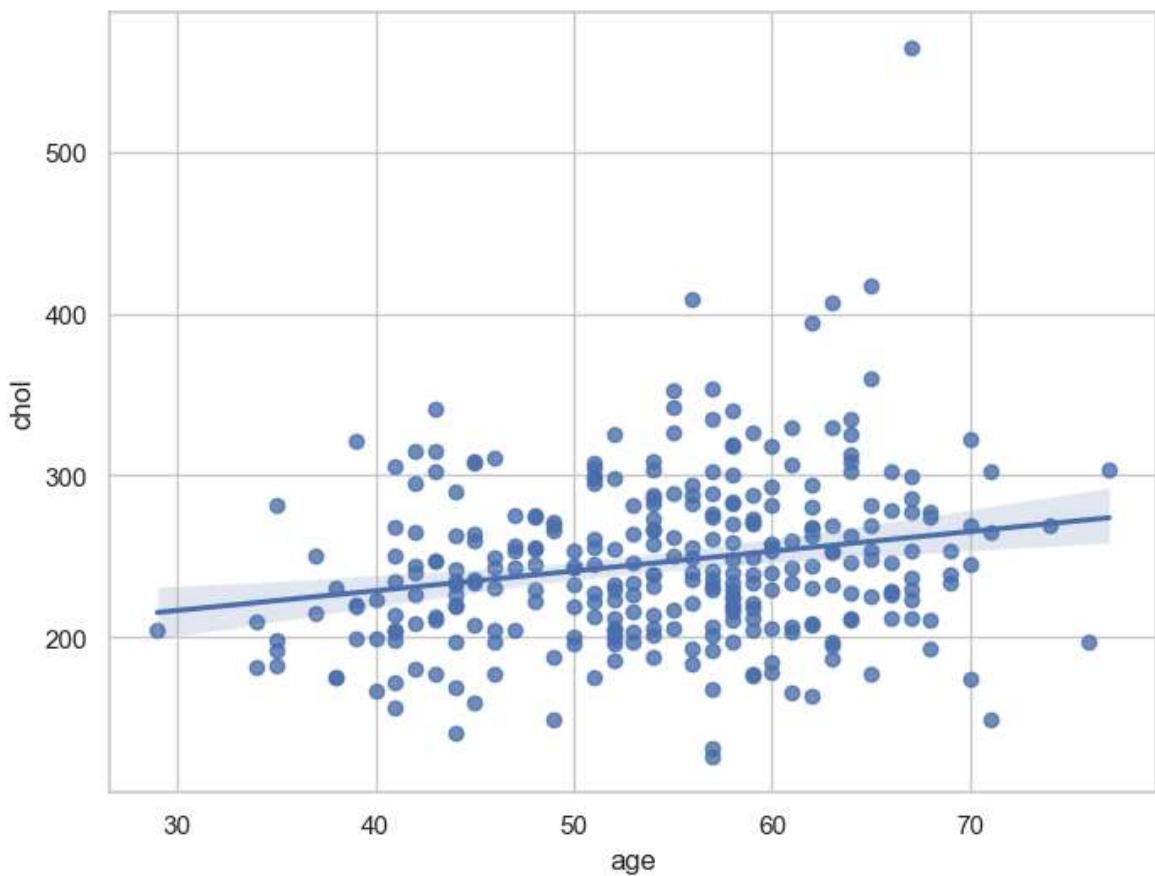


```
In [163]: f,ax=plt.subplots(figsize=(8,6))
ax=sns.scatterplot(x='age',y='chol',data=df)
plt.show()
```

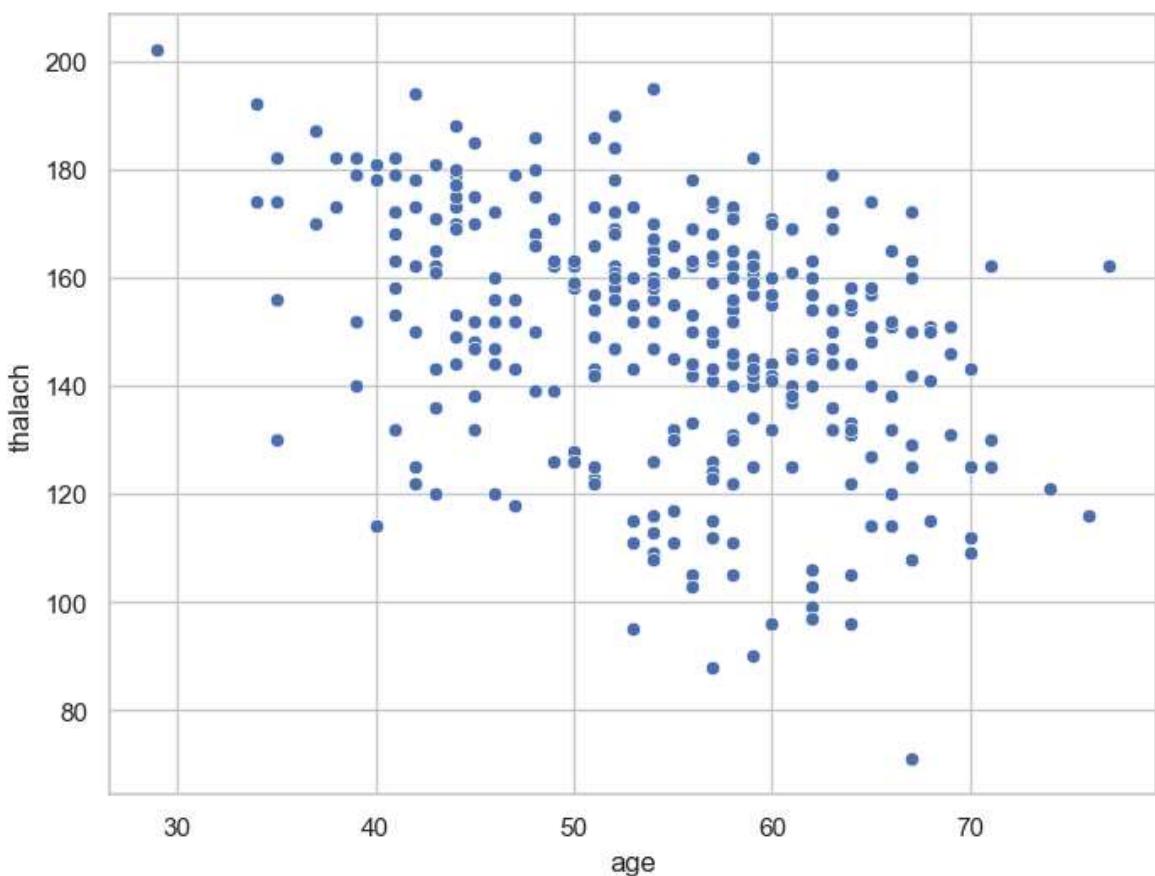


```
In [167]: f,ax=plt.subplots(figsize=(8,6))
ax=sns.regplot(x='age',y='chol',data=df)
```

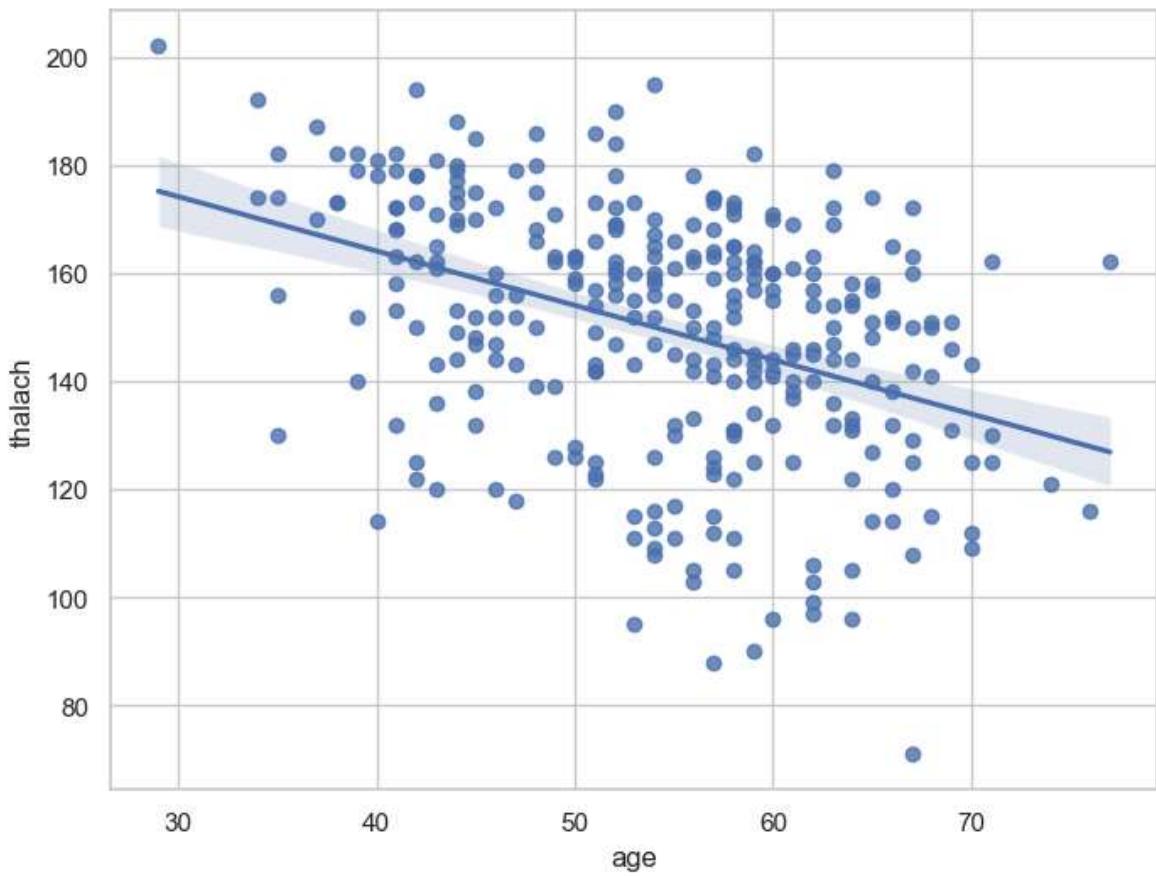
```
plt.show()
```



```
In [171]: f,ax=plt.subplots(figsize=(8,6))  
ax=sns.scatterplot(x='age',y='thalach',data=df)  
plt.show()
```



```
In [175...  
f,ax=plt.subplots(figsize=(8,6))  
ax=sns.regplot(x='age',y='thalach',data=df)  
plt.show()
```



```
In [177... df.isna().sum()
```

```
Out[177... age      0  
sex      0  
cp       0  
trestbps 0  
chol     0  
fbs      0  
restecg  0  
thalach  0  
exang    0  
oldpeak  0  
slope    0  
ca       0  
thal     0  
target   0  
dtype: int64
```

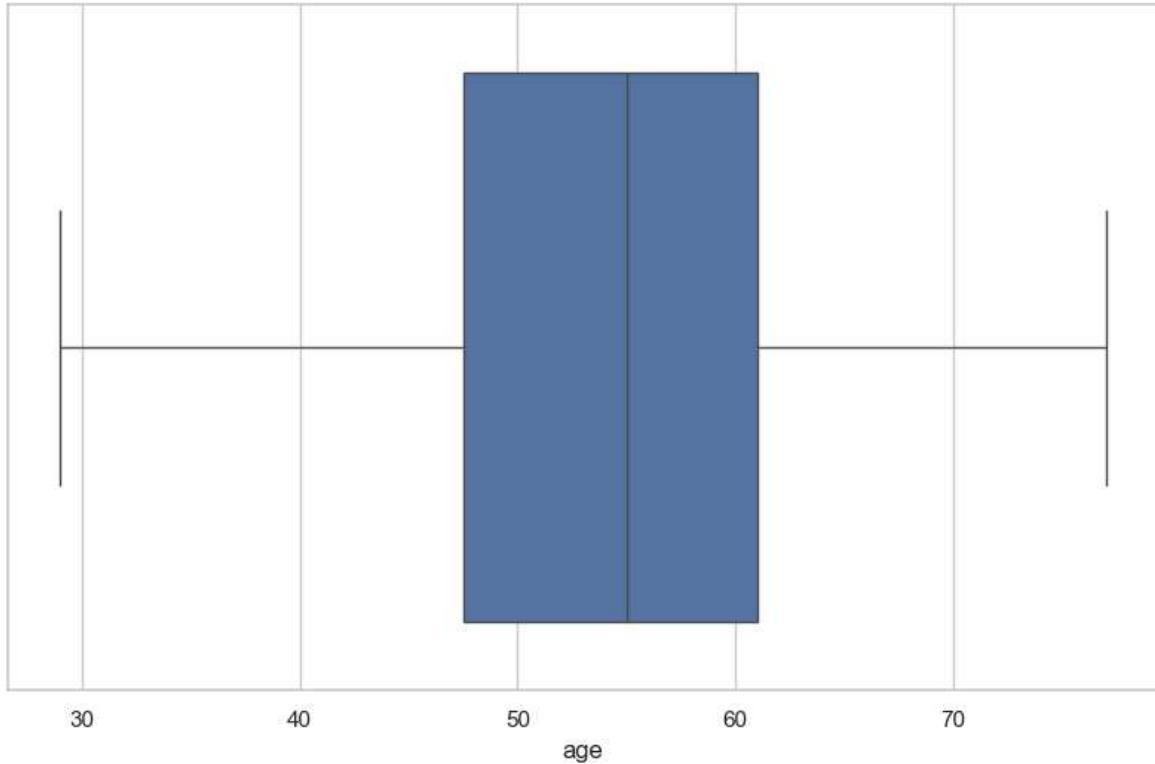
```
In [179... assert pd.notnull(df).all().all()
```

```
In [181... assert(df>=0).all().all()
```

```
In [183... df['age'].describe()
```

```
Out[183...]: 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
```

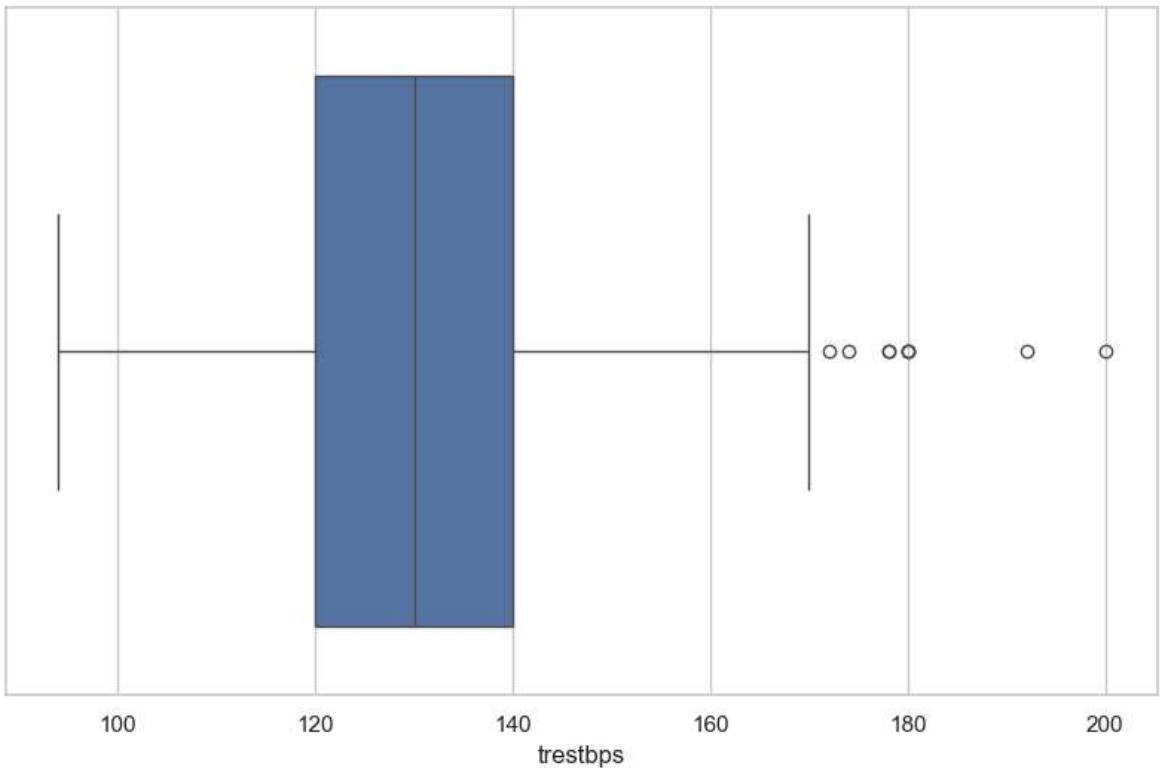
```
In [185...]: f,ax=plt.subplots(figsize=(10,6))
sns.boxplot(x=df['age'])
plt.show()
```



```
In [187...]: df['trestbps'].describe()
```

```
Out[187...]: 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
```

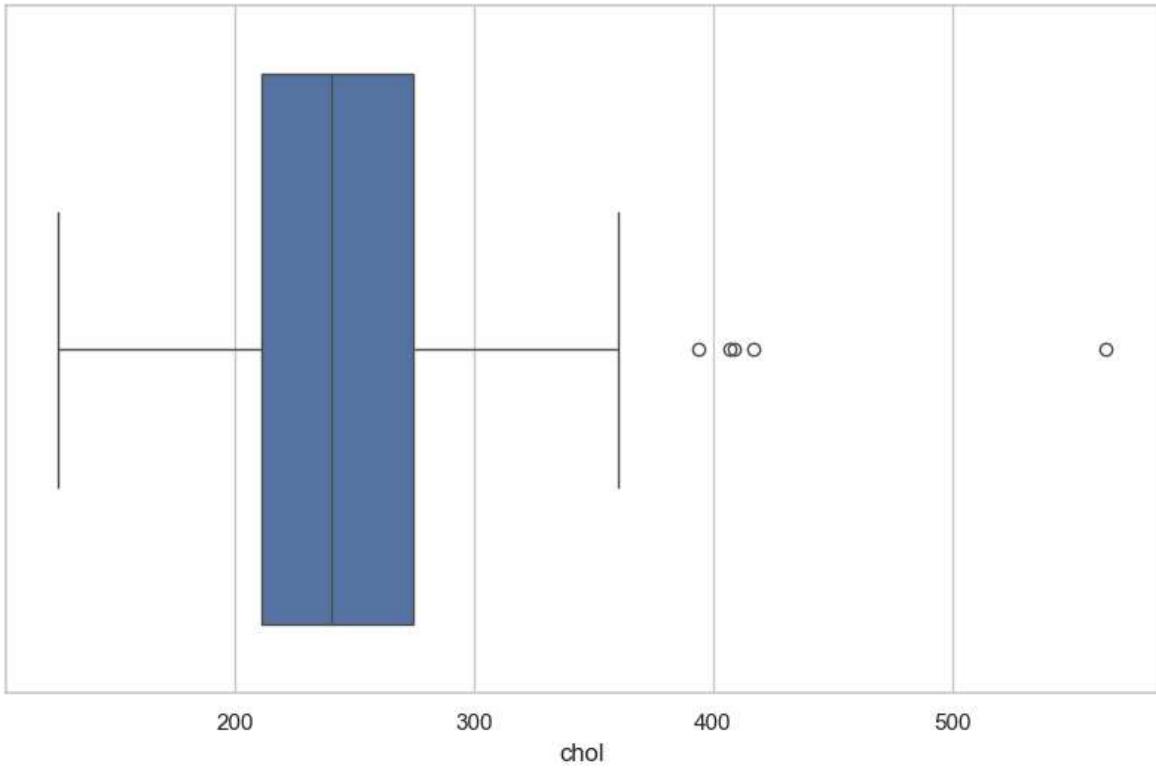
```
In [189...]: f,ax=plt.subplots(figsize=(10,6))
sns.boxplot(x=df['trestbps'])
plt.show()
```



```
In [193...]: df['chol'].describe()
```

```
Out[193...]: 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
```

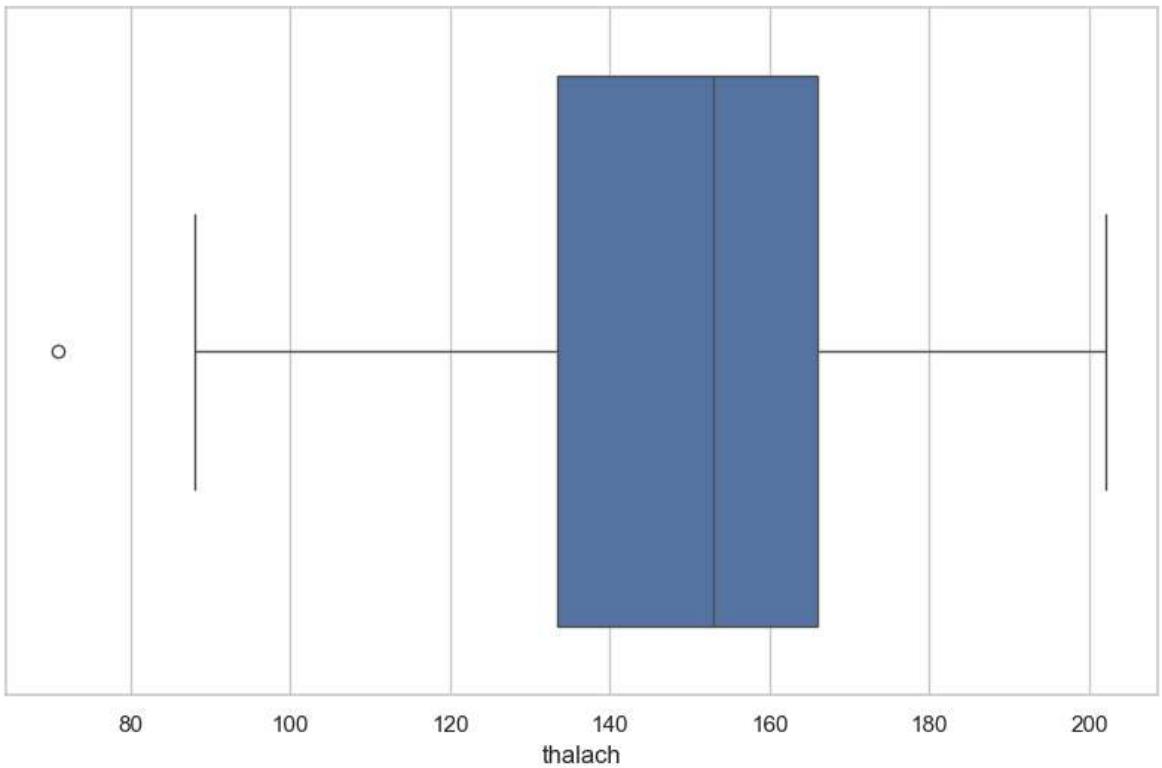
```
In [195...]: f, ax=plt.subplots(figsize=(10,6))
sns.boxplot(x=df['chol'])
plt.show()
```



```
In [197...]: df['thalach'].describe()
```

```
Out[197...]: 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
```

```
In [199...]: f, ax=plt.subplots(figsize=(10,6))
sns.boxplot(x=df['thalach'])
plt.show()
```

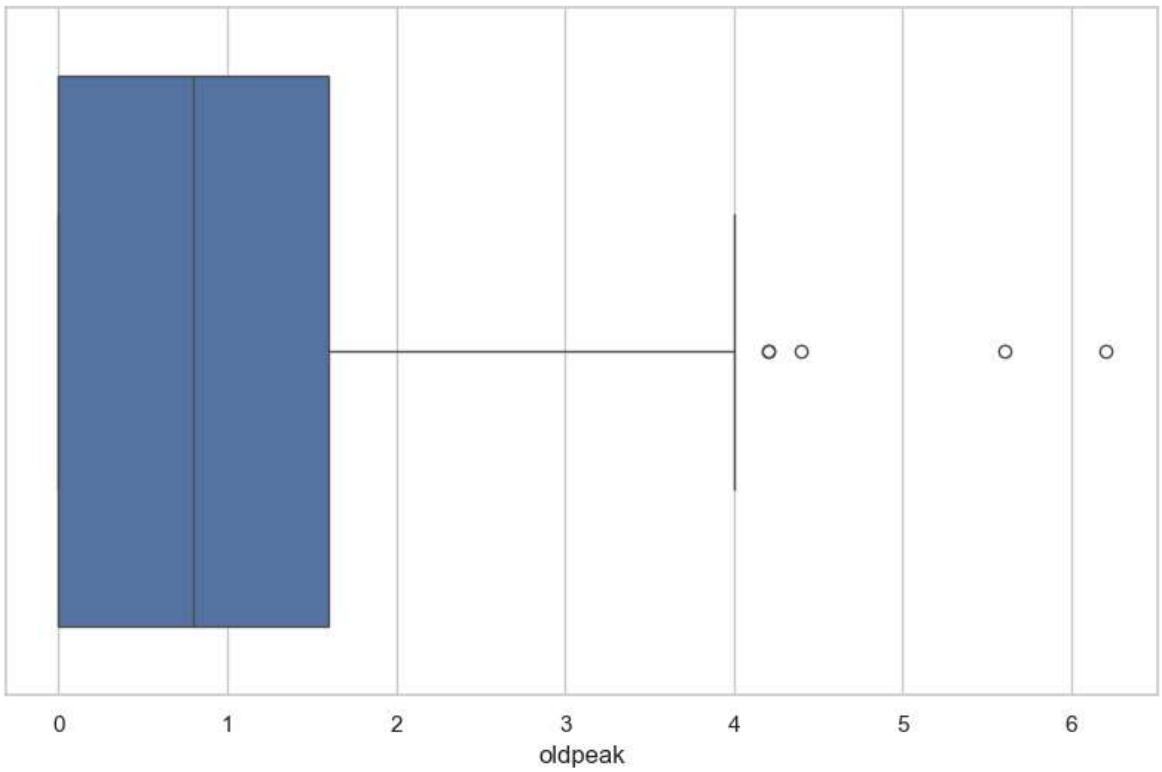


```
In [203...]: df['oldpeak'].describe()
```

```
Out[203...]:
```

	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

```
In [205...]: f, ax=plt.subplots(figsize=(10,6))
sns.boxplot(x=df['oldpeak'])
plt.show()
```



conclusion