

Heart disease or Cardiovascular disease (CVD)

* EDA

* Extensive Analysis + Visualization with Python

```
In [60]: from IPython.display import Image
```

```
# Display the image
img = Image(r"C:\Users\admin\Downloads\Heart.png") # Replace with your actual image
display(img)
```



```
In [3]: import numpy as np # Linear algebra
import pandas as pd # data processing, csv file I/O (e.g. pd.read_csv)
```

- We can see that the input folder contains one input file named heart.csv.

```
In [4]: import seaborn as sns
import matplotlib.pyplot as plt
import scipy.stats as st
%matplotlib inline

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

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

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

import dataset

- I will import the dataset with the usual pandas read_csv() function which is used to import CSV (Comma Separated Value) files.

```
In [7]: df = pd.read_csv(r"C:\DataScience Classnotes\9th- september - seaborn - Seaborn, Ed
df
```

	age	sex	cp	trestbps	chol	fb	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
...
298	57	0	0	140	241	0	1	123	1	0.2	1	0	3
299	45	1	3	110	264	0	1	132	0	1.2	1	0	3
300	68	1	0	144	193	1	1	141	0	3.4	1	2	3
301	57	1	0	130	131	0	1	115	1	1.2	1	1	3
302	57	0	1	130	236	0	0	174	0	0.0	1	1	2

303 rows × 14 columns

Exploratory Data Analysis

The scene has been set up. Now let the actual fun begin.

- Check shape of the dataset
-

It is a good idea to first check the shape of the dataset.

```
In [10]: # 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.

Preview the dataset

```
In [12]: # preview the dataset  
df.head()
```

Out[12]:

	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	0
2	41	0	1	130	204	0	0	172	0	1.4	2	0	2	0
3	56	1	1	120	236	0	1	178	0	0.8	2	0	2	0
4	57	0	0	120	354	0	1	163	1	0.6	2	0	2	0

#Summary of dataset

```
In [13]: 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
```

#Dataset description

- The dataset contains several columns which are as follows-
- age : age in years
- sex : 1 male & 0 female
- cp : chest pain type
- trestbps : resting blood pressure (in mm Hg on admission to the hospital)
- chol : serum cholestrol in mg/dl
- fbs : (fast blood sugar > 120 mg/dl)(1 =true; 0 =false)
- restecg : resting electrocardiographic results
- thalach : maximum heart rate achieved
- exang : exercise induced angina (1 = yes; 0 = no)
- oldpeak : ST depression induced by exercise relative to rest
- slope : the slope of the peak exercise ST segment
- ca : number of major vessels (0-3) colored by flourosopy
- thal : 3 = normal; 6 = fixed defect; 7 = reversable defect
- target : 1 or 0

#Checks the datatypes of columns

```
In [16]: df.dtypes
```

```
Out[16]: 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
```

#Important points about dataset

- `sex` is a character variable. Its data type should be object. But it is encoded as (1 = male; 0 = female). So, its data type is given as `int64`.
- Same is the case with several other variables - `fbs`, `exang` and `target`.
- `fbs` (fasting blood sugar) should be a character variable as it contains only 0 and 1 as values (1 = true; 0 = false). As it contains only 0 and 1 as values, so its data type is given as `int64`.
- `exang` (exercise induced angina) should also be a character variable as it contains only 0 and 1 as values (1 = yes; 0 = no). It also contains only 0 and 1 as values, so its data type is given as `int64`.
- `target` should also be a character variable. But, it also contains 0 and 1 as values. So, its data type is given as `int64`.

#Statistical properties of dataset

```
In [19]: # statistical properties of dataset
          df.describe()
```

Out[19]:

	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	0.528053
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860
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 [20]:

```
# view columns name
df.columns
```

```
Out[20]: Index(['age', 'sex', 'cp', 'trestbps', 'chol', 'fbs', 'restecg', 'thalach',
       'exang', 'oldpeak', 'slope', 'ca', 'thal', 'target'],
      dtype='object')
```

#Univariate Analysis

Analysis of target feature variables

- One feature variable is target.
- It refers to the presence of heart disease in the patient
- It is integer valued as it contains two integers 0 and 1 - (0 stands for absence of heart disease and 1 for presence of heart disease).
- So, in this section, I will analyze the target variable.

#Check the number of unique values in target variable

In [24]:

```
df['target'].nunique()
```

Out[24]:

We can see that there are 2 unique values in the target variable

In [26]:

```
#view the unique values in target variable
df['target'].unique()
```

Out[26]:

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

Comment

- So, the unique values are 1 and 0. (1 stands for presence of heart disease and 0 for absence of heart disease).

#Frequency distribution of target variable

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

```
Out[29]: target
1    165
0    138
Name: count, dtype: int64
```

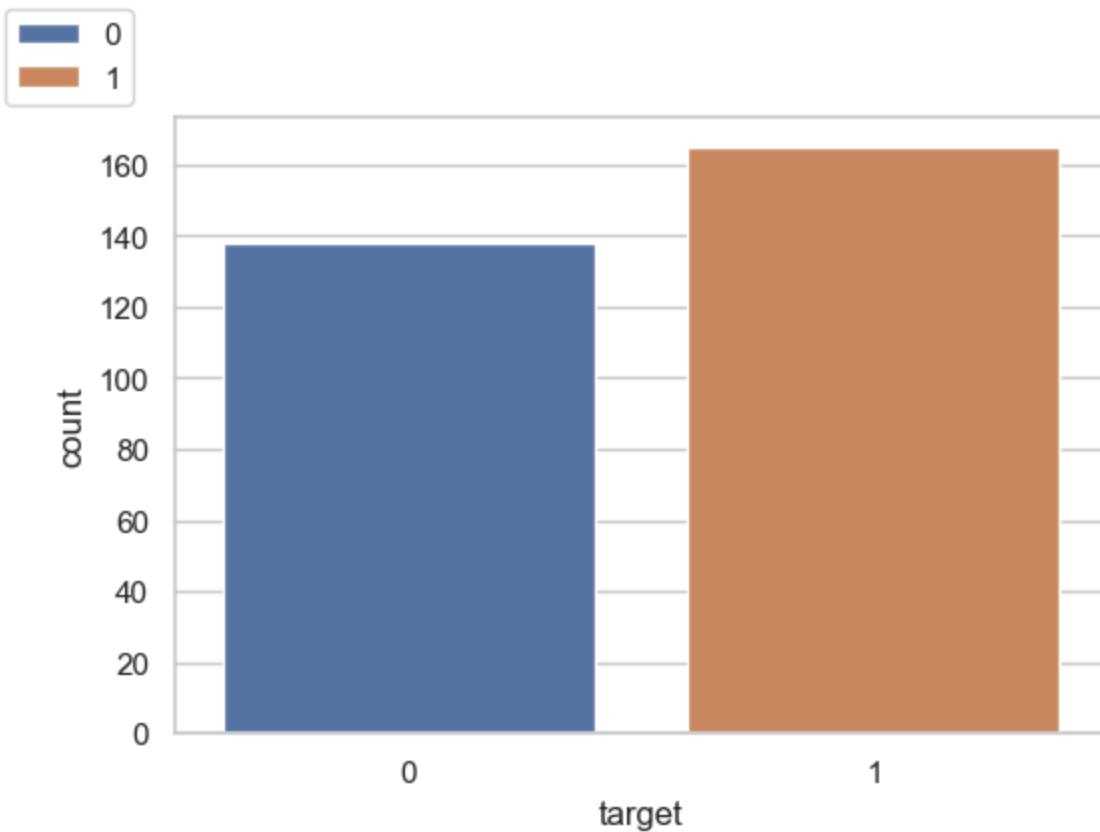
Comment

- 1 stands for presence of heart disease. So, there are 165 patients suffering from heart disease.
- Similarly, 0 stands for absence of heart disease. So, there are 138 patients who do not have any heart disease.
- We can visualize this information below.

Visualize frequency distribution of target variable

```
In [32]: f, ax = plt.subplots(figsize=(6,4))

ax = sns.countplot(x = 'target', data=df ,hue='target')
plt.legend(loc = 'upper left',bbox_to_anchor=(-0.2, 1.2) )
plt.show()
```



Interpretation

- The above plot confirms the findings that -
 - There are 165 patients suffering from heart disease, and
 - There are 138 patients who do not have any heart disease.

Frequency distribution of `target` variable wrt `sex`

```
In [35]: df.groupby('sex')['target'].value_counts()
```

```
Out[35]: sex  target
          0      1       72
                  0       24
          1      0      114
                  1       93
Name: count, dtype: int64
```

Comment

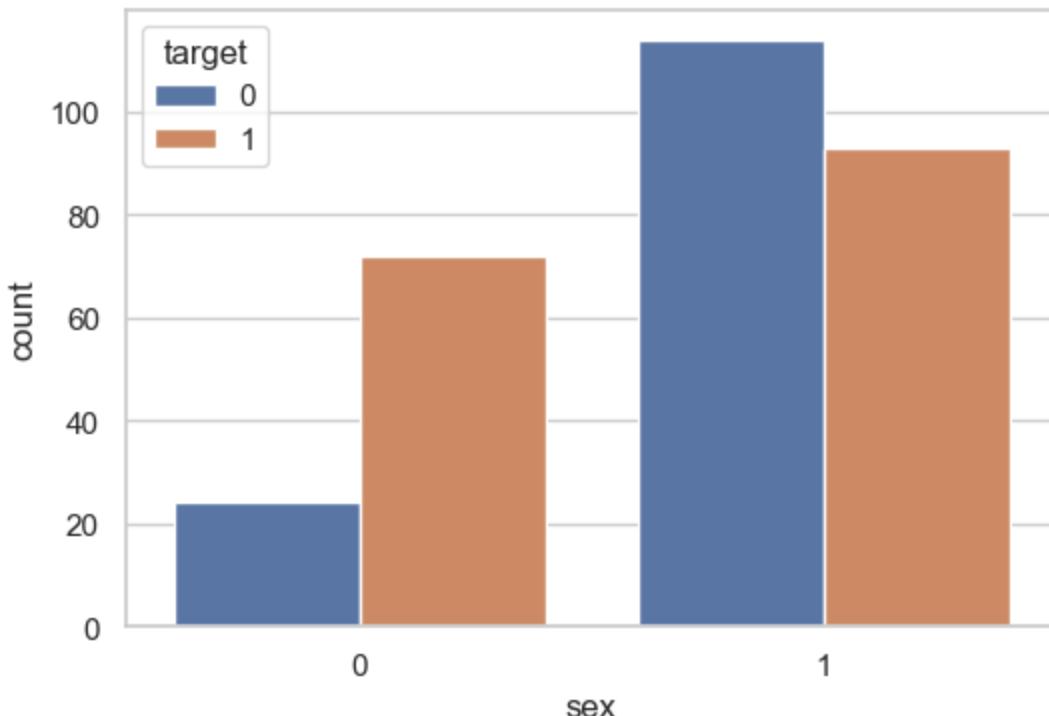
- `sex` variable contains two integer values 1 and 0 : (1 = male; 0 = female).
- `target` variable also contains two integer values 1 and 0 : (1 = Presence of heart disease; 0 = Absence of heart disease)

- So, out of 96 females - 72 have heart disease and 24 do not have heart disease.
- Similarly, out of 207 males - 93 have heart disease and 114 do not have heart disease.
- We can visualize this information below.

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

```
In [38]: #A countplot in Seaborn is a bar plot that shows the number of occurrences of each
#It helps visualize the distribution of categorical data.
f, ax=plt.subplots(figsize=(6,4))

ax = sns.countplot(x='sex', hue='target', data=df)
plt.show()
```

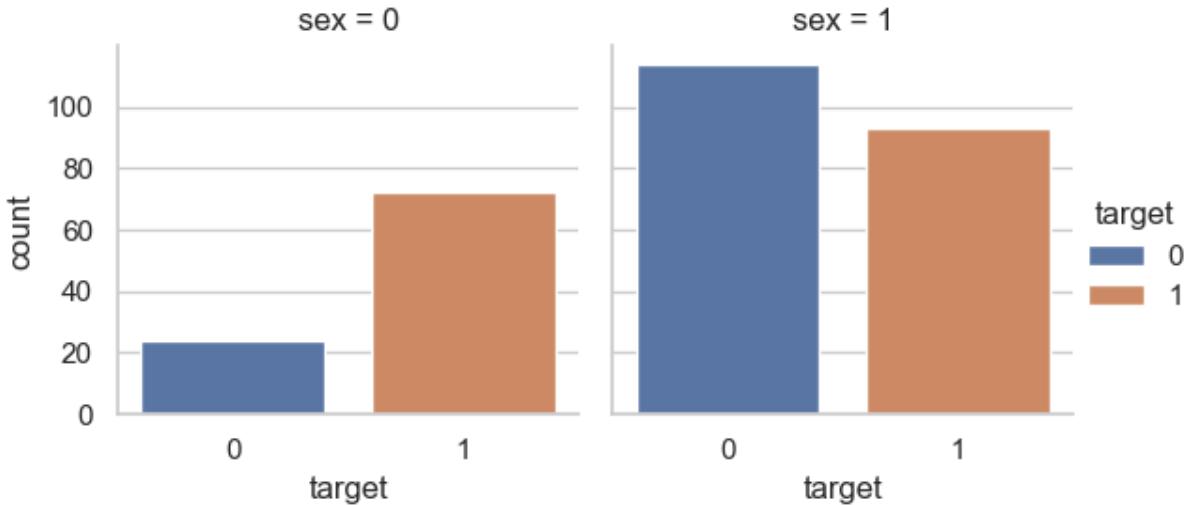


Interpretation

- We can see that the values of `target` variable are plotted wrt `sex` : (1 = male; 0 = female).
- `target` variable also contains two integer values 1 and 0 : (1 = Presence of heart disease; 0 = Absence of heart disease)
- The above plot confirms our findings that -
 - Out of 96 females - 72 have heart disease and 24 do not have heart disease.
 - Similarly, out of 207 males - 93 have heart disease and 114 do not have heart disease.

Alternatively, we can visualize the same information as follows :

```
In [41]: ax = sns.catplot(x='target', col='sex', data = df, kind='count', height=3, aspect=1, h
# col='sex': Creates separate subplots for each unique value in the 'sex' column.
# kind='count': Specifies that it's a count plot.
```

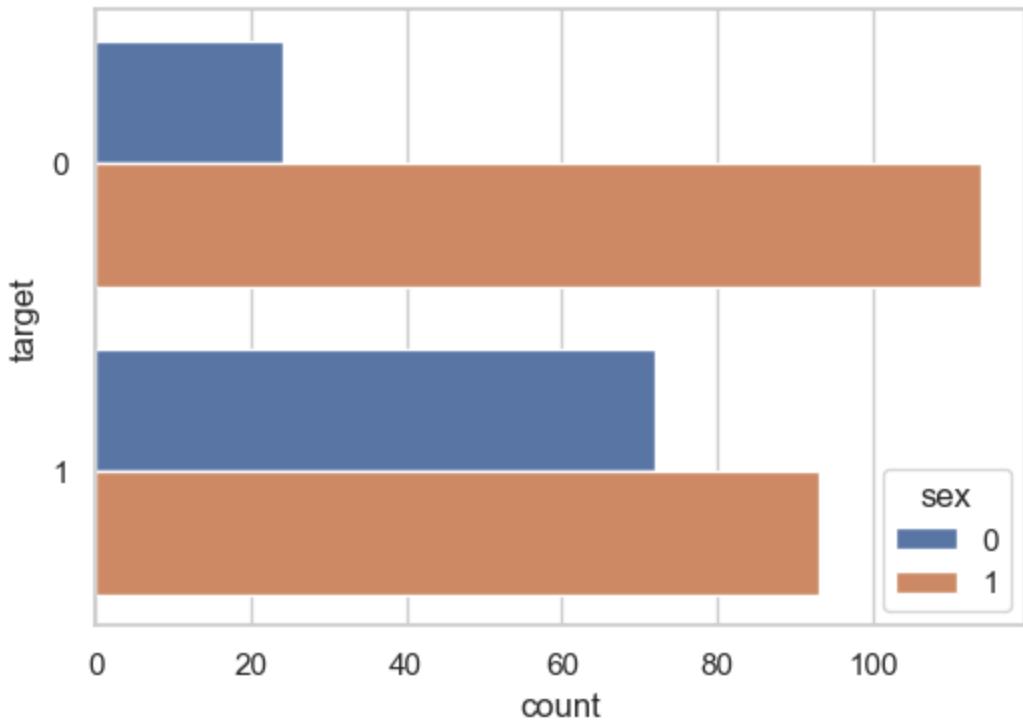


Comment

- The above plot segregate the values of `target` variable and plot on two different columns labelled as (sex = 0, sex = 1).
- I think it is more convinient way of interpret the plots.

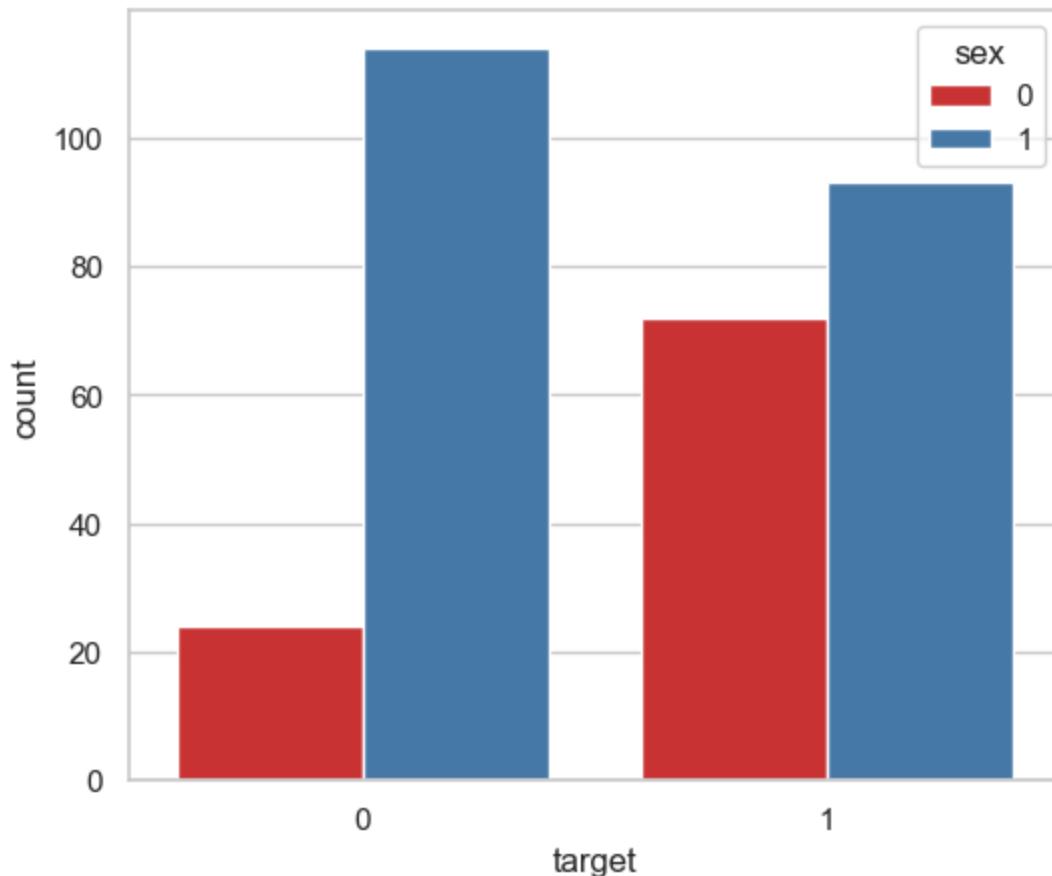
We can plot the bars horizontally as follows :

```
In [44]: f, ax= plt.subplots(figsize=(6,4))
ax = sns.countplot(y='target', hue='sex', data=df)
plt.show()
```



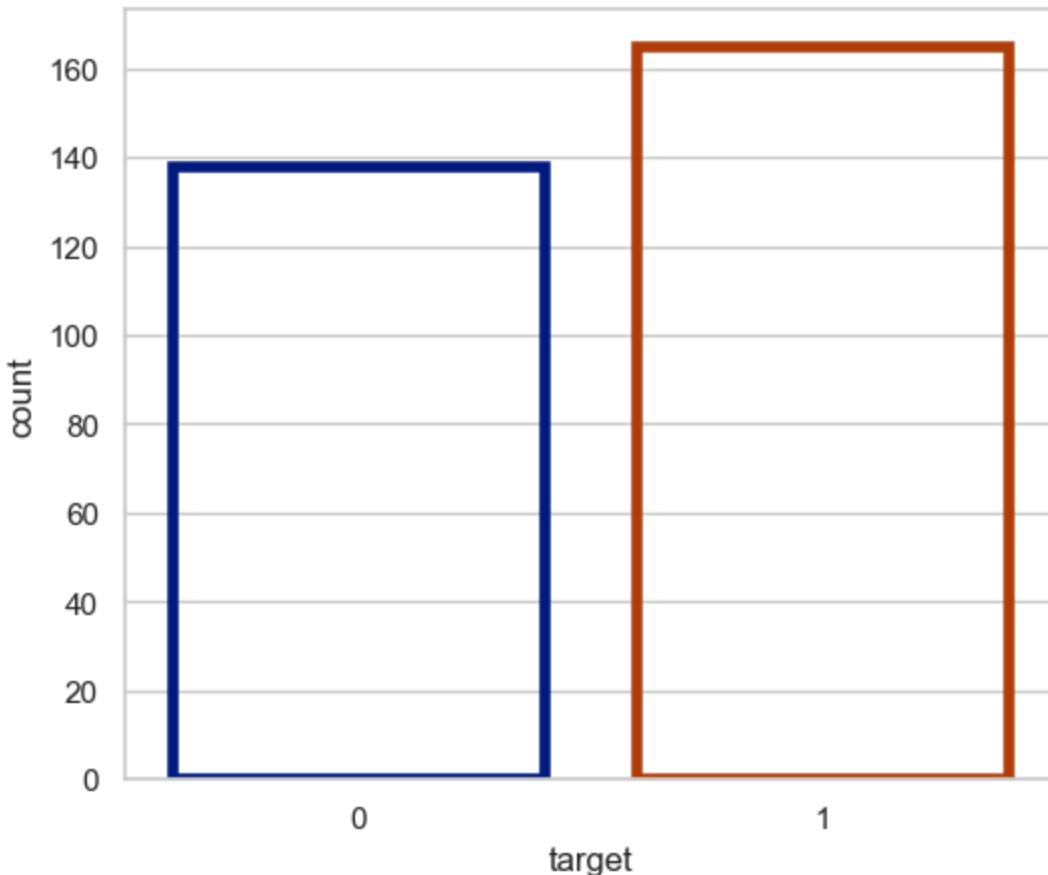
We can use a different color palette as follows :

```
In [46]: f,ax = plt.subplots(figsize=(6,5))
ax = sns.countplot(x='target', hue='sex', data=df, palette="Set1")
plt.show()
```



We can use `plt.bar` keyword arguments for a different look :

```
In [48]: f, ax = plt.subplots(figsize=(6,5))
ax = sns.countplot(x='target', data = df, facecolor=(0, 0, 0, 0), linewidth=4, edgecolor='black')
plt.show()
```



Comment

- I have visualize the `target` values distribution wrt `sex`.
- We can follow the same principles and visualize the `target` values distribution wrt `fbs` (fasting blood sugar) and `exang` (exercise induced angina).

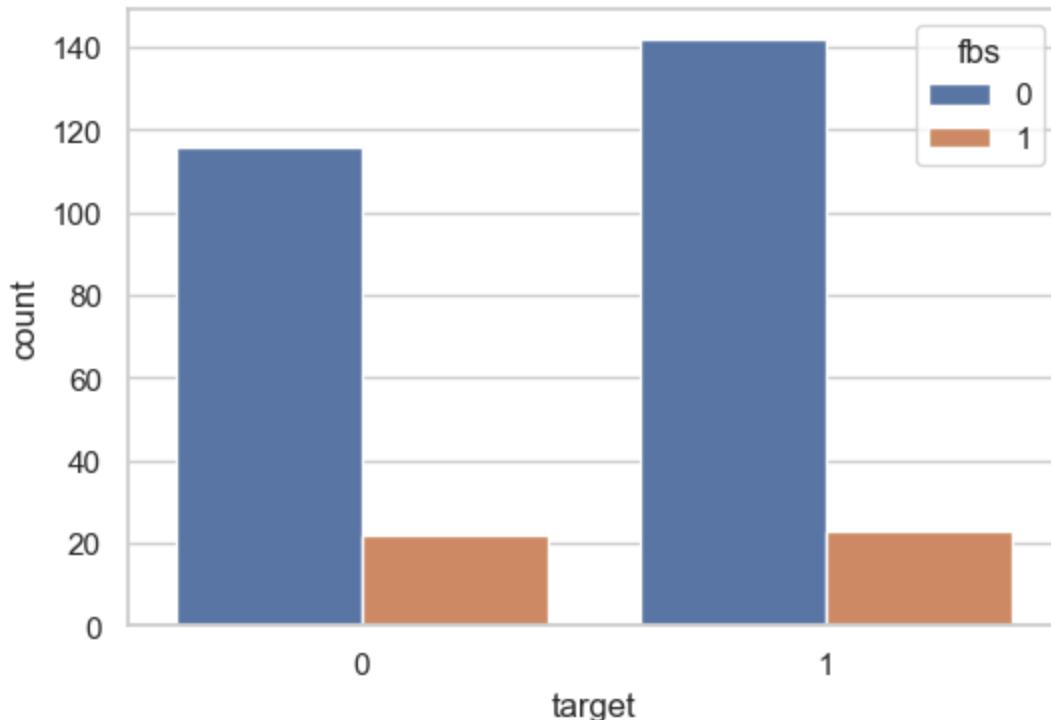
```
In [50]: df['target'].value_counts()
```

```
Out[50]: target
1    165
0    138
Name: count, dtype: int64
```

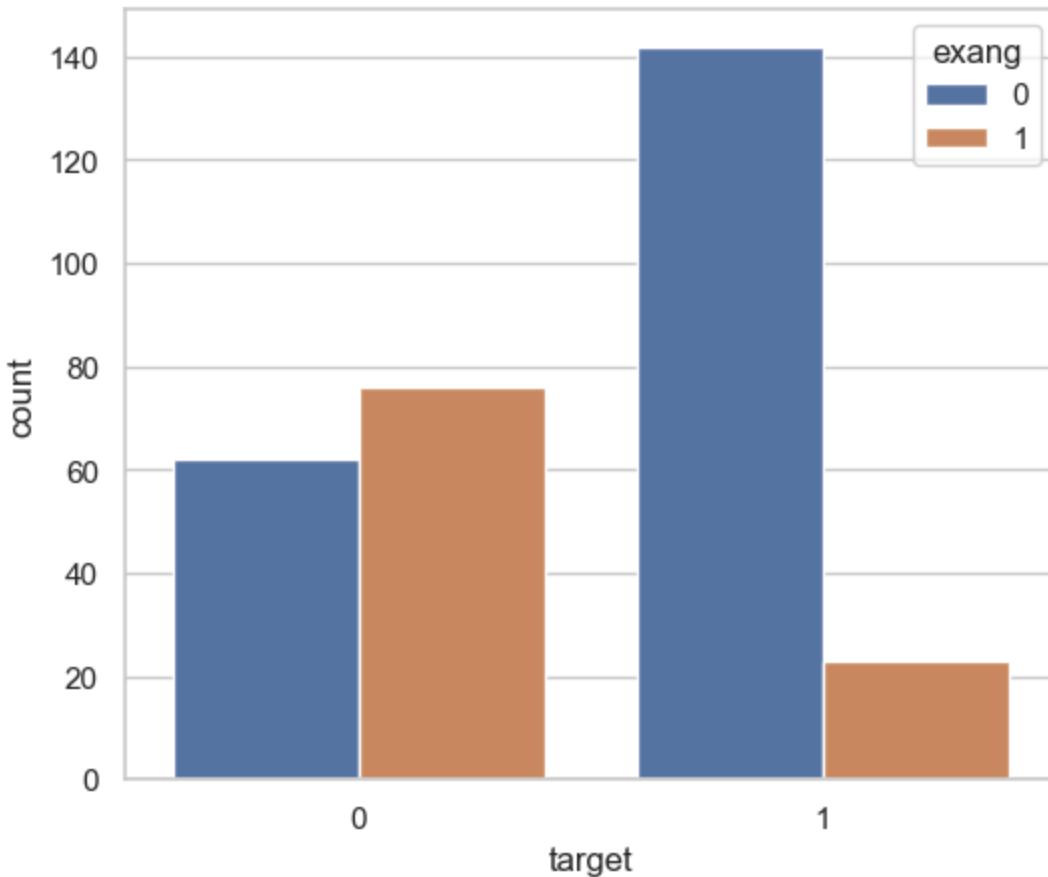
```
In [51]: df.groupby('sex')['target'].value_counts()
```

```
Out[51]: sex  target
          0      72
          0      24
          1     114
          1      93
Name: count, dtype: int64
```

```
In [52]: f,ax = plt.subplots(figsize=(6,4))
ax = sns.countplot(x='target',hue='fbs',data=df)
plt.show()
```



```
In [53]: f,ax = plt.subplots(figsize=(6,5))
ax = sns.countplot(x='target',hue='exang',data=df)
plt.show()
```



In []:

#Findings of Univariate Analysis

Findings of univariate analysis are as follows:-

- Our feature variable of interest is `target`.
- It refers to the presence of heart disease in the patient.
- It is integer valued as it contains two integers 0 and 1 - (0 stands for absence of heart disease and 1 for presence of heart disease).
- `1` stands for presence of heart disease. So, there are 165 patients suffering from heart disease.
- Similarly, `0` stands for absence of heart disease. So, there are 138 patients who do not have any heart disease.
- There are 165 patients suffering from heart disease, and
- There are 138 patients who do not have any heart disease.
- Out of 96 females - 72 have heart disease and 24 do not have heart disease.

- Similarly, out of 207 males - 93 have heart disease and 114 do not have heart disease.

#Bivariate Analysis

Estimate correlation coefficients

Our dataset is very small. So, I will compute the standard correlation coefficient (also called Pearson's r) between every pair of attributes. I will compute it using the `df.corr()` method as follows:-

```
In [57]: correlation = df.corr()
```

The target variable is `target`. So, we should check how each attribute correlates with the `target` variable. We can do it as follows:-

```
In [59]: correlation['target'].sort_values(ascending=False)
```

```
Out[59]: 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
```

Interpretation of correlation coefficient

- The correlation coefficient ranges from -1 to +1.
- When it is close to +1, this signifies that there is a strong positive correlation. So, we can see that there is no variable which has strong positive correlation with `target` variable.
- When it is close to -1, it means that there is a strong negative correlation. So, we can see that there is no variable which has strong negative correlation with `target` variable.
- When it is close to 0, it means that there is no correlation. So, there is no correlation between `target` and `fbs`.

- We can see that the `cp` and `thalach` variables are mildly positively correlated with `target` variable. So, I will analyze the interaction between these features and `target` variable.

Analysis of `target` and `cp` variable

Explore `cp` variable

- `cp` stands for chest pain type.
- First, I will check number of unique values in `cp` variable.

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

```
Out[63]: 4
```

So, there are 4 unique values in `cp` variable. Hence, it is a categorical variable.

Now, I will view its frequency distribution as follows :

```
In [66]: df['cp'].unique()
```

```
Out[66]: array([3, 2, 1, 0], dtype=int64)
```

```
In [67]: df['cp'].value_counts()
```

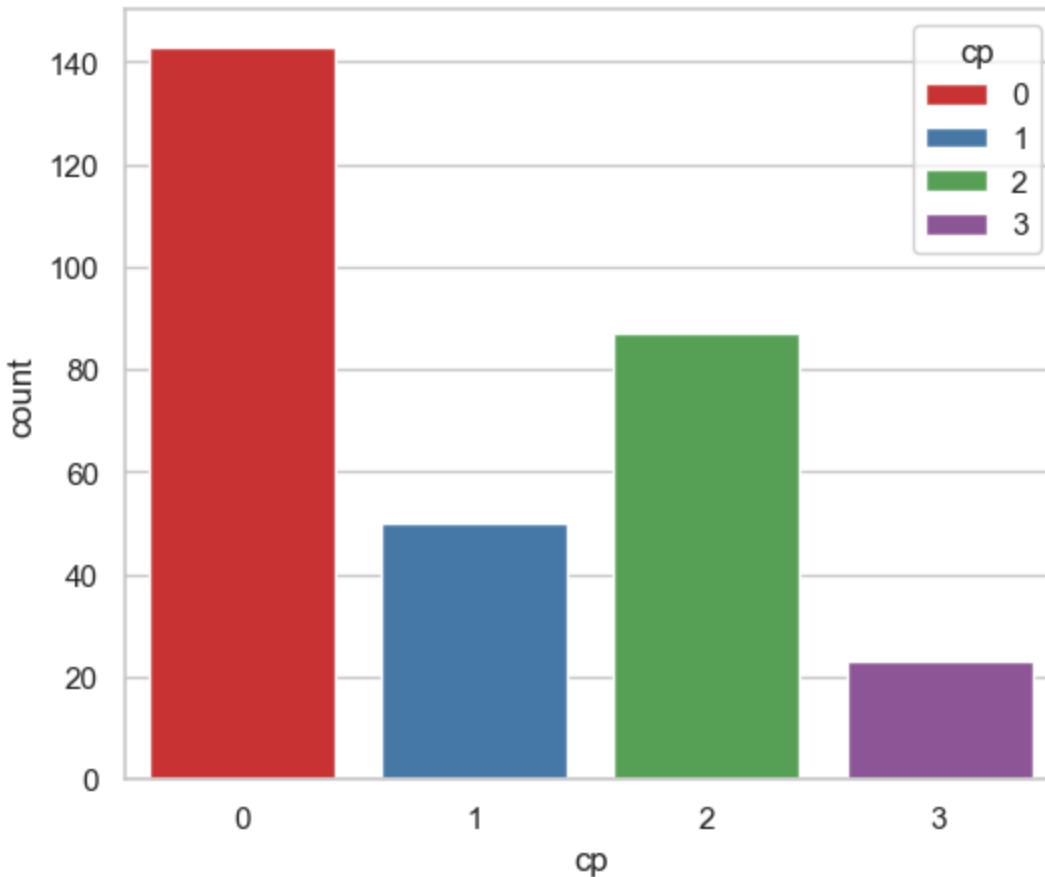
```
Out[67]: cp
0    143
2     87
1     50
3     23
Name: count, dtype: int64
```

Comment

- It can be seen that `cp` is a categorical variable and it contains 4 types of values - 0, 1, 2 and 3.

#Visualize the frequency distribution of `cp` variable

```
In [70]: f, ax= plt.subplots(figsize=(6,5))
ax = sns.countplot(x='cp', data=df, hue='cp', palette="Set1")
plt.show()
```



Frequency distribution of `target` variable wrt `cp`

```
In [72]: df.groupby('cp')['target'].value_counts()
```

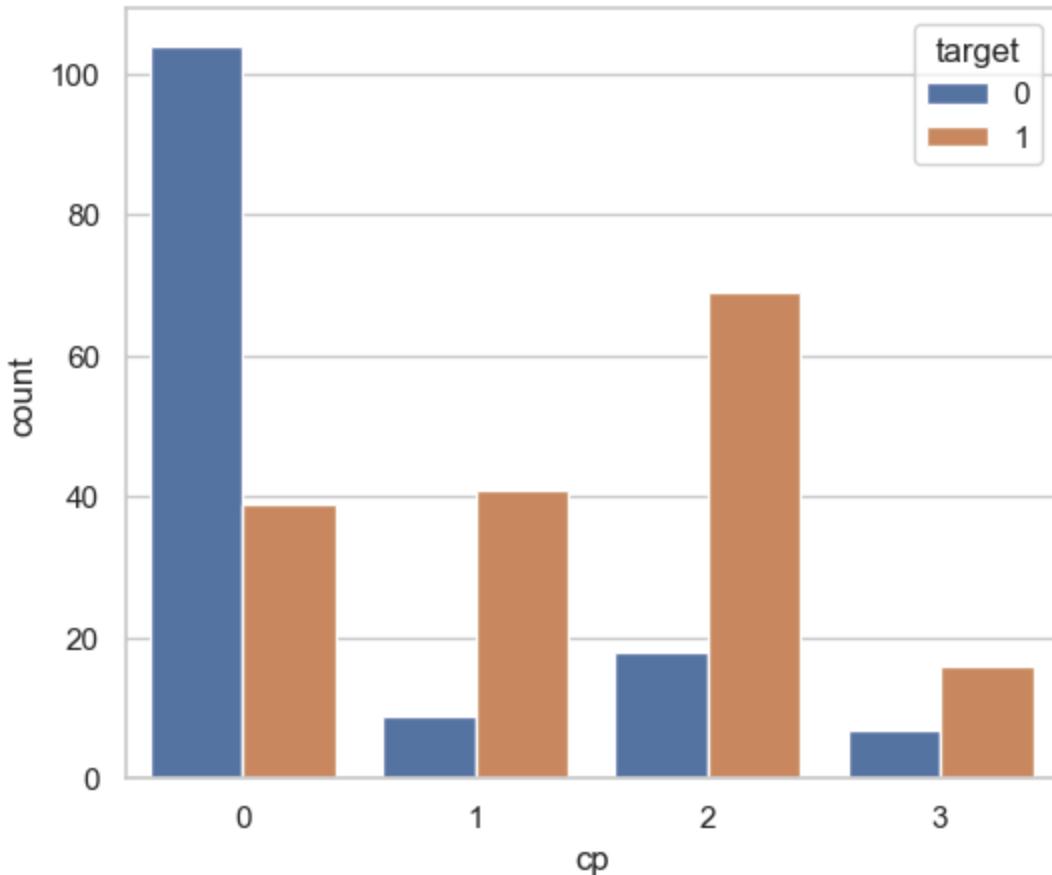
```
Out[72]: cp  target
0   0      104
     1      39
1   1      41
     0      9
2   1      69
     0      18
3   1      16
     0      7
Name: count, dtype: int64
```

Comment

- `cp` variable contains four integer values 0, 1, 2 and 3.
- `target` variable contains two integer values 1 and 0 : (1 = Presence of heart disease; 0 = Absence of heart disease)
- So, the above analysis gives `target` variable values categorized into presence and absence of heart disease and groupby `cp` variable values.
- We can visualize this information below.

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

```
In [75]: f,ax= plt.subplots(figsize=(6,5))
ax = sns.countplot(x='cp',hue='target',data=df)
plt.show()
```

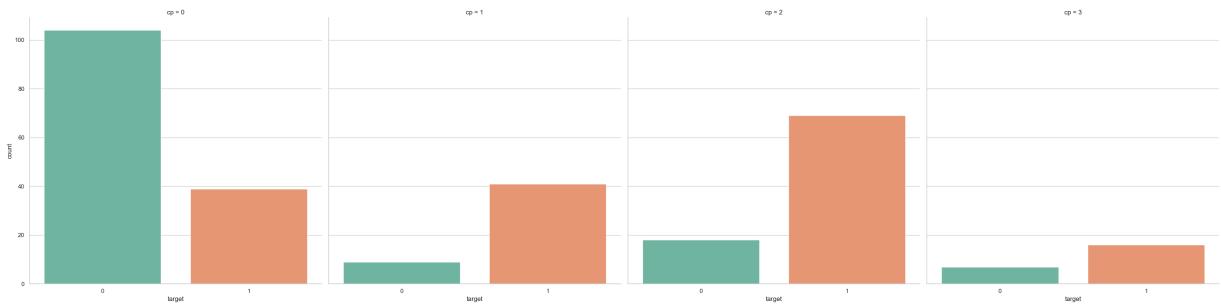


Interpretation

- We can see that the values of `target` variable are plotted wrt `cp`.
- `target` variable contains two integer values 1 and 0 : (1 = Presence of heart disease; 0 = Absence of heart disease)
- The above plot confirms our above findings,

Alternatively, we can visualize the same information as follows :

```
In [78]: ax = sns.catplot(x='target', col='cp', data=df, kind= "count",height=8,aspect=1,pal
```



Analysis of target and thalach variable

Explore thalach variable

- `thalach` stands for maximum heart rate achieved.
- I will check number of unique values in `thalach` variable as follows :

```
In [81]: df['thalach'].nunique()
```

```
Out[81]: 91
```

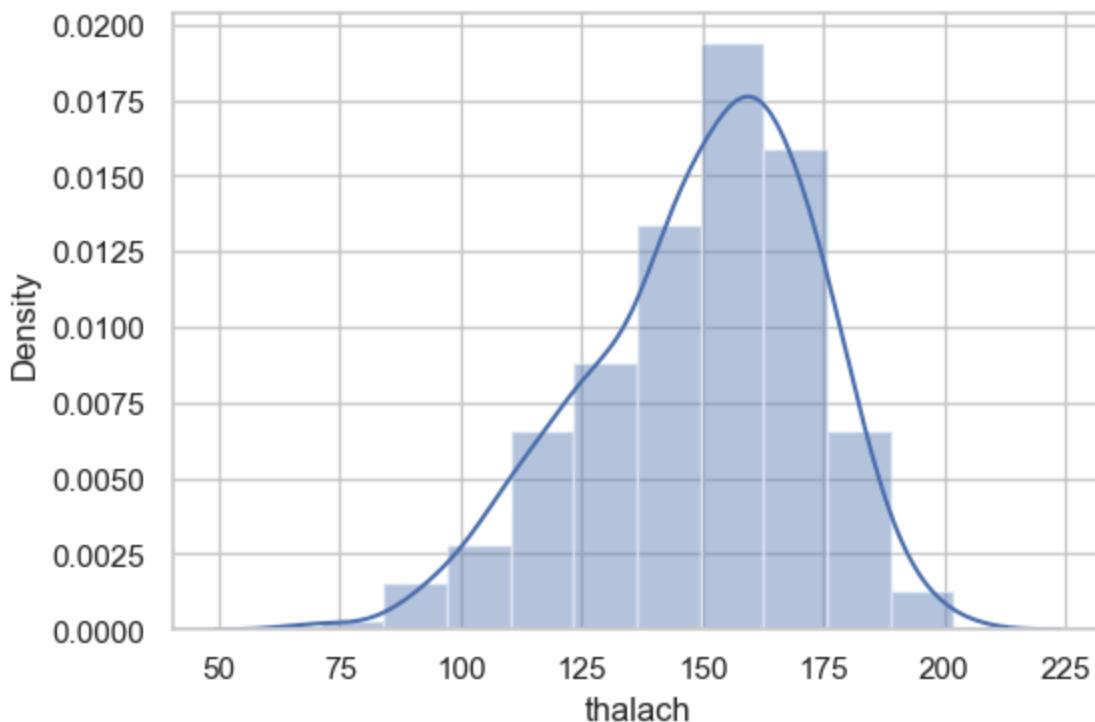
```
In [82]: df['thalach'].unique()
```

```
Out[82]: array([150, 187, 172, 178, 163, 148, 153, 173, 162, 174, 160, 139, 171,
       144, 158, 114, 151, 161, 179, 137, 157, 123, 152, 168, 140, 188,
       125, 170, 165, 142, 180, 143, 182, 156, 115, 149, 146, 175, 186,
       185, 159, 130, 190, 132, 147, 154, 202, 166, 164, 184, 122, 169,
       138, 111, 145, 194, 131, 133, 155, 167, 192, 121, 96, 126, 105,
       181, 116, 108, 129, 120, 112, 128, 109, 113, 99, 177, 141, 136,
       97, 127, 103, 124, 88, 195, 106, 95, 117, 71, 118, 134, 90],
      dtype=int64)
```

- So, number of unique values in `thalach` variable is 91. Hence, it is numerical variable.
- I will visualize its frequency distribution of values as follows :

Visualize the frequency distribution of thalach variable

```
In [85]: f, ax = plt.subplots(figsize=(6,4))
x = df['thalach']
ax = sns.distplot(x, bins=10)
plt.show()
```

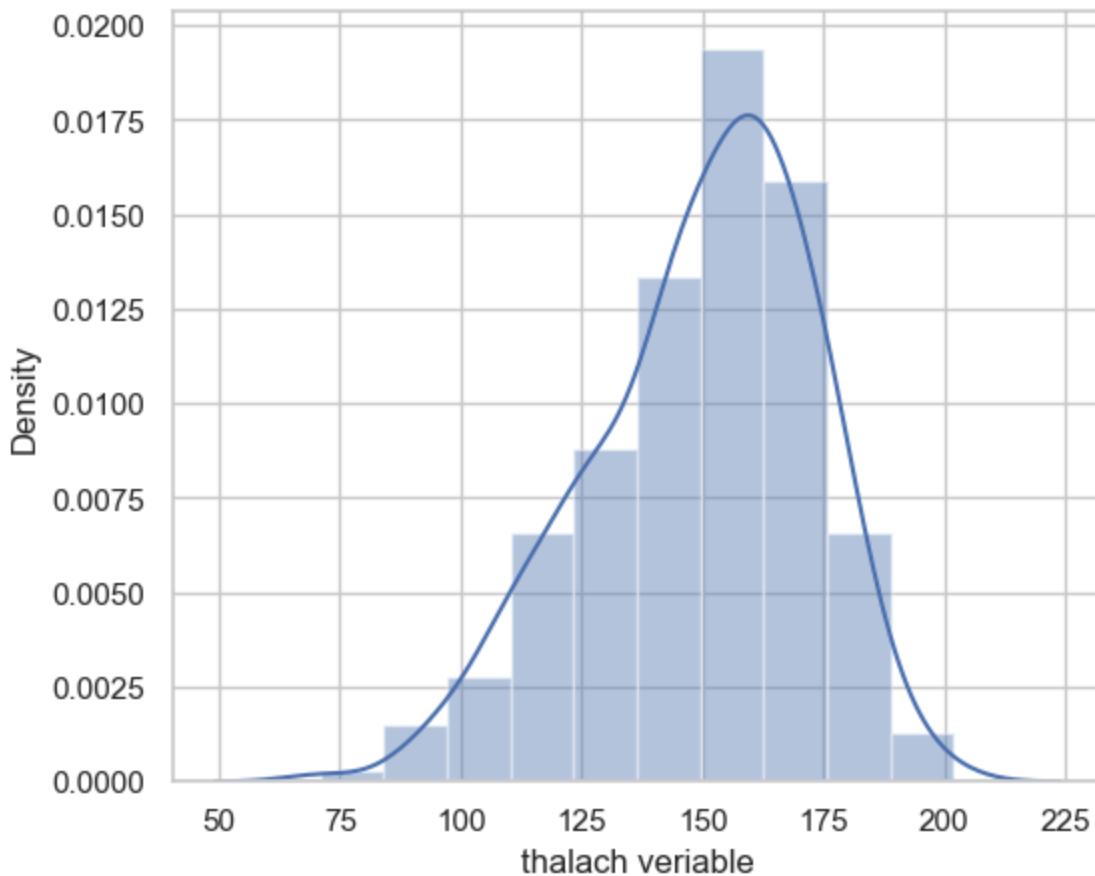


Comment

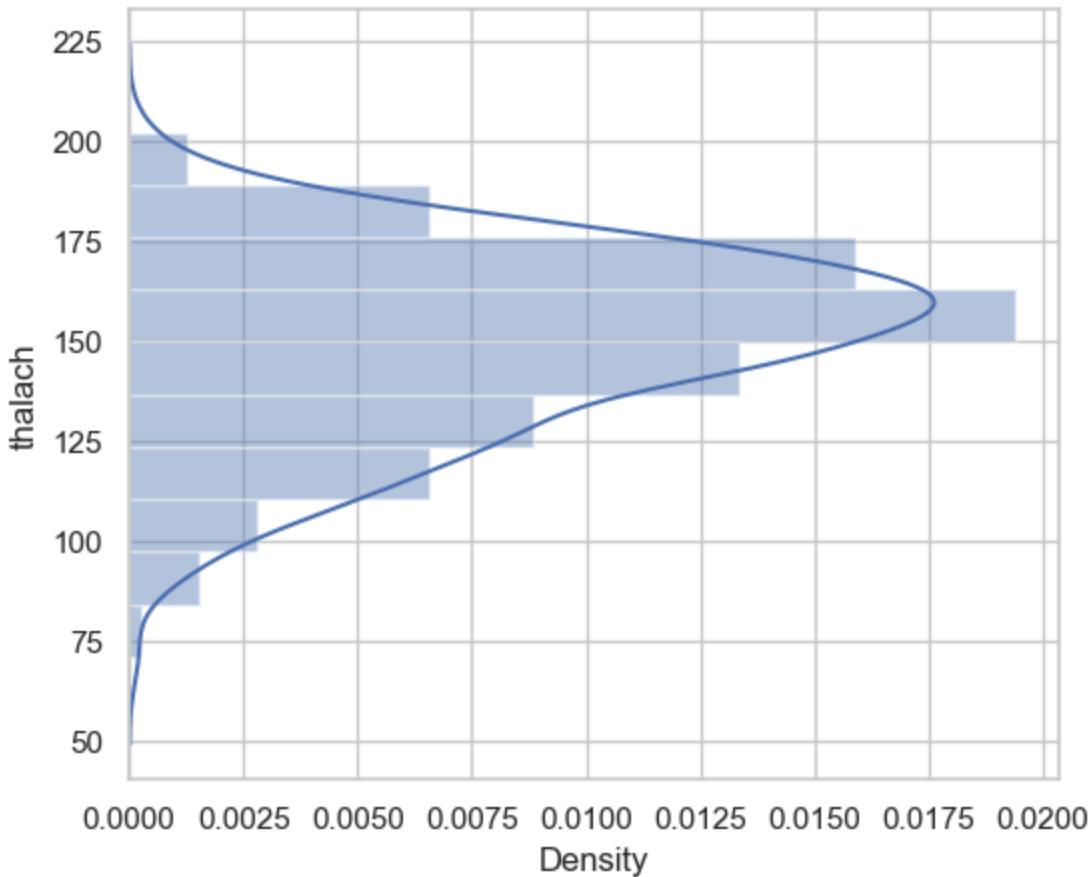
- We can see that the `thalach` variable is slightly negatively skewed.

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

```
In [88]: f ,ax = plt.subplots(figsize=(6,5))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.distplot(x,bins=10)
plt.show()
```



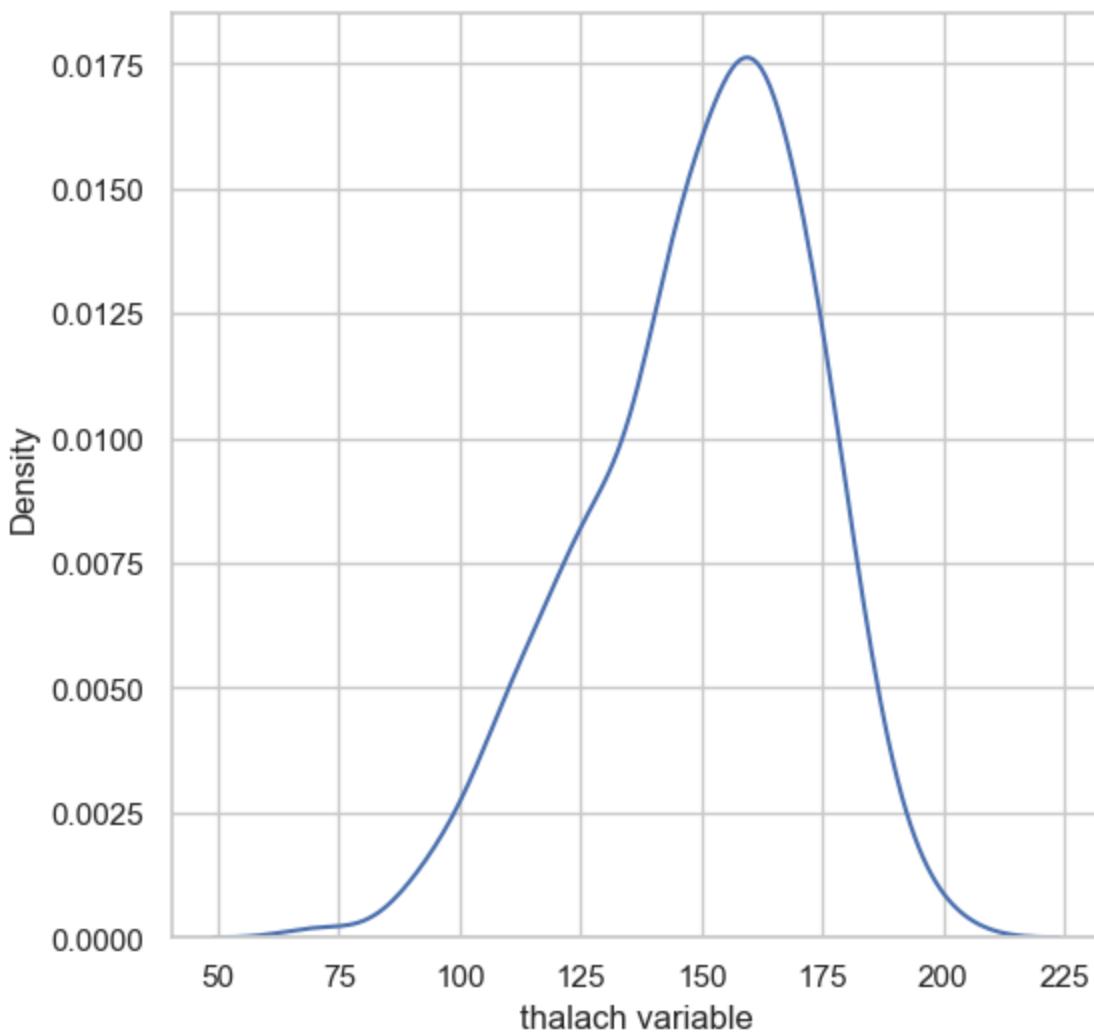
```
In [89]: f, ax = plt.subplots(figsize=(6,5))
x = df['thalach']
ax = sns.distplot(x, bins=10, vertical=True)
plt.show()
```



Seaborn Kernel Density Estimation (KDE) Plot

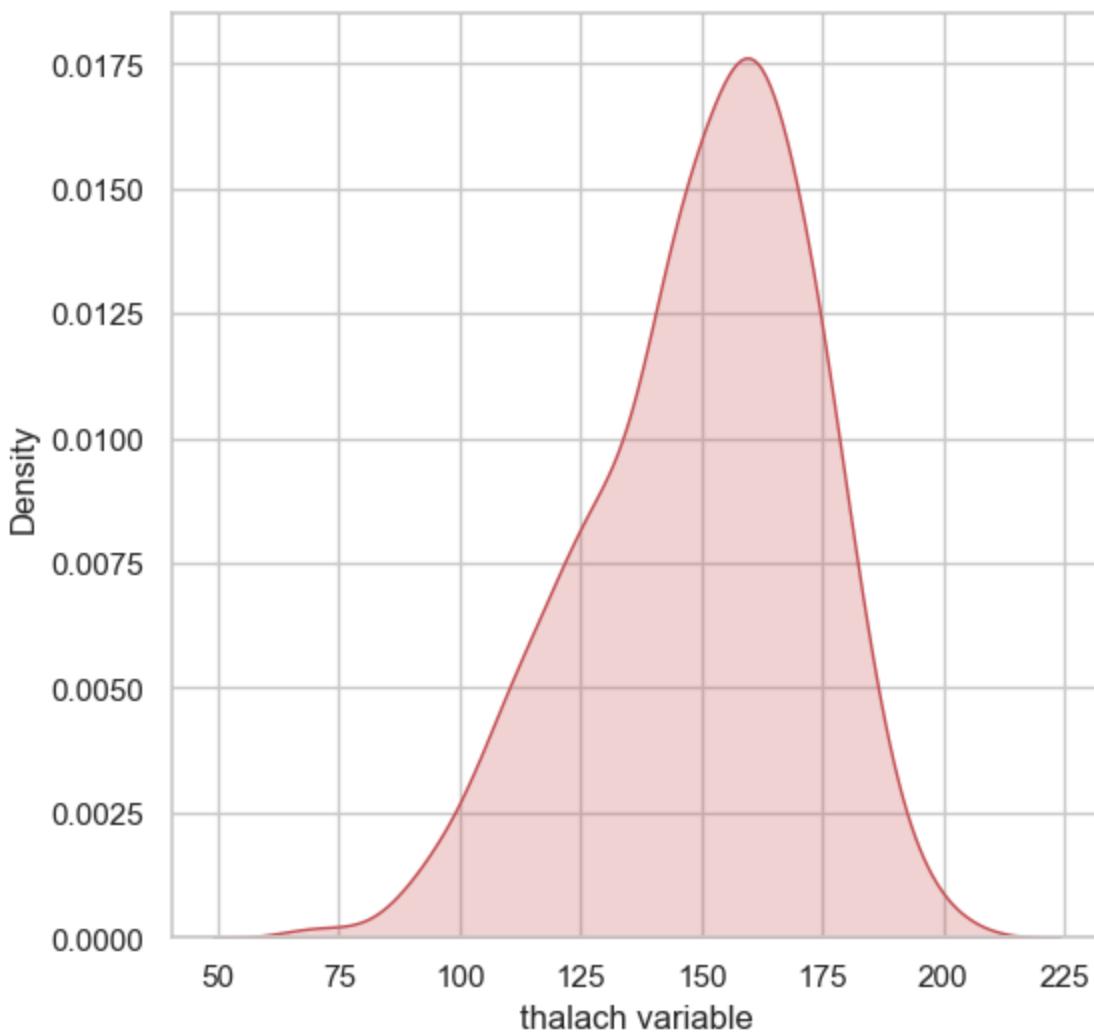
- The kernel density estimate (KDE) plot is a useful tool for plotting the shape of a distribution.
- The KDE plot plots the density of observations on one axis with height along the other axis.
- We can plot a KDE plot as follows :

```
In [91]: f, ax = plt.subplots(figsize=(6,6))
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 [93]: f, ax = plt.subplots(figsize=(6,6))
x = df['thalach']
x = pd.Series(x, name="thalach variable")
ax = sns.kdeplot(x, shade=True, color='r')
plt.show()
```



Histogram

- A histograms represent the distribution of data by forming bins along the range of the data and then drawing bars to show the number of observations that fall in each bin.
- We can plot histogram as follows:
- thalach : maximum heart rate achieved

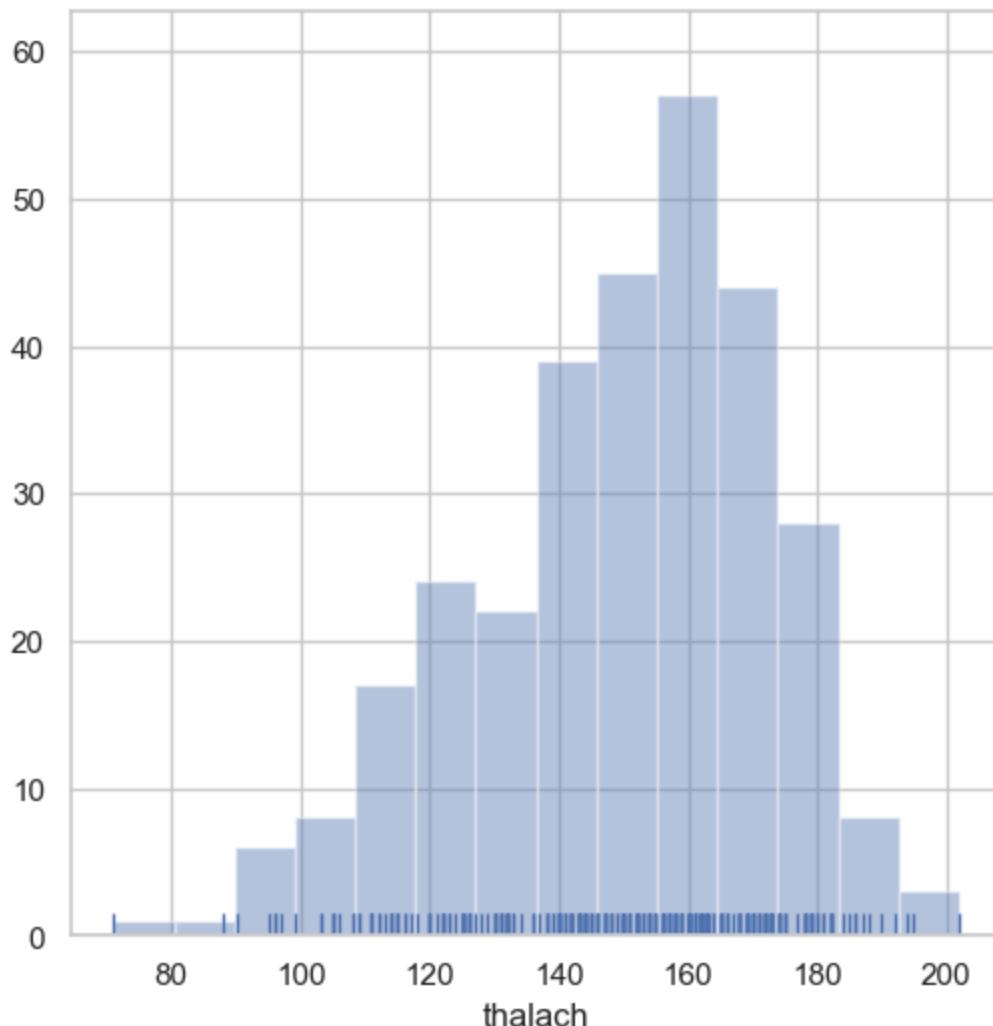
```
In [95]: df['thalach'].unique()
```

```
Out[95]: array([150, 187, 172, 178, 163, 148, 153, 173, 162, 174, 160, 139, 171,
       144, 158, 114, 151, 161, 179, 137, 157, 123, 152, 168, 140, 188,
       125, 170, 165, 142, 180, 143, 182, 156, 115, 149, 146, 175, 186,
       185, 159, 130, 190, 132, 147, 154, 202, 166, 164, 184, 122, 169,
       138, 111, 145, 194, 131, 133, 155, 167, 192, 121, 96, 126, 105,
       181, 116, 108, 129, 120, 112, 128, 109, 113, 99, 177, 141, 136,
       97, 127, 103, 124, 88, 195, 106, 95, 117, 71, 118, 134, 90],
      dtype=int64)
```

```
In [96]: df['thalach'].nunique()
```

```
Out[96]: 91
```

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



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

```
In [100...]: df.groupby('thalach')['target'].value_counts().head(50)
```

Out[100...]:

	thalach	target
71	0	1
88	0	1
90	0	1
95	0	1
96	0	1
	1	1
97	0	1
99	0	1
103	0	2
105	0	2
	1	1
106	0	1
108	0	2
109	0	2
111	0	2
	1	1
112	0	2
113	0	1
114	0	2
	1	1
115	1	2
	0	1
116	0	1
	1	1
117	0	1
118	0	1
120	0	3
121	1	1
122	1	3
	0	1
123	0	1
	1	1
124	0	1
125	0	5
	1	2
126	0	3
	1	1
127	0	1
128	0	1
129	0	1
130	0	3
	1	1
131	0	2
	1	2
132	0	6
	1	1
133	0	1
	1	1
134	0	1
136	0	2

Name: count, dtype: int64

- 71 (thalach) and 0 (target): There is 1 person with a heart rate of 71 and no heart disease.

- 96 (thalach):

target = 0: There is 1 person with no heart disease. target = 1: There is 1 person with heart disease.

- 103 (thalach):

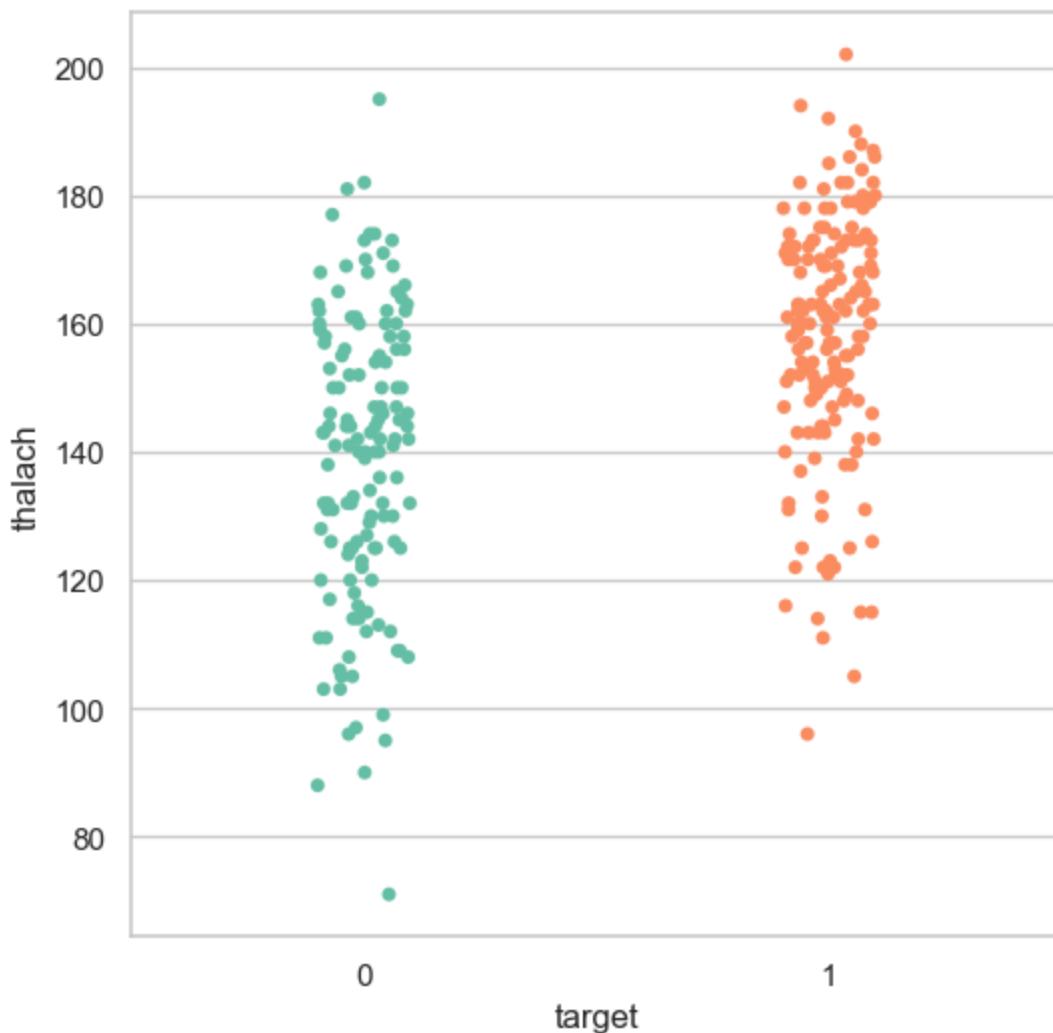
target = 0: There are 2 people with a heart rate of 103 and no heart disease.

- 105 (thalach):

target = 0: 2 people with no heart disease. target = 1: 1 person with heart disease.

In [102...]

```
f,ax = plt.subplots(figsize=(6,6))
sns.stripplot(x='target',y='thalach',data=df,palette='Set2')
plt.show()
```

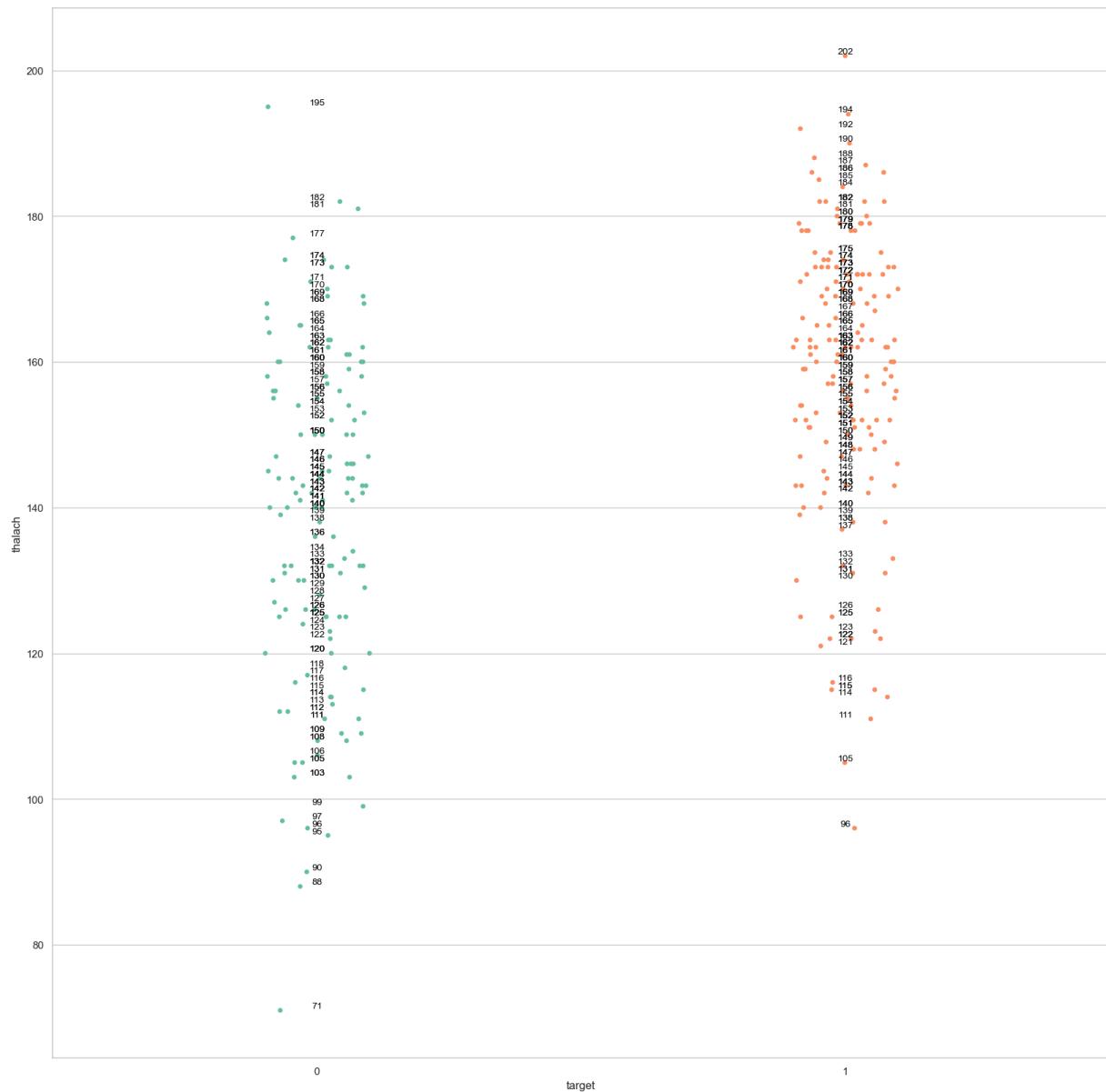


In []:

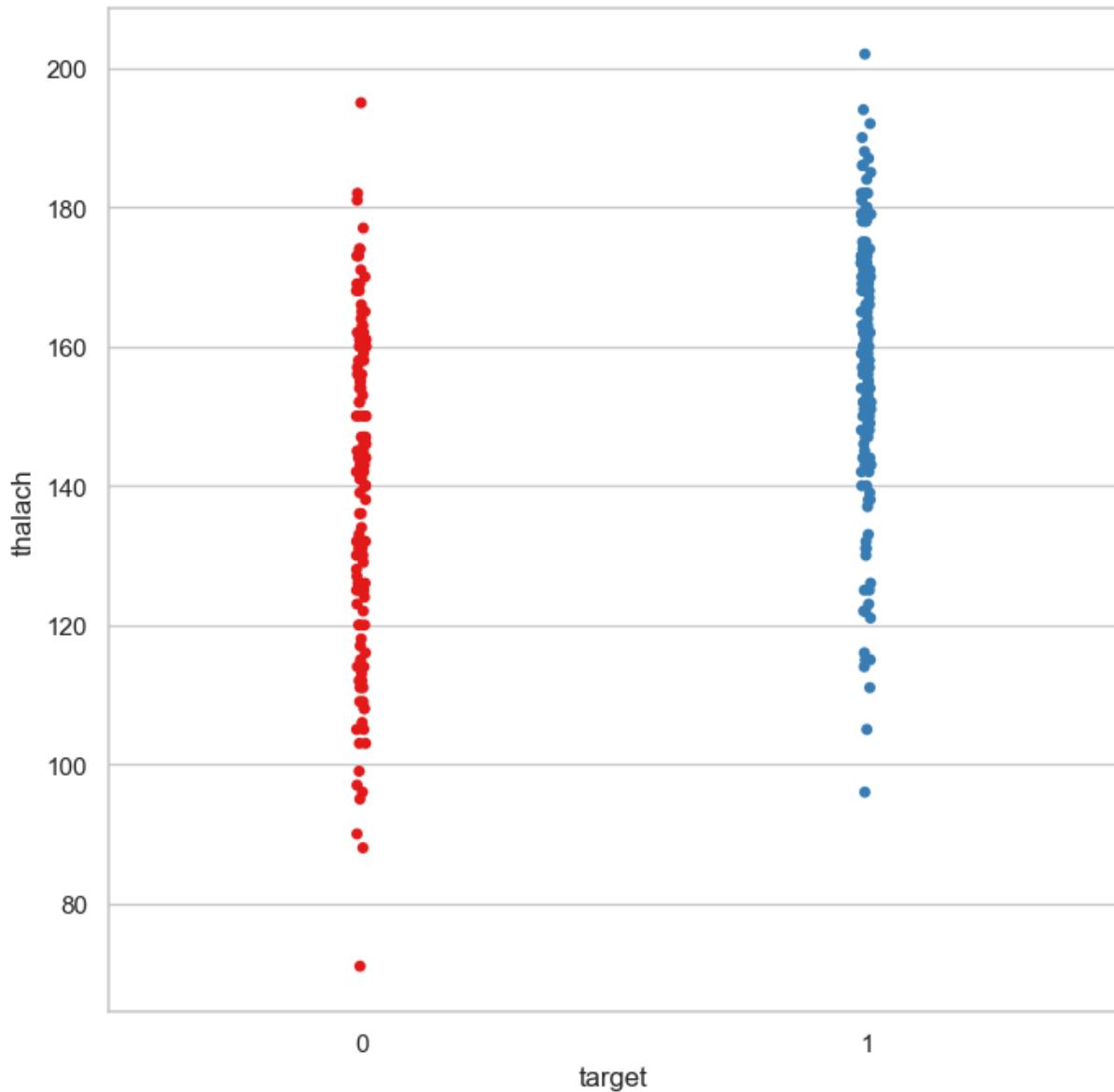
In [103...]

```
f,ax = plt.subplots(figsize=(20,20))
sns.stripplot(x='target',y='thalach',data=df,palette='Set2')
```

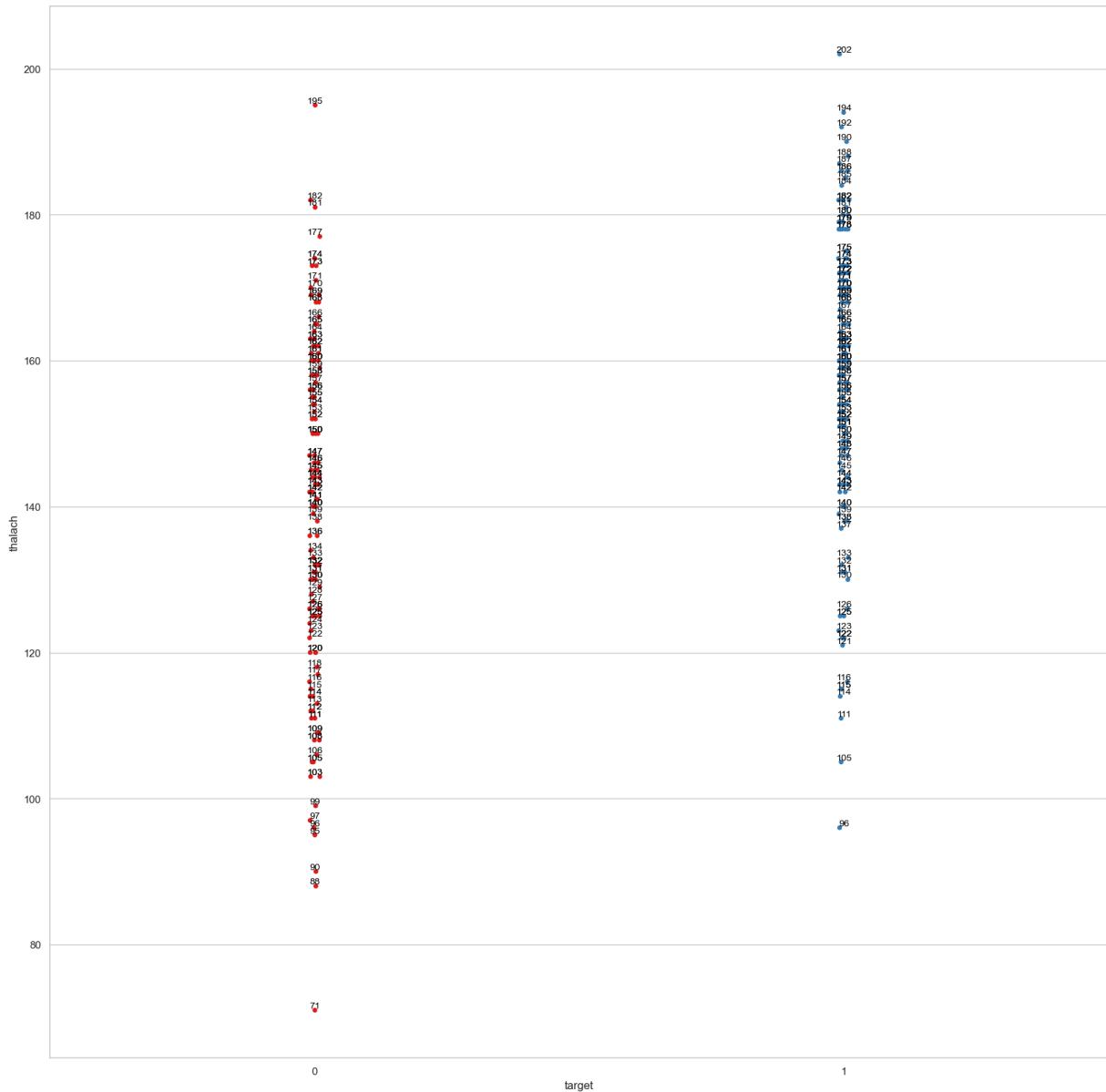
```
for i in range(len(df)):
    ax.text(df['target'].iloc[i], df['thalach'].iloc[i],
            str(df['thalach'].iloc[i]),
            color='black', ha='center', va='bottom', fontsize=10)
plt.show()
```



```
In [104...]: f, ax=plt.subplots(figsize=(8,8))
sns.stripplot(x='target',y='thalach',data = df,palette='Set1',jitter=0.01)
plt.show()
```

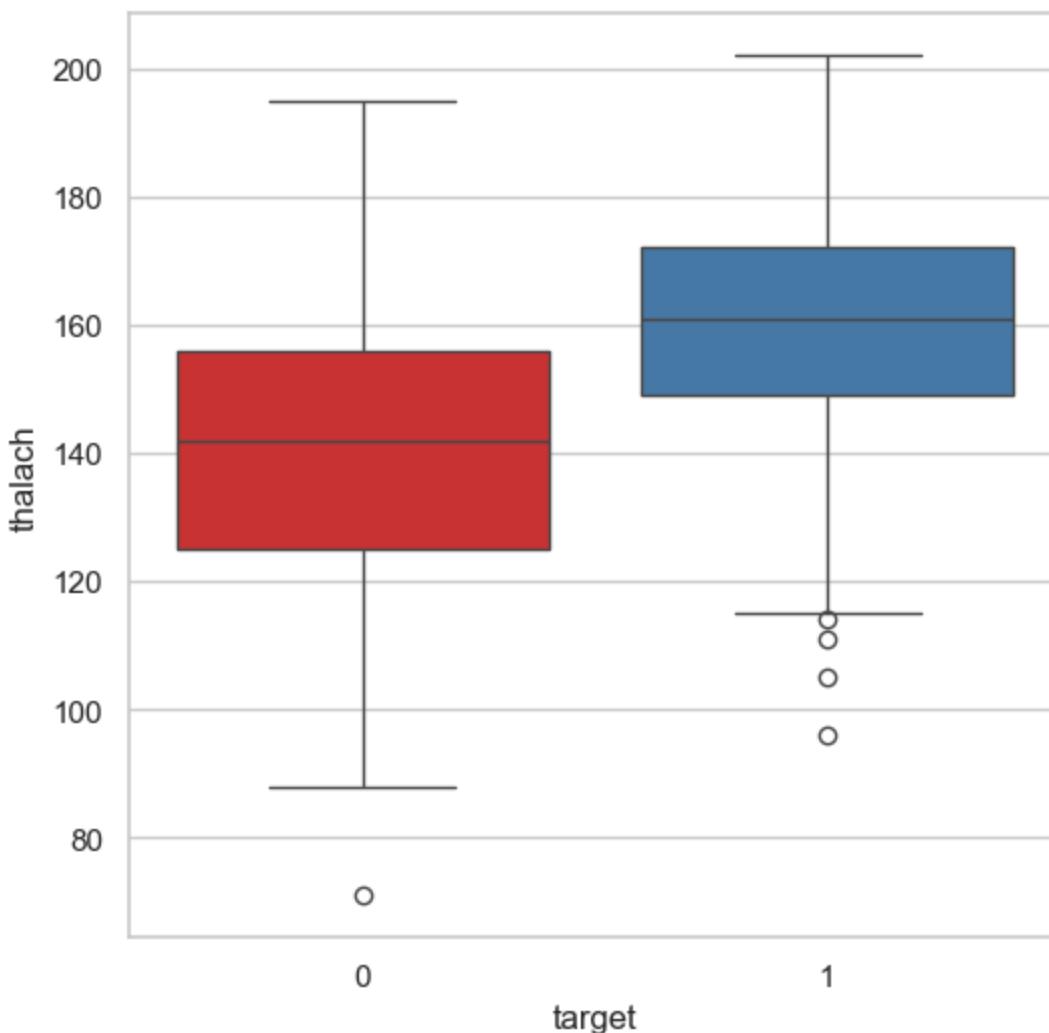


```
In [105]:  
f, ax = plt.subplots(figsize=(20,20))  
sns.stripplot(x='target',y='thalach',data = df,palette='Set1',jitter=0.01)  
for i in range(len(df)):  
    ax.text(df['target'].iloc[i], df['thalach'].iloc[i],  
            str(df['thalach'].iloc[i]),  
            color='black', ha='center', va='bottom', fontsize=10)  
plt.show()
```



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

```
In [107]:  
f,ax = plt.subplots(figsize=(6,6))  
sns.boxplot(x='target',y='thalach',data=df,palette='Set1')  
plt.show()
```



Interpretation

The above boxplot confirms our finding that people suffering from heart disease (target = 1) have relatively higher heart rate (thalach) as compared to people who are not suffering from heart disease (target = 0).

Findings of Bivariate Analysis

Findings of Bivariate Analysis are as follows –

- There is no variable which has strong positive correlation with `target` variable.
- There is no variable which has strong negative correlation with `target` variable.
- There is no correlation between `target` and `fbs`.
- The `cp` and `thalach` variables are mildly positively correlated with `target` variable.
- We can see that the `thalach` variable is slightly negatively skewed.

- The people suffering from heart disease (target = 1) have relatively higher heart rate (thalach) as compared to people who are not suffering from heart disease (target = 0).
- The people suffering from heart disease (target = 1) have relatively higher heart rate (thalach) as compared to people who are not suffering from heart disease (target = 0).

#Multivariate analysis

- The objective of the multivariate analysis is to discover patterns and relationships in the dataset.

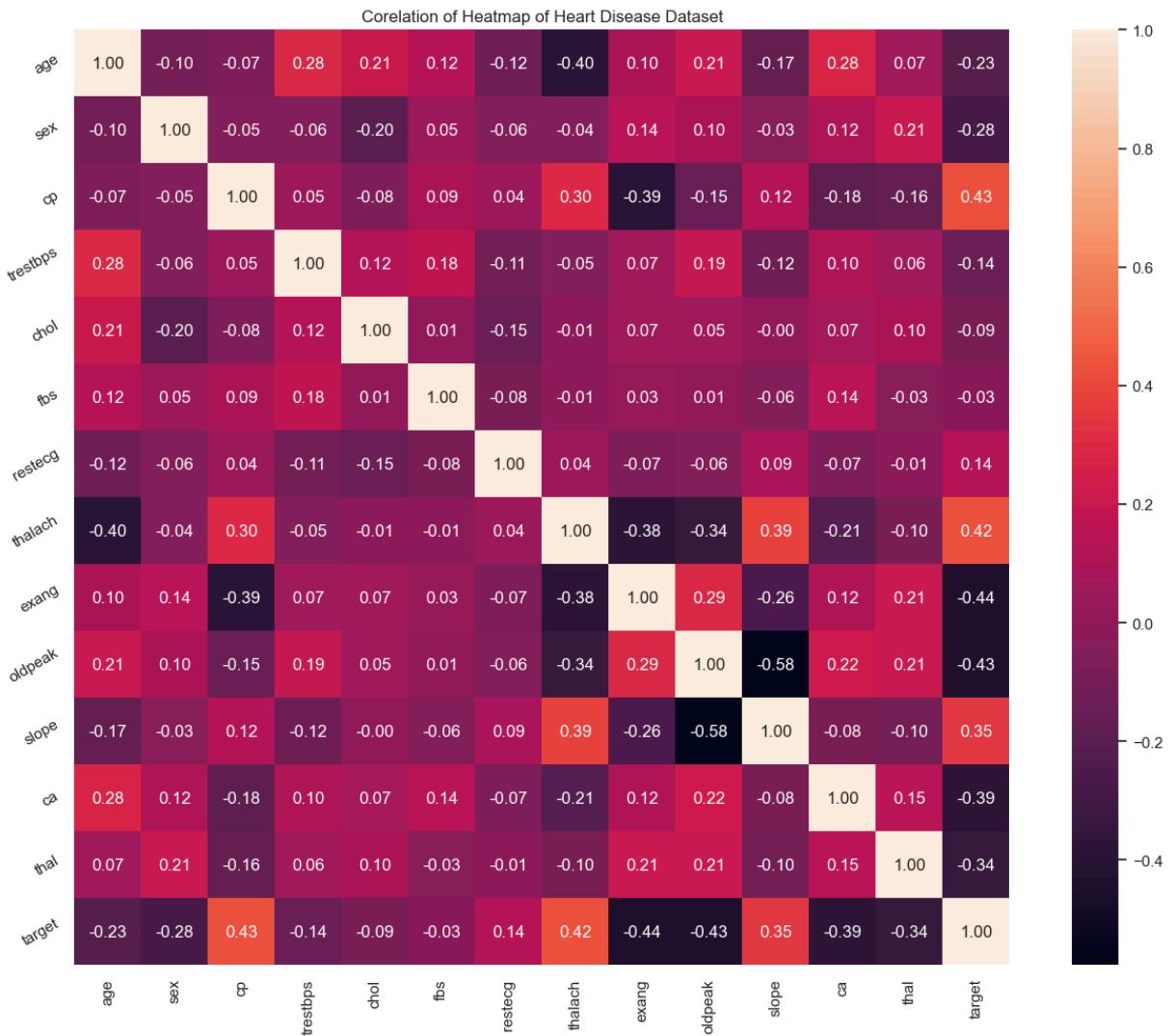
Discover patterns and relationships

- An important step in EDA is to discover patterns and relationships between variables in the dataset.
- I will use `heat map` and `pair plot` to discover the patterns and relationships in the dataset.
- First of all, I will draw a `heat map`.

#Heat map

In [113...]

```
plt.figure(figsize=(16,12))
plt.title("Corelation of 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()
```



- correlation: This should be a DataFrame or matrix containing the correlation values.
- `square=True` : This makes each cell square-shaped.
- `annot=True` : This displays the actual correlation values inside the heatmap cells.
- `fmt=' .2f '` : This controls the format of the annotations, limiting them to two decimal places.
- `linecolor='white'` : Adds white grid lines between the cells.

Interpretation

From the above correlation heat map, we can conclude that :-

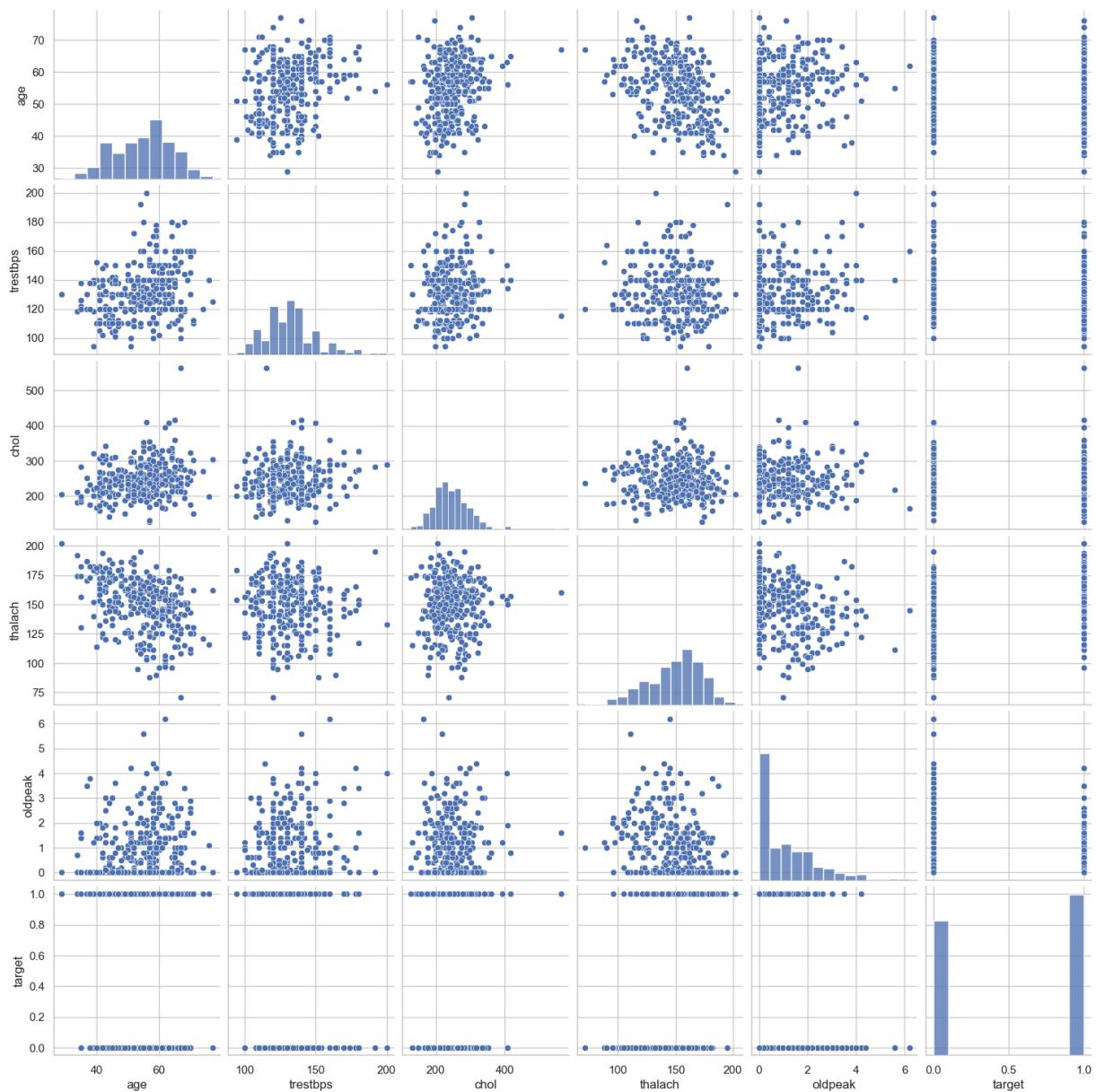
- `target` and `cp` variable are mildly positively correlated (correlation coefficient = 0.43).
- `target` and `thalach` variable are also mildly positively correlated (correlation coefficient = 0.42).
- `target` and `slope` variable are weakly positively correlated (correlation coefficient = 0.35).

- `target` and `exang` variable are mildly negatively correlated (correlation coefficient = -0.44).
- `target` and `oldpeak` variable are also mildly negatively correlated (correlation coefficient = -0.43).
- `target` and `ca` variable are weakly negatively correlated (correlation coefficient = -0.39).
- `target` and `thal` variable are also weakly negatively correlated (correlation coefficient = -0.34).

#Pair Plot

In [117...]

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



- df[num_var]: DataFrame df se specified numerical columns ko select karta hai.
- kind='scatter': Off-diagonal subplots par scatter plots banata hai, jo pairs of variables ke beech relation dikhata hai.
- diag_kind='hist': Diagonal plots par histograms banata hai jo har variable ki distribution ko show karta hai.

Comment

- I have defined a variable `num_var`. Here `age`, `trestbps`, `chol`, `thalach` and `oldpeak` are numerical variables and `target` is the categorical variable.
- So, I will check relationships between these variables.

Pairplot Kya Dikhata Hai:

- Scatter Plots: Diagonal ke bahar ke subplots do variables ke beech relationship ko dikhate hain. Jaise, 'age' aur 'thalach' ke beech ka plot.
- Histograms: Diagonal par histograms har variable ki distribution ko dikhate hain. Jaise, 'age' ke diagonal plot par age distribution dikhegi.

Analysis of Age and other variables

Check the number of unique values in `age` variable

```
In [123...]: df['age'].nunique()
```

```
Out[123...]: 41
```

#view statistical summary of age variable

```
In [125...]: df['age'].describe()
```

```
Out[125...]: 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
```

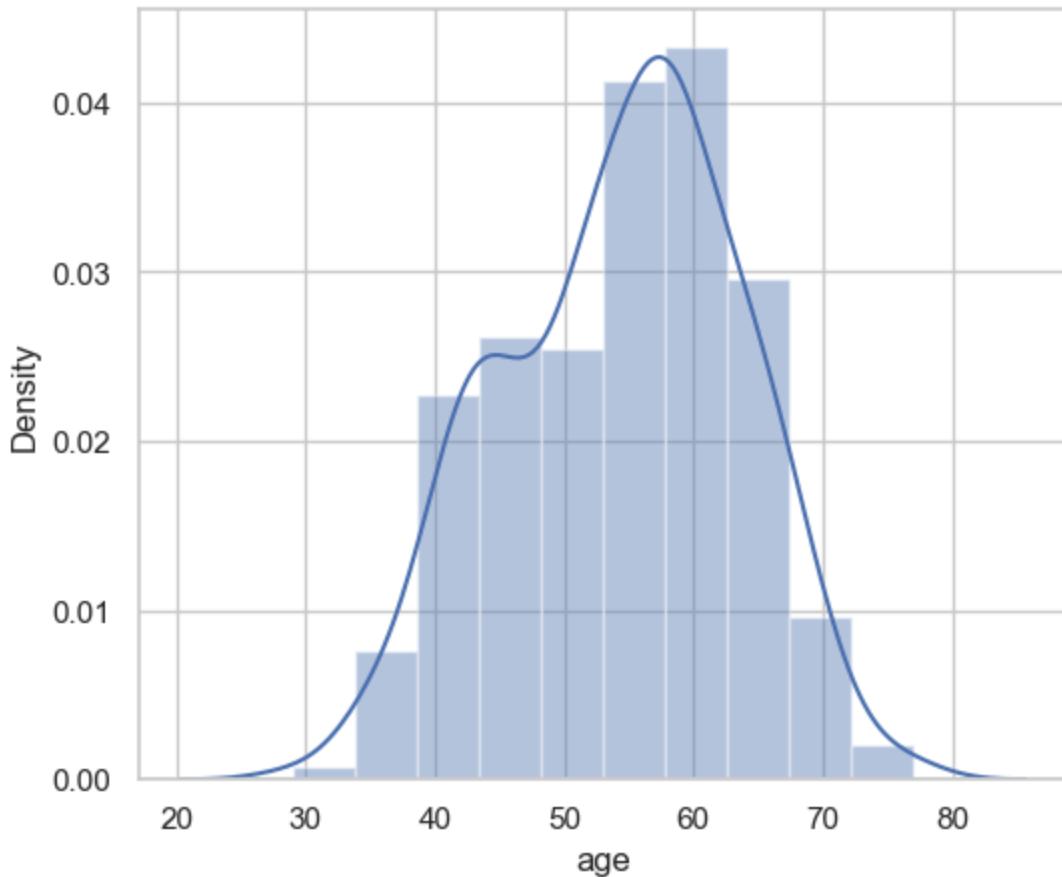
Interpretation

- The mean value of the `age` variable is 54.37 years.
- The minimum and maximum values of `age` are 29 and 77 years.

Plot the distribution of `age` variable

Now, I will plot the distribution of `age` variable to view the statistical properties.

```
In [128...]: f, ax = plt.subplots(figsize=(6,5))
x = df['age']
ax = sns.distplot(x,bins=10)
plt.show()
```



Interpretation

- The `age` variable distribution is approximately normal.

Analyze `age` and `target` variable

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

```
In [132...]: df['age'].unique()
```

```
Out[132...]: array([63, 37, 41, 56, 57, 44, 52, 54, 48, 49, 64, 58, 50, 66, 43, 69, 59,
       42, 61, 40, 71, 51, 65, 53, 46, 45, 39, 47, 62, 34, 35, 29, 55, 60,
       67, 68, 74, 76, 70, 38, 77], dtype=int64)
```

```
In [133...]: df['age'].nunique()
```

```
Out[133...]: 41
```

```
In [134...]: df.groupby('age')['target'].value_counts()
```

```
Out[134...]:
```

	age	target
29	1	1
34	1	2
35	0	2
	1	2
37	1	2
	..	
70	1	1
71	1	3
74	1	1
76	1	1
77	0	1

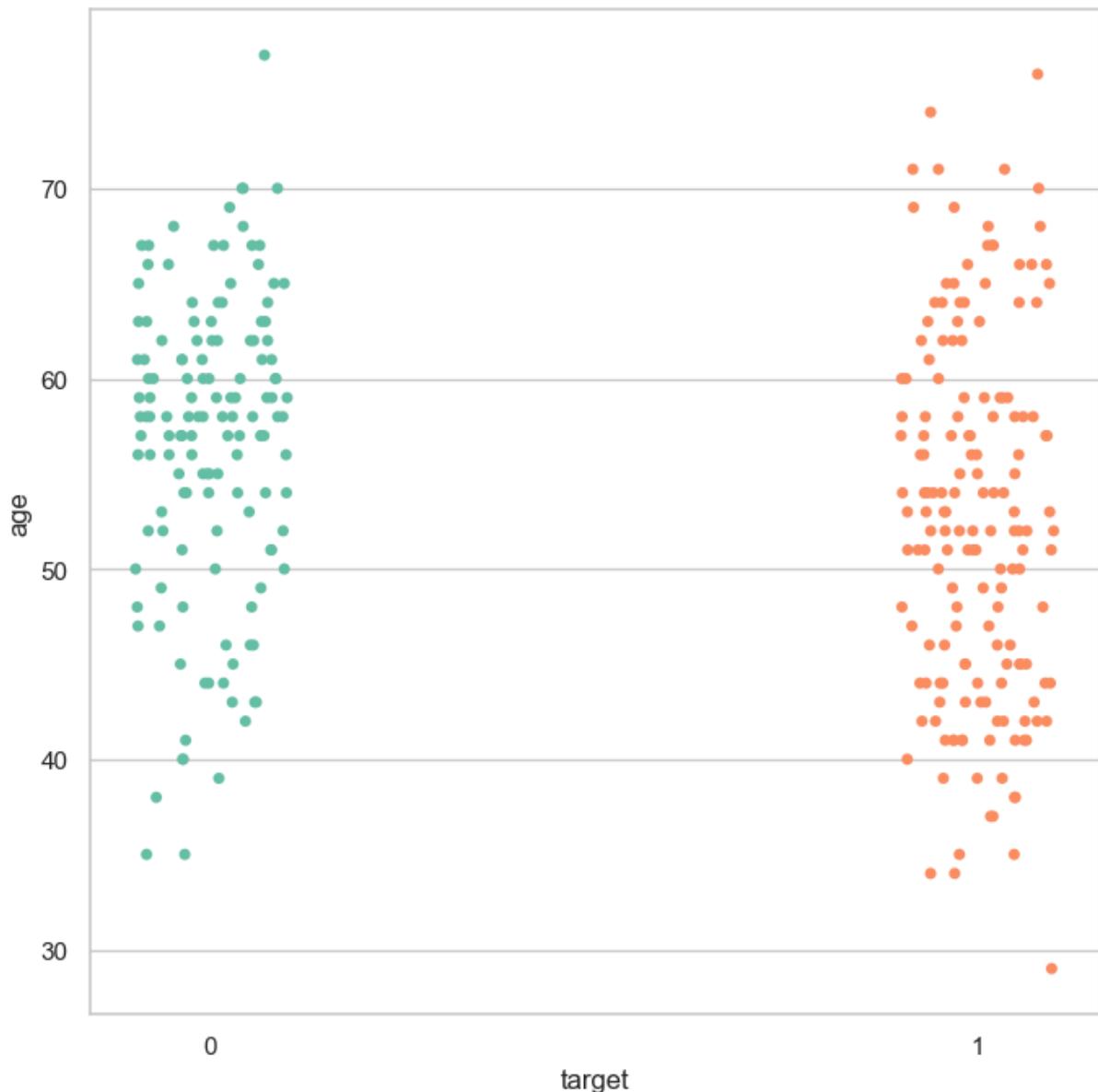
Name: count, Length: 75, dtype: int64

- 29 1 1: For age 29, there is 1 occurrence where the target is 1.
- 34 1 2: For age 34, there are 2 occurrences where the target is 1.
- 35 0 2: For age 35, there are 2 occurrences where the target is 0.
- 35 1 2: For age 35, there are 2 occurrences where the target is 1.

```
In [136...]:
```

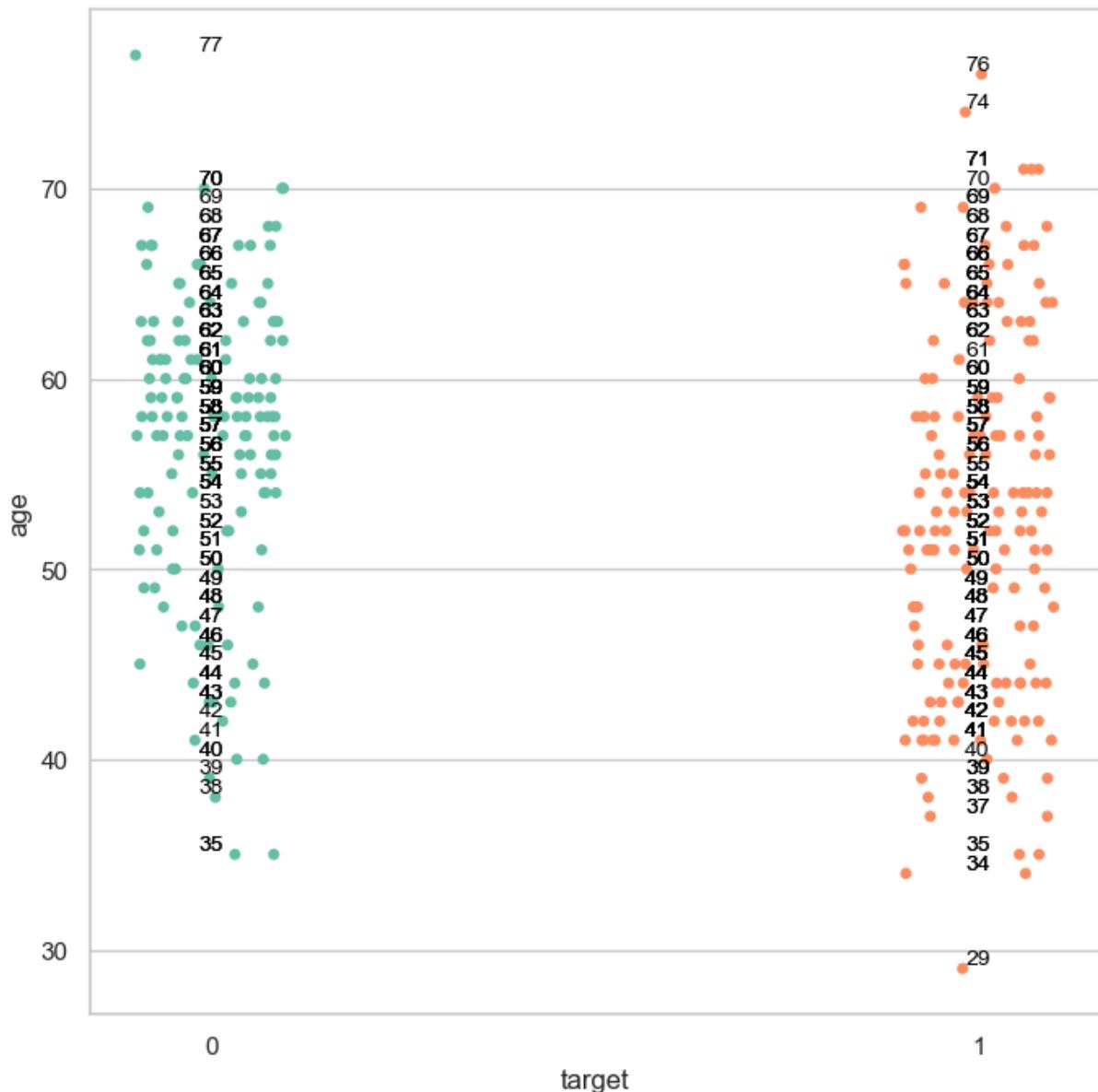
```
f,ax = plt.subplots(figsize=(8,8))
sns.stripplot(x='target',y='age',data= df, palette='Set2')
plt.plot()
```

```
Out[136...]: []
```



```
In [137...]:  
f,ax = plt.subplots(figsize=(8,8))  
sns.stripplot(x='target',y='age',data= df, palette='Set2')  
for i in range(len(df)):  
    ax.text(df['target'].iloc[i], df['age'].iloc[i],  
            str(df['age'].iloc[i]),  
            color='black', ha='center', va='bottom', fontsize=10)  
plt.plot()
```

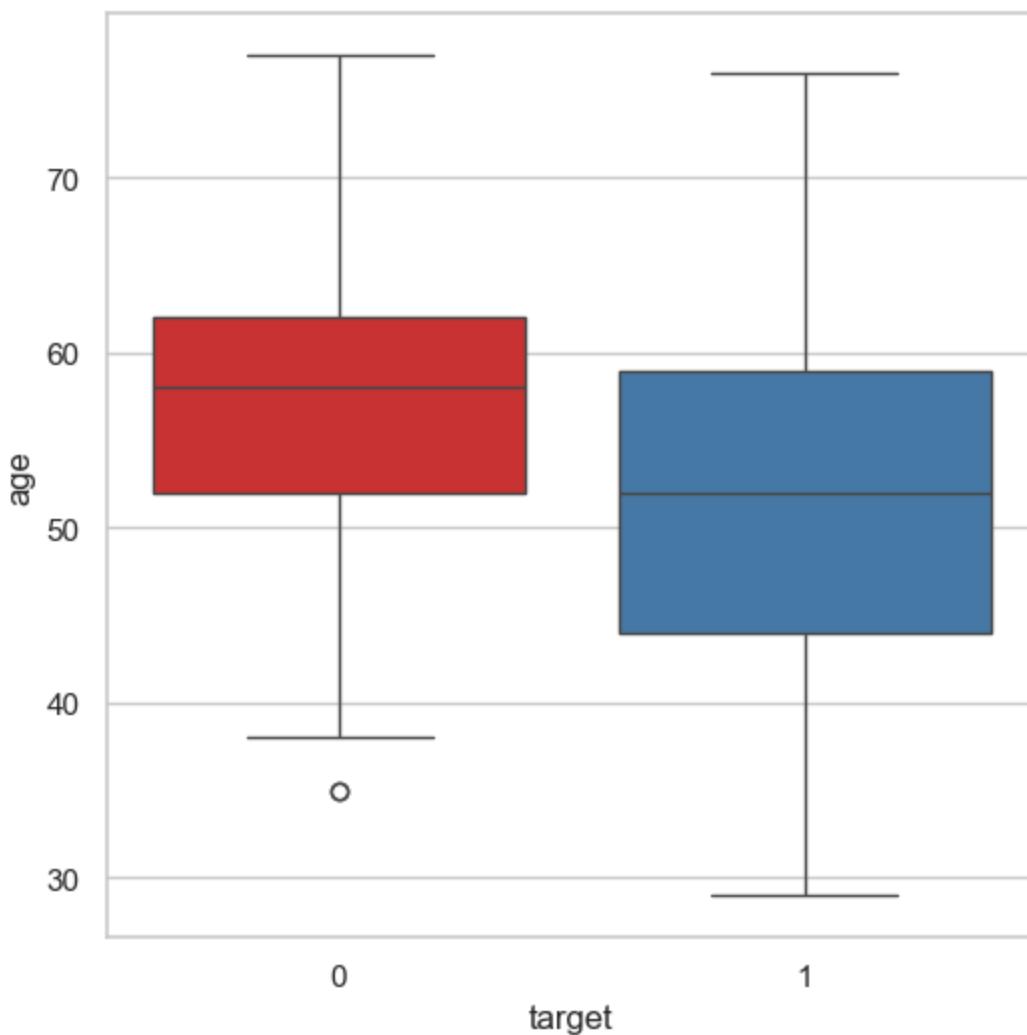
```
Out[137...]: []
```



Interpretation

- We can see that the people suffering from heart disease (target = 1) and people who are not suffering from heart disease (target = 0) have comparable ages.

```
In [139...]: f, ax = plt.subplots(figsize=(6,6))
sns.boxplot(x='target', y='age', data=df, palette='Set1')
plt.show()
```



Interpretation

- The above boxplot tells two different things :
 - The mean age of the people who have heart disease is less than the mean age of the people who do not have heart disease.
 - The dispersion or spread of age of the people who have heart disease is greater than the dispersion or spread of age of the people who do not have heart disease.

Analyze `age` and `trestbps` variable

`trestbps` : resting blood pressure (in mm Hg on admission to the hospital)

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

In [144...]

```
df.groupby('age')['trestbps'].value_counts()
```

```
Out[144...]:
```

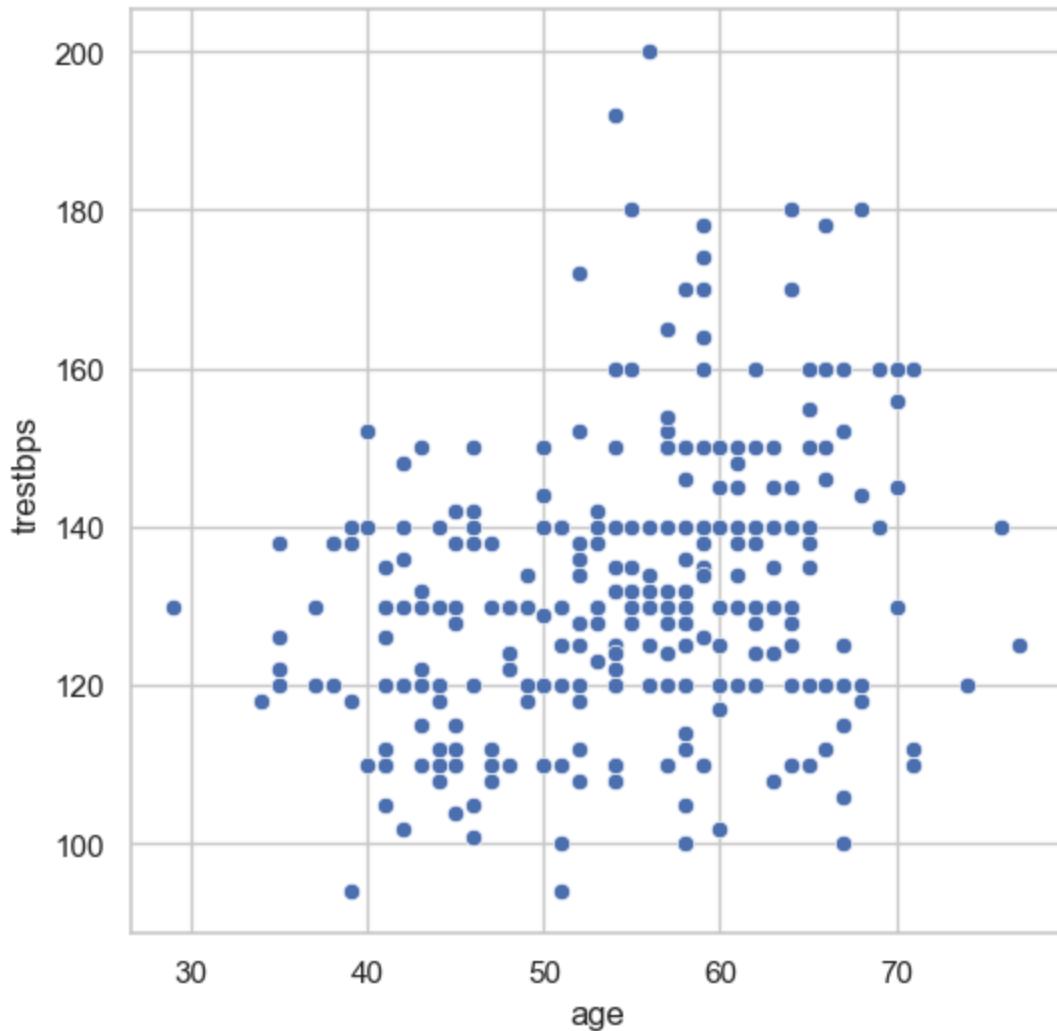
age	trestbps	
29	130	1
34	118	2
35	120	1
	122	1
	126	1
	..	
71	112	1
	160	1
74	120	1
76	140	1
77	125	1

Name: count, Length: 242, dtype: int64

- 29 130 1: At age 29, there is 1 person with a resting blood pressure of 130 mm Hg.
- 34 118 2: At age 34, there are 2 people with a resting blood pressure of 118 mm Hg.
- 35 120 1: At age 35, there is 1 person with a resting blood pressure of 120 mm Hg.
- 35 122 1: At age 35, there is 1 person with a resting blood pressure of 122 mm Hg.

```
In [146...]:
```

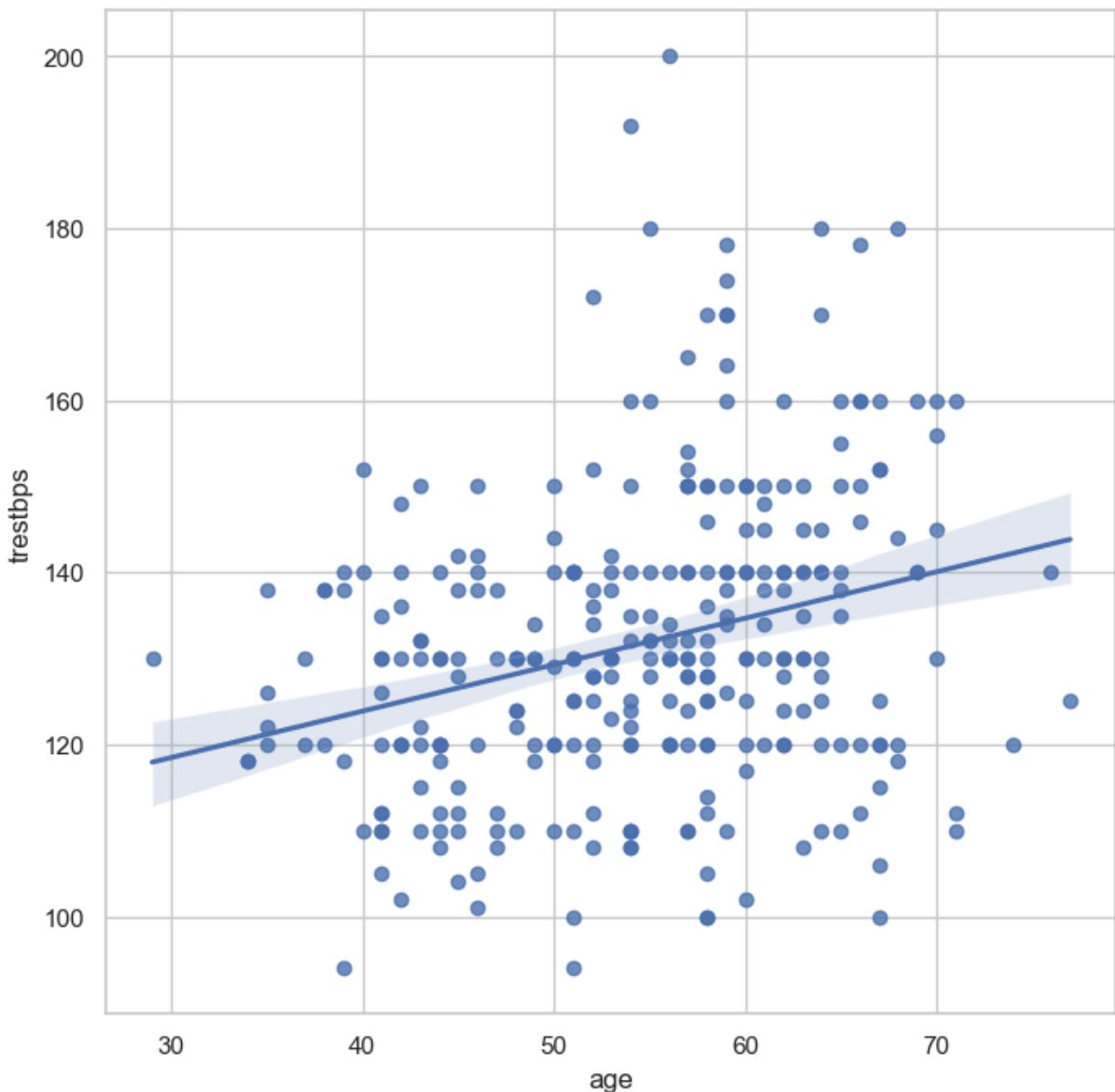
```
f,ax = plt.subplots(figsize=(6,6))
sns.scatterplot(x='age', y='trestbps', data=df, palette='Set1')
plt.show()
```



Interpretation

- The above scatter plot shows that there is no correlation between `age` and `trestbps` variable.

```
In [148]:  
f, ax = plt.subplots(figsize=(8,8))  
sns.regplot(x='age',y='trestbps',data=df)  
plt.show()
```



Interpretation

- The above line shows that linear regression model is not good fit to the data.

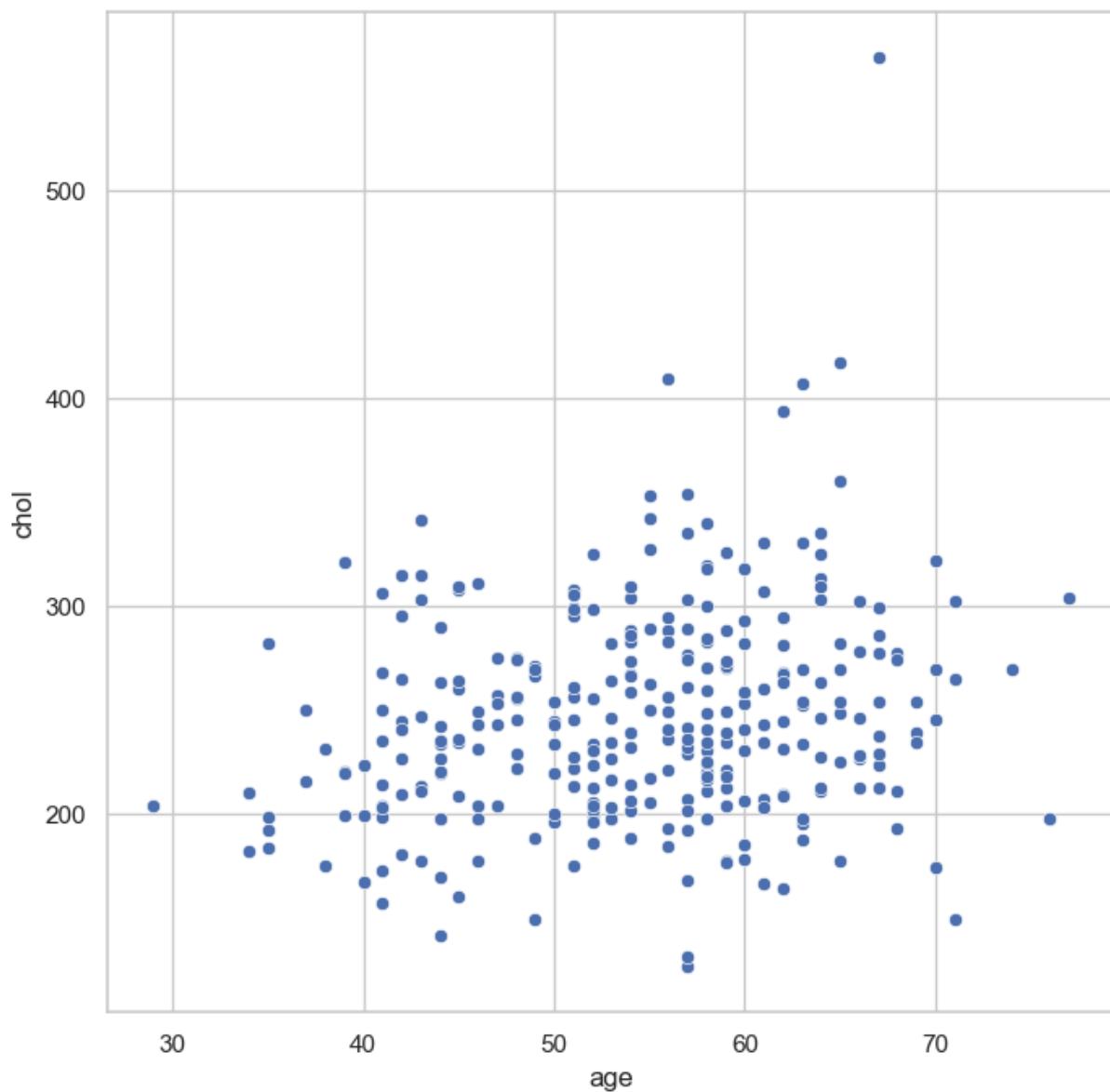
Analyze age and chol variable

- chol : serum cholesterol in mg/dl

```
In [152]: df.groupby('age')['chol'].value_counts()
```

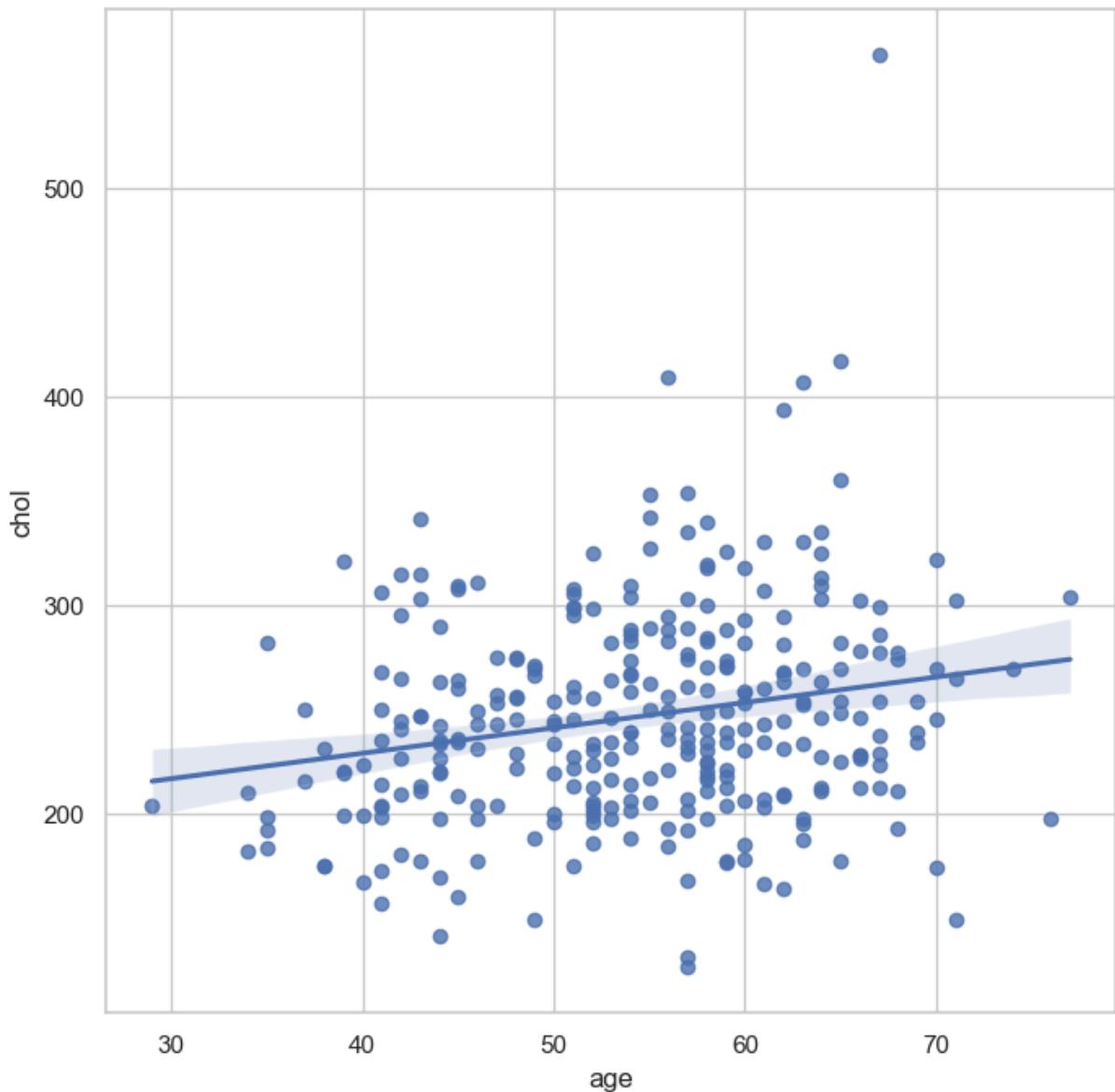
```
Out[152...]: age  chol
29    204    1
34    182    1
      210    1
35    183    1
      192    1
      ..
71    265    1
      302    1
74    269    1
76    197    1
77    304    1
Name: count, Length: 298, dtype: int64
```

```
In [153...]: f, ax = plt.subplots(figsize=(8,8))
sns.scatterplot(x='age',y='chol', data=df)
plt.show()
```



```
In [154...]: f, ax = plt.subplots(figsize=(8,8))
sns.regplot(x='age',y='chol',data=df)
```

```
plt.show()
```



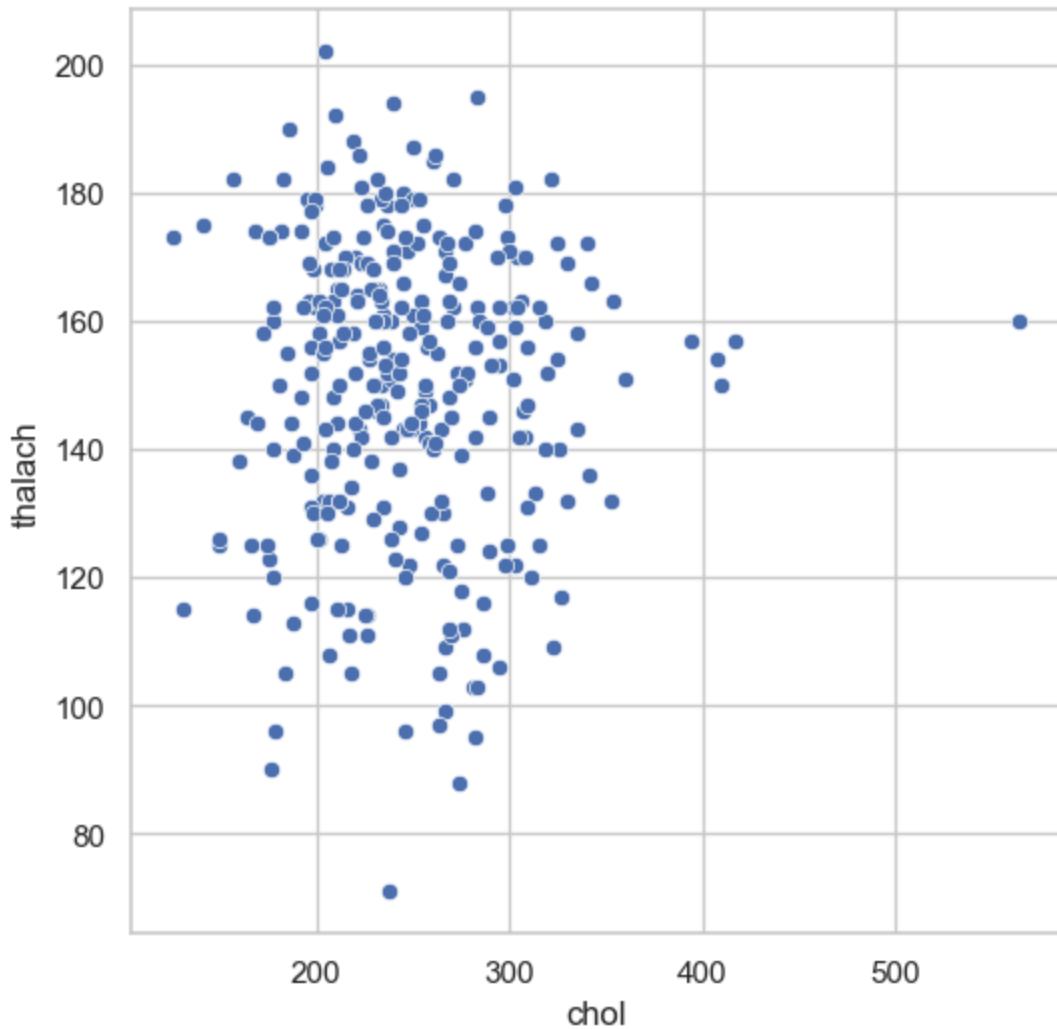
Interpretation

- The above plot confirms that there is a slightly positive correlation between `age` and `chol` variables.

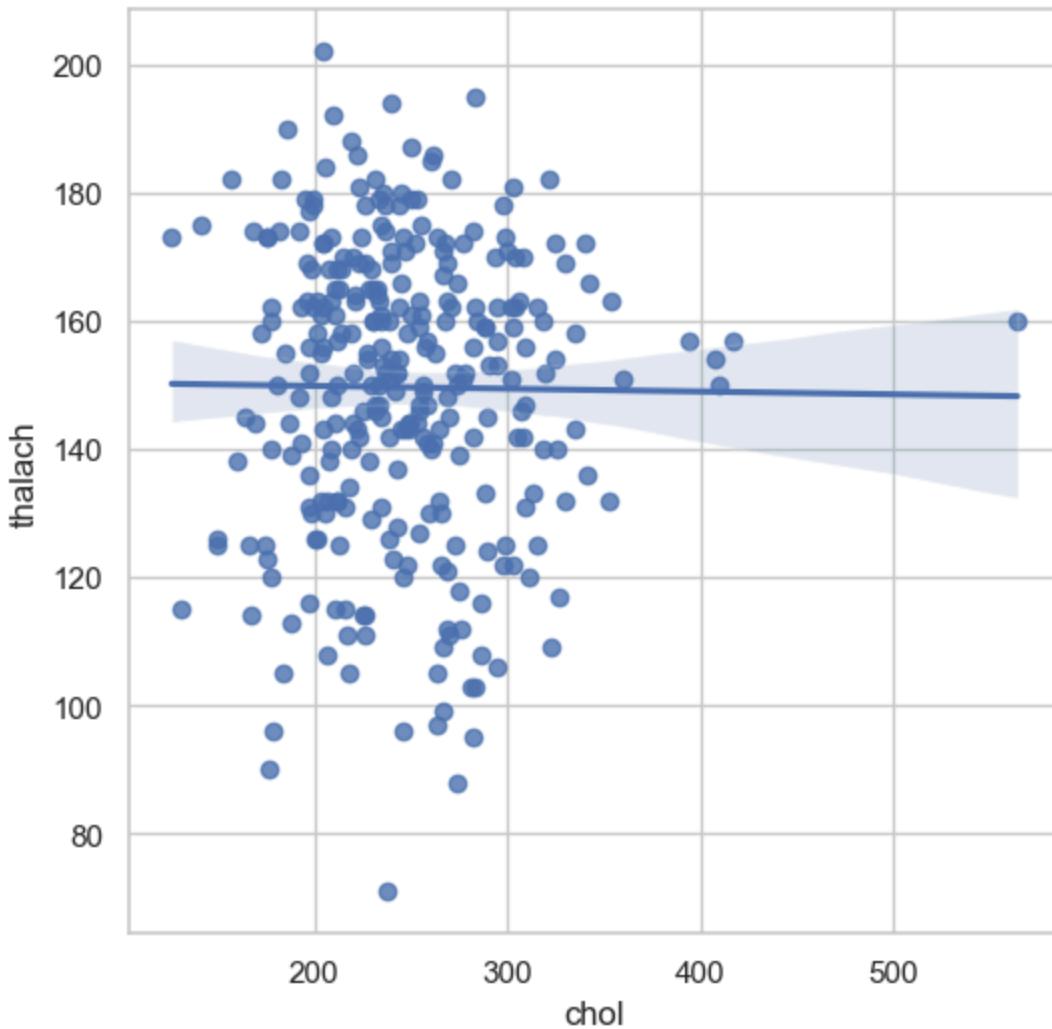
Analyze `chol` and `thalach` variable

In [157...]

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



```
In [230]:  
f,ax = plt.subplots(figsize=(6,6))  
sns.regplot(x='chol',y='thalach',data=df)  
plt.show()
```



Interpretation

- The above plot shows that there is no correlation between `chol` and `thalach` variable.

#Dealing with missing values

- In Pandas missing data is represented by two values:
 - None**: None is a Python singleton object that is often used for missing data in Python code.
 - NaN** : NaN (an acronym for Not a Number), is a special floating-point value recognized by all systems that use the standard IEEE floating-point representation.
- There are different methods in place on how to detect missing values.

Pandas `isnull()` and `notnull()` functions

- Pandas offers two functions to test for missing data - `isnull()` and `notnull()`. These are simple functions that return a boolean value indicating whether the passed in argument value is in fact missing data.
- Below, I will list some useful commands to deal with missing values.

Useful commands to detect missing values

- **`df.isnull()`**

The above command checks whether each cell in a dataframe contains missing values or not. If the cell contains missing value, it returns True otherwise it returns False.

- **`df.isnull().sum()`**

The above command returns total number of missing values in each column in the dataframe.

- **`df.isnull().sum().sum()`**

It returns total number of missing values in the dataframe.

- **`df.isnull().mean()`**

It returns percentage of missing values in each column in the dataframe.

- **`df.isnull().any()`**

It checks which column has null values and which has not. The columns which has null values returns TRUE and FALSE otherwise.

- **`df.isnull().any().any()`**

It returns a boolean value indicating whether the dataframe has missing values or not. If dataframe contains missing values it returns TRUE and FALSE otherwise.

- **`df.isnull().values.any()`**

It checks whether a particular column has missing values or not. If the column contains missing values, then it returns TRUE otherwise FALSE.

- **`df.isnull().values.sum()`**

It returns the total number of missing values in the dataframe.

In [243...]

```
#check missing values
df.isnull().sum()
```

```
Out[243...]: 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
```

Interpretation

We can see that there are no missing values in the dataset.

#Check with ASSERT statement

- We must confirm that our dataset has no missing values.
- We can write an **assert statement** to verify this.
- We can use an assert statement to programmatically check that no missing, unexpected 0 or negative values are present.
- This gives us confidence that our code is running properly.
- **Assert statement** will return nothing if the value being tested is true and will throw an AssertionError if the value is false.
- **Asserts**
 - assert 1 == 1 (return Nothing if the value is True)
 - assert 1 == 2 (return AssertionError if the value is False)

```
In [250...]: #assert that there are no missing values in the dataframe
           assert pd.notnull(df).all().all()
```

```
In [252...]: #assert all values are greater than or equal to 0
           assert (df >= 0).all().all()
```

Interpretation

- The above two commands do not throw any error. Hence, it is confirmed that there are no missing or negative values in the dataset.

- All the values are greater than or equal to zero.

#Outlier detection

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

`age , trestbps , chol , thalach and oldpeak` variables.

age variable

In [260...]: `df['age'].describe()`

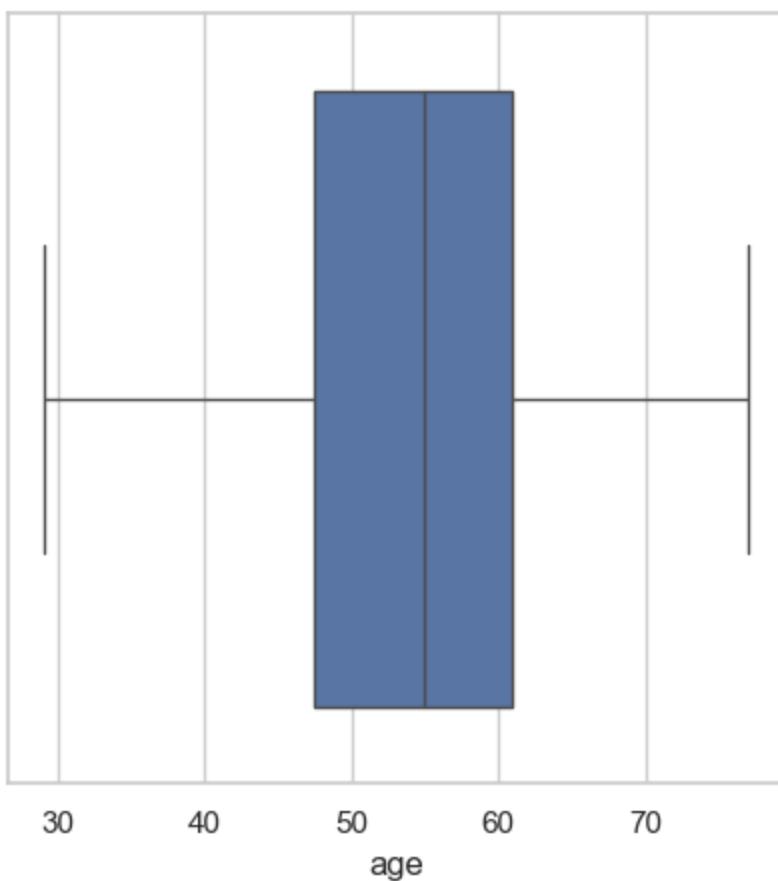
Out[260...]:

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 [269...]:

```
f,ax = plt.subplots(figsize=(5,5))
sns.boxplot(x=df['age'])
plt.show()
```



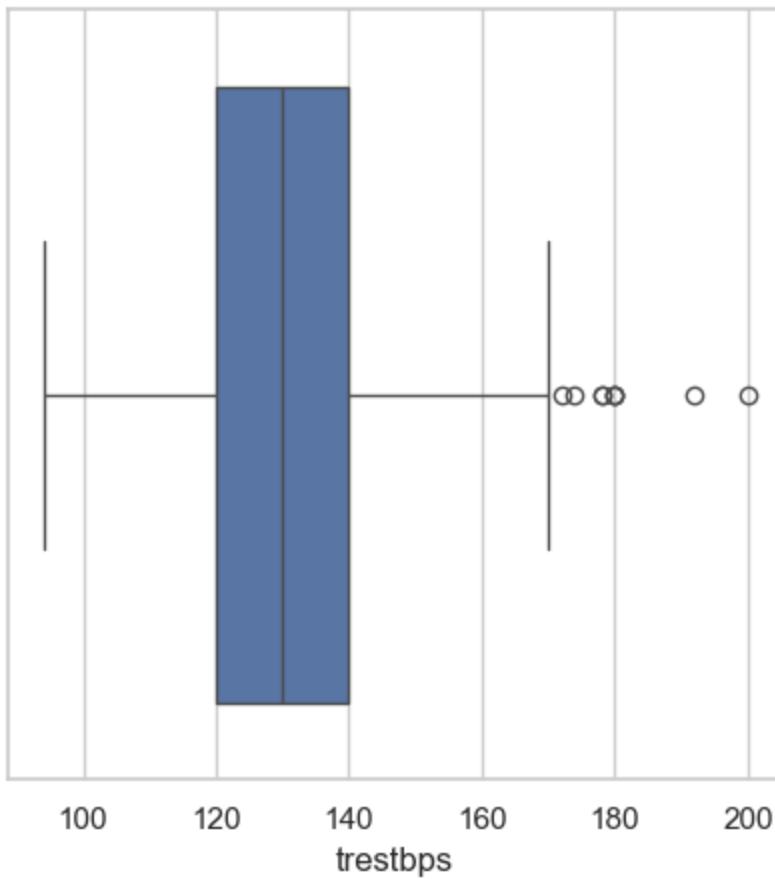
trestbps variable

```
In [272...]: df['trestbps'].describe()
```

```
Out[272...]: 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 [279...]: f, ax = plt.subplots(figsize=(5,5))
sns.boxplot(x=df['trestbps'])
plt.show()
```



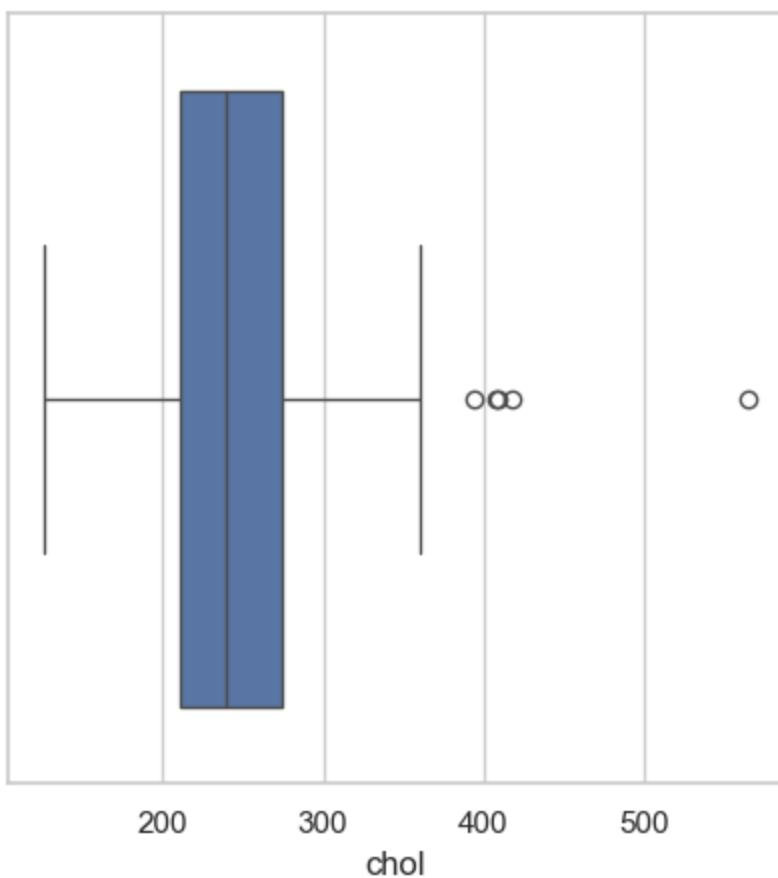
chol variable

```
In [282...]: df['chol'].describe()
```

```
Out[282...]: 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 [296...]: f,ax = plt.subplots(figsize=(5,5))
sns.boxplot(x=df['chol'])
plt.show()
```

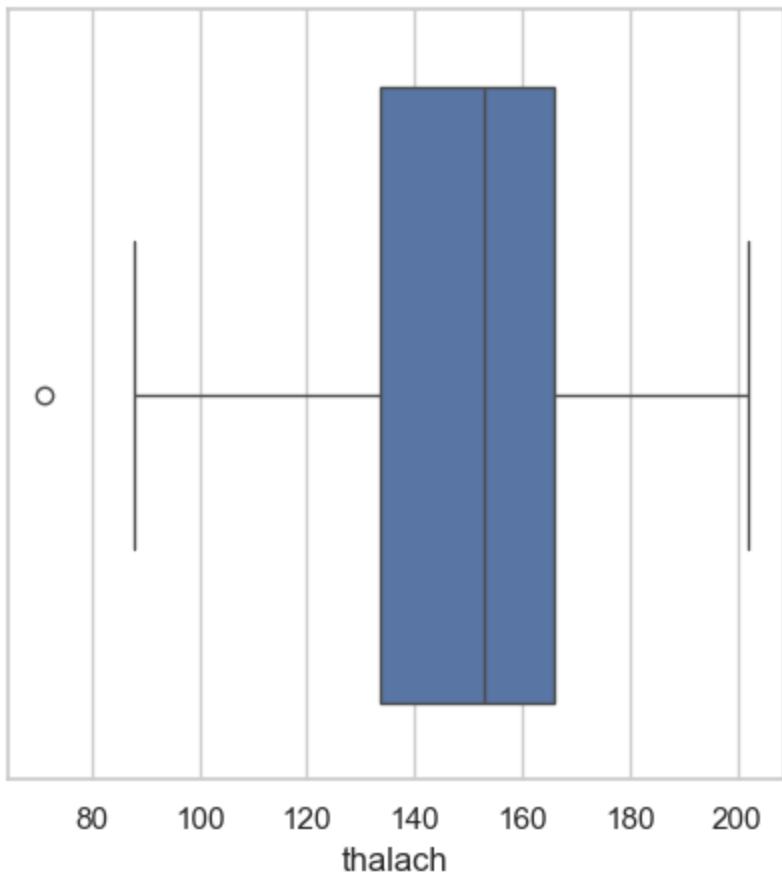


Thalach variable

```
In [299...]: df['thalach'].describe()
```

```
Out[299...]: 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 [301...]: f,ax = plt.subplots(figsize=(5,5))
sns.boxplot(x=df['thalach'])
plt.show()
```

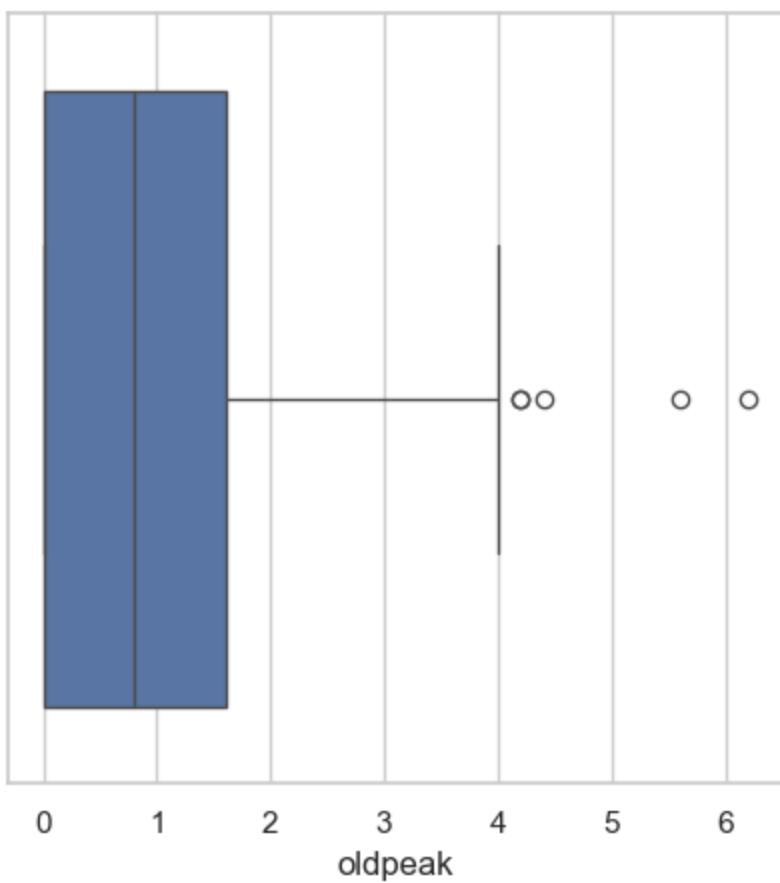


oldpeak variable

```
In [304...]: df['oldpeak'].describe()
```

```
Out[304...]: 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 [308...]: f, ax = plt.subplots(figsize=(5,5))
sns.boxplot(x=df['oldpeak'])
plt.show()
```



Findings

- The `age` variable does not contain any outlier.
- `trestbps` variable contains outliers to the right side.
- `chol` variable also contains outliers to the right side.
- `thalach` variable contains a single outlier to the left side.
- `oldpeak` variable contains outliers to the right side.
- Those variables containing outliers needs further investigation.

In []: