

In []:

```
1 '''Heart Dataset Analysis
2
3 The goal of this dataset is to train a model so that it predicts
4 whether a person is likely to suffer from heart disease (whether the probability is above or below 50%); 
5 however, in this exercise/project, we are simply going to observe and analyze the distribution of the data,
6 search for outliers and missing values, and assess the relationships between features.
7
8 REQUIRED:
9
10 check for outliers, missing values, and the trends and
11 relationships between different features of the dataset to gain a better understanding of the
12 available data and derive useful insights from it.
13
14 ...
15'''
```

In []:

```
1 ...
2 *Loading and Understanding the Data
3
4 *Dealing with Outliers
5
6 *Plotting the Distributions and Relationships Between Specific Features
7
8 *Plotting Distributions and Relationships between Columns with Respect to the Target Column
9
10 *Plotting the Relationship between the Presence of Heart Disease and Maximum Recorded Heart Rate
11
12 *Plotting the Relationship between the Presence of Heart Disease and the Cholesterol Column
13
14 *Observing Correlations with a Heatmap
15
16'''
```

In [424]:

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import seaborn as sns
5 Path = '/Users/tomisin/Dropbox/My Mac (Tomisin-MacBook-Pro.local)/Documents/MY WORKSPACE/Database/heart.csv'
6 df = pd.read_csv(Path)
```

In [425]:

```
1 df=df.rename(columns={'age':'AGE','trestbps':'TRestBPS', 'sex':'SEX','restecg':'restECG', 'chol':'CHOL', 'thalach':'Thalach', 'oldpeak':'OldPeak', 'slope':'SLOPE', 'ca':'CA', 'thal':'THAL'})
```

In [426]:

```
1 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 [427]: 1 df.describe()

Out[427]:

	AGE	SEX	CP	TRestBPS	CHOL	FBS	restECG	Thalach	EXANG	OldPeak	SLOPE
count	303.000000	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	1.399340
std	9.082101	0.466011	1.032052	17.538143	51.830751	0.356198	0.525860	22.905161	0.469794	1.161075	0.616226
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	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	1.000000
50%	55.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000
75%	61.000000	1.000000	2.000000	140.000000	274.500000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	202.000000	1.000000	6.200000	2.000000

In [428]:

```
1  
2  
3 #INFORMATION ABOUT THE DATASET  
4  
5 ...  
6  
7  
8 *** df.info() shows the structure of this dataset is one that has 303 samples, a total of 14 variables  
9 with all non-categorical features  
10  
11 *** df.describe() shows if the data in each variable normally distributed (gaussian curve) or not.  
12 The data in this data set is generally skewed except for the CP feature.  
13  
14 *** There are no missing values in this dataset  
15  
16  
17  
18 ...
```

Out[428]: '\n\n*** df.info() shows the structure of this dataset is one that has 303 samples, a total of 14 variables\nwith all non-categorical features\n\n*** df.describe() shows if the data in each variable normally distributed (gaussian curve) or not.\n\nThe data in this data set is generally skewed except for the CP feature.\n\n\n\n'

In [429]:

```
1  '''There are no missing values in this dataset. df.info(), df.describe() tells me so'''
2
3
4 #Dealing with Outliers
5
6
7  '''There are certainly outliers in some of the features.'''
8 #In 'age', there is no outlier, minimum and maximum ages are 29 and 77 respectively.
9 #Data sample was taken amongst adults especially the older ones.
10
11 #In 'sex', there is/are no outlier, sex can either be 1 or 0 (male or female).
12
13 #In 'cp', there are no outliers, 4 data points in the density plot
14
15 #In 'trestbps', the presence of outliers indicating values between 170 till 200(max value)
16 #does not cause skewness in the data. The density plot using matplotlib appears quite normally distributed
17 #except for the hunches at these values which represent the outliers.
18
19 #In 'chol', presence of outliers indicating values between 350 till 564(max value).
20 #does not cause skewness in the data.
21 #The density plot using matplotlib appears quite normally distributed
22 #except for the hunches at these values which represent the outliers.
23
24 #In 'fbs', we have more zeros (250) than ones (53) making 1 an outlier,
25 #The ones should not be removed from this column in my opinion.
26 #The zeros and ones definitely infer something. Hopefully, the right readings were collected and reported.
27
28 #In 'restecg', presence of no outlier, 3 data points in the density plot
29
30 #In 'thalach', presence of an outliers indicating values between 71 and 120.
31 #Should not be removed assuming the right readings were collected and reported
32
33 #In 'exang', there are no outliers, 2 data points in the density plot.
34
35 #In 'oldpeak', presence of outliers between 4 and 6.2.
36
37 #In 'slope', there are no outliers, 3 data points in the density plot
38
39 #In 'ca', presence of outliers between 3 and 4.
40 #4 data points in the density plot.
41 #The outliers should not be removed assuming the right readings were collected and reported.
```

```
42  
43 #In 'thal', presence of outlier representing the 0 input  
44 #The zeros should not be removed assuming the right readings were collected and reported.
```

Out[429]: 'There are certainly outliers in some of the features.'

In [430]:

```
1 #UNIVARATE ANALYSIS
2
3 ...
4
5 *** For 'AGE' variable, the distribution shows participants from the ages 29-77.
6 It is likely that both young and old adults within this age range are prone to develop a heart disease
7 even though we do not have so many participants from ages (29-40).
8
9 *** For 'SEX' variable, there are more males than females.
10 Either of these genders can be prone to have a heart disease
11
12 *** For 'CP' variable, majority of the participants (246) have an angina-related chest pain.
13 I would want to think that angina is one of the symptom of a diseased heart.
14 This could also mean that some of those with non-angina chest pain may also have a diseased heart.
15
16 *** For 'TRestBPS' variable, A total of 243 participants have either
17 borderline high blood pressure or high blood pressure which already exceeds
18 the total number of participants with a diseased heart (165).
19 This could also mean that some of those with normal BP may also have a diseased heart
20 due to the influence of another variable in this dataset.
21
22 *** For 'CHOL', a total of 252 participants have either
23 borderline high levels or high levels of cholesterol in their blood.
24 There is a likelihood that most of those with a diseased heart have excess cholesterol in their blood.
25 However, if some or most of the rest of the participants with normal cholesterol levels have a diseased heart
26 it may be due to the influence of another variable in this dataset.
27
28 *** For 'FBS', 258 out of the 303 participants have sugar in their blood.
29 Majority of reported diseased heart participants will have high amounts sugar in their blood.
30 However if majority of the remaining 45 participants have a diseased heart,
31 Depending on the proportion of those with a normal blood sugar levels who have a diseased heart
32 FBS may not be an important feature in the model design.
33
34 *** For 'restECG', we have like a 1:1 of those with a normal and an abnormal ECG.
35 Depending on the proportion of those with a normal ECG who have a diseased heart
36 restECG may not be an important feature in the model design.
37
38 *** For 'Thalach', an occurrence of varying heart rates at its maximum ranging from 77-202.
39 Higher maximum heart rates in a patient would mean the presence of Thalassemia disease.
40 Depending on the proportion of those with a normal Thalach who have a diseased heart
41 Thalach may not be an important feature in the model design.
```

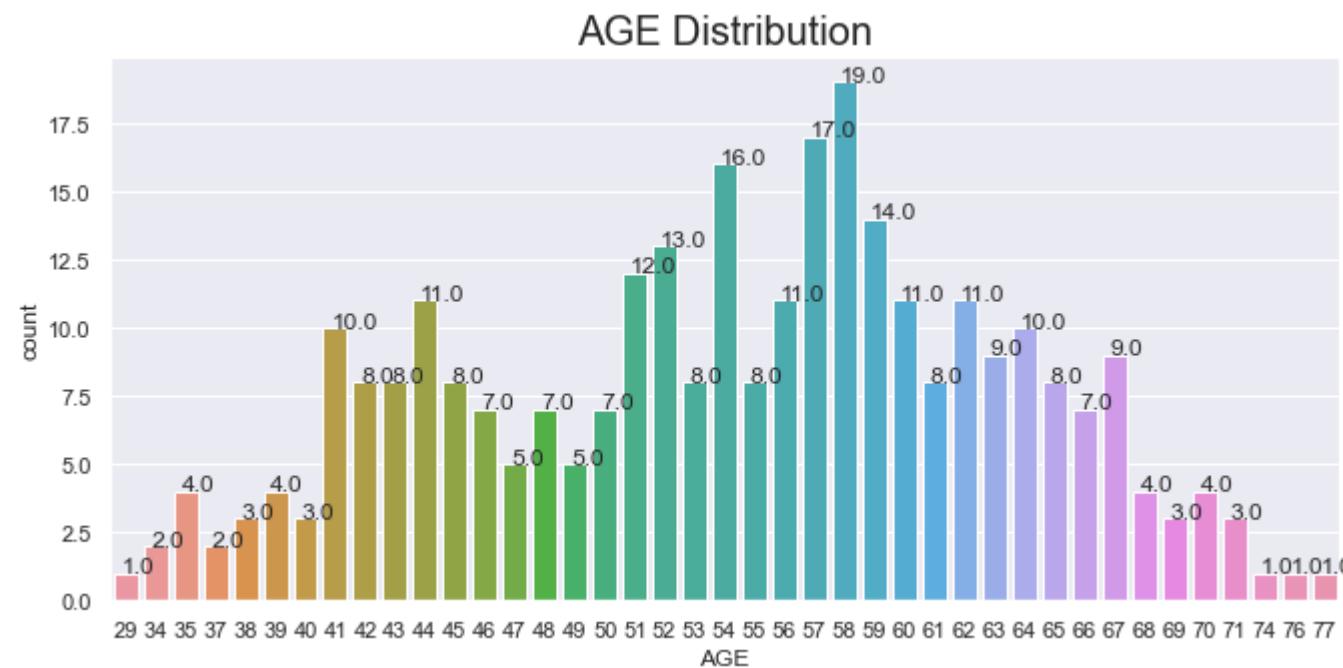
```
42
43 *** For 'EXANG', one of these two types of angina (non-induced and induced angina)
44 are seen in all the participants. Either of them can be a sign of a heart disease.
45 This variable may likely be an important one in model design.
46
47 *** For 'CA', 175 participants have no coloured artery, the rest have at least 1.
48 Depending on the proportion of those with a no colored artery who have a diseased heart
49 CA may not be an important feature in the model design.
50
51 *** For 'THAL', 283 people have fixed and reversible defect, the remaining 20 have no defect.
52 This would mean that some of those who have a fixed and reversible defect have a heart disease.
53
54 *** For 'TARGET', 165 have a heart disease and 138 do not. More like 1:!.
55 With most of the features having majority of the participant on the extremities (positive or negative)
56 May be difficult to decide which of them to keep or do away with
57
58
59
60
61
62 ...
```

Out[430]: "
*** For 'AGE' variable, the distribution shows participants from the ages 29-77. It is likely that both young and old adults within this age range and prone to develop a heart disease even though we do not have so many participants from ages (29-40).
*** For 'SEX' variable, there are more males than females. Either of these genders can be prone to have a heart disease
*** For 'CP' variable, majority of the participants (246) have an angina-related chest pain. It would want to think that angina is one of the symptom of a diseased heart.
This could also mean that some of those with non-angina chest pain may also have a diseased heart.
*** For 'TRestBPS' variable, A total of 243 participants have either borderline high blood pressure or high blood pressure which already exceeds the total number of participants with a diseased heart (165).
This could also mean that some of those with normal BP may also have a diseased heart due to the influence of another variable in this dataset.
*** For 'CHOL', a total of 252 participants have either borderline high levels or high levels of cholesterol in their blood.
There is a likelihood that most of those with a diseased heart have excess cholesterol in their blood.
However, if some or most of the rest of the participants with normal cholesterol levels have a diseased heart, it may be due to the influence of another variable in this dataset.
*** For 'FBS', 258 out of the 303 participants have sugar in their blood.
Majority of reported diseased heart participants will have high amounts sugar in their blood.
However if majority of the remaining 45 participants have a diseased heart, Depending on the proportion of those with a normal blood sugar levels who have a diseased heart
FBS may not be an important feature in the model design.
*** For 'restECG', we have like a 1:1 of those with a normal and an abnormal ECG.
Depending on the proportion of those with a normal ECG who have a diseased heart
restECG may not be an important feature in the model design.
...
*** For 'Thalach', an occurrence of varying heart rates at its maximum ranging from 77-202.
Higher max

imum heart rates in a patient would mean the presence of Thalassemia disease. \nDepending on the proportion o f those with a normal Thalach who have a diseased heart\nThalach may not be an important feature in the model design. \n\n*** For 'EXANG', one of these two types of angina (non-induced and induced angina) \nare seen in all the participants. Either of them can be a sign of a heart disease. \nThis variable may likely be an impor tant one in model design.\n\n*** For 'CA', 175 participants have no coloured artery, the rest have at least 1. \nDepending on the proportion of those with a no colored artery who have a diseased heart\nCA may not be a n important feature in the model design.\n\n*** For 'THAL', 283 people have fixed and reversible defect, the remaining 20 have no defect. \nThis would mean that some of those who have a fixed and reversible defect have a heart disease.\n\n*** For 'TARGET', 165 have a heart disease and 138 do not. More like 1:!. \nWith most of the features having majority of the participant on the extremities (positive or negative)\nMay be difficult t o decide which of them to keep or do away with\n\n\n\n\n"

In [433]:

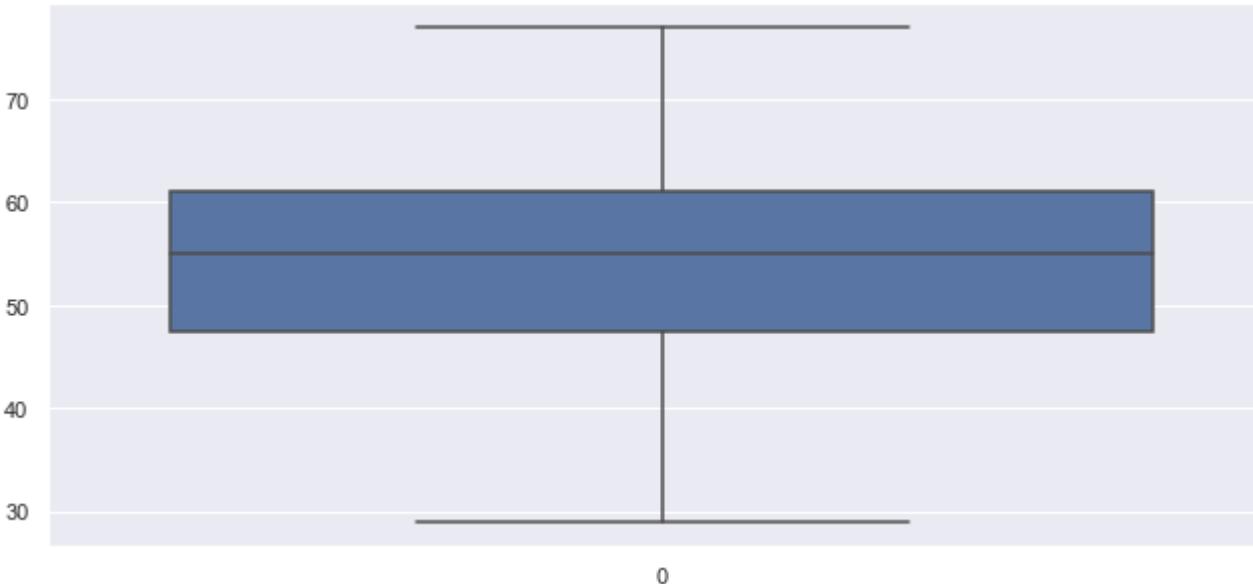
```
1 ax = sns.countplot(x="AGE", data=df)
2 ax.set_title('AGE Distribution', fontsize=20)
3
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
10 sns.set(rc={'figure.figsize':(11,5)})
```



In [434]:

```
1 import seaborn as sb  
2 sb.boxplot(data=df['AGE'])
```

Out[434]: <AxesSubplot:>



In [435]:

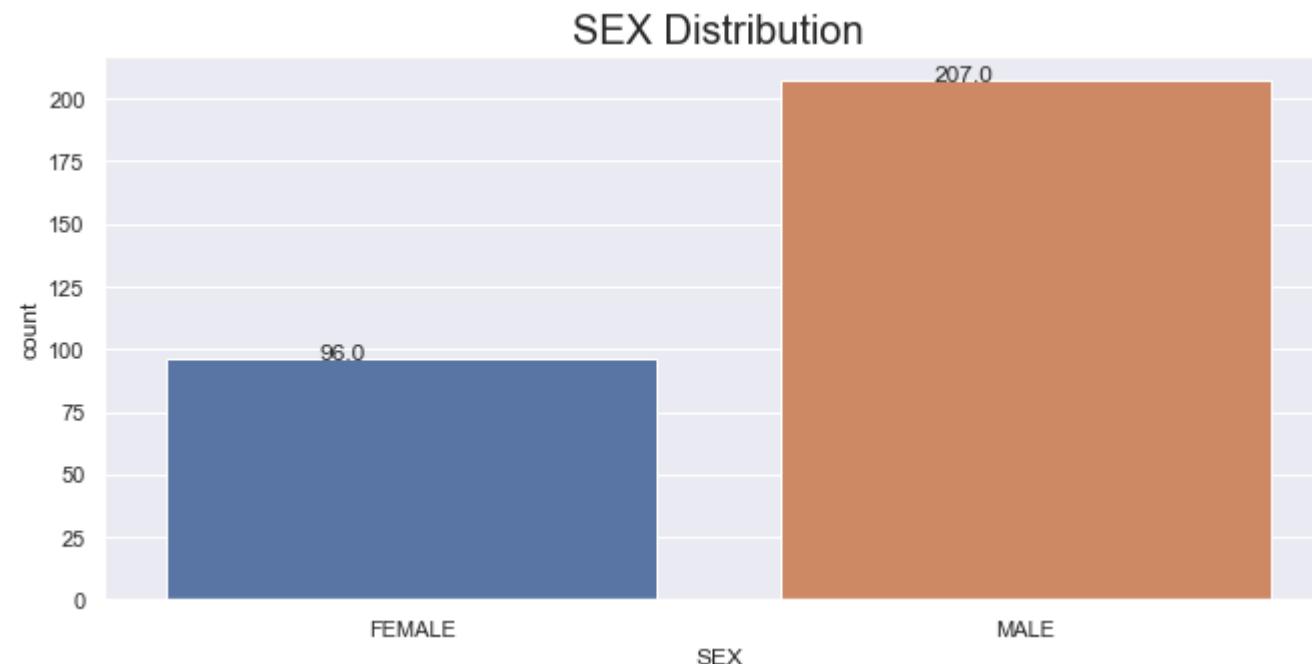
```
1 import seaborn as sb  
2 sb.boxplot(data=df['SEX'])
```

Out[435]: <AxesSubplot:>



In [436]:

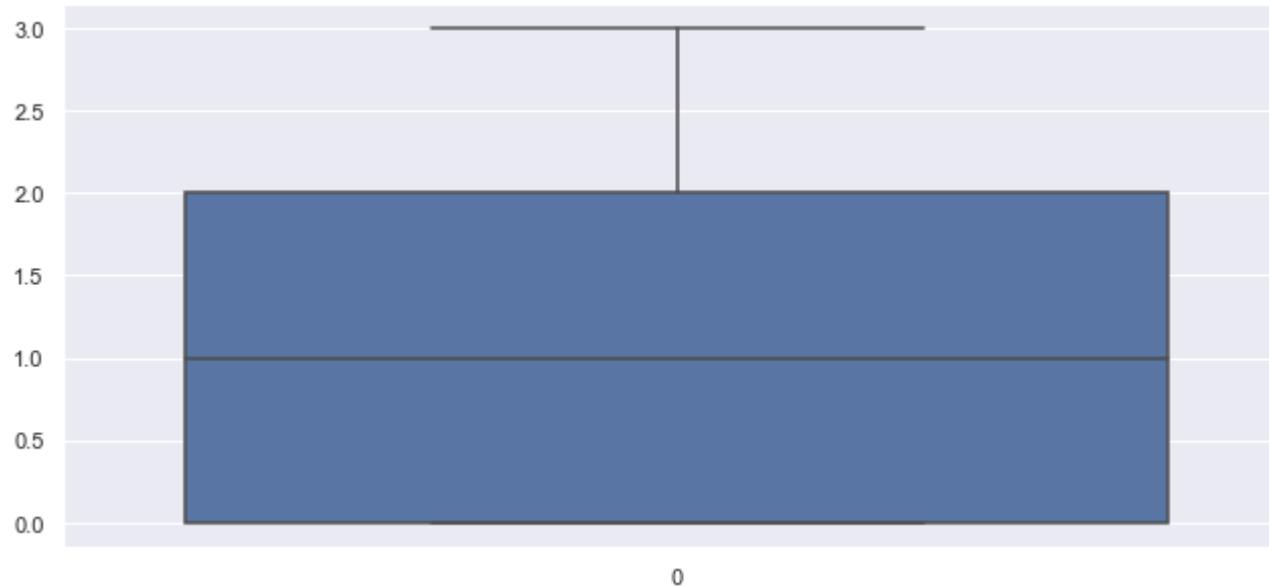
```
1 ax = sns.countplot(x="SEX", data=df)
2 ax.set_title('SEX Distribution', fontsize=20)
3 ax.set_xticklabels(['FEMALE', 'MALE'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7
8 for p in ax.patches:
9     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
10
```



In [437]:

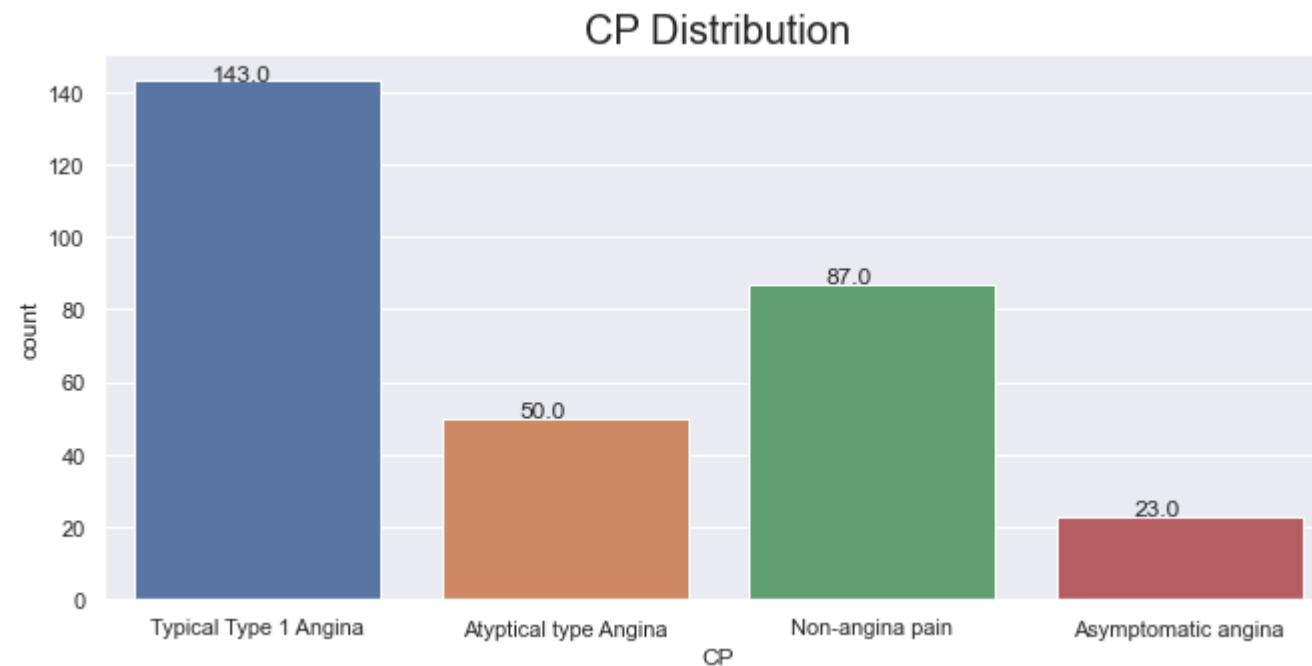
```
1 import seaborn as sb  
2 sb.boxplot(data=df['CP'])
```

Out[437]: <AxesSubplot:>



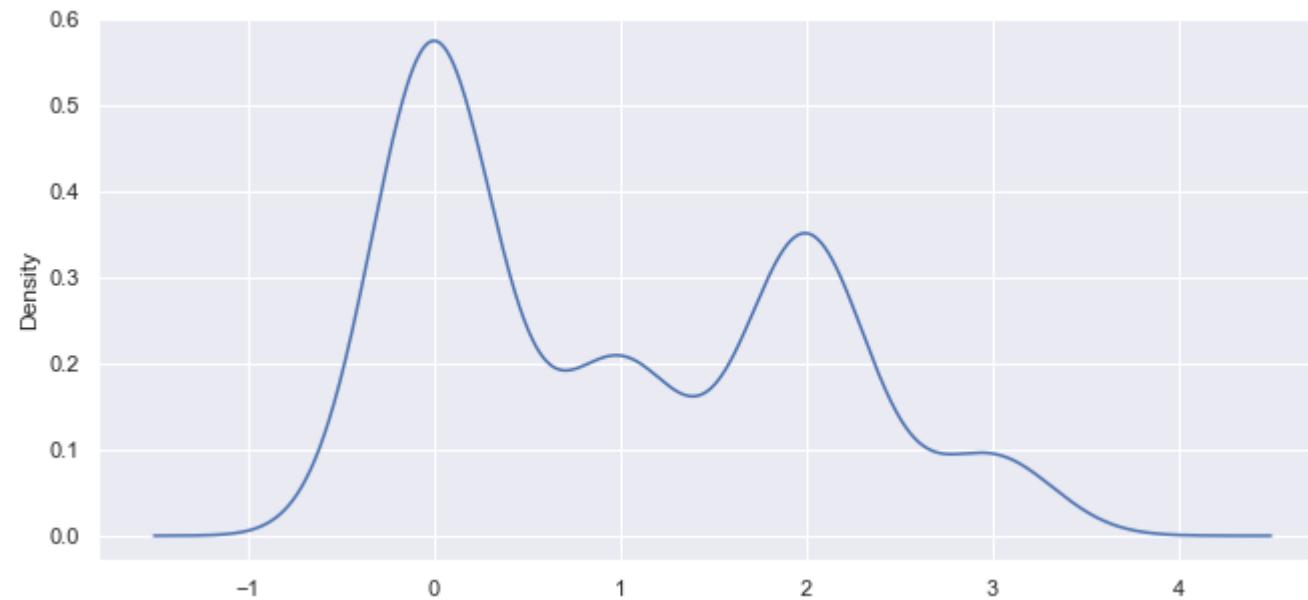
In [438]:

```
1 ax = sns.countplot(x="CP", data=df)
2 ax.set_title('CP Distribution', fontsize=20)
3 ax.set_xticklabels(['Typical Type 1 Angina','Atypical type Angina','Non-angina pain','Asymptomatic angina'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
```



```
In [439]: 1 df['CP'].plot(kind = 'density')
```

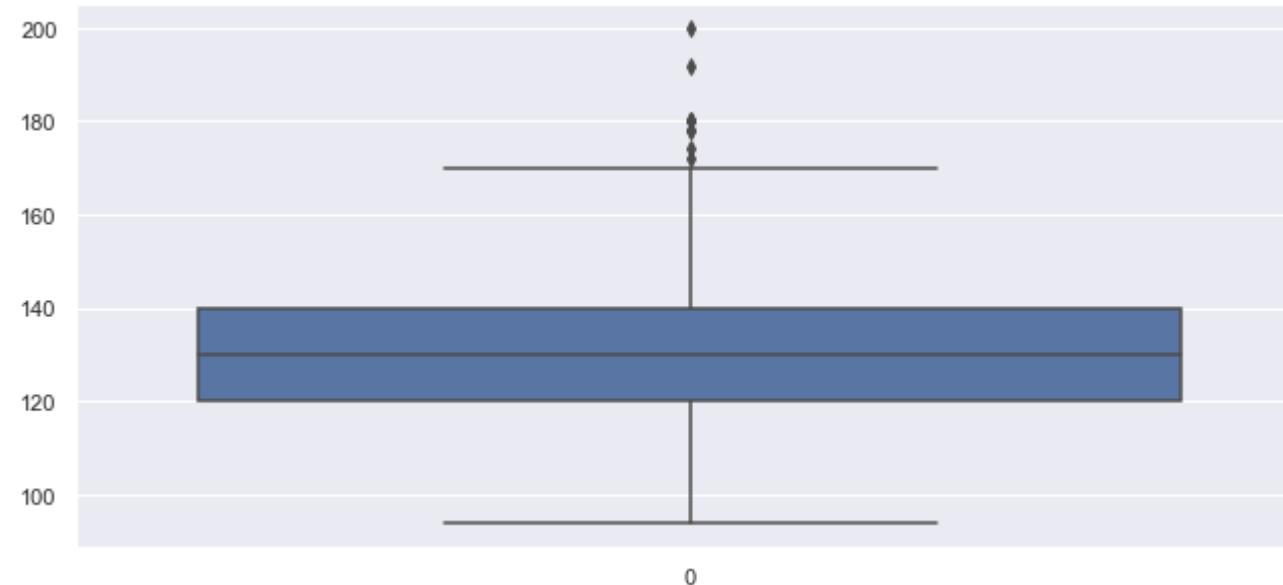
```
Out[439]: <AxesSubplot:ylabel='Density'>
```



In [440]:

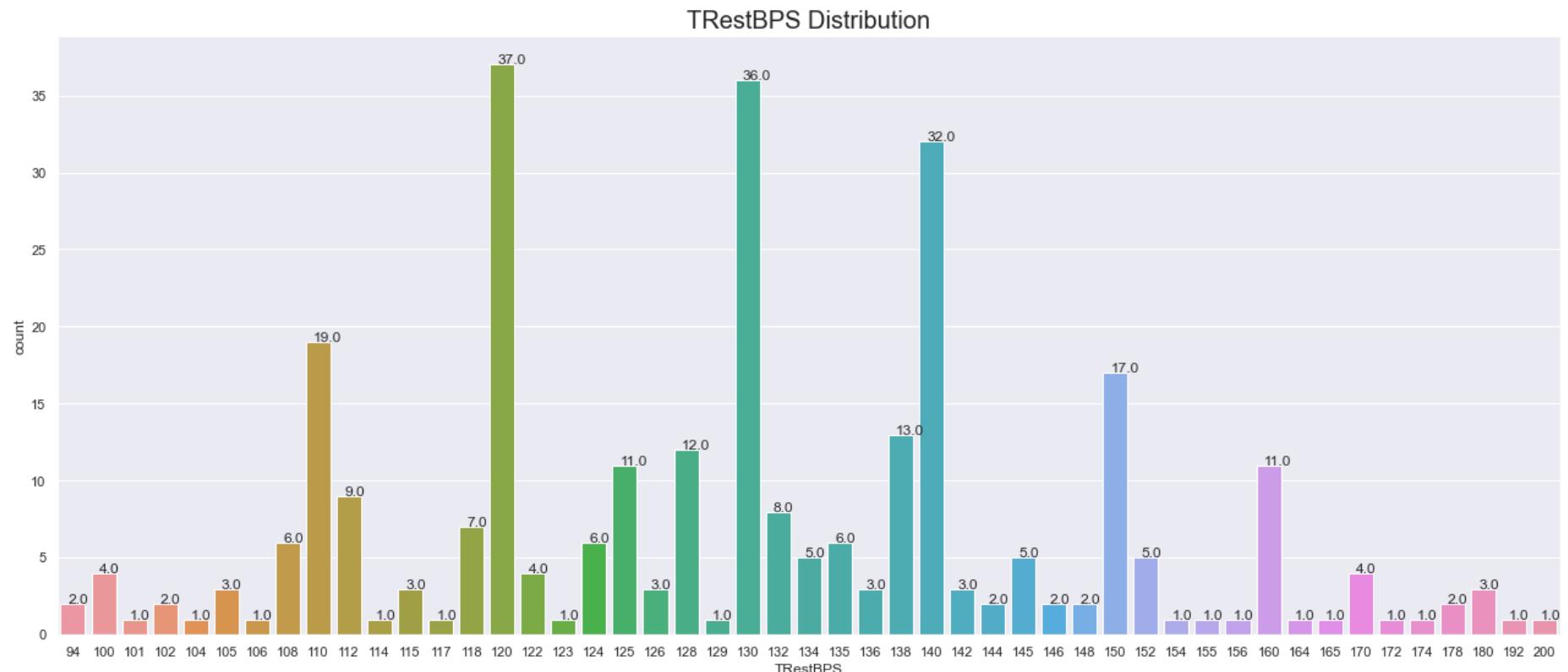
```
1 import seaborn as sb  
2 sb.boxplot(data=df['TRestBPS'])  
3  
4
```

Out[440]: <AxesSubplot:>



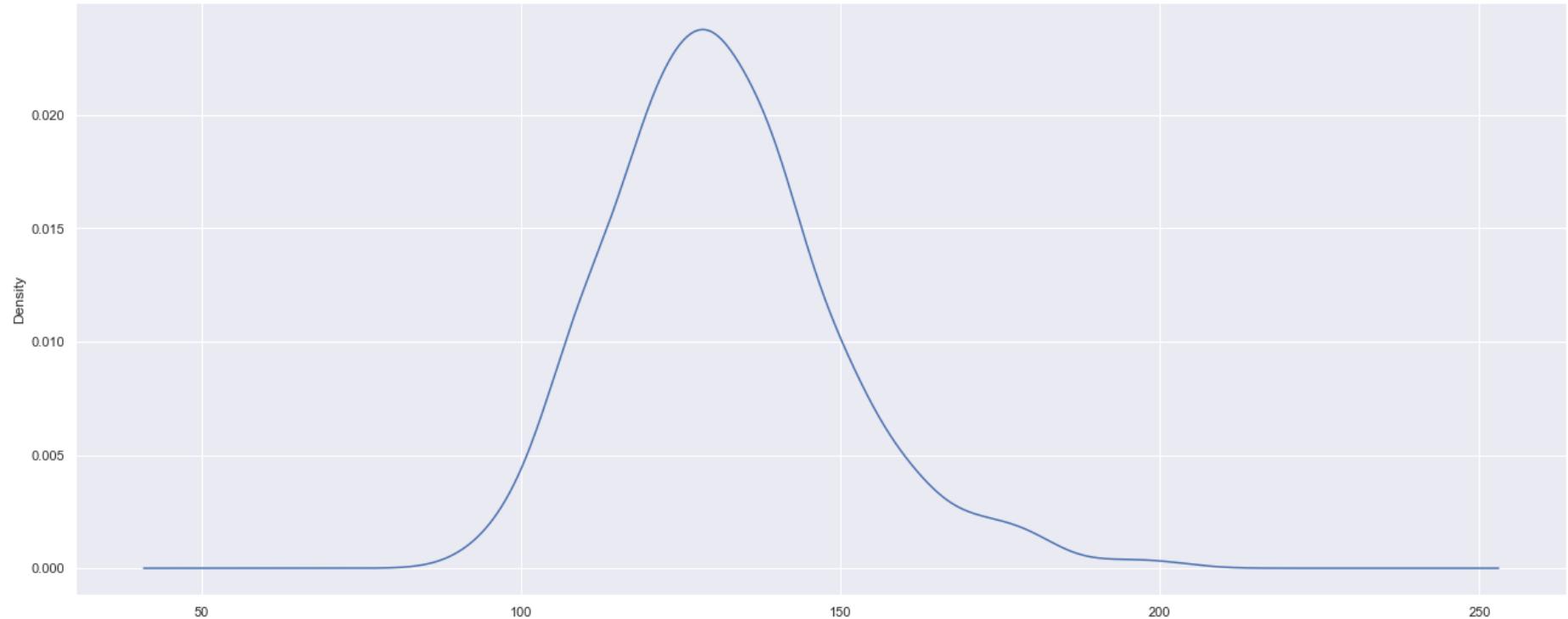
In [443]:

```
1 ax = sns.countplot(x="TRestBPS", data=df)
2 ax.set_title('TRestBPS Distribution', fontsize=20)
3 '''for p in ax.patches:
4     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
5
6 for p in ax.patches:
7     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
8 sns.set(rc={'figure.figsize':(22,9)})
9
```



```
In [444]: 1 df[ 'TRestBPS' ].plot(kind = 'density')
```

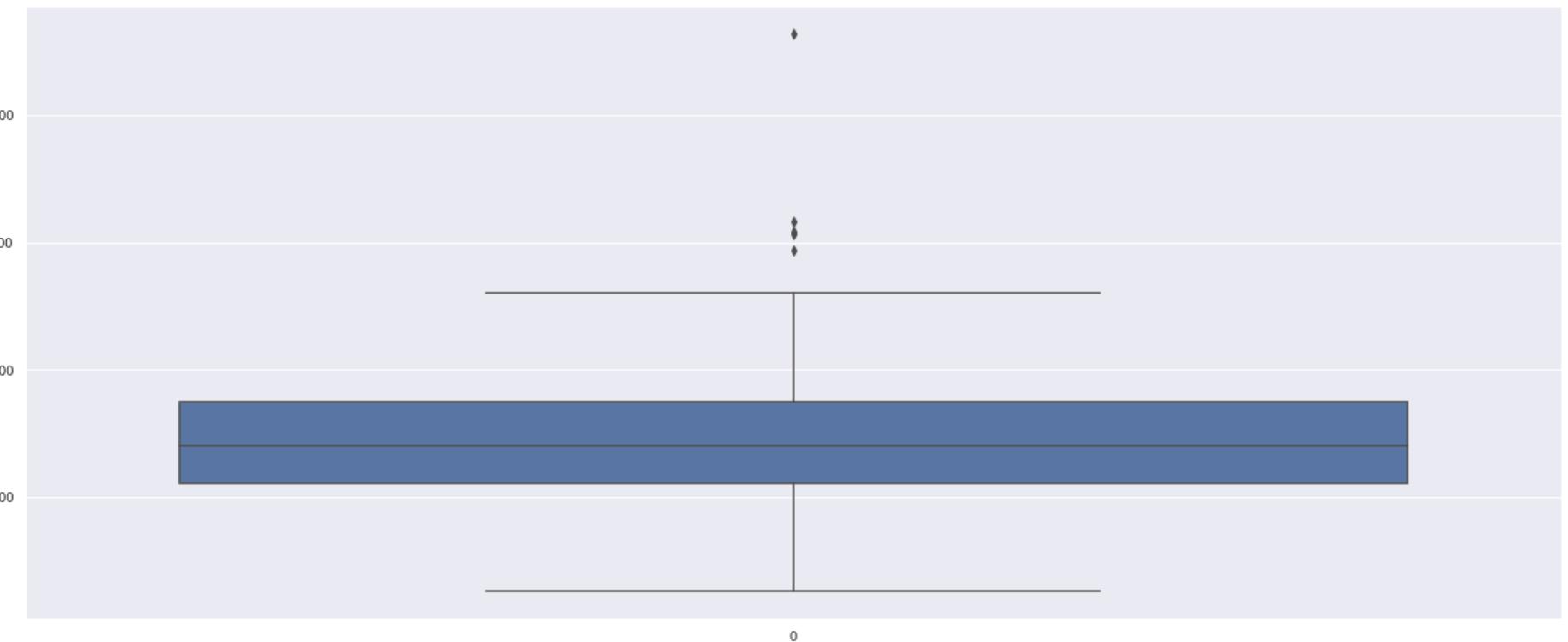
```
Out[444]: <AxesSubplot:ylabel='Density'>
```



In [445]:

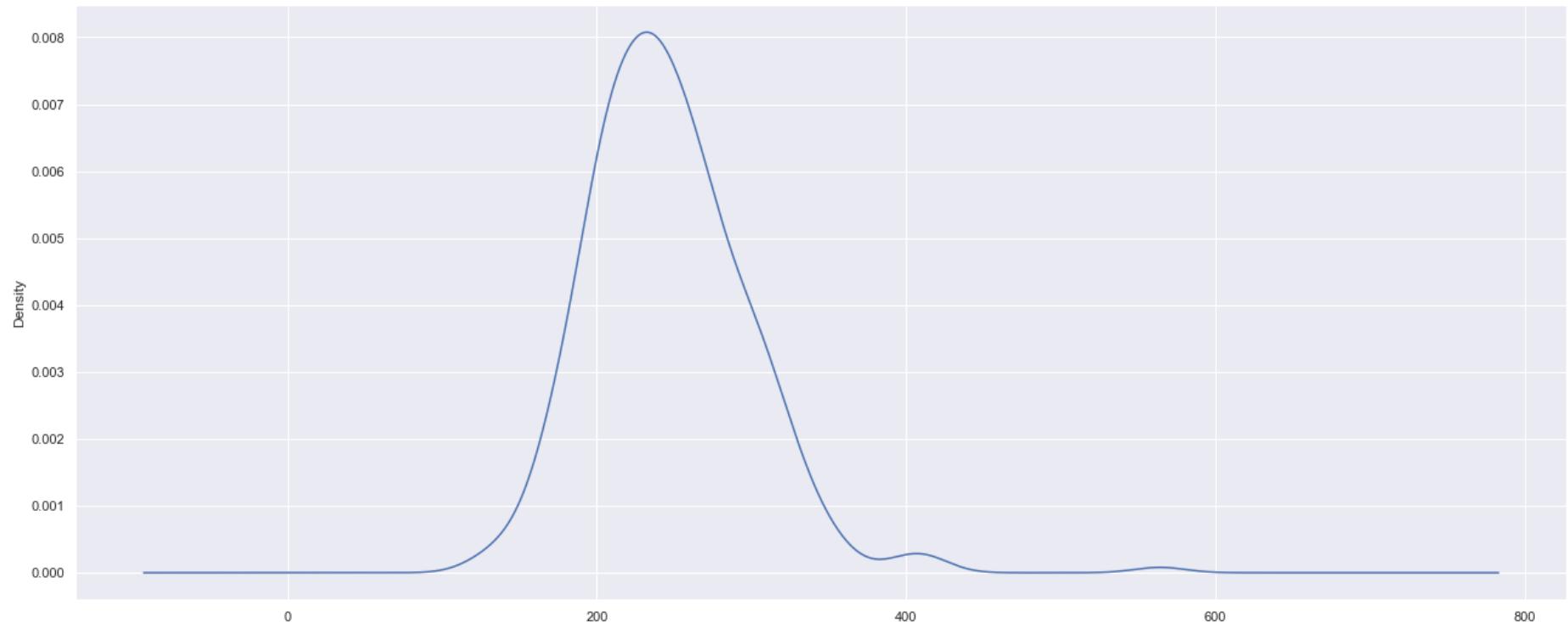
```
1 import seaborn as sb  
2 sb.boxplot(data=df['CHOL'])  
3  
4
```

Out[445]: <AxesSubplot:>



```
In [446]: 1 df[ 'CHOL' ].plot(kind = 'density')
```

```
Out[446]: <AxesSubplot:ylabel='Density'>
```

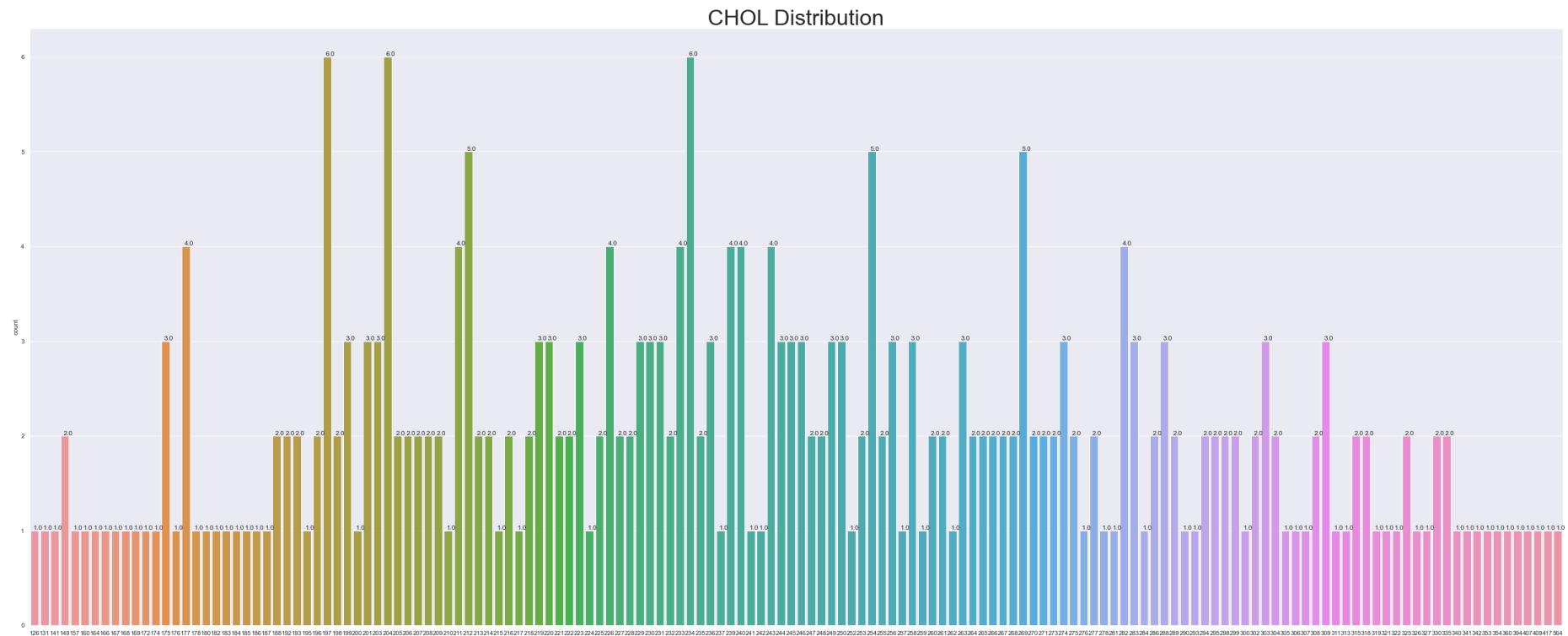


In [448]:

```

1 ax = sns.countplot(x="CHOL", data=df)
2 ax.set_title('CHOL Distribution', fontsize=40)
3 '''for p in ax.patches:
4     ax.annotate('%.1f' % p.get_height(), (p.get_x() + 0.1, p.get_height() + 50))'''
5
6 for p in ax.patches:
7     ax.annotate('%.1f' % p.get_height(), (p.get_x() + 0.25, p.get_height() + 0.01))
8 sns.set(rc={'figure.figsize':(50,20)})
9

```



In [449]:

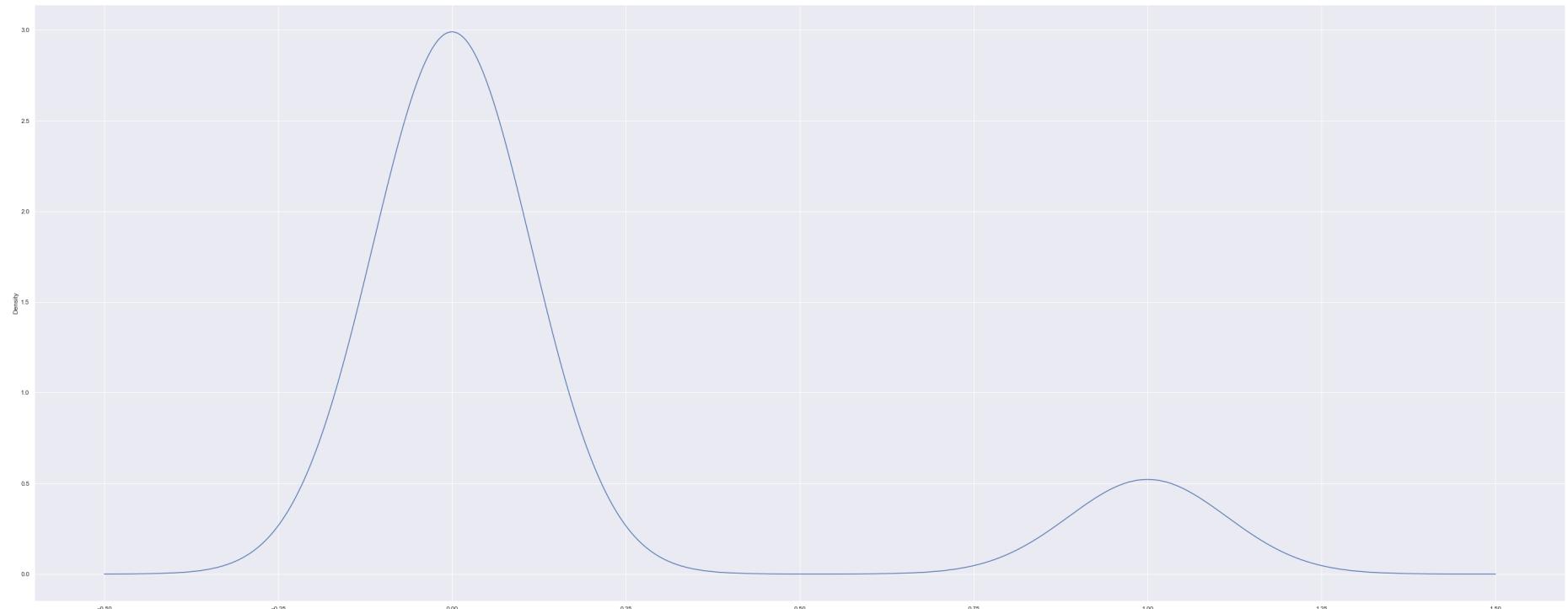
```
1 import seaborn as sb  
2 sb.boxplot(data=df['FBS'])
```

Out[449]: <AxesSubplot:>



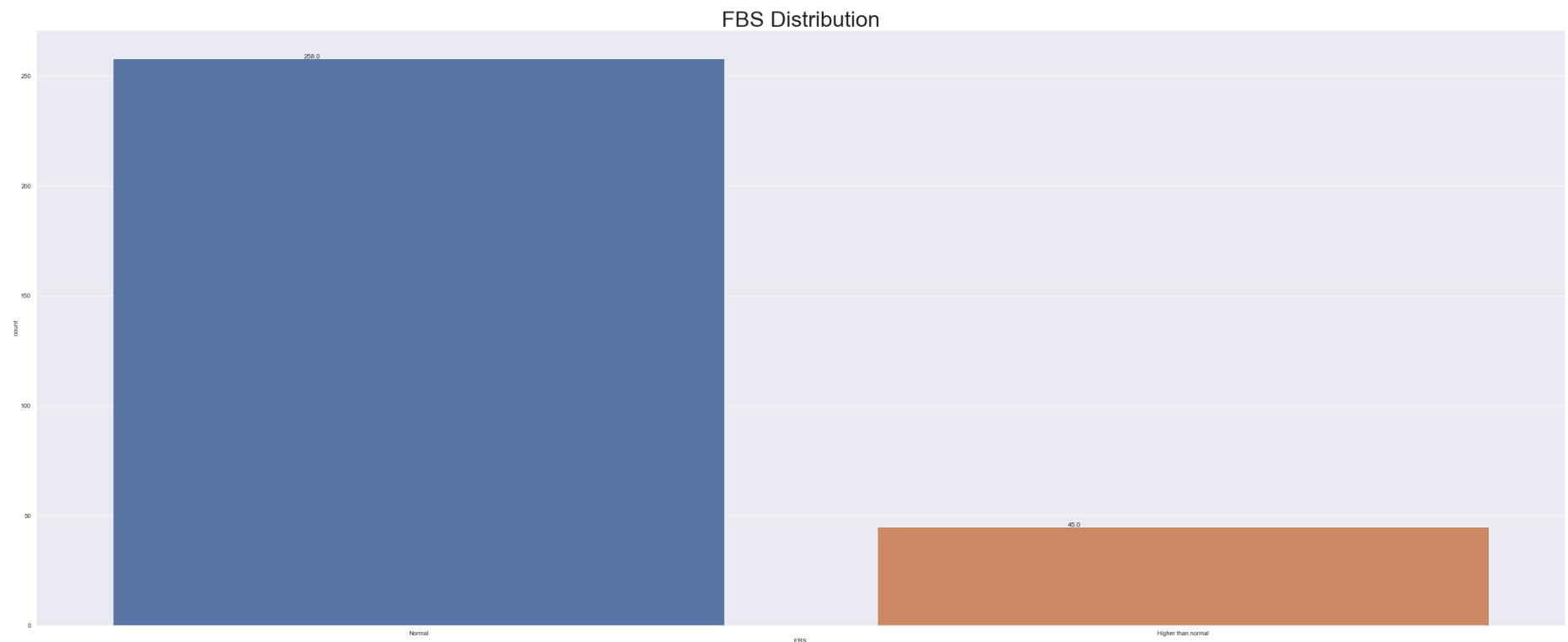
```
In [450]: 1 df[ 'FBS' ].plot(kind = 'density')
```

```
Out[450]: <AxesSubplot:ylabel='Density'>
```



In [455]:

```
1 ax = sns.countplot(x="FBS", data=df)
2 ax.set_title('FBS Distribution', fontsize=40)
3 ax.set_xticklabels(['Normal', 'Higher than normal'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7
8 for p in ax.patches:
9     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
10
```

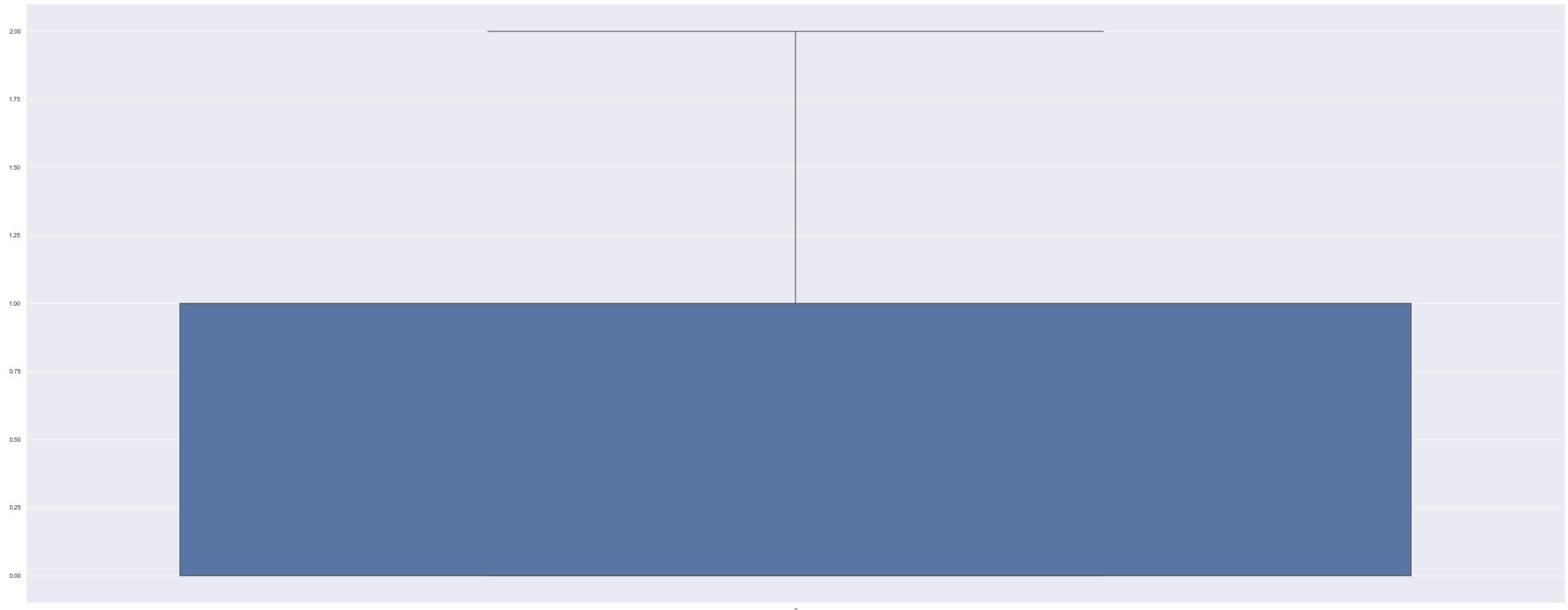


```
In [313]: 1 df.columns
```

```
Out[313]: Index(['AGE', 'SEX', 'CP', 'TRestBPS', 'CHOL', 'FBS', 'restECG', 'Thalach',
   'EXANG', 'OldPeak', 'SLOPE', 'CA', 'THAL', 'TARGET'],
  dtype='object')
```

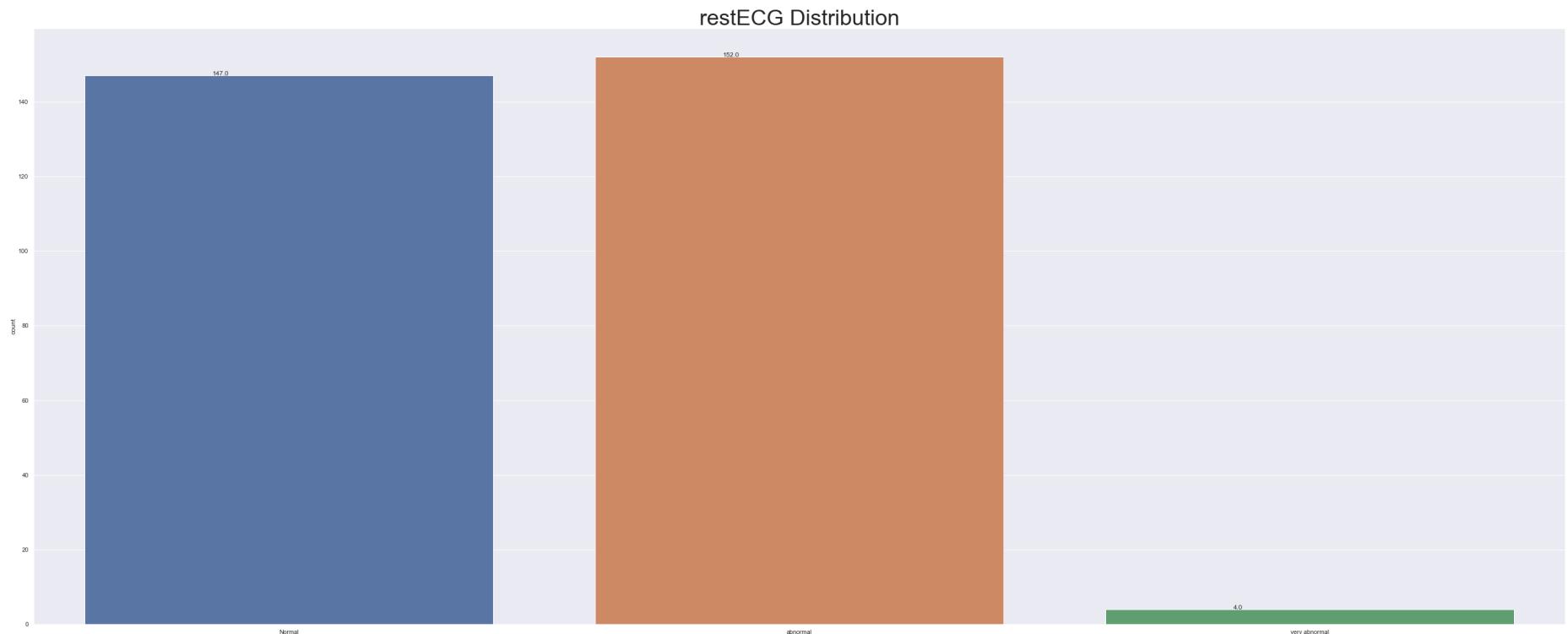
```
In [314]: 1 import seaborn as sb
2 sb.boxplot(data=df['restECG'])
```

```
Out[314]: <AxesSubplot:>
```



In [457]:

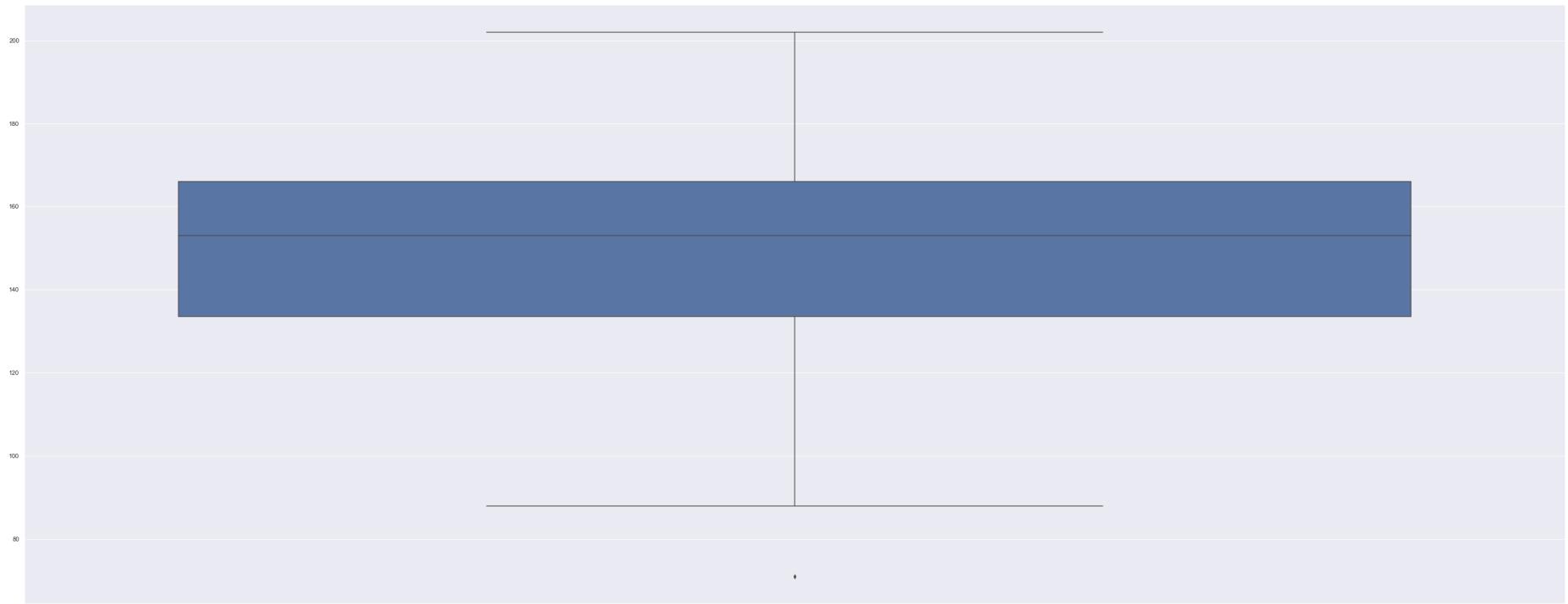
```
1 ax = sns.countplot(x="restECG", data=df)
2 ax.set_title('restECG Distribution', fontsize=40)
3 ax.set_xticklabels(['Normal', 'abnormal', 'very abnormal'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
```



In [458]:

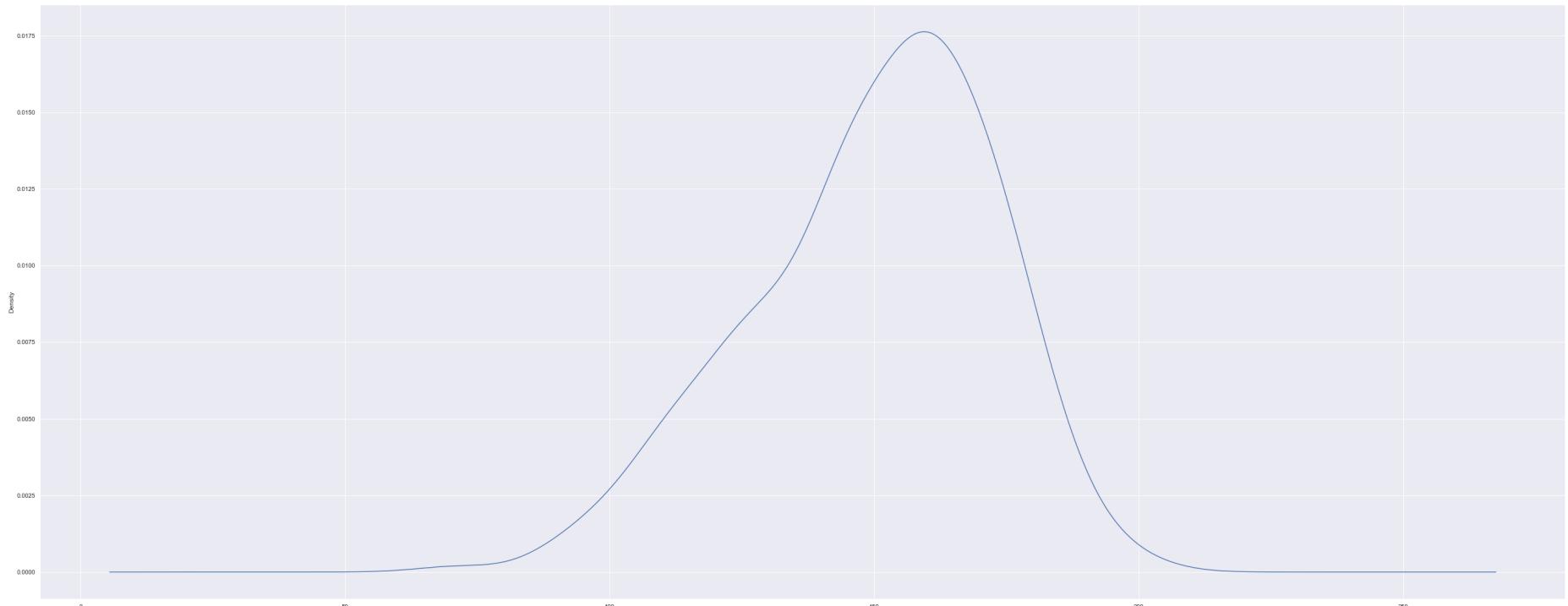
```
1 import seaborn as sb  
2 sb.boxplot(data=df['Thalach'])
```

Out[458]: <AxesSubplot:>



```
In [459]: 1 df['Thalach'].plot(kind = 'density')
```

```
Out[459]: <AxesSubplot:ylabel='Density'>
```

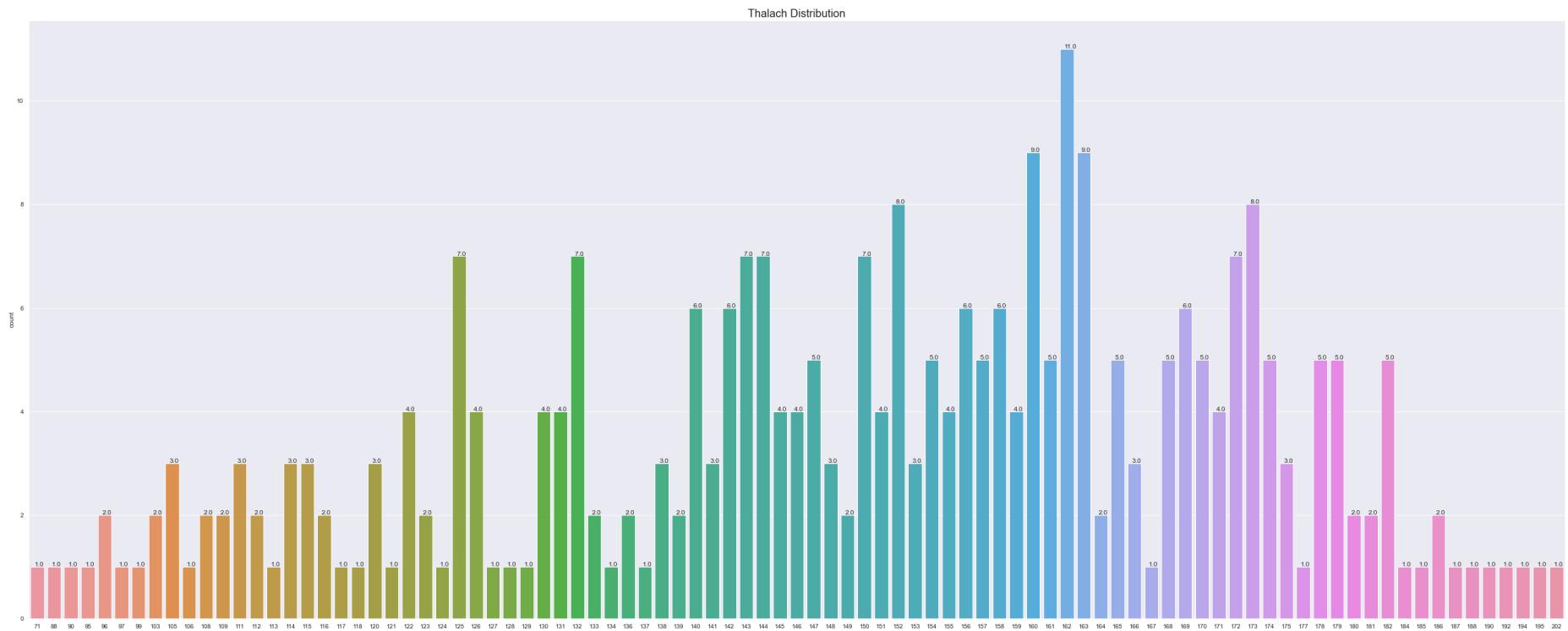


In [460]:

```

1 ax = sns.countplot(x="Thalach", data=df)
2 ax.set_title('Thalach Distribution', fontsize=20)
3 '''for p in ax.patches:
4     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
5
6 for p in ax.patches:
7     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
8

```



In [461]:

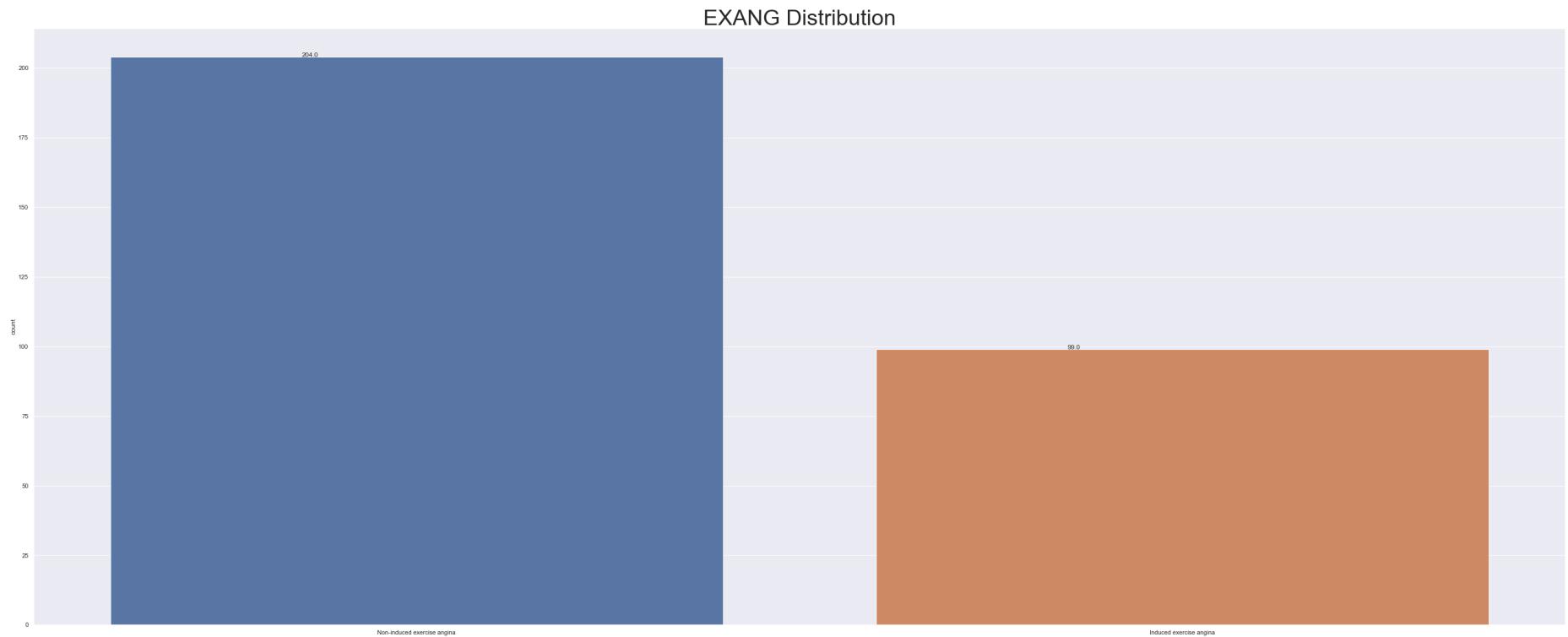
```
1 import seaborn as sb  
2 sb.boxplot(data=df['EXANG'])
```

Out[461]: <AxesSubplot:>



In [462]:

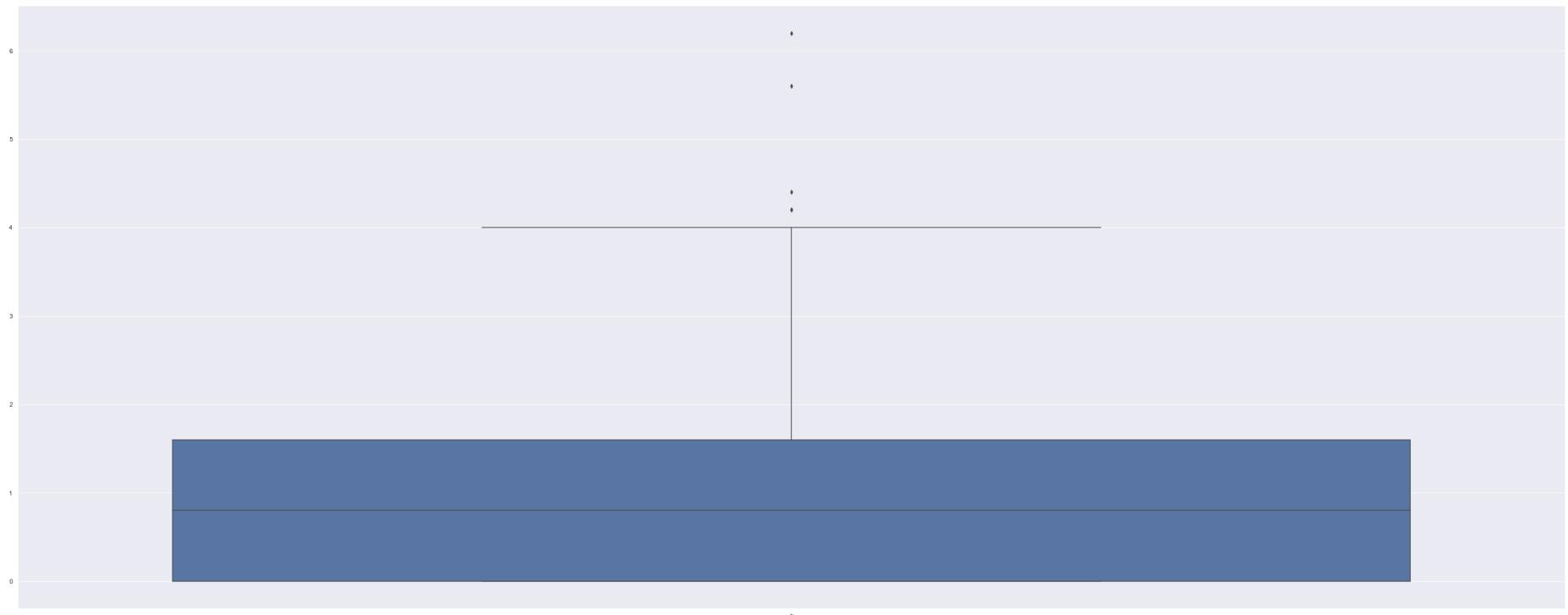
```
1 ax = sns.countplot(x="EXANG", data=df)
2 ax.set_title('EXANG Distribution', fontsize=40)
3 ax.set_xticklabels(['Non-induced exercise angina', 'Induced exercise angina'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
```



In [463]:

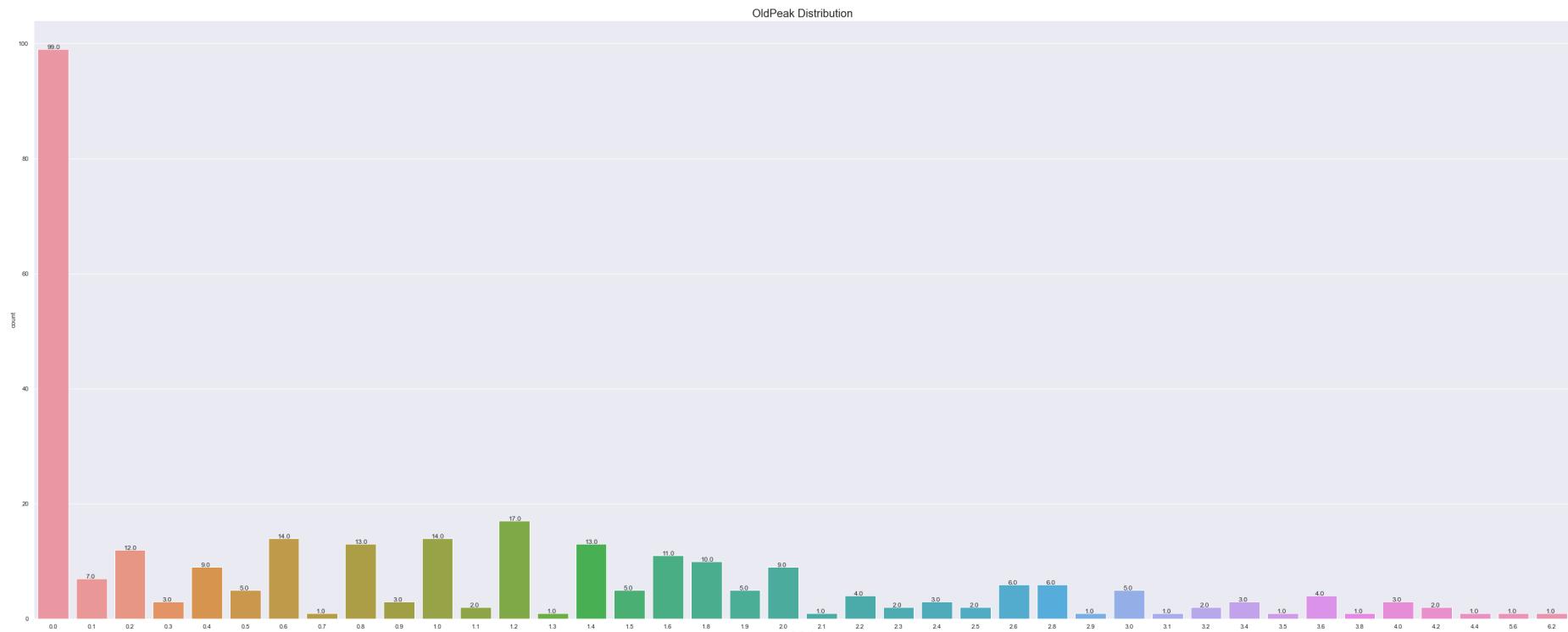
```
1 import seaborn as sb  
2 sb.boxplot(data=df['OldPeak'])
```

Out[463]: <AxesSubplot:>



In [464]:

```
1 ax = sns.countplot(x="OldPeak", data=df)
2 ax.set_title('OldPeak Distribution', fontsize=20)
3 '''for p in ax.patches:
4     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
5
6 for p in ax.patches:
7     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
8
```



In [465]: 1 df.head()

Out[465]:

	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

In [466]:

```
1 import seaborn as sb  
2 sb.boxplot(data=df['SLOPE'])
```

Out[466]: <AxesSubplot:>



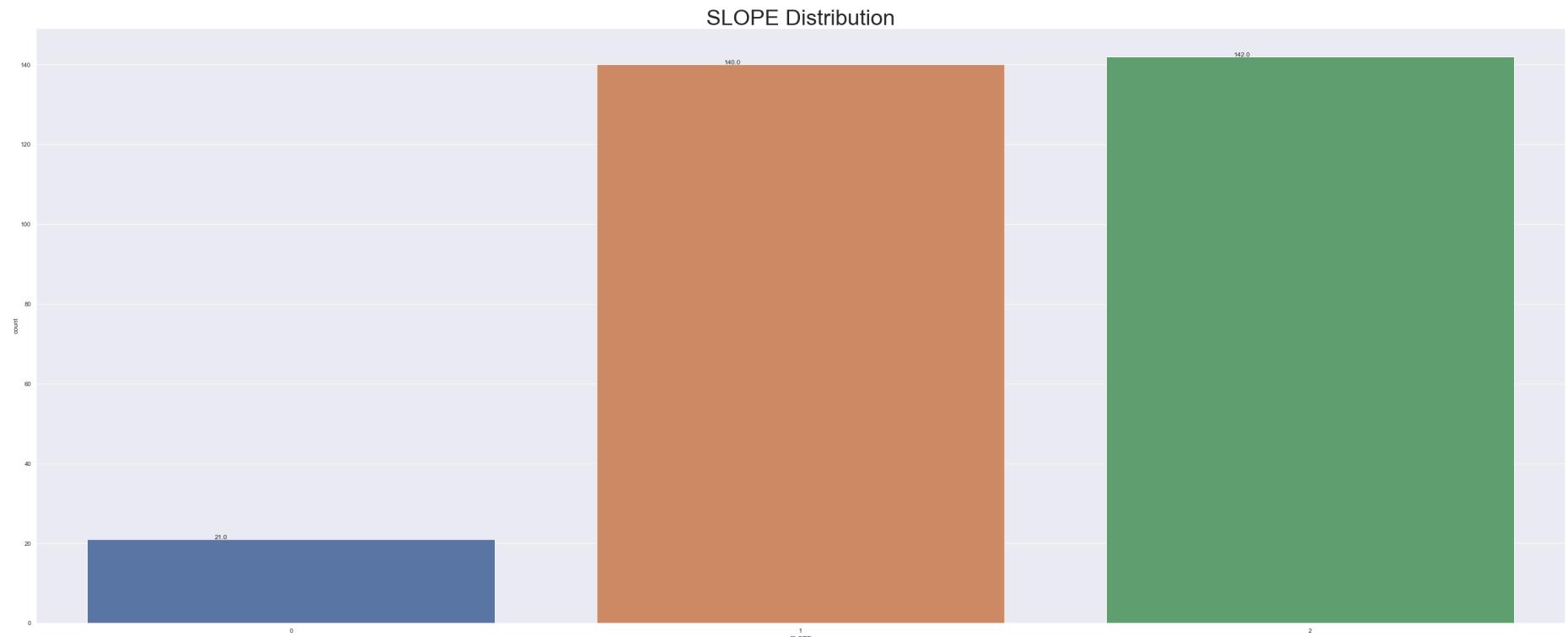
```
In [467]: 1 df[ 'SLOPE' ].plot(kind='density')
```

```
Out[467]: <AxesSubplot:ylabel='Density'>
```



In [469]:

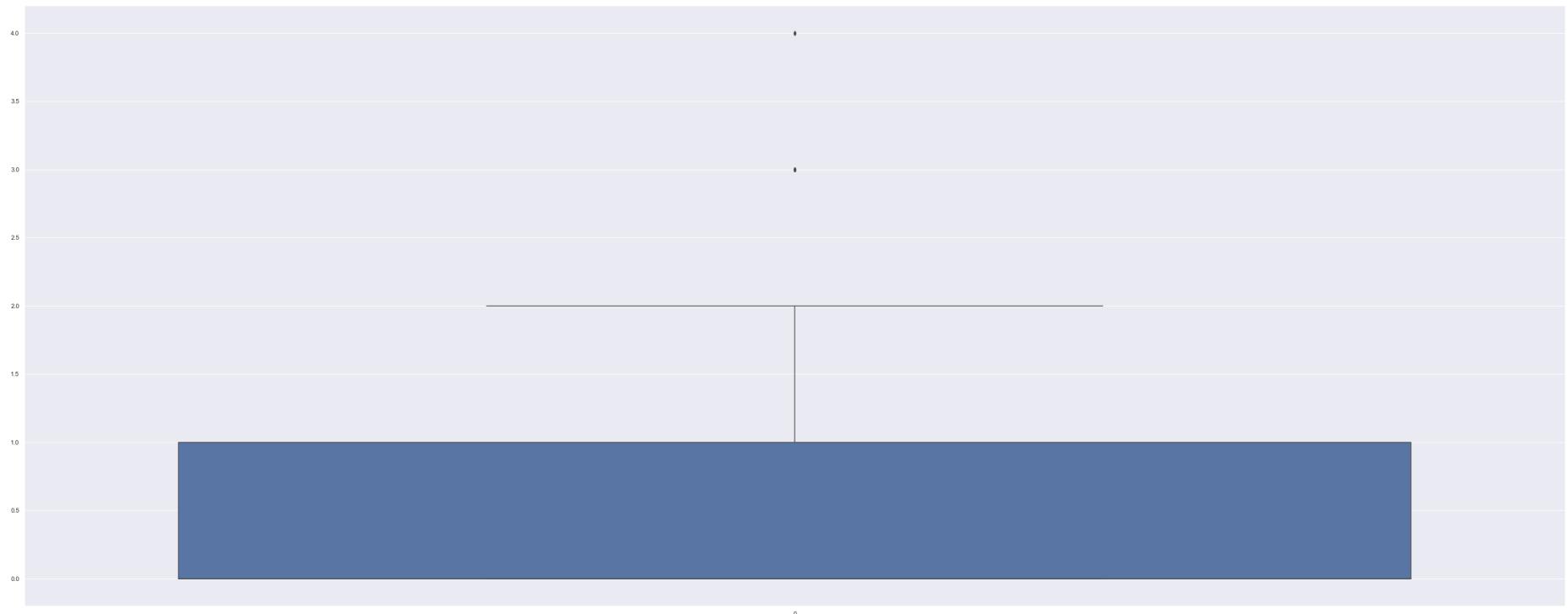
```
1 ax = sns.countplot(x="SLOPE", data=df)
2 ax.set_title('SLOPE Distribution', fontsize=40)
3 ax.set_xticklabels(['0','1','2'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
```



In [470]:

```
1 import seaborn as sb  
2 sb.boxplot(data=df['CA'])
```

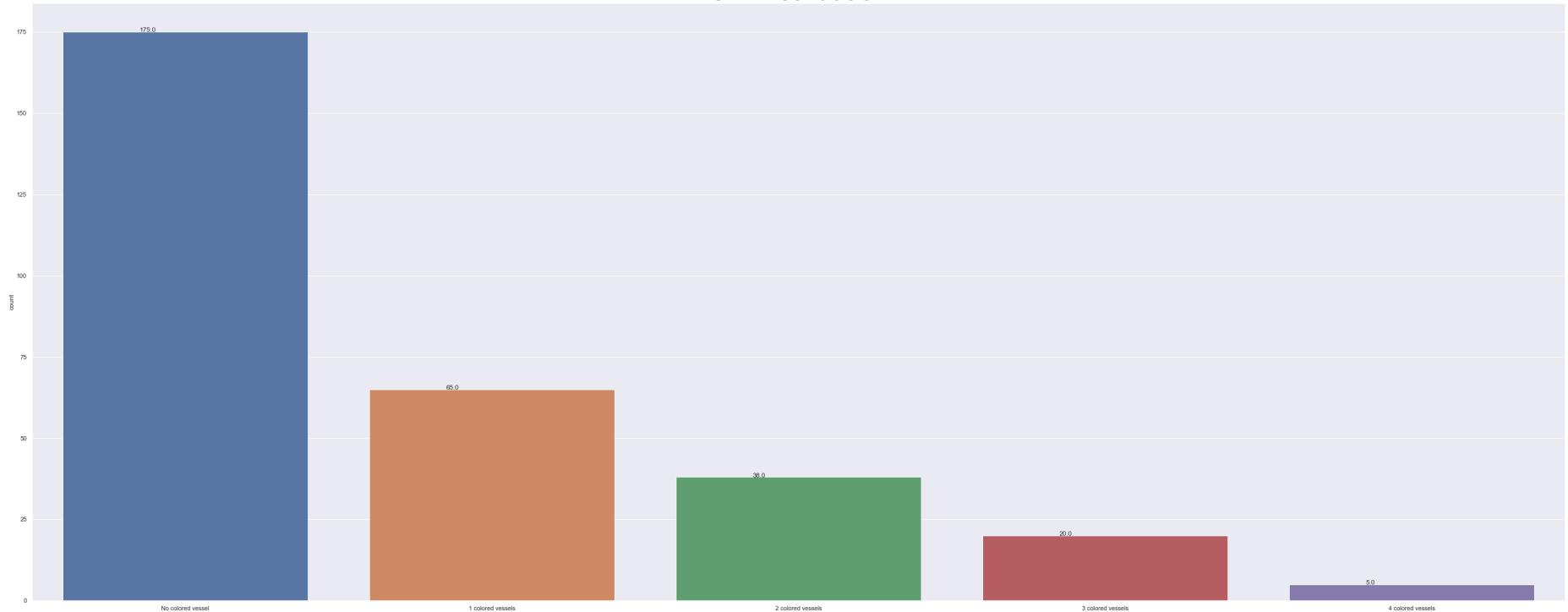
Out[470]: <AxesSubplot:>



In [475]:

```
1 ax = sns.countplot(x="CA", data=df)
2 ax.set_title('CA Distribution', fontsize=50)
3 ax.set_xticklabels(['No colored vessel', '1 colored vessels', '2 colored vessels', '3 colored vessels', '4 colored vessels'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7
8 for p in ax.patches:
9     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x() + 0.25, p.get_height() + 0.01))
```

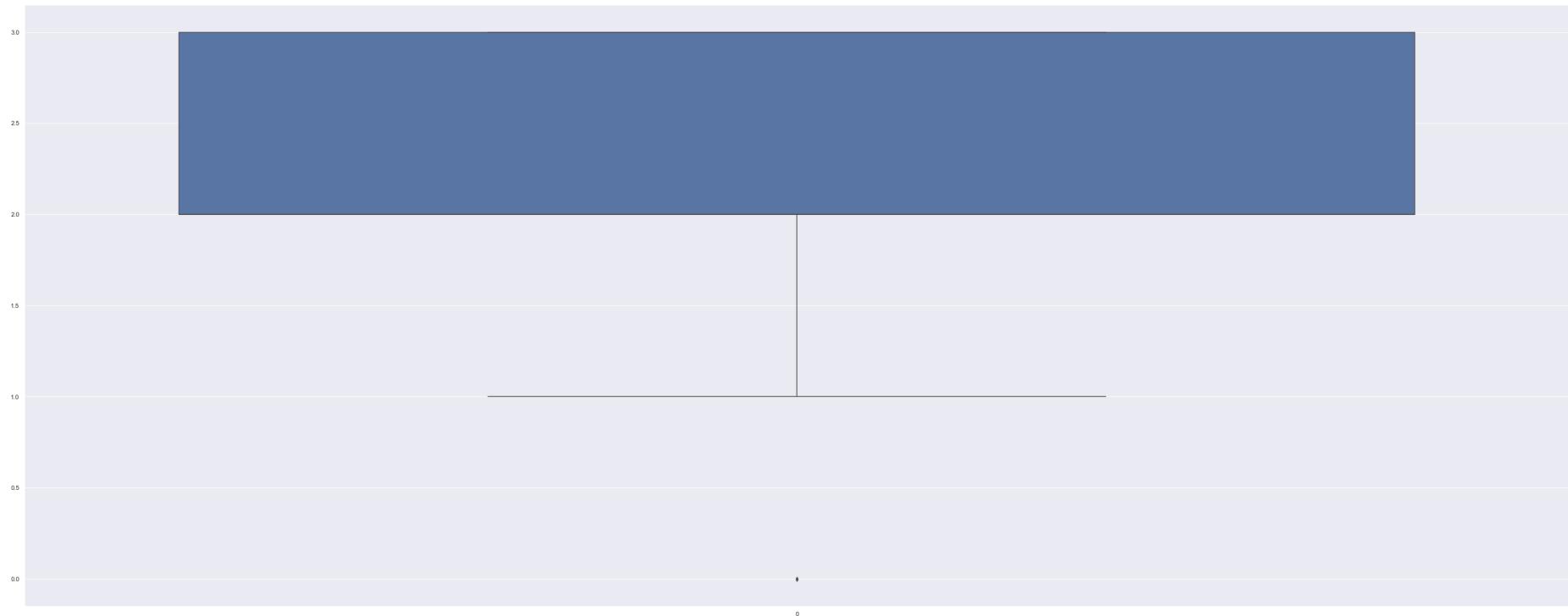
CA Distribution



In [476]:

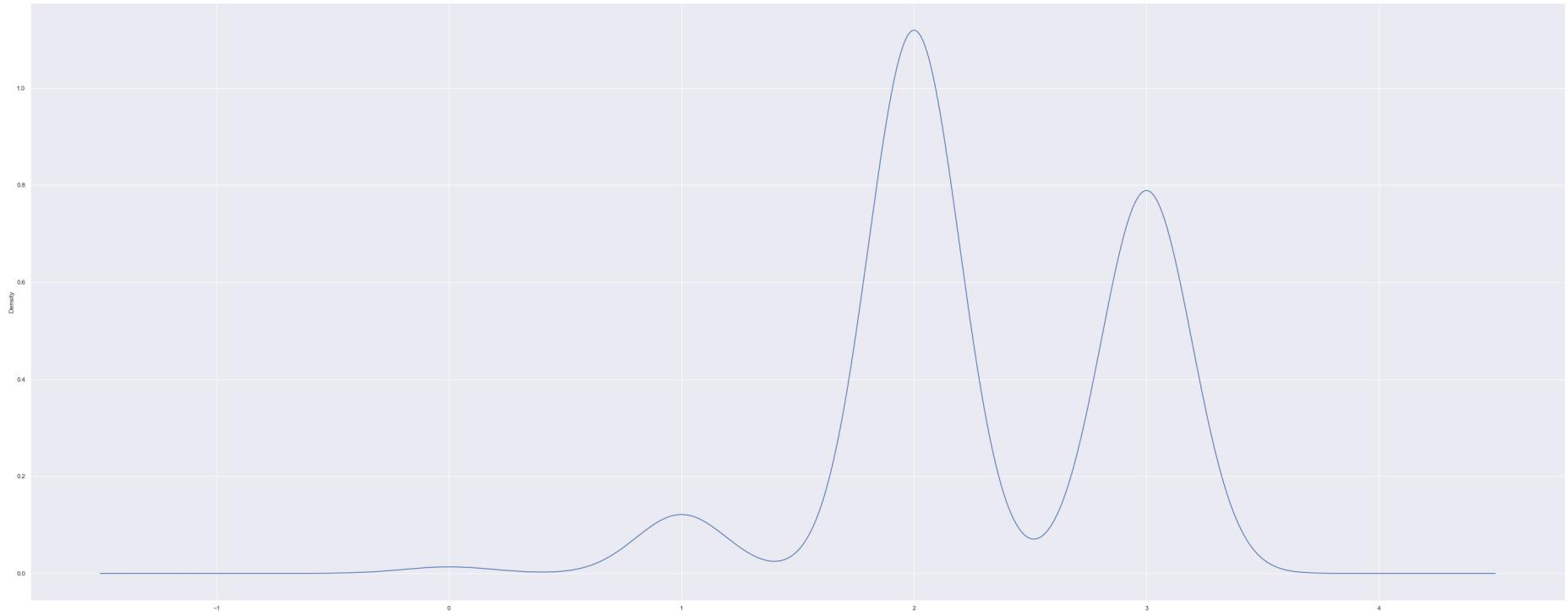
```
1 import seaborn as sb  
2 sb.boxplot(data=df['THAL'])
```

Out[476]: <AxesSubplot:>



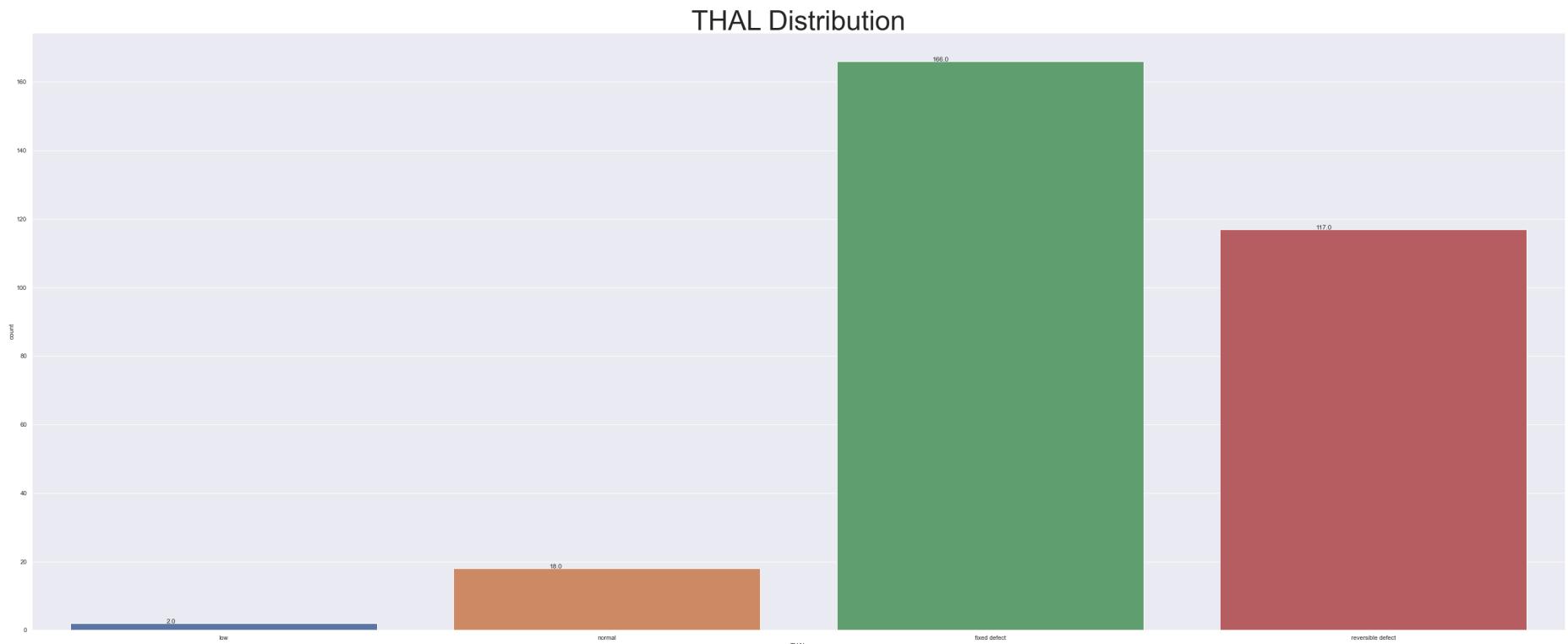
```
In [477]: 1 df[ 'THAL' ].plot(kind='density')
```

```
Out[477]: <AxesSubplot:ylabel='Density'>
```



In [478]:

```
1 ax = sns.countplot(x="THAL", data=df)
2 ax.set_title('THAL Distribution', fontsize=50)
3 ax.set_xticklabels(['low', 'normal','fixed defect','reversible defect'])
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in ax.patches:
8     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
9
```



In [479]:

```
1 ax = sns.countplot(x="TARGET", data=df)
2 ax.set_title('TARGET Distribution', fontsize=50)
3 '''for p in ax.patches:
4     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
5
6 for p in ax.patches:
7     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
8
9
```



In [480]: 1 df.head()

Out[480]:

	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

1. age
2. sex
3. chest pain type (4 values)
4. resting blood pressure
5. serum cholestoral in mg/dl
6. fasting blood sugar > 120 mg/dl
7. resting electrocardiographic results (values 0,1,2)
8. maximum heart rate achieved
9. exercise induced angina
10. oldpeak = ST depression induced by exercise relative to rest
11. the slope of the peak exercise ST segment
12. number of major vessels (0-3) colored by flourosopy
13. thal: 0 = normal; 1 = fixed defect; 2 = reversable defect

The names and social security numbers of the patients were recently removed from the database, replaced with dummy values.

S.No	Attribute Name	Description	Range of Values
1	Age	Age of the person in years	29 to 79
2	Sex	Gender of the person [1: Male, 0: Female]	0, 1
3	Cp	Chest pain type [1-Typical Type 1 Angina 2- Atypical Type Angina 3-Non-angina pain 4-Asymptomatic)	1, 2, 3, 4
4	Trestbps	Resting Blood Pressure in mm Hg	94 to 200
5	Chol	Serum cholesterol in mg/dl	126 to 564
6	Fbs	Fasting Blood Sugar in mg/dl	0, 1
7	Restecg	Resting Electrocardiographic Results	0, 1, 2
8	Thalach	Maximum Heart Rate Achieved	71 to 202
9	Exang	Exercise Induced Angina	0, 1
10	OldPeak	ST depression induced by exercise relative to rest	1 to 3
11	Slope	Slope of the Peak Exercise ST segment	1, 2, 3
12	Ca	Number of major vessels colored by fluoroscopy	0 to 3
13	Thal	3 – Normal, 6 – Fixed Defect, 7 –	3, 6, 7

14	Num	reversible defect	Class Attribute	0 or 1
----	-----	-------------------	-----------------	--------

1. age
2. sex
3. chest pain type (4 values)
4. resting blood pressure
5. serum cholestorol in mg/dl
6. fasting blood sugar > 120 mg/dl
7. resting electrocardiographic results (values 0,1,2)
8. maximum heart rate achieved
9. exercise induced angina
10. oldpeak = ST depression induced by exercise relative to rest
11. the slope of the peak exercise ST segment
12. number of major vessels (0-3) colored by flourosopy
13. that: 0 = normal; 1 = fixed defect; 2 = reversible defect

The names and social security numbers of the patients were recently removed from the database, replaced with dummy values.

```
In [481]: 1 minAge = min(df['AGE'])
2 print(minAge)
3 maxAge = max(df['AGE'])
4 print(maxAge)
5 '''29-40 = 0
6     41-52 = 1
7     53-64 = 2
8     64-77 = 3'''
```

29
77

Out[481]: '29-40 = 0\n 41-52 = 1\n 53-64 = 2\n 64-77 = 3'

```
In [482]: 1 df.AGE.unique()
```

Out[482]: 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])

```
In [483]: 1 fill = (df.AGE == 29) | (df.AGE == 34) | (df.AGE == 35) | (df.AGE == 38) | (df.AGE == 39) | (df.AGE == 40) | (df
2 df.loc[fill, 'AGE'] = 0
```

```
In [484]: 1 fill = (df.AGE == 41) | (df.AGE == 42) | (df.AGE == 44) | (df.AGE == 43) | (df.AGE == 50) | (df.AGE == 45) | (df
2 df.loc[fill, 'AGE'] = 1
```

```
In [485]: 1 fill = (df.AGE == 53) | (df.AGE == 54) | (df.AGE == 55) | (df.AGE == 57) | (df.AGE == 59) | (df.AGE == 56) | (df.AG
2 df.loc[fill, 'AGE'] = 2
```

```
In [486]: 1 fill = (df.AGE == 65) | (df.AGE == 66) | (df.AGE == 67) | (df.AGE == 68) | (df.AGE == 69) | (df.AGE == 70) | (df.AG
2 df.loc[fill, 'AGE'] = 3
```

In [487]:

```
1 minBP = min(df[ 'TRestBPS' ])
2 print(minBP)
3
4 maxBP = max(df[ 'TRestBPS' ])
5 print(maxBP)
6
7 '''90-119 = Normal
8     120-139 = Borderline HBP
9     140-200 = High BP'''
```

```
94
200
```

```
Out[487]: '90-119 = Normal\n    120-139 = Borderline HBP\n    140-200 = High BP'
```

In [488]:

```
1 fill = (df.TRestBPS == 94) | (df.TRestBPS == 110) | (df.TRestBPS == 104) | (df.TRestBPS == 105) | (df.TRestBPS == 106)
2 df.loc[fill, 'TRestBPS'] = 0
```

In [489]:

```
1 fill = (df.TRestBPS == 120) | (df.TRestBPS == 130) | (df.TRestBPS == 135) | (df.TRestBPS == 125) | (df.TRestBPS == 120)
2 df.loc[fill, 'TRestBPS'] = 1
```

In [490]:

```
1 fill = (df.TRestBPS == 140) | (df.TRestBPS == 145) | (df.TRestBPS == 172) | (df.TRestBPS == 150) | (df.TRestBPS == 155)
2 df.loc[fill, 'TRestBPS'] = 2
```

In [491]:

```
1 df.TRestBPS.unique()
```

```
Out[491]: array([2, 1, 0])
```

In []:

```
1
```

In [492]:

```
1 minChol = min(df[ 'CHOL' ])
2 print(minChol)
3 maxChol = max(df[ 'CHOL' ])
4 print(maxChol)
5
6 '''126-200 = Normal
7     201-239 = Borderline high
8     240-564 = high'''
```

126
564

Out[492]: '126-200 = Normal\n 201-239 = Borderline high\n 240-564 = high'

In [493]:

```
1 fill = (df.CHOL == 126) | (df.CHOL == 192) | (df.CHOL == 199) | (df.CHOL == 168) | (df.CHOL == 175) | (df.CHOL ==
2 df.loc[fill, 'CHOL'] = 0
```

In [494]:

```
== 203) | (df.CHOL == 206) | (df.CHOL == 209) | (df.CHOL == 217) | (df.CHOL == 218) | (df.CHOL == 224) | (df.CHOL == 2
DL'] = 1
```

In [495]:

```
1 fill = (df.CHOL == 241) | (df.CHOL == 242) | (df.CHOL == 244) | (df.CHOL == 294) | (df.CHOL == 246) | (df.C
2 df.loc[fill, 'CHOL'] = 2
```

In []:

1

In [496]:

```
1 sorted(df.CHOL.unique())
```

Out[496]: [0, 1, 2]

In [497]:

```

1 minThalach = min(df['Thalach'])
2 print(minThalach)
3
4 maxThalach = max(df['Thalach'])
5 print(maxThalach)
6
7 '''71-100 = low Thalach
8 101-130 = normal Thalach
9 131-160 = borderline high Thalach
10 161-202 = high Thalach'''
11

```

71
202

Out[497]: '71-100 = low Thalach\n101-130 = normal Thalach\n131-160 = borderline high Thalach\n161-202 = high Thalach'

In [498]:

```

1 fill = (df.Thalach == 71) | (df.Thalach == 95) | (df.Thalach == 88) | (df.Thalach == 97) | (df.Thalach == 99) |
2 df.loc[fill, 'Thalach'] = 0

```

In [499]:

```

1 fill = (df.Thalach == 106) | (df.Thalach == 103) | (df.Thalach == 127) | (df.Thalach == 105) | (df.Thalach ==
2 df.loc[fill, 'Thalach'] = 1

```

In [500]:

```

1 fill = (df.Thalach == 131) | (df.Thalach == 132) | (df.Thalach == 133) | (df.Thalach == 134) | (df.Thalach ==
2 df.loc[fill, 'Thalach'] = 2

```

In [501]:

```

1 fill = (df.Thalach == 161) | (df.Thalach == 162) | (df.Thalach == 163) | (df.Thalach == 164) | (df.Thalach == 16
2 df.loc[fill, 'Thalach'] = 3

```

In [502]:

```
1 sorted(df.Thalach.unique())
```

Out[502]: [0, 1, 2, 3]

In [503]:

```
1 df.columns
```

Out[503]: Index(['AGE', 'SEX', 'CP', 'TRestBPS', 'CHOL', 'FBS', 'restECG', 'Thalach',
'EXANG', 'OldPeak', 'SLOPE', 'CA', 'THAL', 'TARGET'],
dtype='object')

In [504]:

```
1 #BIVARATE ANALYSIS
2
3
4 '''I will be plotting distributions between these columns and the target column:
5 'AGE', 'SEX', 'CP', 'TRestBPS', 'CHOL', 'FBS', 'restECG', 'Thalach',
6 'EXANG', 'OldPeak', 'SLOPE', 'CA', 'THAL' '''
7
8 #For AGE,
9 '''we have a majority of the younger adults (A total of 84 out of 120 adults,
10 between 29 and 52; around 70%)
11 have a diseased heart as we don't expect whereas,
12 the reverse is the case for the older adults (A total of 81 out of 183;
13 between ages 53 and 77 around 44.2% ), have a diseased heart
14 Based on this outcome, heart disease is more common in younger adults.'''
15
16 #SEX
17 '''This study was carried out on fewer females than men but yet
18 there is a higher percentage of heart disease in females than in males which is very much significant.
19 This amounts to like 7 out of 9 (75% of the women) and 9 out of 20 (45% of the men) respectively
20 I believe there is a relationship between these two, we could infer that
21 females are likely more prone to a heart disease than males for whatever reason.
22 Studies have shown that women have smaller arteries than men, this could be a reason'''
23
24 #CP
25 '''The proportion of those with a diseased heart who also have one form of chest pain or the other
26 is higher than those with a healthy heart but it appears that
27 even those who do not have pain in their chest(are asymptomatic) and those who had pain not angina-induced
28 were considered to have a heart disease. Therefore, There is a relationship between these two columns'''
29
30 '''For TRestBPS,
31 we see that elevated blood pressure has no influence on heart disease development.
32 Even individuals who have normal blood pressure tend to have a higher proportion of heart disease development.
33 Safe to say that high resting blood pressure has NO influence on heart disease development.'''
34
35
36 '''For CHOL,
37 There is almost an equal proportion of those who have a diseased heart
38 who also have normal or borderline high levels of blood cholesterol.
39 Even a 1:1 ratio of diseased and healthy heart in those who have higher blood cholesterol levels.
40 No increase in percentage of those with a diseased heart from normal to high cholesterol levels
41 and vice versa means no relationship between these two features'''
```

```
42
43     '''FBS
44 The number of those with high blood sugar way less than those with normal blood sugar.
45 For the relationship between FBS and target, proportion of those with a diseased heart
46 who also have a normal blood sugar level is even higher than those with high blood sugar levels
47 such that we cannot tell if individuals with higher than normal blood sugar levels
48 or with normal blood sugar levels will more likely have a diseased heart.
49 For me, I do not think there is a relationship between these two
50 and having diabetes is not a yardstick for having a heart disease.'''
51
52 #For restECG,
53 '''Proportion of those with a diseased heart increases (from like 46.2% to 63%)
54 as we go from normal to abnormal restECG conditions.
55 If only we had more people with very abnormal ECG conditions, we may likely see an increase too
56 I think there is a relationship between these two columns
57 and restECG will be an important feature for the model design'''
58
59 #For Thalach,
60 '''Proportion of those with a diseased heart begins to increase with increase in their maximum heart rates.
61 I think there is a relationship between these two columns
62 and Thalach will be an important feature for the model design'''
63
64 #For EXANG,
65 '''There is a higher proportion (of those who have a diseased heart who also have non-induced exercise angina)
66 than those with induced angina. Both type of anginas are not good and can be a sign of a failing heart.
67 I think there is a relationship between these two columns'''
68
69 #For SLOPE,
70 '''Proportion of those with a diseased heart begins to increase
71 with a much significant increase in their slope measurements(from flat-sloping to down-sloping).
72 I think there is a relationship between these two columns
73 and SLOPE will be an important feature for the model design'''
74
75
76 '''For CA feature, I expect a majority of those with 2, 3 or 4 coloured vessels
77 to have a diseased heart but the reverse is the case.
78 Even the majority of those with no coloured vessel at all have a diseased heart
79 I think there is no relationship between these two columns and will not be an import feature for the model'''
80
81 #THAL,
82 '''There is a higher proportion (of those who have a diseased heart who also have fixed-defect thal)
83 than those with reversible defect thal. Both type of thal are not good and can be a sign of a failing heart'''
```

```
84 I think there is also a relationship between these two columns'''  
85  
86  
87 '''The features in blue highlights all have a relationship with TARGET(outcome variable)'''
```

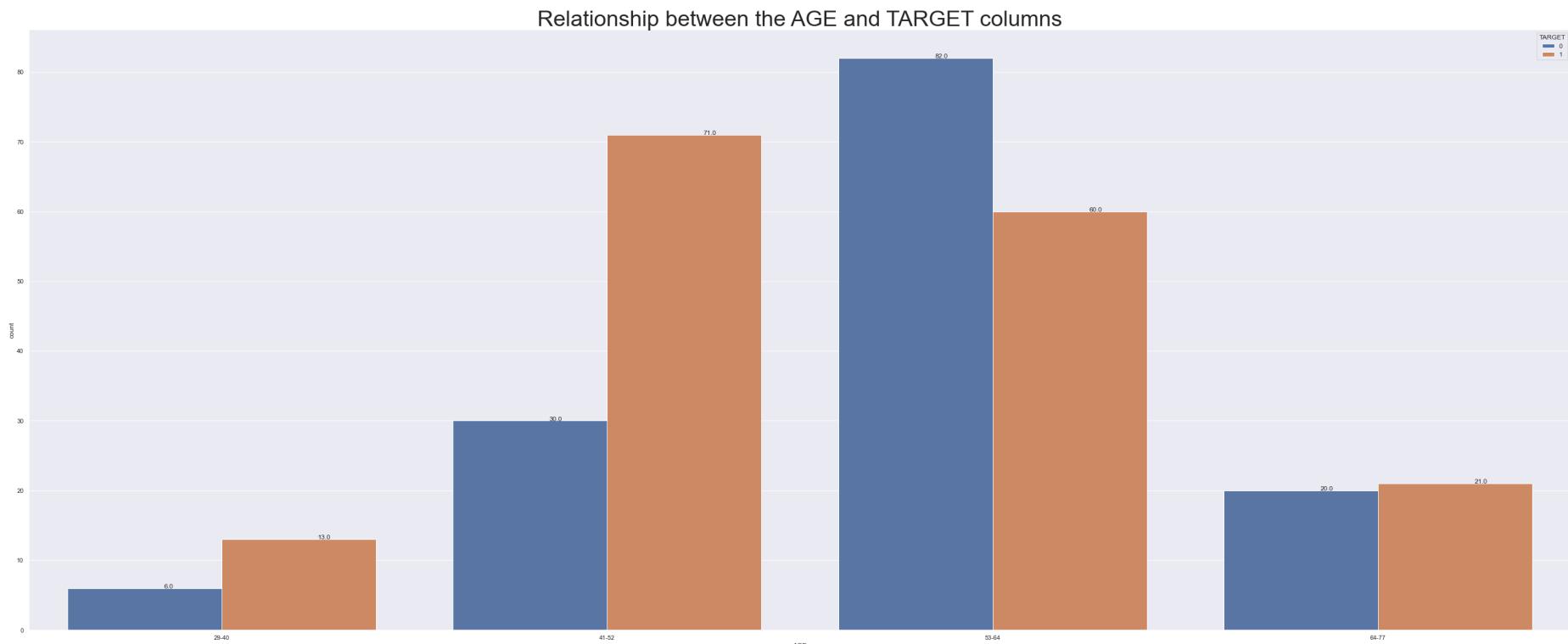
```
Out[504]: 'The features in blue highlights all have a relationship with TARGET(outcome variable)'
```

```
In [505]: 1 df['TARGET'].value_counts()
```

```
Out[505]: 1    165  
0    138  
Name: TARGET, dtype: int64
```

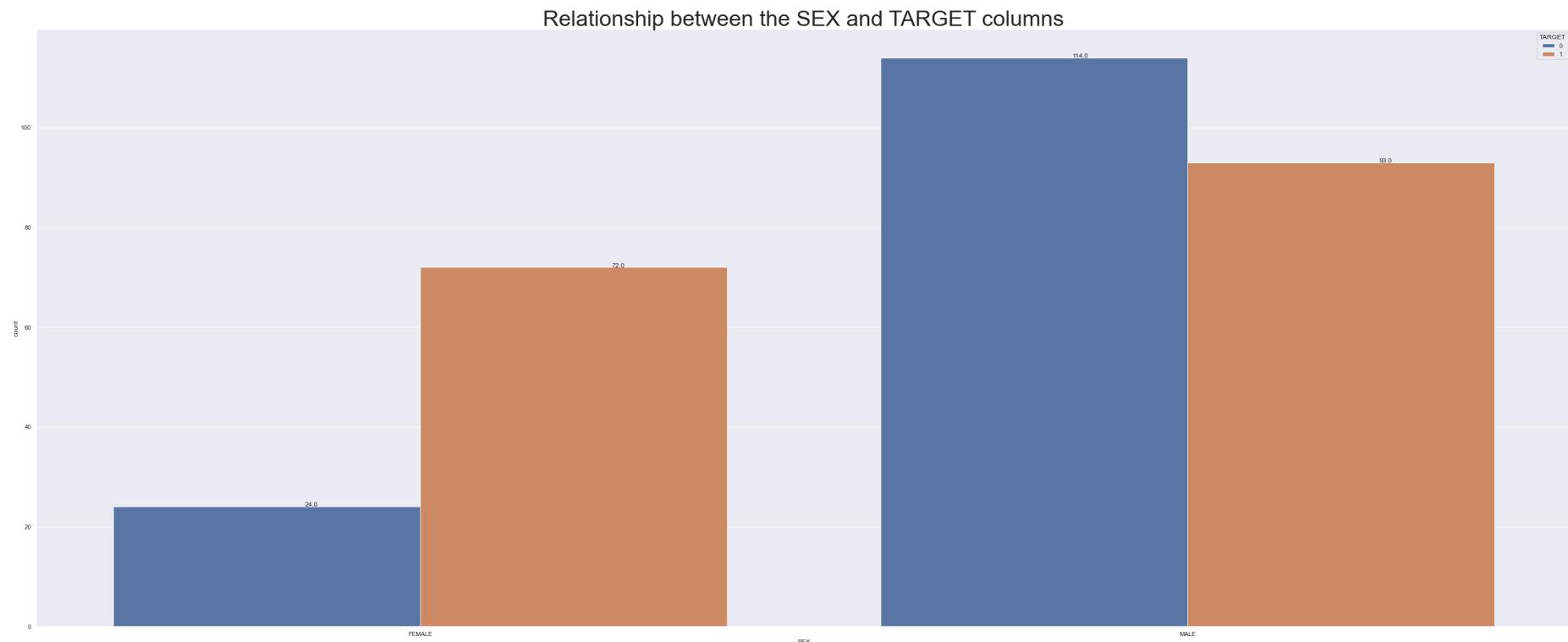
In [506]:

```
1 edu = sns.countplot(x='AGE', hue='TARGET', data=df)
2 edu.set_xticklabels(['29-40', '41-52', '53-64', '64-77'])
3 edu.set_title('Relationship between the AGE and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



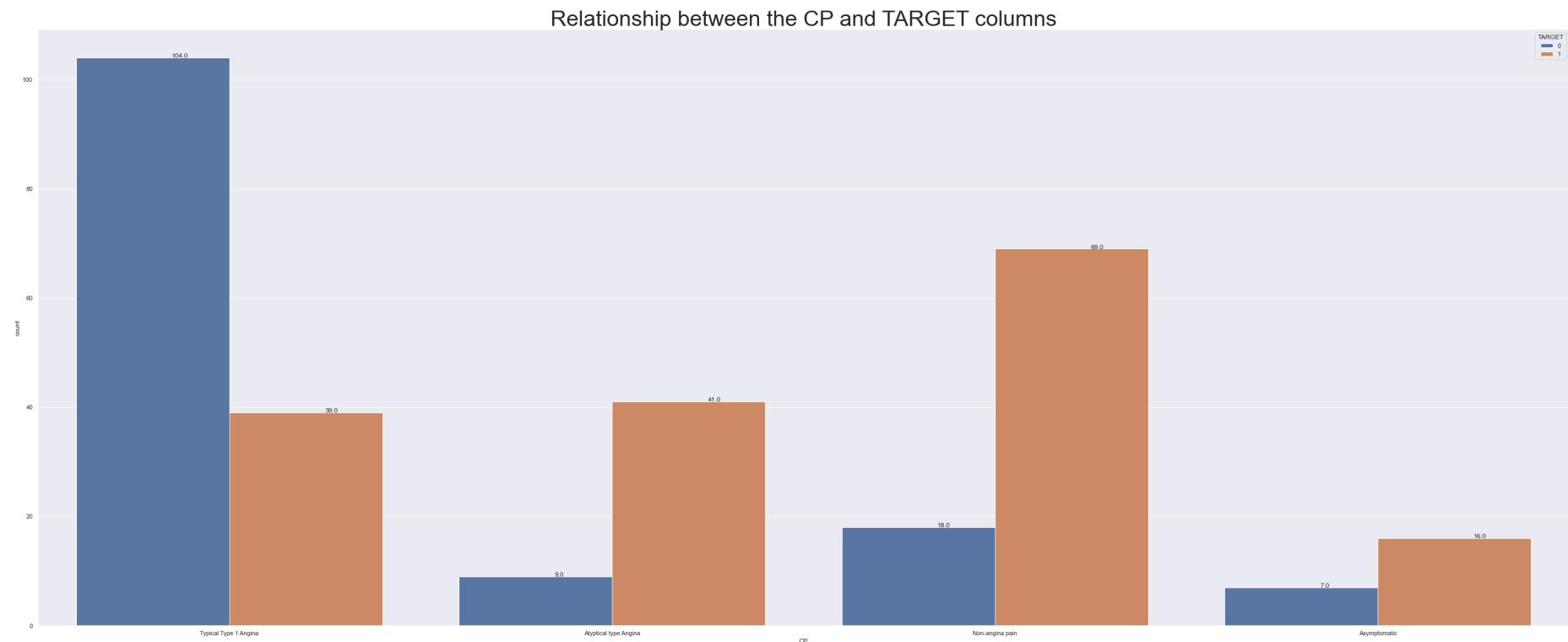
In [507]:

```
1 edu = sns.countplot(x='SEX', hue='TARGET', data=df)
2 edu.set_xticklabels(['FEMALE', 'MALE'])
3 edu.set_title('Relationship between the SEX and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



In [508]:

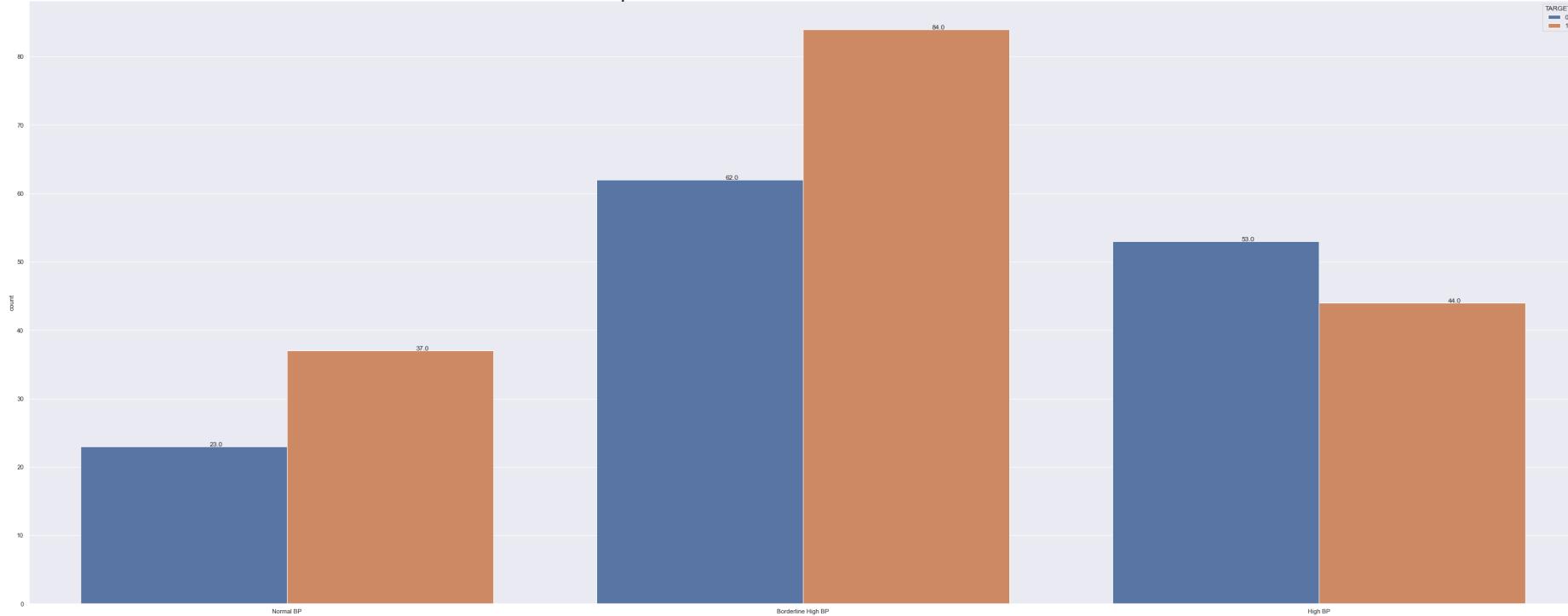
```
1 edu = sns.countplot(x='CP', hue='TARGET', data=df)
2 edu.set_xticklabels(['Typical Type 1 Angina', 'Atypical type Angina', 'Non-angina pain', 'Asymptomatic'])
3 edu.set_title('Relationship between the CP and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



In [509]:

```
1 edu = sns.countplot(x='TRestBPS', hue='TARGET', data=df)
2 edu.set_xticklabels(['Normal BP','Borderline High BP','High BP'])
3 edu.set_title('Relationship between the TRestBPS and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```

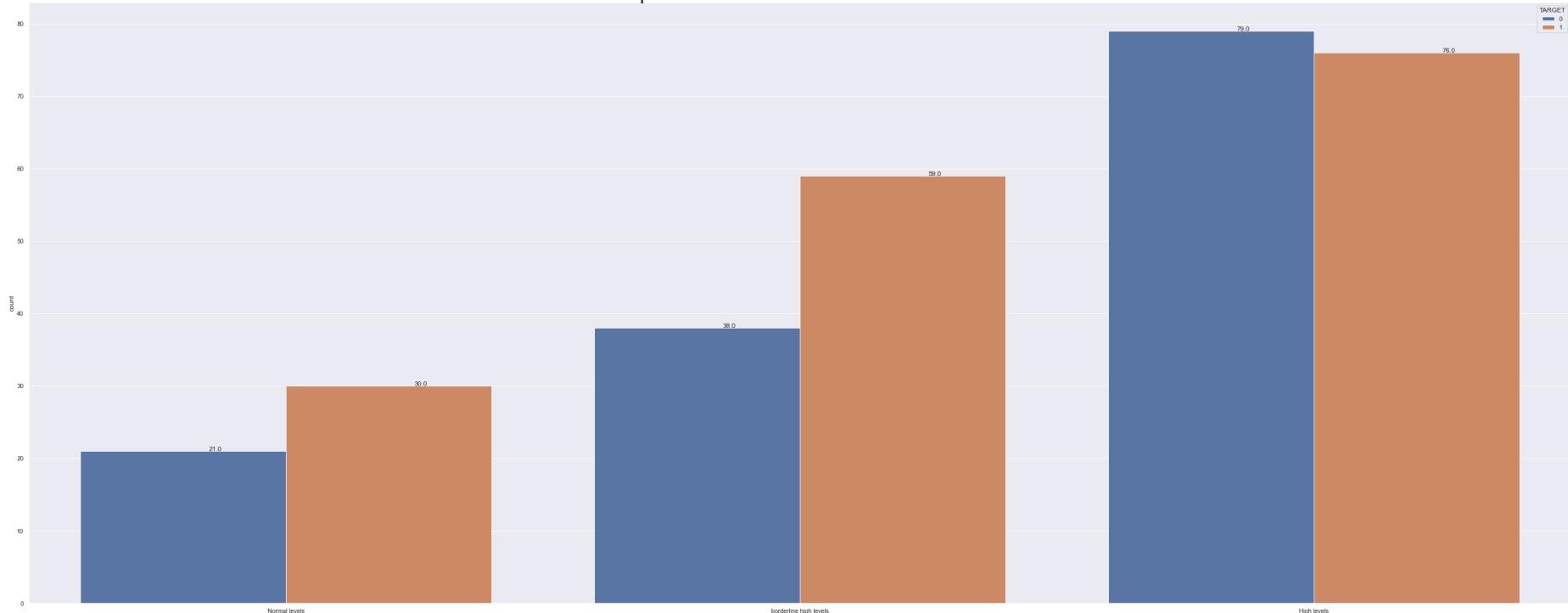
Relationship between the TRestBPS and TARGET columns



In [510]:

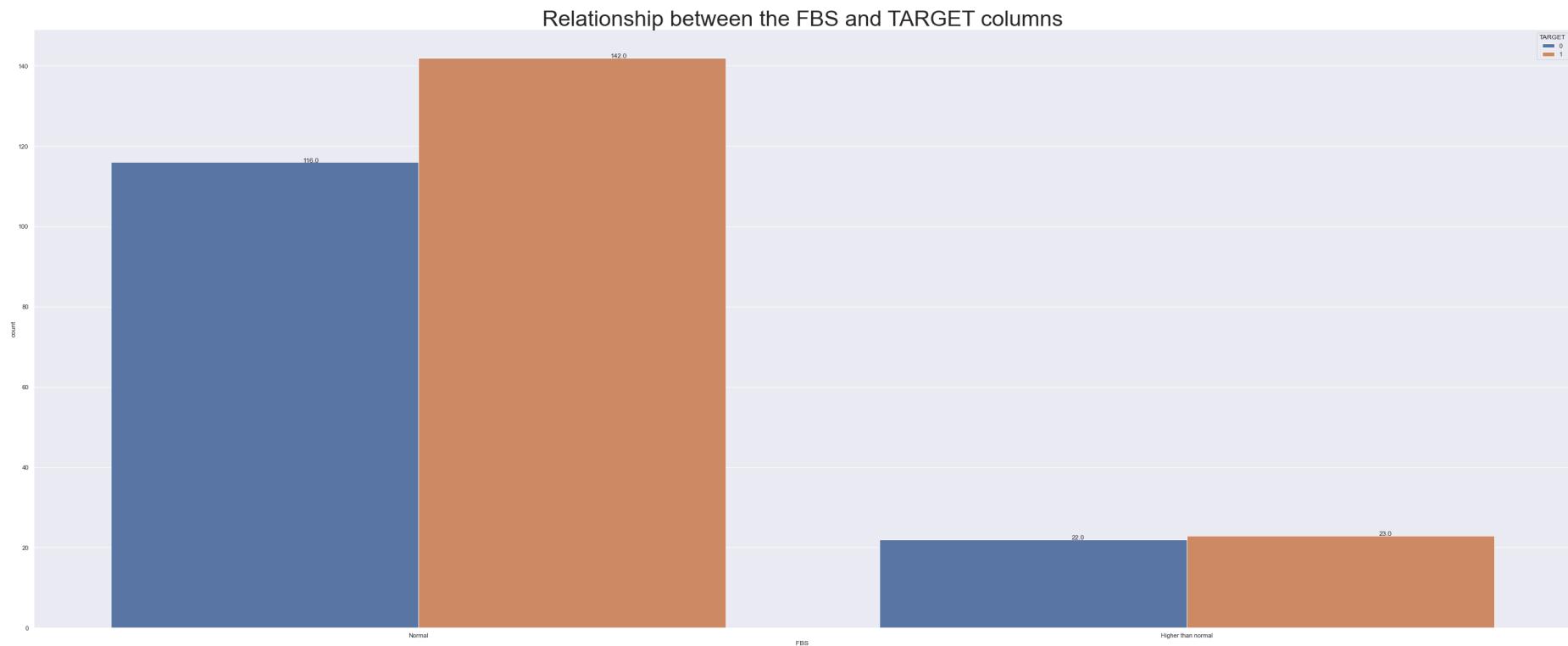
```
1 edu = sns.countplot(x='CHOL', hue='TARGET', data=df)
2 edu.set_xticklabels(['Normal levels','borderline high levels','High levels'])
3 edu.set_title('Relationship between the CHOL and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```

Relationship between the CHOL and TARGET columns



In [511]:

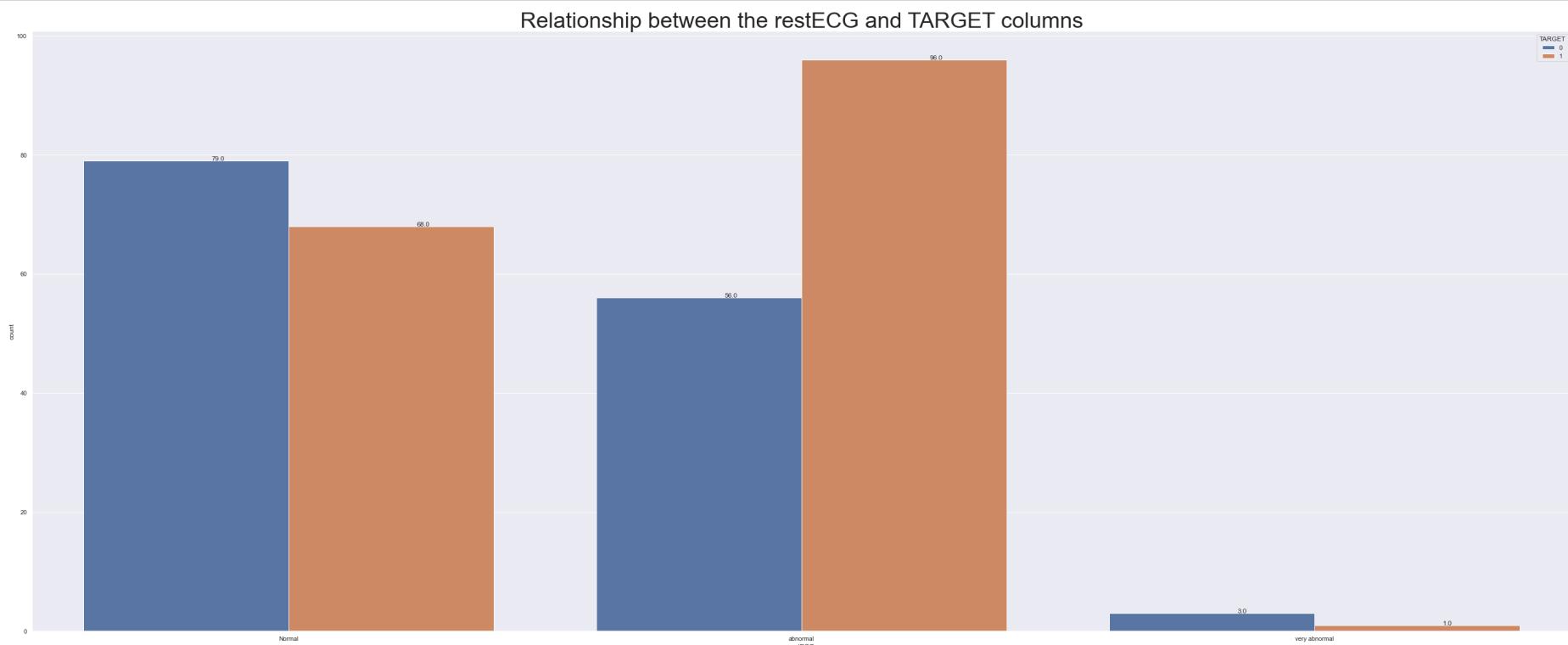
```
1 edu = sns.countplot(x='FBS', hue='TARGET', data=df)
2 edu.set_xticklabels(['Normal', 'Higher than normal'])
3 edu.set_title('Relationship between the FBS and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



```
In [512]: 1 df['FBS'].value_counts()
```

```
Out[512]: 0    258  
1     45  
Name: FBS, dtype: int64
```

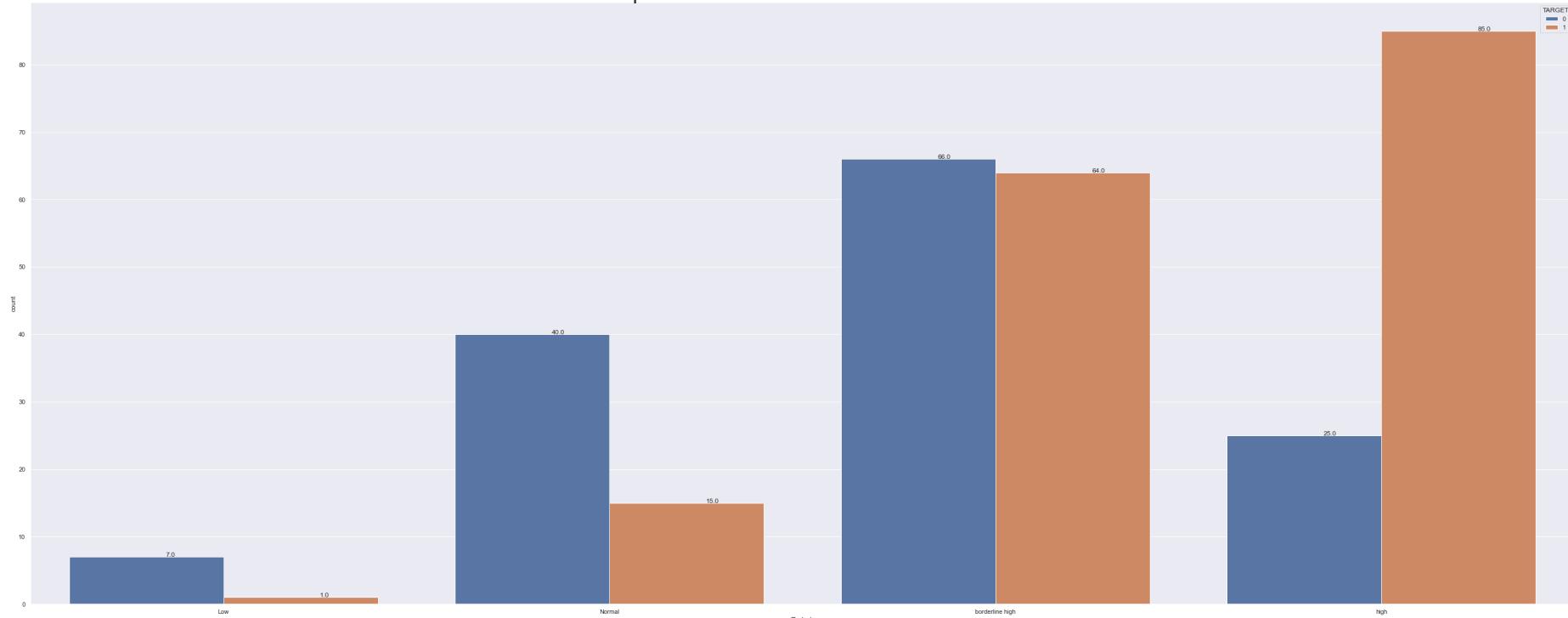
```
In [513]: 1 edu = sns.countplot(x='restECG', hue='TARGET', data=df)  
2 edu.set_xticklabels(['Normal', 'abnormal', 'very abnormal'])  
3 edu.set_title('Relationship between the restECG and TARGET columns', fontsize=40)  
4 '''for p in ax.patches:  
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''  
6  
7 for p in edu.patches:  
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x() + 0.25, p.get_height() + 0.01))
```



In [514]:

```
1 edu = sns.countplot(x='Thalach', hue='TARGET', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```

Relationship between the Thalach and TARGET columns



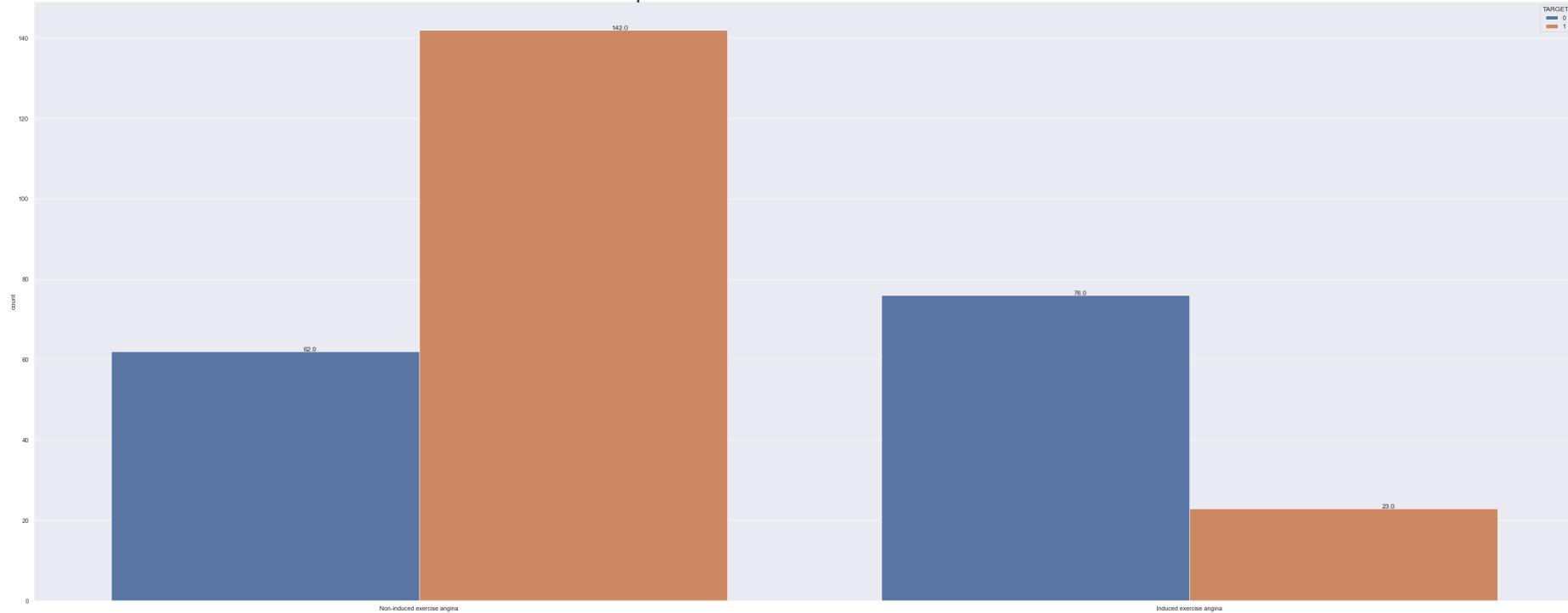
In []:

1

In [515]:

```
1 edu = sns.countplot(x='EXANG', hue='TARGET', data=df)
2 edu.set_xticklabels(['Non-induced exercise angina', 'Induced exercise angina'])
3 edu.set_title('Relationship between the EXANG and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7
8 for p in edu.patches:
9     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```

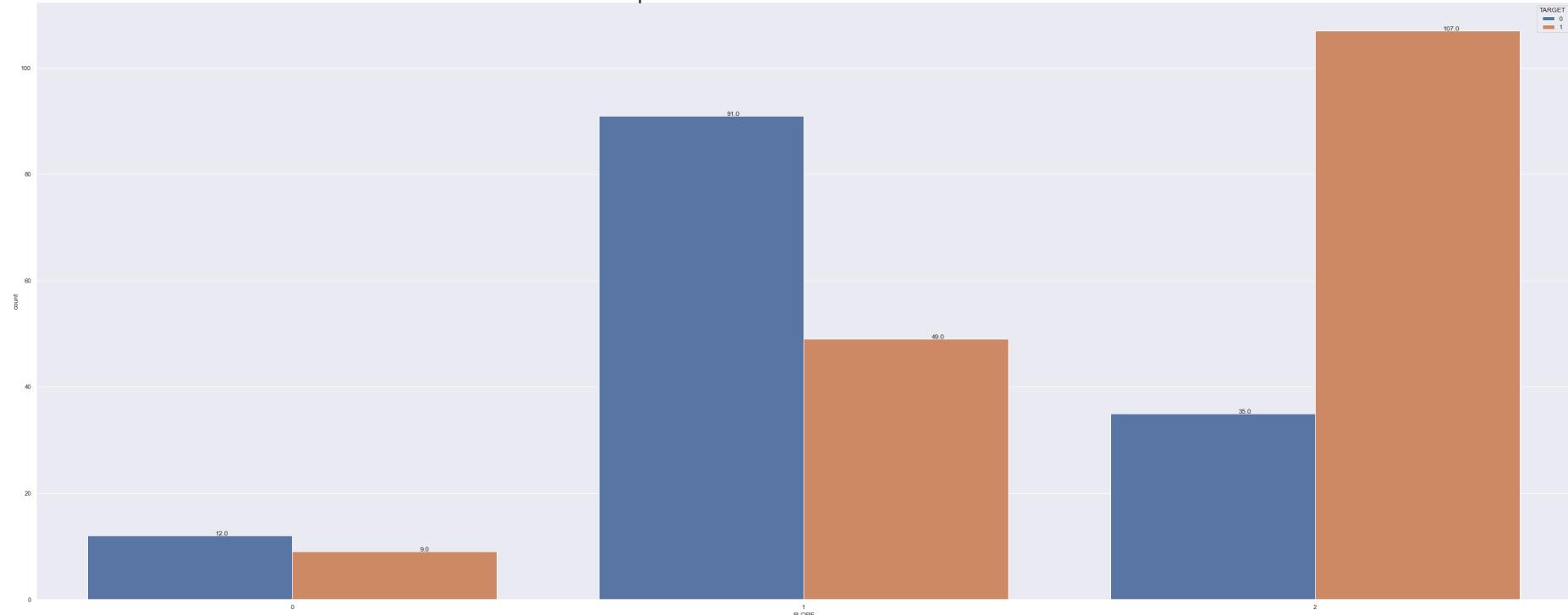
Relationship between the EXANG and TARGET columns



In [516]:

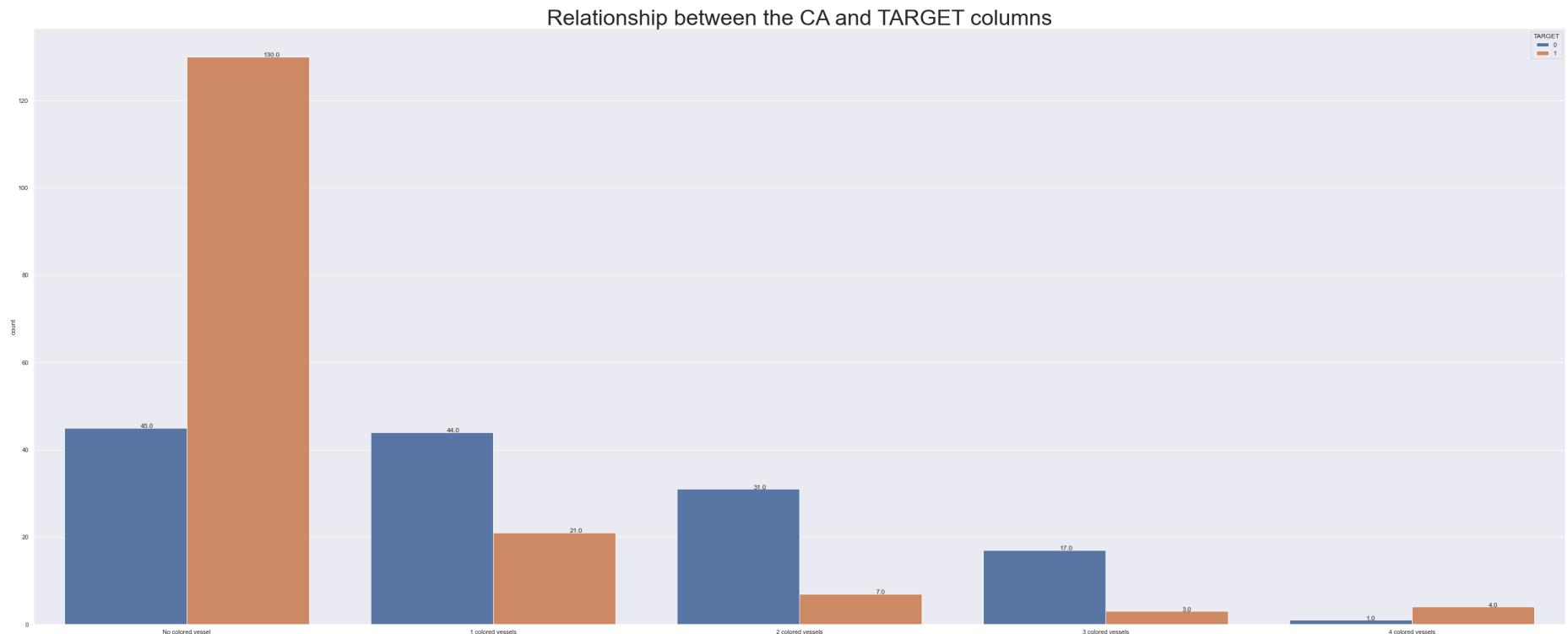
```
1 edu = sns.countplot(x='SLOPE', hue='TARGET', data=df)
2 edu.set_xticklabels(['0','1','2'])
3 edu.set_title('Relationship between the SLOPE and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```

Relationship between the SLOPE and TARGET columns



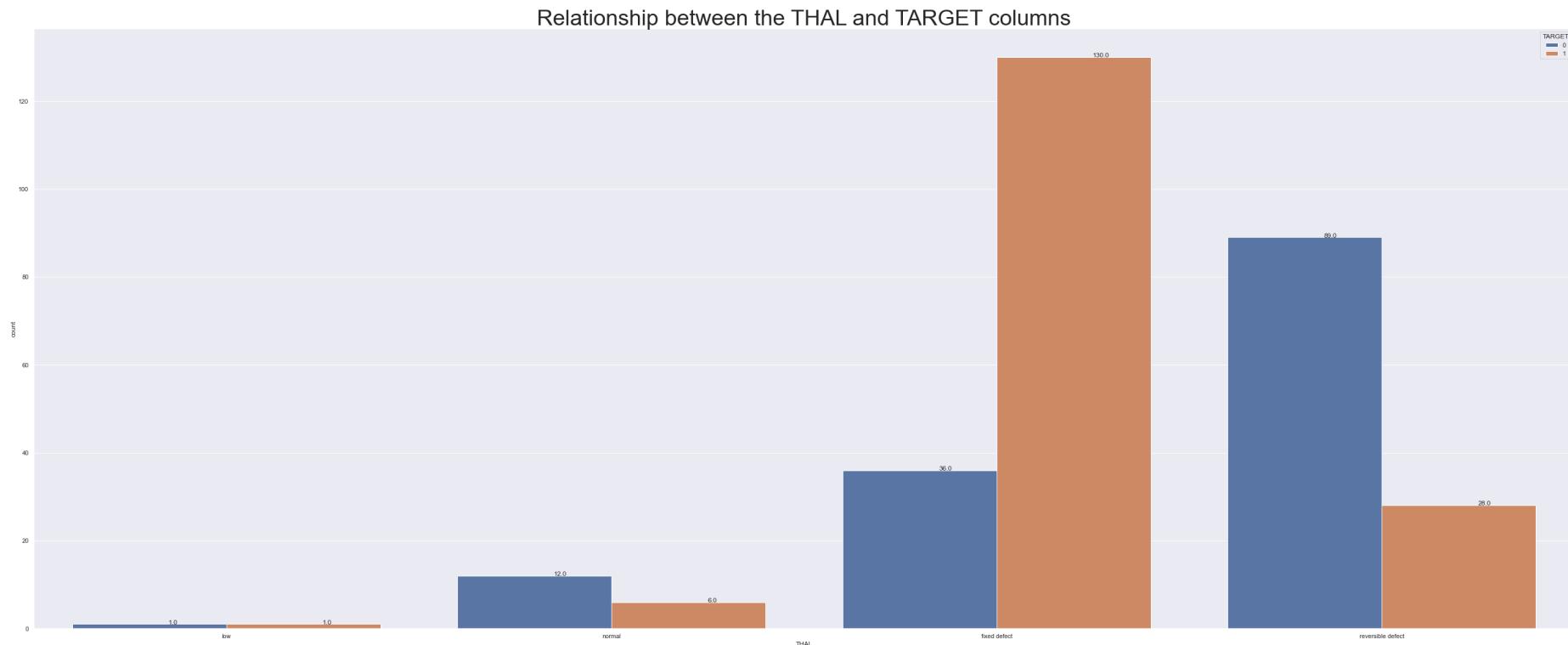
In [517]:

```
1 edu = sns.countplot(x='CA', hue='TARGET', data=df)
2 edu.set_xticklabels(['No colored vessel','1 colored vessels','2 colored vessels','3 colored vessels','4 colored vessels'])
3 edu.set_title('Relationship between the CA and TARGET columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7
8 for p in edu.patches:
9     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



In [518]:

```
1
2 edu = sns.countplot(x='THAL', hue='TARGET', data=df)
3 edu.set_xticklabels(['low', 'normal','fixed defect','reversible defect'])
4 edu.set_title('Relationship between the THAL and TARGET columns', fontsize=40)
5 '''for p in ax.patches:
6     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
7
8 for p in edu.patches:
9     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.25, p.get_height()+0.01))
```



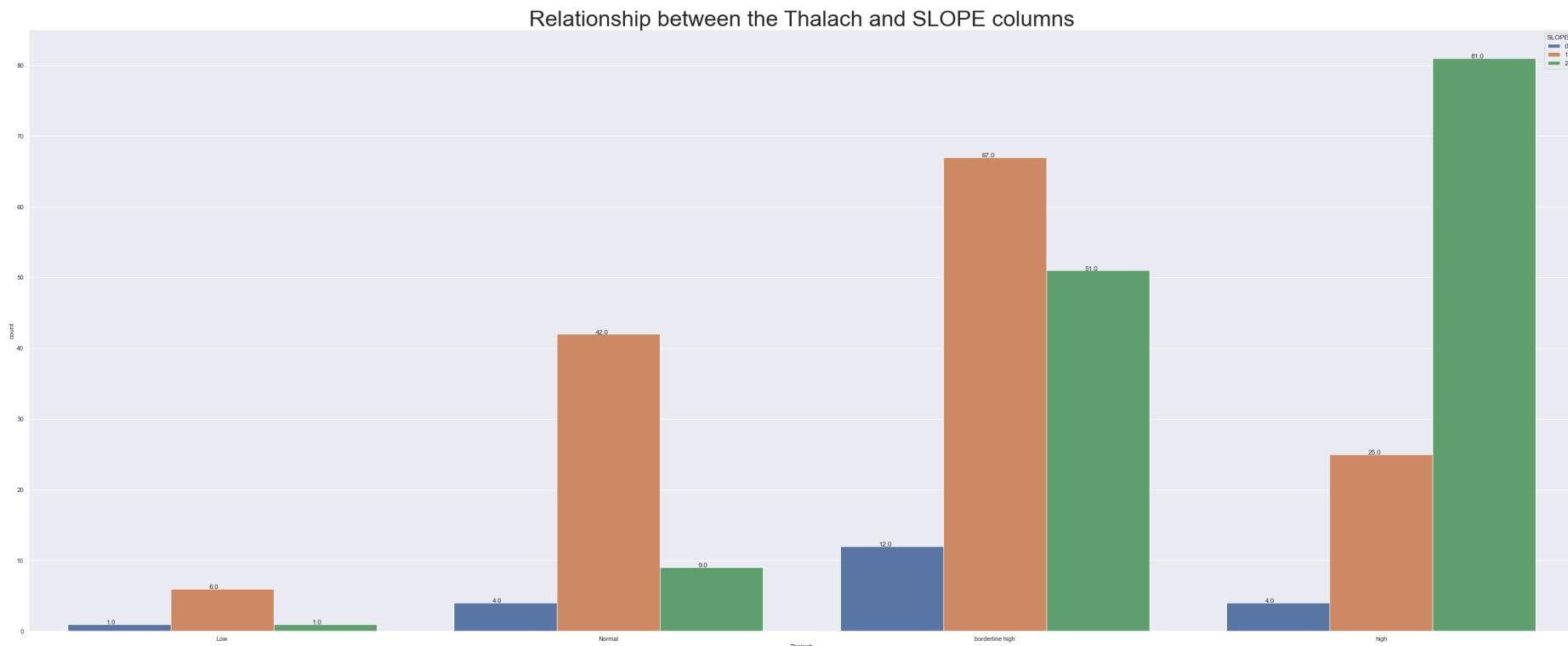
```
In [519]: plotting the Distributions and Relationships Between Specific Features
plots between "Thalach" and 'slope'
    "Thalach" and 'CP'
    'Age' and 'TRestBPS'
    'Age' and 'CA'
    'Thalach' and 'CHOL'
    'Thalach' and 'CA'
    'Thalach' and 'restECG'
```

The plots between specific variables I have made below all have a relationship as one variable increases, the other also increases while for AGE and CA, 'Thalach' and 'CHOL' and 'Thalach' and 'CA' the reverse is the case

```
Out[519]: "The plots between specific variables I have made below all have a relationship as one variable increases, \nthe other also increases while for AGE and CA, 'Thalach' and 'CHOL' and 'Thalach' and 'CA' the reverse is the case "
```

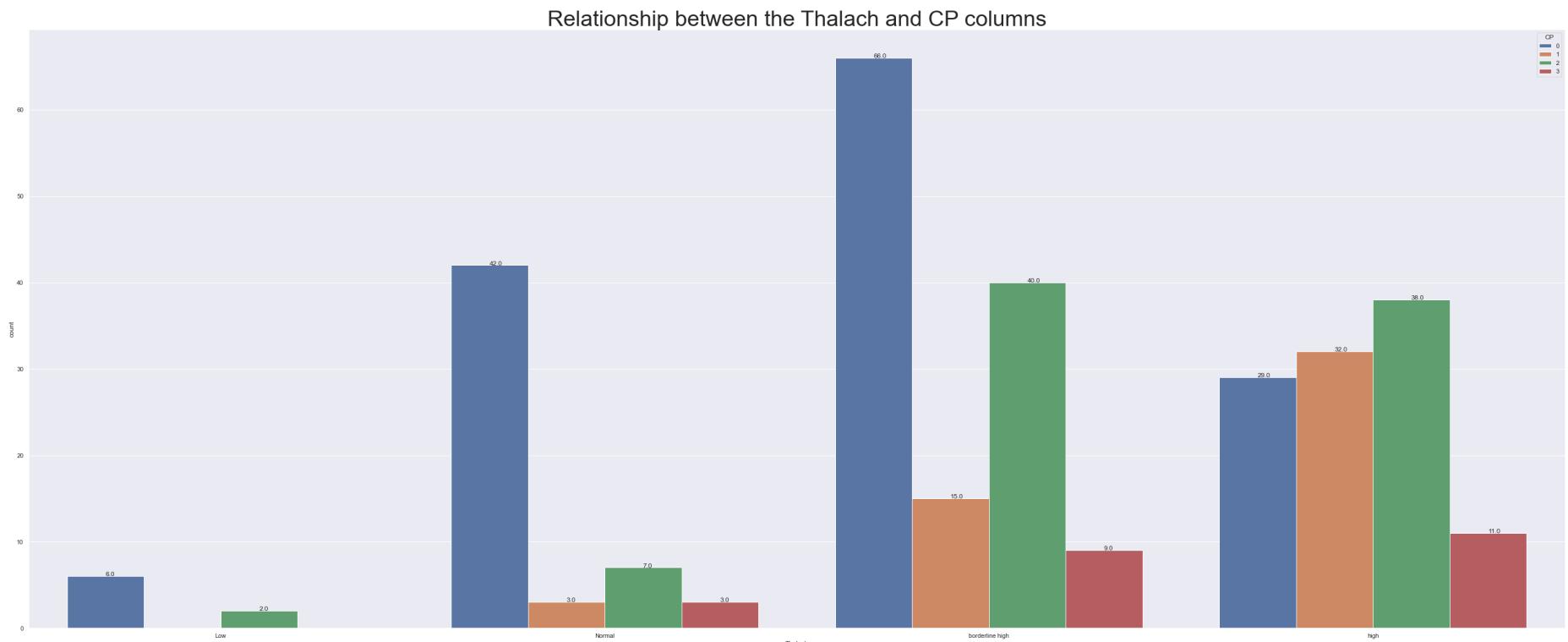
In [520]:

```
1 edu = sns.countplot(x='Thalach', hue='SLOPE', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and SLOPE columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```



In [521]:

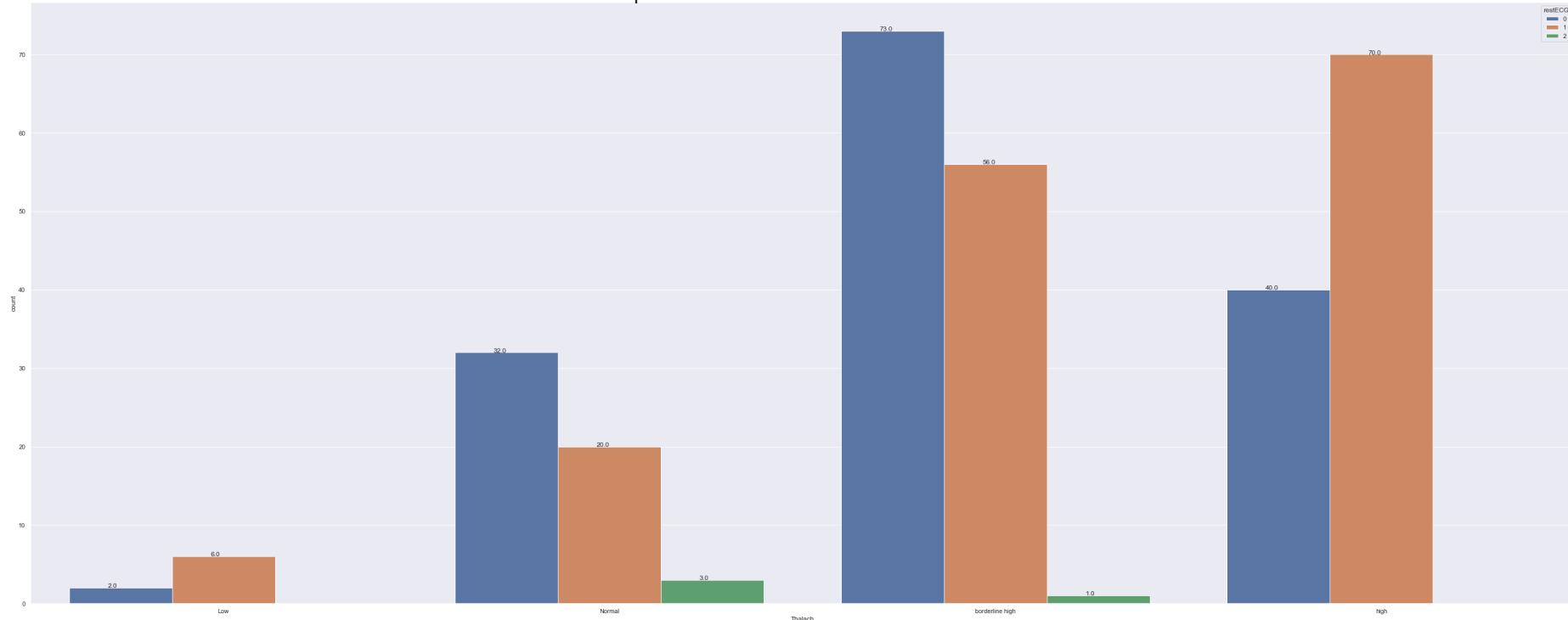
```
1 edu = sns.countplot(x='Thalach', hue='CP', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and CP columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```



In [522]:

```
1 edu = sns.countplot(x='Thalach', hue='restECG', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and restECG columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```

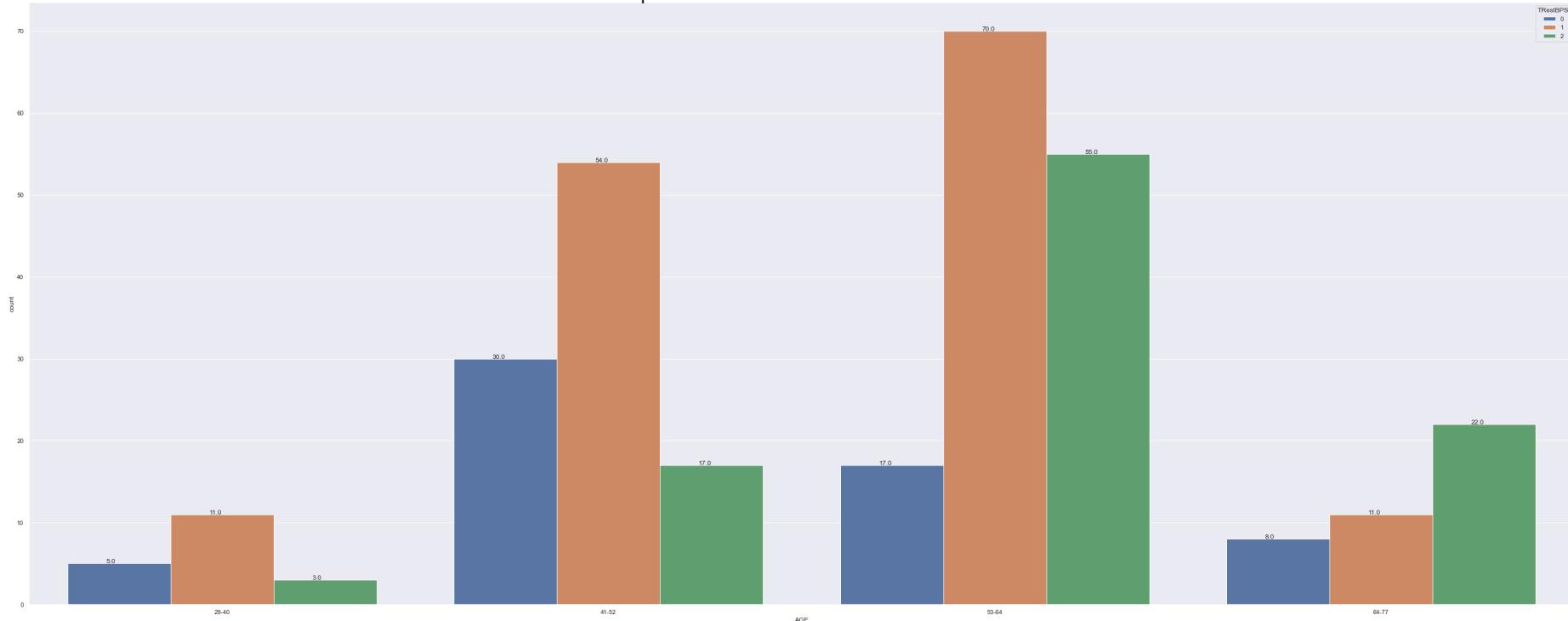
Relationship between the Thalach and restECG columns



In [523]:

```
1 edu = sns.countplot(x='AGE', hue='TRestBPS', data=df)
2 edu.set_xticklabels(['29-40', '41-52', '53-64', '64-77'])
3 edu.set_title('Relationship between the AGE and TRestBPS columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```

Relationship between the AGE and TRestBPS columns

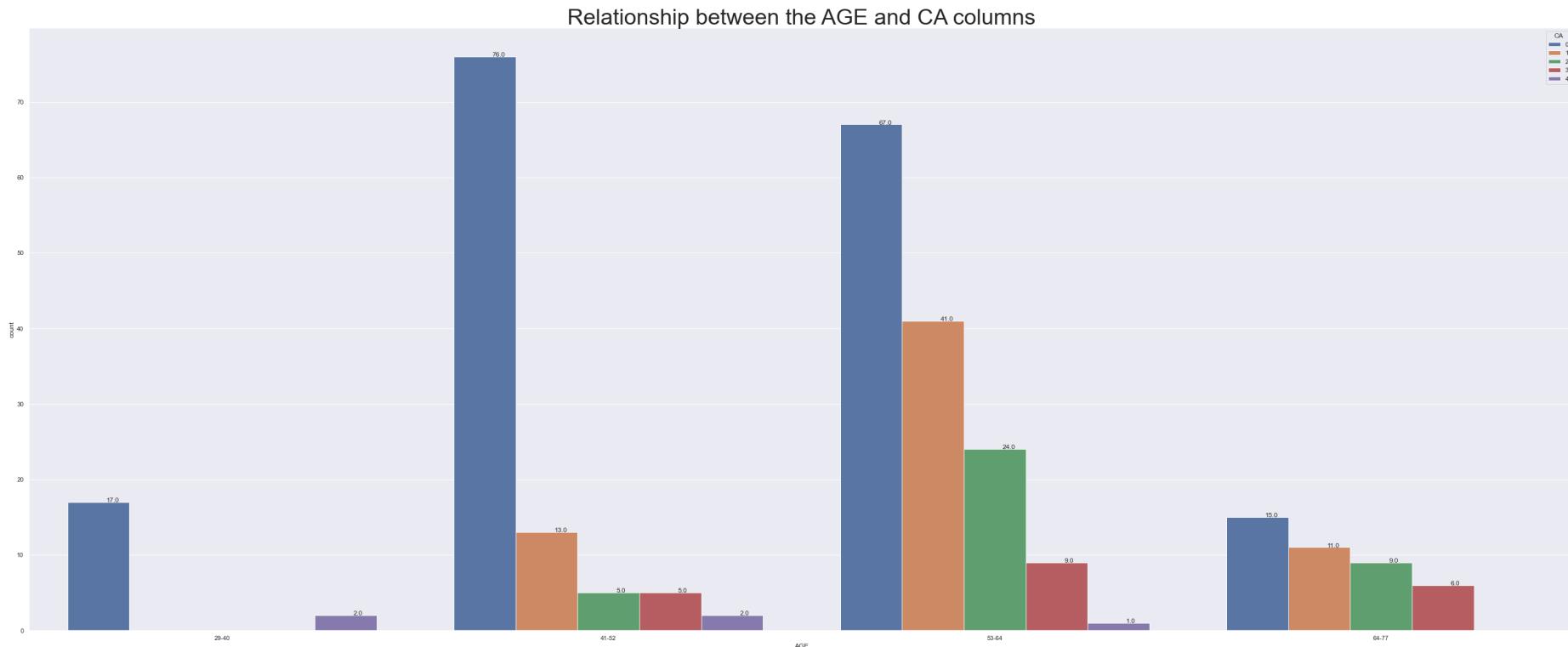


```
In [ ]: 1
```

```
In [524]: 1 df.FBS.unique()
```

```
Out[524]: array([1, 0])
```

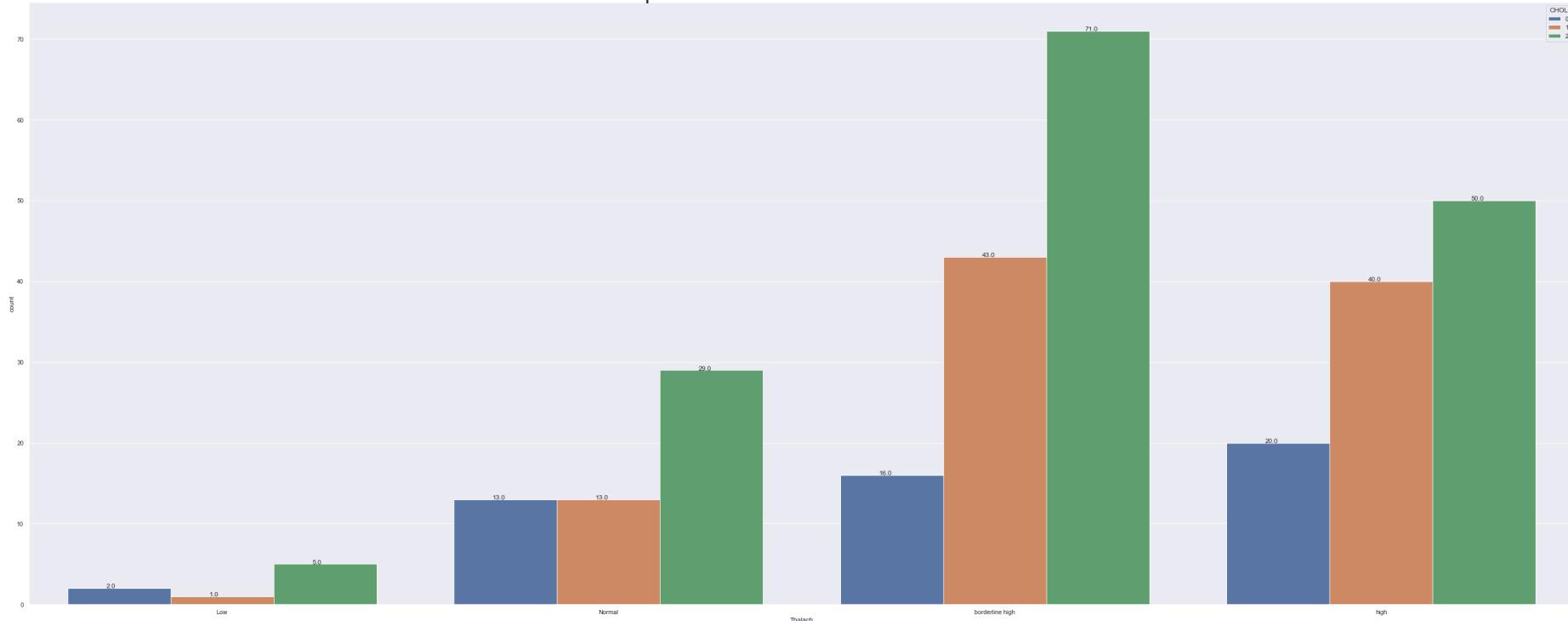
```
In [525]: 1 edu = sns.countplot(x='AGE', hue='CA', data=df)
2 edu.set_xticklabels(['29-40', '41-52', '53-64', '64-77'])
3 edu.set_title('Relationship between the AGE and CA columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```



In [526]:

```
1 edu = sns.countplot(x='Thalach', hue='CHOL', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and CHOL columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```

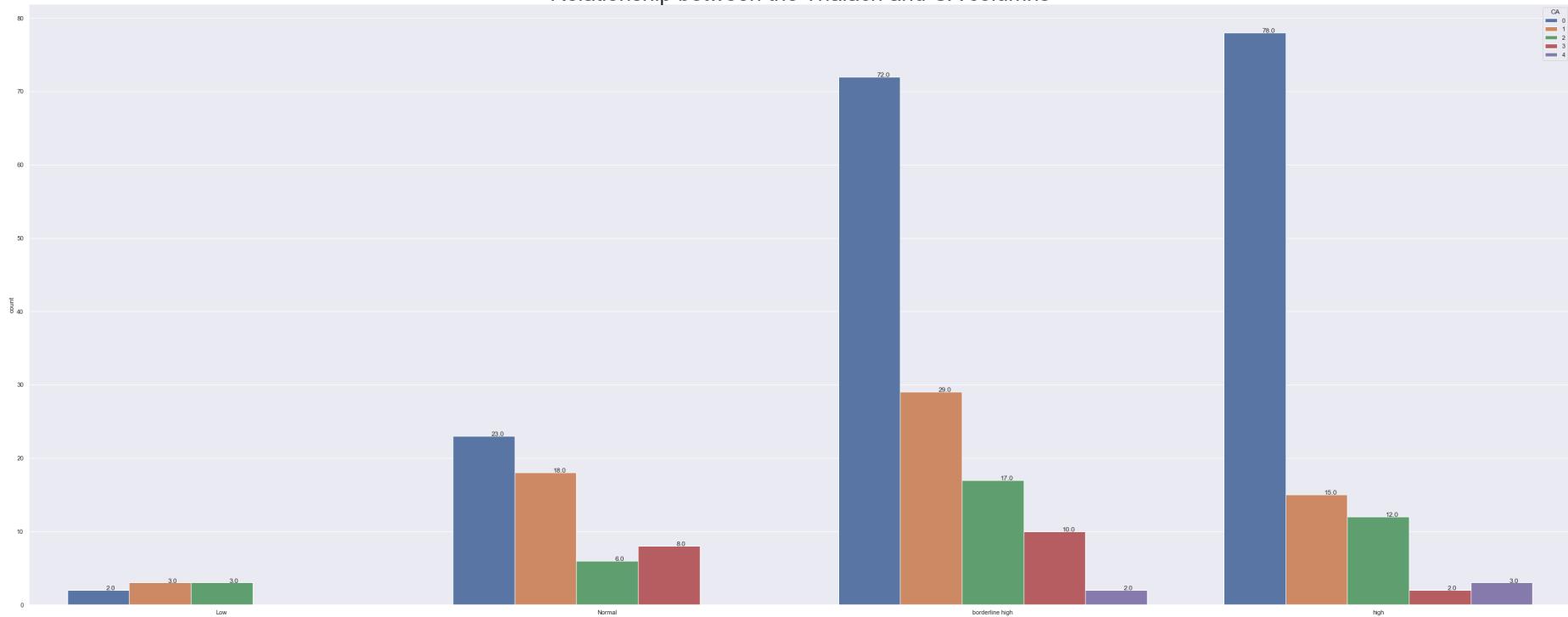
Relationship between the Thalach and CHOL columns



In [527]:

```
1 edu = sns.countplot(x='Thalach', hue='CA', data=df)
2 edu.set_xticklabels(['Low', 'Normal', 'borderline high', 'high'])
3 edu.set_title('Relationship between the Thalach and CA columns', fontsize=40)
4 '''for p in ax.patches:
5     ax.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.1, p.get_height()+50))'''
6
7 for p in edu.patches:
8     edu.annotate('{:.1f}'.format(p.get_height()), (p.get_x()+0.10, p.get_height()+0.01))
```

Relationship between the Thalach and CA columns



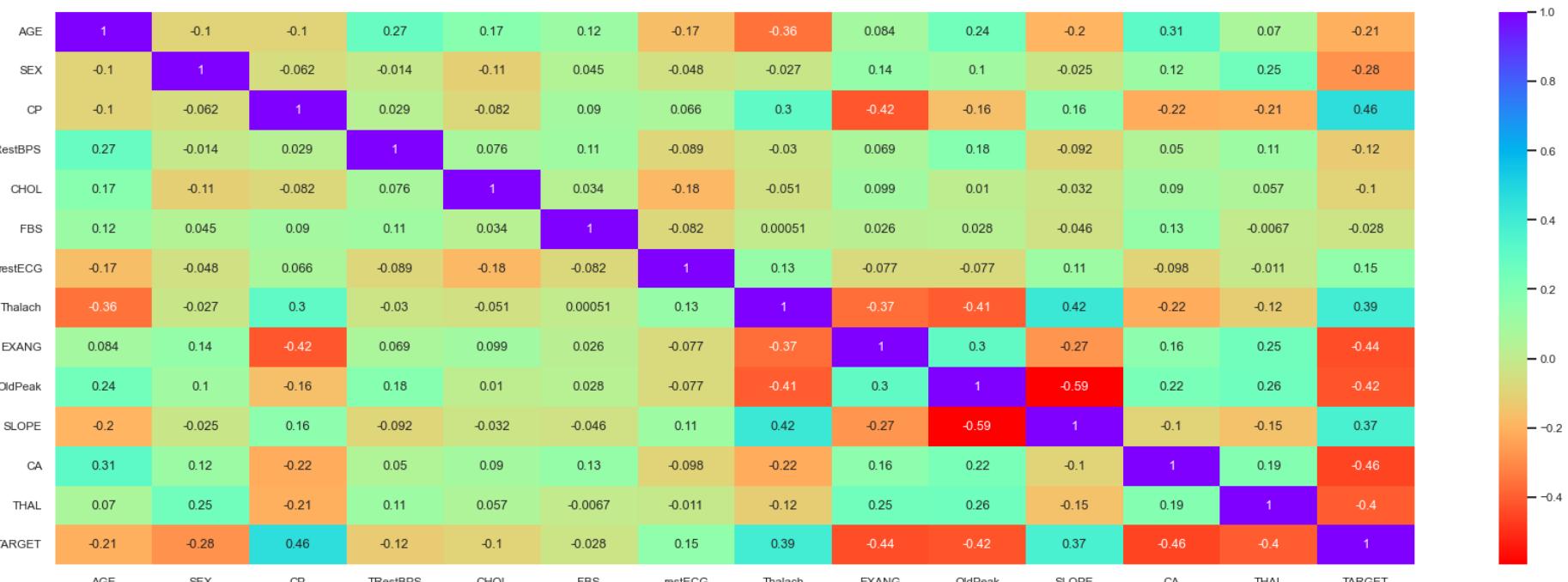
In [528]:

```

1 #Evaluating the Correlation between Columns Using a Heatmap
2
3 #Plot the heatmap for all the columns in the DataFrame (other than the ID column) by using sns.heatmap and
4 #keep the figure size as 30,10 for better visibility:
5
6 sns.set(rc={'figure.figsize':(30,10)})
7 sns.set_context("talk", font_scale=0.7)
8
9
10 #Use Spearman as the method parameter to compute Spearman's rank correlation coefficient:
11
12 sns.heatmap(df.iloc[:,0:14].corr(method='spearman'),cmap='rainbow_r', annot=True)

```

Out[528]: <AxesSubplot:>



In [529]:

```
1 #FINAL NOTES
2
3 '''A younger female adult (between ages 29-52) with chest pain (angina), thalassemia disease and also with
4
5 Features having a strong relationship with the outcome(TARGET) according to their decreasing rank are:
6
7 'CP' - 0.46
8 'Thalach' - 0.43
9 'SLOPE' - 0.37
10 'restECG' - 0.15
11 ''
12'''
```

Out[529]: "A younger female adult (between ages 29-52) with chest pain (angina), thalassemia disease and also with an abnormal resting ECG..... is at more risk of having a heart disease.\n\n\nFeatures having a strong relationship with the outcome(TARGET) according to their decreasing rank are:\n\n'CP' - 0.46\n'Thalach' - 0.43\n'SLOPE' - 0.37\n'restECG' - 0.15"

In []:

1

In []:

1