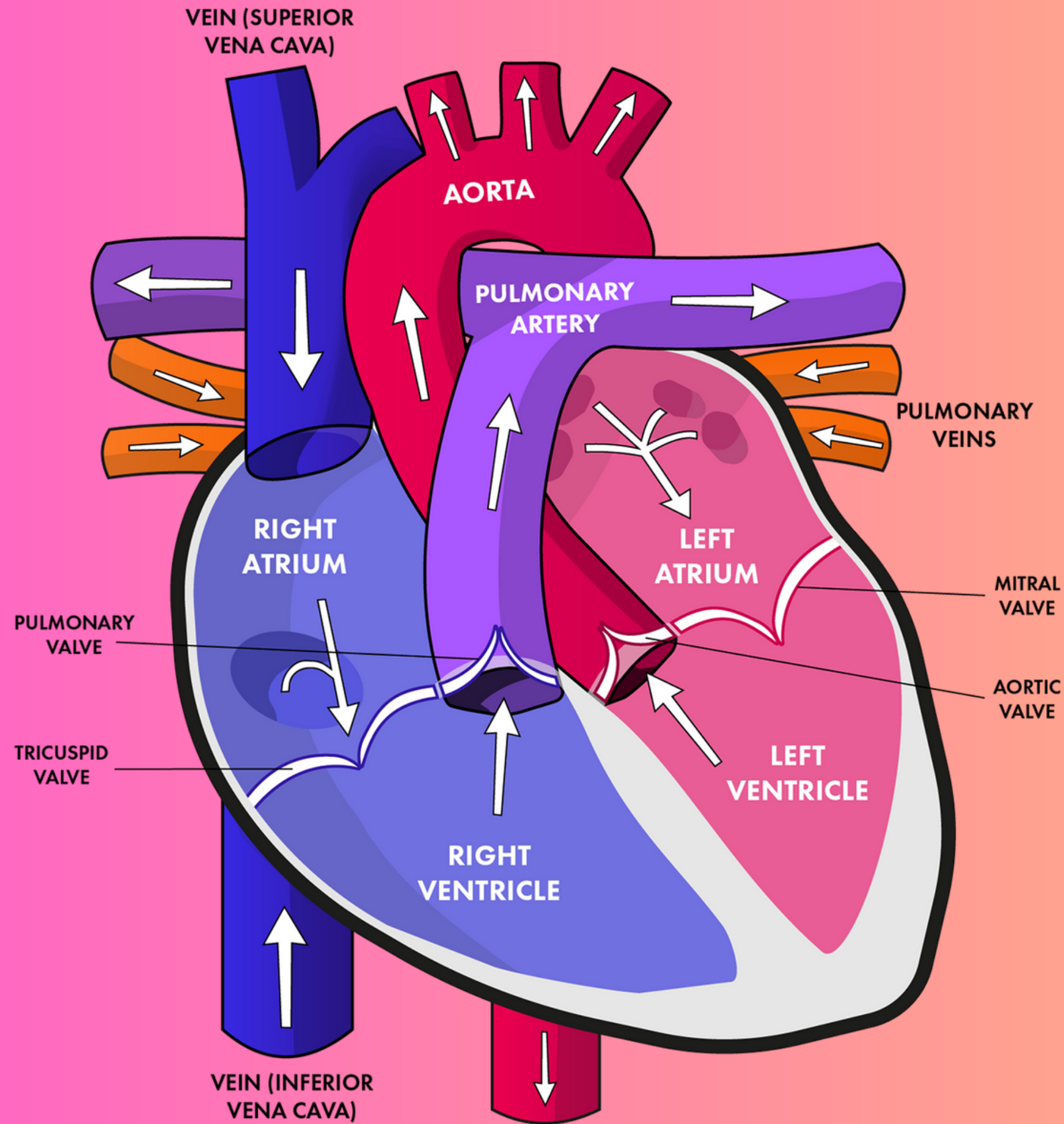


Team Project

Heart Disease Prediction



Context:

Cardiovascular diseases (CVDs) are the number 1 cause of death globally, taking an estimated 17.9 million lives each year, which accounts for 31% of all deaths worldwide. Four out of 5 CVD deaths are due to heart attacks and strokes, and one-third of these deaths occur prematurely in people under 70 years of age. Heart failure is a common event caused by CVDs and this dataset contains 13 features that can be used to predict a possible heart disease.

People with cardiovascular disease or who are at high cardiovascular risk (due to the presence of one or more risk factors such as hypertension, diabetes, hyperlipidaemia or already established disease) need early detection and management wherein a machine learning model can be of great help.

Data Contains

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 cholestoral in mg/dl

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 - have disease or not (1=yes, 0=no)

[illegible]

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  -
 0   age         294 non-null    float64
 1   sex         300 non-null    float64
 2   cp          297 non-null    float64
 3   trestbps    299 non-null    float64
 4   chol        297 non-null    float64
 5   fbs         299 non-null    float64
 6   restecg     300 non-null    float64
 7   thalach     291 non-null    float64
 8   exang       298 non-null    float64
 9   oldpeak     300 non-null    float64
10   slope       297 non-null    float64
11   ca          298 non-null    float64
12   thal        300 non-null    float64
13   target      303 non-null    int64
dtypes: float64(13), int64(1)
memory usage: 33.3 KB
```

Data Information

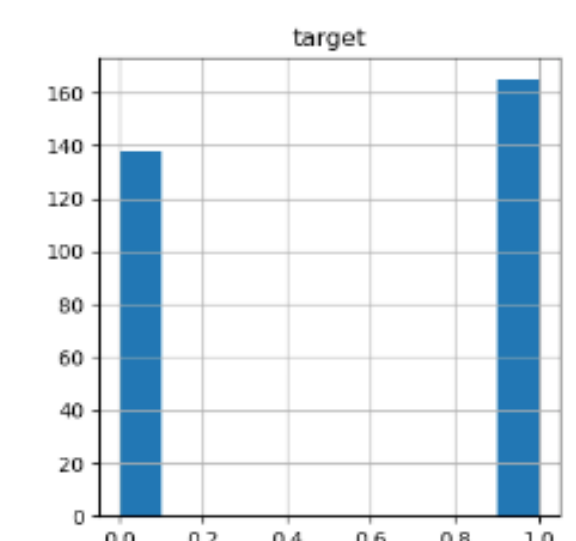
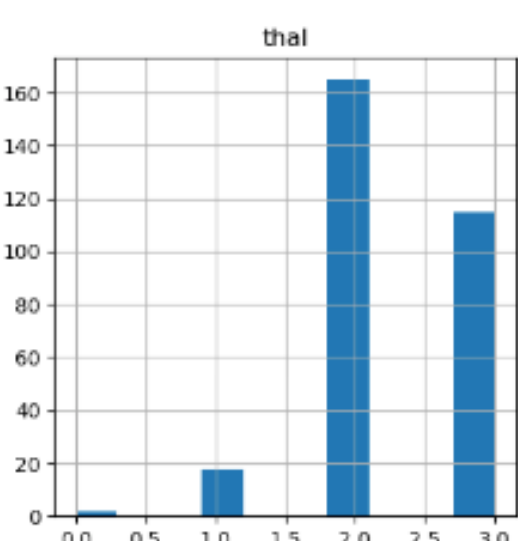
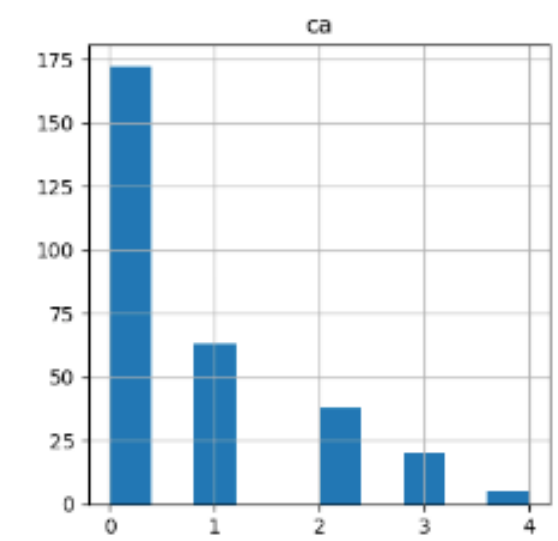
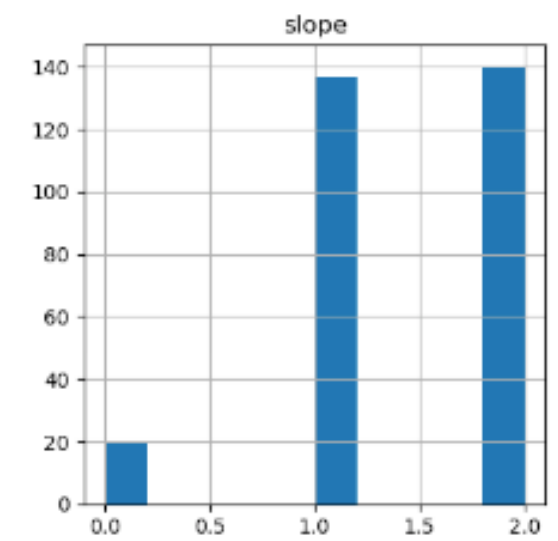
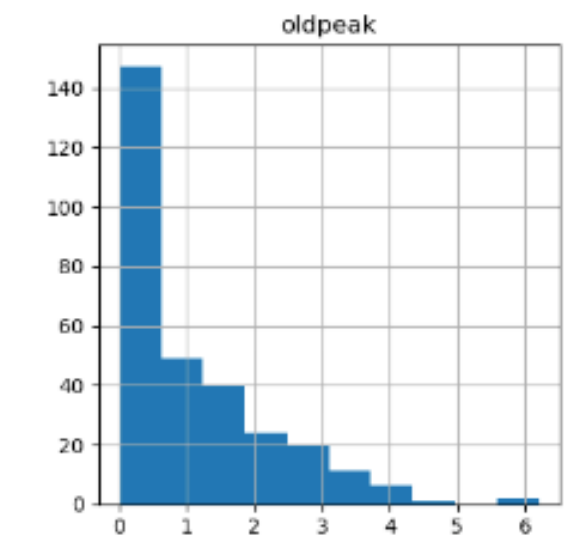
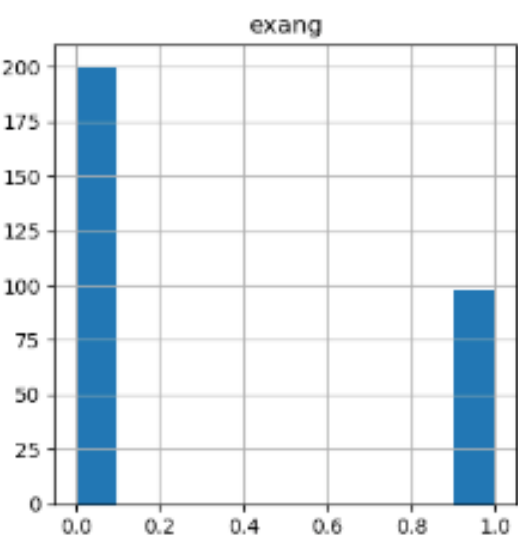
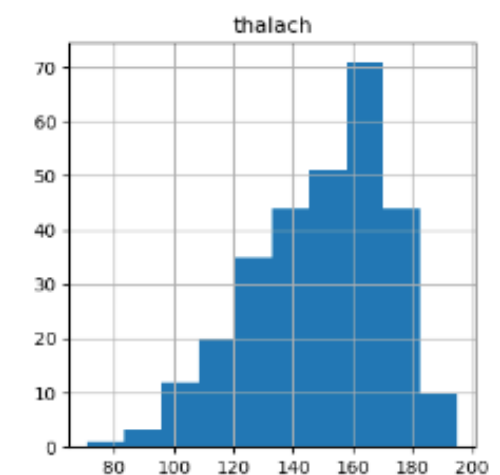
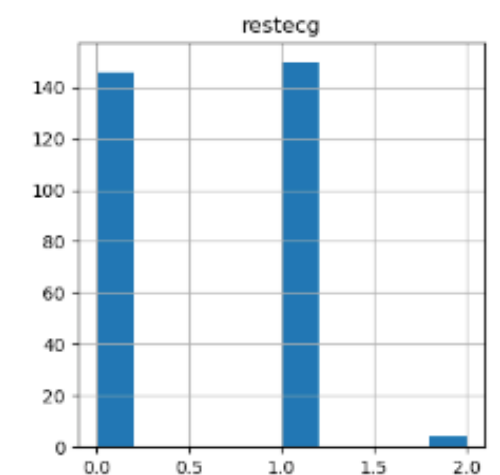
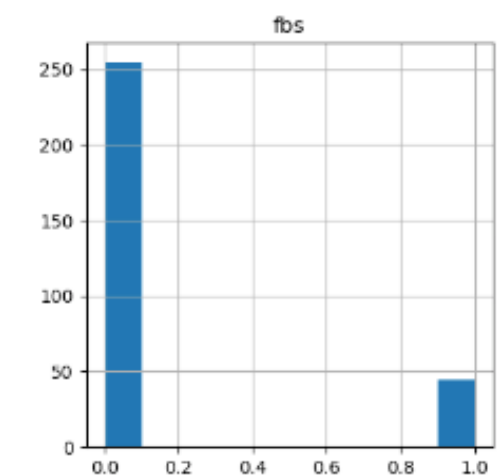
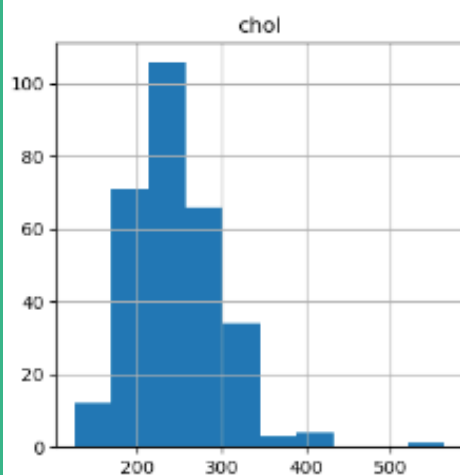
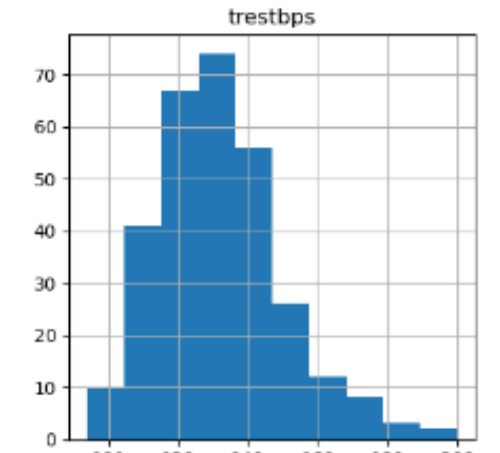
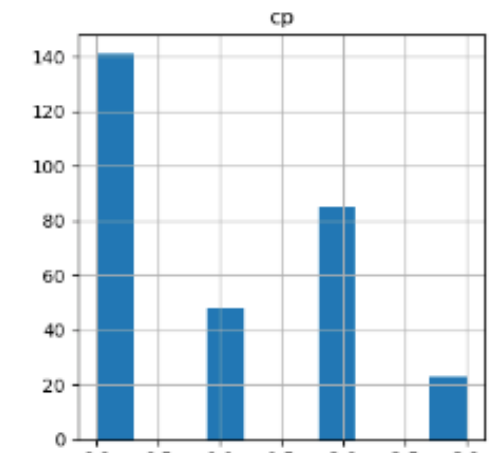
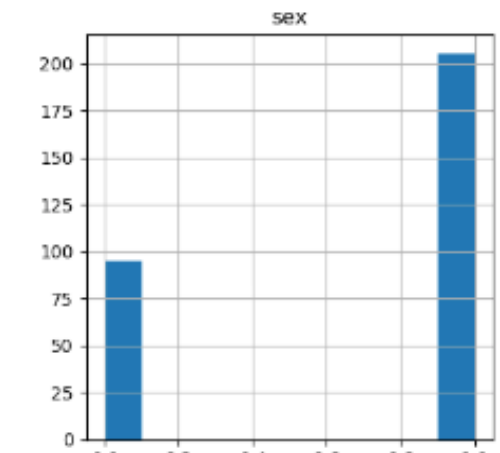
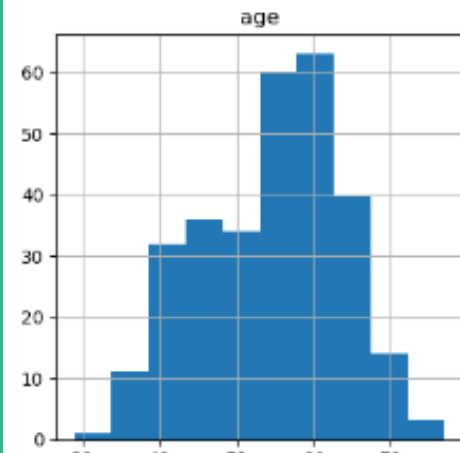
	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal	target
0	63.0	1.0	3.0	145.0	233.0	1.0	0.0	150.0	0.0	2.3	0.0	0.0	1.0	1
1	37.0	1.0	2.0	130.0	250.0	0.0	1.0	187.0	0.0	3.5	0.0	0.0	2.0	1
2	41.0	0.0	1.0	130.0	204.0	0.0	0.0	NaN	0.0	1.4	2.0	0.0	2.0	1
3	56.0	1.0	1.0	120.0	236.0	0.0	1.0	178.0	0.0	0.8	2.0	NaN	2.0	1
4	57.0	0.0	0.0	120.0	354.0	0.0	1.0	163.0	1.0	0.6	2.0	0.0	2.0	1
5	57.0	NaN	0.0	140.0	192.0	0.0	1.0	148.0	0.0	0.4	1.0	0.0	1.0	1
6	56.0	0.0	1.0	140.0	294.0	0.0	0.0	NaN	0.0	1.3	1.0	0.0	2.0	1
7	44.0	1.0	1.0	120.0	263.0	0.0	1.0	173.0	0.0	0.0	2.0	0.0	3.0	1
8	52.0	1.0	2.0	172.0	199.0	1.0	1.0	162.0	0.0	0.5	2.0	0.0	3.0	1
9	57.0	1.0	2.0	150.0	168.0	0.0	1.0	174.0	0.0	1.6	2.0	0.0	2.0	1

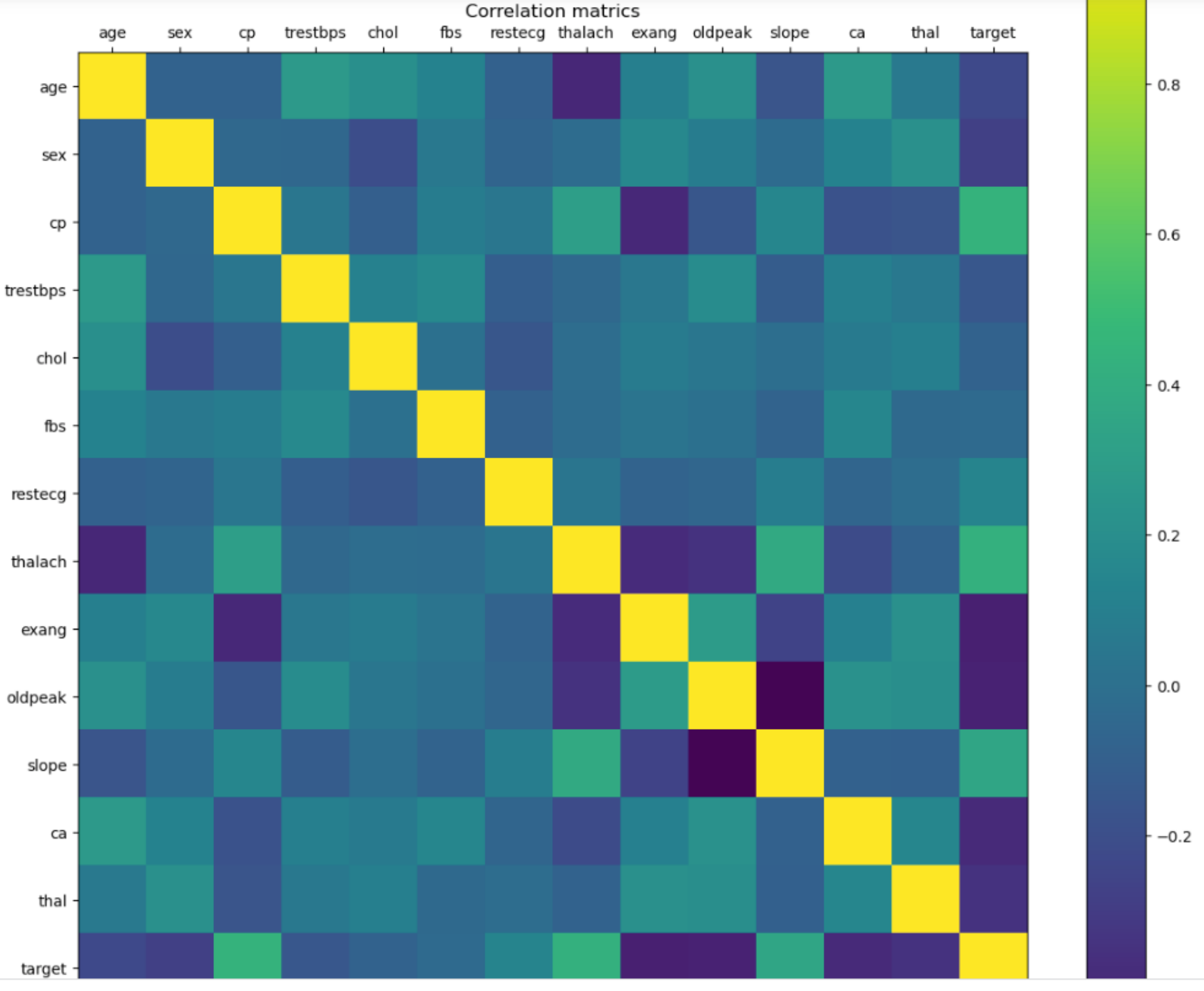
1	##### Types of Variable
2	dataset.dtypes
age float64	
sex float64	
cp float64	
trestbps float64	
chol float64	
fbs float64	
restecg float64	
thalach float64	
exang float64	
oldpeak float64	
slope float64	
ca float64	
thal float64	
target int64	
dtype: object	

All variables in dataset numerical

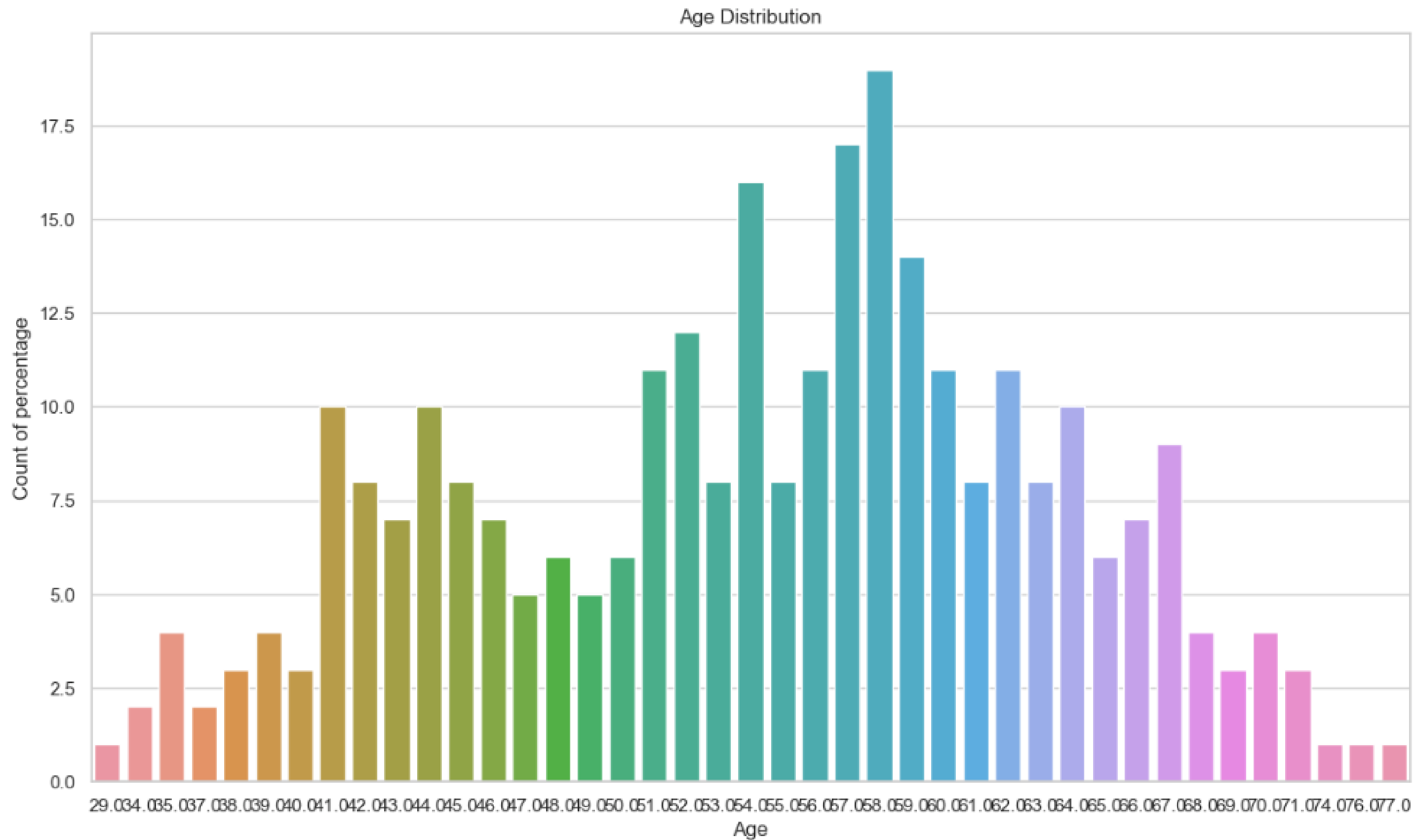
Data Describe

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	
count	294.000000	300.000000	297.000000	299.000000	297.000000	299.000000	300.000000	291.000000	298.000000	300.000000	297.000000	298.000000	300.00
mean	54.394558	0.683333	0.966330	131.505017	246.084175	0.147157	0.526667	149.505155	0.328859	1.046000	1.404040	0.734899	2.31
std	9.106168	0.465953	1.035947	17.502516	52.016723	0.354856	0.526192	22.824574	0.470589	1.163729	0.613792	1.028315	0.61
min	29.000000	0.000000	0.000000	94.000000	126.000000	0.000000	0.000000	71.000000	0.000000	0.000000	0.000000	0.000000	0.00
25%	47.250000	0.000000	0.000000	120.000000	211.000000	0.000000	0.000000	133.500000	0.000000	0.000000	1.000000	0.000000	2.00
50%	56.000000	1.000000	1.000000	130.000000	240.000000	0.000000	1.000000	153.000000	0.000000	0.800000	1.000000	0.000000	2.00
75%	61.000000	1.000000	2.000000	140.000000	274.000000	0.000000	1.000000	166.000000	1.000000	1.600000	2.000000	1.000000	3.00
max	77.000000	1.000000	3.000000	200.000000	564.000000	1.000000	2.000000	195.000000	1.000000	6.200000	2.000000	4.000000	3.00



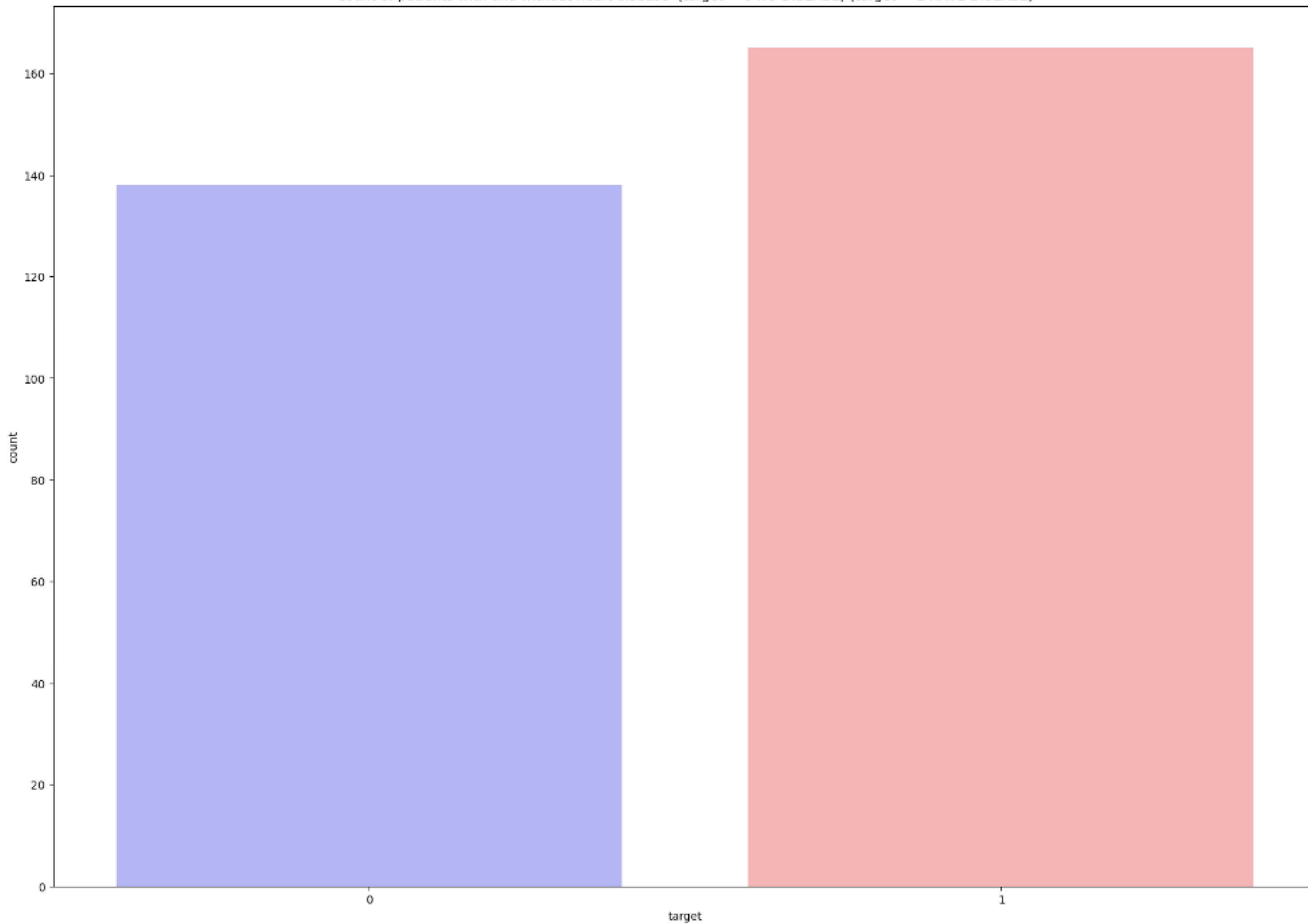


Correlation Matrices



Count plot histogram for the feature "Age" to find out Age Distribution

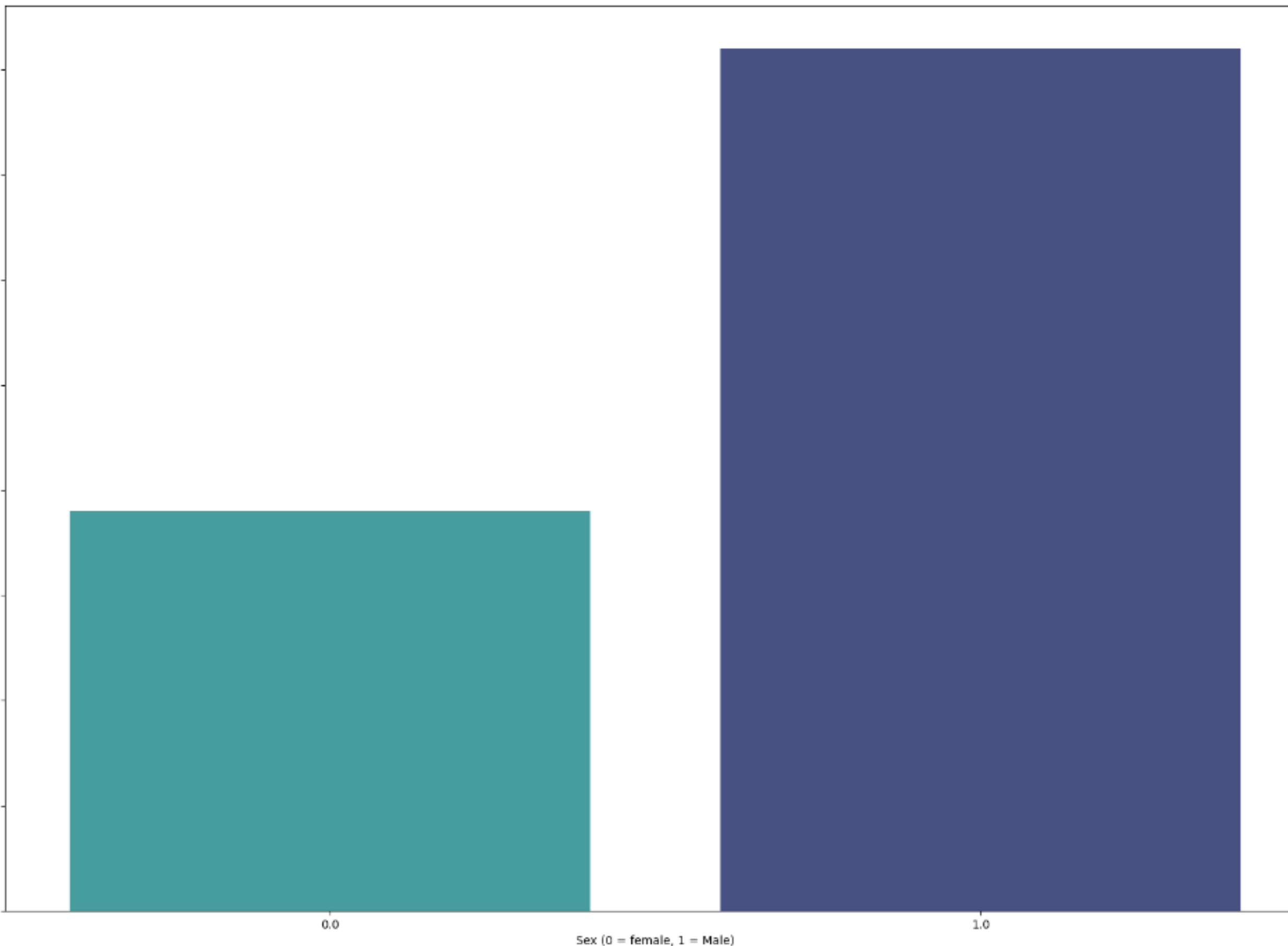
count of patients with and without heart disease (target = 0 NO DISEASE) (target = 1 HAVE DISEASE)



Target Count
0 = No Disease
1 = Have Disease

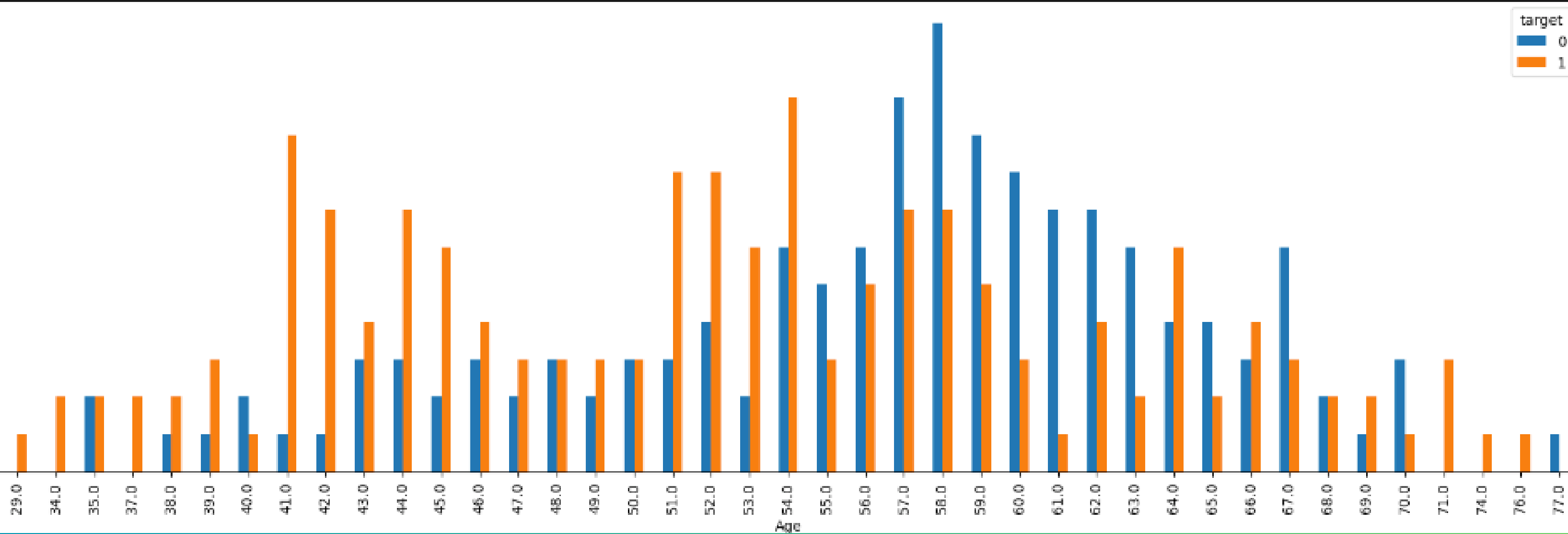
Percentage of Patients Haven't Heart Disease: 45.54%

Percentage of Patients Have Heart Disease: 54.46%



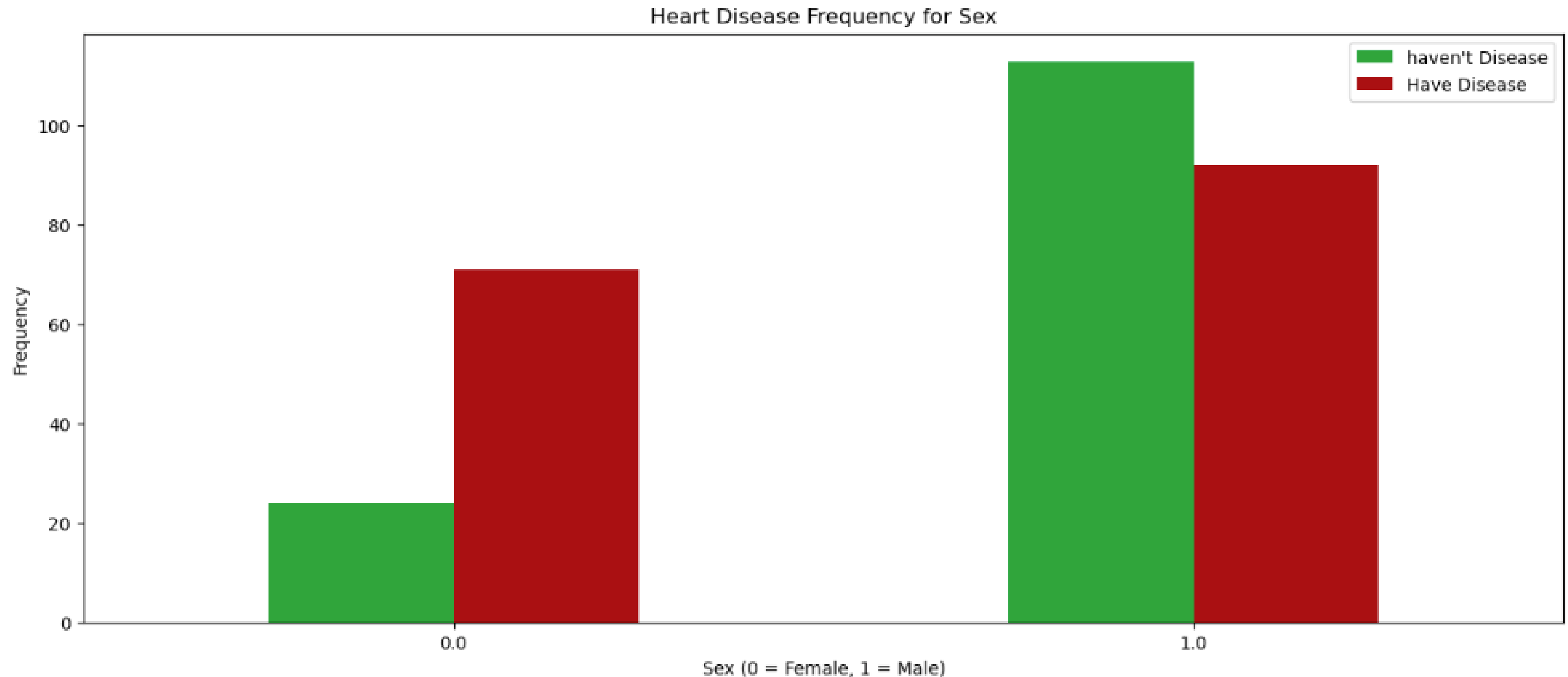
**Count of
male and
female**

heart Disease Frequency for Ages



**In this graphs we can see
the frequency of heart
disease by Age**

Frequency of heart by Sex



Impute Missing Values

Impute missing values:

EndTailImputer

```
1 endTailImpute = EndTailImputer(  
2     imputation_method = 'gaussian',  
3     tail = 'right'  
4 )  
5  
6 endTailImpute.fit(X_train)
```

EndTailImputer()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

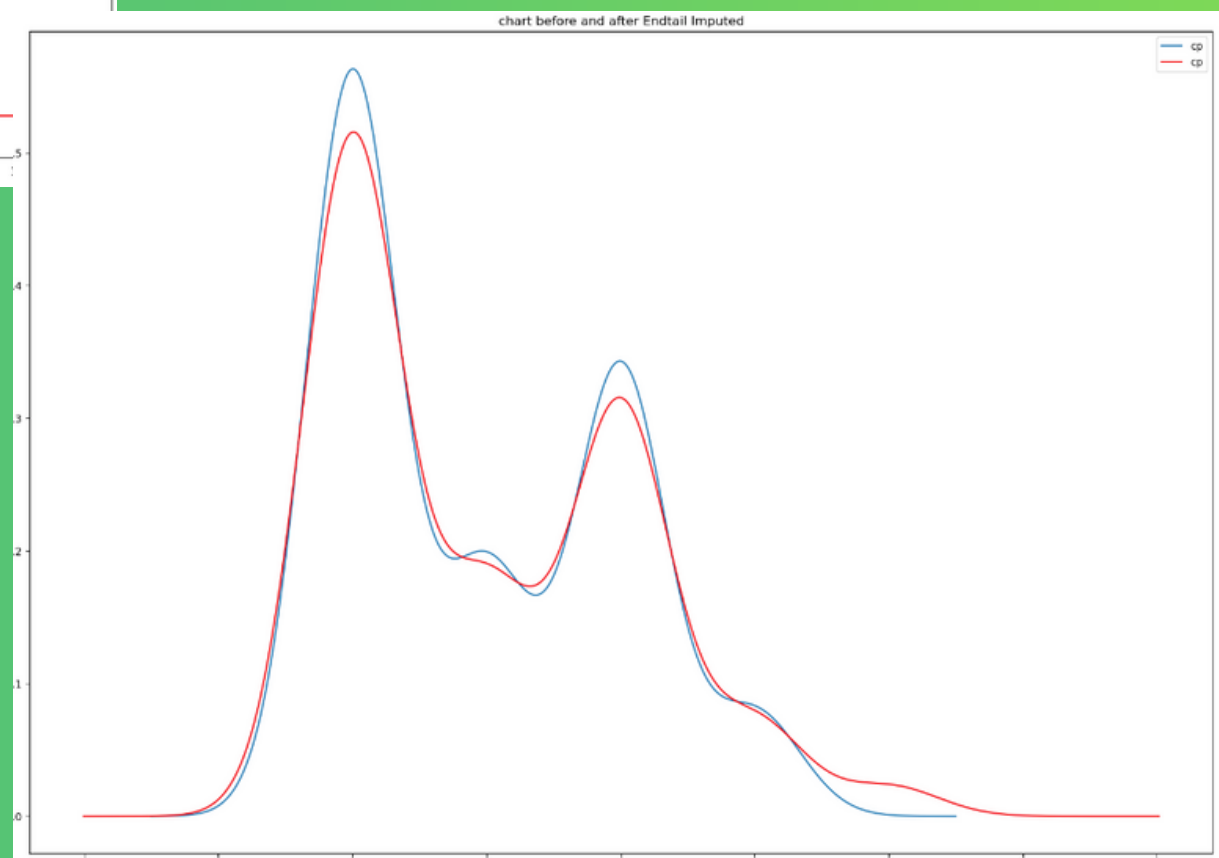
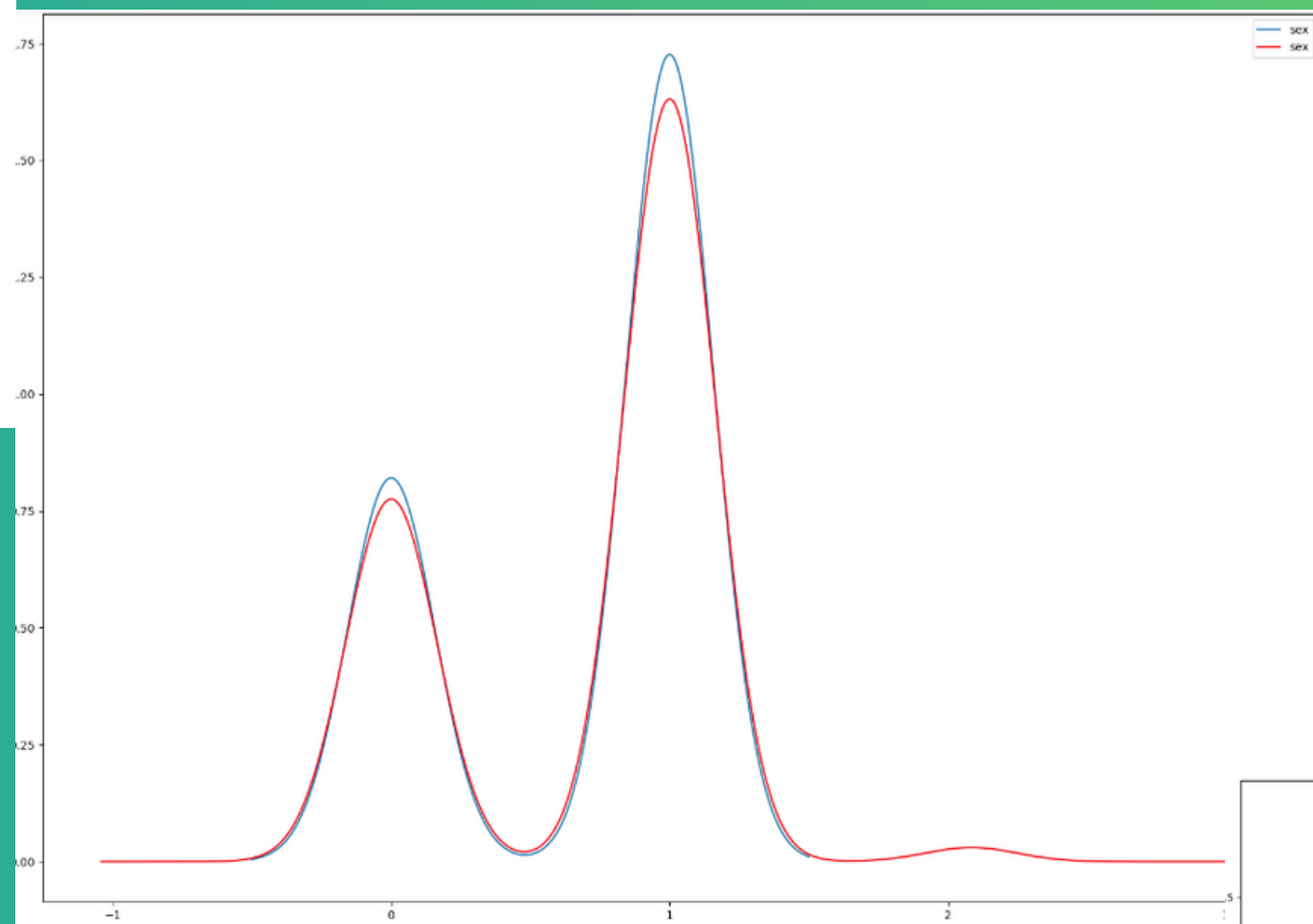
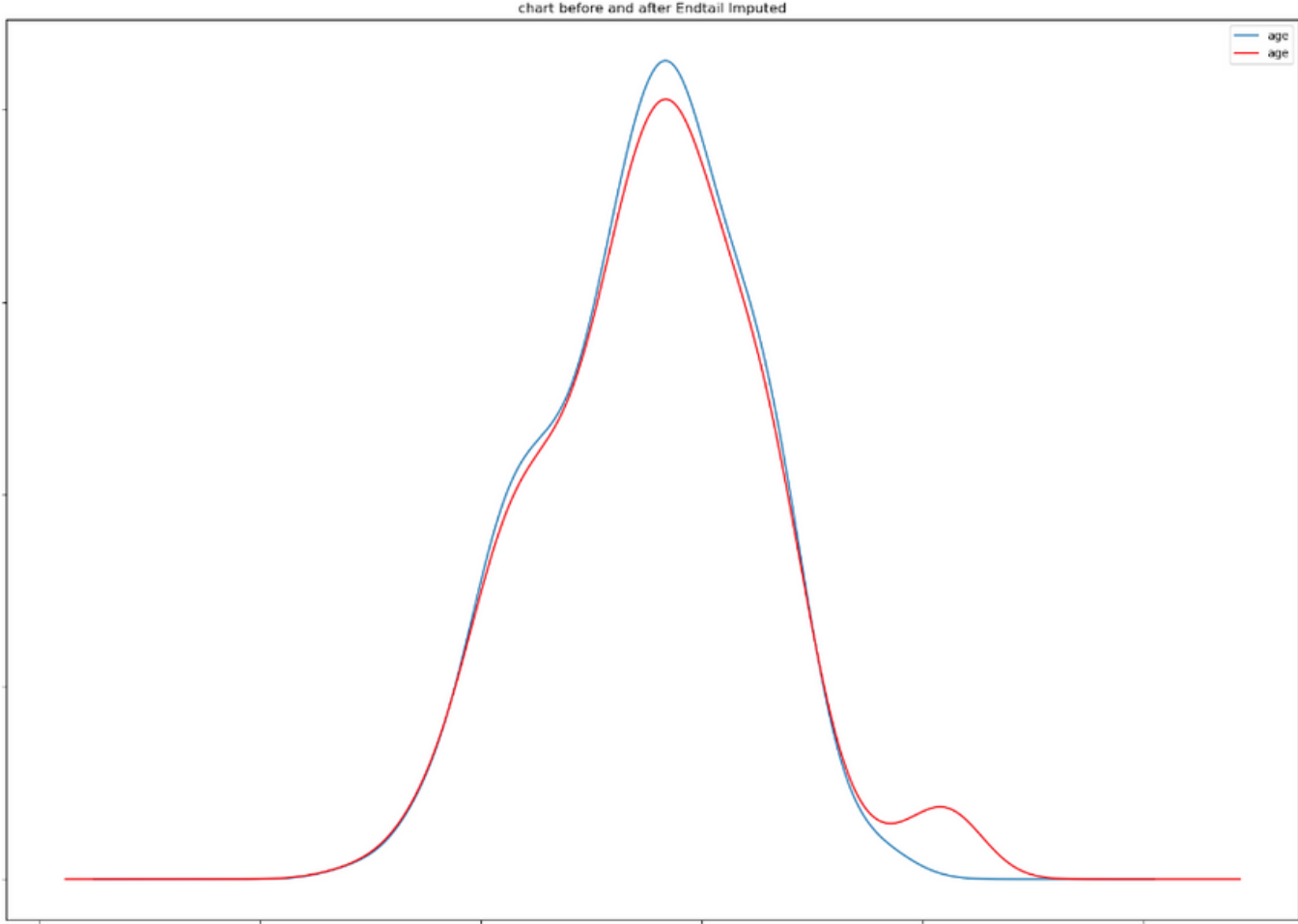
EndTail Imputer

```
1 endTailImpute.variables_
```

```
[ 'age',  
  'sex',  
  'cp',  
  'trestbps',  
  'chol',  
  'fbs',  
  'restecg',  
  'thalach',  
  'exang',  
  'oldpeak',  
  'slope',  
  'ca',  
  'thal' ]
```

```
1 endTailImpute.imputer_dict_
```

```
{ 'age': 82.18659080810232,  
  'sex': 2.0826960696630894,  
  'cp': 4.009833273341641,  
  'trestbps': 186.22514160368743,  
  'chol': 402.3367006620612,  
  'fbs': 1.228433700824096,  
  'restecg': 2.082977941606012,  
  'thalach': 216.35208813034708,  
  'exang': 1.7843474629507259,  
  'oldpeak': 4.542578580266698,  
  'slope': 3.2423135429602743,  
  'ca': 3.6070212966202324,  
  'thal': 4.177213690998737 }
```



Arbitrary Number Imputer

```
1 imputer = ArbitraryNumberImputer(imputer_dict = {
2     'age': 99,
3     'sex': -1,
4     'cp': -1,
5     'trestbps': 99,
6     'chol': 222,
7     'fbs': -1,
8     'restecg': -1,
9     'thalach': 99,
10    'exang': -1,
11    'oldpeak': -1,
12    'slope': -1,
13    'ca': 2,
14    'thal': 3
15 })
16 imputer.fit(X_train)
```

Arbitrary Number Imputer

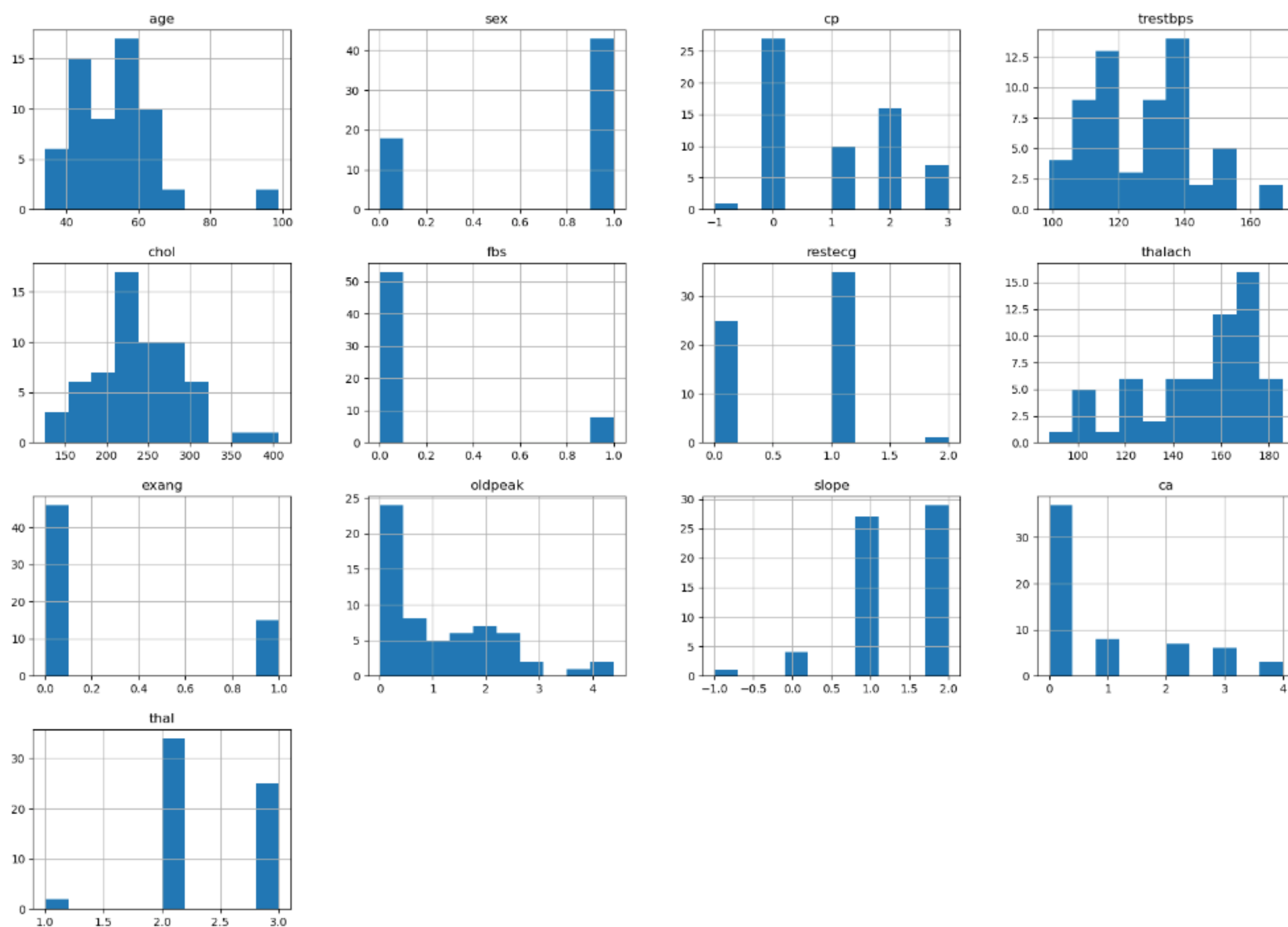
```
1 imputer.imputer_dict_

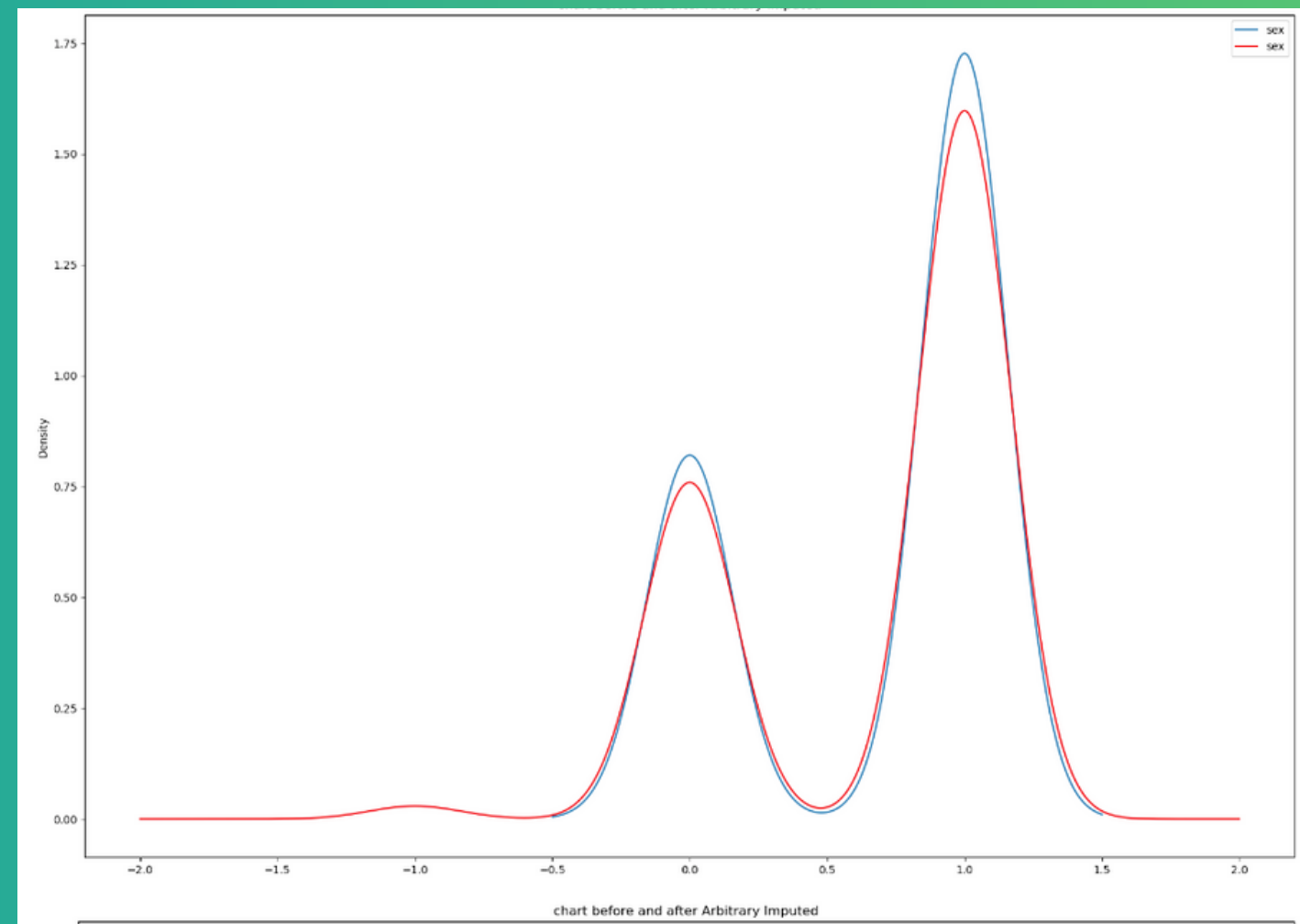
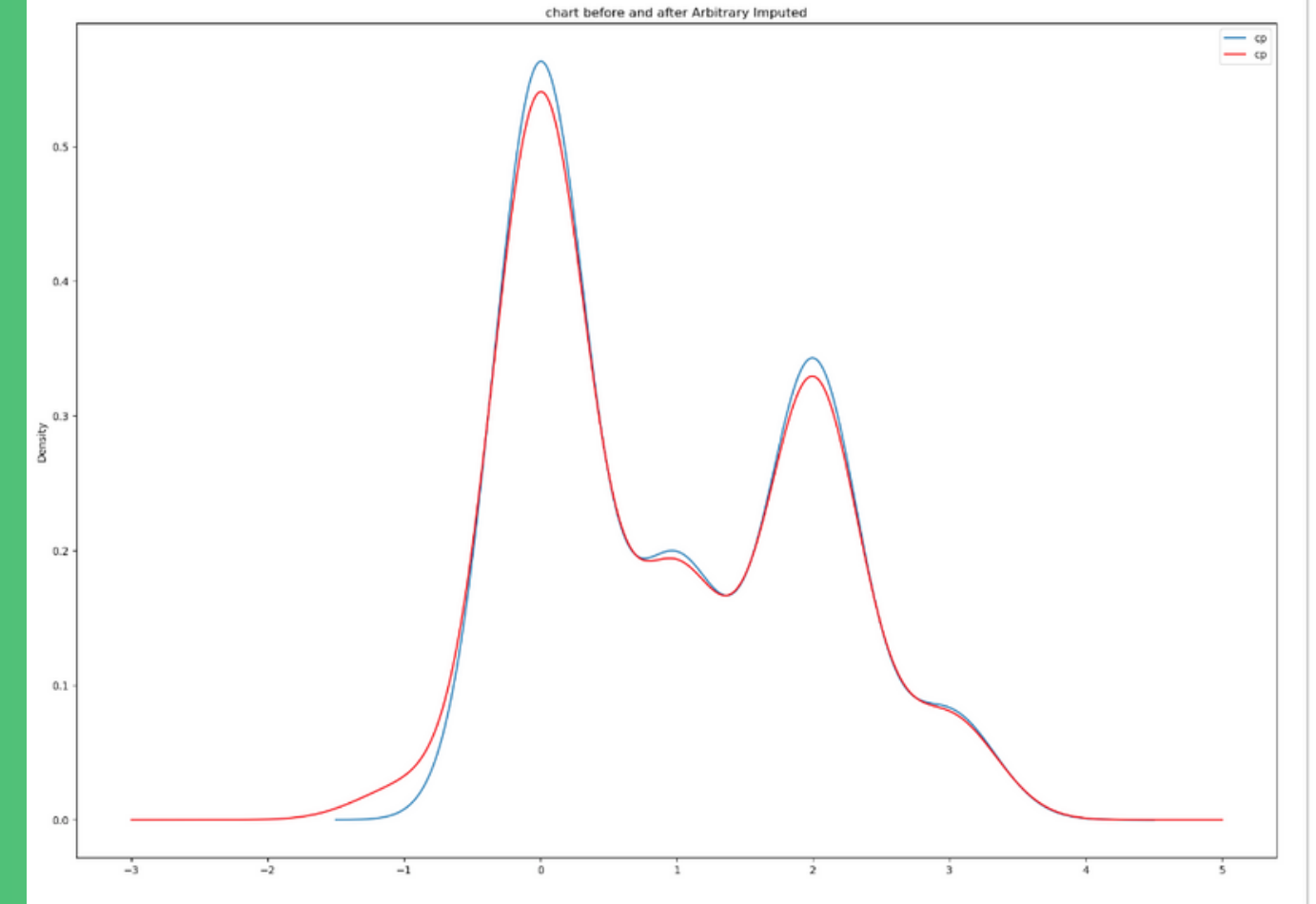
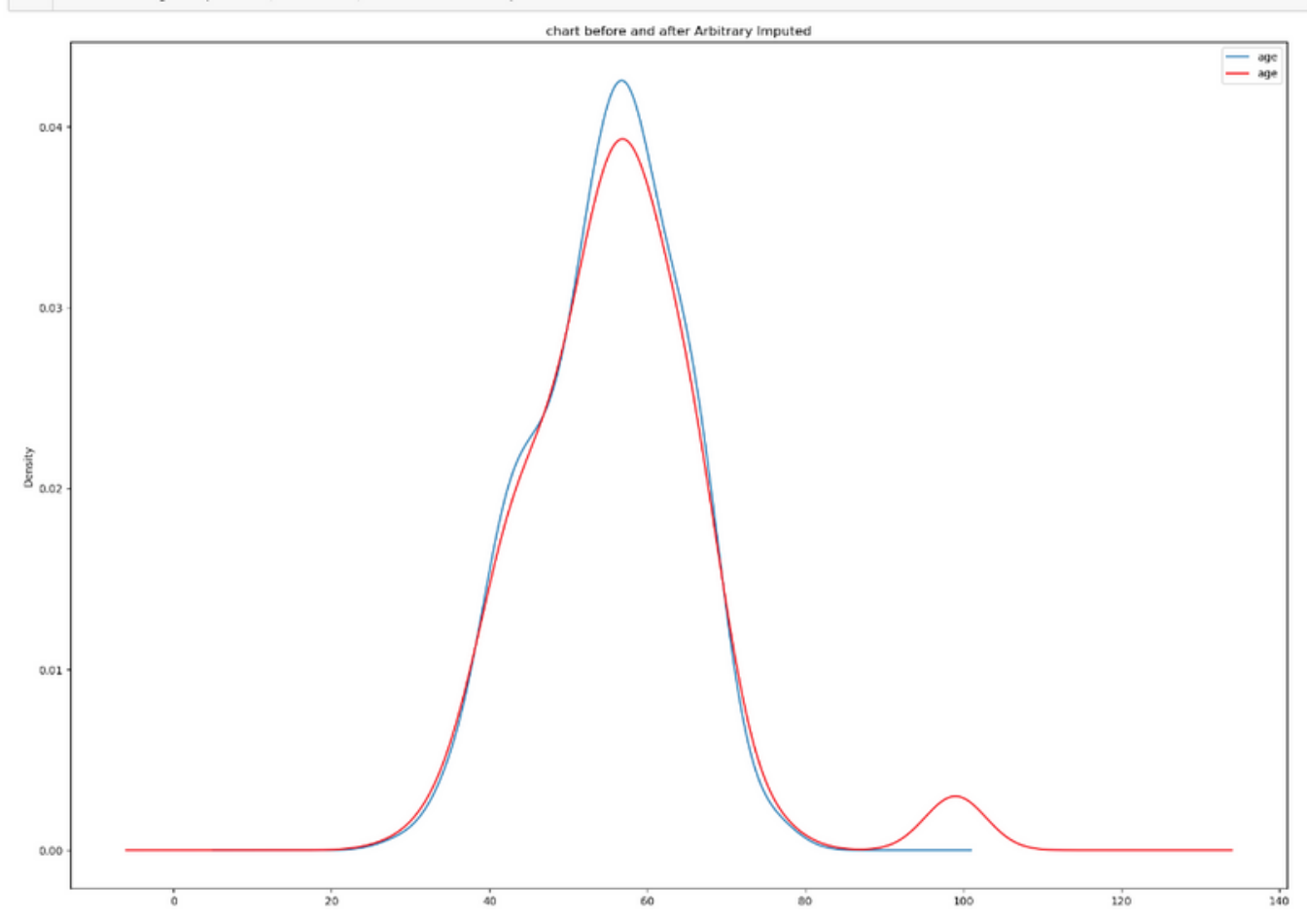
{'age': 99,
 'sex': -1,
 'cp': -1,
 'trestbps': 99,
 'chol': 222,
 'fbs': -1,
 'restecg': -1,
 'thalach': 99,
 'exang': -1,
 'oldpeak': -1,
 'slope': -1,
 'ca': 2,
 'thal': 3}
```

```
1 tmp_arbTrain = imputer.transform(X_train)
2 tmp_arbTest = imputer.transform(X_test)
3 tmp_arbTrain[imputer.variables_].isnull().mean()
```

```
age      0.0
sex      0.0
cp       0.0
trestbps 0.0
chol     0.0
fbs      0.0
restecg  0.0
thalach  0.0
exang    0.0
oldpeak  0.0
slope    0.0
ca       0.0
thal     0.0
dtype: float64
```

Histogram after Arbitrary Number Imputer





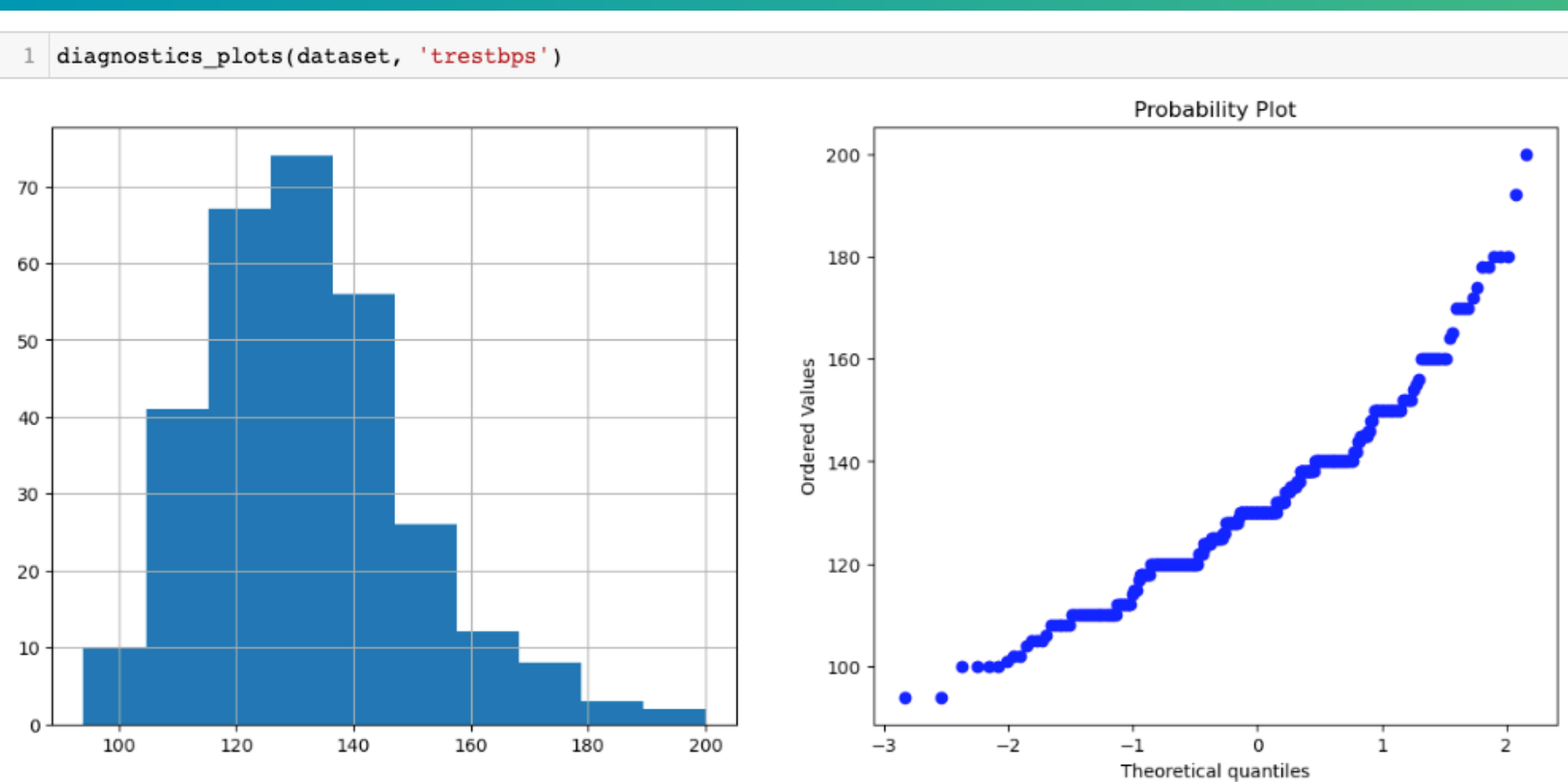
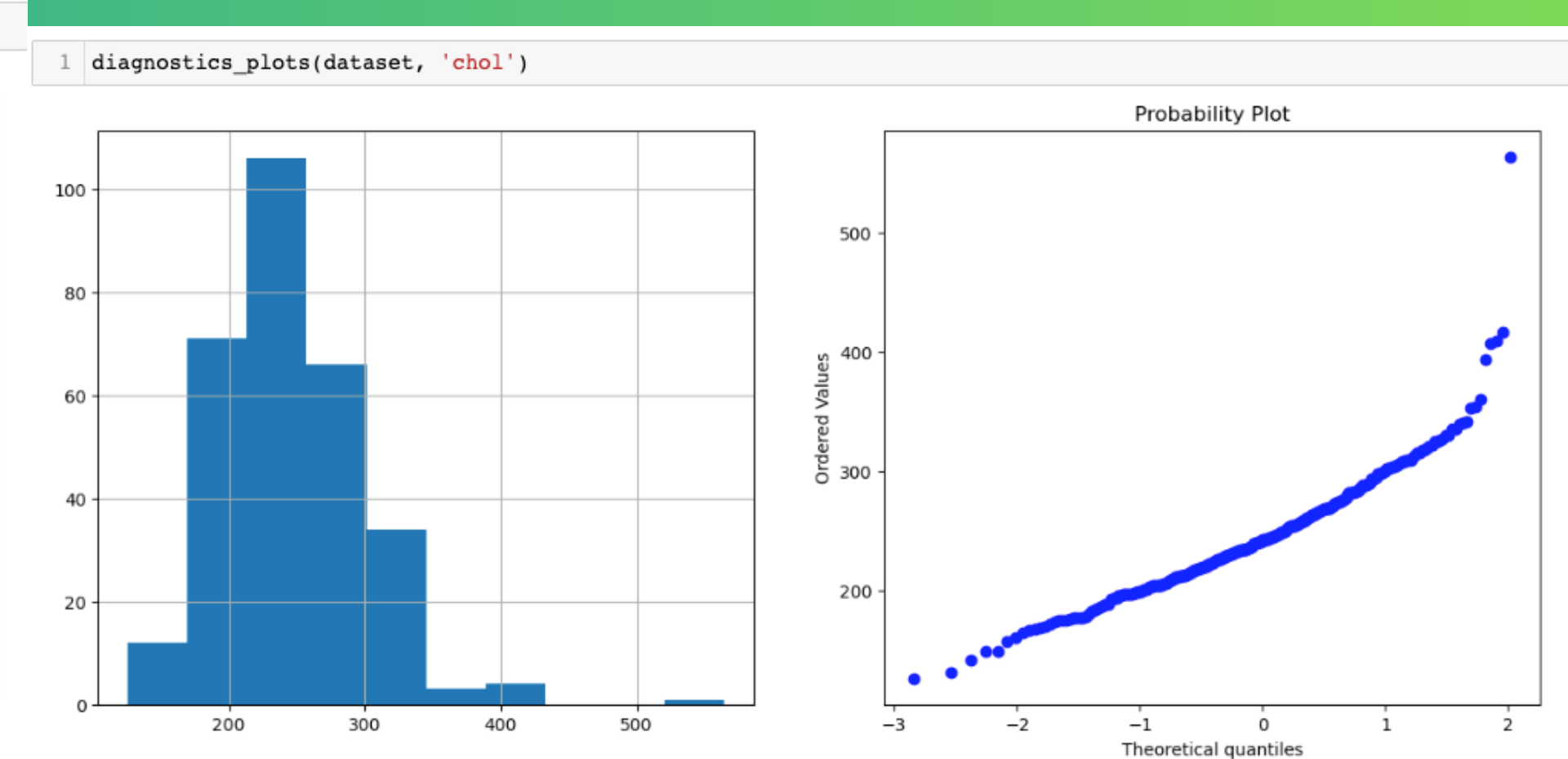
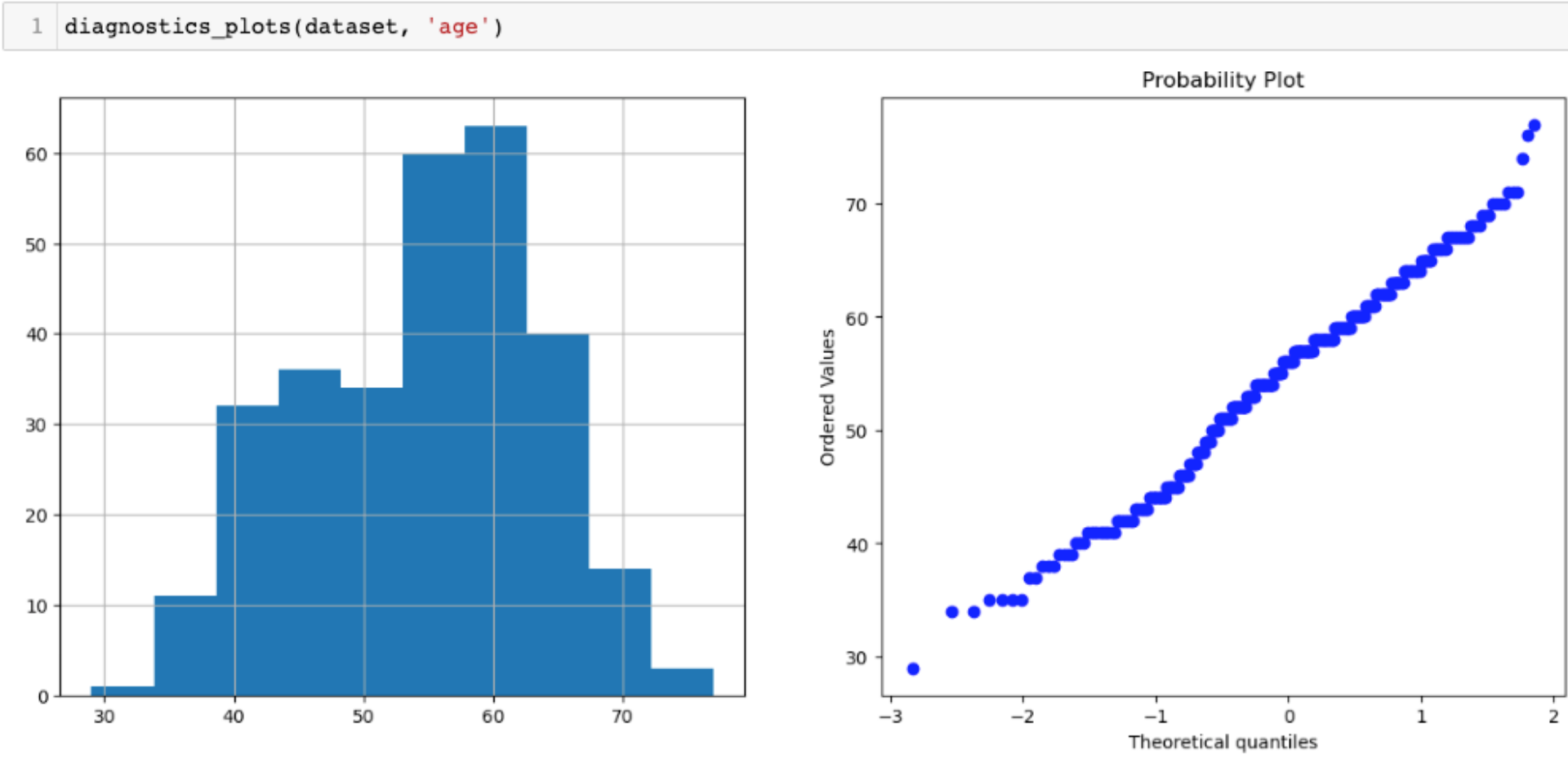
Transformation with:

Log Transformer

Exponential Transformer

Yeo-Johnson Transformer

```
def diagnostics_plots(dataset, variable):  
  
    plt.figure(figsize = (15, 6))  
  
    plt.subplot(1, 2, 1)  
    dataset[variable].hist()  
  
    plt.subplot(1, 2, 2)  
    stats.probplot(dataset[variable], dist = 'norm', plot = plt)  
  
    plt.show()
```

Before
Transformation

```
1 # Log Transform
2 lt = vt.LogCpTransformer(variables = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak'])
3 lt.fit(imputer_MeanMedian)
```

LogCpTransformer(variables=['age', 'trestbps', 'chol', 'thalach', 'oldpeak'])

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.

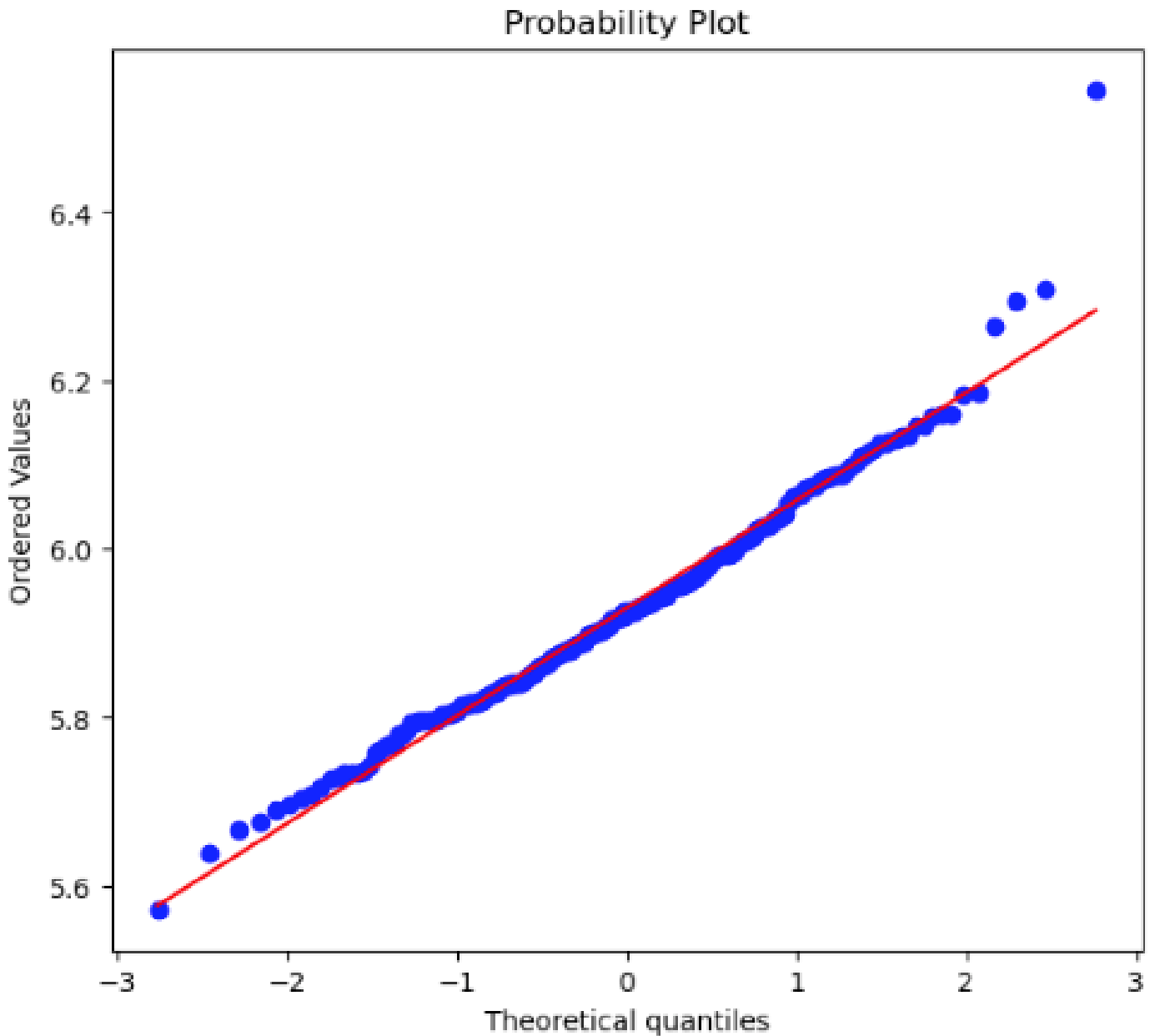
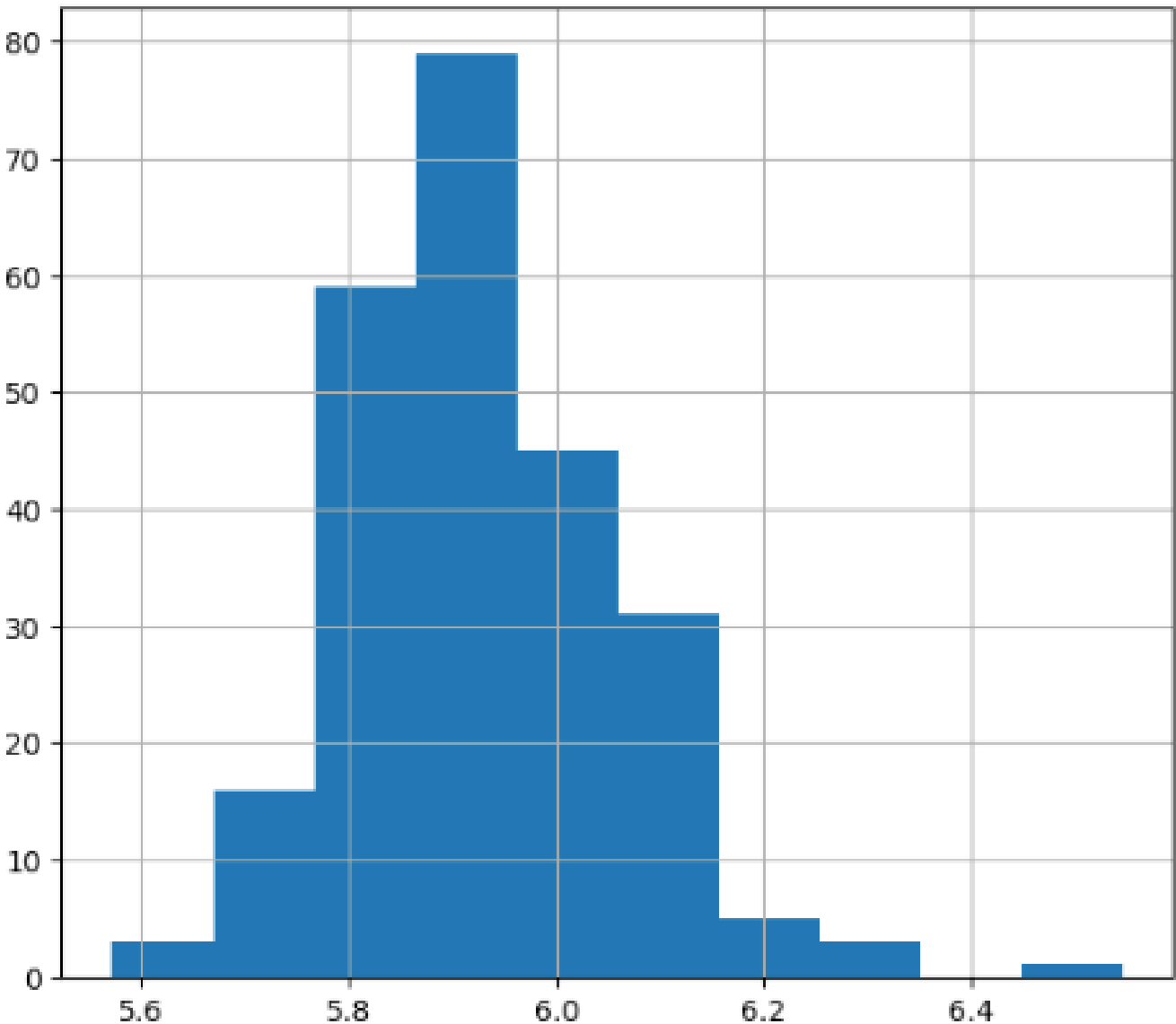
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
1 lt.variables_
```

['age', 'trestbps', 'chol', 'thalach', 'oldpeak']

```
1 data_tf = lt.transform(imputer_MeanMedian)
```

```
1 diagnostics_plots(data_tf, 'chol')
```



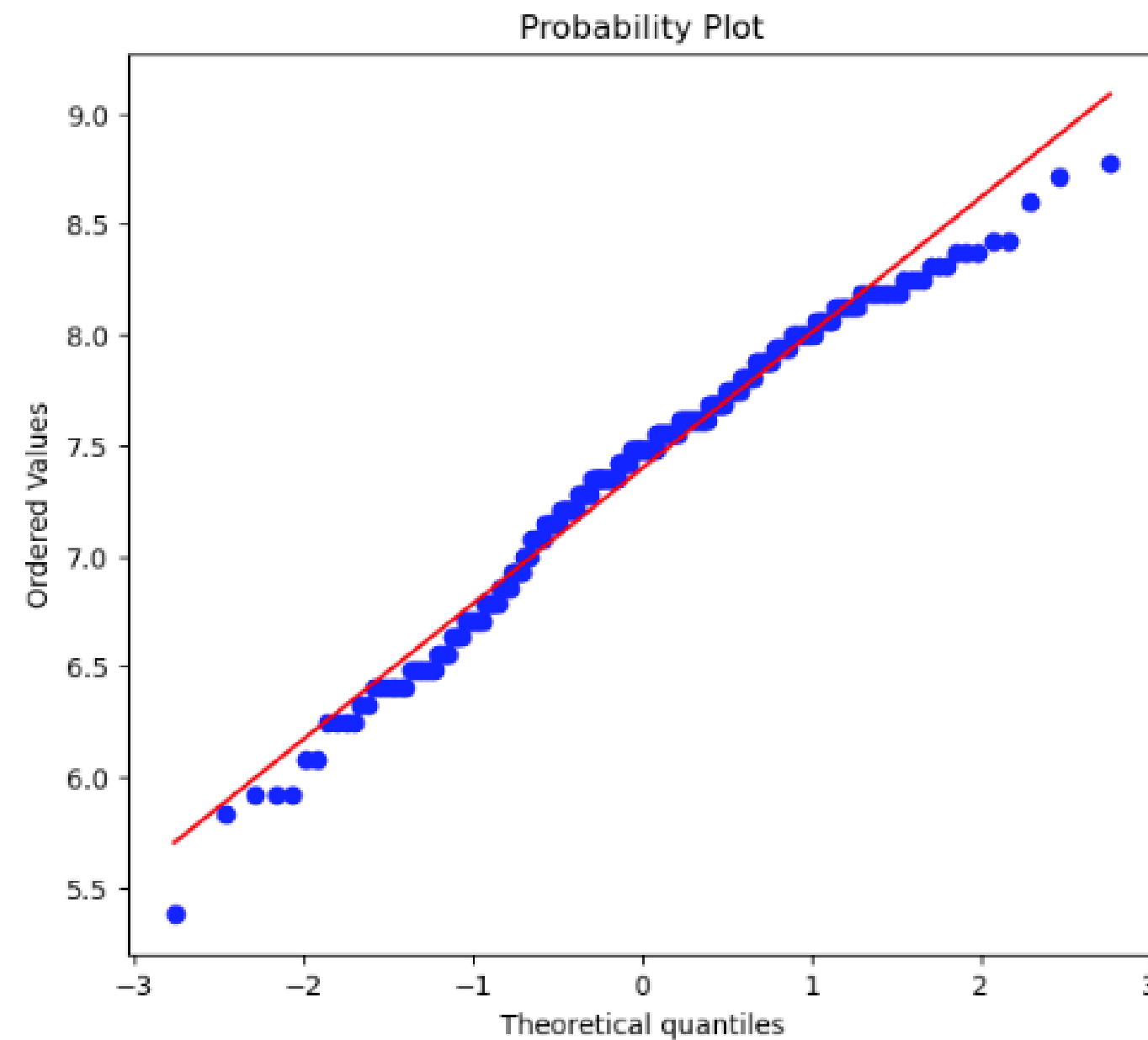
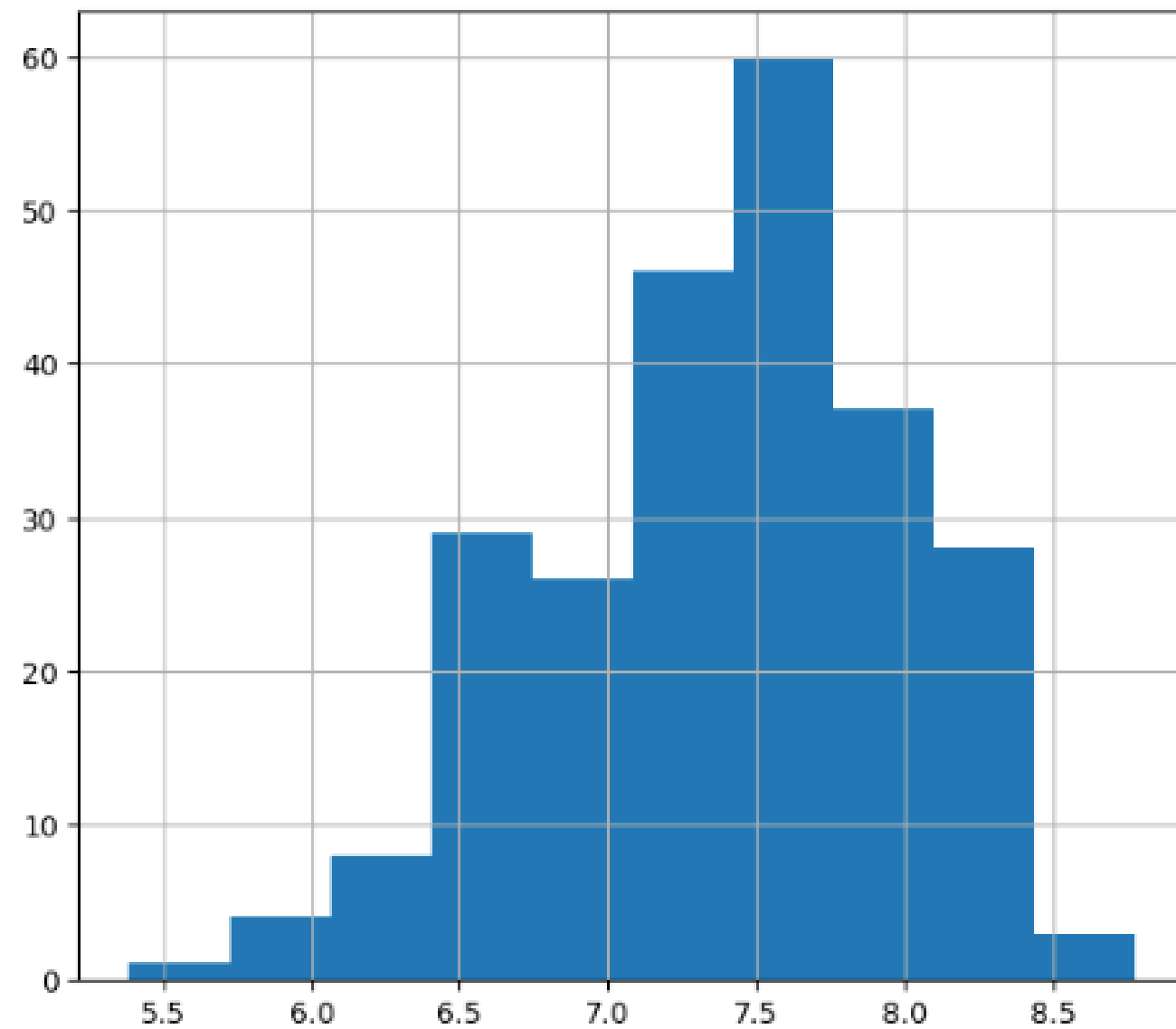
Log Transformer

```
1 #Exponential Transofrmer
2 et = vt.PowerTransformer(variables = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak'])
3 et.fit(imputer_MeanMedian)
```

PowerTransformer(variables=['age', 'trestbps', 'chol', 'thalach', 'oldpeak'])
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
1 data_tf = et.transform(imputer_MeanMedian)
```

```
1 diagnostics_plots(data_tf, 'age')
```



Exponential Transformer

```

1 #Yeo-Johnson
2 yjt = vt.YeoJohnsonTransformer(variables = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak'])
3 yjt.fit(imputer_MeanMedian)

```

```

YeoJohnsonTransformer(variables=['age', 'trestbps', 'chol', 'thalach',
                                   'oldpeak'])

```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```

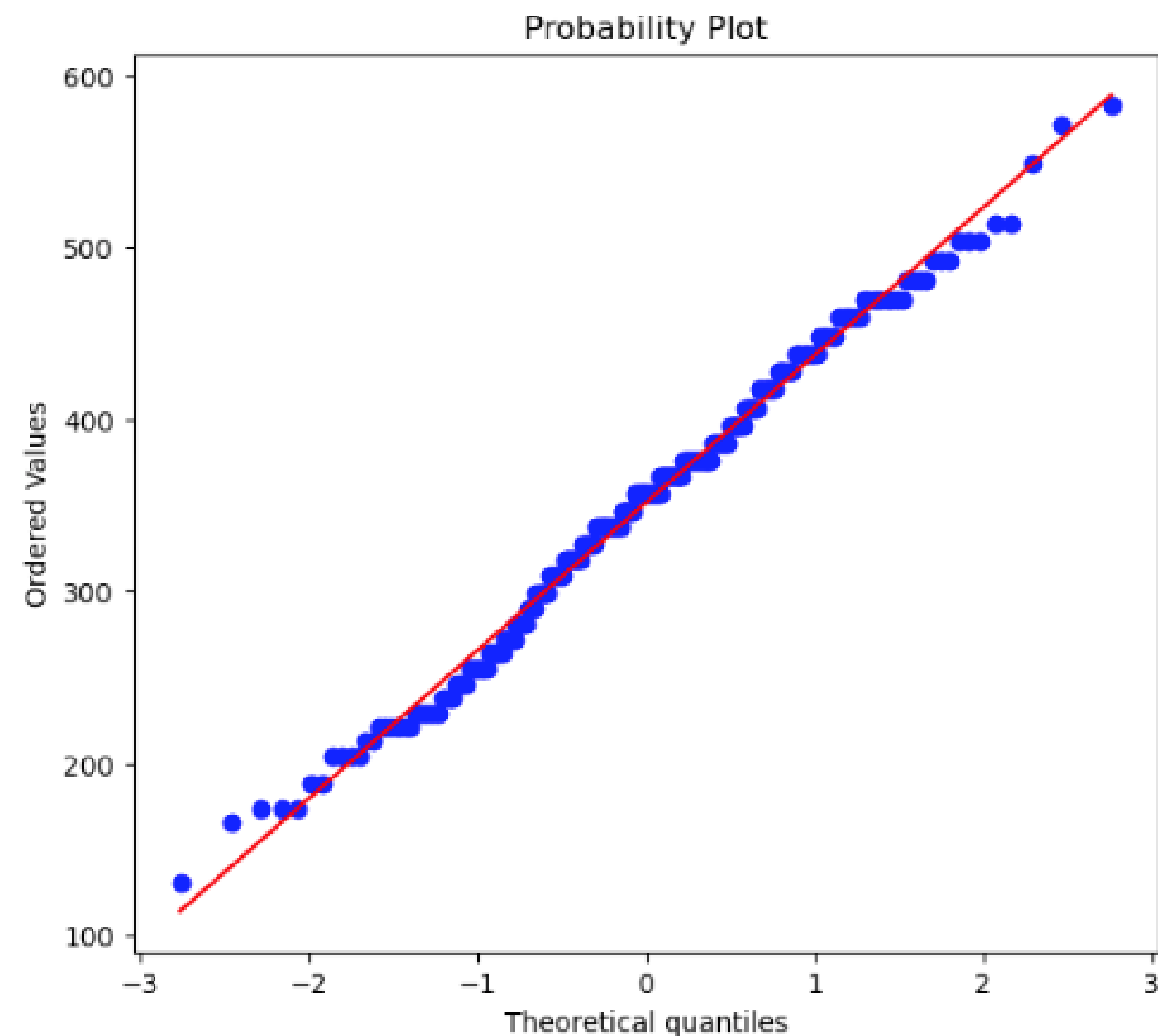
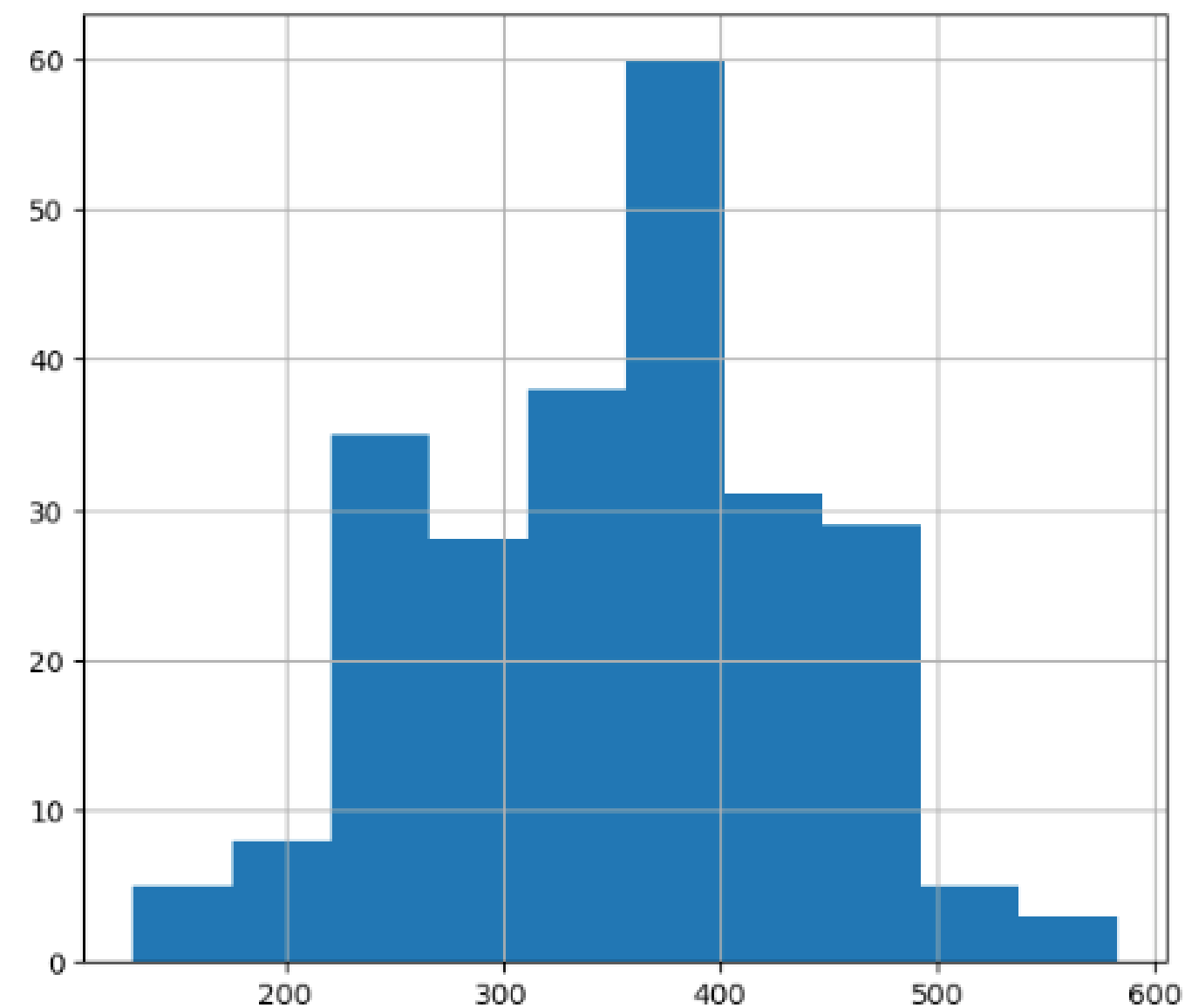
1 data_tf = yjt.transform(imputer_MeanMedian)

```

```

1 diagnostics_plots(data_tf, 'age')

```



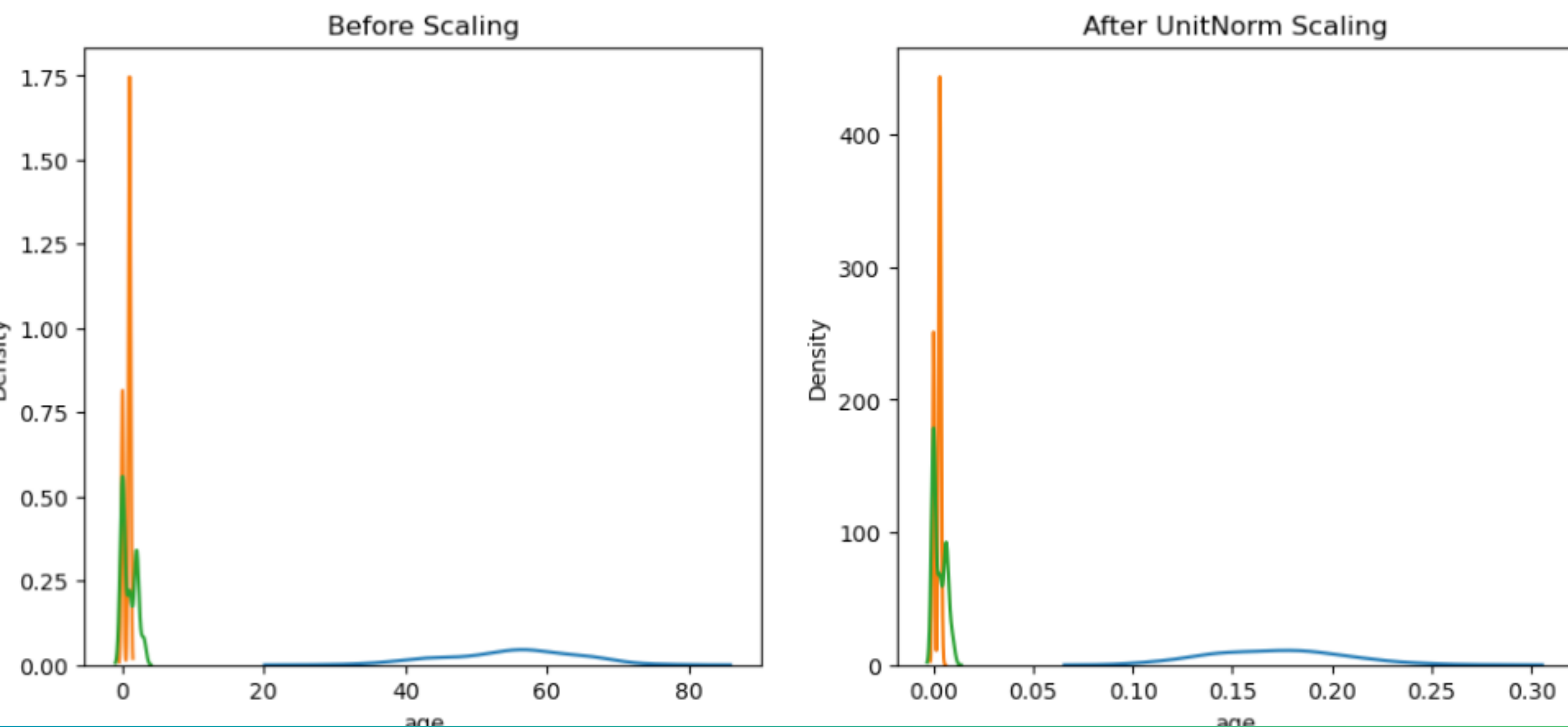
Yeo-Johnson Transformer

Scaling with Unit Norm, Robust Scaler


```
fig, (ax1, ax2) = plt.subplots(ncols = 2, figsize = (12, 5))
#before scaling
ax1.set_title('Before Scaling')
sns.kdeplot(imputer_MeanMedian['age'], ax = ax1)
sns.kdeplot(imputer_MeanMedian['sex'], ax = ax1)
sns.kdeplot(imputer_MeanMedian['cp'], ax = ax1)

#after scaling
ax2.set_title('After UnitNorm Scaling')
sns.kdeplot(X_train_scaled['age'], ax = ax2)
sns.kdeplot(X_train_scaled['sex'], ax = ax2)
sns.kdeplot(X_train_scaled['cp'], ax = ax2)
```

axesSubplot: title={'center': 'After UnitNorm Scaling'}, xlabel='age', ylabel='Density'>



Robust Scaler

```
1 scaler = RobustScaler()  
2  
3 scaler.fit(imputer_MeanMedian)  
4  
5 X_train_scaled = scaler.transform(imputer_MeanMedian)  
6 X_test_scaled = scaler.transform(imputer_MeanMedianTest)
```

```
1 X_train_scaled = pd.DataFrame(X_train_scaled, columns = imputer_MeanMedian.columns)  
2 X_test_scaled = pd.DataFrame(X_test_scaled, columns = imputer_MeanMedianTest.columns)
```

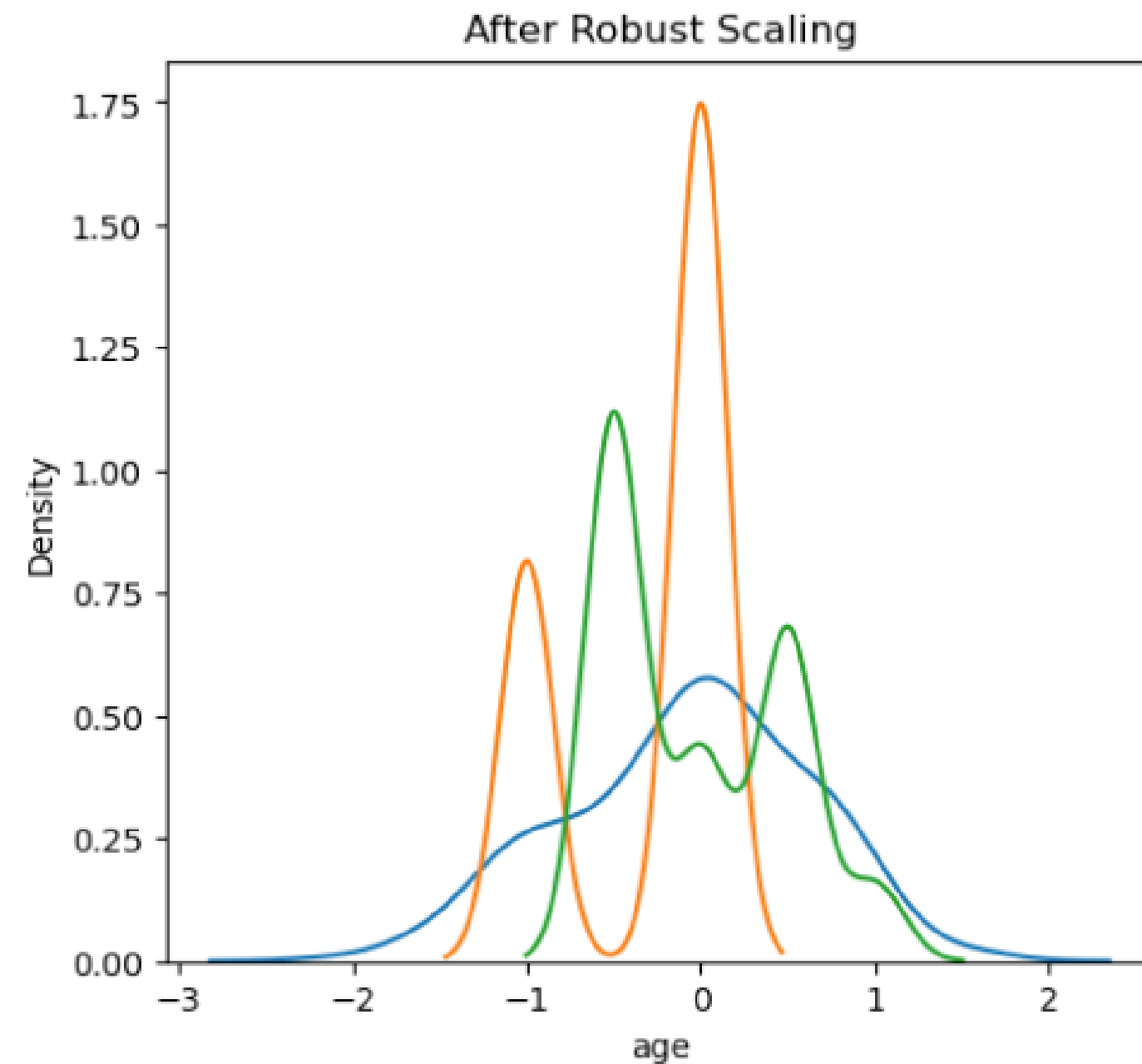
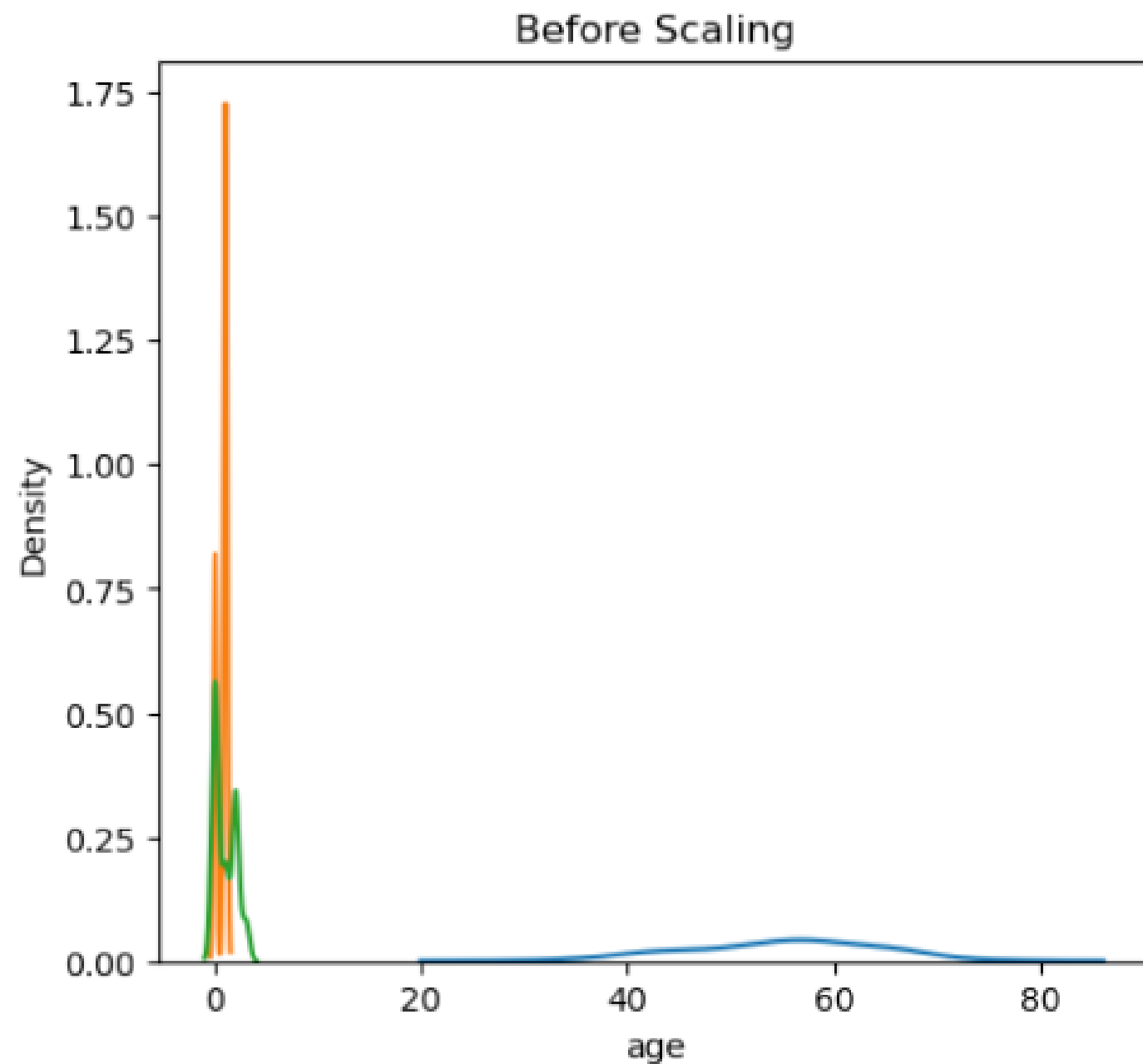
```
1 X_train_scaled
```

	age	sex	cp	trestbps	chol	fbs	restecg	thalach	exang	oldpeak	slope	ca	thal
0	1.019608	0.0	0.5	0.5	0.178138	0.0	0.0	-0.188034	0.0	0.7500	0.0	3.0	1.0
1	0.627451	-1.0	-0.5	2.5	1.327935	0.0	1.0	0.000000	1.0	-0.5000	1.0	0.0	0.0
2	-0.862745	-1.0	0.0	0.0	-0.145749	0.0	0.0	0.803419	0.0	-0.1250	0.0	0.0	0.0
3	0.000000	-1.0	0.5	0.5	2.817814	1.0	0.0	0.188034	0.0	0.0000	0.0	1.0	0.0
4	0.784314	0.0	0.0	1.5	0.048583	0.0	1.0	-1.076923	1.0	-0.5000	0.0	3.0	-1.0
...
237	0.392157	0.0	0.5	0.0	0.000000	1.0	1.0	-0.495726	1.0	0.1250	0.0	0.0	0.0
238	0.941176	0.0	0.5	-0.6	0.550607	0.0	1.0	-0.017094	0.0	0.1250	1.0	1.0	1.0
239	-0.862745	-1.0	0.0	-0.9	-1.344130	0.0	1.0	-0.461538	0.0	0.0000	0.0	0.0	0.0
240	-0.078431	0.0	-0.5	0.0	-0.421053	0.0	1.0	-1.384615	1.0	3.0000	-1.0	0.0	1.0
241	0.156863	0.0	-0.5	-1.5	-0.145749	0.0	1.0	0.153846	0.0	-0.4375	1.0	1.0	1.0

242 rows x 13 columns

Robust Scaler

```
1 fig, (ax1, ax2) = plt.subplots(ncols = 2, figsize =(12, 5))
2
3 #before scaling
4
5 ax1.set_title('Before Scaling')
6 sns.kdeplot(X_train['age'], ax = ax1)
7 sns.kdeplot(X_train['sex'], ax = ax1)
8 sns.kdeplot(X_train['cp'], ax = ax1)
9
10 #after Scaling
11 ax2.set_title('After Robust Scaling')
12 sns.kdeplot(X_train_scaled['age'], ax = ax2)
13 sns.kdeplot(X_train_scaled['sex'], ax = ax2)
14 sns.kdeplot(X_train_scaled['cp'], ax = ax2)
15
16 plt.show()
```



Discretization

with

EqualFrequency,

EqualWidthDisc

```
1 # EqualFrequencyDiscretizer
2 disc = EqualFrequencyDiscretiser(
3     q = 10,
4     variables = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
5 )
6
7 disc.fit(imputer_MeanMedian)
```

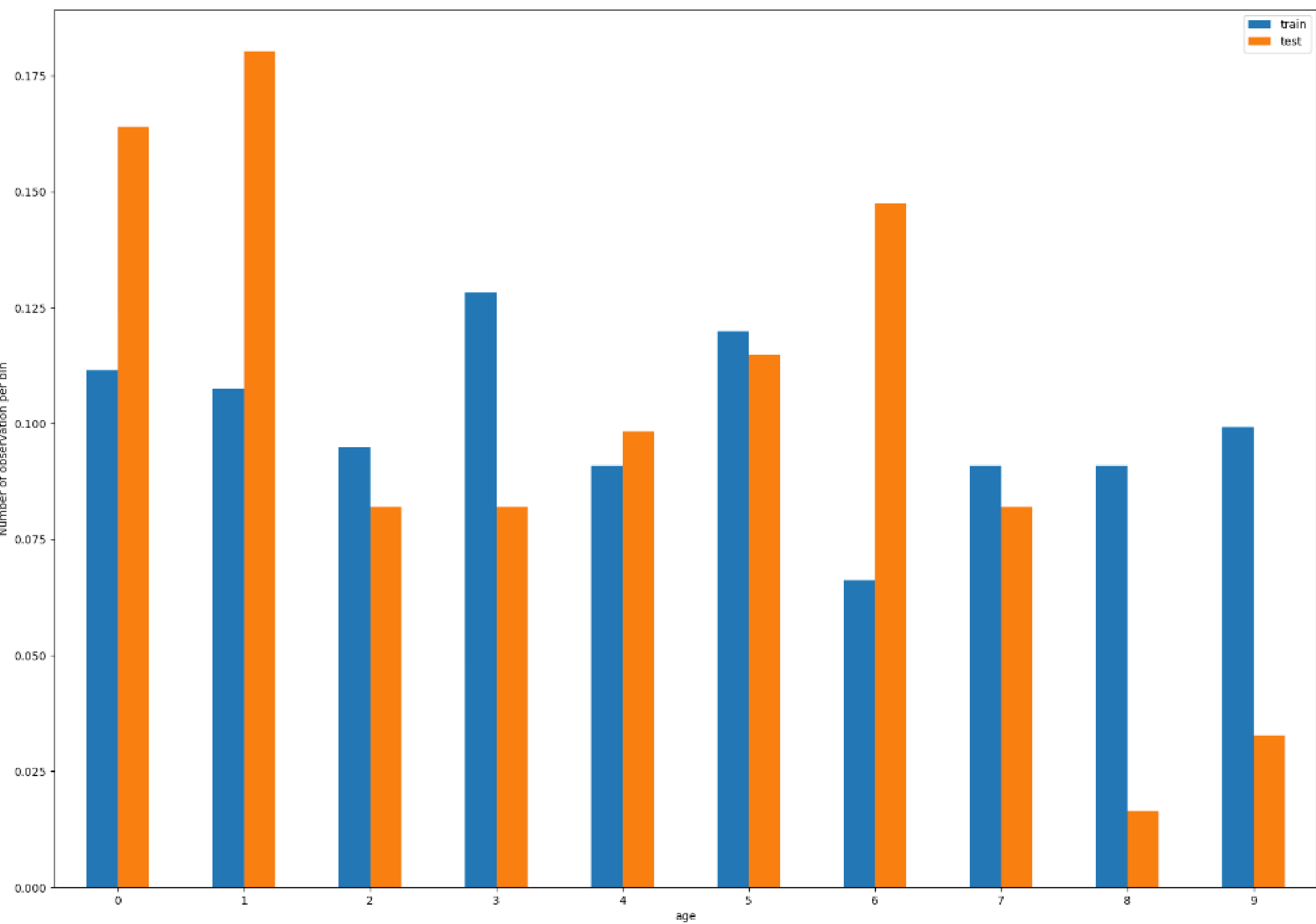
```
EqualFrequencyDiscretiser(variables=['age', 'trestbps', 'chol', 'thalach',
                                     'oldpeak'])
```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```
1 disc.binner_dict_
```

```
{'age': [-inf, 42.0, 47.0, 51.0, 54.0, 56.0, 58.0, 60.0, 63.0, 66.0, inf],
'trestbps': [-inf,
110.0,
120.0,
122.0,
128.0,
130.0,
132.0,
140.0,
144.8,
154.9,
inf],
'chol': [-inf,
196.0,
206.0,
218.3,
229.4,
243.0,
252.60000000000002,
266.70000000000005,
283.0,
309.0,
inf],
'thalach': [-inf,
116.0,
130.2,
140.0,
145.4,
151.5,
156.0,
162.0,
166.8,
174.9,
inf],
```

Equal Frequency




```

1 #EqualWidthDiscretizer
2 disc = EqualWidthDiscretiser(
3     bins = 10,
4     variables = ['age', 'trestbps', 'chol', 'thalach', 'oldpeak']
5 )
6
7 disc.fit(imputer_MeanMedian)

```

```

EqualWidthDiscretiser(variables=['age', 'trestbps', 'chol', 'thalach',
                                'oldpeak'])

```

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

```

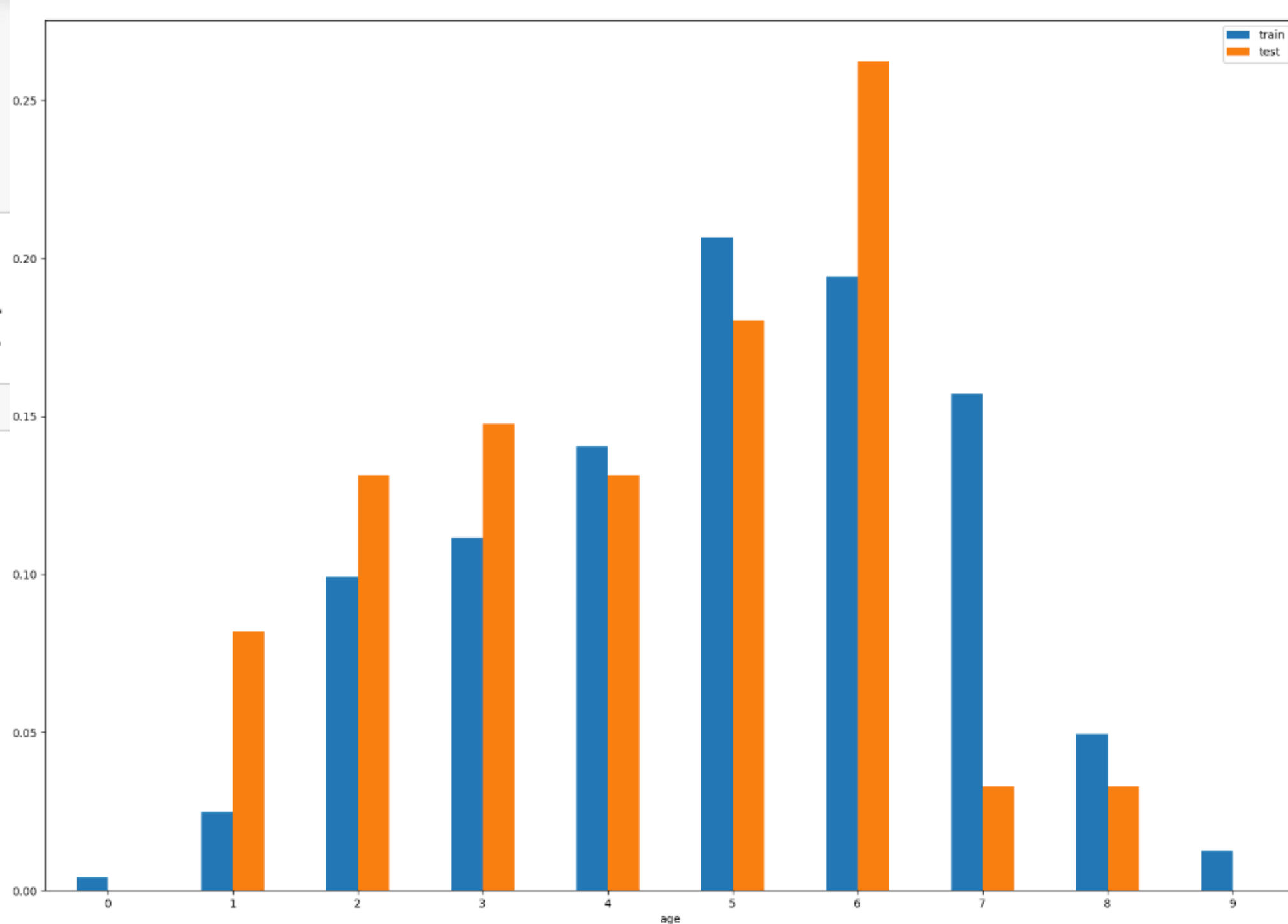
1 disc.binner_dict_

```

```

{'age': [-inf,
 33.8,
 38.6,
 43.4,
 48.2,
 53.0,
 57.8,
 62.6,
 67.4,
 72.19999999999999,
 inf],
'trestbps': [-inf,
 104.6,
 115.2,
 125.8,
 136.4,
 147.0,
 157.6,
 168.2,
 178.8,
 189.39999999999998,
 inf],
'chol': [-inf,
 174.3,
 217.6,
 260.9,
 304.2,
 347.5,
 390.79999999999995,
 434.09999999999997,
 477.4,
 520.7,
 inf],
'thalach': [-inf,

```



EqualWidth Discretization

Model Results

1	metrics_score(X_train, X_test, y_train, y_test, best_DT_model)				
Roc-auc score					
Train set: 82.039%					
Test set: 81.944%					
Accuracy score					
Train set: 75.583%					
Test set: 85.238%					
Standart deviation					
Train set: 7.773%					
Test set: 17.236%					
Classification report					
	precision	recall	f1-score	support	
0	0.72	0.72	0.72	25	
1	0.81	0.81	0.81	36	
accuracy			0.77	61	
macro avg	0.76	0.76	0.76	61	
weighted avg	0.77	0.77	0.77	61	

Best Results

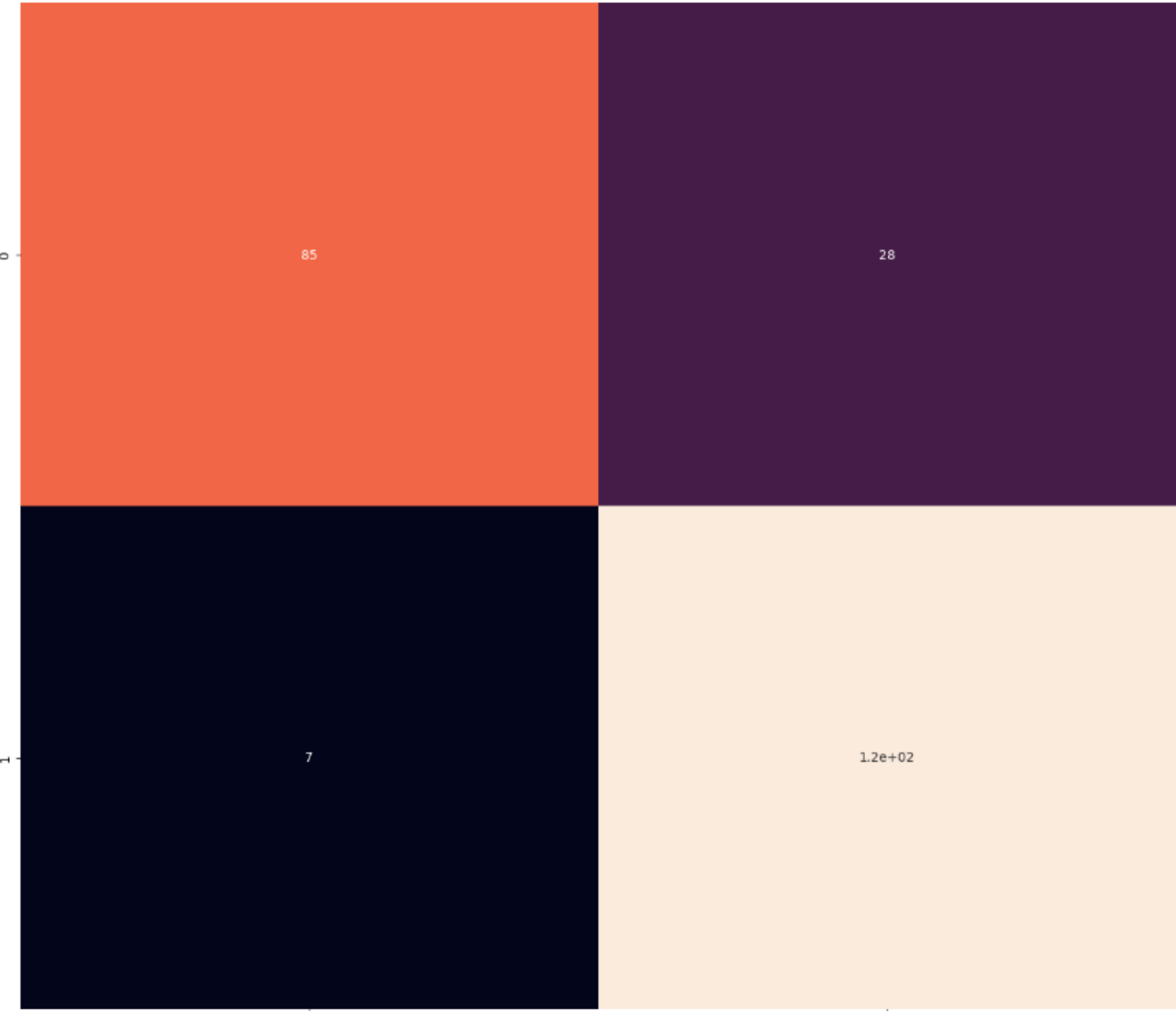
Decision

Tree

Model

```
1 y_pred_train = best_DT_model.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

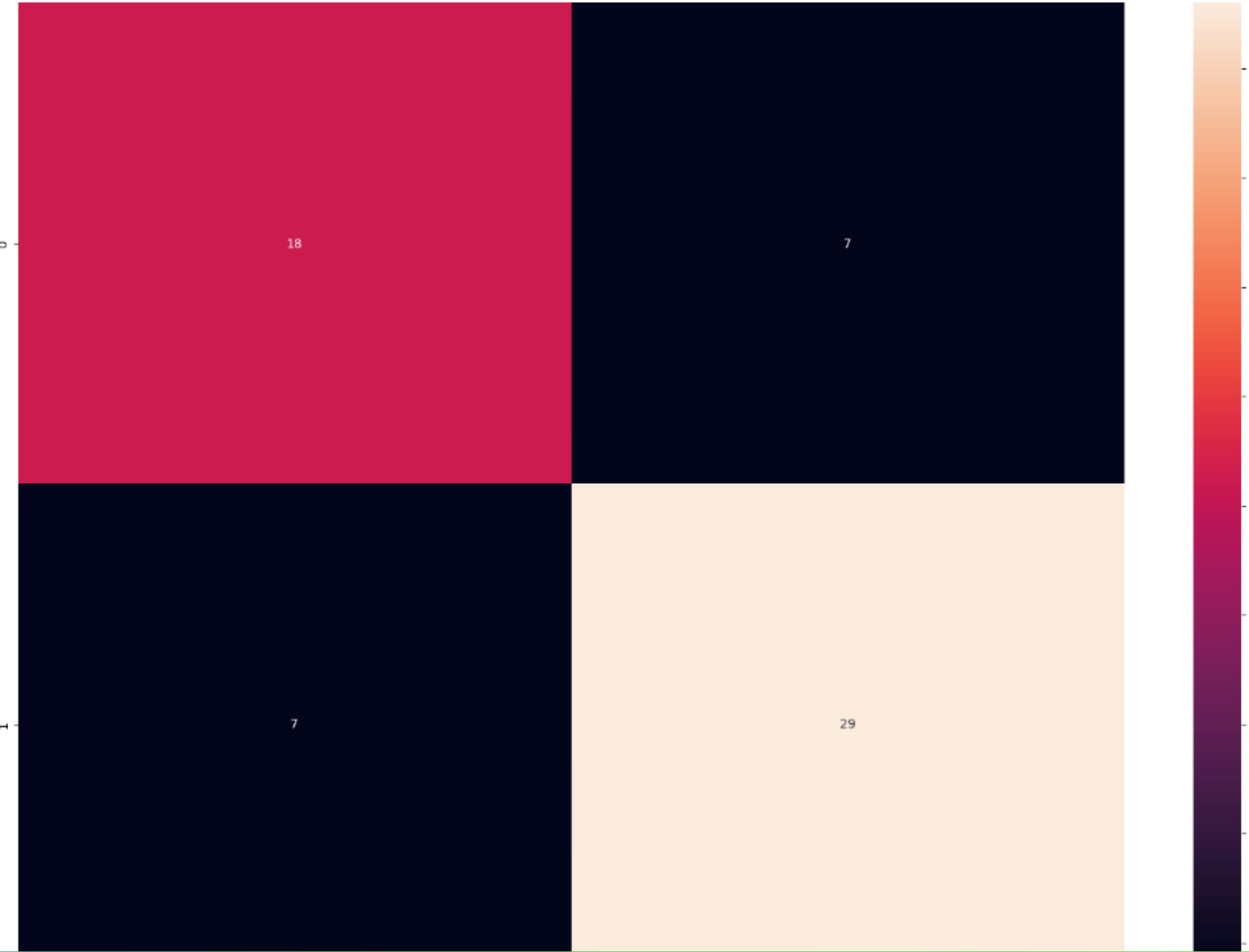
<AxesSubplot: >



Train

```
1 y_pred_test = best_DT_model.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```

<AxesSubplot: >



Test

Confusion Matrices

Decision Tree

Best Results Random Forest Classifier

```
1 metrics_score(X_train, X_test, y_train, y_test, rf_best_model)
```

Roc-auc score

```
Train set: 90.170%
```

```
Test set: 90.556%
```

Accuracy score

```
Train set: 81.417%
```

```
Test set: 85.000%
```

Standart deviation

```
Train set: 5.101%
```

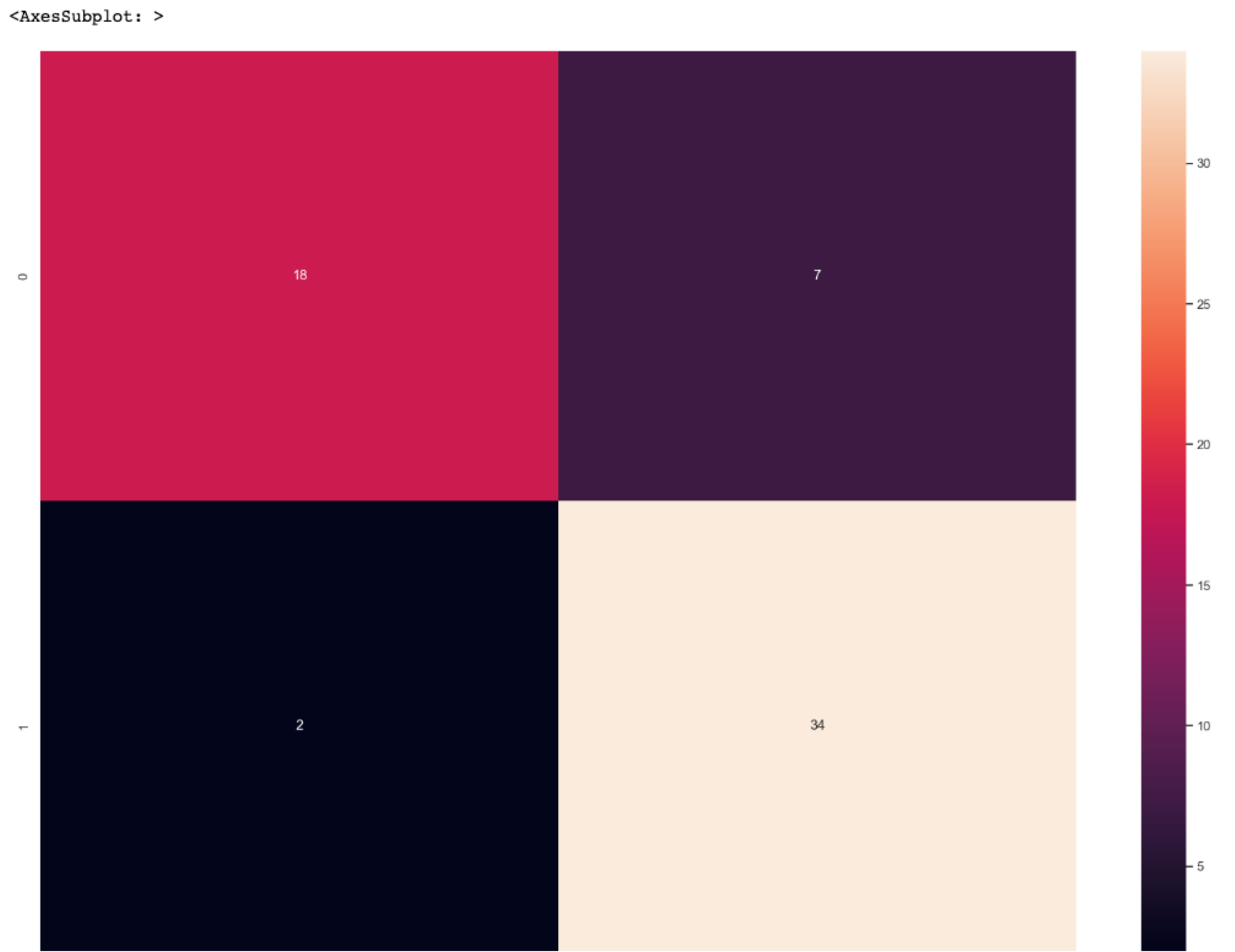
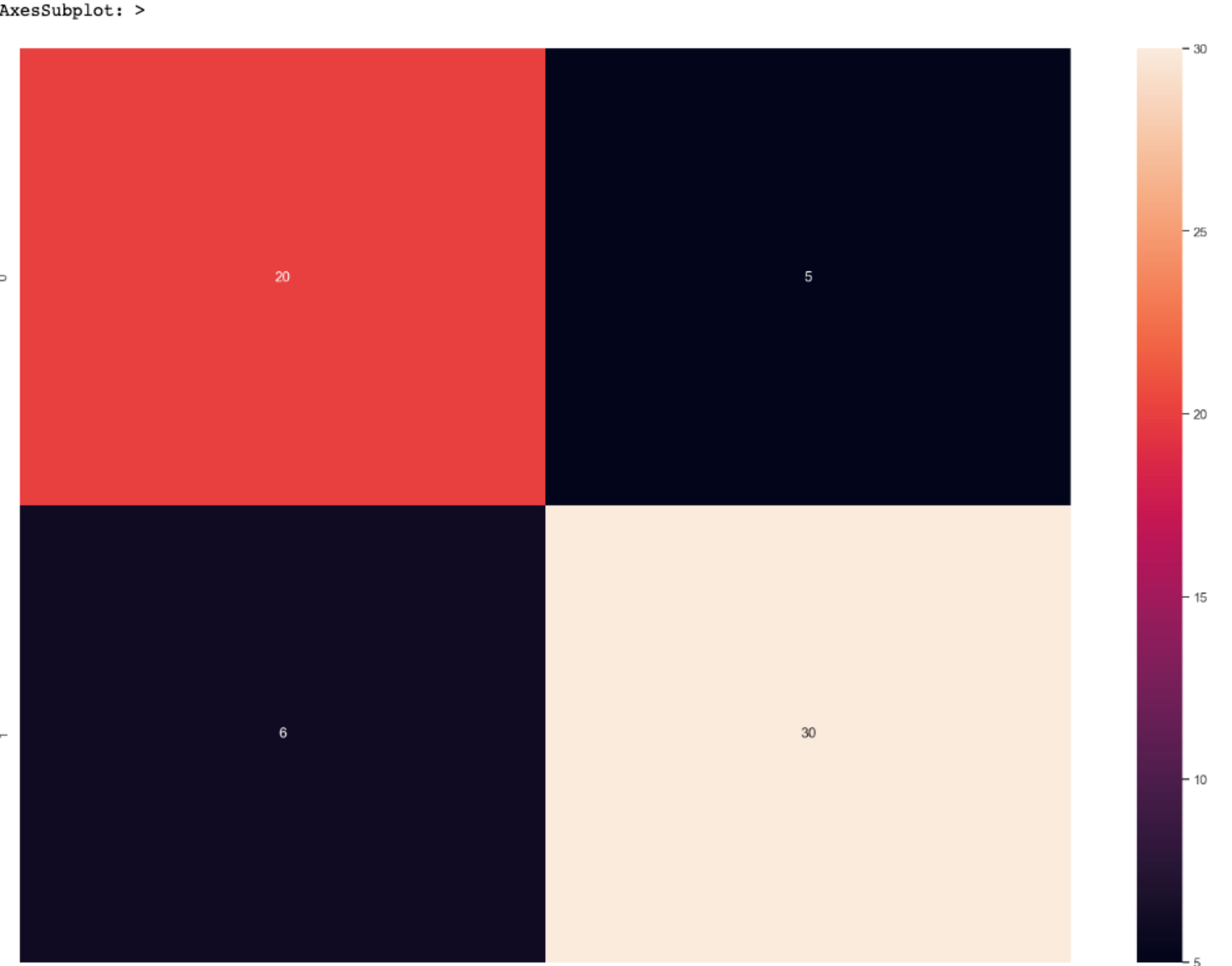
```
Test set: 12.877%
```

Classification report

	precision	recall	f1-score	support
0	0.90	0.72	0.80	25
1	0.83	0.94	0.88	36
accuracy			0.85	61
macro avg	0.86	0.83	0.84	61
weighted avg	0.86	0.85	0.85	61

```
1 y_pred_train = rf_best_model.predict(X_train)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```

```
1 y_pred_test = rf_best_model.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```



Train

Test

Confusion Matrices
Random Forest

```
1 metrics_score(X_train, X_test, y_train, y_test, best_model_KNNeighbors)
```

Roc-auc score

Train set: 87.669%
Test set: 89.444%

Accuracy score

Train set: 80.567%
Test set: 76.667%

Standart deviation

Train set: 7.313%
Test set: 10.869%

Classification report

	precision	recall	f1-score	support
0	0.81	0.68	0.74	25
1	0.80	0.89	0.84	36
accuracy			0.80	61
macro avg	0.80	0.78	0.79	61
weighted avg	0.80	0.80	0.80	61

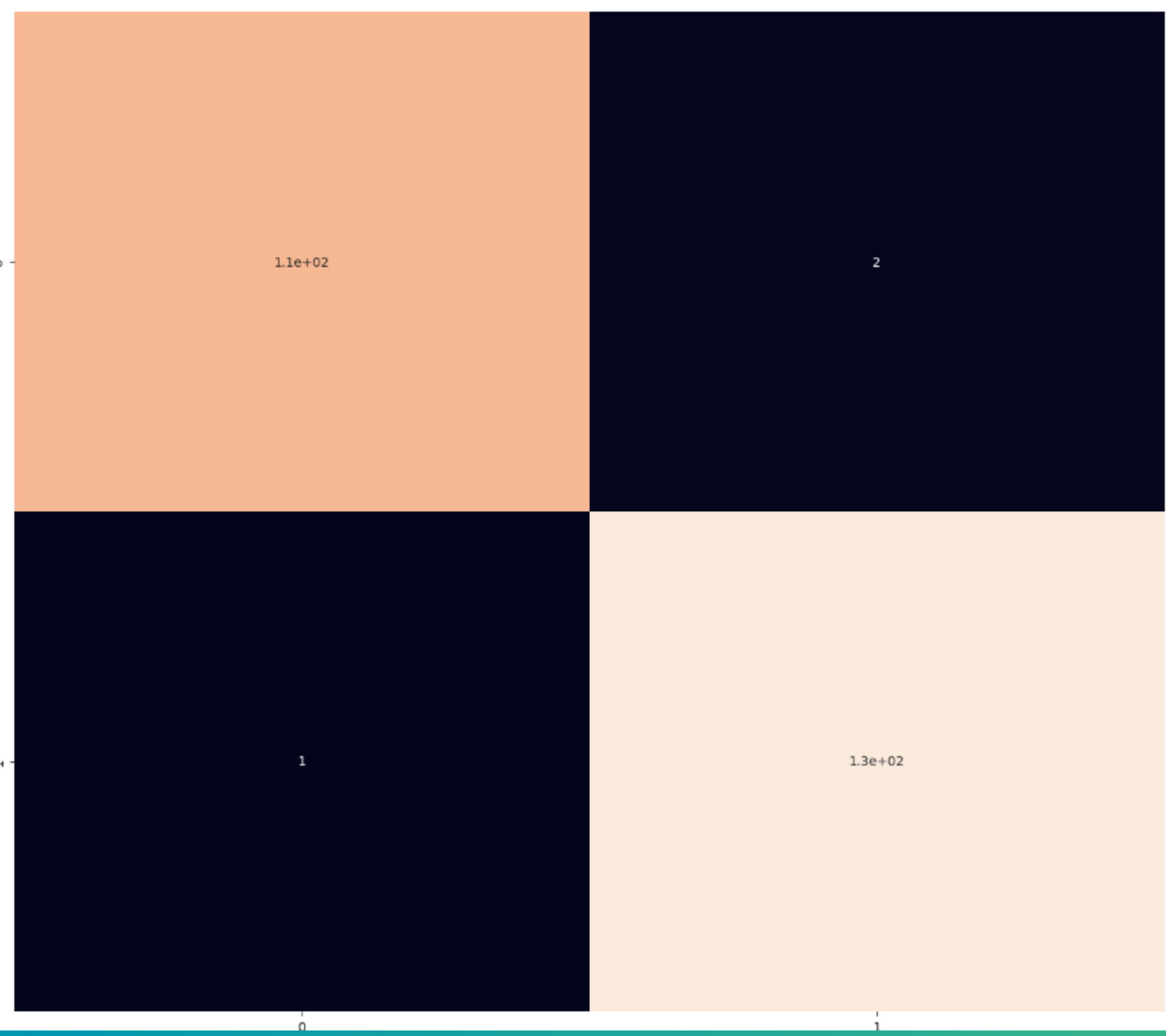
Best Results

KNNeighbors

Classifier

```
1 y_pred_train = best_model_KNNNeighbors.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

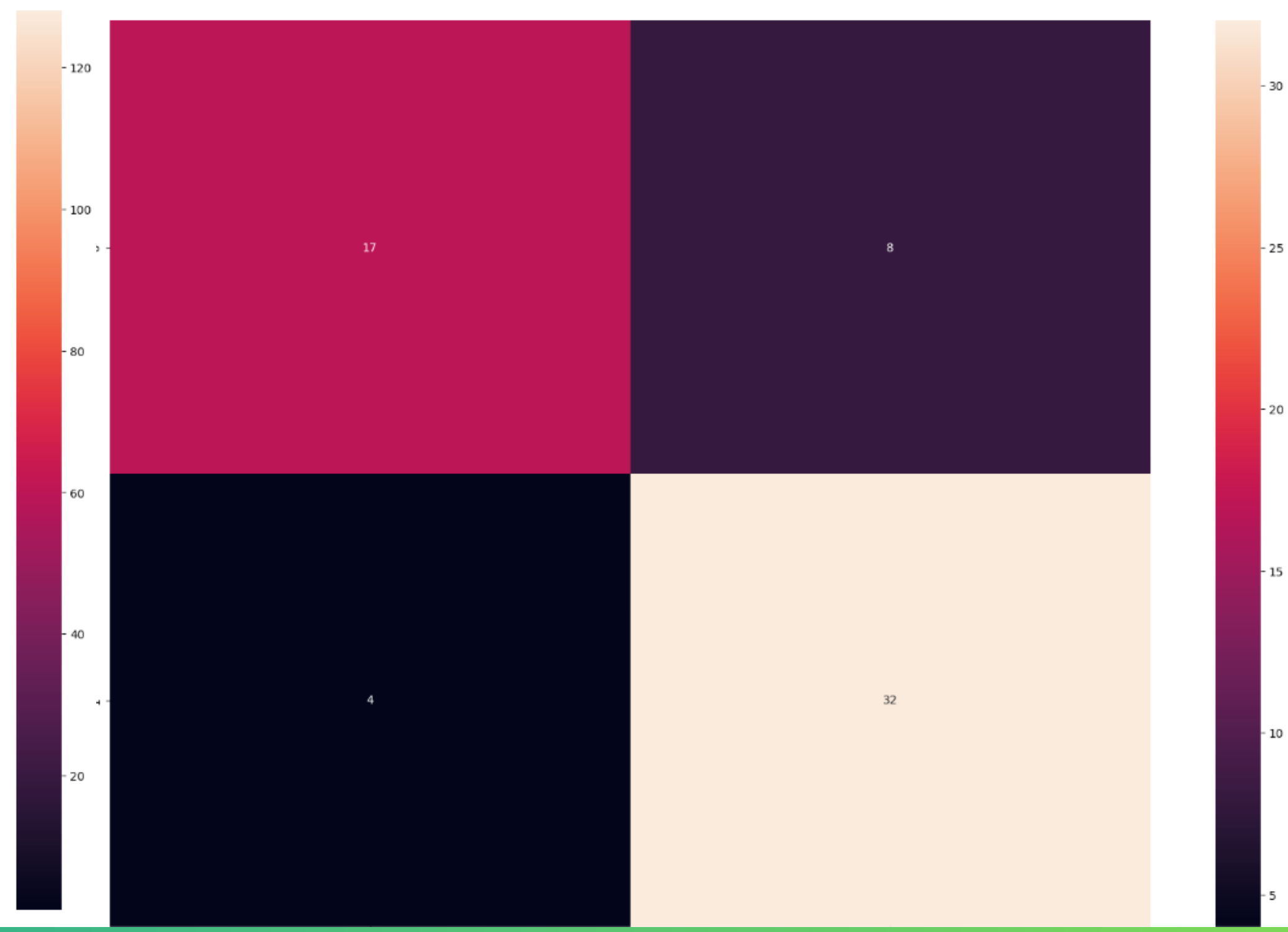
AxesSubplot: >



Train

```
1 y_pred_test = best_model_KNNNeighbors.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```

AxesSubplot: >



Test

Confusion Matrices

KNNNeighbors Classifier

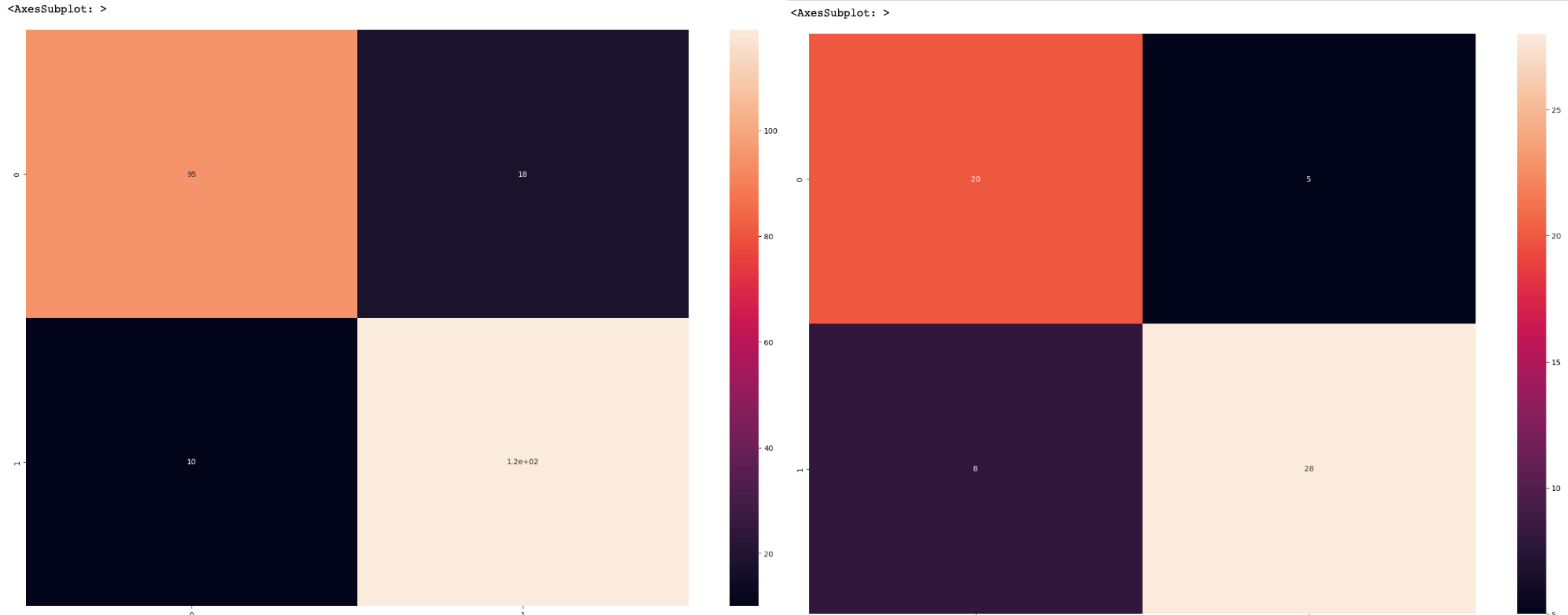
1	metrics_score(X_train, X_test, y_train, y_test, best_model_SVC)				
Roc-auc score					
Train set: 89.604%					
Test set: 89.722%					
Accuracy score					
Train set: 82.217%					
Test set: 86.667%					
Standart deviation					
Train set: 7.627%					
Test set: 10.833%					
Classification report					
	precision	recall	f1-score	support	
0	0.71	0.80	0.75	25	
1	0.85	0.78	0.81	36	
accuracy			0.79	61	
macro avg	0.78	0.79	0.78	61	
weighted avg	0.79	0.79	0.79	61	

Best Results

Support Vector Classifier

```
1 y_pred_train = best_model_SVC.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

```
1 y_pred_test = best_model_SVC.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```



Train

Test

Confusion Matrics

Support Vector Classifier

```
1 metrics_score(X_train, X_test, y_train, y_test, best_model_logit)
```

Roc-auc score

Train set: 90.058%

Test set: 89.444%

Accuracy score

Train set: 82.633%

Test set: 83.333%

Standart deviation

Train set: 6.117%

Test set: 11.167%

Classification report

	precision	recall	f1-score	support
0	0.77	0.80	0.78	25
1	0.86	0.83	0.85	36
accuracy			0.82	61
macro avg	0.81	0.82	0.81	61
weighted avg	0.82	0.82	0.82	61

Best Results

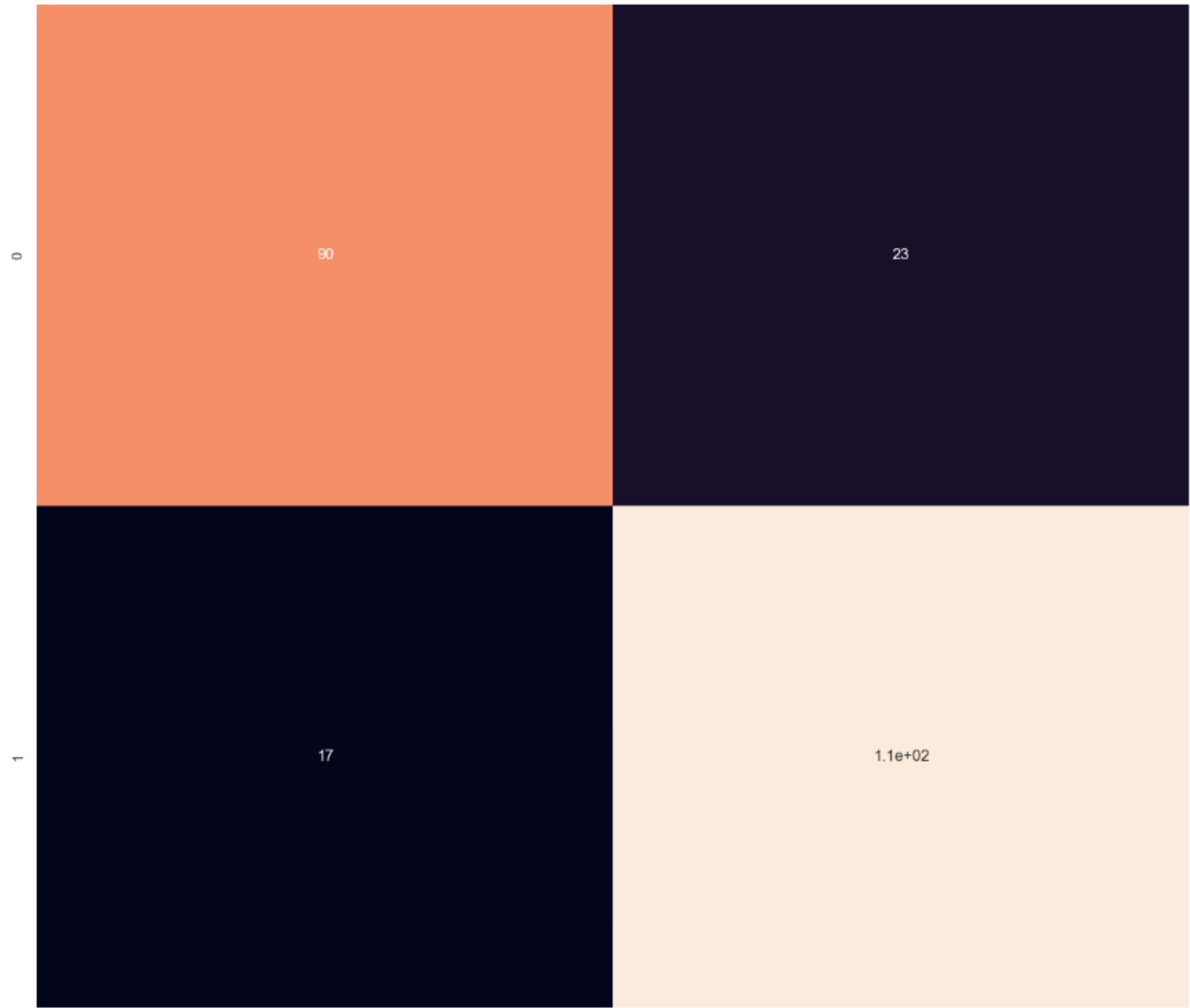
Logistic

Regression

```
1 y_pred_test = best_model_logit.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

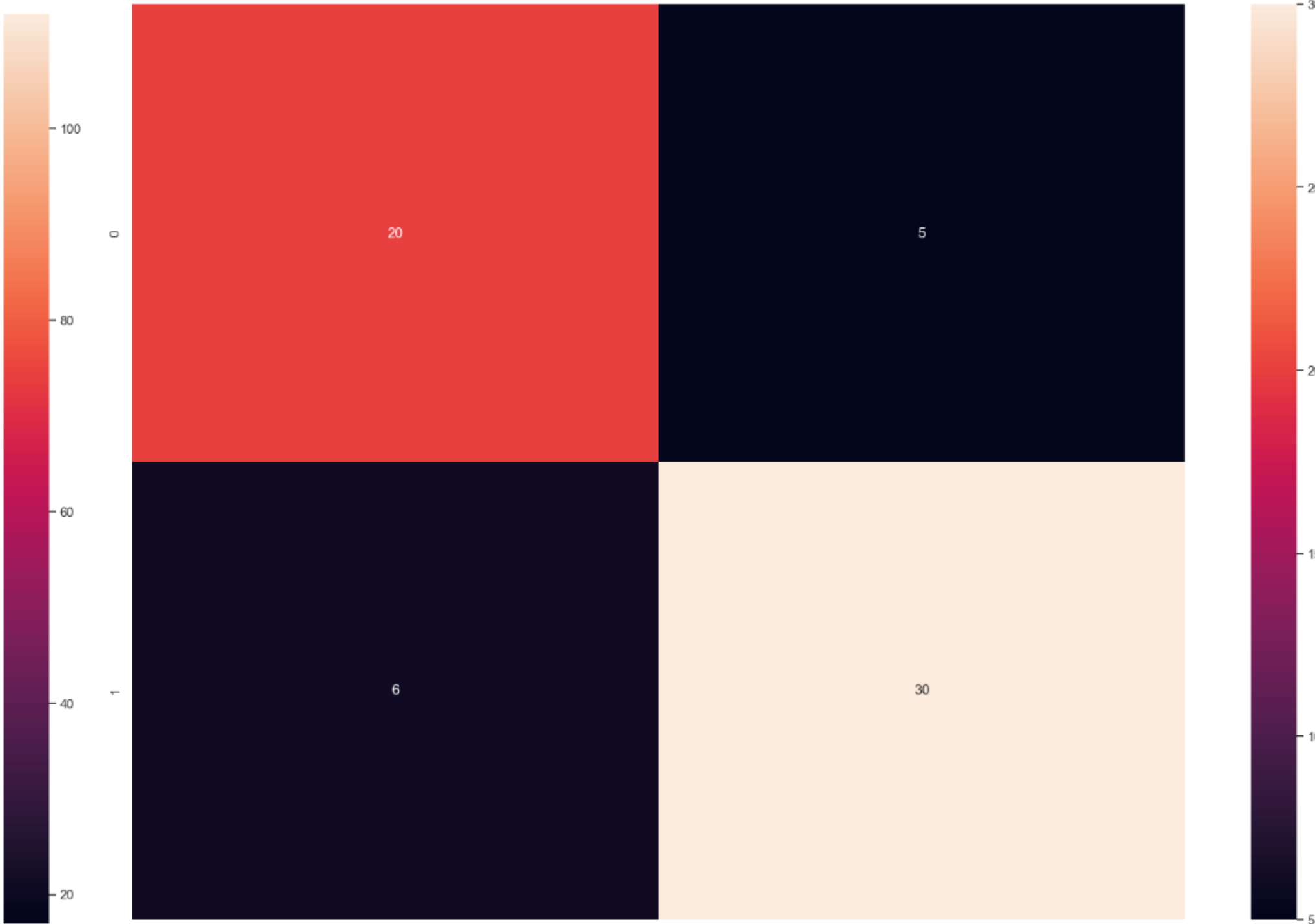
```
1 y_pred_test = best_model_logit.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```

<AxesSubplot: >



Train

<AxesSubplot: >



Test

Confusion Matrices

Logistic Regression

1

metrics_score(X_train, X_test, y_train, y_test, best_model_ada)

Roc-auc score

Train set: 89.704%

Test set: 90.278%

Accuracy score

Train set: 83.083%

Test set: 86.667%

Standart deviation

Train set: 5.739%

Test set: 14.540%

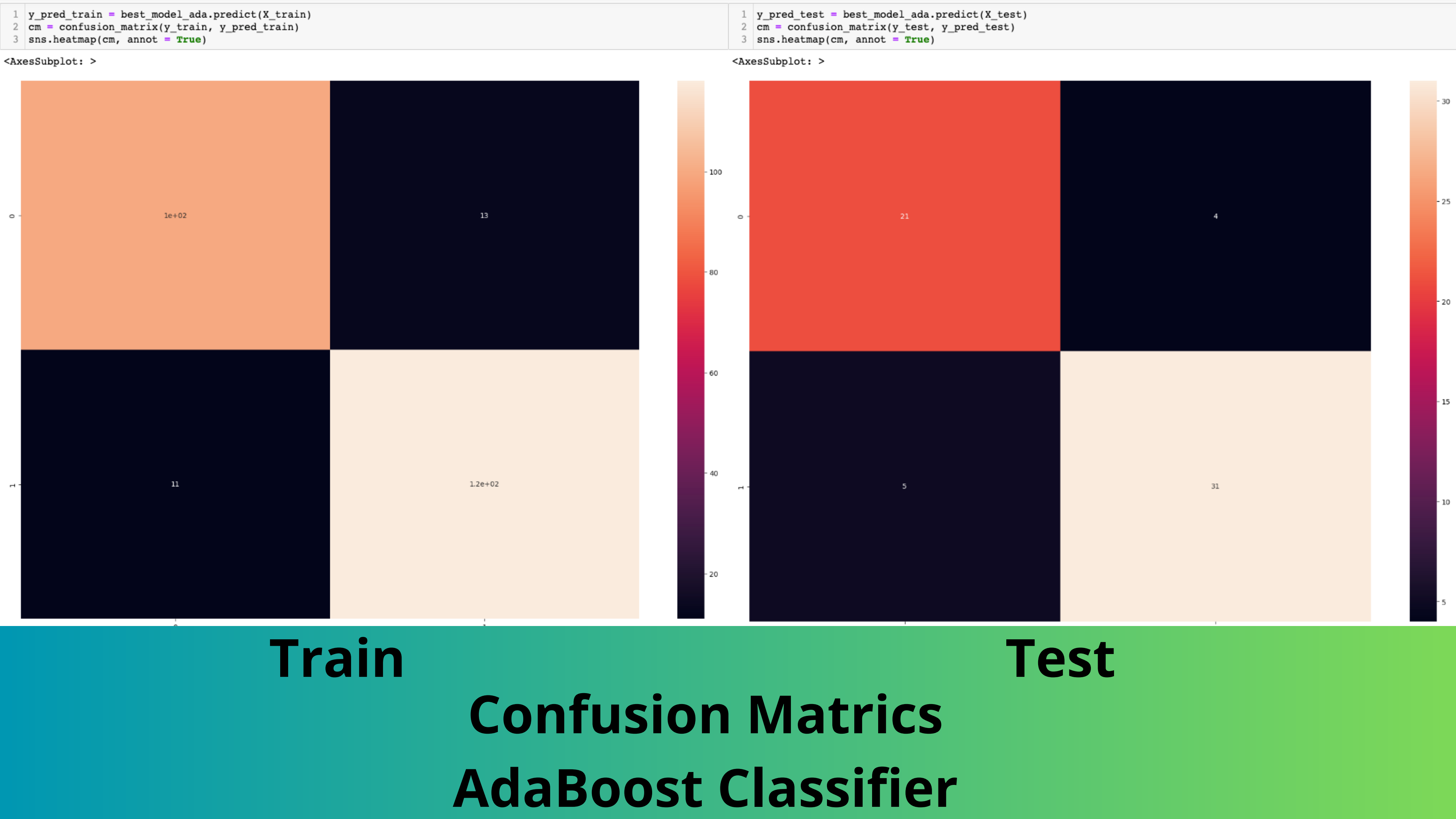
Classification report

	precision	recall	f1-score	support
0	0.81	0.84	0.82	25
1	0.89	0.86	0.87	36
accuracy			0.85	61
macro avg	0.85	0.85	0.85	61
weighted avg	0.85	0.85	0.85	61

Best Results

AdaBoost

Classifier



Roc-auc score

Train set: 87.068%

Test set: 90.694%

Accuracy score

Train set: 80.200%

Test set: 83.333%

Standart deviation

Train set: 5.586%

Test set: 12.851%

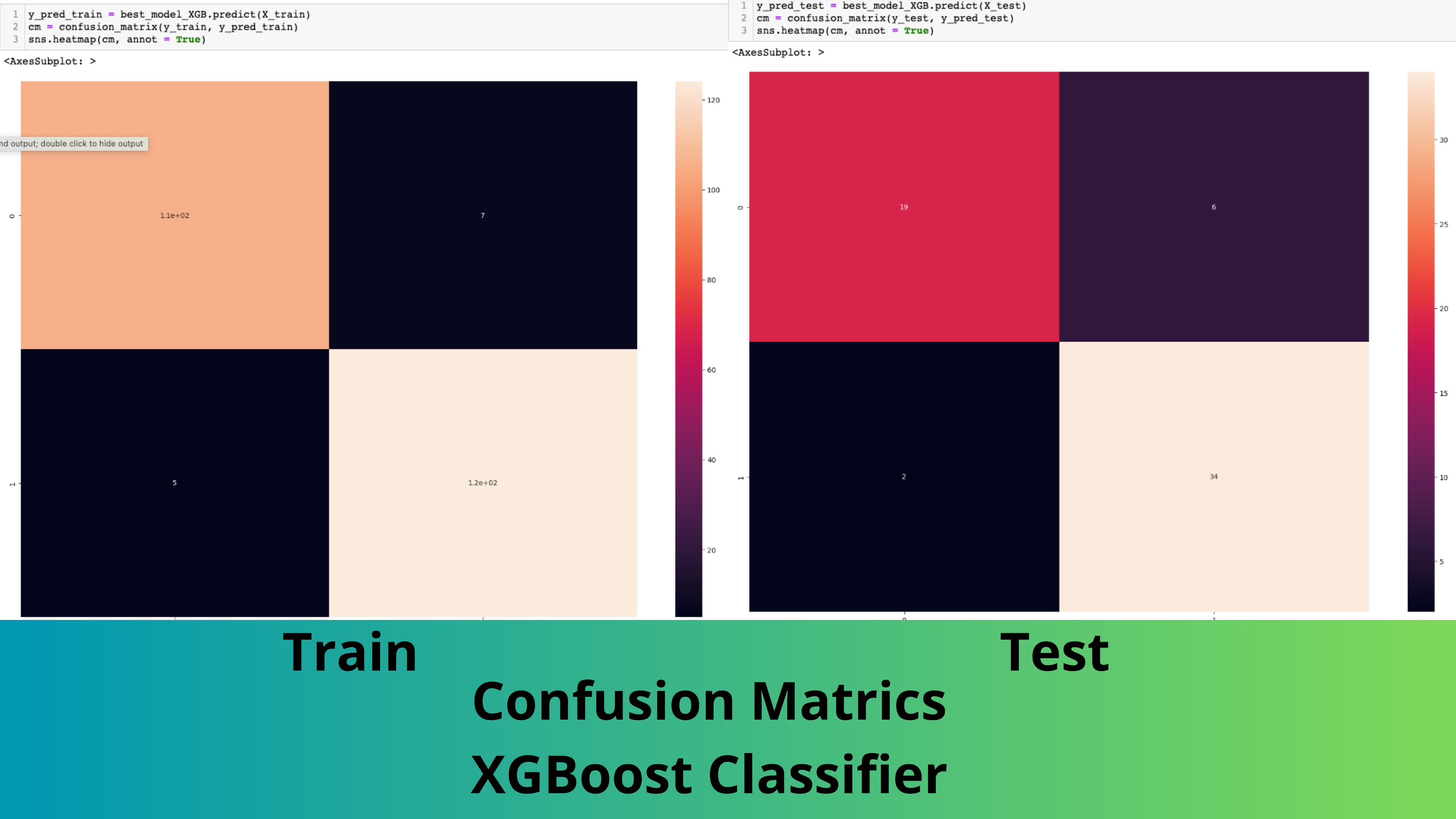
Classification report

	precision	recall	f1-score	support
0	0.90	0.76	0.83	25
1	0.85	0.94	0.89	36
accuracy			0.87	61
macro avg	0.88	0.85	0.86	61
weighted avg	0.87	0.87	0.87	61

Best Results

XGBoost

Classifier



Roc-auc score

Train set: 88.946%
Test set: 94.167%

Accuracy score

Train set: 81.333%
Test set: 85.000%

Standart deviation

Train set: 5.485%
Test set: 9.497%

Classification report

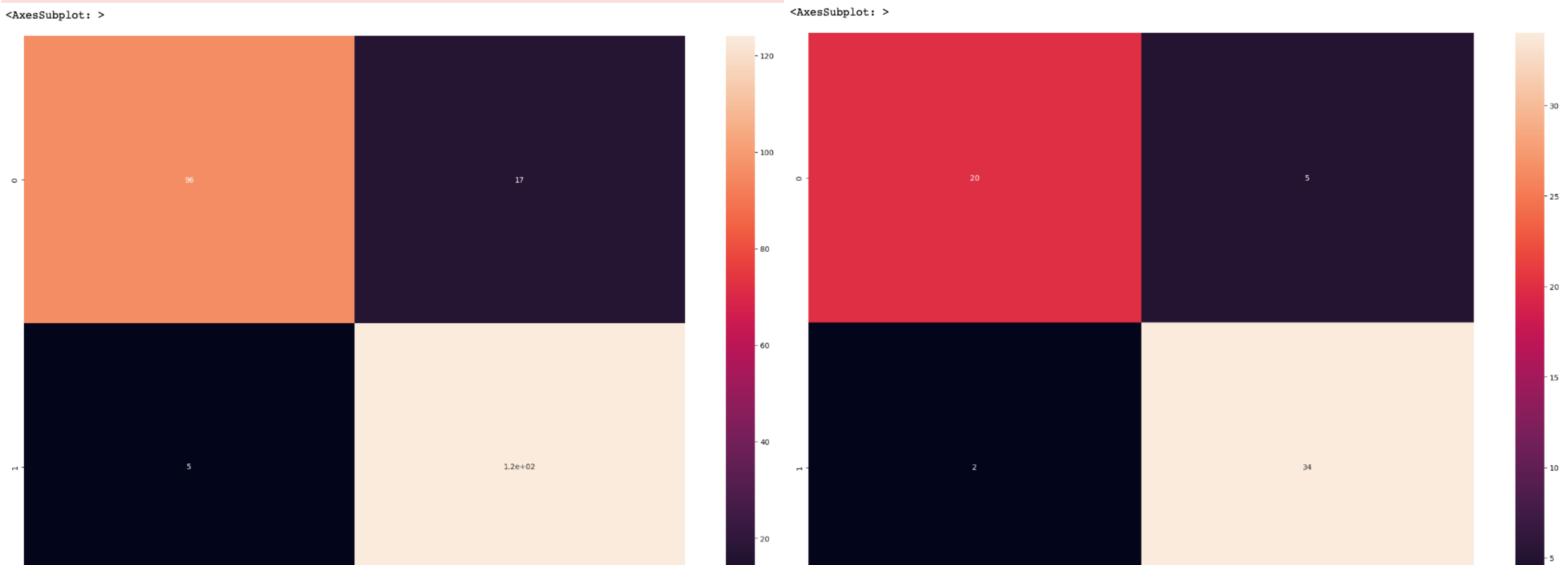
	precision	recall	f1-score	support
0	0.91	0.80	0.85	25
1	0.87	0.94	0.91	36
accuracy			0.89	61
macro avg	0.89	0.87	0.88	61
weighted avg	0.89	0.89	0.88	61

Best Results

Bagging
Classifier

```
1 y_pred_train = best_model_Bagg.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

```
[Parallel(n_jobs=8)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=8)]: Done 3 out of 8 | elapsed: 0.0s remaining: 0.0s
[Parallel(n_jobs=8)]: Done 8 out of 8 | elapsed: 0.0s finished
```



Train

Test

Confusion Matrices

Bagging Classifier

1	metrics_score(X_train, X_test, y_train, y_test, best_model_Extra)
---	---

Roc-auc score

Train set: 90.420%

Test set: 91.806%

Accuracy score

Train set: 82.183%

Test set: 85.000%

Standart deviation

Train set: 4.041%

Test set: 10.476%

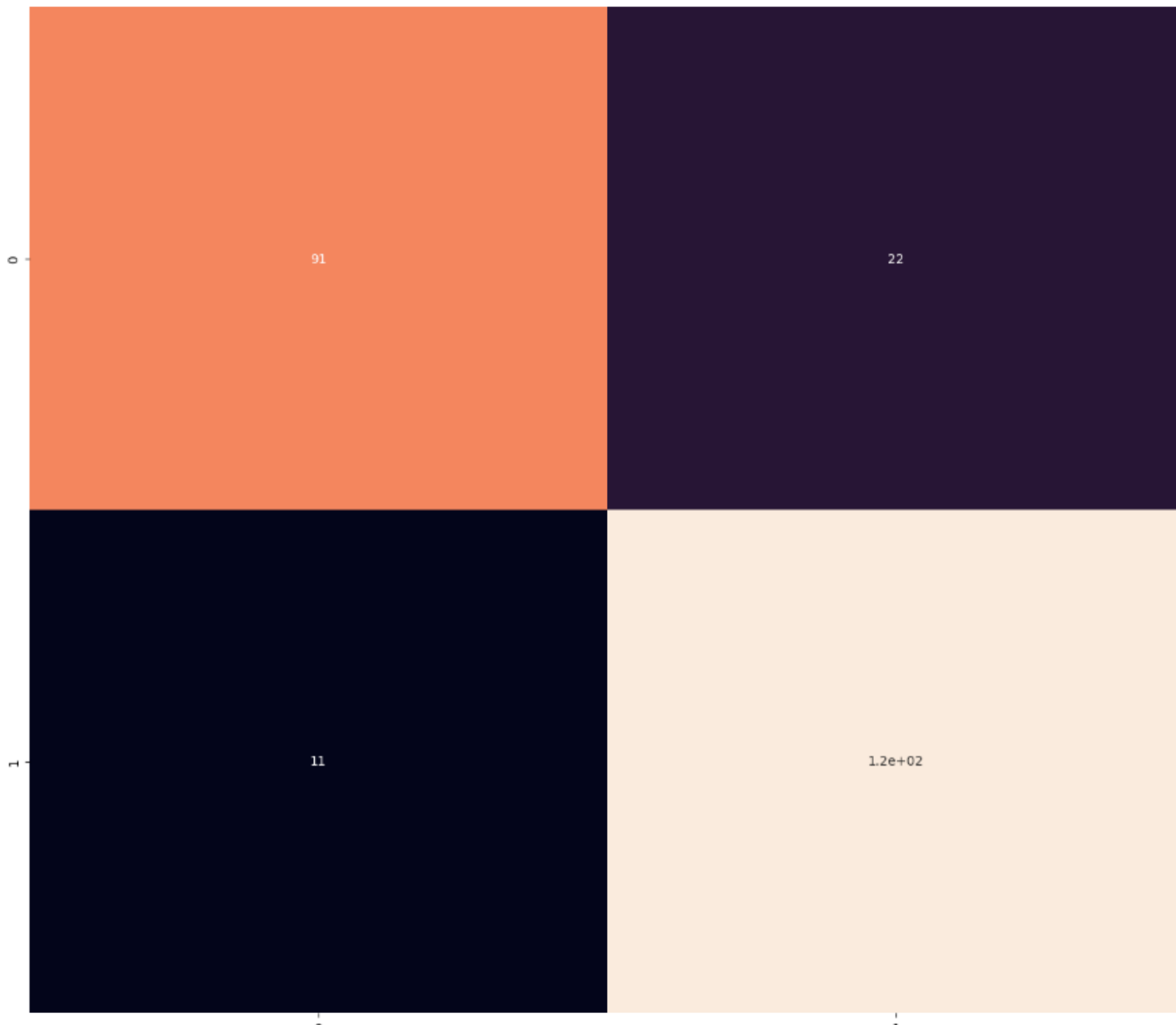
Classification report

	precision	recall	f1-score	support
0	0.86	0.72	0.78	25
1	0.82	0.92	0.87	36
accuracy			0.84	61
macro avg	0.84	0.82	0.83	61
weighted avg	0.84	0.84	0.83	61

Best Results ExtraTree Classifier

```
1 y_pred_train = best_model_Extra.predict(X_train)
2 cm = confusion_matrix(y_train, y_pred_train)
3 sns.heatmap(cm, annot = True)
```

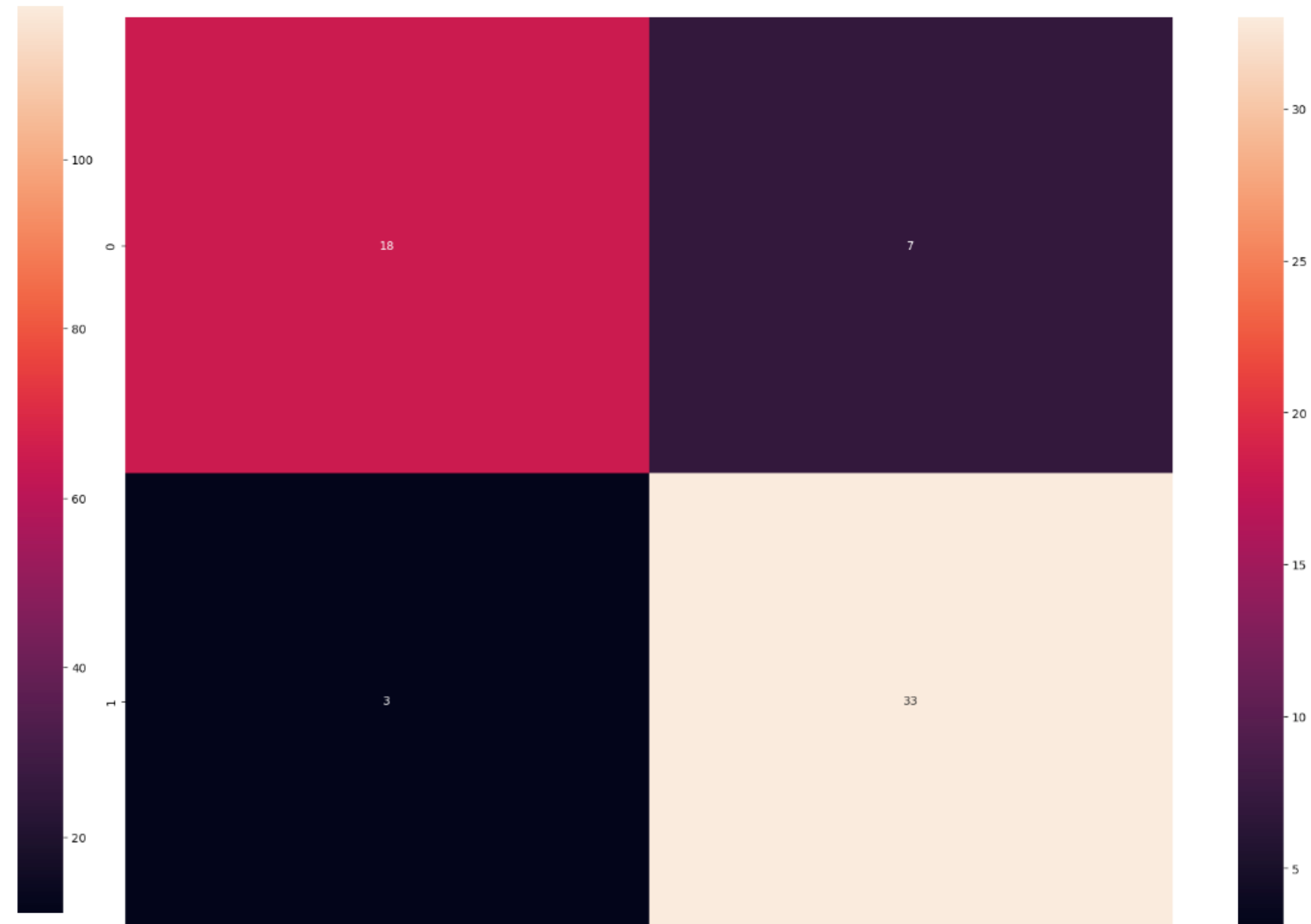
<AxesSubplot: >



Train

```
1 y_pred_test = best_model_Extra.predict(X_test)
2 cm = confusion_matrix(y_test, y_pred_test)
3 sns.heatmap(cm, annot = True)
```

<AxesSubplot: >



Test

Confusion Matrices
ExtraTree Classifier

1

metrics_score(X_train, X_test, y_train, y_test, best_model_NB)

Roc-auc score

Train set: 88.903%

Test set: 92.917%

Accuracy score

Train set: 80.967%

Test set: 80.000%

Standart deviation

Train set: 5.608%

Test set: 9.470%

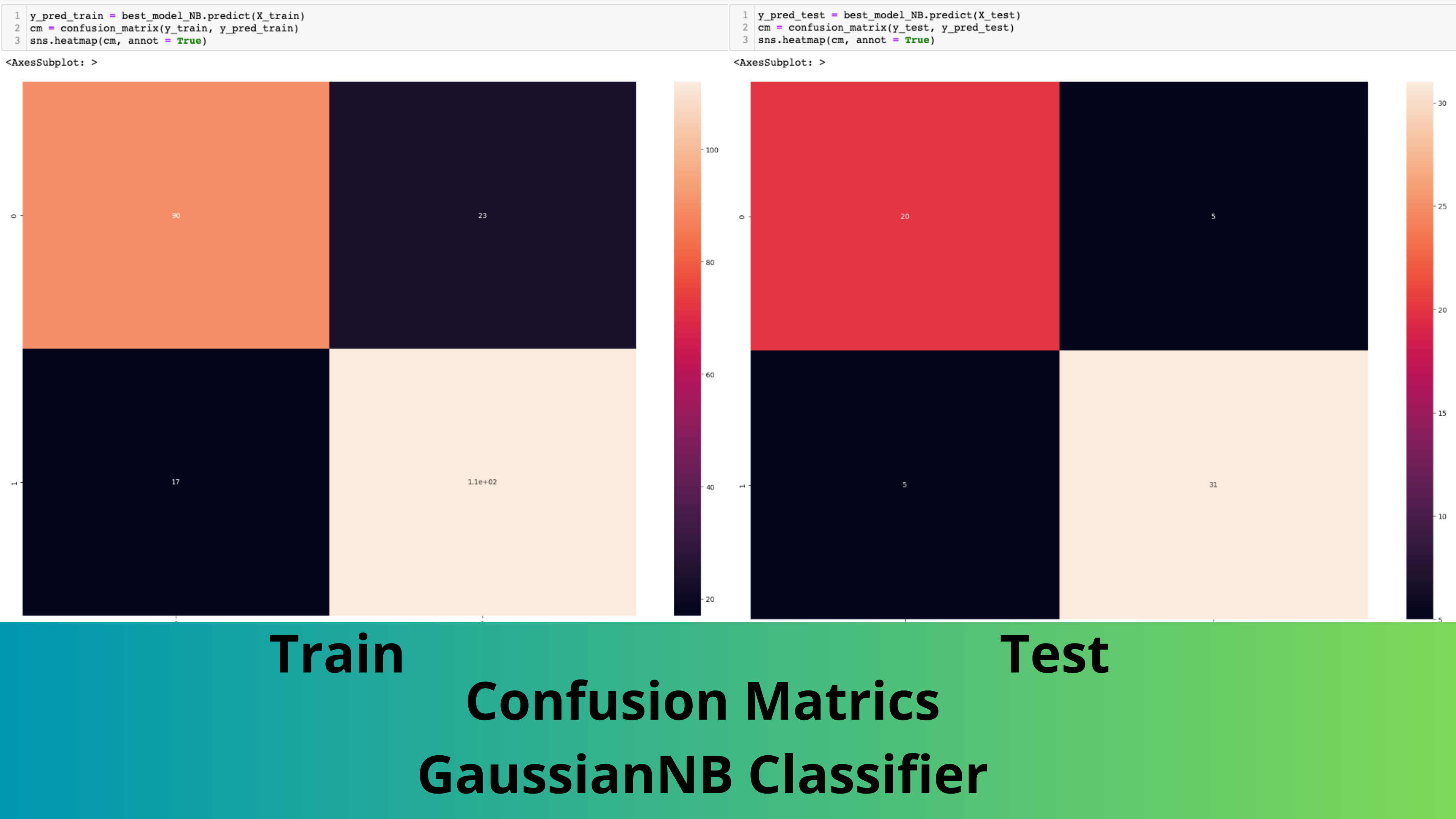
Classification report

	precision	recall	f1-score	support
0	0.80	0.80	0.80	25
1	0.86	0.86	0.86	36
accuracy			0.84	61
macro avg	0.83	0.83	0.83	61
weighted avg	0.84	0.84	0.84	61

Best Results

GaussianNB

Classifier



All Best Results

```
Decision Tree by accuracy: 77%  
Random Forest by accuracy: 85%  
KNeighbor Classifier by accuracy: 80%  
SVC by accuracy: 79%  
Ada Boost by accuracy: 85%  
Bagging Classifier by accuracy: 89%  
Extra Tree Classifier by accuracy: 84%  
XGBoost Classifier by accuracy: 83%  
GaussianNB Classifier by accuracy: 84%  
Logistic Regression by accuracy: 82%
```

Conclusion

1 Goal summary:

find the best machine learning model, in this project we solved the problem of binary classification

Where there is heart disease or not

2 Main Results

as a results of our project, not bad scores were achived for the following metrics such as Accuracy,

roc_auc score, precision, recall, f1-score and Standart Deviation

3 Models

Of all the algorithms, Bagging Classifier give us the highest accuracy score