

# AMEO 2015 analysis

Our job here is to find out which factors affect the future salary of an individual. We will find out the insights after processing the data. The dataset has been obtained from Aspiring Minds from the Aspiring Mind Employment Outcome 2015 (AMEO). The study is primarily limited only to students with engineering disciplines.

```
In [183]: #installing packages
import numpy as np
import pandas as pd
import matplotlib as plt
import seaborn as sb
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import r2_score
data_train, data_test = train_test_split(
    from sklearn.model_selection import cross_val_score
from sklearn.model_selection import cross_val_score
from sklearn.metrics import mean_squared_error
from sklearn.metrics import classification_report
from sklearn.pipeline import Pipeline

from sklearn.preprocessing import PolynomialFeatures
from sklearn.preprocessing import LabelEncoder

from scipy.stats import boxcox
from datetime import datetime
from datetime import date
import statsmodels.api as sm
from statsmodels.formula.api import ols
from sklearn.linear_model import LinearRegression
from sklearn.decomposition import PCA

from scipy.stats import chi2_contingency
from sklearn import metrics
from sklearn.preprocessing import OneHotEncoder
from sklearn.compose import ColumnTransformer

import larrypredict
from larrypredict.supervised import LazyRegressor
```

Out[2]:

ID	Salary	DOJ	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	...	English	Logical	Quant	Domain	Marks
0	203097	420000	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	...	515	585	525	0.635979
1	1027655	300000	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	...	560	555	620	-1.000000
2	947847	300000	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	...	590	435	380	0.356536
3	87291	600000	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	...	565	455	565	0.765674
4	606655	270000	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	...	380	580	470	0.229482

5 rows x 30 columns

```
In [3]: #examining the dataset
ameo.info()
```

Out[3]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1961 entries, 0 to 1960
Data columns (total 30 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   ID                    1961 non-null  int64
1   Salary               1961 non-null  float64
2   DOJ                  1961 non-null  object
3   DOL                  1961 non-null  object
4   Designation          1961 non-null  object
5   JobCity              1961 non-null  object
6   Gender               1961 non-null  object
7   DOB                  1961 non-null  object
8   Sgradyr              1961 non-null  object
9   percentage           1961 non-null  float64
10  board                1961 non-null  object
11  CollegeID            1961 non-null  int64
12  Collegetier           1961 non-null  int64
13  Degree               1961 non-null  object
14  Specialization        1961 non-null  object
15  collegeGPA            1961 non-null  float64
16  collegeCityID         1961 non-null  int64
17  collegeCityTier       1961 non-null  int64
18  CollegeState          1961 non-null  object
19  CGradYear             1961 non-null  object
20  English               1961 non-null  object
21  Logical               1961 non-null  int64
22  Quant                1961 non-null  int64
23  Domain               1961 non-null  float64
24  Marks_obtained        1935 non-null  float64
25  conscientiousness     1961 non-null  float64
26  agreeableness         1961 non-null  float64
27  extraversion          1961 non-null  float64
28  neuroticism           1961 non-null  float64
29  openness_to_experience 1961 non-null  float64
dtypes: float64(10), int64(12), object(8)
memory usage: 459.7+ KB
```

```
In [4]: #examining the shape
ameo.shape
```

Out[4]:

```
(1961, 30)
```

## Step 2: Data Pre processing

In this step, we are going to clean our dataset. We are going to look for null values and replace them with mean and mode. We are going to modify some variables if it is necessary and change datatypes for better analysis. We will also remove outliers from the dataset. Outliers hamper the machine learning algorithms and hence they have to be removed.

```
In [5]: # column name cleaning
ameo.columns = ameo.columns.str.replace(' ', '')
ameo.columns
```

Out[5]:

```
Index(['ID', 'Salary', 'DOJ', 'DOL', 'Designation', 'JobCity', 'Gender', 'DOB', 'Sgradyr', 'percentage', 'board', 'CollegeID', 'Collegetier', 'Degree', 'Specialization', 'collegeGPA', 'collegeCityID', 'CollegeCityTier', 'CollegeState', 'CGradYear', 'English', 'Logical', 'Quant', 'Domain', 'Marks_obtained', 'conscientiousness', 'agreeableness', 'extraversion', 'neuroticism', 'openness_to_experience'],
      dtype='object')
```

```
In [6]: #checking null values
ameo.isna().sum(axis=0)
```

Out[6]:

```
ID                0
Salary            0
DOJ               0
DOL               0
Designation       0
JobCity           0
Gender            0
DOB              0
Sgradyr           0
percentage        0
board             0
CollegeID         0
Collegetier       0
Degree            0
Specialization    0
collegeGPA        0
collegeCityID     0
CollegeCityTier   0
CollegeState      0
CGradYear         0
English           0
Logical           0
Quant             0
Domain            0
Marks_obtained    26
conscientiousness 0
agreeableness     0
extraversion      0
neuroticism       0
openness_to_experience 0
dtype: int64
```

```
In [149]: #changing data type to right format
ameo.DOB=ameo.DOB.astype(object)
ameo.DOB=ameo.DOB.astype(object)
ameo.Sgradyr=ameo.Sgradyr.astype(object)
ameo.CGradYear=ameo.CGradYear.astype(object)
ameo.Salary=ameo.Salary.astype(float)
```

ameo.dtypes

```
AttributeError                                Traceback (most recent call last)
Cell In [149], line 6
      5 ameo.Sgradyr=ameo.Sgradyr.astype(object)
      6 ameo.CGradYear=ameo.CGradYear.astype(object)
----> 7 ameo.Salary=ameo.Salary.astype(float)
      8 ameo.dtypes

File ~\anaconda3\lib\site-packages\pandas\core\generic.py:5583, in NDFrame._get_attr_(self, name)
    5578 if (
    5579     name not in self._internal_names_set
    5580     and name not in self._metadata
    5581     and name not in self._accessors
    5582 ):
    5583     return self[name]
--> 5583 return Object._getattribute_(self, name)

AttributeError: 'DataFrame' object has no attribute 'Salary'
```

```
In [8]: #replacing missing values with numpy nan
ameo.replace("",np.nan,inplace = True)
ameo.head()
```

Out[8]:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1961 entries, 0 to 1960
Data columns (total 30 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   ID                    1961 non-null  int64
1   Salary               1961 non-null  float64
2   DOJ                  1961 non-null  object
3   DOL                  1961 non-null  object
4   Designation          1961 non-null  object
5   JobCity              1961 non-null  object
6   Gender               1961 non-null  object
7   DOB                  1961 non-null  object
8   Sgradyr              1961 non-null  object
9   percentage           1961 non-null  float64
10  board                1961 non-null  object
11  CollegeID            1961 non-null  int64
12  Collegetier           1961 non-null  int64
13  Degree               1961 non-null  object
14  Specialization        1961 non-null  object
15  collegeGPA            1961 non-null  float64
16  collegeCityID         1961 non-null  int64
17  collegeCityTier       1961 non-null  int64
18  CollegeState          1961 non-null  object
19  CGradYear             1961 non-null  object
20  English               1961 non-null  object
21  Logical               1961 non-null  int64
22  Quant                1961 non-null  int64
23  Domain               1961 non-null  float64
24  Marks_obtained        1935 non-null  float64
25  conscientiousness     1961 non-null  float64
26  agreeableness         1961 non-null  float64
27  extraversion          1961 non-null  float64
28  neuroticism           1961 non-null  float64
29  openness_to_experience 1961 non-null  float64
dtypes: float64(10), int64(8), object(12)
memory usage: 459.7+ KB
None
```

```
In [9]: #gathering mean values of numerical variables
ameo.ameo[ameo['Marks_obtained'].astype("float").mean(axis=0)
print("The average marks of a student in our Test is",avg_Marks_obtained,"points.")
```

```
#replacing numerical missing values with mean
ameo['Marks_obtained'].replace(np.nan, avg_Marks_obtained, inplace=True)
```

The average marks of a student in our Test is 678.6428364047317 points.

```
In [10]: #creating a compact variable-apptitude score
ameo['apptitude_score']=ameo['English']+ameo['Logical']+ameo['Quant']/3
ameo['apptitude_score']
ameo.head()
```

Out[10]:

ID	Salary	DOJ	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	...	Logical	Quant	Domain	Marks obtained	
0	203097	420000.0	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	...	585	525	0.635979	4.0
1	1027655	300000.0	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	...	555	620	-1.000000	8.0
2	947847	300000.0	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	...	435	380	0.356536	4.0
3	87291	600000.0	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	...	455	565	0.765674	8.0
4	606655	270000.0	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	...	580	470	0.229482	4.0

5 rows x 31 columns

```
In [11]: #removing unnecessary variables
ameo=ameo.drop(['English','Logical','Quant','CollegeCityID','CollegeState','ID'],axis=1)
ameo.head()
```

Out[11]:

Salary	DOJ	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	board	...	CollegeCityTier	CGradYear	Domain	M
0	420000.0	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	cbse	...	0	2011	0.635979
1	300000.0	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	cbse	...	0	2014	-1.000000
2	300000.0	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	cbse	...	1	2014	0.356536
3	600000.0	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	cbse	...	0	2010	0.765674
4	270000.0	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	cbse	...	0	2013	0.229482

5 rows x 25 columns

```
In [12]: #having a look at the new variables
ameo.dtypes
```

Out[12]:

```
Salary                float64
DOJ                   object
DOL                   object
Designation           object
JobCity               object
Gender                object
DOB                  object
Sgradyr               object
percentage             float64
board                 object
CollegeID             int64
Collegetier           int64
Degree                object
Specialization         object
collegeGPA            float64
CollegeCityTier       int64
CGradYear             object
Domain                float64
Marks_obtained        float64
conscientiousness     float64
agreeableness         float64
extraversion          float64
neuroticism           float64
apptitude_score       float64
dtype: object
```

```
In [13]: #Separate categorical and numerical columns
cat_cols = list(ameo.select_dtypes(include=['int64','object']).columns)
num_cols = list(ameo.select_dtypes(include=['float64']).columns)
```

Categorical columns: ['DOJ', 'DOL', 'Designation', 'JobCity', 'Gender', 'DOB', 'Sgradyr', 'board', 'CollegeID', 'Collegetier', 'Degree', 'Specialization', 'CollegeCityTier', 'CGradYear']  
Numerical columns: ['Salary', 'percentage', 'collegeGPA', 'Domain', 'Marks\_obtained', 'conscientiousness', 'agreeableness', 'extraversion', 'neuroticism', 'openness\_to\_experience', 'apptitude\_score']

## Univariate analysis

We are using boxplot,histogram and plot to find the outliers,extreme values and distribution of the numerical variables.

```
In [14]: #checking distribution for Salary
ameo['Salary'].plot.density()
```

Out[14]:

```
In [15]: #checking outliers using boxplot Salary
sb.boxplot(ameo['Salary'])
```

Out[15]:

```
In [16]: #outlier treatment for variable Salary
ameo['Salary'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

Out[16]:

```
0.10    110000.0
0.25    170000.0
0.50    250000.0
0.70    325000.0
0.90    450000.0
0.95    500000.0
0.99    847000.0
Name: Salary, dtype: float64
```

```
In [17]: #seeing limits
sal_HB=ameo['Salary'] > 500000.0.copy()
```

```
In [18]: #putting them at one place
print(len(sal_HB))
```

93

```
In [19]: #removing outliers
for x in ameo['Salary']:
    if x > 500000.0:
        ameo['Salary'].replace(x,np.nan,inplace=True)
```

```
In [28]: #checking if outliers gone
sb.boxplot(ameo['Salary'])
```

Out[28]:

```
In [21]: #replacing with median values
medsal=ameo['Salary'].median()
ameo['Salary'].replace(np.nan, medsal, inplace=True)
```

```
In [22]: #perform Shapiro-Wilk test
stat, p = stats.shapiro(ameo['Salary'])
```

Print output results:  
alpha = 0.05  
if p > alpha: Data looks Gaussian (fail to reject H0)  
else: print("Data does not look Gaussian (reject H0)")  
Data does not look Gaussian (reject H0)

```
In [23]: #Doing box-cox test to determine transform method
data = ameo['Salary']
```

```
# plot histogram of data
plt.hist(data, bins=50)
plt.show()
```

```
# apply Box-Cox transformation to data
data_transformed, lambda_ = boxcox(data)
```

```
# plot histogram of transformed data
plt.hist(data_transformed, bins=50)
plt.show()
```

```
# print optimal lambda value
print(lambda_)
```

0.49932772633940914

```
In [24]: #log transforming the variable
ameo['Salary']= np.log(ameo['Salary'])
```

```
# plot histogram of data
plt.hist(data, bins=50)
plt.show()
```

```
# apply log transformation
ameo['Salary']= np.log(ameo['Salary'])
```

```
fig, axs = plt.subplots(1, 2, figsize=(10, 4))
axs[0].hist(data, bins=50)
axs[0].set_title('Original data')
```

```
axs[1].hist(ameo['Salary'], bins=50)
axs[1].set_title('log-transformed data')
```

```
plt.show()
```

Original data

log-transformed data

0.49932772633940914

```
In [25]: #removing the original
ameo=ameo.drop(['Salary'],axis=1)
```

```
Out[25]:
```

DOJ	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	board	CollegeID	...	CGradYear	Domain	Marks obtained	
0	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	cbse	1141	...	2011	0.635979	445.0
1	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	cbse	5086	...	2014	-1.000000	872.0
2	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	cbse	314	...	2014	0.356536	445.0
3	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	cbse	403	...	2010	0.765674	822.0
4	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	cbse	2665	...	2013	0.229482	420.0

5 rows x 25 columns

Remarks

The variable salary was right skewed and had about 93 outliers which was removed and have been replaced with median. The boxcox value was 0.49 which required a log transformation to transform the right skewness of the variable.

```
In [26]: #checking distribution for percentage
ameo['percentage'].plot.density()
```

Out[26]:

```
In [27]: #checking outliers using boxplot percentage
sb.boxplot(ameo['percentage'])
```

Out[27]:

Remarks

The variable percentage was normally distributed and had about no outliers hence no transformation was needed.

```
In [28]: #checking distribution for collegeGPA
ameo['collegeGPA'].plot.density()
```

Out[28]:

```
In [29]: #checking outliers using boxplot collegeGPA
sb.boxplot(ameo['collegeGPA'])
```

Out[29]:

```
In [30]: #outlier treatment for variable collegeGPA
ameo['collegeGPA'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

Out[30]:

```
0.10    61.8
0.25    66.2
0.50    71.7
0.70    75.5
0.90    81.0
0.95    84.0
0.99    89.0
Name: collegeGPA, dtype: float64
```

```
In [31]: #seeing limits
gpa_HB=ameo['collegeGPA'] > 80.0.copy()
gpa_HB=ameo['collegeGPA'] > 80.0.copy()
```

```
In [32]: #putting them at one place
print(len(gpa_HB))
```

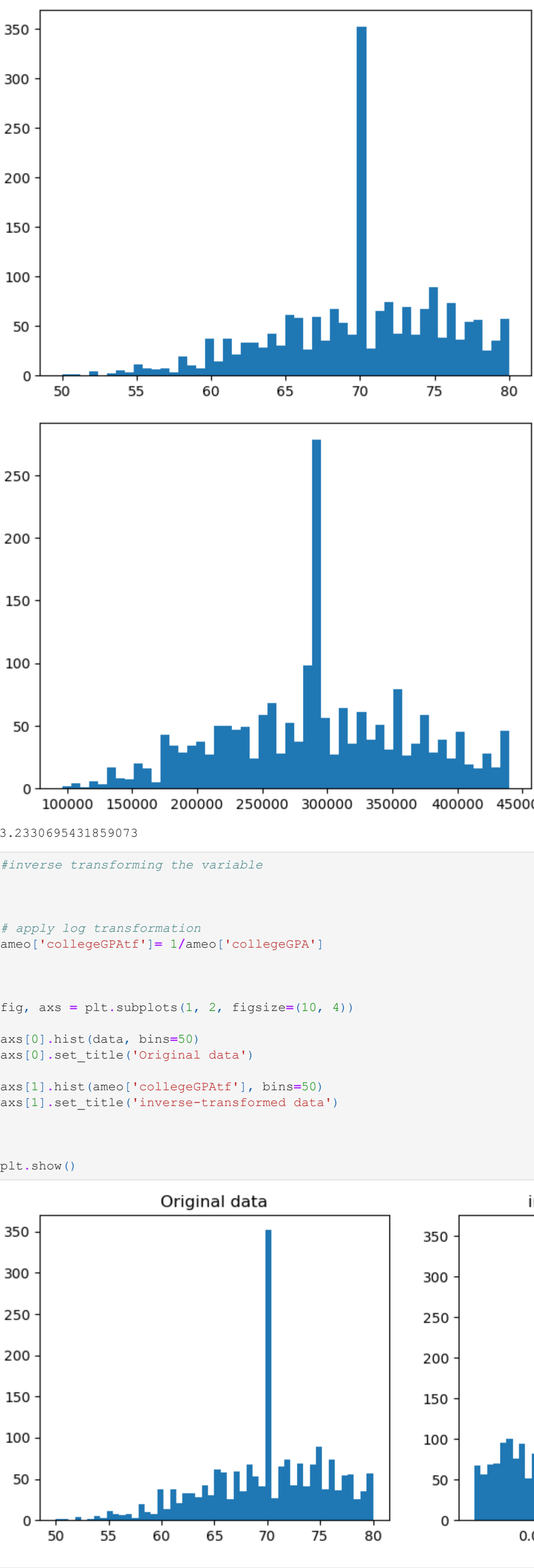
234

```
In [32]: #removing outliers
for x in ameo['collegeGPA']:
    if x > 80.0 or x < 50.0:
        ameo['collegeGPA'].replace(x,np.nan,inplace=True)
```

```
In [33]: #checking if outliers gone
sb.boxplot(ameo['collegeGPA'])
```

Out[33]:





```
3.2330695431859073
```

```
In [37]: #Inverse transforming the variable
```

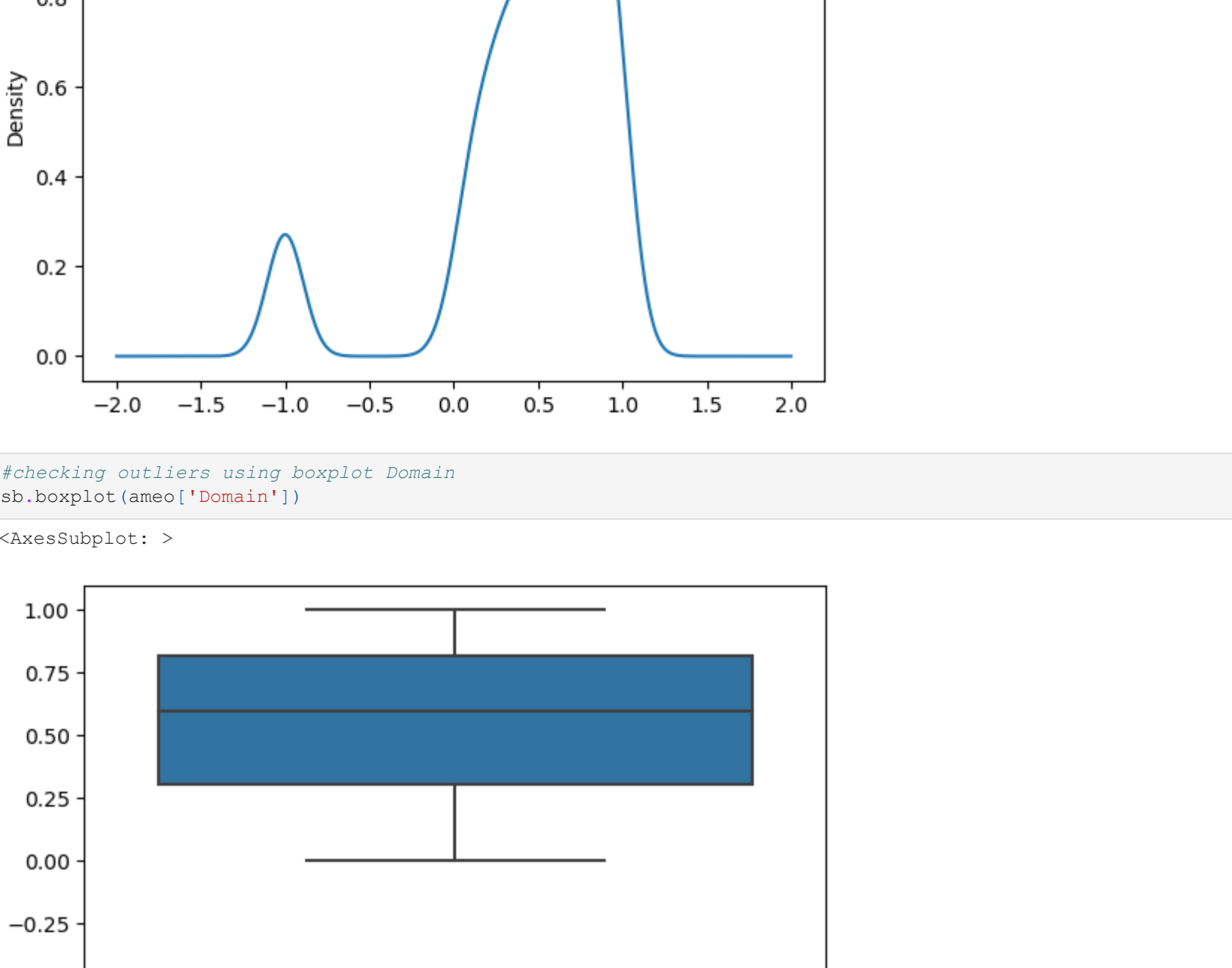
```
# apply log transformation
ameo['collegeGPA_tf'] = 1/ameo['collegeGPA']

fig, axs = plt.subplots(1, 2, figsize=(10, 4))

axs[0].hist(data, bins=50)
axs[0].set_title('Original data')

axs[1].hist(ameo['collegeGPA_tf'], bins=50)
axs[1].set_title('inverse-transformed data')

plt.show()
```



```
#Removing the original
ameo=ameo.drop(['collegeGPA'],axis=1)
ameo.head()
```

	DOI	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	board	CollegeID	...	Domain	Marks_obtained	conscientious
0	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	cbse	1141	..	0.635979	445.0	0
1	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	cbse	5086	..	-1.000000	872.0	-0
2	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	cbse	314	..	0.356536	445.0	1
3	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	cbse	403	..	0.765674	822.0	0
4	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	cbse	2665	..	0.229482	420.0	0

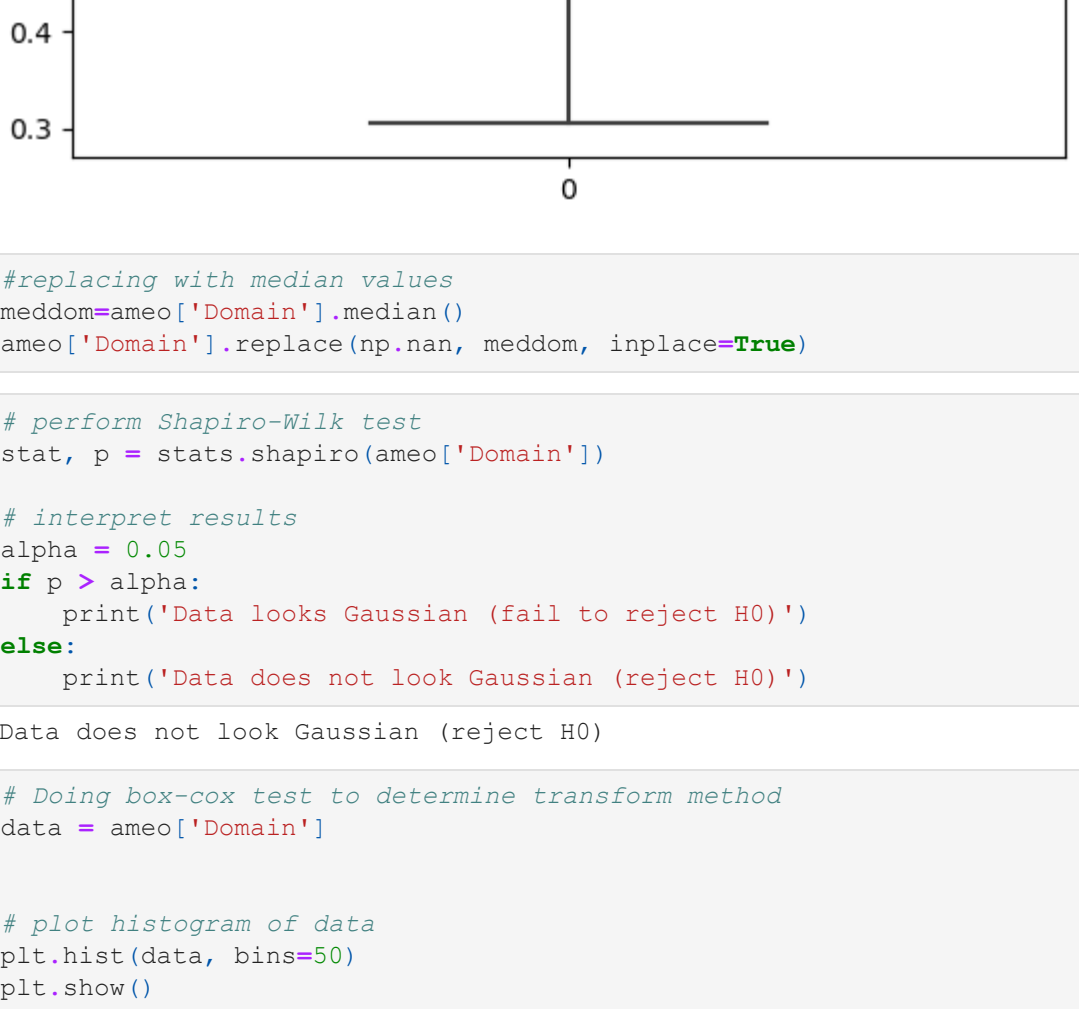
5 rows x 25 columns

### Remarks

The variable GPA was left skewed and had about 243 outliers which was removed and have been replaced with median. The boxcox value was 3 which required an inverse transformation to transform the left skewness of the variable.

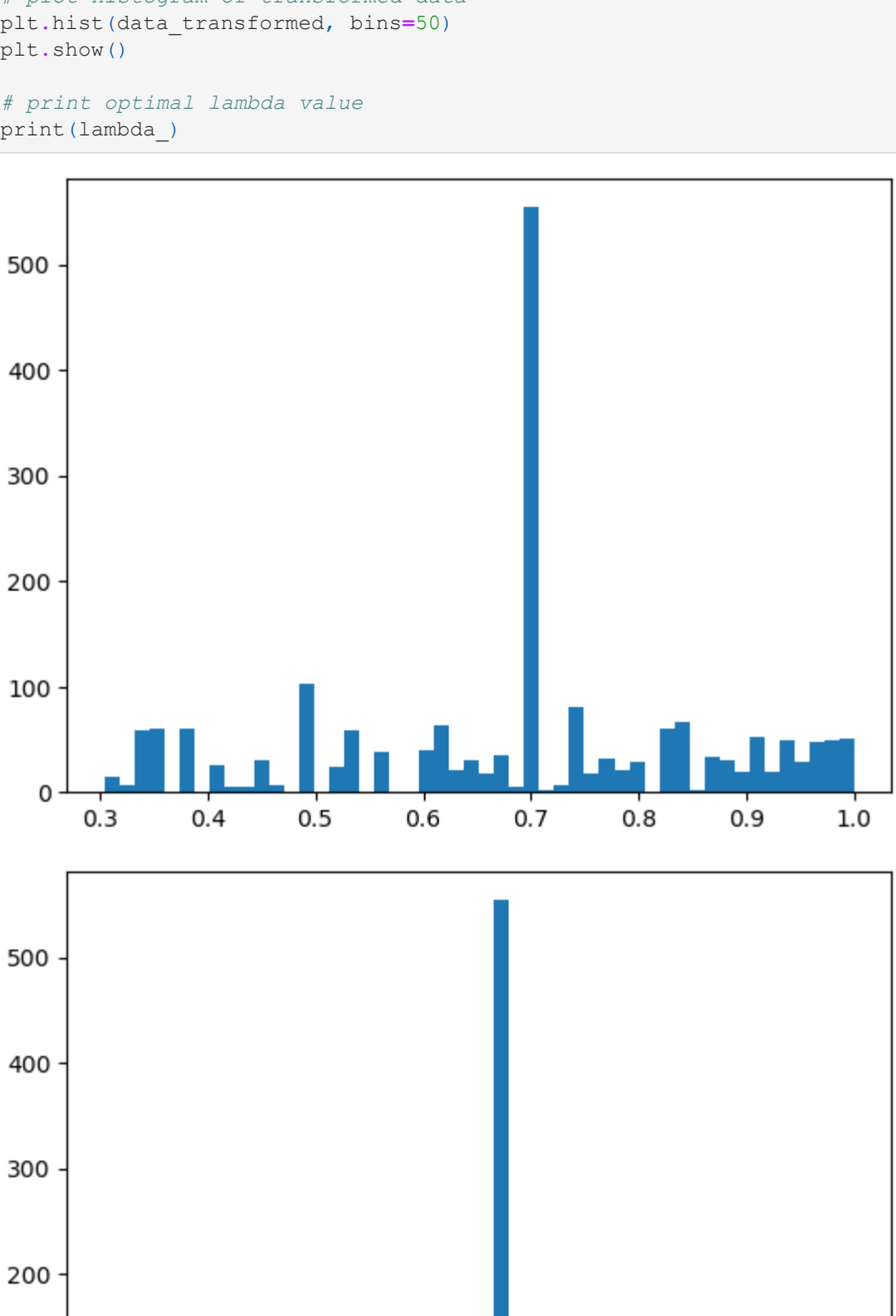
```
In [39]: #checking distribution for Domain
ameo['Domain'].plot.density()
```

```
Out[39]: <AxesSubplot: ylabel='Density'>
```



```
In [48]: #checking outliers using boxplot Domain
sb.boxplot(ameo['Domain'])
```

```
Out[48]: <AxesSubplot: >
```



```
#Outlier treatment for variable Domain
ameo['Domain'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

```
0.10    0.062221
0.25    0.308401
0.50    0.598281
0.70    0.765674
0.90    0.942117
0.95    0.974396
0.99    0.995614
Name: Domain, dtype: float64
```

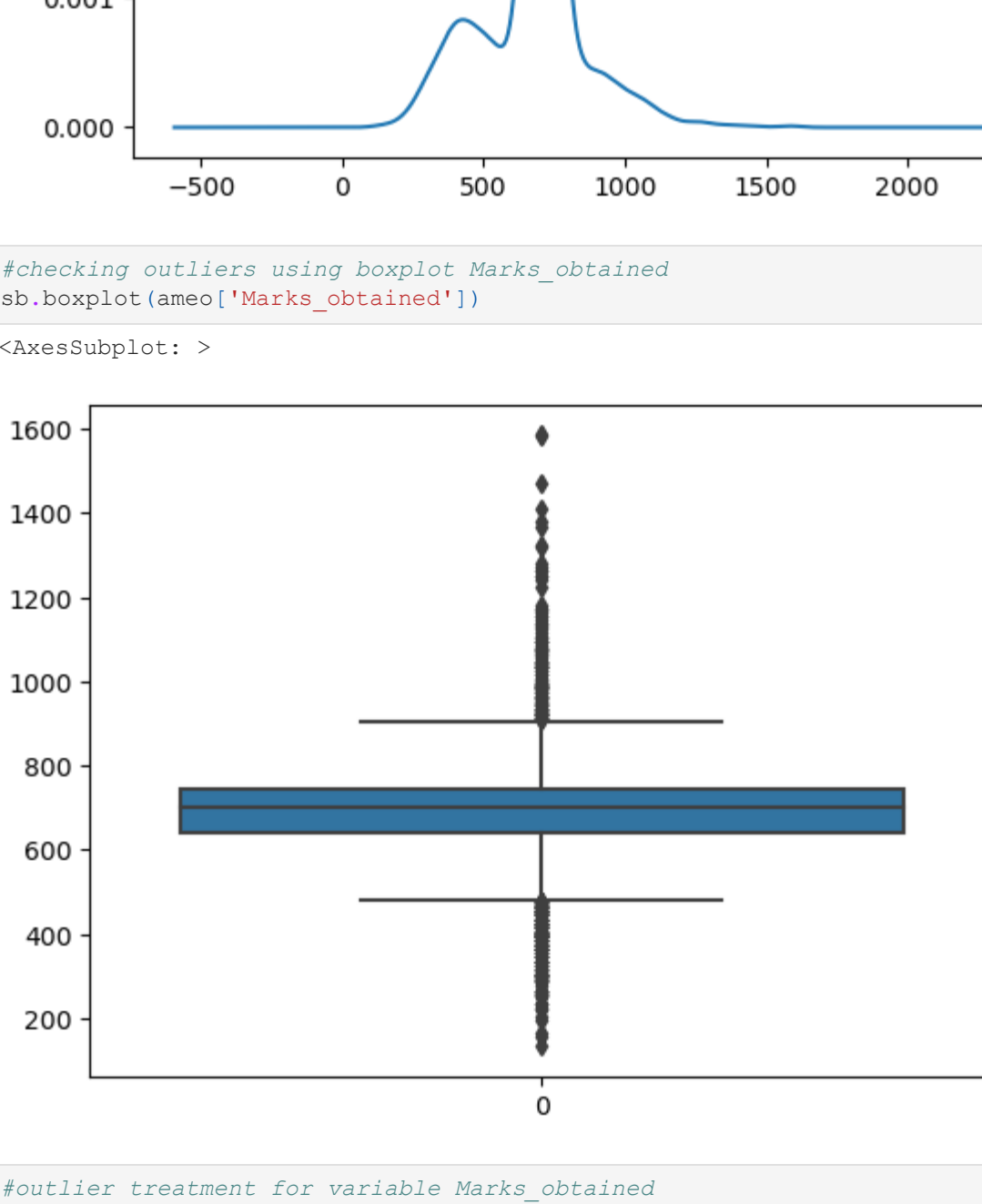
```
#setting limits
dom_lb=ameo['Domain'] < 0.3).copy()
dom_ub=ameo['Domain'] > 0.974396
print(len(dom_lb))
print(len(dom_ub))
```

```
482
```

```
#Removing outliers
for x in ameo['Domain']:
    if x < 0.3:
        ameo['Domain'].replace(x,np.nan,inplace=True)
```

```
In [44]: #checking if outliers gone
sb.boxplot(ameo['Domain'])
```

```
Out[44]: <AxesSubplot: >
```



```
#Replacing with median values
meddom=ameo['Domain'].median()
ameo['Domain'].replace(np.nan, meddom, inplace=True)
```

```
In [46]: # perform Shapiro-Wilk test
stat, p = stats.shapiro(ameo['Domain'])
```

```
# interpret results
alpha = 0.05
if p > alpha:
    print('Data looks Gaussian (fail to reject H0)')
else:
    print('Data does not look Gaussian (reject H0)')
```

Data does not look Gaussian (reject H0)

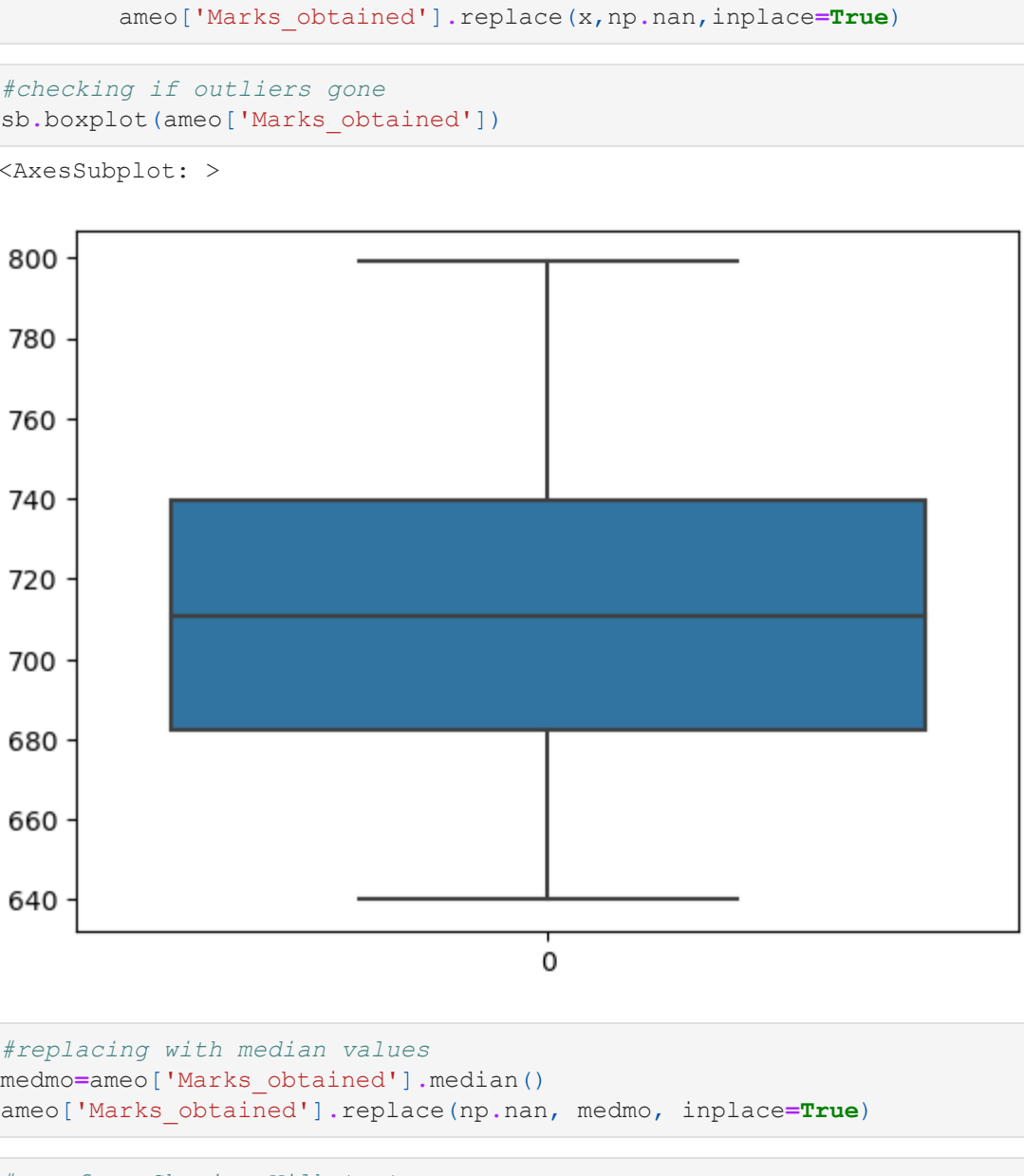
```
# Doing box-cox test to determine transform method
data = ameo['Domain']
```

```
# plot histogram of data
plt.hist(data, bins=50)
plt.show()
```

```
# apply Box-Cox transformation to data
data_transformed, lambda_ = boxcox(data)
```

```
# plot histogram of transformed data
plt.hist(data_transformed, bins=50)
plt.show()
```

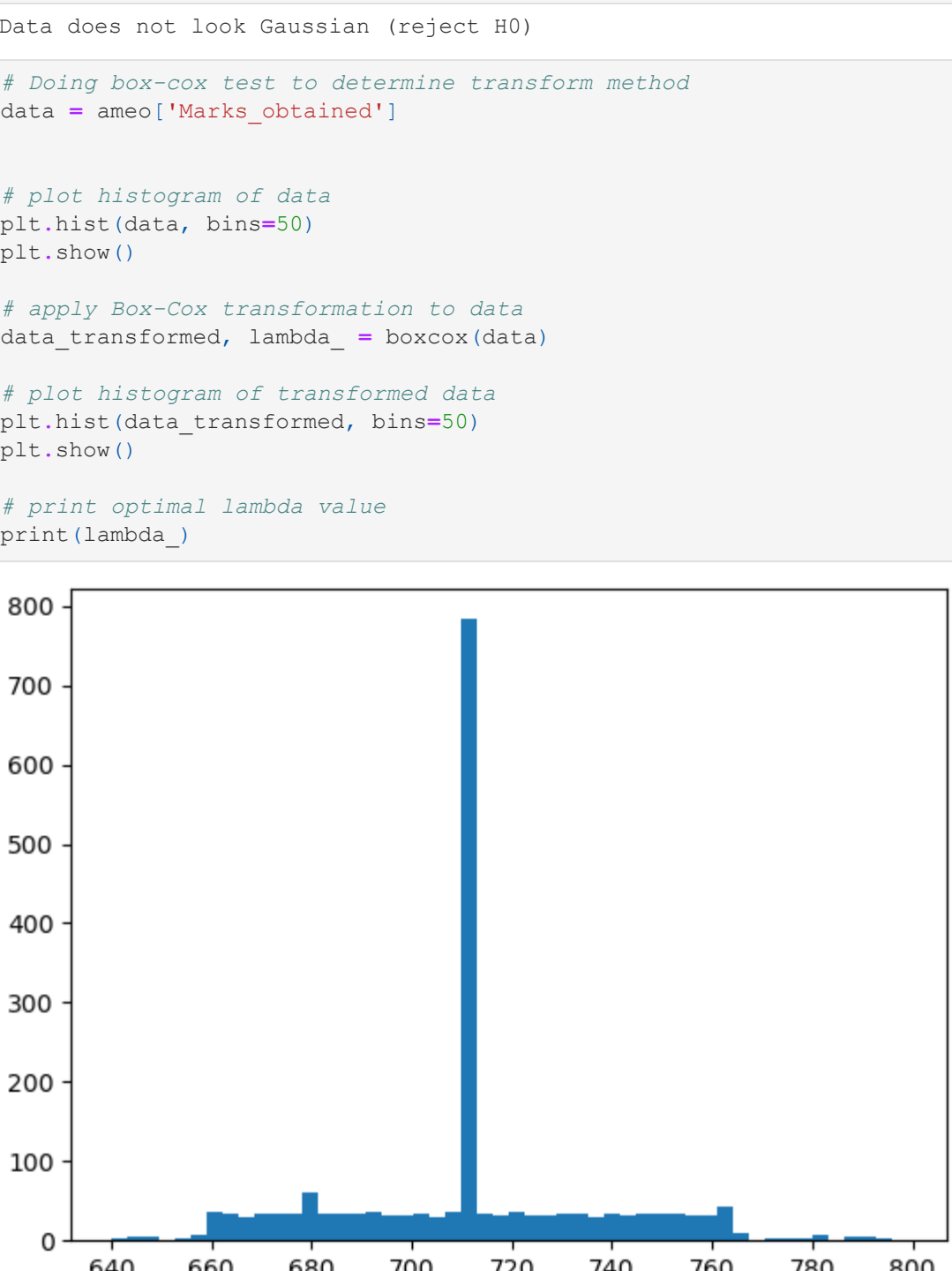
```
# print optimal lambda value
print(lambda_)
```



The variable Domain was left skewed and had about 482 outliers which was removed and have been replaced with median. The boxcox value was 1.2 which required no transformation.

```
In [48]: #checking distribution for Marks_obtained
ameo['Marks_obtained'].plot.density()
```

```
Out[48]: <AxesSubplot: ylabel='Density'>
```



```
In [49]: #checking outliers using boxplot Marks_obtained
sb.boxplot(ameo['Marks_obtained'])
```

```
Out[49]: <AxesSubplot: >
```



```
#Outlier treatment for variable Marks_obtained
ameo['Marks_obtained'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

```
0.10    415.000000
0.25    639.000000
0.50    699.941772
0.70    736.940991
0.90    853.000000
0.95    976.000000
0.99   1164.600000
Name: Marks_obtained, dtype: float64
```

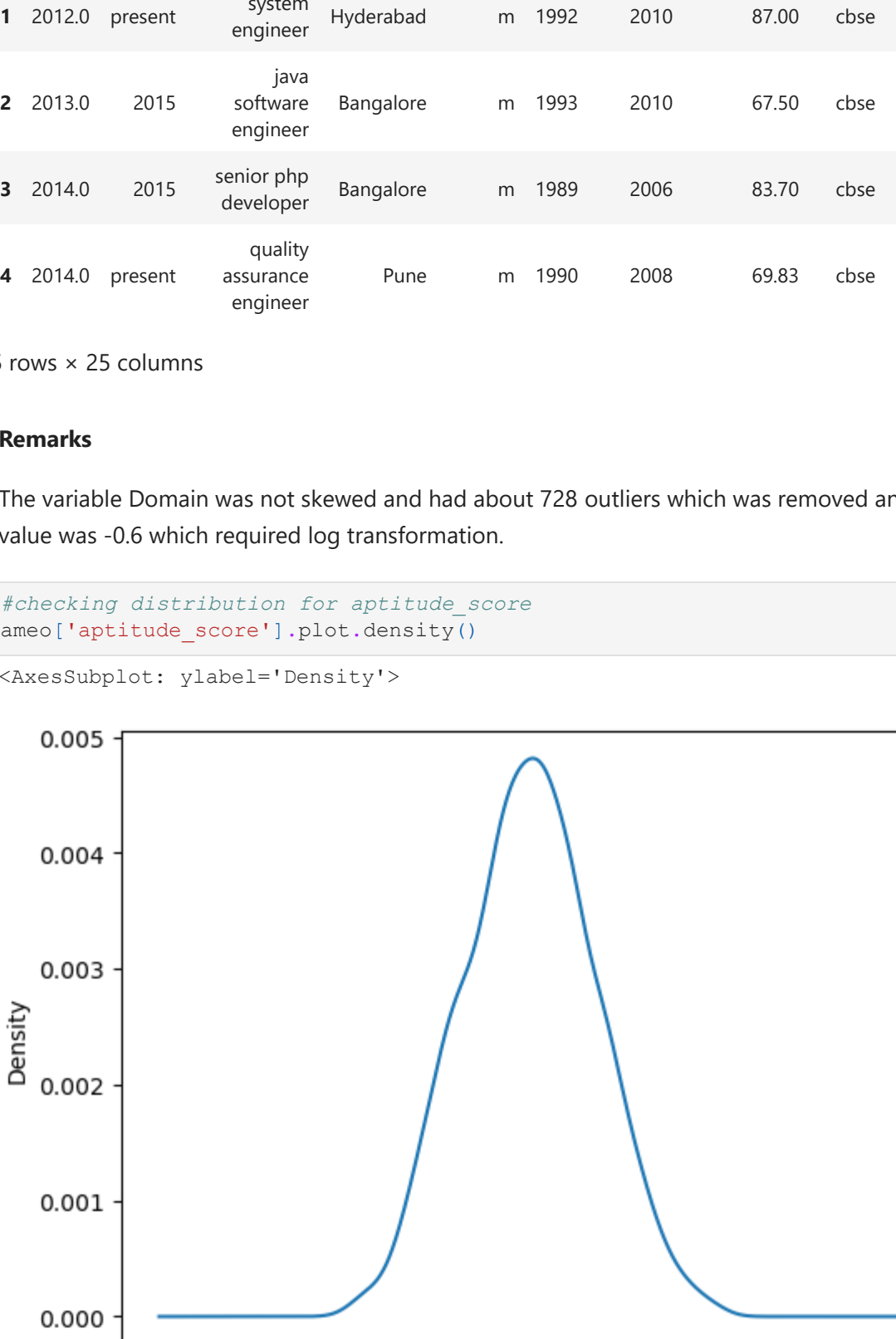
```
#setting limits
marks_lb=ameo['Marks_obtained'] > 820.0).copy()
marks_ub=ameo['Marks_obtained'] < 640.0).copy()
print(len(marks_lb))
print(len(marks_ub))
```

```
230
492
```

```
#Removing outliers
for x in ameo['Marks_obtained']:
    if x > 820.0 or x < 640.0:
        ameo['Marks_obtained'].replace(x,np.nan,inplace=True)
```

```
In [53]: #checking if outliers gone
sb.boxplot(ameo['Marks_obtained'])
```

```
Out[53]: <AxesSubplot: >
```



```
#Replacing with median values
medmarks=ameo['Marks_obtained'].median()
ameo['Marks_obtained'].replace(np.nan, medmarks, inplace=True)
```

```
In [55]: # perform Shapiro-Wilk test
stat, p = stats.shapiro(ameo['Marks_obtained'])
```

```
# interpret results
alpha = 0.05
if p > alpha:
    print('Data looks Gaussian (fail to reject H0)')
else:
    print('Data does not look Gaussian (reject H0)')
```

Data does not look Gaussian (reject H0)

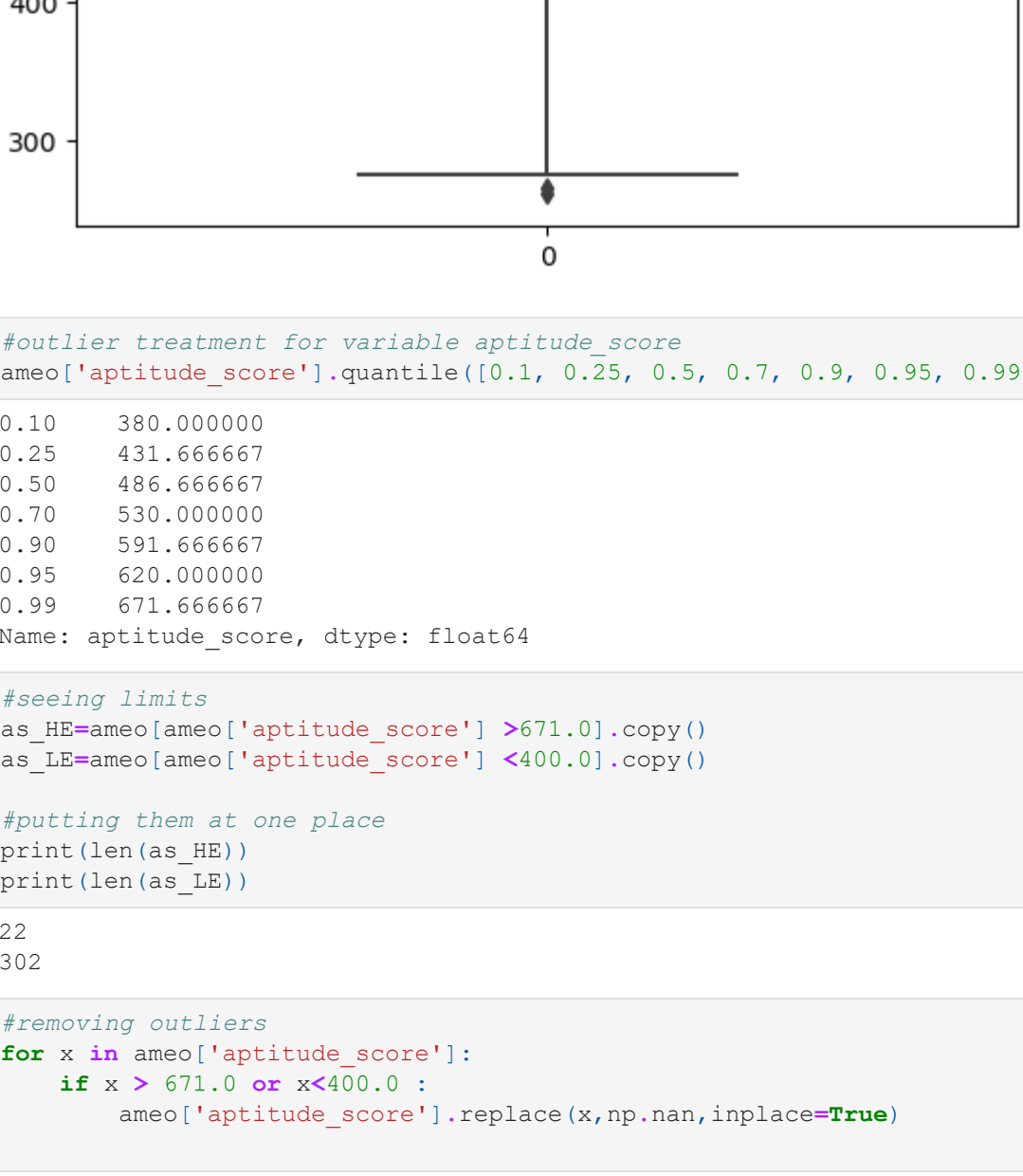
```
# Doing box-cox test to determine transform method
data = ameo['Marks_obtained']
```

```
# plot histogram of data
plt.hist(data, bins=50)
plt.show()
```

```
# apply Box-Cox transformation to data
data_transformed, lambda_ = boxcox(data)
```

```
# plot histogram of transformed data
plt.hist(data_transformed, bins=50)
plt.show()
```

```
# print optimal lambda value
print(lambda_)
```



```
In [57]: #log transforming the variable
```

```
# apply log transformation
ameo['Marks_obtained_tf'] = np.log(ameo['Marks_obtained'])
```

```
fig, axs = plt.subplots(1, 2, figsize=(10, 4))
```

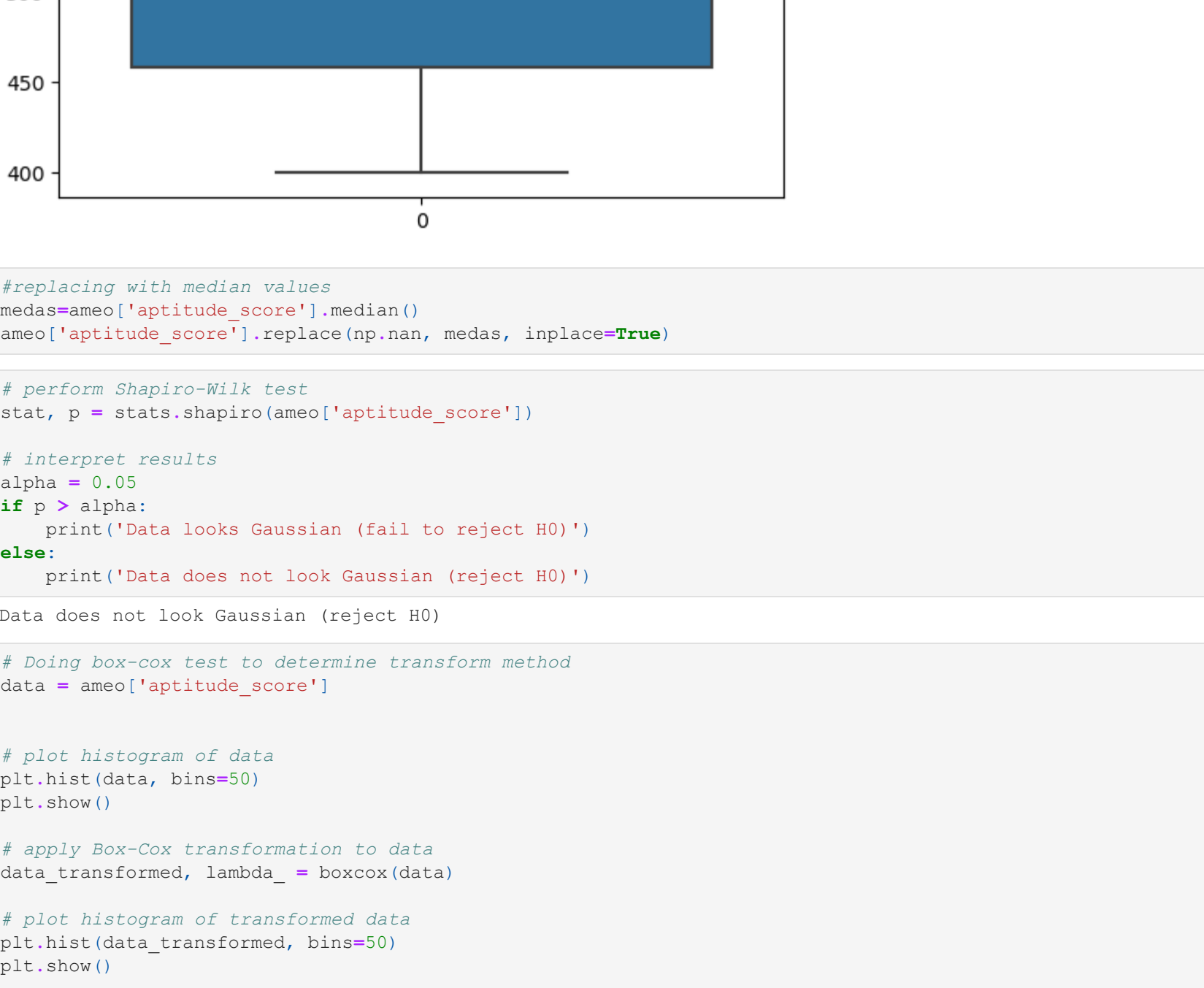
```
axs[0].hist(data, bins=50)
```

```
axs[0].set_title('Original data')
```

```
axs[1].hist(ameo['Marks_obtained_tf'], bins=50)
```

```
axs[1].set_title('log-transformed data')
```

```
plt.show()
```



```
#Removing the original
ameo=ameo.drop(['Marks_obtained'],axis=1)
ameo.head()
```

	DOI	DOL	Designation	JobCity	Gender	DOB	Sgradyr	percentage	board	CollegeID	...	Domain	conscientiousness	agreeableness
0	2012.0	present	senior quality engineer	Bangalore	f	1990	2007	95.80	cbse	1141	..	0.635979	0.9737	0.81
1	2012.0	present	system engineer	Hyderabad	m	1992	2010	87.00	cbse	5086	..	0.694479	-0.3027	-0.62
2	2013.0	2015	java software engineer	Bangalore	m	1993	2010	67.50	cbse	314	..	0.356536	1.7081	-0.10
3	2014.0	2015	senior php developer	Bangalore	m	1989	2006	83.70	cbse	403	..	0.765674	0.0464	0.81
4	2014.0	present	quality assurance engineer	Pune	m	1990	2008	69.83	cbse	2665	..	0.694479	0.1282	0.54

5 rows x 25 columns

### Remarks

The variable Domain was not skewed and had about 728 outliers which was removed and have been replaced with median. The boxcox value was -0.6 which required log transformation.

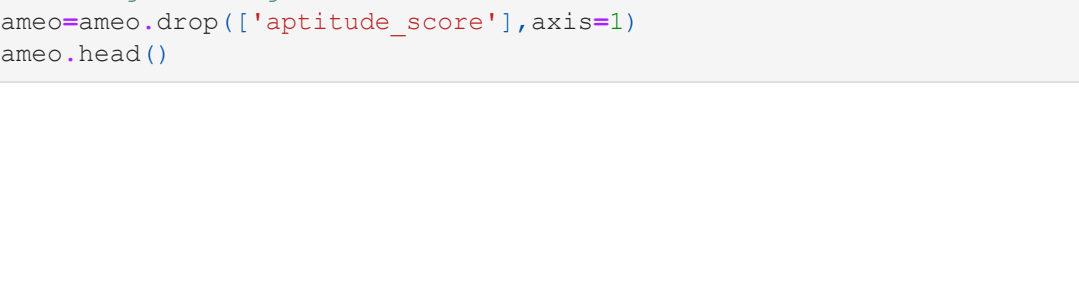
```
In [59]: #checking distribution for aptitude_score
ameo['aptitude_score'].plot.density()
```

```
Out[59]: <AxesSubplot: ylabel='Density'>
```



```
#checking outliers using boxplot aptitude_score
sb.boxplot(ameo['aptitude_score'])
```

```
Out[60]: <AxesSubplot: >
```



```
#Outlier treatment for variable aptitude_score
ameo['aptitude_score'].quantile([0.1, 0.25, 0.5, 0.7, 0.9, 0.95, 0.99])
```

```
0.10    380.000000
0.25    431.666667
0.50    486.666667
0.70    530.000000
0.90    591.666667
0.95    620.000000
0.99    671.666667
Name: aptitude_score, dtype: float64
```

```
#setting limits
as_lb=ameo['aptitude_score'] > 671.0).copy()
as_ub=ameo['aptitude_score'] < 400.0).copy()
print(len(as_lb))
print(len(as_ub))
```

```
22
302
```

```
#Removing outliers
for x in ameo['aptitude_score']:
    if x > 671.0 or x < 400.0:
        ameo['aptitude_score'].replace(x,np.nan,inplace=True)
```

```
In [64]: #checking if outliers gone
sb.boxplot(ameo['aptitude_score'])
```

```
Out[64]: <AxesSubplot: >
```



```
#Replacing with median values
medas=ameo['aptitude_score'].median()
ameo['aptitude_score'].replace(np.nan, medas, inplace=True)
```

```
In [66]: # perform Shapiro-Wilk test
stat, p = stats.shapiro(ameo['aptitude_score'])
```

```
# interpret results
alpha = 0.05
if p > alpha:
    print('Data looks Gaussian (fail to reject H0)')
else:
    print('Data does not look Gaussian (reject H0)')
```

Data does not look Gaussian (reject H0)

```
# Doing box-cox test to determine transform method
data = ameo['aptitude_score']
```

```
# plot histogram of data
plt.hist(data, bins=50)
plt.show()
```

```
# apply Box-Cox transformation to data
data_transformed, lambda_ = boxcox(data)
```

```
# plot histogram of transformed data
plt.hist(data_transformed, bins=50)
plt.show()
```

```
# print optimal lambda value
print(lambda_)
```



```
In [68]: #log transforming the variable
```

```
# apply log transformation
ameo['aptitude_score_tf'] = np.log(ameo['aptitude_score'])
```

```
fig, axs = plt.subplots(1, 2, figsize=(10, 4))
```

```
axs[0].hist(data, bins=50)
```

```
axs[0].set_title('Original data')
```

```
axs[1].hist(ameo['aptitude_score_tf'], bins=50)
```

```
axs[1].set_title('log-transformed data')
```

```
plt.show()
```

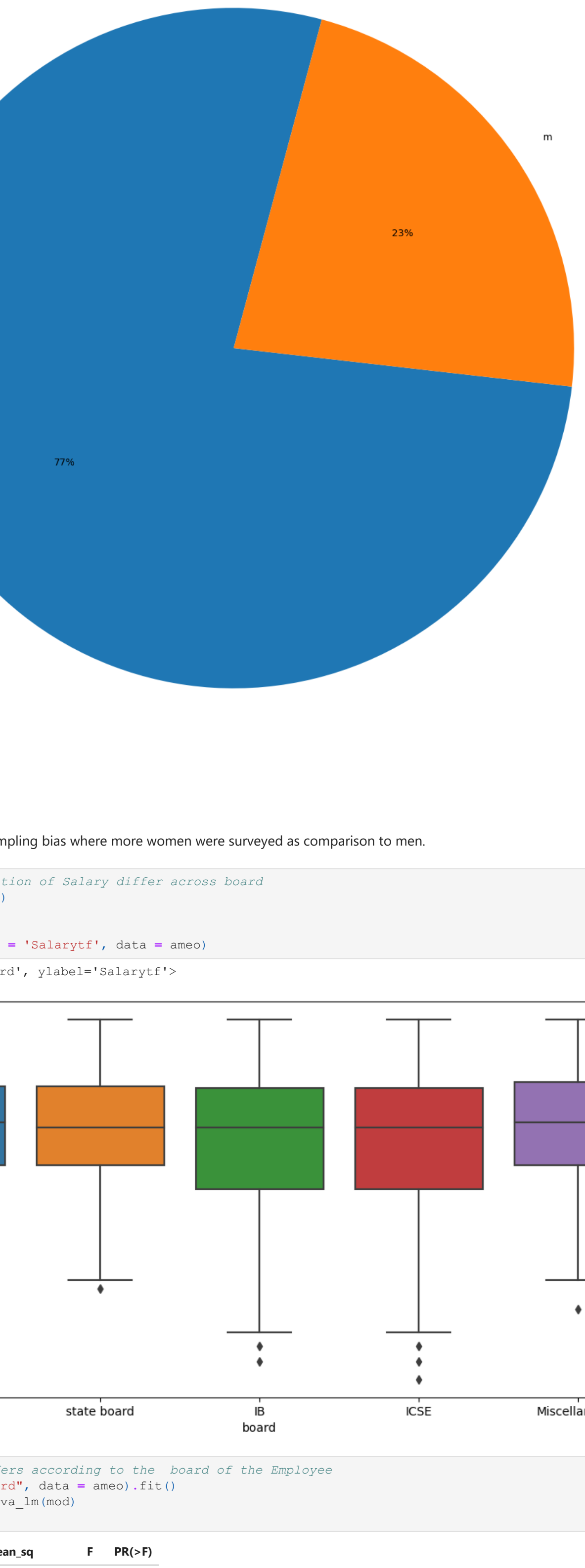


```
#Removing the original
ameo=ameo.drop(['aptitude_score'],axis=1)
ameo.head()
```









Observation

We can confirm that this is a sampling bias where more women were surveyed as comparison to men.

```
In [94]: #Box plot to see distribution of Salary differ across board
plt.figure(figsize=(12,6))

#b.boxplot(x = "board", y = "Salarytf", data = ameo)

Out[94]: <AxesSubplot: xlabel='board', ylabel='Salarytf'>
```



It can be clearly seen that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [97]: #Box plot to see distribution of Salary across CollegeCityTier
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeCityTier", grid = False)

Out[97]: <AxesSubplot: title='center: CollegeCityTier', xlabel='Salarytf'>
```

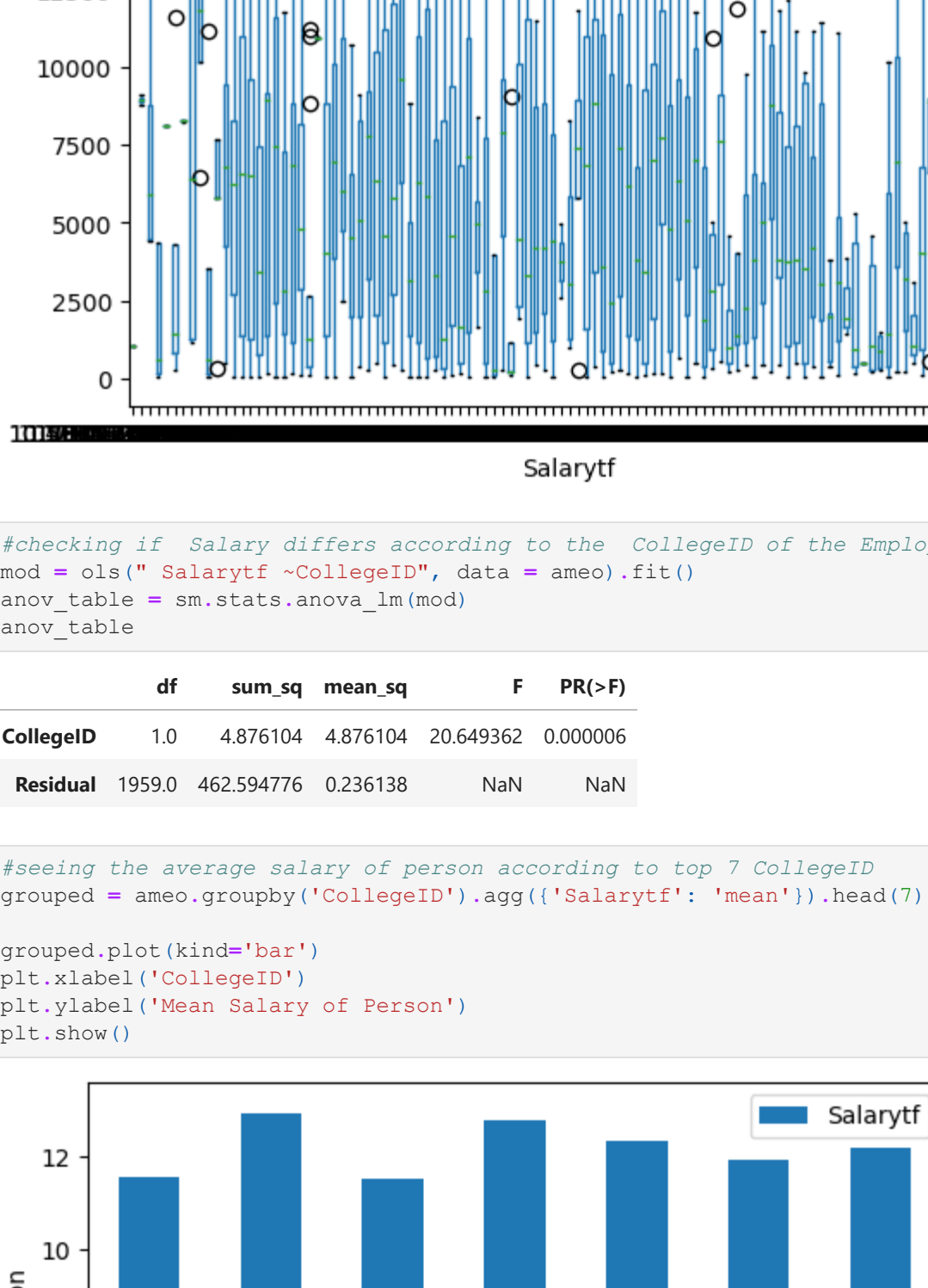


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [97]: #Box plot to see distribution of Salary across CollegeCityTier
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeCityTier", grid = False)

Out[97]: <AxesSubplot: title='center: CollegeCityTier', xlabel='Salarytf'>
```

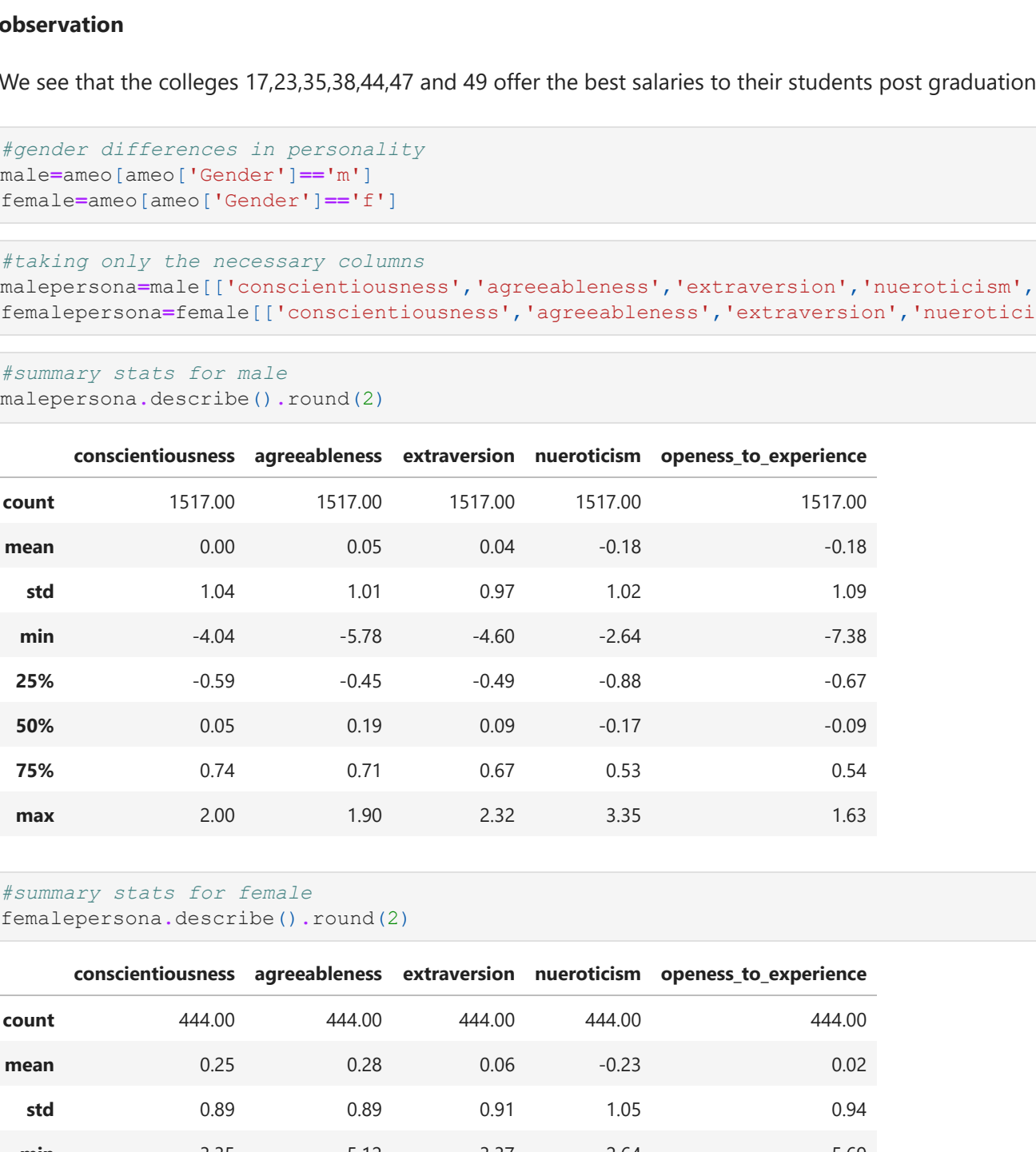


Observation

It can be clearly seen that Tier one cities college allow students to get a higher paying job.

```
In [100]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[100]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```

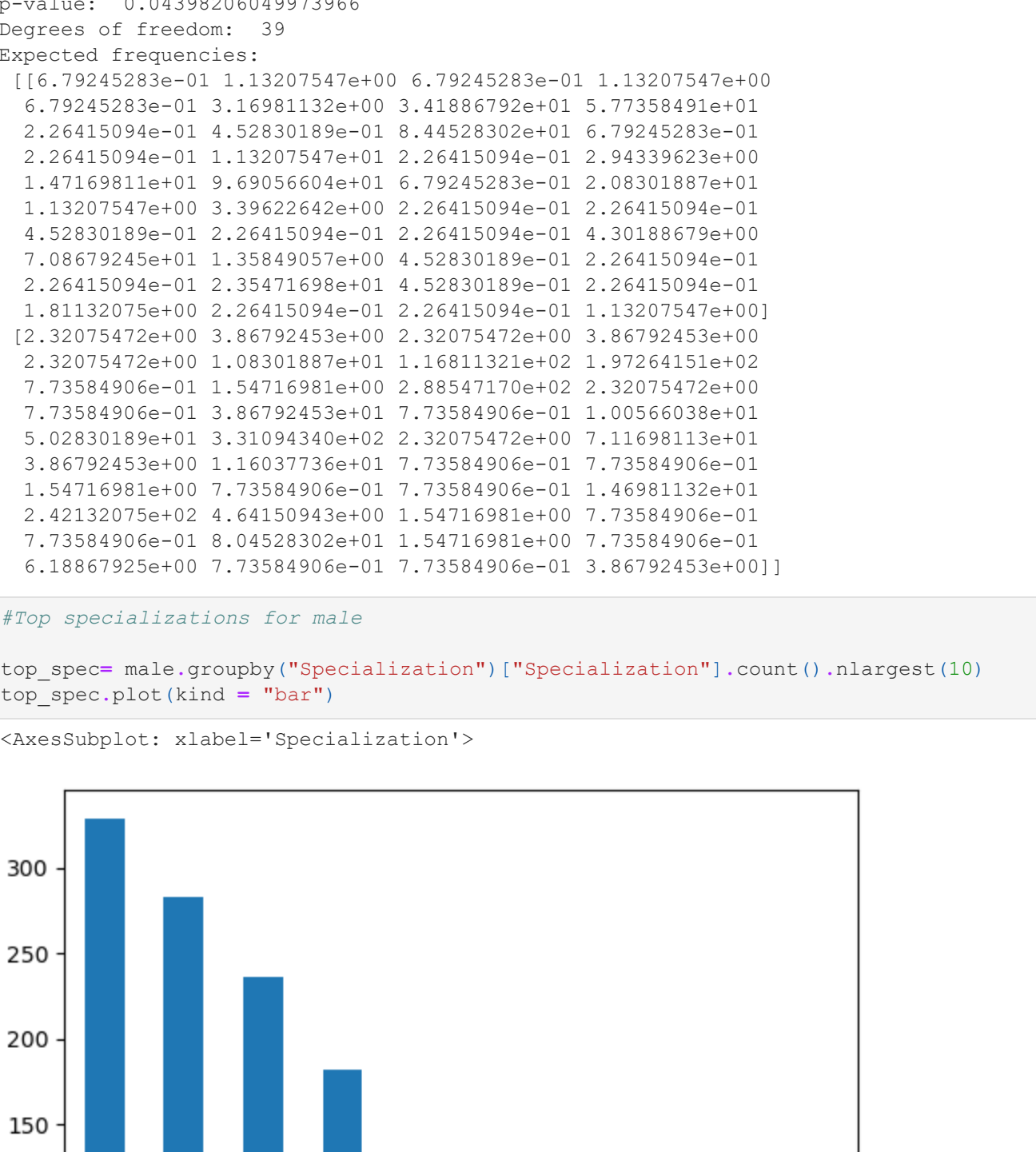


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [101]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[101]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```

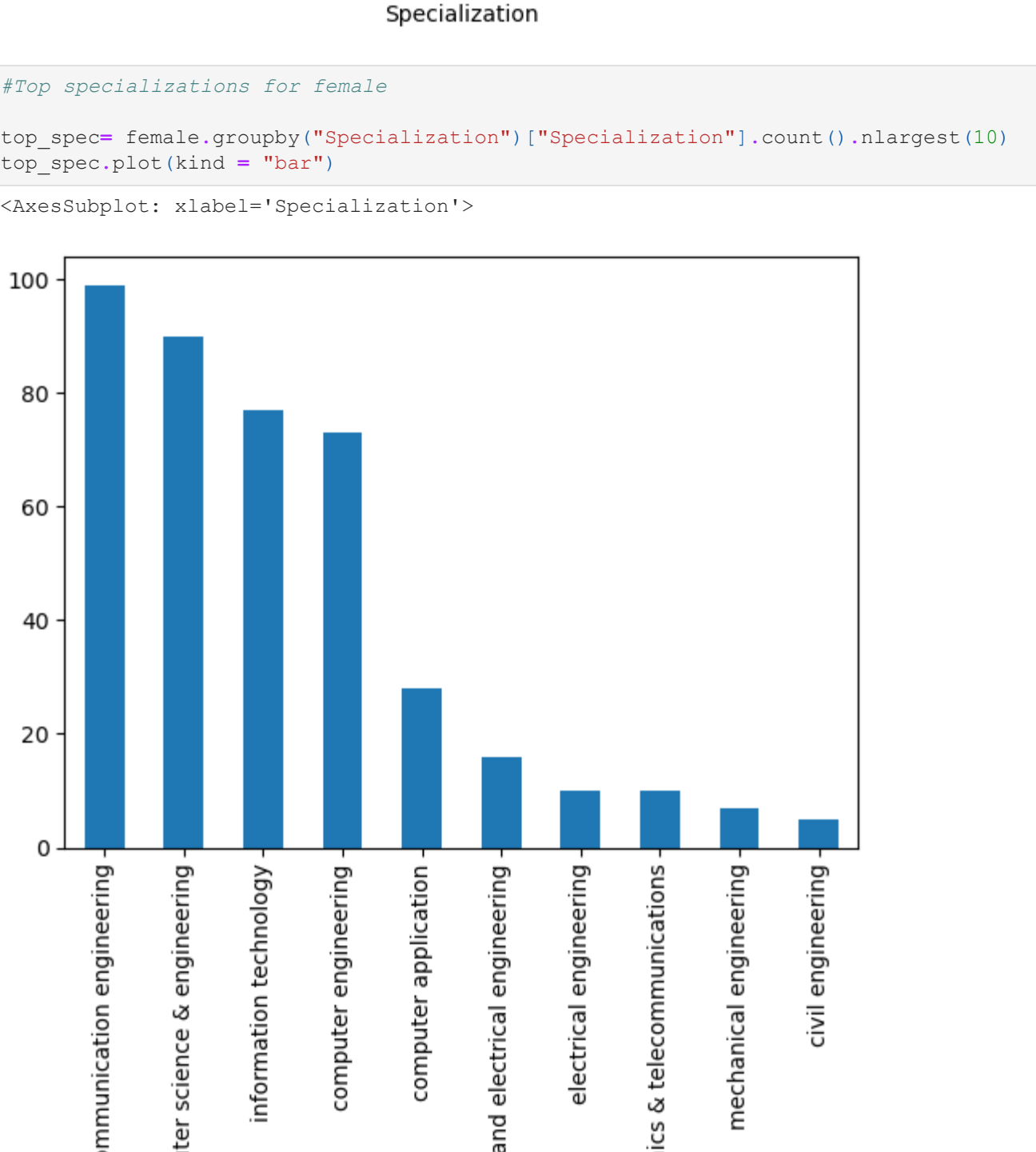


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [102]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[102]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```

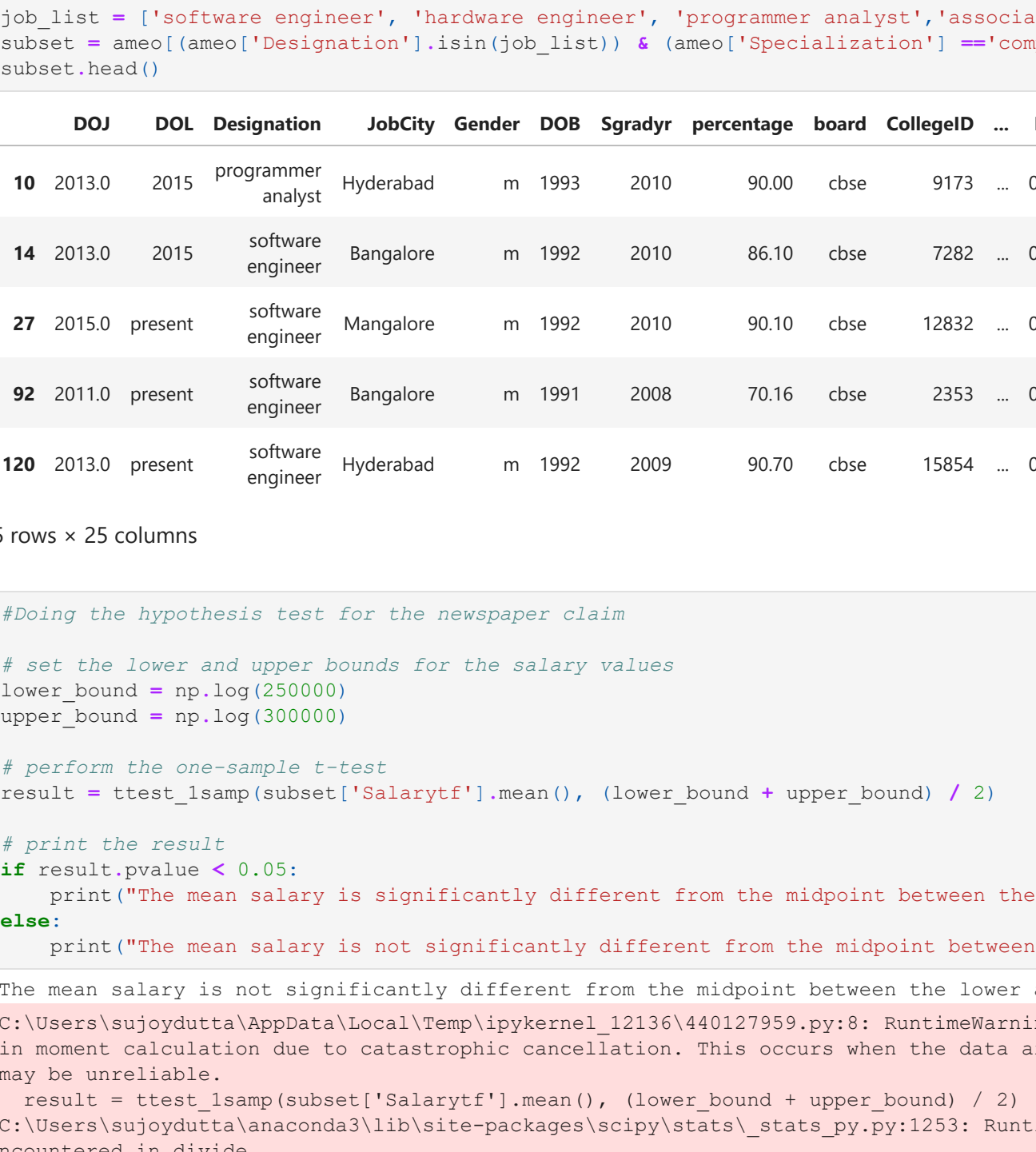


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [103]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[103]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```



Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [104]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[104]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```



Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [105]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[105]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```

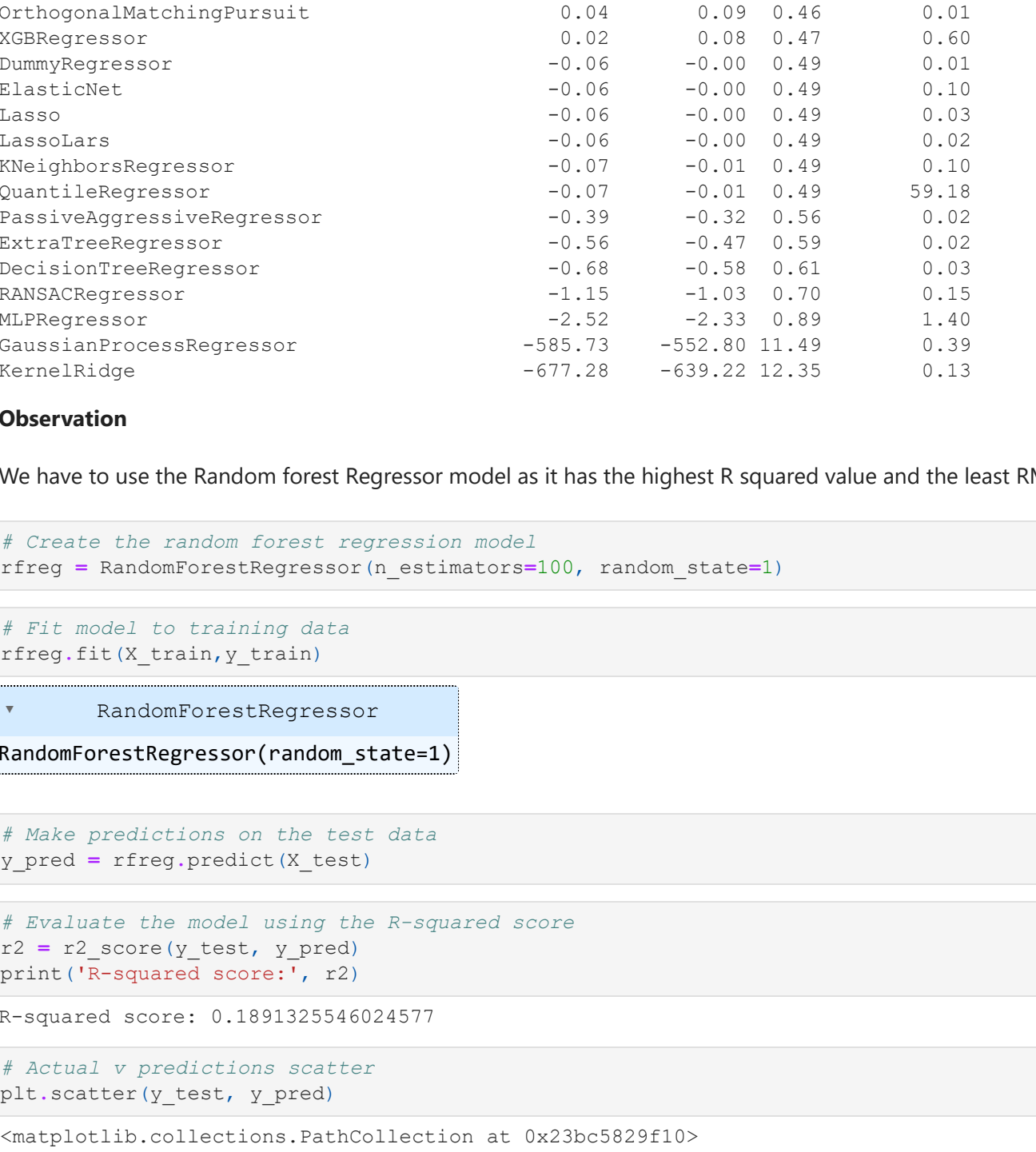


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [106]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[106]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```

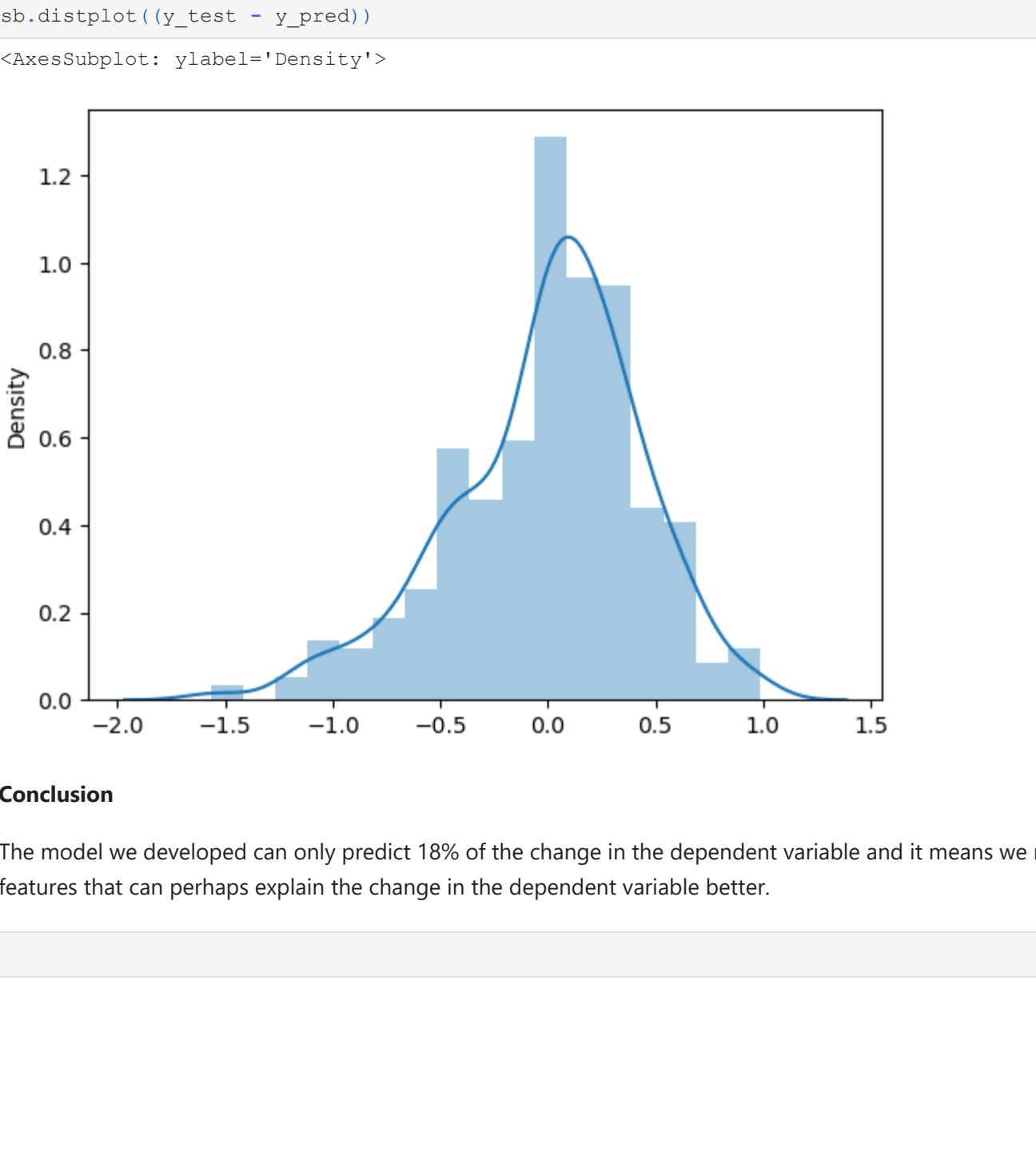


Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.

```
In [107]: #Box plot to see distribution of Salary across CollegeID
plt.figure(figsize=(12,6))
ameo.boxplot(x="Salarytf",column="CollegeID", grid = False)

Out[107]: <AxesSubplot: title='center: CollegeID', xlabel='Salarytf'>
```



Observation

We can see that there is no significant difference in average salary, hence we can say that the board a person studied from is not a deciding factor of the amount of salary they get.