

Problem statement

Build the linear regression model using scikit learn in boston data to predict 'Price' based on other independent variable.

```
In [11]: import numpy as np
import pandas as pd
import scipy.stats as stats
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.datasets import load_boston
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

In [12]: # loading dataset and converting Pandas DataFrame
df = load_boston()
boston = pd.DataFrame(data=df.data, columns=df.feature_names)
boston['PRICE'] = df.target

In [13]: boston.head()
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTCRATIO	B	LSTAT	PRICE
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	15.3	366.90	4.98	24.0
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	366.90	9.14	21.6
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	362.63	4.03	34.7
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	364.63	2.94	33.4
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	366.90	5.33	36.2

```
In [14]: boston.shape
(506, 14)

In [15]: # Checking missing values
boston.isna().sum().sum() # No missing values present
0

In [16]: boston.info()
```

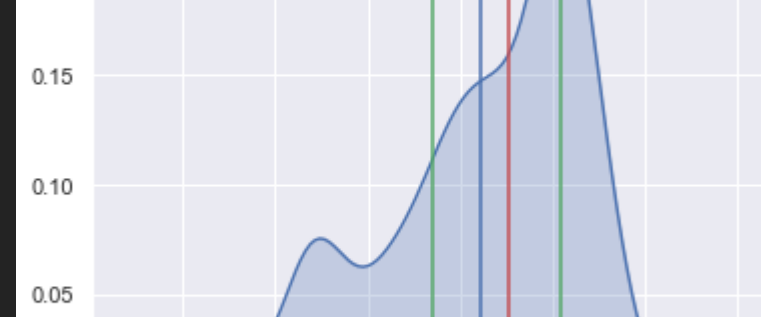
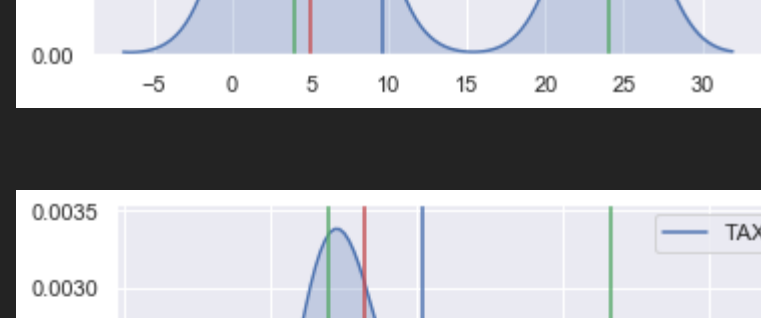
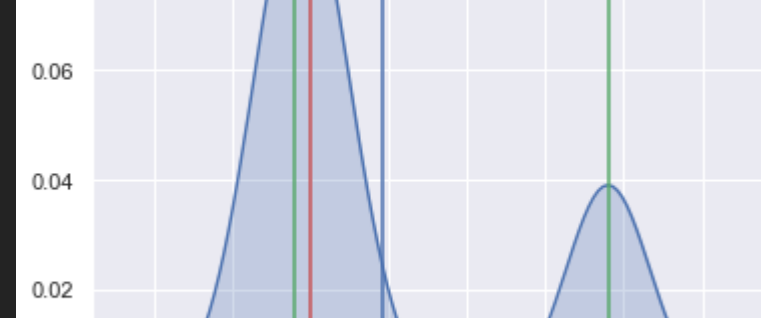
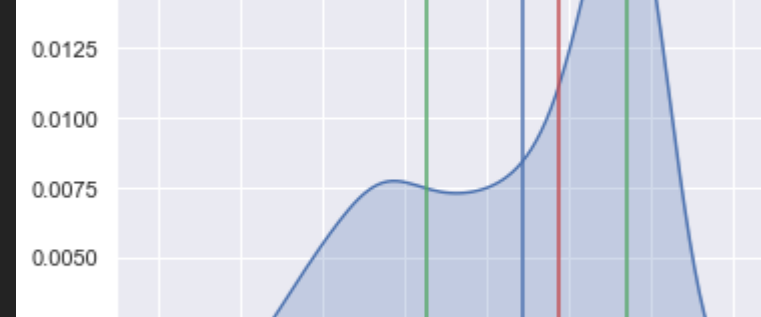
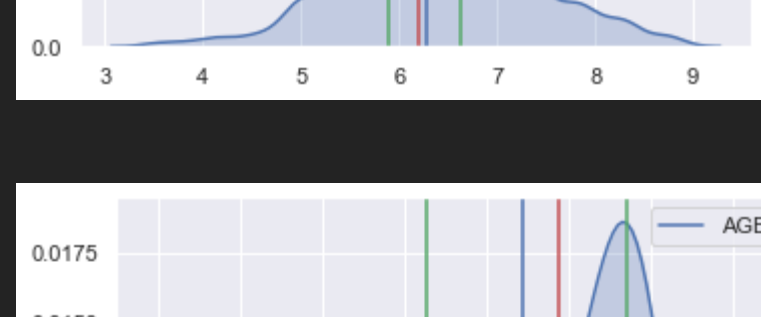
```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 506 entries, 0 to 505
Data columns (total 14 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   CRIM        506 non-null    float64
 1   ZN          506 non-null    float64
 2   INDUS       506 non-null    float64
 3   CHAS        506 non-null    float64
 4   NOX         506 non-null    float64
 5   RM          506 non-null    float64
 6   AGE         506 non-null    float64
 7   DIS         506 non-null    float64
 8   RAD         506 non-null    float64
 9   TAX         506 non-null    float64
10  PTRATIO     506 non-null    float64
11  B           506 non-null    float64
12  LSTAT       506 non-null    float64
13  dtype: float64[14]
memory usage: 55.5 KB
```

Exploratory Data Analysis

```
In [17]: # function for plotting KDE Plot along with mean, median and quantiles
def kde_viz(df, var):
    plt.figure(figsize=(10,1))
    sns.set(color_codes=True)
    sns.kdeplot(df[var], shade=True)

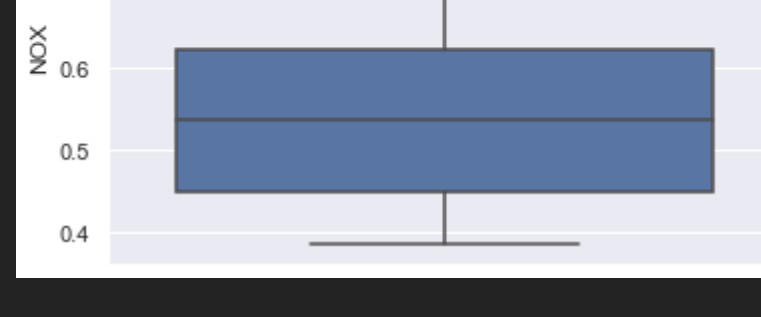
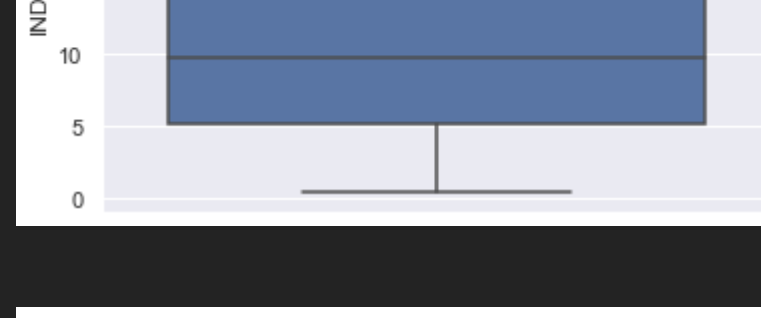
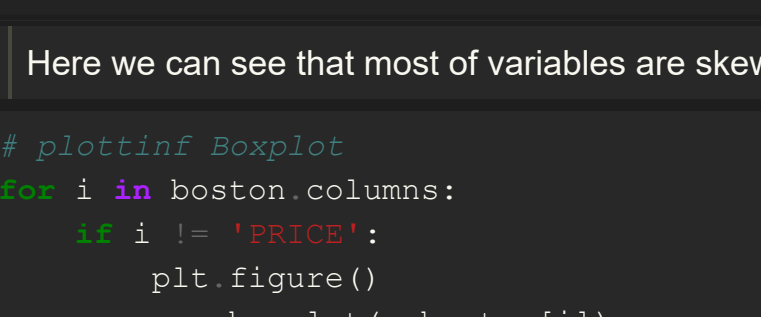
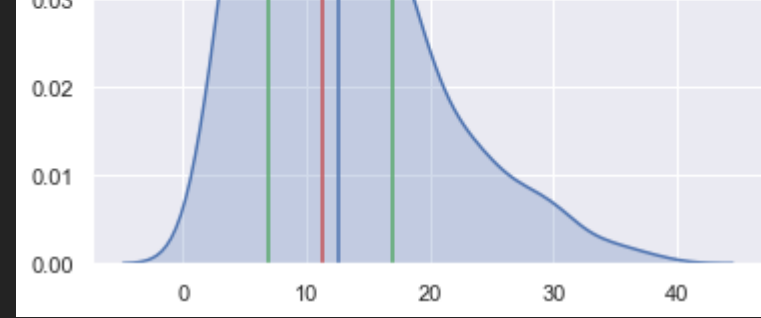
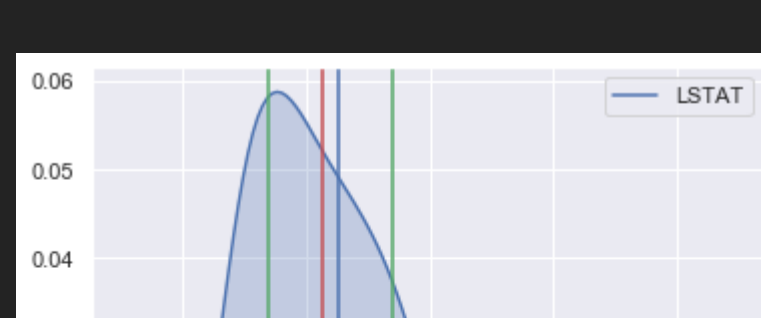
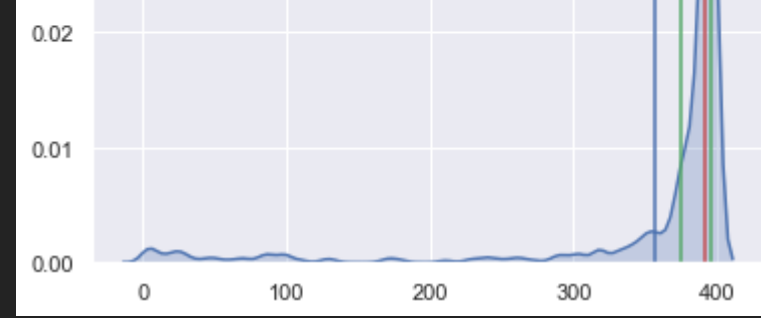
    plt.axvline(df[var].mean())
    plt.axvline(df[var].median(), color='r')
    plt.axvline(df[var].quantile(.25), color='g')
    plt.axvline(df[var].quantile(.75), color='b')
```

```
In [18]: var_col in boston.columns:
var_col != 'CHAS' and col != 'PRICE': # bec CHAS is discrete
    kde_viz(boston, col)
```



Here we can see that most of variables are skewed, non-normalized and have outliers.

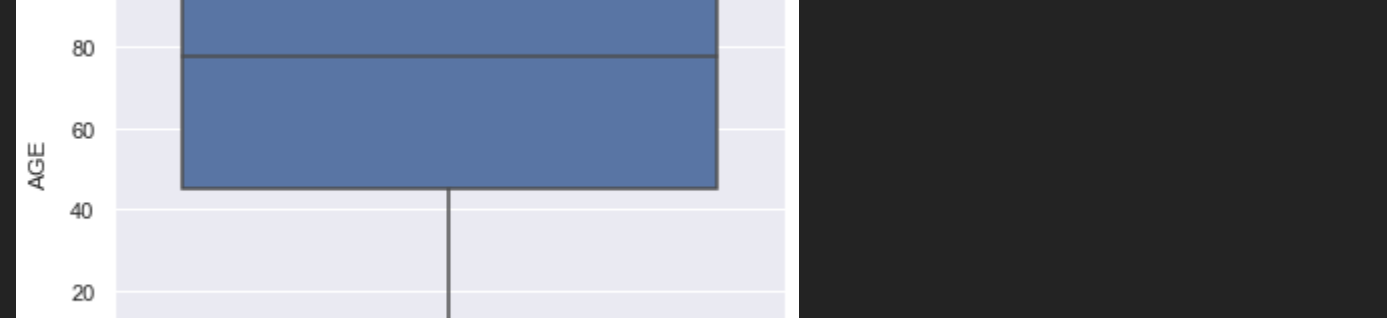
```
In [19]: # plotting Boxplot
var_col in boston.columns:
var_col != 'CHAS':
    plt.figure()
    sns.boxplot(y=boston[col])
```



Here we can see that the CRIM, ZN, RM, B have some outliers so we will handle them later.

```
In [19]: # We have 2 discrete variables. So lets analyse them
fig, ax = plt.subplots(2, 2, figsize=(20, 10))

sns.lineplot(boston['CHAS'], boston['PRICE'], palette='Set1', ax=ax[0][0])
sns.countplot(boston['CHAS'], palette='Set1', ax=ax[0][1])
sns.boxplot(boston['CHAS'], boston['PRICE'], ax=ax[1][0], palette='Set1')
sns.barplot(boston['CHAS'], boston['PRICE'], ax=ax[1][1], palette='Set1')
```



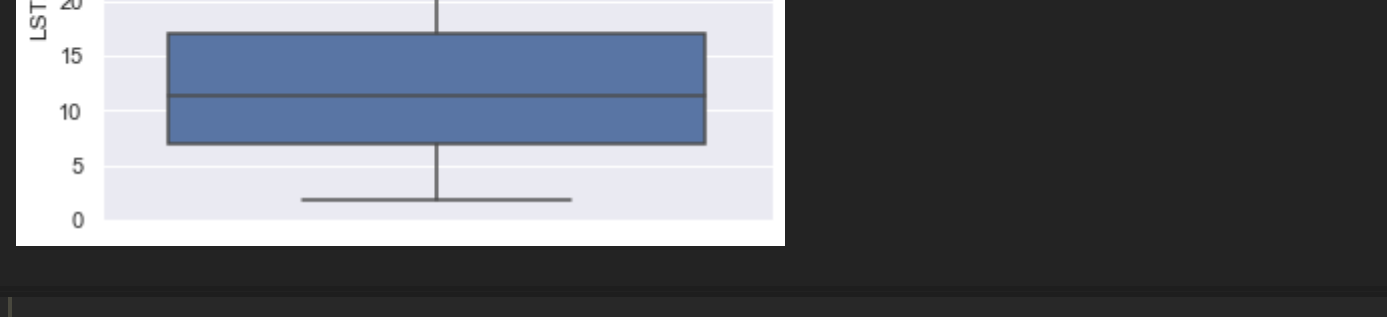
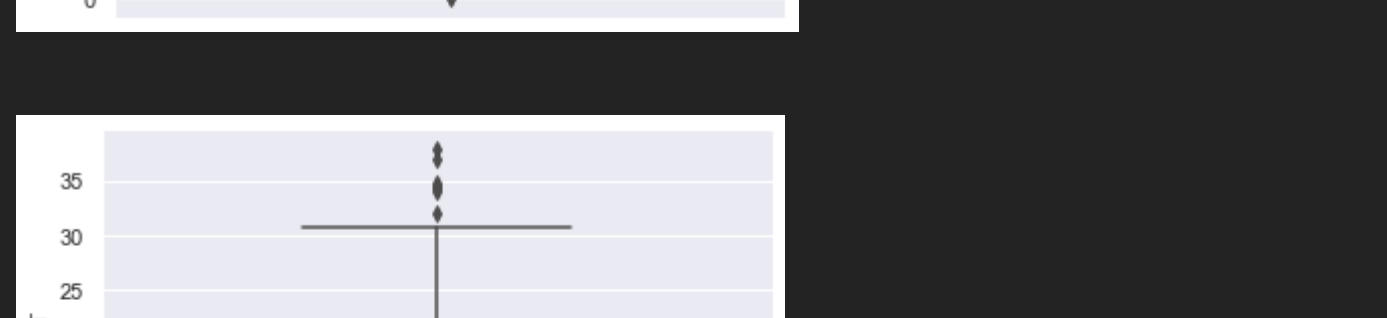
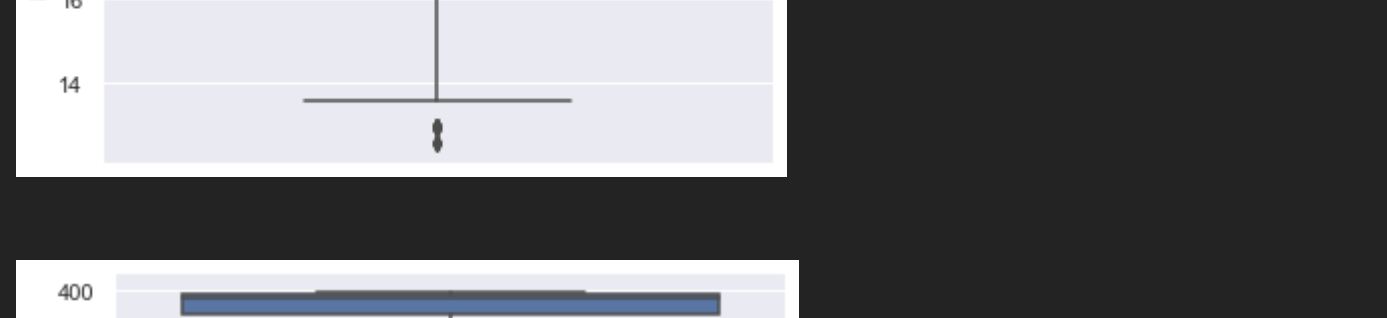
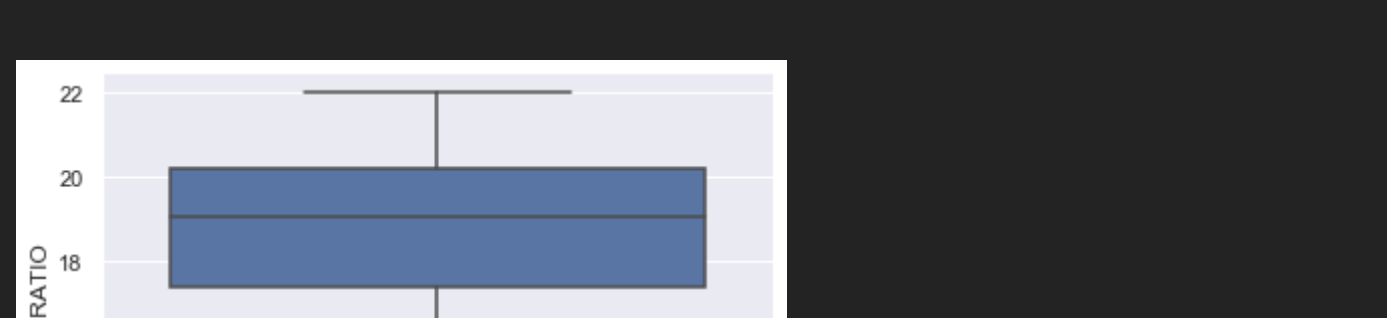
```
<matplotlib.axes._subplots.AxesSubplot at 0a22fac64508>
```



We can see that there is some sort of relation b/w PRICE and the CHAS. PRICE increases if tract bounds the river.

```
In [19]: fig, ax = plt.subplots(2, 2, figsize=(20, 10))

sns.lineplot(boston['RAD'], boston['PRICE'], palette='Blues_d', ax=ax[0][0])
sns.countplot(boston['RAD'], palette='Blues_d', ax=ax[0][1])
sns.boxplot(boston['RAD'], boston['PRICE'], ax=ax[1][0], palette='Blues_d')
sns.barplot(boston['RAD'], boston['PRICE'], ax=ax[1][1], palette='Blues_d')
```



We don't see much of the observations, but we can say that PRICE is greater than 25 for the RAD 2, 3, 5, 7, and 8.

Variable Transformation

We have seen that most of the features have outliers. So we try to transform our data and then see if there are outlier present or not.

```
In [19]: X = boston.drop(['PRICE'], axis=1)
Y = boston['PRICE']

# splitting into train, test
X_tr, X_ts, y_tr, y_ts = train_test_split(X, Y, test_size=0.3, random_state=4)
```

```
In [19]: discrete_cols = [col for col in X.columns if X[col].nunique() < 50]
sns.boxplot(discrete_cols)
```

```
['CRIM', 'INDUS', 'NOX', 'RM', 'AGE', 'DIS', 'TAX', 'B', 'LSTAT']
```

```
In [19]: from feature_engine.variable_transformers import YeoJohnsonTransformer
box_cox_transformer = BoxCoxTransformer
```

```
# fitting data to YeoJohnson Transformer
yco_tf = YeoJohnsonTransformer()
yco_tr = yco_tf.fit_transform(X_tr)
```

```
# grabbing non-zero columns because box-cox transformation does not work with 0 values
non_zero = []
for feature in X_tr.columns:
    if 0 in X_tr[feature].unique():
        pass
    else:
        non_zero.append(feature)
```

```
['CRIM', 'INDUS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT']
```

```
In [19]: # fitting data to BOXCOX Transformer
box_tf = BoxCoxTransformer(variables=non_zero)
box_tr = box_tf.fit_transform(X_tr)
```

```
# transforming Test data
yco_ts = yco_tf.transform(X_ts)
box_ts = box_tf.transform(X_ts)
```

```
In [19]: # diagnostic plots (df, variable):
'''Function that plots Histogram, Q-Q Plot and Boxplot for visualizing the outliers and distribution of the data'''
plt.figure(figsize=(10, 1))
```

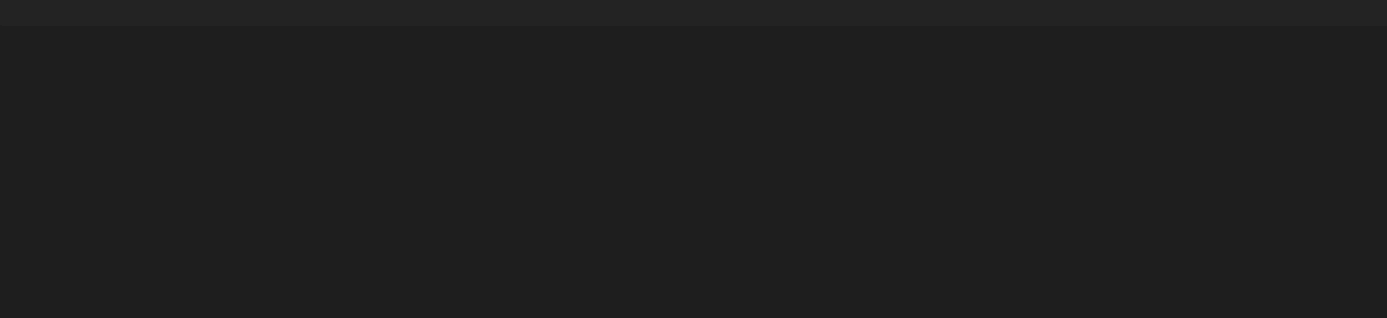
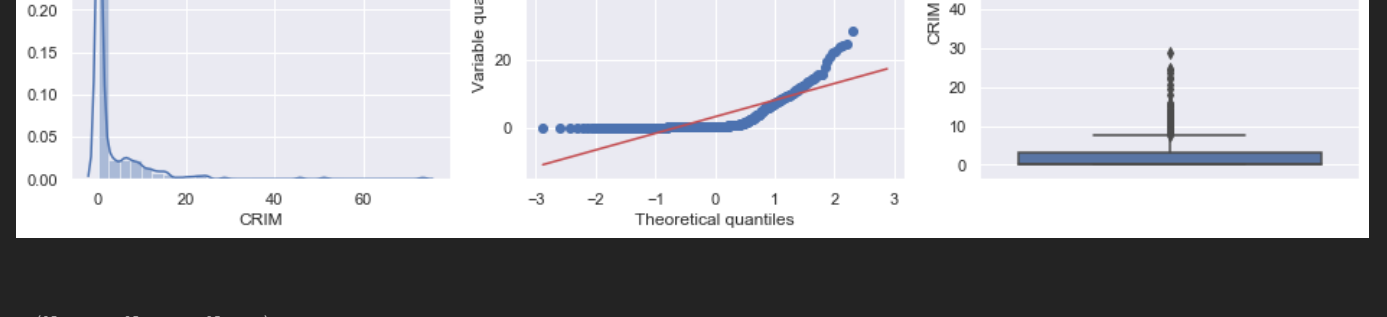
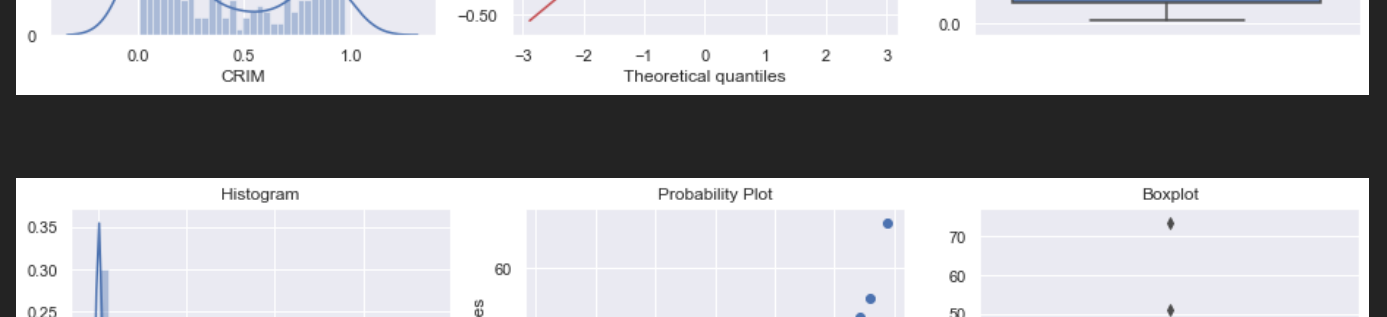
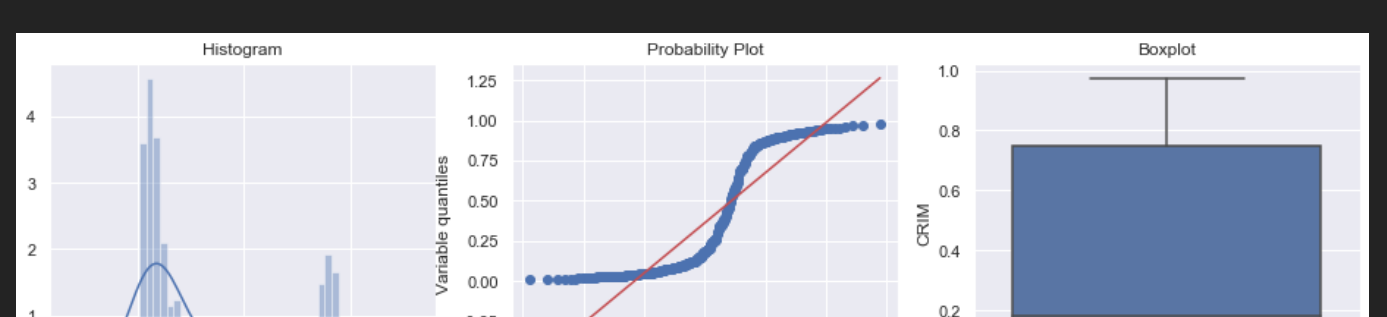
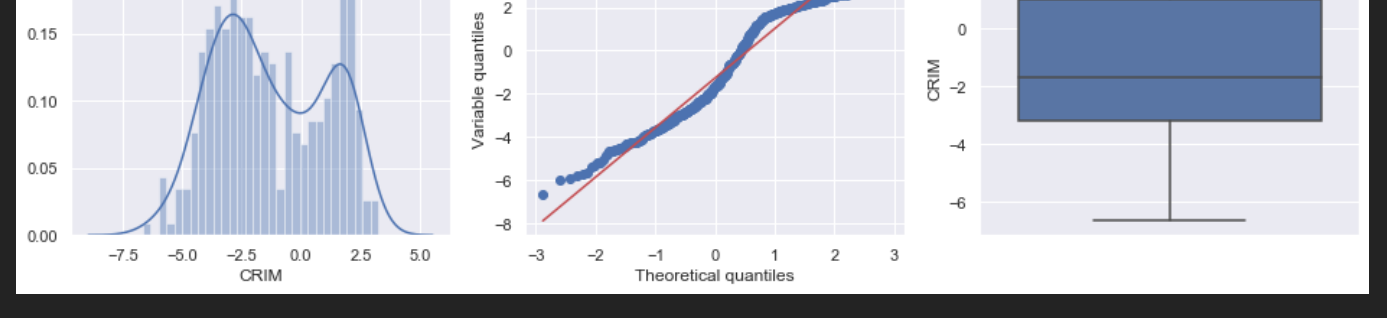
```
# histogram
plt.subplot(1, 3, 1)
sns.distplot(df[variable], bins=30)
plt.title('Histogram')
```

```
# Q-Q plot
plt.subplot(1, 3, 2)
stats.probplot(df[variable], dist='norm', plot=plt)
plt.ylabel('Variable quantiles')
```

```
# boxplot
plt.subplot(1, 3, 3)
sns.boxplot(y=df[variable])
plt.title('Boxplot')
```

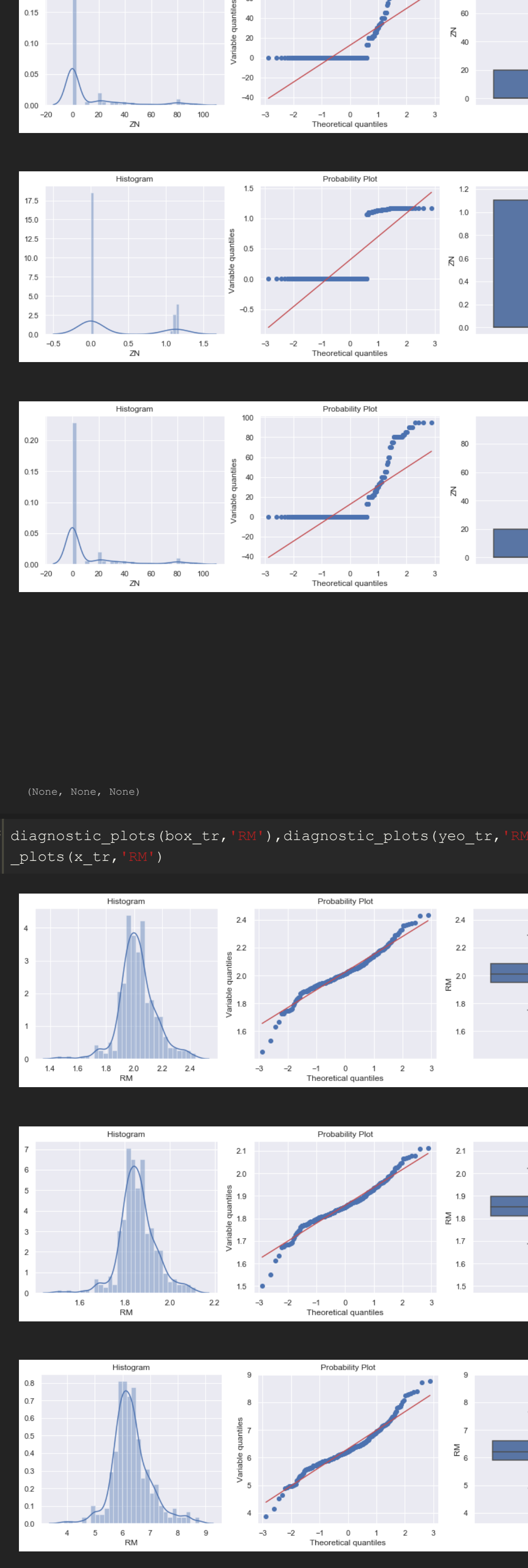
```
plt.show()
```

```
In [19]: # 1. BoxCox 2. YeoJohnson 3. Original Distribution
diagnostic_plots(box_tr, 'CRIM'), diagnostic_plots(yco_tr, 'CRIM'), diagnostic_plots(X_tr, 'CRIM')
```



14 (19)

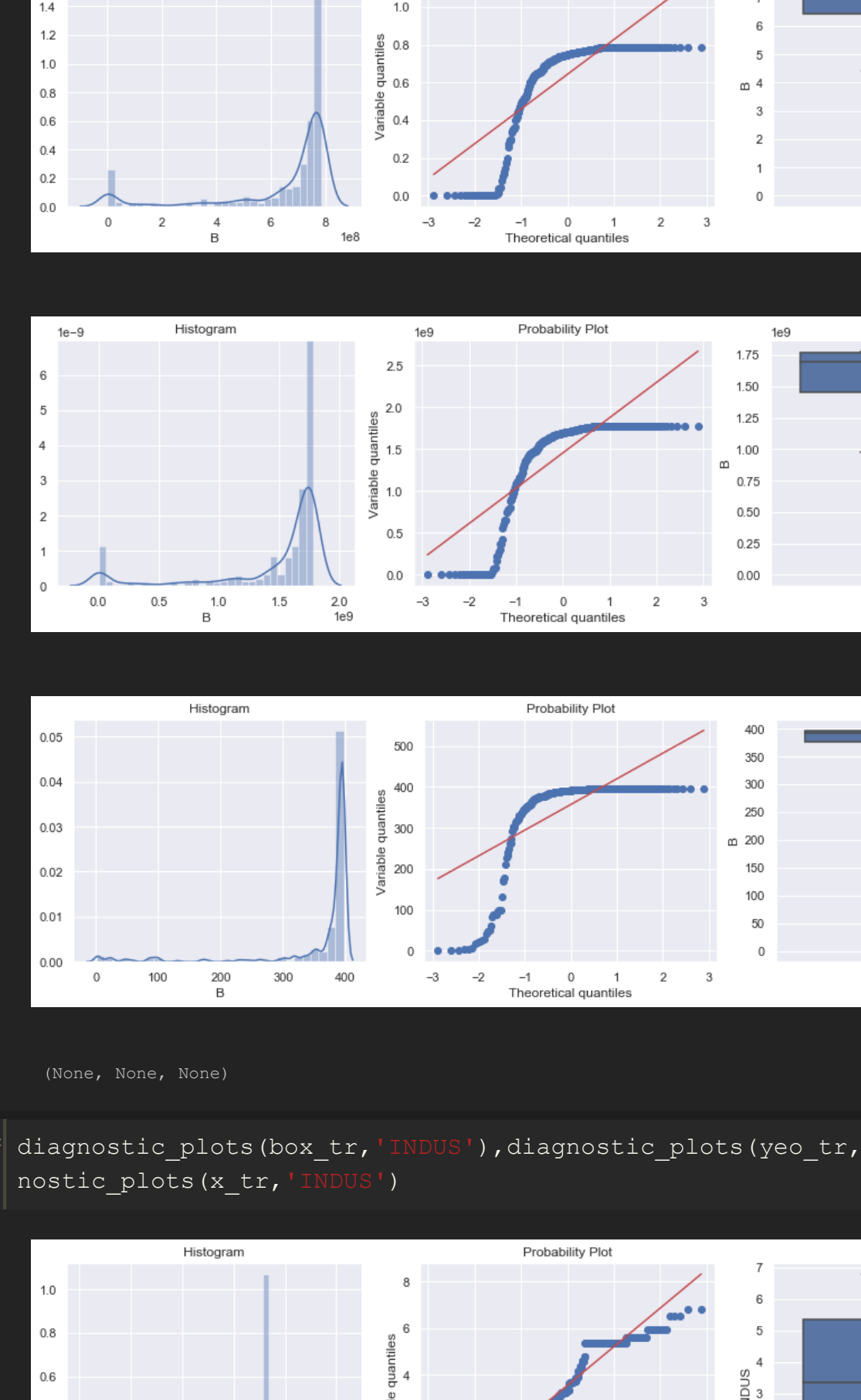
diagnostic_plots(box_tr,'RM'),diagnostic_plots(yeo_tr,'RM'),diagnostic_plots(x_tr,'RM')



(None, None, None)

14 (20)

diagnostic_plots(box_tr,'RM'),diagnostic_plots(yeo_tr,'RM'),diagnostic_plots(x_tr,'RM')



(None, None, None)

14 (21)

diagnostic_plots(box_tr,'B'),diagnostic_plots(yeo_tr,'B'),diagnostic_plots(x_tr,'B')



(None, None, None)

14 (22)

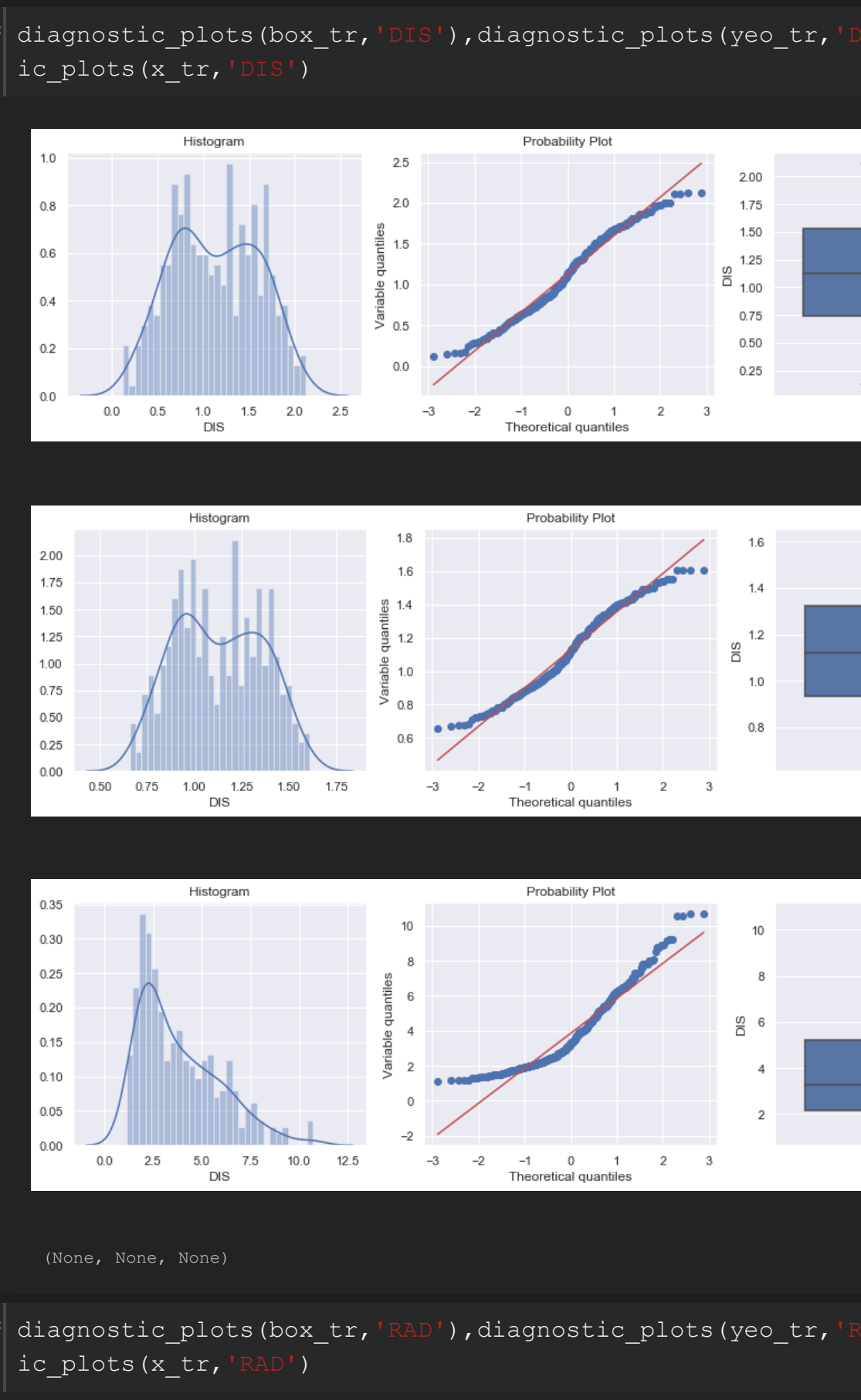
diagnostic_plots(box_tr,'INDUS'),diagnostic_plots(yeo_tr,'INDUS'),diagnostic_plots(x_tr,'INDUS')



(None, None, None)

14 (23)

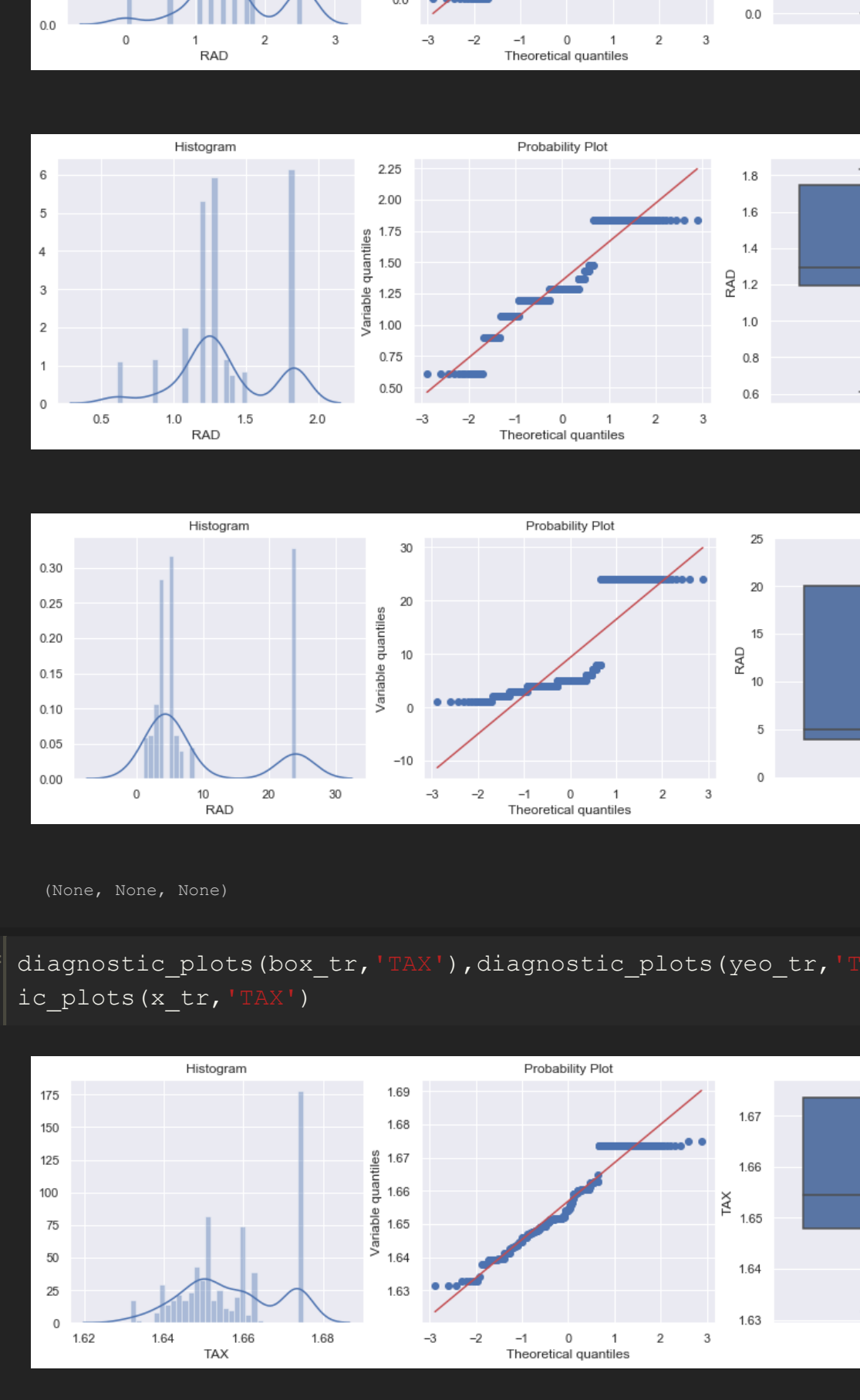
diagnostic_plots(box_tr,'NOX'),diagnostic_plots(yeo_tr,'NOX'),diagnostic_plots(x_tr,'NOX')



(None, None, None)

14 (24)

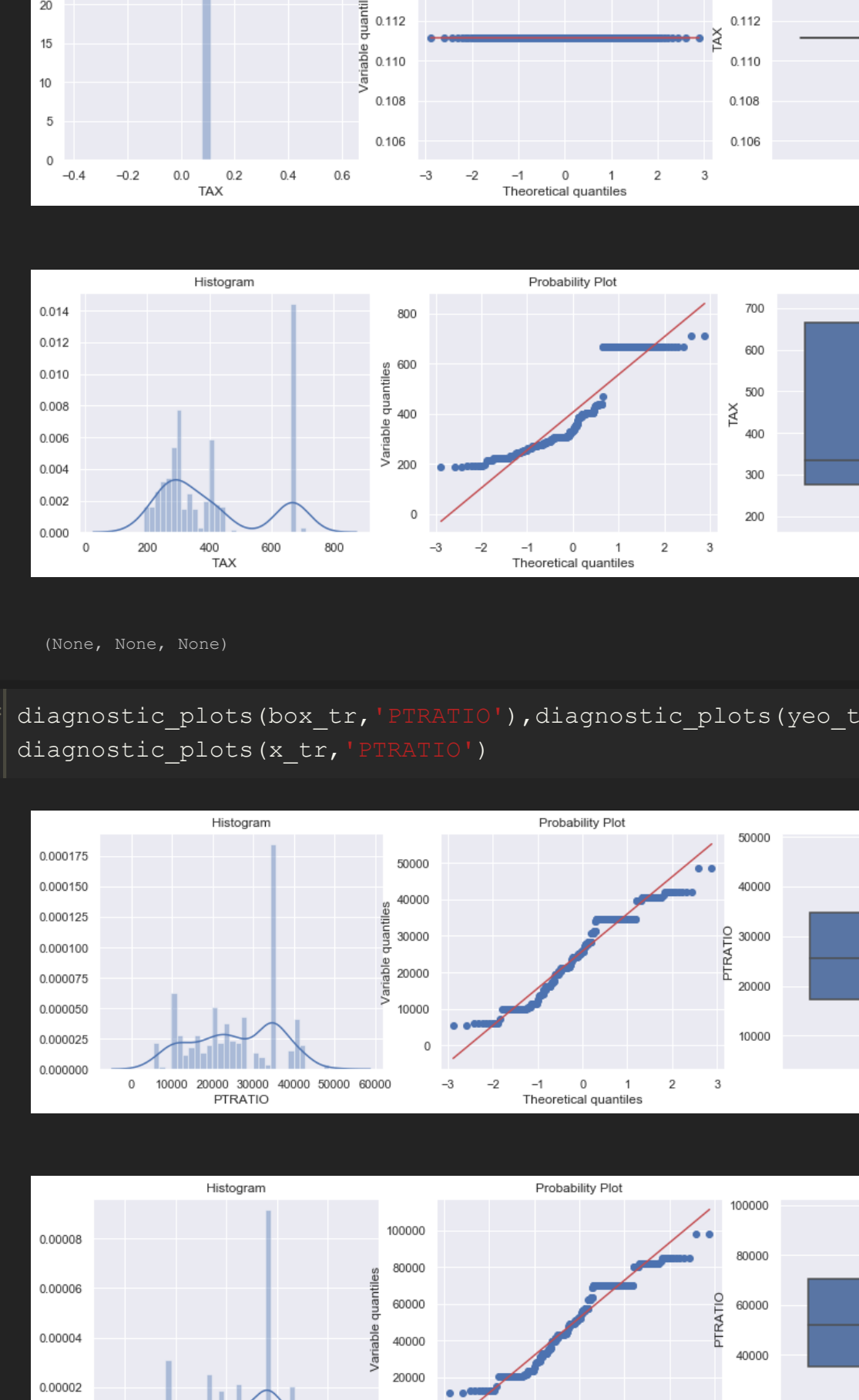
diagnostic_plots(box_tr,'AGE'),diagnostic_plots(yeo_tr,'AGE'),diagnostic_plots(x_tr,'AGE')



(None, None, None)

14 (25)

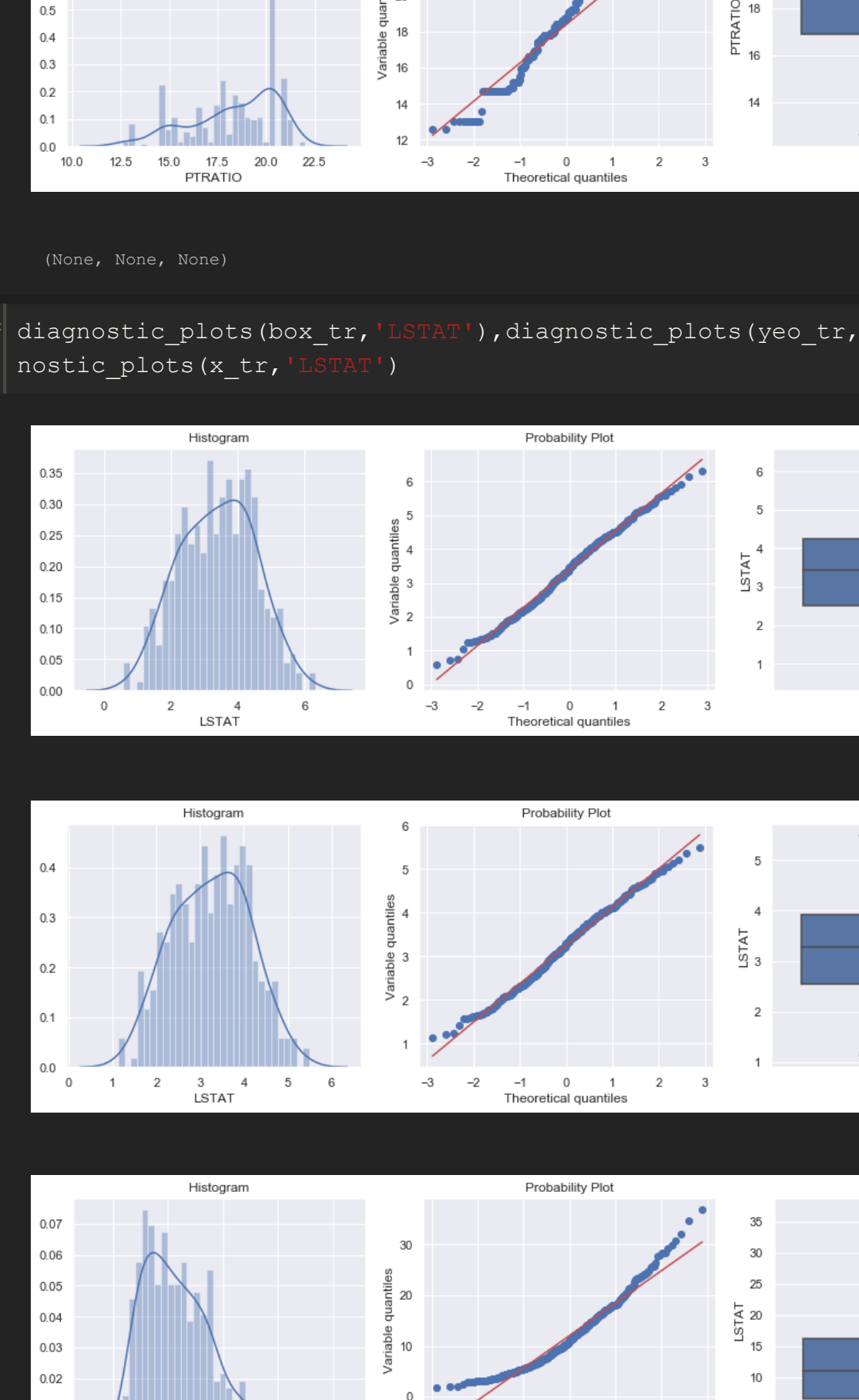
diagnostic_plots(box_tr,'DIS'),diagnostic_plots(yeo_tr,'DIS'),diagnostic_plots(x_tr,'DIS')



(None, None, None)

14 (26)

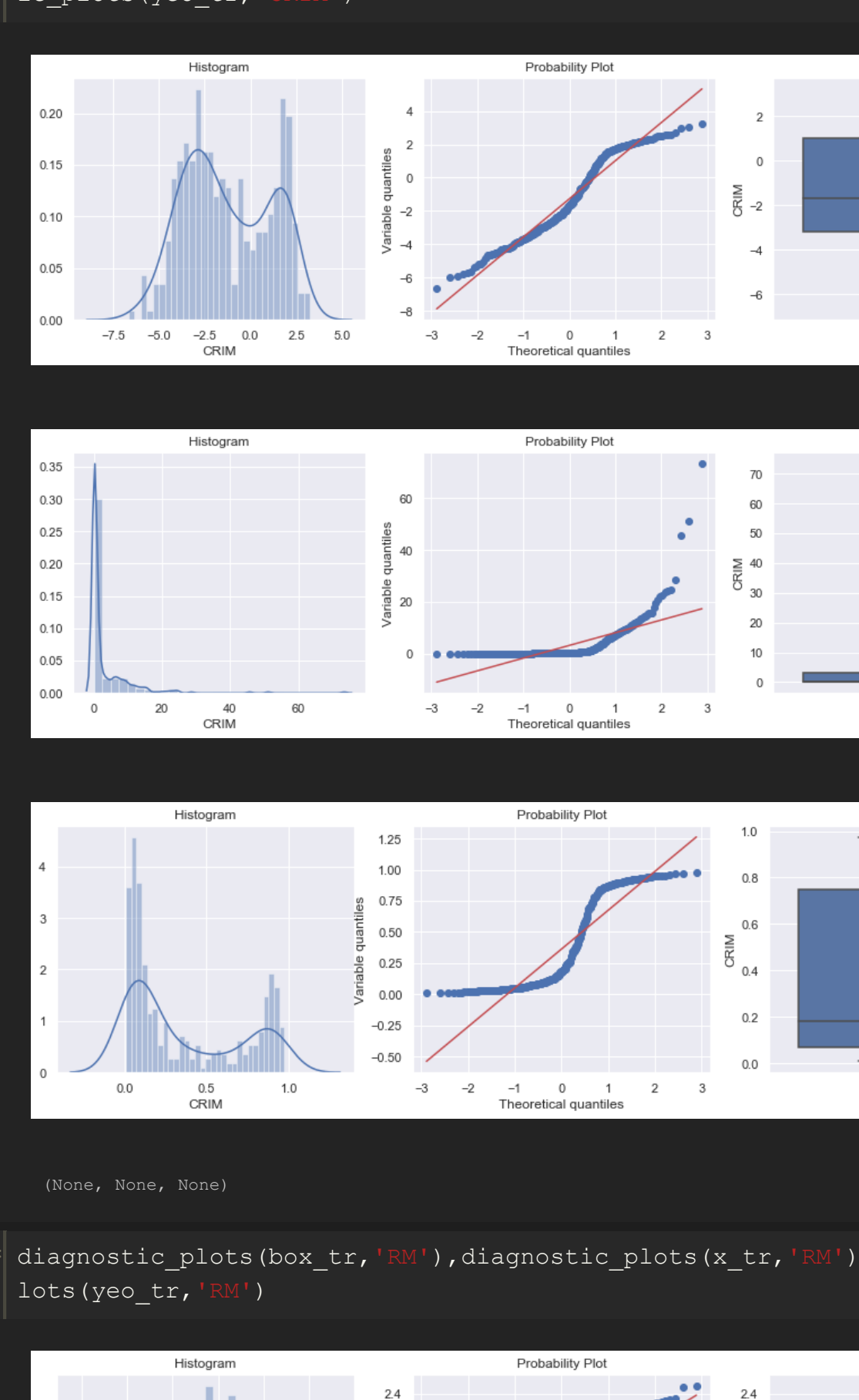
diagnostic_plots(box_tr,'RAD'),diagnostic_plots(yeo_tr,'RAD'),diagnostic_plots(x_tr,'RAD')



(None, None, None)

14 (27)

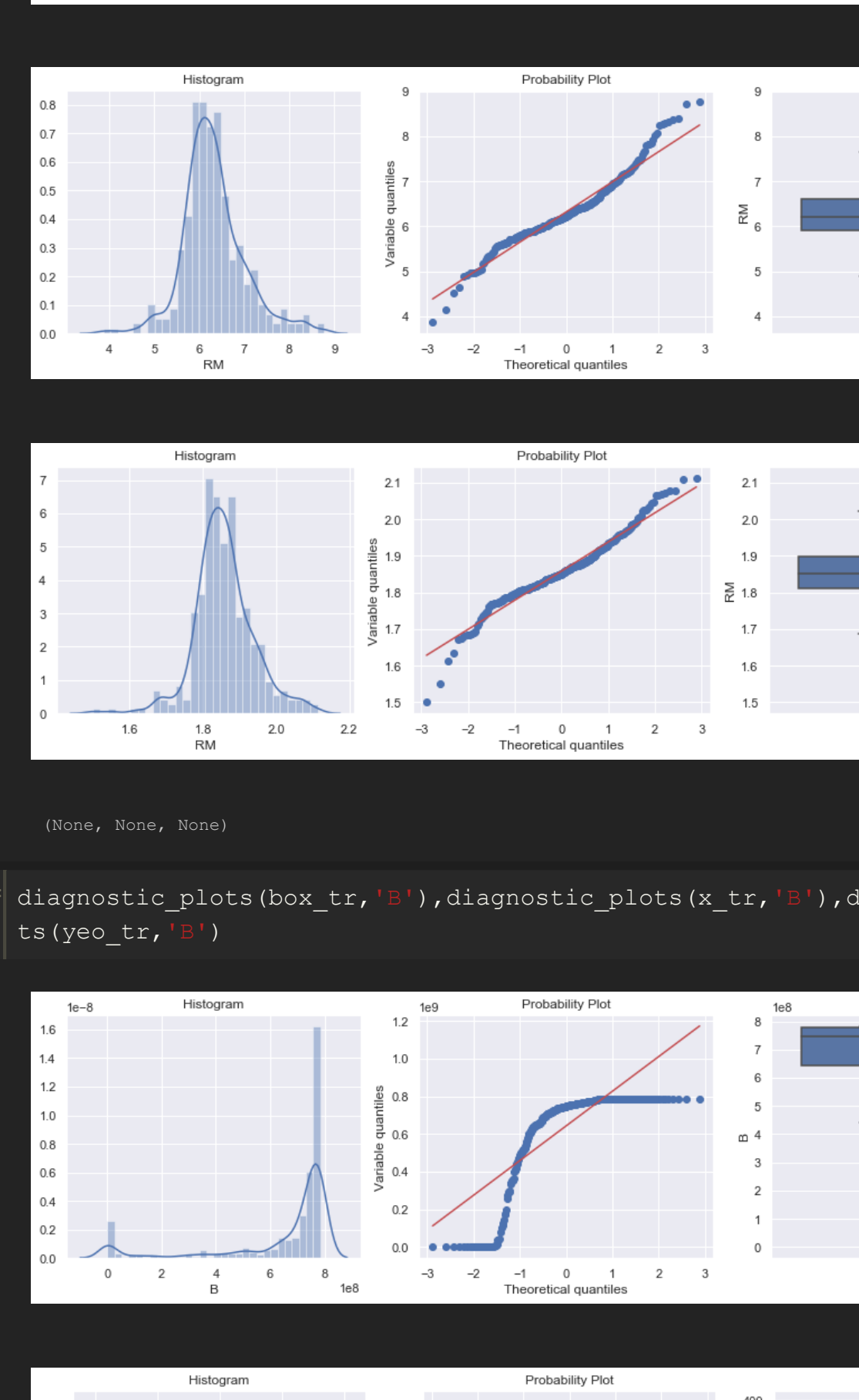
diagnostic_plots(box_tr,'TAX'),diagnostic_plots(yeo_tr,'TAX'),diagnostic_plots(x_tr,'TAX')



(None, None, None)

14 (28)

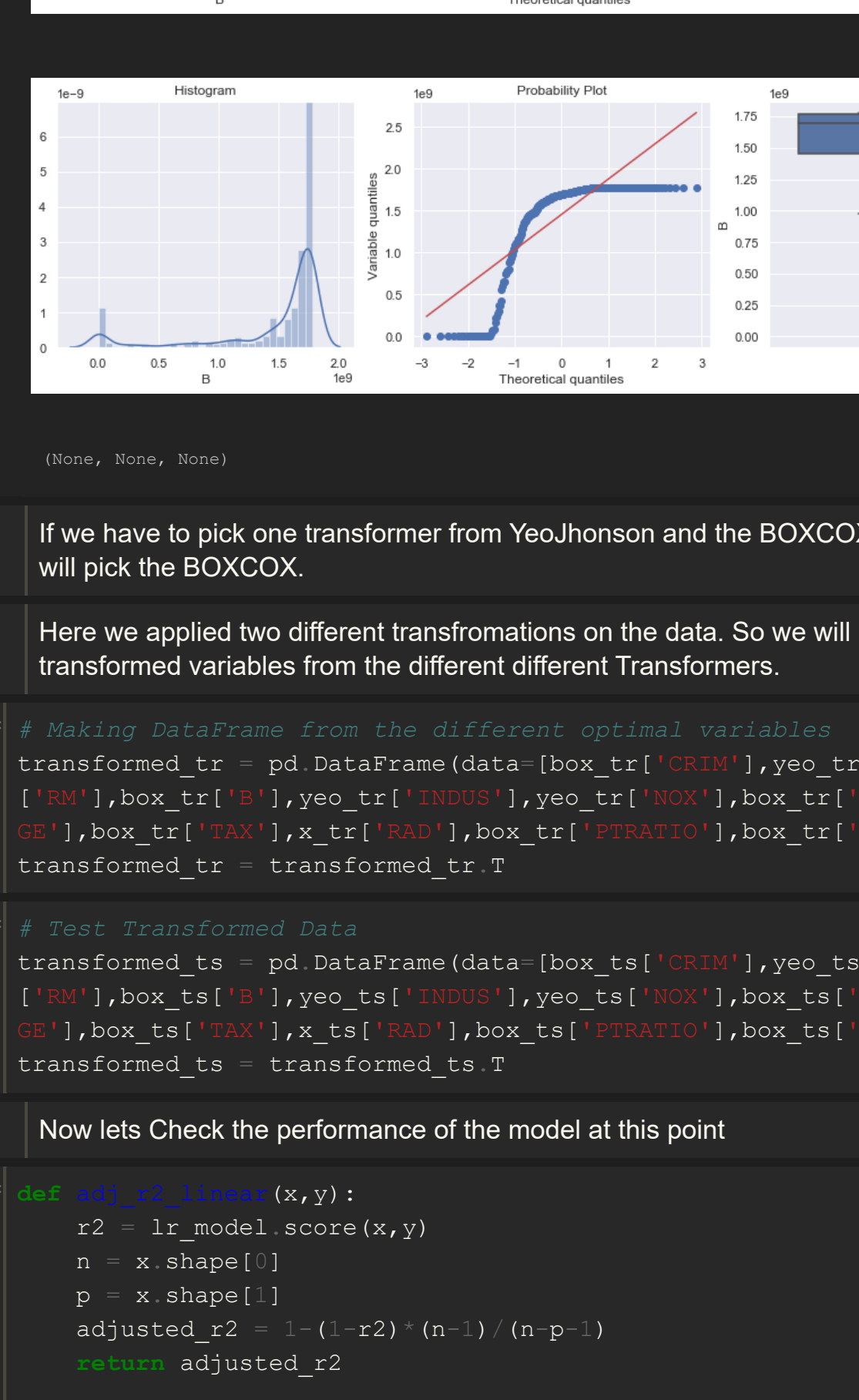
diagnostic_plots(box_tr,'PTRATIO'),diagnostic_plots(yeo_tr,'PTRATIO'),diagnostic_plots(x_tr,'PTRATIO')



(None, None, None)

14 (29)

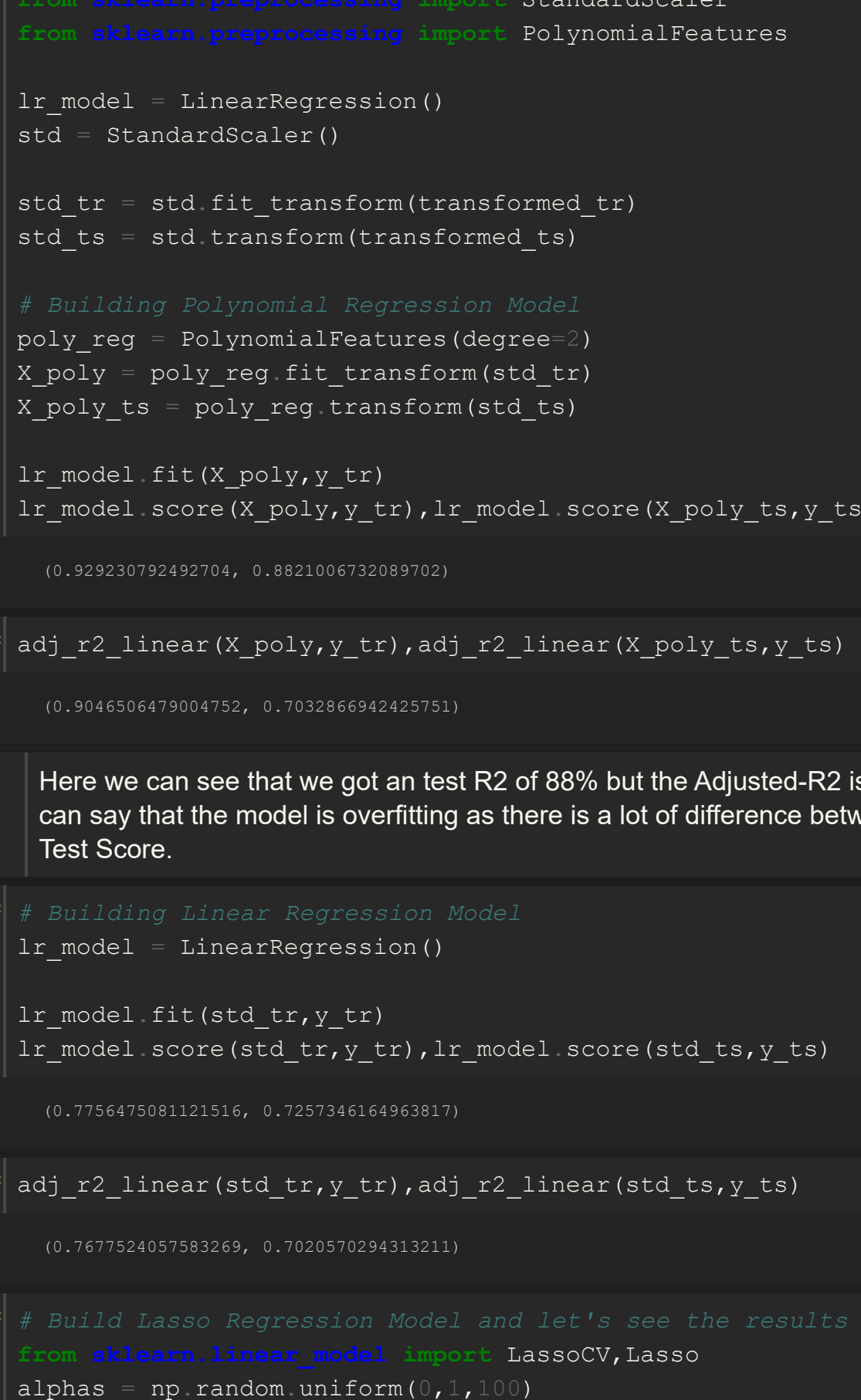
diagnostic_plots(box_tr,'LSTAT'),diagnostic_plots(yeo_tr,'LSTAT'),diagnostic_plots(x_tr,'LSTAT')



(None, None, None)

14 (30)

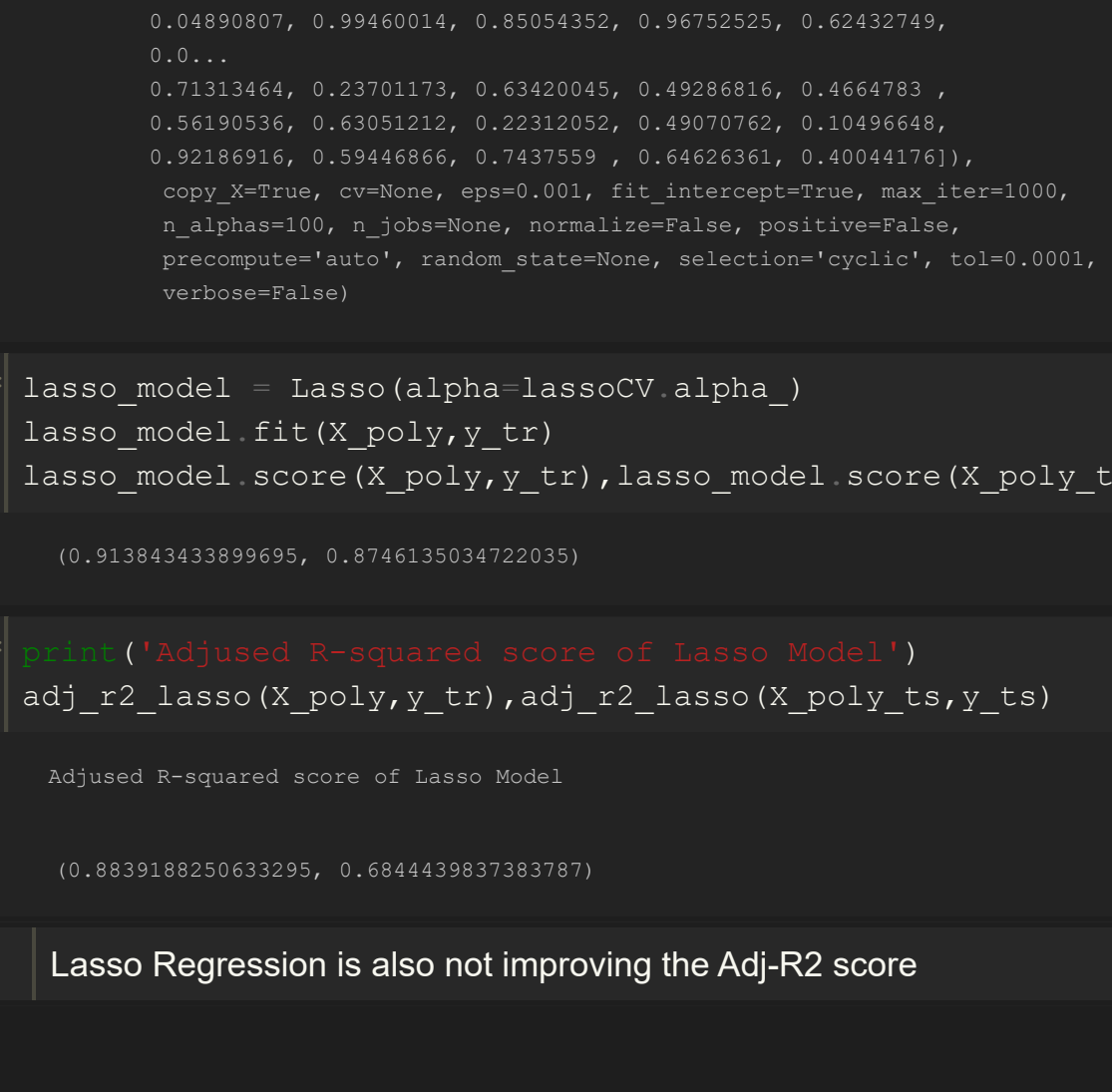
diagnostic_plots(box_tr,'CRIM'),diagnostic_plots(x_tr,'CRIM'),diagnostic_plots(yeo_tr,'CRIM')



(None, None, None)

14 (31)

diagnostic_plots(box_tr,'B'),diagnostic_plots(x_tr,'B'),diagnostic_plots(yeo_tr,'B')



(None, None, None)

If we have to pick one transformer from 'Yeohonson' and the 'BOXCOX' Transformer We will pick the 'BOXCOX'.

Here we applied two different transformations on the data. So we will pick the best transformed variables from the different different Transformers.

14 (32)

Making DataFrame from the different optimal variables

```
transformed_tr = pd.DataFrame(data=[box_tr['CRIM'],yeo_tr['LSTAT'],box_tr['RM'],box_tr['B'],yeo_tr['INDUS'],yeo_tr['NOX'],box_tr['DIS'],x_tr['AGE'],box_tr['TAX'],x_tr['RAD'],box_tr['PTRATIO'],box_tr['LSTAT']])
transformed_tr = transformed_tr.T
```

14 (33)

Test Transformed Data

```
transformed_ts = pd.DataFrame(data=[box_ts['CRIM'],yeo_ts['LSTAT'],box_ts['RM'],box_ts['B'],yeo_ts['INDUS'],yeo_ts['NOX'],box_ts['DIS'],x_ts['AGE'],x_ts['TAX'],x_ts['RAD'],box_ts['PTRATIO'],box_ts['LSTAT']])
transformed_ts = transformed_ts.T
```

Now lets Check the performance of the model at this point

14 (34)

```
def adj_r2_linear(x,y):
    r2 = lr_model.score(x,y)
    n = x.shape[0]
    p = x.shape[1]
    adjusted_r2 = 1-(1-r2)*(n-1)/(n-p-1)
    return adjusted_r2
```

```
def adj_r2_lasso(x,y):
    r2 = lasso_model.score(x,y)
    n = x.shape[0]
    p = x.shape[1]
    adjusted_r2 = 1-(1-r2)*(n-1)/(n-p-1)
    return adjusted_r2
```

14 (35)

```
from sklearn.preprocessing import StandardScaler
from sklearn.preprocessing import PolynomialFeatures
```

```
lr_model = LinearRegression()
std = StandardScaler()
```

```
std_tr = std.fit_transform(transformed_tr)
std_ts = std.transform(transformed_ts)
```

```
# Building Polynomial Regression Model
poly_reg = PolynomialFeatures(degree=2)
X_poly = poly_reg.fit_transform(std_tr)
X_poly_ts = poly_reg.transform(std_ts)
```

```
lr_model.fit(X_poly,y_tr)
lr_model.score(X_poly,y_tr),lr_model.score(X_poly_ts,y_ts)
```

```
(0.929320792492704, 0.8821046732089702)
```

14 (36)

```
adj_r2_linear(X_poly,y_tr),adj_r2_linear(X_poly_ts,y_ts)
```

```
(0.804060479004732, 0.70328664222701)
```

Here we can see that we got an test R2 of 88% but the Adjusted-R2 is quite low. So we can say that the model is overfitting as there is a lot of difference between the Train and Test Score.

14 (37)

```
# Building Linear Regression Model
```

```
lr_model = LinearRegression()
```

```
lr_model.fit(std_tr,y_tr)
```

```
lr_model.score(std_tr,y_tr),lr_model.score(std_ts,y_ts)
```

```
(0.7786475081121516, 0.7257346164963817)
```

14 (38)

```
adj_r2_linear(std_tr,y_tr),adj_r2_linear(std_ts,y_ts)
```

```
(0.7677924057583269, 0.7020570294313211)
```

14 (39)

```
# Build Lasso Regression Model and let's see the results
```

```
from sklearn.linear_model import LassoCV,Lasso
```

```
alphas = np.random.uniform(0,1,100)
```

```
lassoCV = LassoCV(alphas=alphas)
```

```
lassoCV.fit(X_poly,y_tr)
```

```
LassoCV(alphas=array([2.70051868, 0.77219045, 0.49536695, 0.42488879, 0.40054149, 0.0319799, 0.81613055, 0.20722039, 0.50866339, 0.06050549, 0.99400399, 0.73097373, 0.83872439, 0.94585336, 0.01741711, 0.36281821, 0.06700706, 0.67434539, 0.89305136, 0.29370701, 0.8896868, 0.43650769, 0.78814879, 0.57801069, 0.92197273, 0.04390807, 0.99460516, 0.15050492, 0.96752525, 0.5042745, 0.0, ..., 0.7113464, 0.23701173, 0.03220045, 0.49286826, 0.4664703, 0.36190916, 0.60031212, 0.23312052, 0.49070702, 0.10194648, 0.92186916, 0.59446866, 0.7437559, 0.64626361, 0.40644166], copy_X=True, cv=None, eps=0.001, fit_intercept=True, max_iter=1000, n_alpha=100, n_jobs=None, normalize=False, positive=False, precompute='auto', random_state=None, selection='cyclic', tol=0.0001, verbose=False)
```

14 (40)

```
lasso_model = Lasso(alpha=lassoCV.alpha_)
```

```
lasso_model.fit(X_poly,y_tr)
```

```
lasso_model.score(X_poly,y_tr),lasso_model.score(X_poly_ts,y_ts)
```

```
(0.81383433899695, 0.874613534722035)
```

14 (41)

```
print('Adjusted R-squared score of Lasso Model')
```

```
adj_r2_lasso(X_poly,y_tr),adj_r2_lasso(X_poly_ts,y_ts)
```

```
Adjusted R-squared score of Lasso Model
```

```
(0.8839188250632295, 0.684439837383787)
```

Lasso Regression is also not improving the Adj-R2 score

