

Name : Shivam Indrabhan Borse

Roll No : 21119

Subject: Software Laboratory III (DATA SCIENCE)

Assignment No : 04

Problem statement:

Data Analytics I

Create a Linear Regression Model using Python/R to predict home prices using Boston Housing Dataset (<https://www.kaggle.com/c/boston-housing>). The Boston Housing dataset contains information about various houses in Boston through different parameters. There are 506 samples and 14 feature variables in this dataset.

CODE :

```
----- Assignment NO : 04 -----

Data Analytics I Create a Linear Regression Model using Python/R to predict home prices using Boston Housing Dataset (https://www.kaggle.com/c/boston-housing). The Boston Housing dataset contains information about various houses in Boston through different parameters. There are 506 samples and 14 feature variables in this dataset.

>--
import numpy as np
import matplotlib.pyplot as plt
from sklearn.datasets import load_boston
import pandas as pd
import seaborn as sns
%matplotlib inline

[10] ✓ 0.0s Python

+ Code + Markdown

boston_dataset = load_boston()
boston_dataset.keys()

[11] ✓ 0.0s Python

... Output exceeds the size limit. Open the full output data in a text editor
c:\Users\ASUS\AppData\Local\Programs\Python\Python39\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated.

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source::
```

```
boston = pd.DataFrame(boston_dataset.data, columns=boston_dataset.feature_names)
boston.head()
```

[12] ✓ 0.0s

```
CRIM  ZN  INDUS  CHAS  NOX  RM  AGE  DIS  RAD  TAX  PTRATIO  B  LSTAT
0  0.00632  18.0  2.31  0.0  0.538  6.575  65.2  4.0900  1.0  296.0  15.3  396.90  4.98
1  0.02731  0.0  7.07  0.0  0.469  6.421  78.9  4.9671  2.0  242.0  17.8  396.90  9.14
2  0.02729  0.0  7.07  0.0  0.469  7.185  61.1  4.9671  2.0  242.0  17.8  392.83  4.03
3  0.03237  0.0  2.18  0.0  0.458  6.998  45.8  6.0622  3.0  222.0  18.7  394.63  2.94
4  0.06905  0.0  2.18  0.0  0.458  7.147  54.2  6.0622  3.0  222.0  18.7  396.90  5.33
```

```
boston.isnull().sum()
```

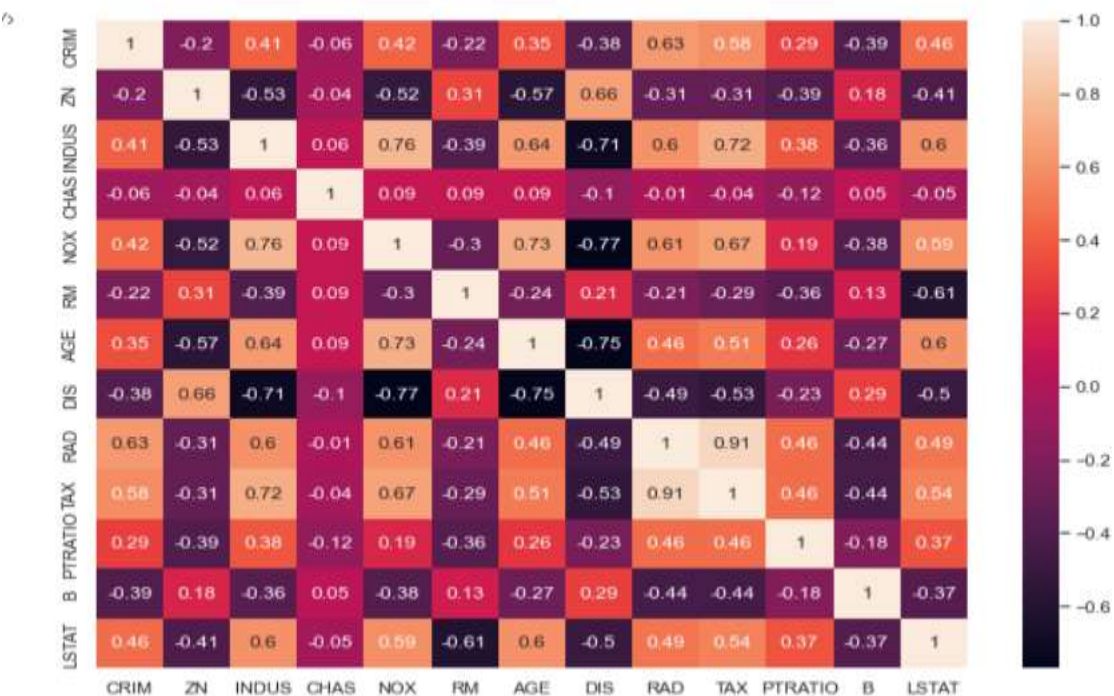
[13] ✓ 0.0s

```
CRIM  0
ZN  0
INDUS  0
CHAS  0
NOX  0
RM  0
AGE  0
DIS  0
RAD  0
TAX  0
PTRATIO  0
B  0
LSTAT  0
dtype: int64
```

```
correlation_matrix = boston.corr().round(2)
sns.set(rc={'figure.figsize':(11.7,8.27)})
sns.heatmap(data=correlation_matrix, annot=True)
```

[14] ✓ 2.8s

<AxesSubplot:>



```
X = pd.DataFrame(np.c_[boston['LSTAT'], boston['RM']], columns = ['LSTAT', 'RM'])
Y = boston['TAX']
```

[15] ✓ 0.0s

```
from sklearn.model_selection import train_test_split
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, random_state=5)
print(X_train.shape)
print(X_test.shape)
print(Y_train.shape)
print(Y_test.shape)
```

[16] ✓ 0.1s

```
... (404, 2)
... (102, 2)
... (404,)
... (102,)
```

```
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

lin_model = LinearRegression()
lin_model.fit(X_train, Y_train)
```

[17] ✓ 0.1s

```
... LinearRegression()
```

```
print("The model performance for training set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
print("\n")

# model evaluation for testing set

y_test_predict = lin_model.predict(X_test)
# root mean square error of the model
rmse = (np.sqrt(mean_squared_error(Y_test, y_test_predict)))

# r-squared score of the model
r2 = r2_score(Y_test, y_test_predict)

print("The model performance for testing set")
print("-----")
print('RMSE is {}'.format(rmse))
print('R2 score is {}'.format(r2))
```

✓ 0.1s

```
The model performance for training set
-----
RMSE is 141.50944983061052
R2 score is 0.2915976508661784

The model performance for testing set
-----
RMSE is 138.96850882629893
R2 score is 0.32167241065376195
```



--> SHIVAM BORSE