

# Boston\_linear\_reg

May 11, 2022

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.

```
[4]: # Importing Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt

# Importing Data
from sklearn.datasets import load_boston
boston = load_boston()
```

```
[5]: boston.data.shape
```

```
[5]: (506, 13)
```

```
[6]: boston.feature_names
```

```
[6]: array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD',
        'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U7')
```

```
[7]: data = pd.DataFrame(boston.data)
data.columns = boston.feature_names

data.head(10)
```

```
[7]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	\
0	0.00632	18.0	2.31	0.0	0.538	6.575	65.2	4.0900	1.0	296.0	
1	0.02731	0.0	7.07	0.0	0.469	6.421	78.9	4.9671	2.0	242.0	
2	0.02729	0.0	7.07	0.0	0.469	7.185	61.1	4.9671	2.0	242.0	
3	0.03237	0.0	2.18	0.0	0.458	6.998	45.8	6.0622	3.0	222.0	
4	0.06905	0.0	2.18	0.0	0.458	7.147	54.2	6.0622	3.0	222.0	
5	0.02985	0.0	2.18	0.0	0.458	6.430	58.7	6.0622	3.0	222.0	
6	0.08829	12.5	7.87	0.0	0.524	6.012	66.6	5.5605	5.0	311.0	
7	0.14455	12.5	7.87	0.0	0.524	6.172	96.1	5.9505	5.0	311.0	
8	0.21124	12.5	7.87	0.0	0.524	5.631	100.0	6.0821	5.0	311.0	

```
9 0.17004 12.5 7.87 0.0 0.524 6.004 85.9 6.5921 5.0 311.0
```

```

PTRATIO      B  LSTAT
0    15.3 396.90  4.98
1    17.8 396.90  9.14
2    17.8 392.83  4.03
3    18.7 394.63  2.94
4    18.7 396.90  5.33
5    18.7 394.12  5.21
6    15.2 395.60 12.43
7    15.2 396.90 19.15
8    15.2 386.63 29.93
9    15.2 386.71 17.10

```

```
[8]: data['Price'] = boston.target
data.head(10)
```

```

[8]:      CRIM      ZN  INDUS  CHAS    NOX     RM   AGE     DIS  RAD    TAX  \
0  0.00632  18.0    2.31   0.0  0.538  6.575  65.2  4.0900  1.0  296.0
1  0.02731   0.0    7.07   0.0  0.469  6.421  78.9  4.9671  2.0  242.0
2  0.02729   0.0    7.07   0.0  0.469  7.185  61.1  4.9671  2.0  242.0
3  0.03237   0.0    2.18   0.0  0.458  6.998  45.8  6.0622  3.0  222.0
4  0.06905   0.0    2.18   0.0  0.458  7.147  54.2  6.0622  3.0  222.0
5  0.02985   0.0    2.18   0.0  0.458  6.430  58.7  6.0622  3.0  222.0
6  0.08829  12.5    7.87   0.0  0.524  6.012  66.6  5.5605  5.0  311.0
7  0.14455  12.5    7.87   0.0  0.524  6.172  96.1  5.9505  5.0  311.0
8  0.21124  12.5    7.87   0.0  0.524  5.631 100.0  6.0821  5.0  311.0
9  0.17004  12.5    7.87   0.0  0.524  6.004  85.9  6.5921  5.0  311.0

```

```

PTRATIO      B  LSTAT  Price
0    15.3 396.90  4.98   24.0
1    17.8 396.90  9.14   21.6
2    17.8 392.83  4.03   34.7
3    18.7 394.63  2.94   33.4
4    18.7 396.90  5.33   36.2
5    18.7 394.12  5.21   28.7
6    15.2 395.60 12.43   22.9
7    15.2 396.90 19.15   27.1
8    15.2 386.63 29.93   16.5
9    15.2 386.71 17.10   18.9

```

```
[9]: data.describe()
```

```

[9]:      CRIM      ZN      INDUS      CHAS      NOX      RM  \
count  506.000000  506.000000  506.000000  506.000000  506.000000  506.000000
mean     3.613524  11.363636  11.136779    0.069170    0.554695    6.284634
std     8.601545  23.322453    6.860353    0.253994    0.115878    0.702617

```

min	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000
25%	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500
50%	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500
75%	3.677083	12.500000	18.100000	0.000000	0.624000	6.623500
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000

	AGE	DIS	RAD	TAX	PTRATIO	B \
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	68.574901	3.795043	9.549407	408.237154	18.455534	356.674032
std	28.148861	2.105710	8.707259	168.537116	2.164946	91.294864
min	2.900000	1.129600	1.000000	187.000000	12.600000	0.320000
25%	45.025000	2.100175	4.000000	279.000000	17.400000	375.377500
50%	77.500000	3.207450	5.000000	330.000000	19.050000	391.440000
75%	94.075000	5.188425	24.000000	666.000000	20.200000	396.225000
max	100.000000	12.126500	24.000000	711.000000	22.000000	396.900000

	LSTAT	Price
count	506.000000	506.000000
mean	12.653063	22.532806
std	7.141062	9.197104
min	1.730000	5.000000
25%	6.950000	17.025000
50%	11.360000	21.200000
75%	16.955000	25.000000
max	37.970000	50.000000

```
[10]: data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 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  Price       506 non-null    float64
```

dtypes: float64(14)  
memory usage: 55.5 KB

```
[11]: # Input Data
x = boston.data
# Output Data
y = boston.target
# splitting data to training and testing dataset.

from sklearn.model_selection import train_test_split

xtrain, xtest, ytrain, ytest = train_test_split(x, y, test_size=0.
↪2, random_state = 0)

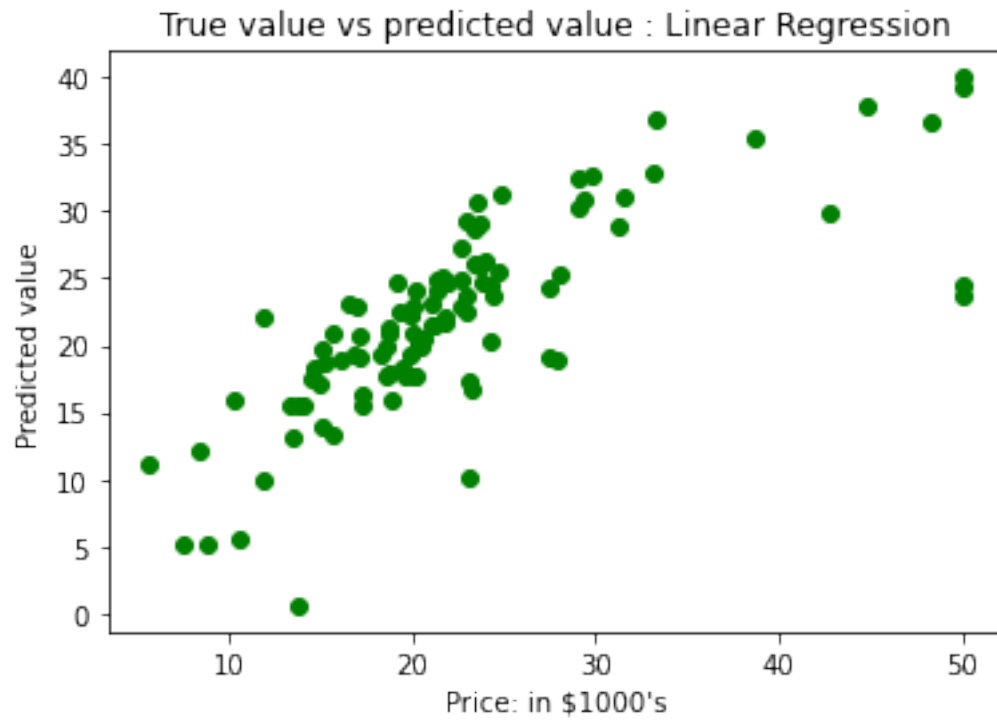
print("xtrain shape : ", xtrain.shape)
print("xtest shape : ", xtest.shape)
print("ytrain shape : ", ytrain.shape)
print("ytest shape : ", ytest.shape)
```

xtrain shape : (404, 13)  
xtest shape : (102, 13)  
ytrain shape : (404,)  
ytest shape : (102,)

```
[12]: # Fitting Multi Linear regression model to training model
from sklearn.linear_model import LinearRegression
regressor = LinearRegression()
regressor.fit(xtrain, ytrain)

# predicting the test set results
y_pred = regressor.predict(xtest)
```

```
[13]: # Plotting Scatter graph to show the prediction
# results - 'ytrue' value vs 'y_pred' value
plt.scatter(ytest, y_pred, c = 'green')
plt.xlabel("Price: in $1000's")
plt.ylabel("Predicted value")
plt.title("True value vs predicted value : Linear Regression")
plt.show()
```



```
[14]: from sklearn.metrics import mean_squared_error

mse = mean_squared_error(ytest, y_pred)
print("Mean Square Error : ", mse)
```

Mean Square Error : 33.44897999767661