

KUMAR GAURAV 20122065**ML FINAL EXAM ON 02 JUNE 2021**

- to build a linear regression model and explain each components Implement using python*

```
import matplotlib.pyplot as plt
import numpy as np
from sklearn import datasets, linear_model, metrics
```

load the boston dataset

```
boston = datasets.load_boston(return_X_y=False)
```

```
boston
```

```
[2.7290e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9283e+02,
 4.0300e+00],
...,
[6.0760e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
 5.6400e+00],
[1.0959e-01, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9345e+02,
 6.4800e+00],
[4.7410e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
 7.8800e+00]]),
'feature_names': array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS',
  'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U7'),
'filename': '/usr/local/lib/python3.7/dist-packages/sklearn/datasets/data/boston_
'target': array([24. , 21.6, 34.7, 33.4, 36.2, 28.7, 22.9, 27.1, 16.5, 18.9, 15.
 18.9, 21.7, 20.4, 18.2, 19.9, 23.1, 17.5, 20.2, 18.2, 13.6, 19.6,
 15.2, 14.5, 15.6, 13.9, 16.6, 14.8, 18.4, 21. , 12.7, 14.5, 13.2,
 13.1, 13.5, 18.9, 20. , 21. , 24.7, 30.8, 34.9, 26.6, 25.3, 24.7,
 21.2, 19.3, 20. , 16.6, 14.4, 19.4, 19.7, 20.5, 25. , 23.4, 18.9,
 35.4, 24.7, 31.6, 23.3, 19.6, 18.7, 16. , 22.2, 25. , 33. , 23.5,
 19.4, 22. , 17.4, 20.9, 24.2, 21.7, 22.8, 23.4, 24.1, 21.4, 20. ,
 20.8, 21.2, 20.3, 28. , 23.9, 24.8, 22.9, 23.9, 26.6, 22.5, 22.2,
 23.6, 28.7, 22.6, 22. , 22.9, 25. , 20.6, 28.4, 21.4, 38.7, 43.8,
 33.2, 27.5, 26.5, 18.6, 19.3, 20.1, 19.5, 19.5, 20.4, 19.8, 19.4,
 21.7, 22.8, 18.8, 18.7, 18.5, 18.3, 21.2, 19.2, 20.4, 19.3, 22. ,
 20.3, 20.5, 17.3, 18.8, 21.4, 15.7, 16.2, 18. , 14.3, 19.2, 19.6,
 23. , 18.4, 15.6, 18.1, 17.4, 17.1, 13.3, 17.8, 14. , 14.4, 13.4,
 15.6, 11.8, 13.8, 15.6, 14.6, 17.8, 15.4, 21.5, 19.6, 15.3, 19.4,
 17. , 15.6, 13.1, 41.3, 24.3, 23.3, 27. , 50. , 50. , 50. , 22.7,
 25. , 50. , 23.8, 23.8, 22.3, 17.4, 19.1, 23.1, 23.6, 22.6, 29.4,
 23.2, 24.6, 29.9, 37.2, 39.8, 36.2, 37.9, 32.5, 26.4, 29.6, 50. ,
 32. , 29.8, 34.9, 37. , 30.5, 36.4, 31.1, 29.1, 50. , 33.3, 30.3,
 34.6, 34.9, 32.9, 24.1, 42.3, 48.5, 50. , 22.6, 24.4, 22.5, 24.4,
 20. , 21.7, 19.3, 22.4, 28.1, 23.7, 25. , 23.3, 28.7, 21.5, 23. ,
 26.7, 21.7, 27.5, 30.1, 44.8, 50. , 37.6, 31.6, 46.7, 31.5, 24.3,
 31.7, 41.7, 48.3, 29. , 24. , 25.1, 31.5, 23.7, 23.3, 22. , 20.1,
 22.2, 23.7, 17.6, 18.5, 24.3, 20.5, 24.5, 26.2, 24.4, 24.8, 29.6,
 42.8, 21.9, 20.9, 44. , 50. , 36. , 30.1, 33.8, 43.1, 48.8, 31. ,
```

```

36.5, 22.8, 30.7, 50. , 43.5, 20.7, 21.1, 25.2, 24.4, 35.2, 32.4,
32. , 33.2, 33.1, 29.1, 35.1, 45.4, 35.4, 46. , 50. , 32.2, 22. ,
20.1, 23.2, 22.3, 24.8, 28.5, 37.3, 27.9, 23.9, 21.7, 28.6, 27.1,
20.3, 22.5, 29. , 24.8, 22. , 26.4, 33.1, 36.1, 28.4, 33.4, 28.2,
22.8, 20.3, 16.1, 22.1, 19.4, 21.6, 23.8, 16.2, 17.8, 19.8, 23.1,
21. , 23.8, 23.1, 20.4, 18.5, 25. , 24.6, 23. , 22.2, 19.3, 22.6,
19.8, 17.1, 19.4, 22.2, 20.7, 21.1, 19.5, 18.5, 20.6, 19. , 18.7,
32.7, 16.5, 23.9, 31.2, 17.5, 17.2, 23.1, 24.5, 26.6, 22.9, 24.1,
18.6, 30.1, 18.2, 20.6, 17.8, 21.7, 22.7, 22.6, 25. , 19.9, 20.8,
16.8, 21.9, 27.5, 21.9, 23.1, 50. , 50. , 50. , 50. , 50. , 13.8,
13.8, 15. , 13.9, 13.3, 13.1, 10.2, 10.4, 10.9, 11.3, 12.3, 8.8,
7.2, 10.5, 7.4, 10.2, 11.5, 15.1, 23.2, 9.7, 13.8, 12.7, 13.1,
12.5, 8.5, 5. , 6.3, 5.6, 7.2, 12.1, 8.3, 8.5, 5. , 11.9,
27.9, 17.2, 27.5, 15. , 17.2, 17.9, 16.3, 7. , 7.2, 7.5, 10.4,
8.8, 8.4, 16.7, 14.2, 20.8, 13.4, 11.7, 8.3, 10.2, 10.9, 11. ,
9.5, 14.5, 14.1, 16.1, 14.3, 11.7, 13.4, 9.6, 8.7, 8.4, 12.8,
10.5, 17.1, 18.4, 15.4, 10.8, 11.8, 14.9, 12.6, 14.1, 13. , 13.4,
15.2, 16.1, 17.8, 14.9, 14.1, 12.7, 13.5, 14.9, 20. , 16.4, 17.7,
19.5, 20.2, 21.4, 19.9, 19. , 19.1, 19.1, 20.1, 19.9, 19.6, 23.2,
29.8, 13.8, 13.3, 16.7, 12. , 14.6, 21.4, 23. , 23.7, 25. , 21.8,
20.6, 21.2, 19.1, 20.6, 15.2, 7. , 8.1, 13.6, 20.1, 21.8, 24.5,
23.1, 19.7, 18.3, 21.2, 17.5, 16.8, 22.4, 20.6, 23.9, 22. , 11.9]})}

```

defining feature matrix(X) and response vector(y)

```

X = boston.data
y = boston.target

```

There are 506 rows and total 14 columns where, In x there are 13 columns for feature , In y there are one columns for target

There are 4 keys in the bunch ['data', 'target', 'feature_names', 'DESCR'] as mentioned above. The data has 506 rows and 13 feature variable.

```
X.shape
```

```
(506, 13)
```

```
y.shape
```

```
(506,)
```

splitting X and y into training and testing sets

```
print(boston.DESCR)
```

```
.. _boston_dataset:
```

```
Boston house prices dataset
```

```
-----
```

****Data Set Characteristics:****

:Number of Instances: 506

:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute

:Attribute Information (in order):

- CRIM per capita crime rate by town
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS proportion of non-retail business acres per town
- CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
- NOX nitric oxides concentration (parts per 10 million)
- RM average number of rooms per dwelling
- AGE proportion of owner-occupied units built prior to 1940
- DIS weighted distances to five Boston employment centres
- RAD index of accessibility to radial highways
- TAX full-value property-tax rate per \$10,000
- PTRATIO pupil-teacher ratio by town
- B $1000(Bk - 0.63)^2$ where Bk is the proportion of blacks by town
- LSTAT % lower status of the population
- MEDV Median value of owner-occupied homes in \$1000's

:Missing Attribute Values: None

:Creator: Harrison, D. and Rubinfeld, D.L.

This is a copy of UCI ML housing dataset.

<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>

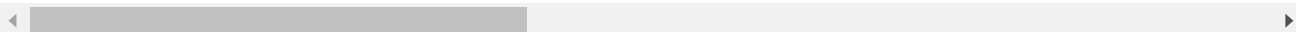
This dataset was taken from the StatLib library which is maintained at Carnegie Mellon

The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic prices and the demand for clean air', J. Environ. Economics & Management, vol.5, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression diagnostics ...', Wiley, 1980. N.B. Various transformations are used in the table on pages 244-261 of the latter.

The Boston house-price data has been used in many machine learning papers that address problems.

.. topic:: References

- Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and
- Quinlan,R. (1993). Combining Instance-Based and Model-Based Learning. In Proceed



```
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.4,
                                                    random_state=1)
```

create linear regression object

```
reg = linear_model.LinearRegression()
```

train the model using the training sets

```
reg.fit(X_train, y_train)
```

regression coefficients

```
print('Coefficients: ', reg.coef_)
```

```
Coefficients: [-8.95714048e-02  6.73132853e-02  5.04649248e-02  2.18579583e+00
-1.72053975e+01  3.63606995e+00  2.05579939e-03 -1.36602886e+00
 2.89576718e-01 -1.22700072e-02 -8.34881849e-01  9.40360790e-03
-5.04008320e-01]
```

variance score: 1 means perfect prediction

```
print('Variance score: {}'.format(reg.score(X_test, y_test)))
```

```
Variance score: 0.7209056672661777
```

setting plot style

```
plt.style.use('fivethirtyeight')
```

To create our model, we must “learn” or estimate the values of regression coefficients b_0 and b_1 . And once we’ve estimated these coefficients, we can use the model to predict responses! In this article, we are going to use the principle of Least Squares . Now consider: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i = h(x_i) + \epsilon_i \rightarrow \epsilon_i = y_i - h(x_i)$
Here, ϵ_i is residual error in i th observation. So, our aim is to minimize the total residual error

```
def estimate_coef(x, y):
    # number of observations/points
    n = np.size(x)

    # mean of x and y vector
    m_x = np.mean(x)
    m_y = np.mean(y)

    # calculating cross-deviation and deviation about x
    SS_xy = np.sum(y*x) - n*m_y*m_x
    SS_xx = np.sum(x*x) - n*m_x*m_x
```

```
# calculating regression coefficients
b_1 = SS_xy / SS_xx
b_0 = m_y - b_1*m_x

return (b_0, b_1)
```

```
def plot_regression_line(x, y, b):
    # plotting the actual points as scatter plot
    plt.scatter(x, y, color = "m",
                marker = "o", s = 30)

    # predicted response vector
    y_pred = b[0] + b[1]*x

    # plotting the regression line
    plt.plot(x, y_pred, color = "g")

    # putting labels
    plt.xlabel('x')
    plt.ylabel('y')

    # function to show plot
    plt.show()
```

Estimated coefficients:

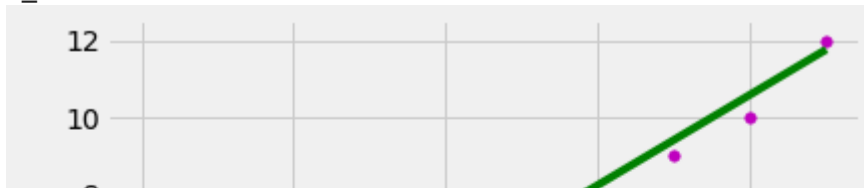
```
def main():
    # observations / data
    x = np.array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])
    y = np.array([1, 3, 2, 5, 7, 8, 8, 9, 10, 12])

    # estimating coefficients
    b = estimate_coef(x, y)
    print("Estimated coefficients:\nb_0 = {} \
        \nb_1 = {}".format(b[0], b[1]))

    # plotting regression line
    plot_regression_line(x, y, b)

if __name__ == "__main__":
    main()
```

Estimated coefficients:
 $b_0 = 1.2363636363636363$
 $b_1 = 1.1696969696969697$



Multiple linear regression Multiple linear regression attempts to model the relationship between two or more features and a response by fitting a linear equation to the observed data

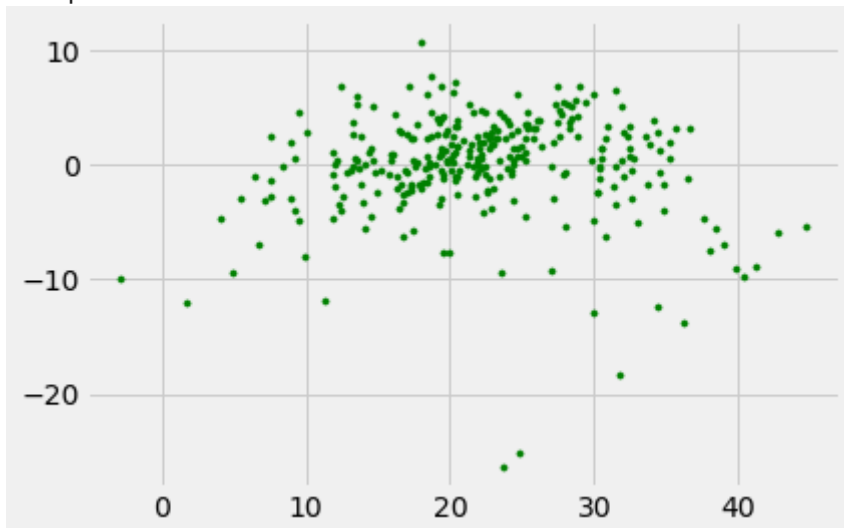


plotting residual errors in training data



```
plt.scatter(reg.predict(X_train), reg.predict(X_train) - y_train,
            color = "green", s = 10, label = 'Train data')
```

<matplotlib.collections.PathCollection at 0x7f2437923ed0>



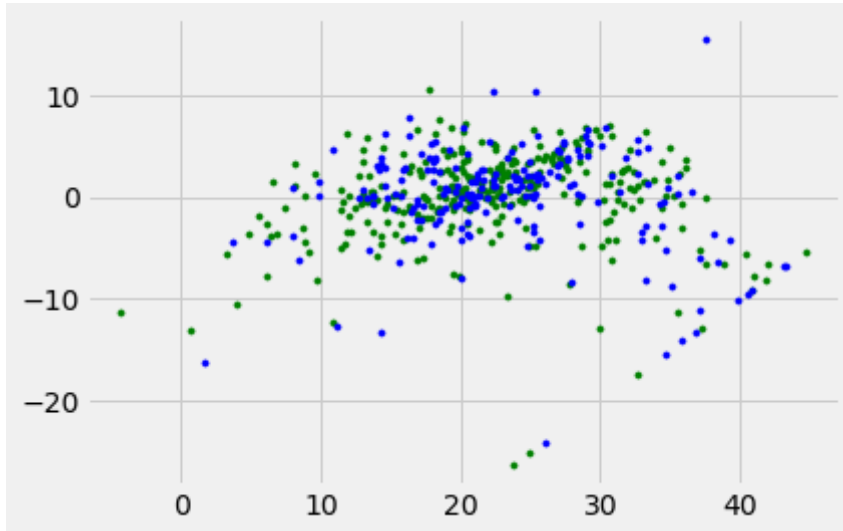
plotting residual errors in test data

```
plt.scatter(reg.predict(X_test), reg.predict(X_test) - y_test,
            color = "blue", s = 10, label = 'Test data')
```

```
<matplotlib.collections.PathCollection at 0x7f2437904dd0>
```



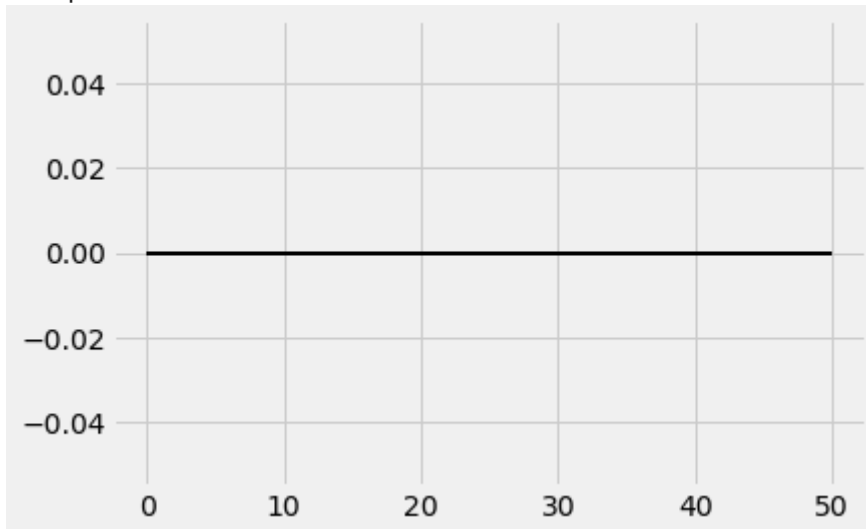
```
plt.scatter(reg.predict(X_train), reg.predict(X_train) - y_train,  
            color = "green", s = 10, label = 'Train data')  
plt.scatter(reg.predict(X_test), reg.predict(X_test) - y_test,  
            color = "blue", s = 10, label = 'Test data')  
plt.show()
```



plotting line for zero residual error

```
plt.hlines(y = 0, xmin = 0, xmax = 50, linewidth = 2)
```

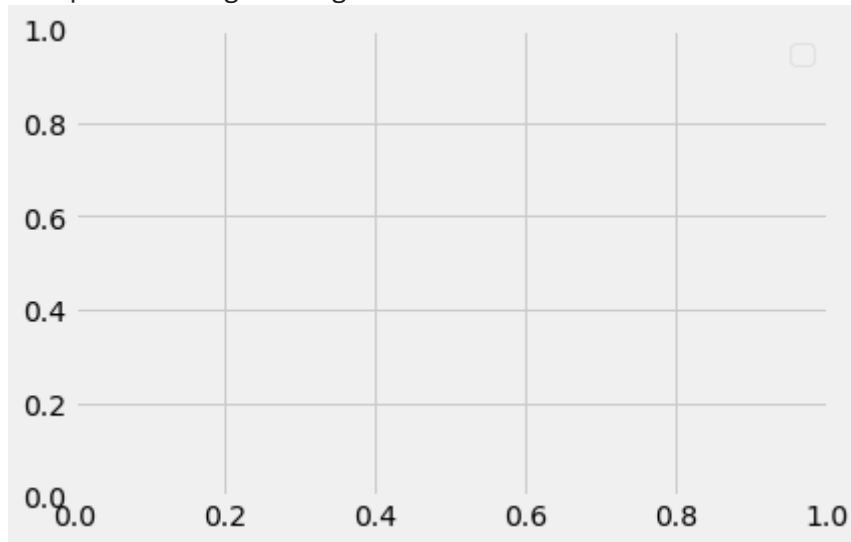
```
<matplotlib.collections.LineCollection at 0x7f243786c2d0>
```



plotting legend

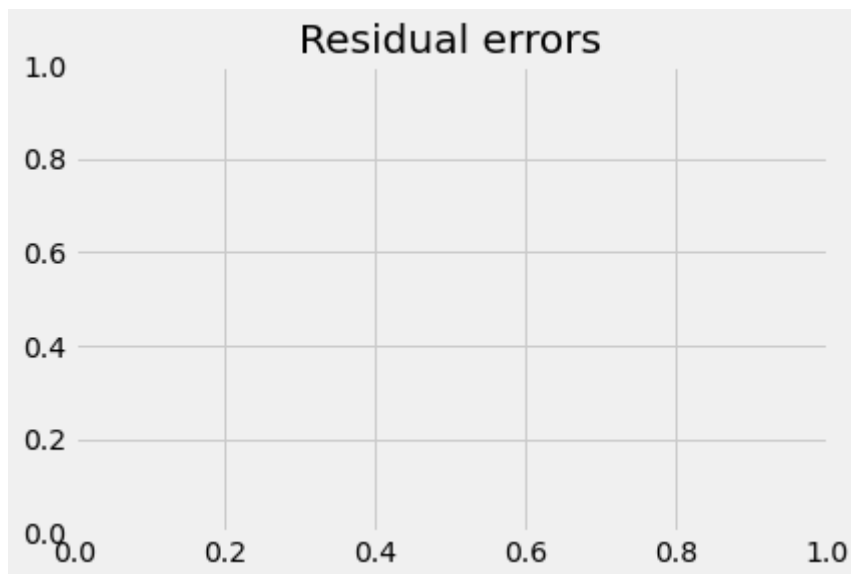
```
plt.legend(loc = 'upper right')
```

No handles with labels found to put in legend.
<matplotlib.legend.Legend at 0x7f243787ac50>

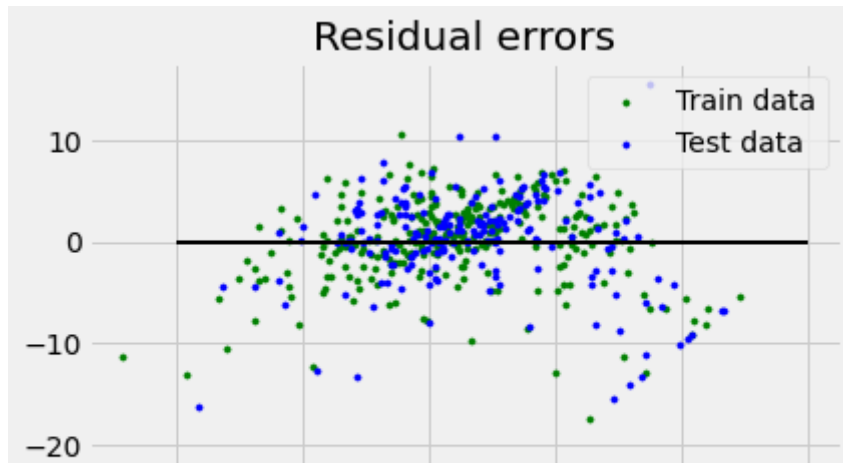


plot title

```
plt.title("Residual errors")  
plt.show()
```



```
plt.scatter(reg.predict(X_train), reg.predict(X_train) - y_train,  
            color = "green", s = 10, label = 'Train data')  
plt.scatter(reg.predict(X_test), reg.predict(X_test) - y_test,  
            color = "blue", s = 10, label = 'Test data')  
plt.hlines(y = 0, xmin = 0, xmax = 50, linewidth = 2)  
plt.legend(loc = 'upper right')  
plt.title("Residual errors")  
plt.show()
```

```
# Mean X and Y
mean_x = np.mean(X)
mean_y = np.mean(y)

# Total number of values
n = len(X)

# Using the formula to calculate m and c
numer = 0
denom = 0
for i in range(n):
    numer += (X[i] - mean_x) * (y[i] - mean_y)
    denom += (X[i] - mean_x) ** 2
m = numer / denom
c = mean_y - (m * mean_x)

# Print coefficients
print(m, c)
```

```
[-6.82675435e-03  1.93400822e-02 -8.65216721e-03  8.33753038e-05
 -9.40446529e-05  1.10197483e-03 -1.22813721e-01  1.09854680e-03
 -8.15772581e-03 -5.07922474e-03 -3.78049638e-03  3.08913388e-03
 -1.44417592e-02] [23.01118408 21.17757004 23.139098  22.52696389 22.53939641 22.455
 31.13885096 22.45582679 23.10445053 22.88872775 22.7977207  22.31633846
 23.54479768]
```

Calculating R2 Score

[+ Code](#)
[+ Text](#)

```
from sklearn.metrics import mean_squared_error
Y_pred = reg.predict(X)

# Calculating R2 Score
r2_score = reg.score(X, y)

print(r2_score)
```

```
0.7406426641094095
```

✓ 0s completed at 4:53 PM

● ✕