

Linear Regression (boston dataset)

December 18, 2022

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import sklearn as sl
from sklearn import linear_model
from sklearn.linear_model import LinearRegression
from sklearn.datasets import load_boston
```

```
[2]: df=load_boston()
```

C:\Users\Deepak\ana-conda-3\lib\site-packages\sklearn\utils\deprecation.py:87:
FutureWarning: Function load_boston is deprecated; `load_boston` is deprecated
in 1.0 and will be removed in 1.2.

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::

```
import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])
target = raw_df.values[1::2, 2]
```

Alternative datasets include the California housing dataset (i.e. :func:`~sklearn.datasets.fetch_california_housing`) and the Ames housing dataset. You can load the datasets as follows::

```
from sklearn.datasets import fetch_california_housing
```

```
housing = fetch_california_housing()
```

for the California housing dataset and::

```
from sklearn.datasets import fetch_openml
housing = fetch_openml(name="house_prices", as_frame=True)
```

for the Ames housing dataset.

```
warnings.warn(msg, category=FutureWarning)
```

```
[3]: df
```

```
[3]: {'data': array([[6.3200e-03, 1.8000e+01, 2.3100e+00, ..., 1.5300e+01,
3.9690e+02,
4.9800e+00],
[2.7310e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9690e+02,
9.1400e+00],
[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]]),
'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]],
'feature_names': array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE',
'DIS', 'RAD',
'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U7'),
'DESCR': ".. _boston_dataset:\n\nBoston house prices
dataset\n-----\n\n**Data Set Characteristics:** \n\n
: Number of Instances: 506 \n\n : Number of Attributes: 13 numeric/categorical
predictive. Median Value (attribute 14) is usually the target.\n\n : Attribute
Information (in order):\n - CRIM per capita crime rate by town\n
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.\n
- INDUS proportion of non-retail business acres per town\n - CHAS
Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)\n -
NOX nitric oxides concentration (parts per 10 million)\n - RM
average number of rooms per dwelling\n - AGE proportion of owner-
occupied units built prior to 1940\n - DIS weighted distances to
five Boston employment centres\n - RAD index of accessibility to
radial highways\n - TAX full-value property-tax rate per $10,000\n
- PTRATIO pupil-teacher ratio by town\n - B 1000(Bk - 0.63)^2
where Bk is the proportion of black people by town\n - LSTAT % lower
status of the population\n - MEDV Median value of owner-occupied
homes in $1000's\n\n : Missing Attribute Values: None\n\n : Creator:

```


[506 rows x 13 columns]

```
[6]: X=dataset  
      y=df.target
```

```
[7]: y
```

```
[7]: 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])

```

```
[8]: from sklearn.model_selection import train_test_split
```

```

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.30, random_state=42)
X_train

```

```

[8]:
      0      1      2      3      4      5      6      7      8      9  \
5      0.02985  0.0  2.18  0.0  0.458  6.430  58.7  6.0622  3.0  222.0
116    0.13158  0.0 10.01  0.0  0.547  6.176  72.5  2.7301  6.0  432.0
45     0.17142  0.0  6.91  0.0  0.448  5.682  33.8  5.1004  3.0  233.0
16     1.05393  0.0  8.14  0.0  0.538  5.935  29.3  4.4986  4.0  307.0
468   15.57570  0.0 18.10  0.0  0.580  5.926  71.0  2.9084 24.0  666.0
..      ...      ...      ...      ...      ...      ...      ...
106    0.17120  0.0  8.56  0.0  0.520  5.836  91.9  2.2110  5.0  384.0
270    0.29916 20.0  6.96  0.0  0.464  5.856  42.1  4.4290  3.0  223.0
348    0.01501 80.0  2.01  0.0  0.435  6.635  29.7  8.3440  4.0  280.0
435   11.16040  0.0 18.10  0.0  0.740  6.629  94.6  2.1247 24.0  666.0
102    0.22876  0.0  8.56  0.0  0.520  6.405  85.4  2.7147  5.0  384.0

      10      11      12
5      18.7  394.12  5.21
116    17.8  393.30 12.04
45     17.9  396.90 10.21
16     21.0  386.85  6.58
468    20.2  368.74 18.13
..      ...      ...
106    20.9  395.67 18.66
270    18.6  388.65 13.00
348    17.0  390.94  5.99
435    20.2  109.85 23.27
102    20.9   70.80 10.63

```

[354 rows x 13 columns]

```
[9]: from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
```

```
[10]: X_train=scaler.fit_transform(X_train)
```

```

[11]: X_test=scaler.transform(X_test)

[12]: from sklearn.linear_model import LinearRegression
      from sklearn.model_selection import cross_val_score

[13]: regression=LinearRegression()
      regression.fit(X_train,y_train)

[13]: LinearRegression()

[14]: mse=cross_val_score(regression,X_train,y_train,scoring='neg_mean_squared_error',cv=10)

[15]: np.mean(mse)

[15]: -25.55066079166079

[16]: reg_pred=regression.predict(X_test)

[17]: reg_pred

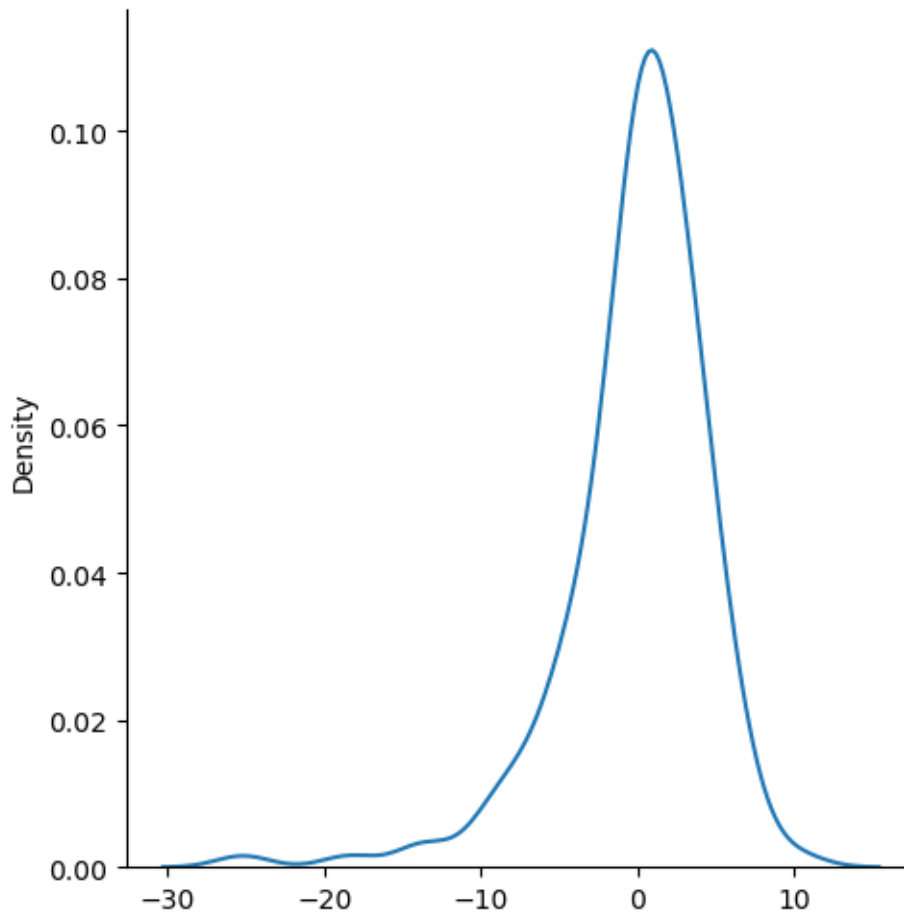
[17]: array([28.64896005, 36.49501384, 15.4111932 , 25.40321303, 18.85527988,
        23.14668944, 17.3921241 , 14.07859899, 23.03692679, 20.59943345,
        24.82286159, 18.53057049, -6.86543527, 21.80172334, 19.22571177,
        26.19191985, 20.27733882,  5.61596432, 40.44887974, 17.57695918,
        27.44319095, 30.1715964 , 10.94055823, 24.02083139, 18.07693812,
        15.934748 , 23.12614028, 14.56052142, 22.33482544, 19.3257627 ,
        22.16564973, 25.19476081, 25.31372473, 18.51345025, 16.6223286 ,
        17.50268505, 30.94992991, 20.19201752, 23.90440431, 24.86975466,
        13.93767876, 31.82504715, 42.56978796, 17.62323805, 27.01963242,
        17.19006621, 13.80594006, 26.10356557, 20.31516118, 30.08649576,
        21.3124053 , 34.15739602, 15.60444981, 26.11247588, 39.31613646,
        22.99282065, 18.95764781, 33.05555669, 24.85114223, 12.91729352,
        22.68101452, 30.80336295, 31.63522027, 16.29833689, 21.07379993,
        16.57699669, 20.36362023, 26.15615896, 31.06833034, 11.98679953,
        20.42550472, 27.55676301, 10.94316981, 16.82660609, 23.92909733,
         5.28065815, 21.43504661, 41.33684993, 18.22211675,  9.48269245,
        21.19857446, 12.95001331, 21.64822797,  9.3845568 , 23.06060014,
        31.95762512, 19.16662892, 25.59942257, 29.35043558, 20.13138581,
        25.57297369,  5.42970803, 20.23169356, 15.1949595 , 14.03241742,
        20.91078077, 24.82249135, -0.47712079, 13.70520524, 15.69525576,
        22.06972676, 24.64152943, 10.7382866 , 19.68622564, 23.63678009,
        12.07974981, 18.47894211, 25.52713393, 20.93461307, 24.6955941 ,
         7.59054562, 19.01046053, 21.9444339 , 27.22319977, 32.18608828,
        15.27826455, 34.39190421, 12.96314168, 21.01681316, 28.57880911,
        15.86300844, 24.85124135,  3.37937111, 23.90465773, 25.81792146,
        23.11020547, 25.33489201, 33.35545176, 20.60724498, 38.4772665 ,
        13.97398533, 25.21923987, 17.80946626, 20.63437371,  9.80267398,
        21.07953576, 22.3378417 , 32.32381854, 31.48694863, 15.46621287,

```

```
16.86242766, 28.99330526, 24.95467894, 16.73633557, 6.12858395,  
26.65990044, 23.34007187, 17.40367164, 13.38594123, 39.98342478,  
16.68286302, 18.28561759])
```

```
[18]: import seaborn as sns  
sns.displot(reg_pred-y_test,kind='kde')
```

```
[18]: <seaborn.axisgrid.FacetGrid at 0x244e5ed5790>
```



```
[19]: from sklearn.metrics import r2_score
```

```
[20]: score=r2_score(reg_pred,y_test)
```

```
[21]: score
```

```
[21]: 0.6693702691495591
```

```
[ ]:
```