

```
1 DSBDA 04
```

```
1 Roll No : 13320
```

```
In [75]: 1 import pandas as pd
          2 import numpy as np
          3 import matplotlib.pyplot as plt
```

```
In [76]: 1 x=np.array([95,85,80,70,60])
          2 y=np.array([85,95,70,65,70])
```

```
In [77]: 1 model= np.polyfit(x, y, 1)
          2 model
```

```
Out[77]: array([ 0.64383562, 26.78082192])
```

```
In [78]: 1 predict = np.poly1d(model)
          2 predict(65)
```

```
Out[78]: 68.63013698630137
```

```
In [79]: 1 y_pred= predict(x)
          2 y_pred
```

```
Out[79]: array([87.94520548, 81.50684932, 78.28767123, 71.84931507, 65.4109589 ])
```

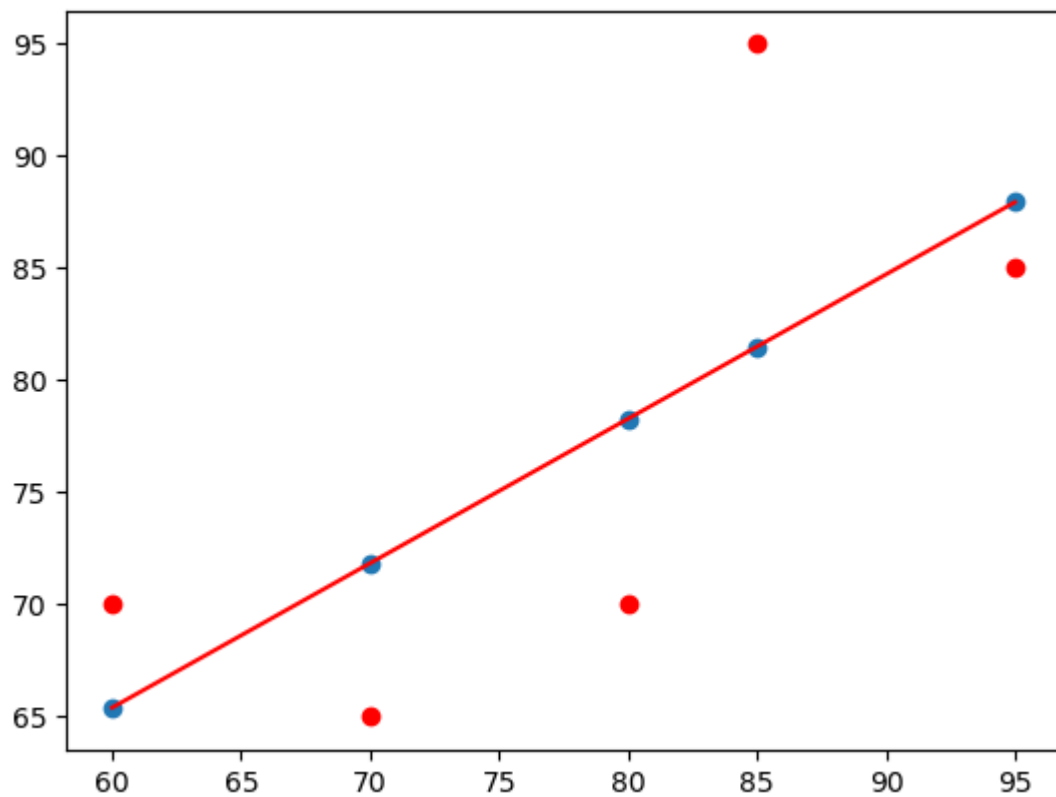
```
In [80]: 1 from sklearn.metrics import r2_score
          2 rs=r2_score(y, y_pred)
          3 print(rs)
          4 print("Accuracy of data = ",rs*100,"%")
```

```
0.4803218090889326
```

```
Accuracy of data = 48.03218090889326 %
```

```
In [81]: 1 y_line = model[1] + model[0]* x
2 plt.plot(x, y_line, c = 'r')
3 plt.scatter(x, y_pred)
4 plt.scatter(x,y,c='r')
5
```

Out[81]: <matplotlib.collections.PathCollection at 0x1e6bb4ab450>



```
In [82]: 1 from sklearn.datasets import load_boston
```

```
-----
-
ImportError                                Traceback (most recent call last)
Cell In[82], line 1
----> 1 from sklearn.datasets import load_boston

File D:\anaconda\Lib\site-packages\sklearn\datasets\__init__.py:157, in __
getattr__(name)
    108 if name == "load_boston":
    109     msg = textwrap.dedent("""
    110         `load_boston` has been removed from scikit-learn since ver
sion 1.2.
    111     (...)
    155         <https://www.researchgate.net/publication/4974606_Hedonic_
housing_prices_and_the_demand_for_clean_air>
    156     """)
--> 157     raise ImportError(msg)
    158 try:
    159     return globals()[name]
```

ImportError:
`load_boston` has been removed from scikit-learn since version 1.2.

The Boston housing prices dataset has an ethical problem: as investigated in [1], the authors of this dataset engineered a non-invertible variable "B" assuming that racial self-segregation had a positive impact on house prices [2]. Furthermore the goal of the research that led to the creation of this dataset was to study the impact of air quality but it did not give adequate demonstration of the validity of this assumption.

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 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.

```
[1] M Carlisle.  
"Racist data destruction?"  
<https://medium.com/@docintangible/racist-data-destruction-113e3eff54a8>  
  
[2] Harrison Jr, David, and Daniel L. Rubinfeld.  
"Hedonic housing prices and the demand for clean air."  
Journal of environmental economics and management 5.1 (1978): 81-102.  
<https://www.researchgate.net/publication/4974606\_Hedonic\_housing\_prices\_and\_the\_demand\_for\_clean\_air>
```

```
In [83]: 1 data_url = "http://lib.stat.cmu.edu/datasets/boston"  
2 raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)  
3 data = np.hstack([raw_df.values[::2, :], raw_df.values[1::2, :2]])  
4 target = raw_df.values[1::2, 2]
```

```
In [84]: 1 print(data)
```

```
[[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]]
```

In [85]:

1 `print(target)`

```
[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]
```

In [86]:

1 `data= pd.DataFrame(data)`

In [87]:

1

data

Out[87]:

	0	1	2	3	4	5	6	7	8	9	10	11	12
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
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99	9.67
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90	9.08
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90	5.64
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45	6.48
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90	7.88

506 rows × 13 columns

In [88]:

1

dt=pd.DataFrame(target)

2

dt

Out[88]:

	0
0	24.0
1	21.6
2	34.7
3	33.4
4	36.2
...	...
501	22.4
502	20.6
503	23.9
504	22.0
505	11.9

506 rows × 1 columns

In [89]:

1

data.isnull().sum()

Out[89]:

00

10

20

30

40

50

60

70

80

90

100

110

120

dtype: int64

In [90]:

1x=data

2y=target

3x

Out[90]:

	0	1	2	3	4	5	6	7	8	9	10	11	12
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
...
501	0.06263	0.0	11.93	0.0	0.573	6.593	69.1	2.4786	1.0	273.0	21.0	391.99	9.67
502	0.04527	0.0	11.93	0.0	0.573	6.120	76.7	2.2875	1.0	273.0	21.0	396.90	9.08
503	0.06076	0.0	11.93	0.0	0.573	6.976	91.0	2.1675	1.0	273.0	21.0	396.90	5.64
504	0.10959	0.0	11.93	0.0	0.573	6.794	89.3	2.3889	1.0	273.0	21.0	393.45	6.48
505	0.04741	0.0	11.93	0.0	0.573	6.030	80.8	2.5050	1.0	273.0	21.0	396.90	7.88

506 rows × 13 columns

In [91]:

1from sklearn.model_selection import train_test_split

2xtrain, xtest, ytrain, ytest =train_test_split(x, y, test_size =0.2,random_state=42)

3#xtrain,ytrain

In [92]:

1import sklearn

2from sklearn.linear_model import LinearRegression

3lm = LinearRegression()

4model=lm.fit(xtrain, ytrain)

In [93]:

1ytrain_pred=lm.predict(xtrain)

2ytest_pred=lm.predict(xtest)

In [94]:

1

dt=pd.DataFrame(ytrain_pred,ytrain)

2

dt

Out[94]:

	0
26.7	32.556927
21.7	21.927095
22.0	27.543826
22.9	23.603188
10.4	6.571910
...	...
18.5	19.494951
36.4	33.326364
19.2	23.796208
16.6	18.458353
23.1	23.249181

404 rows × 1 columns

In [95]:

1

dt=pd.DataFrame(ytest_pred,ytest)

2

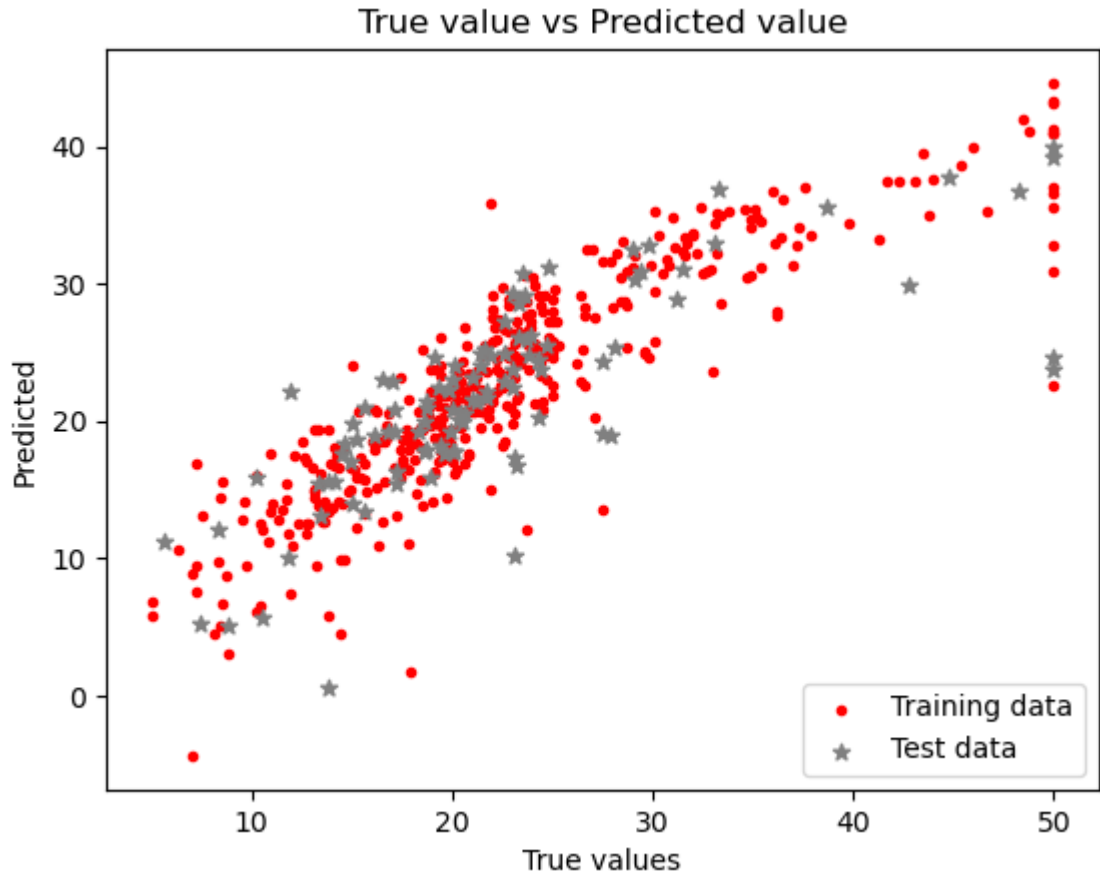
dt

Out[95]:

	0
22.6	24.889638
50.0	23.721411
23.0	29.364999
8.3	12.122386
21.2	21.443823
...	...
24.7	25.442171
14.1	15.571783
18.7	17.937195
28.1	25.305888
19.8	22.373233

102 rows × 1 columns

```
In [96]: 1 plt.scatter(ytrain ,ytrain_pred,c='red',marker='.',label='Training data')
2 plt.scatter(ytest,ytest_pred ,c='grey',marker='*',label='Test data')
3 plt.xlabel('True values')
4 plt.ylabel('Predicted')
5 plt.title("True value vs Predicted value")
6 plt.legend(loc= 'lower right')
7 plt.plot()
8 plt.show()
```



```
In [97]: 1 from sklearn.metrics import mean_squared_error, r2_score
2 mse = mean_squared_error(ytest, ytest_pred)
3 print(mse)
4
```

33.44897999767632

```
In [98]: 1 mse = mean_squared_error(ytrain_pred,ytrain)
2 print(mse)
```

19.326470203585725

```
In [99]: 1 rs=r2_score(ytest,ytest_pred)
2 print(rs)
3 print("Accuracy of test data = ",rs*100,"%")
```

0.5892223849182534

Accuracy of test data = 58.92223849182534 %

```
In [100]: 1 rs=r2_score(ytrain,ytrain_pred)
          2 print(rs)
          3 print("Accuracy of trained data = ",rs*100,"%")
```

0.7730135569264234

Accuracy of trained data = 77.30135569264233 %

```
In [ ]: 1
```