

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import seaborn as sns
        4 import matplotlib.pyplot as plt
        5 %matplotlib inline
        6 import warnings
        7 warnings.filterwarnings('ignore')
```

```
In [2]: 1 from sklearn.datasets import load_boston #" sklearn me load bostn kr k inbult data hota h"
```

```
In [3]: 1 boston = load_boston()
```

In [4]: 1 print(boston)

```
{'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 1
4) is usually the target.\n\n :Attribute Information (in order):\n      - CRIM    per capita crime rate by town\n      - ZN      proportion of residenti
al 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 o
f rooms per dwelling\n      - AGE   proportion of owner-occupied units built prior to 1940\n      - DIS   weighted distances to five Boston empl
oyment 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 o
f the population\n      - MEDV   Median value of owner-occupied homes in $1000's\n\n :Missing Attribute Values: None\n\n :Creator: Harris
on, D. and Rubinfeld, D.L.\n\nThis is a copy of UCI ML housing dataset.\nhttps://archive.ics.uci.edu/ml/machine-learning-databases/housing/\n\n
\nThis dataset was taken from the StatLib library which is maintained at Carnegie Mellon University.\n\nThe Boston house-price data of Harrison,
D. and Rubinfeld, D.L. 'Hedonic\nprices and the demand for clean air', J. Environ. Economics & Management,\nvol.5, 81-102, 1978.  Used in Be
lsley, Kuh & Welsch, 'Regression diagnostics\n...', Wiley, 1980.  N.B. Various transformations are used in the table on\npages 244-261 of the latt
er.\n\nThe Boston house-price data has been used in many machine learning papers that address regression\nproblems.  \n      \n.. topic:: Refer
ences\n\n - Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.\n -
Quinlan,R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Le
arning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.\n", 'filename': 'C:\\Users\\User\\anaconda3\\lib\\site-packages\\sklear
n\\datasets\\data\\boston_house_prices.csv'}

```

In [5]: 1 boston.keys()

Out[5]: dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename'])

```
In [6]: 1 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 14) is usually the target.

: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 $1000(B_k - 0.63)^2$ where B_k is the proportion of black people 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/> (<https://archive.ics.uci.edu/ml/machine-learning-databases/housing/>)

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

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 regression problems.

.. topic:: References

- Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.
- Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.

In [7]:

```
1 print(boston.data) #input feature
```

```
[[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 [8]: 1 `print(boston.target) #outout feature`

```
[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 [9]: 1 `print(boston.feature_names) # sare columns`

```
['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO'
'B' 'LSTAT']
```

In [10]: 1 `## Lets prepare dataframe`

In [11]: 1 `dataset = pd.DataFrame(boston.data, columns = boston.feature_names)`

In [12]: 1 `dataset.head()`

Out[12]:

	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

In [13]: 1 `dataset['Price'] = boston.target # price ko alag se add kia hai`

In [14]: 1 `dataset.head()`

Out[14]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	Price
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	24.0
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	21.6
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	34.7
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	33.4
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	36.2

In [15]: 1 dataset.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
```

In [16]: 1 dataset.describe().T

Out[16]:

	count	mean	std	min	25%	50%	75%	max
CRIM	506.0	3.613524	8.601545	0.00632	0.082045	0.25651	3.677083	88.9762
ZN	506.0	11.363636	23.322453	0.00000	0.000000	0.00000	12.500000	100.0000
INDUS	506.0	11.136779	6.860353	0.46000	5.190000	9.69000	18.100000	27.7400
CHAS	506.0	0.069170	0.253994	0.00000	0.000000	0.00000	0.000000	1.0000
NOX	506.0	0.554695	0.115878	0.38500	0.449000	0.53800	0.624000	0.8710
RM	506.0	6.284634	0.702617	3.56100	5.885500	6.20850	6.623500	8.7800
AGE	506.0	68.574901	28.148861	2.90000	45.025000	77.50000	94.075000	100.0000
DIS	506.0	3.795043	2.105710	1.12960	2.100175	3.20745	5.188425	12.1265
RAD	506.0	9.549407	8.707259	1.00000	4.000000	5.00000	24.000000	24.0000
TAX	506.0	408.237154	168.537116	187.00000	279.000000	330.00000	666.000000	711.0000
PTRATIO	506.0	18.455534	2.164946	12.60000	17.400000	19.05000	20.200000	22.0000
B	506.0	356.674032	91.294864	0.32000	375.377500	391.44000	396.225000	396.9000
LSTAT	506.0	12.653063	7.141062	1.73000	6.950000	11.36000	16.955000	37.9700
Price	506.0	22.532806	9.197104	5.00000	17.025000	21.20000	25.000000	50.0000

```
In [17]: 1 ## check missing value  
        2 dataset.isnull().sum()
```

```
Out[17]: 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  
         Price    0  
         dtype: int64
```

In [18]:

```
1 ## EDA
2 dataset.corr()
```

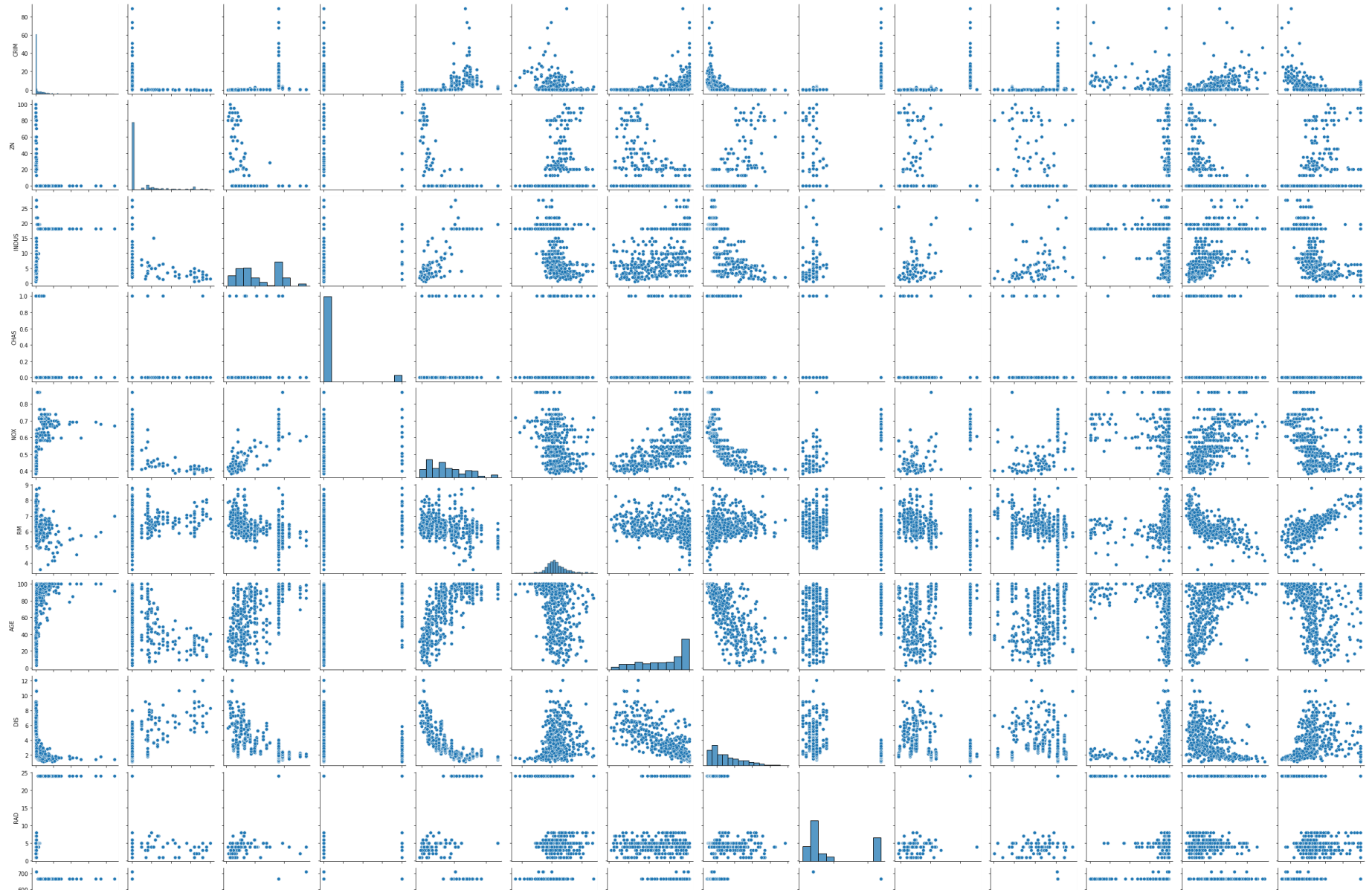
Out[18]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
CRIM	1.000000	-0.200469	0.406583	-0.055892	0.420972	-0.219247	0.352734	-0.379670	0.625505	0.582764	0.289946	-0.385064	0.455621
ZN	-0.200469	1.000000	-0.533828	-0.042697	-0.516604	0.311991	-0.569537	0.664408	-0.311948	-0.314563	-0.391679	0.175520	-0.412995
INDUS	0.406583	-0.533828	1.000000	0.062938	0.763651	-0.391676	0.644779	-0.708027	0.595129	0.720760	0.383248	-0.356977	0.603800
CHAS	-0.055892	-0.042697	0.062938	1.000000	0.091203	0.091251	0.086518	-0.099176	-0.007368	-0.035587	-0.121515	0.048788	-0.053929
NOX	0.420972	-0.516604	0.763651	0.091203	1.000000	-0.302188	0.731470	-0.769230	0.611441	0.668023	0.188933	-0.380051	0.590879
RM	-0.219247	0.311991	-0.391676	0.091251	-0.302188	1.000000	-0.240265	0.205246	-0.209847	-0.292048	-0.355501	0.128069	-0.613808
AGE	0.352734	-0.569537	0.644779	0.086518	0.731470	-0.240265	1.000000	-0.747881	0.456022	0.506456	0.261515	-0.273534	0.602339
DIS	-0.379670	0.664408	-0.708027	-0.099176	-0.769230	0.205246	-0.747881	1.000000	-0.494588	-0.534432	-0.232471	0.291512	-0.496996
RAD	0.625505	-0.311948	0.595129	-0.007368	0.611441	-0.209847	0.456022	-0.494588	1.000000	0.910228	0.464741	-0.444413	0.488676
TAX	0.582764	-0.314563	0.720760	-0.035587	0.668023	-0.292048	0.506456	-0.534432	0.910228	1.000000	0.460853	-0.441808	0.543993
PTRATIO	0.289946	-0.391679	0.383248	-0.121515	0.188933	-0.355501	0.261515	-0.232471	0.464741	0.460853	1.000000	-0.177383	0.374044
B	-0.385064	0.175520	-0.356977	0.048788	-0.380051	0.128069	-0.273534	0.291512	-0.444413	-0.441808	-0.177383	1.000000	-0.366087
LSTAT	0.455621	-0.412995	0.603800	-0.053929	0.590879	-0.613808	0.602339	-0.496996	0.488676	0.543993	0.374044	-0.366087	1.000000
Price	-0.388305	0.360445	-0.483725	0.175260	-0.427321	0.695360	-0.376955	0.249929	-0.381626	-0.468536	-0.507787	0.333461	-0.737661

```
In [19]: 1 plt.figure(figsize = (15,10))  
2 sns.pairplot(dataset)
```

Out[19]: <seaborn.axisgrid.PairGrid at 0x1f009336fa0>

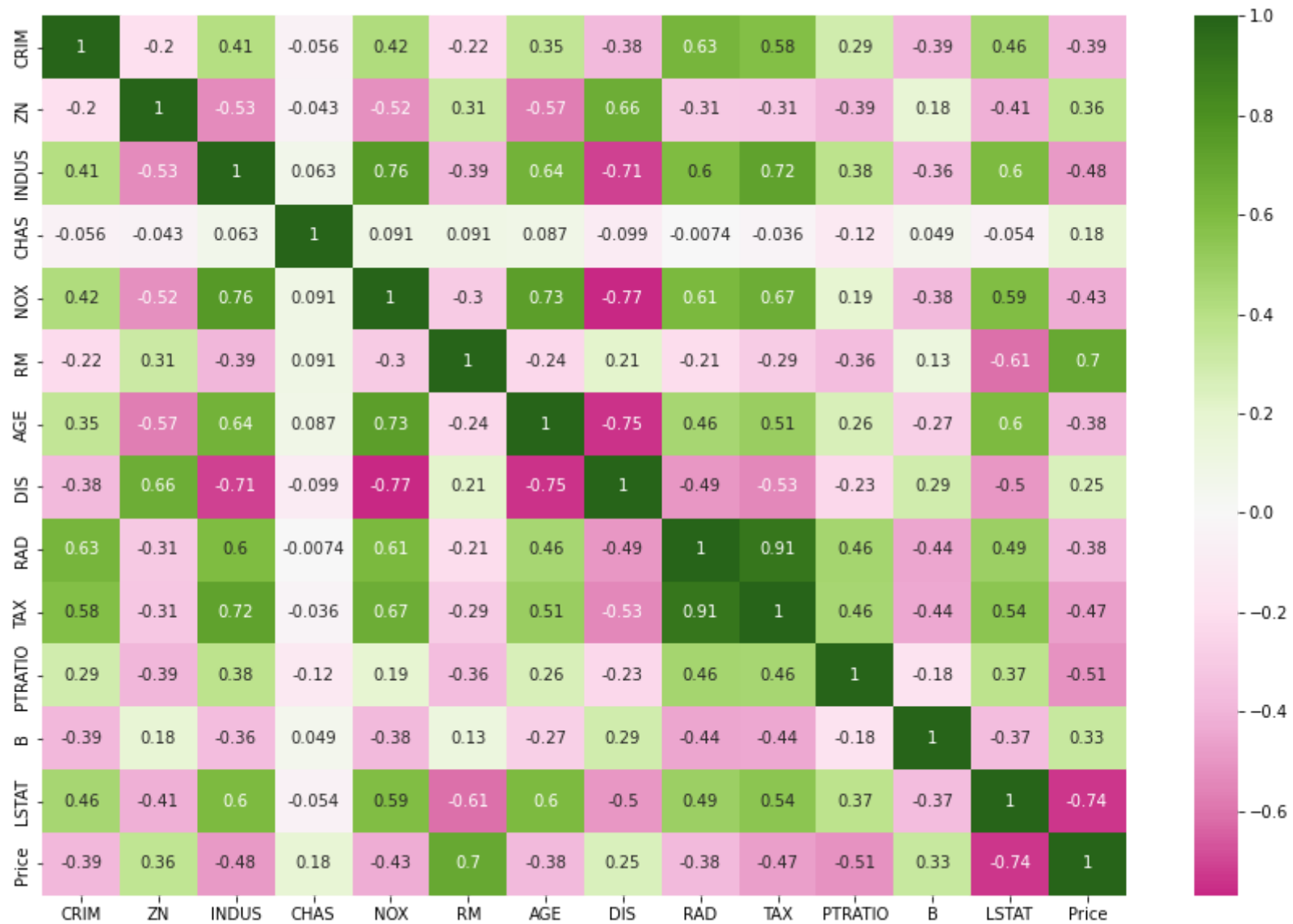
<Figure size 1080x720 with 0 Axes>





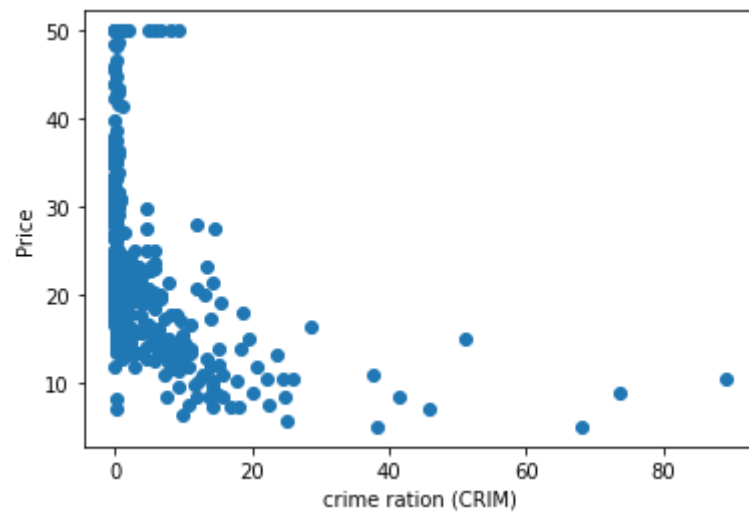
In [20]: 1 *# we usually use pair plot in multi variant analysis*

```
In [21]: 1 plt.figure(figsize = (15,10))
2         sns.heatmap(dataset.corr(), annot=True,center = 0,cmap = 'PiYG')
3         plt.show()
```



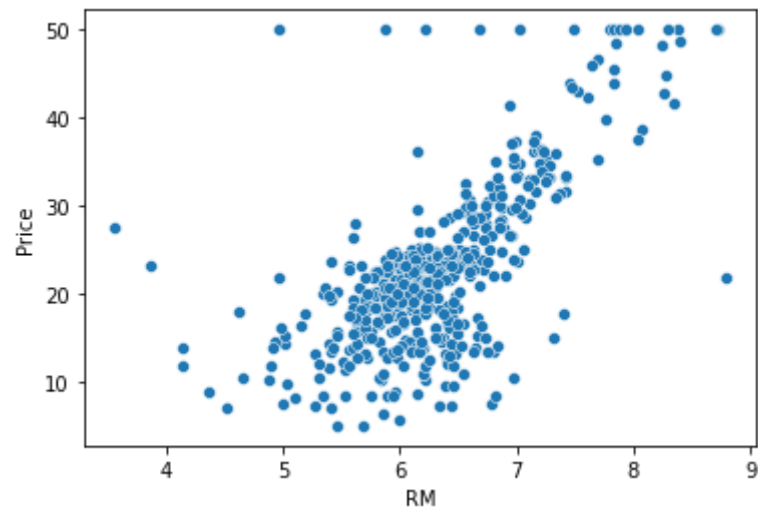
```
In [22]: 1 plt.scatter(dataset['CRIM'],dataset['Price'])  
        2 plt.xlabel('crime ration (CRIM)')  
        3 plt.ylabel('Price')
```

Out[22]: Text(0, 0.5, 'Price')



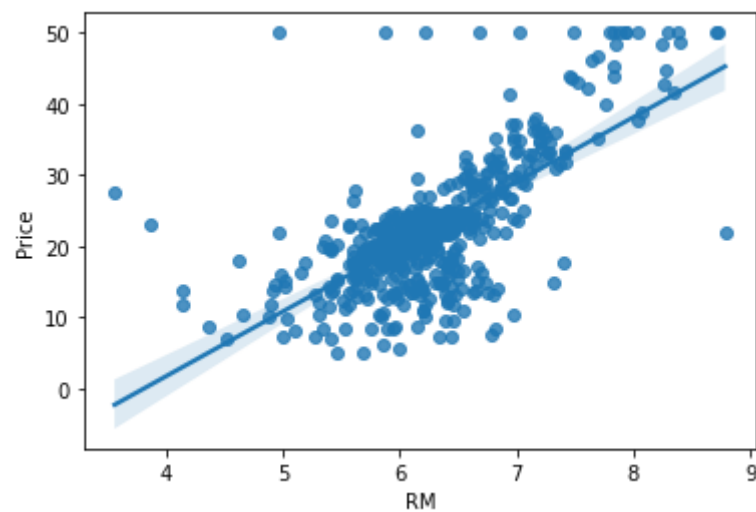

```
In [23]: 1 sns.scatterplot(x=dataset['RM'],y=dataset['Price'])
```

```
Out[23]: <AxesSubplot:xlabel='RM', ylabel='Price'>
```



```
In [24]: 1 sns.regplot(x="RM",y="Price",data=dataset)
```

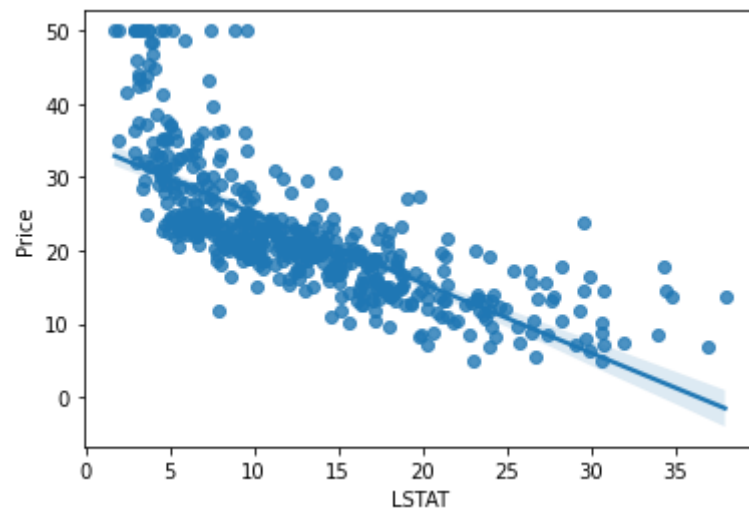
Out[24]: <AxesSubplot:xlabel='RM', ylabel='Price'>



```
In [25]: 1 # shaded region is "Ridge and lasso"
```

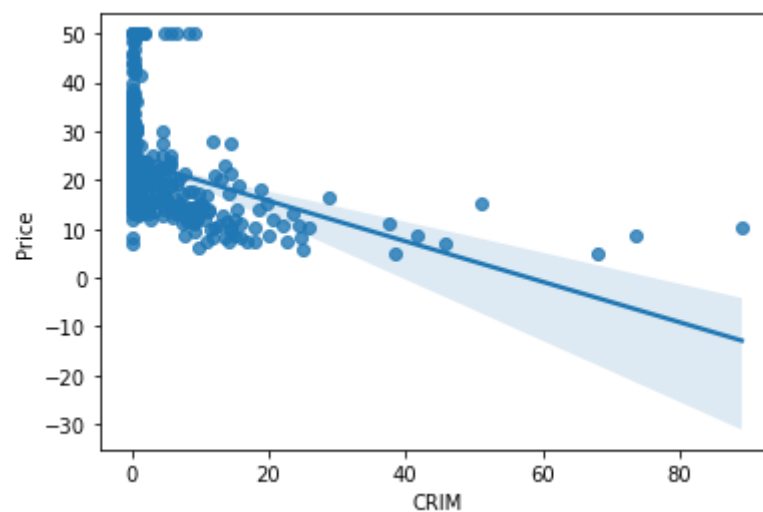
```
In [26]: 1 sns.regplot(x = "LSTAT",y = "Price",data = dataset)
```

Out[26]: <AxesSubplot:xlabel='LSTAT', ylabel='Price'>



```
In [27]: 1 sns.regplot(x="CRIM",y="Price",data = dataset)
```

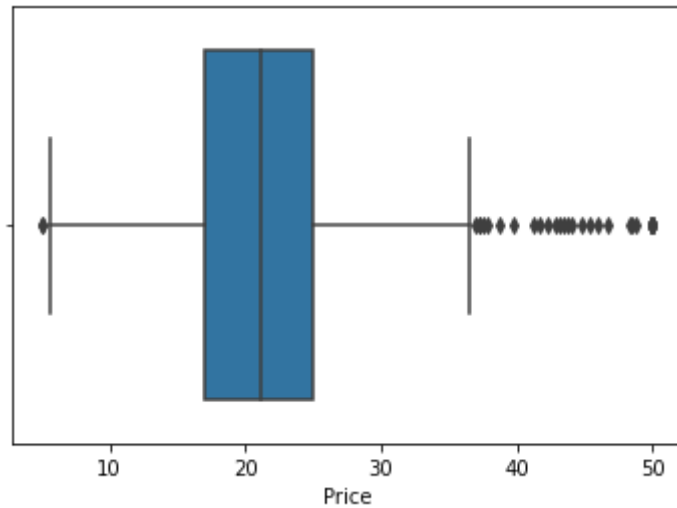
Out[27]: <AxesSubplot:xlabel='CRIM', ylabel='Price'>



```
In [28]: 1 # if there is more concentrated point there it is less shaded and where it is less point more shaded
```

```
In [29]: 1 sns.boxplot(dataset['Price'])
```

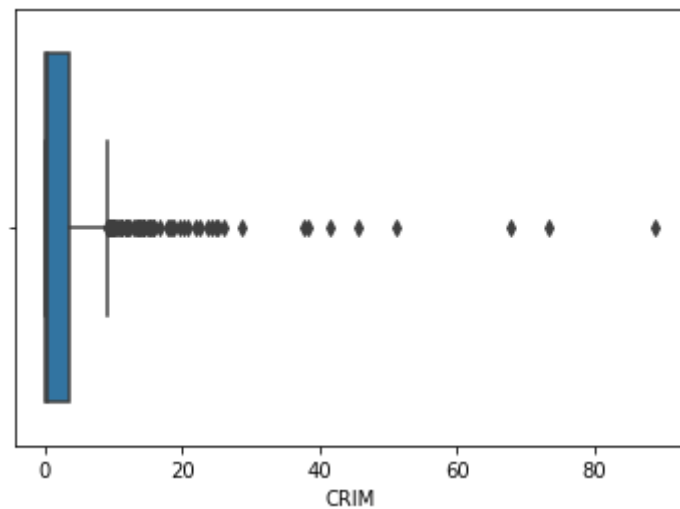
Out[29]: <AxesSubplot:xlabel='Price'>



```
In [30]: 1 # no need to remove outlier because it is dependent feature
```

```
In [31]: 1 sns.boxplot(dataset['CRIM'])
```

```
Out[31]: <AxesSubplot:xlabel='CRIM'>
```



```
In [32]: 1 # not every time need to remove outlier  
2 # suppose heart cancer data , so not every patient have cancer  
3 # therefore we will hypertune the data
```

```
In [33]: 1 ## Independent and dependent feature
```

```
In [34]: 1 X = dataset.iloc[:, :-1]
        2 y = dataset.iloc[:, -1]
```

```
In [35]: 1 X.head()
```

Out[35]:

	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

```
In [36]: 1 y.head()
```

Out[36]:

```
0    24.0
1    21.6
2    34.7
3    33.4
4    36.2
Name: Price, dtype: float64
```

```
In [37]: 1 from sklearn.model_selection import train_test_split
```

```
In [38]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)
```

In [39]: 1 X_train.head()

Out[39]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
478	10.23300	0.0	18.10	0.0	0.614	6.185	96.7	2.1705	24.0	666.0	20.2	379.70	18.03
26	0.67191	0.0	8.14	0.0	0.538	5.813	90.3	4.6820	4.0	307.0	21.0	376.88	14.81
7	0.14455	12.5	7.87	0.0	0.524	6.172	96.1	5.9505	5.0	311.0	15.2	396.90	19.15
492	0.11132	0.0	27.74	0.0	0.609	5.983	83.5	2.1099	4.0	711.0	20.1	396.90	13.35
108	0.12802	0.0	8.56	0.0	0.520	6.474	97.1	2.4329	5.0	384.0	20.9	395.24	12.27

In [40]: 1 X_train.shape

Out[40]: (339, 13)

In [41]: 1 X_test.head()

Out[41]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
173	0.09178	0.0	4.05	0.0	0.510	6.416	84.1	2.6463	5.0	296.0	16.6	395.50	9.04
274	0.05644	40.0	6.41	1.0	0.447	6.758	32.9	4.0776	4.0	254.0	17.6	396.90	3.53
491	0.10574	0.0	27.74	0.0	0.609	5.983	98.8	1.8681	4.0	711.0	20.1	390.11	18.07
72	0.09164	0.0	10.81	0.0	0.413	6.065	7.8	5.2873	4.0	305.0	19.2	390.91	5.52
452	5.09017	0.0	18.10	0.0	0.713	6.297	91.8	2.3682	24.0	666.0	20.2	385.09	17.27

In [42]: 1 X_test.shape

Out[42]: (167, 13)

```
In [43]: 1 y_train.head()
```

```
Out[43]: 478    14.6  
         26     16.6  
         7     27.1  
         492    20.1  
         108    19.8  
         Name: Price, dtype: float64
```

```
In [44]: 1 y_train.shape
```

```
Out[44]: (339,)
```

```
In [45]: 1 y_test.head()
```

```
Out[45]: 173    23.6  
         274    32.4  
         491    13.6  
         72     22.8  
         452    16.1  
         Name: Price, dtype: float64
```

```
In [46]: 1 y_test.shape
```

```
Out[46]: (167,)
```

```
In [47]: 1 # if we do feature scaling it will helpfull in gradient decent calculation
```

feature scaling

```
In [49]: 1 from sklearn.preprocessing import StandardScaler  
         2 scaler = StandardScaler()  
         3 scaler # ye object bnare hai
```

```
Out[49]: StandardScaler()
```

```
In [50]: 1 X_train=scaler.fit_transform(X_train)
```



```
In [51]: 1 X_test = scaler.transform(X_test) # here only transform use not fit transform to avoid dataleakage
```

```
In [ ]: 1 #data leakage ka mtlb toh pta hi hoga tujhe
2 # ab yaha sirf transform kyo kia hai ? ye doubt ata hai
3 # fir transform me kya hoga ki , jese standard scaler hai to fit karega pehle jisme mean or standard deviation nikalega
4 # fir usko us data pe apply kar k transform karenge
5 # ab yaha test data me agar fir kr diya toh uska mean alag hojayega or standard deviation bhi
6 # toh scaling toh kari lekin same level ki scaling nahi hui
7 # toh jo train data se mean , SD aya hai fit se usi ko use kar ke transofrm karenge
8
9 # ab ek or fanda hota h jisme , bolte hai ki pehle data ko split karn ahai fir feature scaling karenge , nito data leakage hojayega .....
10 # ab esa kyo ?
11 # kyoki data scale pehle kr diya or fir split kia toh test data is impacted by train data through transformartion , toh new data ayega toh uski accuracy kam hogi
12
```

```
In [52]: 1 X_train
```

```
Out[52]: array([[ 0.89624872, -0.51060139,  0.98278223, ...,  0.86442095,
  0.24040357,  0.77155612],
 [-0.34895881, -0.51060139, -0.44867555, ...,  1.22118698,
  0.20852839,  0.32248963],
 [-0.41764058,  0.03413008, -0.48748013, ..., -1.36536677,
  0.43481957,  0.92775316],
 ...,
 [-0.43451148,  2.97567999, -1.32968321, ..., -0.56264319,
  0.36745216, -0.90756208],
 [ 1.01703049, -0.51060139,  0.98278223, ...,  0.86442095,
 -2.80977992,  1.50233514],
 [-0.40667333, -0.51060139, -0.38831288, ...,  1.17659123,
 -3.25117205, -0.26046005]])
```

Model Training

```
In [53]: 1 from sklearn.linear_model import LinearRegression
```

```
In [54]: 1 regression = LinearRegression()
```

```
In [55]: 1 regression
```

```
Out[55]: LinearRegression()
```

```
In [56]: 1 regression.fit(X_train,y_train)
```

```
Out[56]: LinearRegression()
```

```
In [57]: 1 ## print the coefficient and the intercept
```

```
In [58]: 1 print(regression.coef_) # 13 coefficient will be there bcoz 13 indepent featurea are there
```

```
[-0.98858032  0.86793276  0.40502822  0.86183791 -1.90009974  2.80813518  
-0.35866856 -3.04553498  2.03276074 -1.36400909 -2.0825356   1.04125684  
-3.92628626]
```

```
In [59]: 1 # or ek sirf intercept hoga , thitha wala jisme x ni hoga  
2 print(regression.intercept_)
```

```
22.970796460176988
```

```
In [60]: 1 # ab agar sare feature ko nahi leta hai or ya zero kr deta hai toh price ki value 22.97 rahegi
```

```
In [61]: 1 # predition on test data  
2 reg_pred = regression.predict(X_test)
```

```
In [62]: 1 reg_pred # ye pridicted values hai
```

```
Out[62]: array([28.53469469, 36.6187006 , 15.63751079, 25.5014496 , 18.7096734 ,
 23.16471591, 17.31011035, 14.07736367, 23.01064388, 20.54223482,
 24.91632351, 18.41098052, -6.52079687, 21.83372604, 19.14903064,
 26.0587322 , 20.30232625, 5.74943567, 40.33137811, 17.45791446,
 27.47486665, 30.2170757 , 10.80555625, 23.87721728, 17.99492211,
 16.02608791, 23.268288 , 14.36825207, 22.38116971, 19.3092068 ,
 22.17284576, 25.05925441, 25.13780726, 18.46730198, 16.60405712,
 17.46564046, 30.71367733, 20.05106788, 23.9897768 , 24.94322408,
 13.97945355, 31.64706967, 42.48057206, 17.70042814, 26.92507869,
 17.15897719, 13.68918087, 26.14924245, 20.2782306 , 29.99003492,
 21.21260347, 34.03649185, 15.41837553, 25.95781061, 39.13897274,
 22.96118424, 18.80310558, 33.07865362, 24.74384155, 12.83640958,
 22.41963398, 30.64804979, 31.59567111, 16.34088197, 20.9504304 ,
 16.70145875, 20.23215646, 26.1437865 , 31.12160889, 11.89762768,
 20.45432404, 27.48356359, 10.89034224, 16.77707214, 24.02593714,
 5.44691807, 21.35152331, 41.27267175, 18.13447647, 9.8012101 ,
 21.24024342, 13.02644969, 21.80198374, 9.48201752, 22.99183857,
 31.90465631, 18.95594718, 25.48515032, 29.49687019, 20.07282539,
 25.5616062 , 5.59584382, 20.18410904, 15.08773299, 14.34562117,
 20.85155407, 24.80149389, -0.19785401, 13.57649004, 15.64401679,
 22.03765773, 24.70314482, 10.86409112, 19.60231067, 23.73429161,
 12.08082177, 18.40997903, 25.4366158 , 20.76506636, 24.68588237,
 7.4995836 , 18.93015665, 21.70801764, 27.14350579, 31.93765208,
 15.19483586, 34.01357428, 12.85763091, 21.06646184, 28.58470042,
 15.77437534, 24.77512495, 3.64655689, 23.91169589, 25.82292925,
 23.03339677, 25.35158335, 33.05655447, 20.65930467, 38.18917361,
 14.04714297, 25.26034469, 17.6138723 , 20.60883766, 9.8525544 ,
 21.06756951, 22.20145587, 32.2920276 , 31.57638342, 15.29265938,
 16.7100235 , 29.10550932, 25.17762329, 16.88159225, 6.32621877,
 26.70210263, 23.3525851 , 17.24168182, 13.22815696, 39.49907507,
 16.53528575, 18.14635902, 25.06620426, 23.70640231, 22.20167772,
 21.22272327, 16.89825921, 23.15518273, 28.69699805, 6.65526482,
 23.98399958, 17.21004545, 21.0574427 , 25.01734597, 27.65461859,
 20.70205823, 40.38214871])
```

```
In [63]: 1 #####
          2 # assumption in Linear regression
```

- 1 Linear regression is an analysis that assesses whether one or more predictor variables explain the dependent (criterion) variable. The regression has five key assumptions:
- 2
- 3 Linear relationship
- 4 Multivariate normality
- 5 No or little multicollinearity
- 6 No auto-correlation
- 7 Homoscedasticity

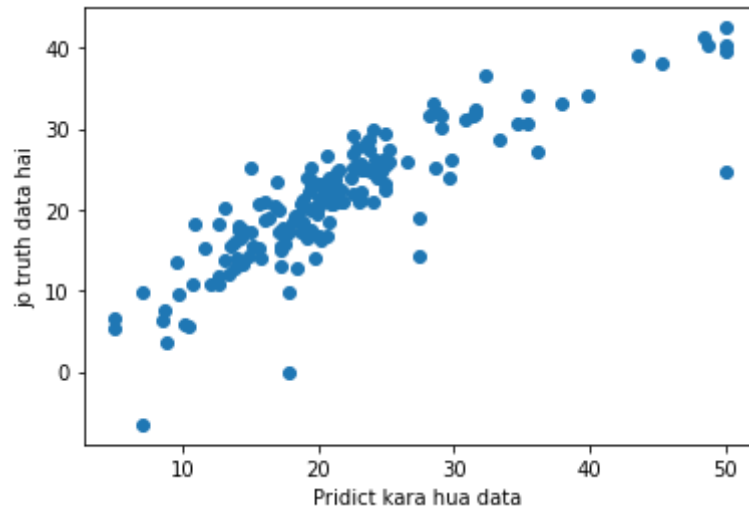
In [64]:

```

1 plt.scatter(y_test,reg_pred) # toh liner graph jesa araha hai toh model sahi banaya hai
2 plt.xlabel('Pridict kara hua data ')
3 plt.ylabel('jo truth data hai')

```

Out[64]: Text(0, 0.5, 'jo truth data hai')



In [65]:

```

1 ## residual
2 residual = y_test-reg_pred

```

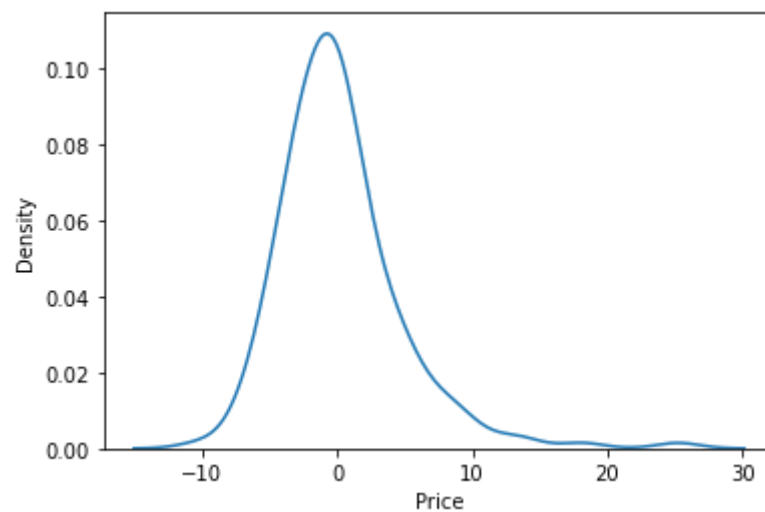
```
In [66]: 1 residual
```

```
Out[66]: 173 -4.934695  
274 -4.218701  
491 -2.037511  
72 -2.701450  
452 -2.609673  
...  
110 0.642557  
321 -1.917346  
265 -4.854619  
29 0.297942  
262 8.417851
```

```
Name: Price, Length: 167, dtype: float64
```

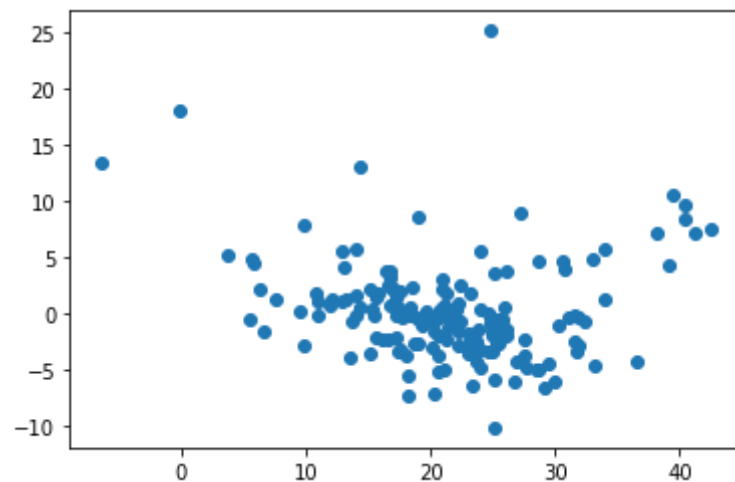
```
In [67]: 1 sns.kdeplot(residual) # ye graph normal distribution araha hai toh model sahi hai .  
2 # thode outlier hai usko lasso ridge se sahi karege
```

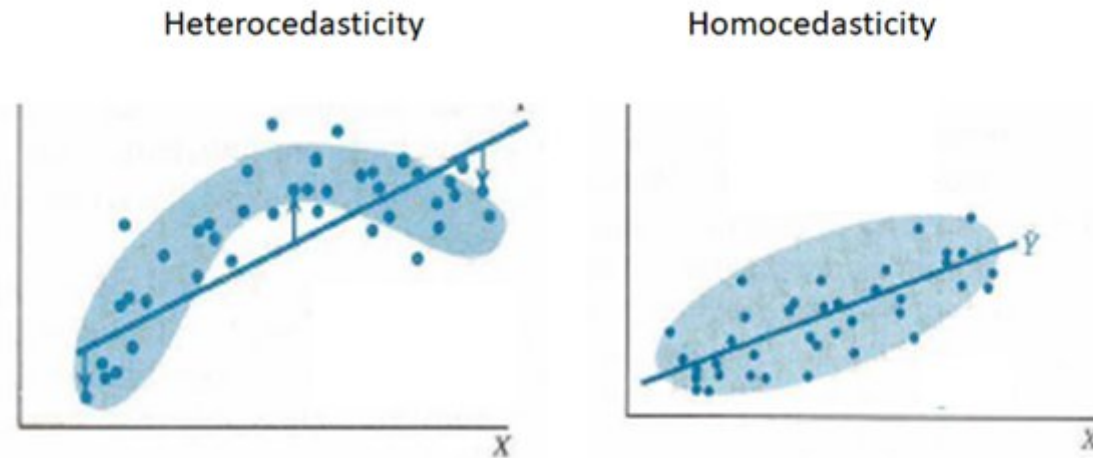
```
Out[67]: <AxesSubplot:xlabel='Price', ylabel='Density'>
```



```
In [68]: 1  ## Scatter plot with prediction and residual  
2  ### UNIFORM DISTRIBUTION  
3  plt.scatter(reg_pred,residual)
```

Out[68]: <matplotlib.collections.PathCollection at 0x1f015782700>





In [69]:

```

1  ## Performance Matrix
2  from sklearn.metrics import mean_squared_error
3  from sklearn.metrics import mean_absolute_error
4  print(mean_squared_error(y_test,reg_pred))
5  print(mean_absolute_error(y_test,reg_pred))
6  print(np.sqrt(mean_squared_error(y_test,reg_pred)))

```

```

20.724023437339753
3.148255754816832
4.552364598463062

```

In [70]:

```

1  ## R square and adjusted R square
2  from sklearn.metrics import r2_score
3  score = r2_score(y_test,reg_pred)
4  print(score)

```

```

0.7261570836552476

```

In [71]:

```

1  ## adjusted R square
2  1 - (1-score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1)

```

Out[71]: 0.7028893848808568

```
In [72]: 1 ## now need to do with ridge, lasso , elastic
```

```
In [73]: 1 from sklearn.linear_model import Ridge  
2 ridge = Ridge()
```

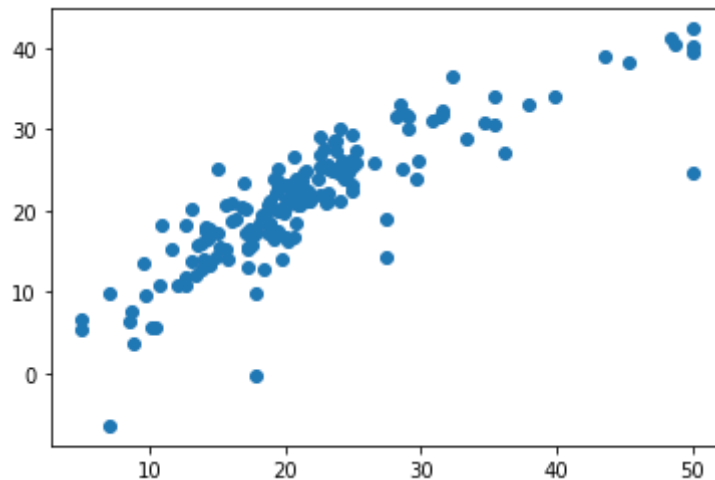
```
In [74]: 1 ridge.fit(X_train,y_train)
```

Out[74]: Ridge()

```
In [75]: 1 ridge_predict = ridge.predict(X_test)
```

```
In [76]: 1 plt.scatter(y_test,ridge_predict)
```

Out[76]: <matplotlib.collections.PathCollection at 0x1f0157f3220>



```
In [77]: 1 ## R square and adjusted R square  
2 from sklearn.metrics import r2_score  
3 score = r2_score(y_test,ridge_predict)  
4 print(score)
```

0.7257819060246209


```
In [78]: 1  ## adjusted R square  
2  1 - (1-score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1)
```

Out[78]: 0.7024823294123338

```
In [79]: 1  from sklearn.linear_model import Lasso
```

```
In [80]: 1  lasso = Lasso(alpha=0.1)
```

```
In [81]: 1  lasso
```

Out[81]: Lasso(alpha=0.1)

```
In [82]: 1  lasso.fit(X_train,y_train)
```

Out[82]: Lasso(alpha=0.1)

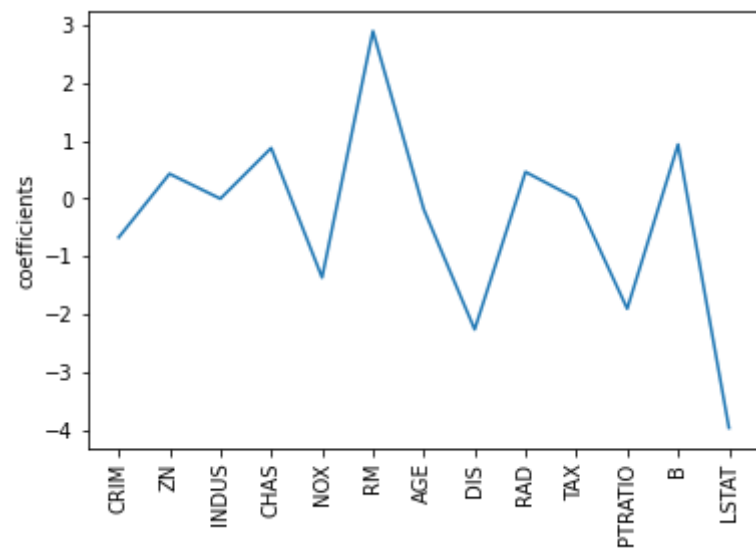
```
In [83]: 1  lasso_pred = lasso.predict(X_test)
```

```
In [84]: 1  from sklearn.metrics import r2_score
```

```
In [85]: 1  r2_score(y_test,lasso_pred)
```

Out[85]: 0.7112387502289164

```
In [86]: 1 plt.plot(X.columns,lasso.coef_)
2         plt.xticks(rotation=90)
3         plt.ylabel("coefficients")
4         plt.show()
```



```
In [87]: 1 # it shows which feature is most important to predict y value
```

```
In [ ]: 1
```

In []:

1