

```
In [1]: 1 import pandas as pd  
2 import numpy as np  
3 import matplotlib.pyplot as plt  
4 %matplotlib inline
```

## Lets load the Boston House Pricing Dataset

```
In [2]: 1 from sklearn.datasets import load_boston
```

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

```
In [4]: 1 boston.keys()
```

```
Out[4]: dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename'])
```

In [5]:

```
1 ## Lets check the description of the dataset
2 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  $1000(Bk - 0.63)^2$  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/> (<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 [6]: 1 print(boston.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 [7]: 1 `print(boston.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. 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 [8]: 1 `print(boston.feature_names)`

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

## Preparing The Dataset

In [9]: 1 `dataset=pd.DataFrame(boston.data,columns=boston.feature_names)`

In [10]: 1 dataset.head()

Out[10]:

	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 [11]: 1 dataset['Price']=boston.target

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.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 [14]:

```
1 ## Summarizing The Stats of the data
2 dataset.describe()
```

Out[14]:

	<b>CRIM</b>	<b>ZN</b>	<b>INDUS</b>	<b>CHAS</b>	<b>NOX</b>	<b>RM</b>	<b>AGE</b>
<b>count</b>	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
<b>mean</b>	3.613524	11.363636	11.136779	0.069170	0.554695	6.284634	68.574901
<b>std</b>	8.601545	23.322453	6.860353	0.253994	0.115878	0.702617	28.148861
<b>min</b>	0.006320	0.000000	0.460000	0.000000	0.385000	3.561000	2.900000
<b>25%</b>	0.082045	0.000000	5.190000	0.000000	0.449000	5.885500	45.025000
<b>50%</b>	0.256510	0.000000	9.690000	0.000000	0.538000	6.208500	77.500000
<b>75%</b>	3.677083	12.500000	18.100000	0.000000	0.624000	6.623500	94.075000
<b>max</b>	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000

In [15]:

```
1 ## Check the missing Values
2 dataset.isnull().sum()
```

Out[15]:

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 [16]:

```
1 ### Exploratory Data Analysis
2 ## Correlation
3 dataset.corr()
```

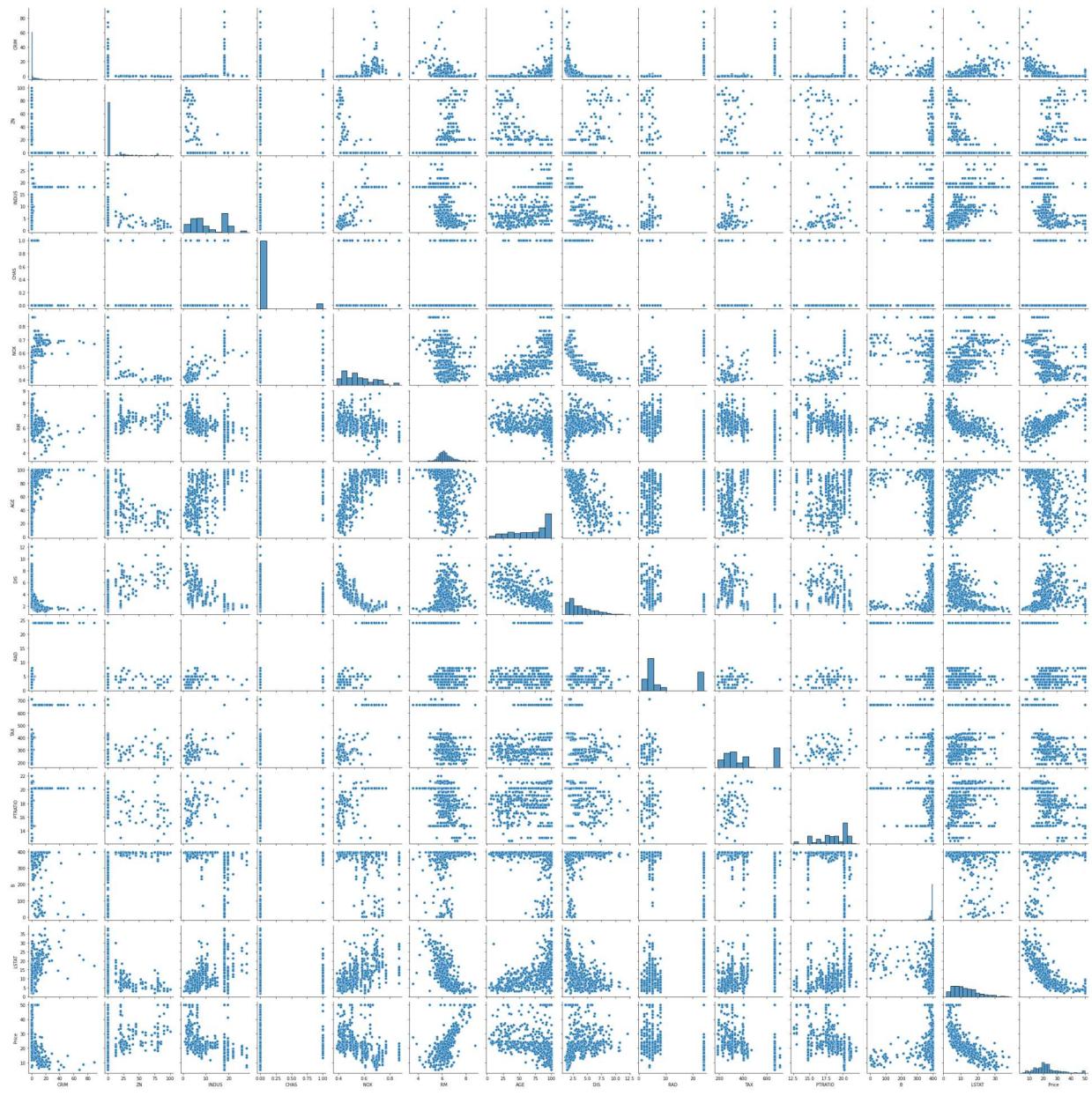
Out[16]:

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

In [17]:

```
1 import seaborn as sns  
2 sns.pairplot(dataset)
```

Out[17]: &lt;seaborn.axisgrid.PairGrid at 0x2344ef6de08&gt;



## Analyzing The Correlated Features

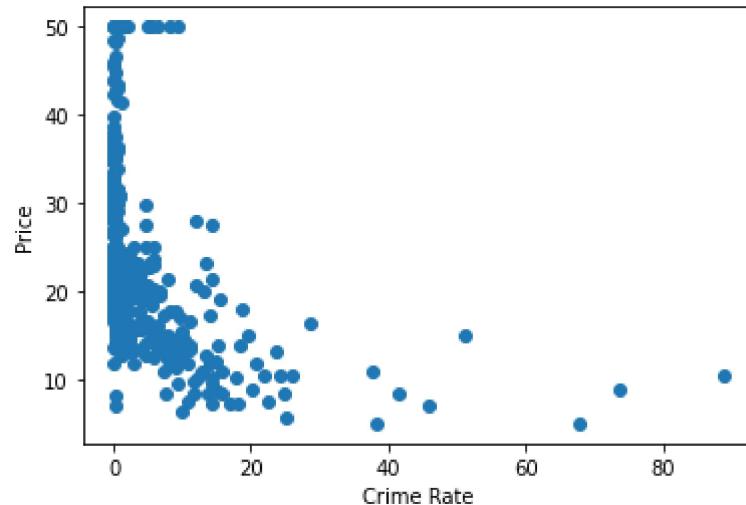
In [18]: 1 dataset.corr()

Out[18]:

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	
CRIM	1.000000	-0.200469	0.406583	-0.055892	0.420972	-0.219247	0.352734	-0.379670	0.
ZN	-0.200469	1.000000	-0.533828	-0.042697	-0.516604	0.311991	-0.569537	0.664408	-0.
INDUS	0.406583	-0.533828	1.000000	0.062938	0.763651	-0.391676	0.644779	-0.708027	0.
CHAS	-0.055892	-0.042697	0.062938	1.000000	0.091203	0.091251	0.086518	-0.099176	-0.
NOX	0.420972	-0.516604	0.763651	0.091203	1.000000	-0.302188	0.731470	-0.769230	0
RM	-0.219247	0.311991	-0.391676	0.091251	-0.302188	1.000000	-0.240265	0.205246	-0.
AGE	0.352734	-0.569537	0.644779	0.086518	0.731470	-0.240265	1.000000	-0.747881	0.
DIS	-0.379670	0.664408	-0.708027	-0.099176	-0.769230	0.205246	-0.747881	1.000000	-0.
RAD	0.625505	-0.311948	0.595129	-0.007368	0.611441	-0.209847	0.456022	-0.494588	1.
TAX	0.582764	-0.314563	0.720760	-0.035587	0.668023	-0.292048	0.506456	-0.534432	0.
PTRATIO	0.289946	-0.391679	0.383248	-0.121515	0.188933	-0.355501	0.261515	-0.232471	0.
B	-0.385064	0.175520	-0.356977	0.048788	-0.380051	0.128069	-0.273534	0.291512	-0.
LSTAT	0.455621	-0.412995	0.603800	-0.053929	0.590879	-0.613808	0.602339	-0.496996	0.
Price	-0.388305	0.360445	-0.483725	0.175260	-0.427321	0.695360	-0.376955	0.249929	-0.

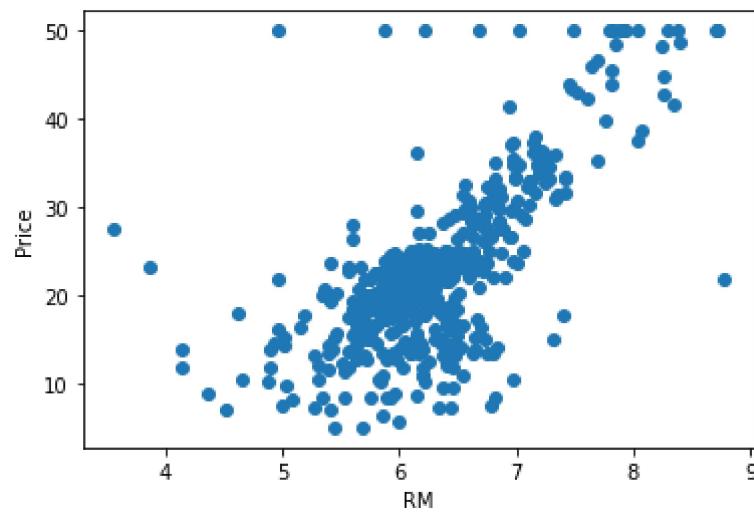
```
In [19]: 1 plt.scatter(dataset['CRIM'],dataset['Price'])  
2 plt.xlabel("Crime Rate")  
3 plt.ylabel("Price")
```

Out[19]: Text(0, 0.5, 'Price')



```
In [20]: 1 plt.scatter(dataset['RM'],dataset['Price'])  
2 plt.xlabel("RM")  
3 plt.ylabel("Price")
```

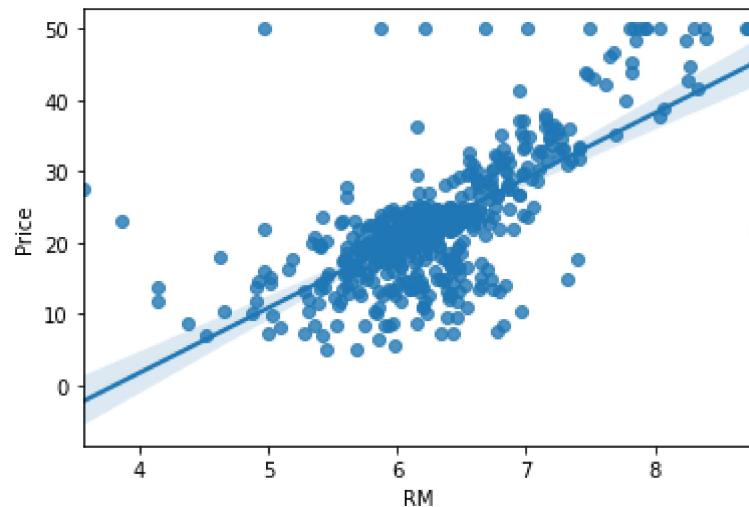
Out[20]: Text(0, 0.5, 'Price')



In [21]:

```
1 import seaborn as sns  
2 sns.regplot(x="RM",y="Price",data=dataset)
```

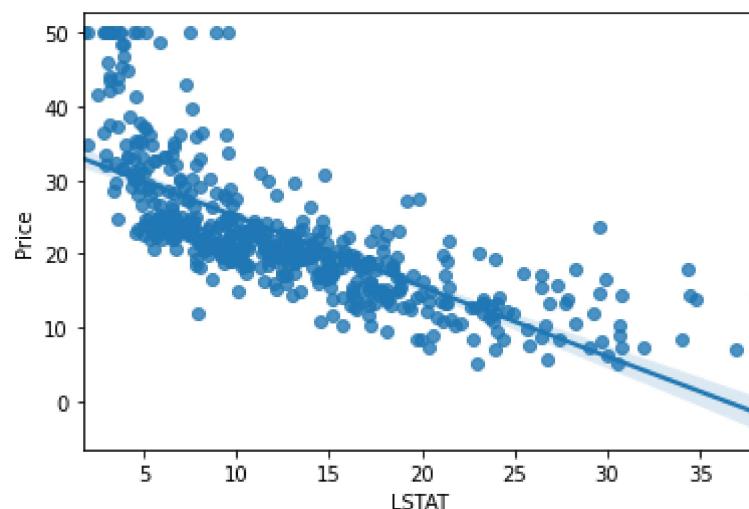
Out[21]: &lt;AxesSubplot:xlabel='RM', ylabel='Price'&gt;



In [22]:

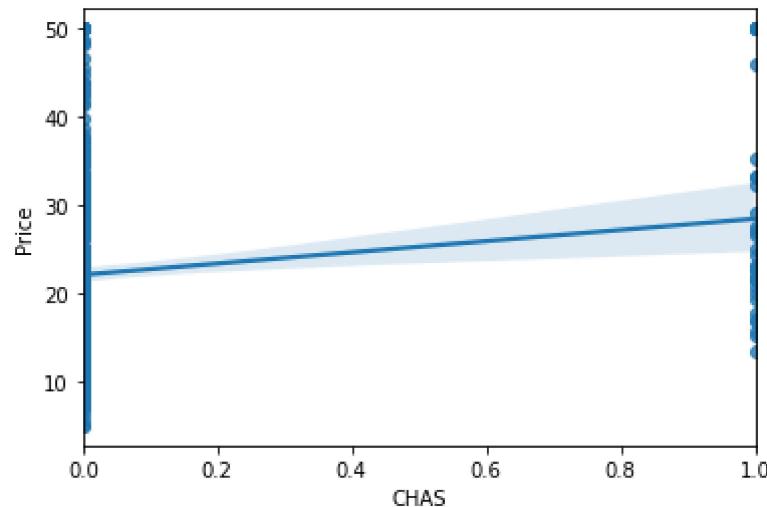
```
1 sns.regplot(x="LSTAT",y="Price",data=dataset)
```

Out[22]: &lt;AxesSubplot:xlabel='LSTAT', ylabel='Price'&gt;



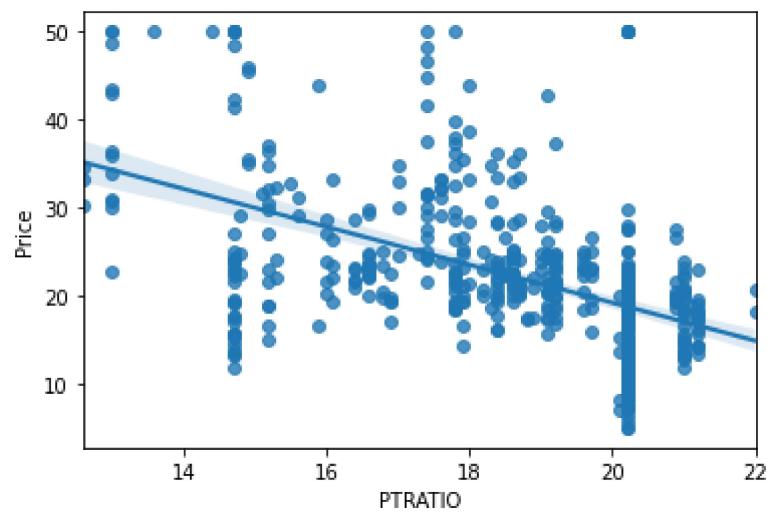
```
In [23]: 1 sns.regplot(x="CHAS",y="Price",data=dataset)
```

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



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

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



```
In [25]: 1 ## Independent and Dependent features
2
3 X=dataset.iloc[:, :-1]
4 y=dataset.iloc[:, -1]
```

```
In [26]: 1 X.head()
```

Out[26]:

	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 [27]: 1 y
```

Out[27]:

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

Name: Price, Length: 506, dtype: float64

```
In [28]: 1 ##Train Test Split
2 from sklearn.model_selection import train_test_split
3 X_train,X_test,y_train,y_test=train_test_split(X,y,test_size=0.3,random_stat
```

In [29]: 1 X\_train

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSI
5	0.02985	0.0	2.18	0.0	0.458	6.430	58.7	6.0622	3.0	222.0	18.7	394.12	5
116	0.13158	0.0	10.01	0.0	0.547	6.176	72.5	2.7301	6.0	432.0	17.8	393.30	12
45	0.17142	0.0	6.91	0.0	0.448	5.682	33.8	5.1004	3.0	233.0	17.9	396.90	10
16	1.05393	0.0	8.14	0.0	0.538	5.935	29.3	4.4986	4.0	307.0	21.0	386.85	6
468	15.57570	0.0	18.10	0.0	0.580	5.926	71.0	2.9084	24.0	666.0	20.2	368.74	18
...	...	...	...	...	...	...	...	...	...	...	...	...	...
106	0.17120	0.0	8.56	0.0	0.520	5.836	91.9	2.2110	5.0	384.0	20.9	395.67	18
270	0.29916	20.0	6.96	0.0	0.464	5.856	42.1	4.4290	3.0	223.0	18.6	388.65	13
348	0.01501	80.0	2.01	0.0	0.435	6.635	29.7	8.3440	4.0	280.0	17.0	390.94	5
435	11.16040	0.0	18.10	0.0	0.740	6.629	94.6	2.1247	24.0	666.0	20.2	109.85	23
102	0.22876	0.0	8.56	0.0	0.520	6.405	85.4	2.7147	5.0	384.0	20.9	70.80	10

354 rows × 13 columns

In [30]: 1 X\_test

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSI
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
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
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
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
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
...	...	...	...	...	...	...	...	...	...	...	...	...	...
441	9.72418	0.0	18.10	0.0	0.740	6.406	97.2	2.0651	24.0	666.0	20.2	385.96	19
23	0.98843	0.0	8.14	0.0	0.538	5.813	100.0	4.0952	4.0	307.0	21.0	394.54	19
225	0.52693	0.0	6.20	0.0	0.504	8.725	83.0	2.8944	8.0	307.0	17.4	382.00	4
433	5.58107	0.0	18.10	0.0	0.713	6.436	87.9	2.3158	24.0	666.0	20.2	100.19	16
447	9.92485	0.0	18.10	0.0	0.740	6.251	96.6	2.1980	24.0	666.0	20.2	388.52	16

152 rows × 13 columns

In [31]: 1 ## Standardize the dataset

```
2 from sklearn.preprocessing import StandardScaler
3 scaler=StandardScaler()
```

```
In [32]: 1 X_train=scaler.fit_transform(X_train)
```

```
In [33]: 1 X_test=scaler.transform(X_test)
```

```
In [34]: 1 import pickle
2 pickle.dump(scaler,open('scaling.pkl','wb'))
```

```
In [34]: 1 X_train
```

```
Out[34]: array([[-0.41425879, -0.50512499, -1.29214218, ... , 0.18727079,
       0.39651419, -1.01531611],
      [-0.40200818, -0.50512499, -0.16208345, ... , -0.21208981,
       0.3870674 , -0.05366252],
      [-0.39721053, -0.50512499, -0.60948856, ... , -0.16771641,
       0.42854113, -0.31132373],
      ... ,
      [-0.41604586,  3.03838247, -1.3166773 , ... , -0.56707702,
       0.35987906, -0.90549329],
      [ 0.92611293, -0.50512499,  1.00549958, ... ,  0.8528718 ,
       -2.87841346,  1.52750437],
      [-0.39030549, -0.50512499, -0.37135358, ... ,  1.16348561,
       -3.32828832, -0.25218837]])
```

```
In [35]: 1 X_test
```

```
Out[35]: array([[-0.406801 , -0.50512499, -1.02225586, ... , -0.74457062,
       0.41241246, -0.47605794],
      [-0.41105674,  1.26662874, -0.68165068, ... , -0.30083661,
       0.42854113, -1.25185755],
      [-0.4051199 , -0.50512499,  2.39678516, ... ,  0.8084984 ,
       0.35031705,  0.79535229],
      ... ,
      [-0.35439903, -0.50512499, -0.71195877, ... , -0.38958342,
       0.25688594, -1.09697922],
      [ 0.25423453, -0.50512499,  1.00549958, ... ,  0.8528718 ,
       -2.98970133,  0.53487511],
      [ 0.77732457, -0.50512499,  1.00549958, ... ,  0.8528718 ,
       0.33199949,  0.56585078]])
```

## Model Training

```
In [36]: 1 from sklearn.linear_model import LinearRegression
```

```
In [37]: 1 regression=LinearRegression()
```

```
In [38]: 1 regression.fit(X_train,y_train)
```

```
Out[38]: LinearRegression()
```

```
In [39]: 1 ## print the coefficients and the intercept  
2 print(regression.coef_)
```

```
[-1.10834602  0.80843998  0.34313466  0.81386426 -1.79804295  2.913858  
-0.29893918 -2.94251148  2.09419303 -1.44706731 -2.05232232  1.02375187  
-3.88579002]
```

```
In [40]: 1 print(regression.intercept_)
```

```
23.01581920903955
```

```
In [41]: 1 ## on which parameters the model has been trained  
2 regression.get_params()
```

```
Out[41]: {'copy_X': True,  
          'fit_intercept': True,  
          'n_jobs': None,  
          'normalize': False,  
          'positive': False}
```

```
In [42]: 1 ##### Prediction With Test Data  
2 reg_pred=regression.predict(X_test)
```

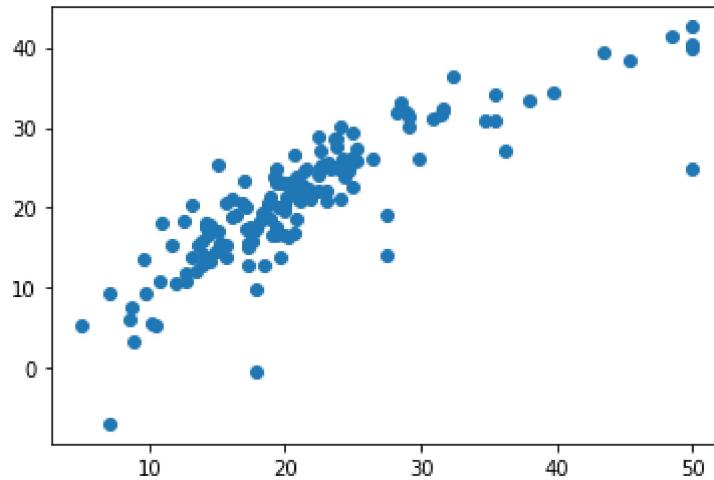
```
In [43]: 1 reg_pred
```

```
Out[43]: 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])
```

## Assumptions

```
In [44]: 1 ## plot a scatter plot for the prediction  
2 plt.scatter(y_test,reg_pred)
```

Out[44]: <matplotlib.collections.PathCollection at 0x21810ccb248>



```
In [45]: 1 ## Residuals  
2 residuals=y_test-reg_pred
```

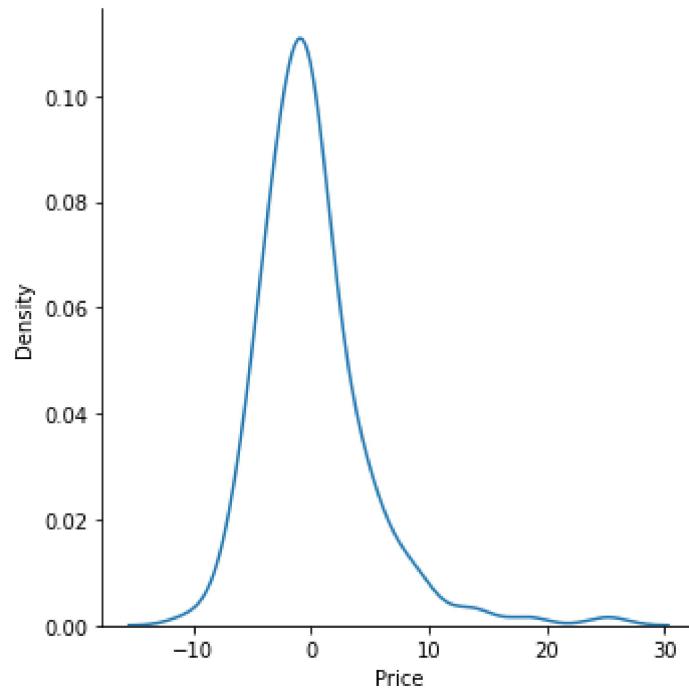
```
In [46]: 1 residuals
```

```
Out[46]: 173    -5.048960  
274    -4.095014  
491    -1.811193  
72     -2.603213  
452    -2.755280  
...  
441    -0.303672  
23     1.114059  
225    10.016575  
433    -2.382863  
447    -5.685618  
Name: Price, Length: 152, dtype: float64
```

In [47]:

```
1 ## Plot this residuals  
2  
3 sns.displot(residuals,kind="kde")
```

Out[47]: &lt;seaborn.axisgrid.FacetGrid at 0x21810fde788&gt;



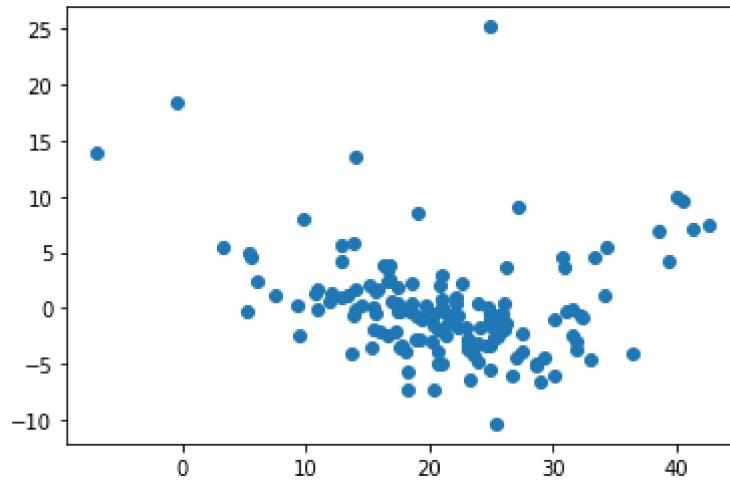
In [48]:

```

1 ## Scatter plot with respect to prediction and residuals
2 ## uniform distribution
3 plt.scatter(reg_pred,residuals)

```

Out[48]: &lt;matplotlib.collections.PathCollection at 0x21810facec8&gt;



In [49]:

```

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

```

3.162709871457405  
 21.517444231177205  
 4.6386899261728205

## R square and adjusted R square

Formula

$$R^2 = 1 - \frac{SSR}{SST}$$

$R^2$  = coefficient of determination SSR = sum of squares of residuals SST = total sum of squares

```
In [50]: 1 from sklearn.metrics import r2_score
2 score=r2_score(y_test,reg_pred)
3 print(score)
```

0.7112260057484933

In [ ]:

1

$$\text{Adjusted R2} = 1 - [(1-R2)*(n-1)/(n-k-1)]$$

where:

R2: The R2 of the model n: The number of observations k: The number of predictor variables

```
In [51]: 1 #display adjusted R-squared
2 1 - (1-score)*(len(y_test)-1)/(len(y_test)-X_test.shape[1]-1)
```

Out[51]: 0.6840226584639311

## New Data Prediction

```
In [52]: 1 boston.data[0].reshape(1,-1)
```

Out[52]: array([[6.320e-03, 1.800e+01, 2.310e+00, 0.000e+00, 5.380e-01, 6.575e+00, 6.520e+01, 4.090e+00, 1.000e+00, 2.960e+02, 1.530e+01, 3.969e+02, 4.980e+00]])

```
In [53]: 1 ##transformation of new data
2 scaler.transform(boston.data[0].reshape(1,-1))
```

Out[53]: array([[-0.41709233, 0.29216419, -1.27338003, -0.28154625, -0.16513629, 0.34715902, -0.13030059, 0.15267164, -0.97798895, -0.66777595, -1.32142483, 0.42854113, -1.04769976]])

```
In [54]: 1 regression.predict(scaler.transform(boston.data[0].reshape(1,-1)))
```

Out[54]: array([30.08649576])

## Pickling The Model file For Deployment

```
In [55]: 1 import pickle
```

```
In [56]: 1 pickle.dump(regression,open('regmodel.pkl','wb'))
```

```
In [57]: 1 pickled_model=pickle.load(open('regmodel.pkl','rb'))
```

```
In [58]: 1 ## Prediction  
2 pickled_model.predict(scaler.transform(boston.data[0].reshape(1,-1)))
```

```
Out[58]: array([30.08649576])
```

```
In [ ]: 1
```