

Our Target here to check how the different features effect the pricing of houses in "Boston"

```
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import seaborn as sns
import warnings
warnings.filterwarnings(action='ignore', category=FutureWarning)
```

load boston data set

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

In [3]: boston = load_boston()

In [4]: type(boston)

Out[4]: sklearn.utils.Bunch

In [5]: boston.keys()

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

check description of dataset

```
In [6]: print(boston.DESCR)

.. _bostondataset:

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
   - INDUS     proportion of non-retail business acres per town
   - 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         lower status of the population
   - LSTAT     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 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, Pt.1-02, 1978. Used in Belsley, Kuh & Welsch, "Regression diagnostics
...", Wiley, 1988. N.B. Various attributes are missing in the table on
rows 244-243 of the latter.

The Boston house-price data has been used in many machine learning papers that address regression
problems:
.. top101:References
- Belsley, Kuh & Welsch, "Regression diagnostics: Identifying Influential Data and Sources of Collinearity", Wiley, 1980. 244-243.
- Quinlan R. (1993). Combining Instance-Based and Model-Based Methods. In Proceedings on the Tenth International Conference of Machine Learning, 236-243,
University of Massachusetts, Amherst. Morgan Kaufmann.
```

```
In [7]: print(boston.data)

[[6.2206e+01 1.8800e+01 2.1308e+00 ... 1.5380e+01 3.9690e+02 4.9800e+00]
 [2.1310e+02 8.0800e+00 7.0700e+00 ... 1.7800e+01 3.9690e+02 9.4400e+00]
 [2.7290e+02 9.0800e+00 7.0700e+00 ... 1.7800e+01 3.9283e+02 4.6300e+00]
 ...
 [0.6700e+02 0.0800e+00 1.1930e+01 ... 2.1080e+01 3.9690e+02 5.4400e+00]
 [1.0950e+01 0.0800e+00 1.1930e+01 ... 2.1080e+01 3.9345e+02 6.4800e+00]
 [4.7430e+02 0.0800e+00 1.1930e+01 ... 2.1080e+01 3.9690e+02 7.4800e+00]]

In [8]: print(boston.target) # which is price of house.

[24. 21.6 34.7 33.4 36.2 28.7 22.9 27.1 16.5 18.9 15. 18.6 21.7 20.4
 18.2 18.9 23.1 17.5 20.2 18.2 23.6 19.8 15.2 14.5 15.6 23.9 16.9 14.8
 18.4 21. 12.7 14.5 13.2 13.1 13.5 18.9 20. 21. 24.7 28.8 34.9 26.6
 24.7 21.2 21.2 19.2 19. 16.6 21.4 19.6 18.9 20. 25. 23.1 18.9 25.
 24.7 31.6 23.3 19.6 18.7 16. 22.2 25. 33. 23.5 19.0 22. 17.4 20.9
 24.1 21.7 22.8 23.4 24.1 21.4 28. 20.8 21.2 26.3 28. 23.9 24.8 22.9
 24.6 26.6 22.5 22.2 23.6 28.7 22.6 22. 22.9 25. 20.6 20.4 21.4 20.7
 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
 48.7 18.5 18.3 21.2 19.2 20.4 19.3 22. 20.2 20.6 17.5 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 6.11 8.3 15.6 14.6 17.8 15.4 21.5 19.5 18.6 15.3 19.4
 37. 15.6 13.1 41.3 24.3 23.3 27. 50. 50. 50. 22.7 25. 16.9 23.6
 23.8 22.3 17.4 10.1 23.1 23.6 22.6 24.8 23.2 24.6 29.9 37.2 39.8 36.2
 27. 13.6 20.4 29.8 60. 32. 28. 24.9 37. 30.6 36. 29.6 21.9 18.1 36.
 33.8 36.3 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. 28.7 21.7 27.5 50.
 44.8 36.6 37.6 31.4 46.7 36.1 51.3 33.7 41.9 48.3 29.9 24. 26.1 33.5
 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. 47.8 21.9 20.9 44. 59. 36. 30.1 33.8 43.1 48.8 31. 36.5 22.8
 98.7 50. 43.5 20.7 21.1 25.2 24.4 35.2 32.4 32. 33.2 23.1 29.1 35.1
 45. 43.4 46. 50. 32.2 22. 28.1 23.2 23.2 24.8 28.5 37.3 27.9 23.9
 28.7 20.7 21.2 20.9 32.5 27.7 29. 19.8 20.6 21.2 19.1 28.6 15.2 7.
 22.8 28.3 16.1 12.1 19.4 21.6 23.8 16.2 17.8 19.8 23.1 21. 23.8 23.1
 20.8 15.25. 24.6 23. 22.2 19.3 22.6 19.8 17.1 19.4 22.2 20.7 21.1
 22.9 15.5 20.6 19. 18.7 37.7 16.5 23.9 33.2 17.5 17.2 23.1 24.5 20.6
 22.9 24.1 18.6 30.1 18.2 28.6 17.8 21.7 22.7 22.6 25. 19.9 20.8 16.8
 21.5 21.5 23.1 15. 19.6 50. 50. 50. 12.5 12.3 12.5. 13.9 3.5
 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.
 11.2 21.9 17.2 27.5 15. 17.2 21.9 16.3 7. 7.2 7.5 10.4 8.8 8.4
 18.7 14.2 20.8 13.4 11.7 8.3 18.2 18.9 11. 9.5 14.5 14.1 16.1 14.3
 22.9 24.1 18.6 30.1 18.2 28.6 17.8 21.7 22.7 22.6 25. 19.9 20.8 16.8
 14.1 13. 13.4 15.2 16.1 17.8 14.9 14.1 12.7 13.5 14.9 28. 16.4 17.7
 19.5 20.6 21.4 19.9 19. 19.1 19.1 20.1 10.9 19.6 23.2 28.6 13.8 13.3
 16. 13. 14.6 21.4 23. 37. 7.25. 21.8 20.6 21.2 19.1 28.6 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 28.6 23.9
 22. 11.9]
```

preparing our data set

```
In [9]: dataset = pd.DataFrame(boston.data,columns =boston.feature_names) #to add the column name we use boston.feature_names

In [10]: dataset.head()

Out[10]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT
0	0.06632	18.0	2.31	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	0.469	6.421	78.9	4.9671	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33

```
In [11]: dataset['price']=boston.target

In [12]: dataset.head()

Out[12]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
0	0.06632	18.0	2.31	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	0.469	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	0.469	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	0.458	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	0.458	7.147	54.2	6.0622	3.0	222.0	18.7	396.90	5.33	36.2

```
In [13]: 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: 95.5 KB
```

summarizing stats of data

```
In [14]: dataset.describe()

Out[14]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
count	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000	506.000000
mean	3.612634	11.363636	11.136779	0.069170	0.554965	6.294034	68.174901	3.769543	8.549407	406.23154	18.459534	356.67022	12.653063	22.538606
std	8.601545	23.322663	6.660353	0.253994	0.115878	0.702617	28.148861	2.106710	8.707259	166.537116	2.164946	91.249684	7.141062	1.910704
min	0.006200	0.000000	0.040000	0.000000	0.399000	3.561000	2.900000	1.126900	1.000000	187.000000	12.600000	0.320000	1.730000	5.000000
25%	0.037600	0.000000	5.190000	0.000000	0.449000	5.895000	4.050000	2.100175	5.000000	279.000000	17.400000	331.000000	6.950000	17.025000
50%	0.056100	0.000000	9.690000	0.000000	0.538000	6.236000	7.050000	3.207405	5.000000	330.000000	19.400000	391.400000	11.300000	21.200000
75%	3.677033	12.500000	18.100000	0.000000	0.624000	6.623500	94.075000	5.188425	24.000000	666.000000	20.200000	982.200000	16.950000	25.000000
max	88.976200	100.000000	27.740000	1.000000	0.871000	8.780000	100.000000	12.126500	24.000000	711.000000	22.000000	986.900000	37.970000	50.000000

check the missing values

```
In [15]: 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
```

Exploratory data analysis

Correlation

```
In [16]: dataset.corr()

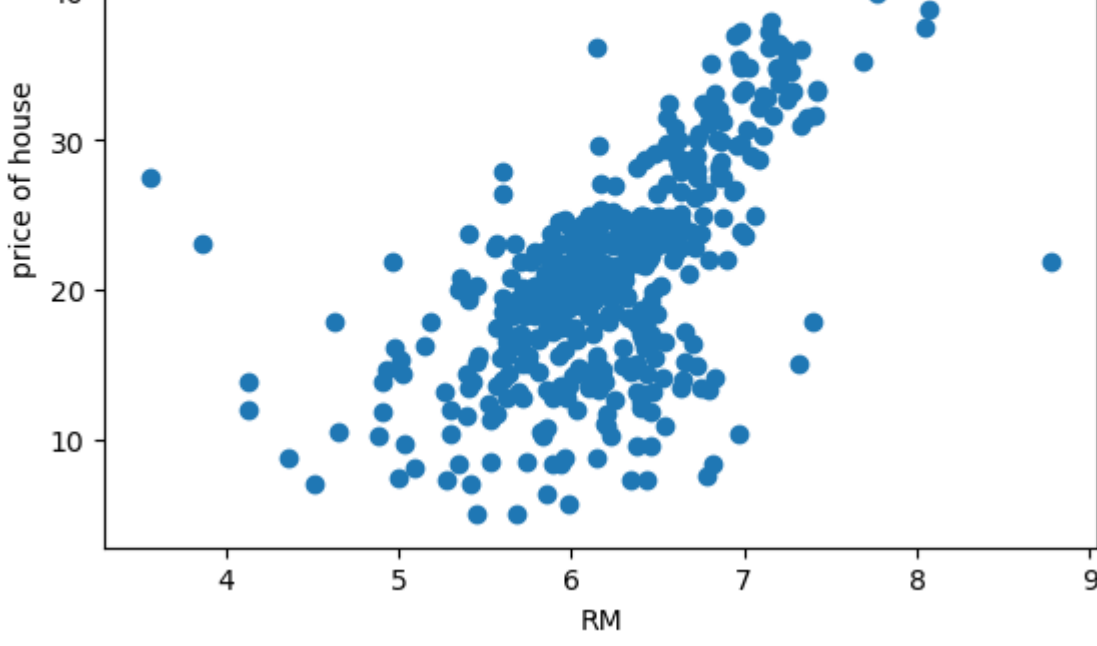
Out[16]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
CRIM	1.000000	-0.200469	0.406583	-0.055892	0.420972	-0.219247	0.352734	-0.379470	0.625505	0.582764	0.289946	-0.385064	0.455621	-0.388405
ZN	-0.200469	1.000000	-0.533828	0.062938	0.163691	-0.393176	0.644778	-0.108027	0.595129	0.720780	0.363248	-0.358977	-0.437280	-0.187589
INDUS	0.406583	-0.533828	1.000000	0.062938	0.091203	-0.092525	0.094158	0.094178	-0.007388	0.038587	-0.121315	0.048788	-0.059559	0.175280
CHAS	-0.055892	-0.062938	0.062938	0.991000	0.000000	-0.362000	0.732470	-0.399200	0.614441	0.008023	0.089023	0.380951	0.908070	-0.423721
NOX	0.420972	0.163691	0.091203	0.000000	1.000000	-0.240205	0.705246	-0.209847	0.200468	-0.355501	0.130669	-0.162080	0.095360	-0.085360
RM	-0.219247	-0.393176	-0.092525	-0.362000	-0.240205	1.000000	-0.460622	-0.494688	0.500000	0.910228	0.464741	0.444413	0.486761	-0.481636
AGE	0.352734	0.644778	0.094158	0.732470	0.705246	-0.460622	1.000000	-0.147881	0.456022	0.505456	0.261515	0.273534	0.052330	-0.378905
DIS	-0.379470	-0.108027	-0.007388	-0.099176	-0.209847	-0.494688	-0.147881	1.000000	-0.494688	0.534432	-0.232471	0.291512	-0.069990	0.249959
RAD	0.625505	-0.311948	0.500000	-0.007388	0.611441	-0.209847	0.456022	-0.494688	1.000000	0.910228	0.464741	0.444413	0.486761	-0.481636
TAX	0.582764	0.720780	0.038587	-0.007388	0.668023	-0.365501	0.273534	0.291512	-0.444413	-0.441908	0.534432	-0.232471	0.291512	-0.069990
PTRATIO	0.289946	-0.358977	0.363248	-0.121515	0.138853	-0.355501	0.061515	-0.232471	0.464741	0.462953	1.000000	0.177383	0.374044	-0.507787
B	-0.385064	0.455621	-0.358977	0.048788	-0.380951	0.129009	-0.273534	0.291512	-0.444413	-0.441908	-0.177383	1.000000	-0.369007	0.730003
LSTAT	0.435661	-0.412095	0.053030	-0.053030	0.506879	-0.613808	0.602339	-0.486986	0.488676	0.543993	0.374044	-0.369007	1.000000	-0.737663
price	-0.388305	0.380445	-0.483725	0.176200	-0.427321	0.065300	-0.378905	0.248920	-0.381626	-0.460536	-0.507787	0.333401	-0.737663	1.000000

analyse relationship b/w few of our fetures with target column

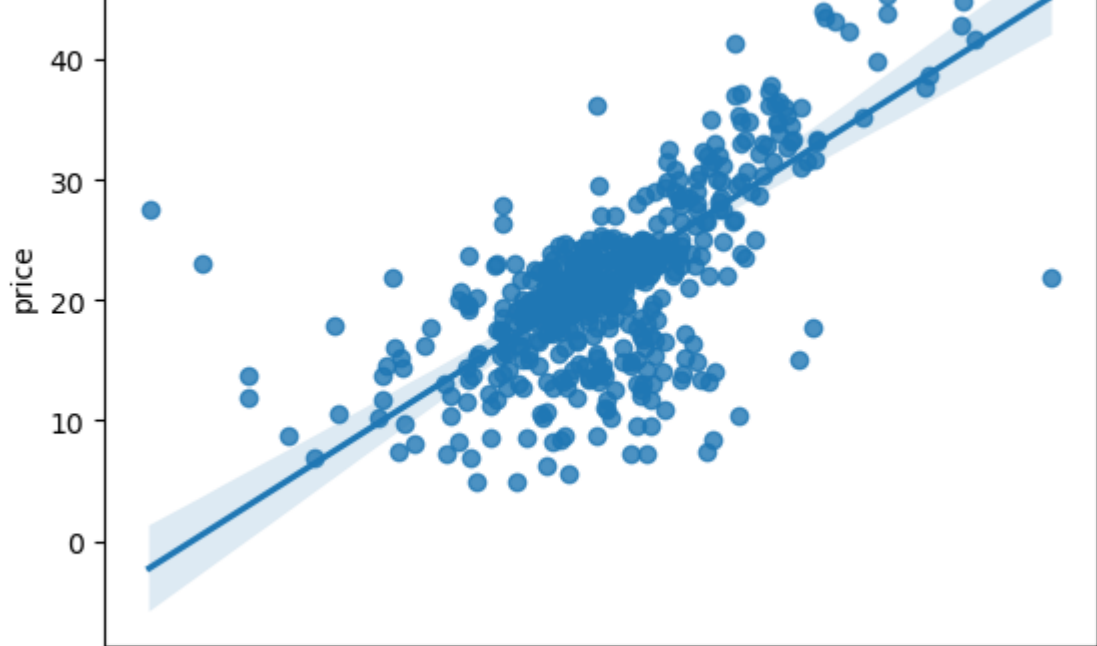
```
In [17]: plt.scatter(dataset['CRIM'],dataset['price'])
plt.xlabel('CRIM')
plt.ylabel('price of house') # AS WE CAN SEE AS CRIMINATE INCREASE PRICE OF HOUSE WILL DECREASE

Out[17]: Text(0, 0.5, 'price of house')
```



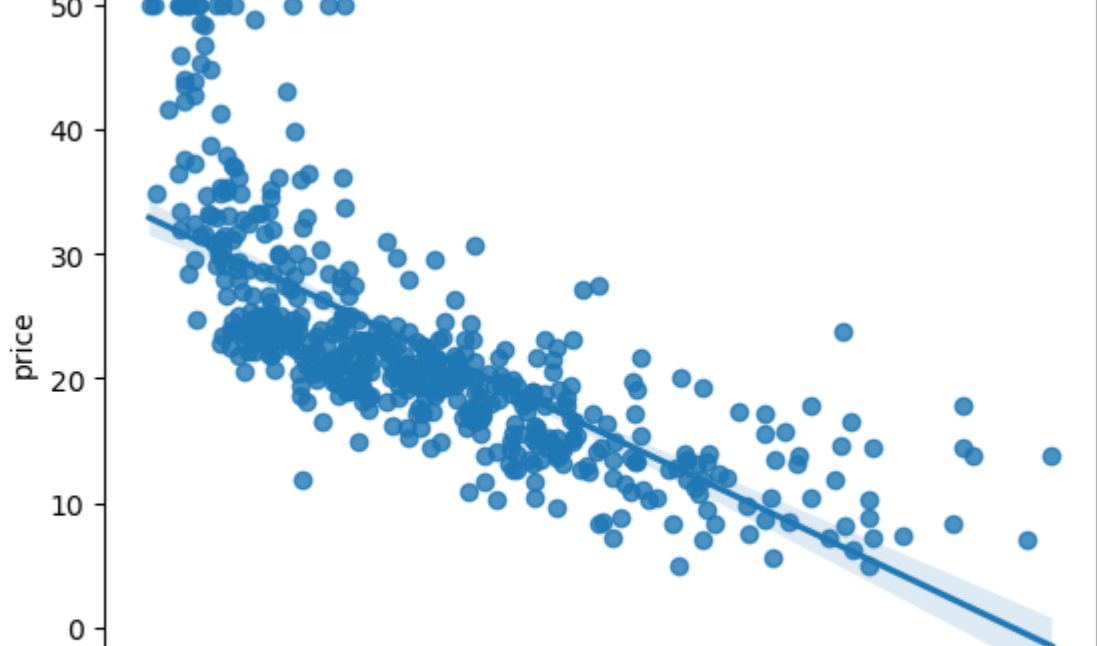
```
In [18]: plt.scatter(dataset['RM'],dataset['price'])
plt.xlabel('RM')
plt.ylabel('price of house') # AS WE CAN SEE AS CRIMINATE INCREASE PRICE OF HOUSE WILL DECREASE

Out[18]: Text(0, 0.5, 'price of house')
```



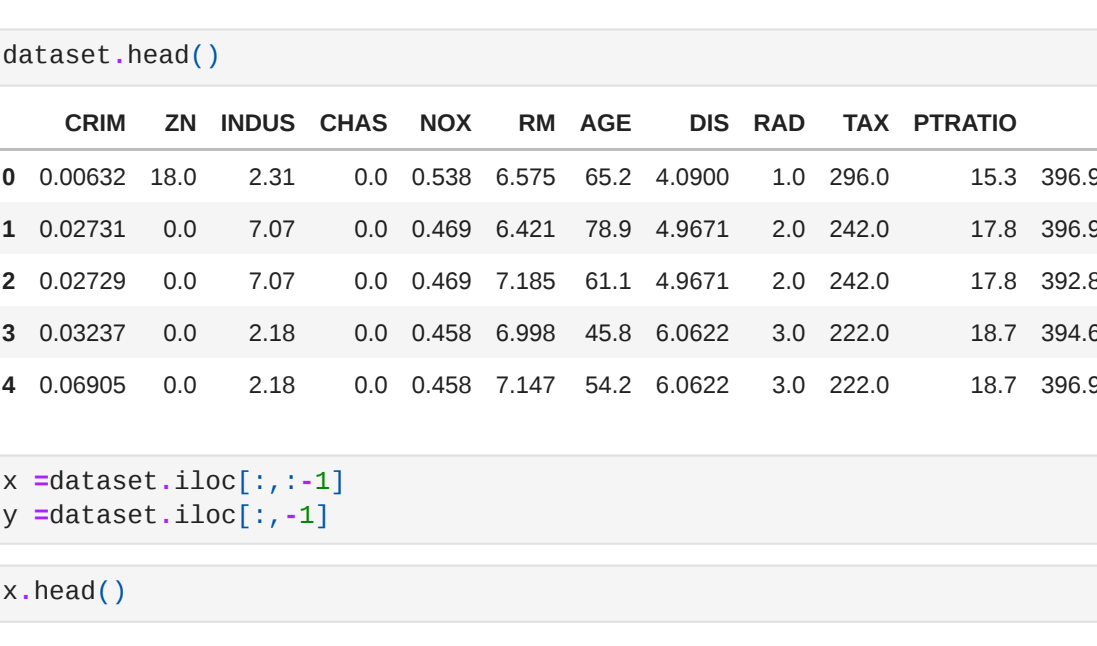
```
In [19]: sns.regplot(x='RM',y='price',data=dataset) #HERE WE CAN SEE POSITIVE CORRELATION BETWEEN AVG ROOM AND PRICE

Out[19]: <AxesSubplot:xlabel='RM', ylabel='price'>
```



```
In [20]: sns.regplot(x='LSTAT',y='price',data=dataset) #HERE WE HAVE NEGATIVE RELATION

Out[20]: <AxesSubplot:xlabel='LSTAT', ylabel='price'>
```



Now its time to split our dataset into dependent and independent features

```
In [21]: dataset.head()

Out[21]:
```

	CRIM	ZN	INDUS	CHAS	NOX	RM	AGE	DIS	RAD	TAX	PTRATIO	B	LSTAT	price
0	0.06632	18.0	2.31	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	0.469	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	0.469	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	0.458	6.998	45.8	6.0622	3.0	222.0	18.7	394.63	2.94	33.4
4	0.06905	0.0	2.											