# boston house prices

March 5, 2022

```
[1]: import pandas as pd
     import numpy as np
     from sklearn import linear_model
     import sklearn
     from sklearn.utils import
     import matplotlib.pyplot as plt
     import seaborn as sns
     import warnings
     import random
     import os
```

```
[2]: #To create machine learning models easily and make predictions.
     from sklearn.datasets import load_boston
     dataset = load_boston()
```

```
[3]: print("[INFO] keys : {}".format(dataset.keys()))
```

```
[INFO] keys : dict_keys(['data', 'target', 'feature_names', 'DESCR',
'filename'])
```

```
[4]: dataset.data
```

```
[4]: 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]])
```

```
[5]: dataset.target
```

```
[5]: 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])
```

```
[6]: dataset.feature_names
```

```
[6]: array(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD',
             'TAX', 'PTRATIO', 'B', 'LSTAT'], dtype='<U7')
```

```
[7]: print("[INFO] dataset summary", dataset.DESCR)
```

[INFO] dataset summary .. _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/


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.

```
[8]: dataset.filename
```

```
[8]: 'C:\\Users\\nomaniqbal\\anaconda3\\lib\\site-
     packages\\sklearn\\datasets\\data\\boston_house_prices.csv'
```

```
[9]: df=pd.DataFrame(dataset.data)
     df
```

[9]:

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 \ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.00632 | 18.0 | 2.31 | 0.0 | 0.538 | 6.575 | 65.2 | 4.0900 | 1.0 | 296.0 | 15.3 |
| 1 | 0.02731 | 0.0 | 7.07 | 0.0 | 0.469 | 6.421 | 78.9 | 4.9671 | 2.0 | 242.0 | 17.8 |
| 2 | 0.02729 | 0.0 | 7.07 | 0.0 | 0.469 | 7.185 | 61.1 | 4.9671 | 2.0 | 242.0 | 17.8 |
| 3 | 0.03237 | 0.0 | 2.18 | 0.0 | 0.458 | 6.998 | 45.8 | 6.0622 | 3.0 | 222.0 | 18.7 |
| 4 | 0.06905 | 0.0 | 2.18 | 0.0 | 0.458 | 7.147 | 54.2 | 6.0622 | 3.0 | 222.0 | 18.7 |
| .. | … | … | … | … | … | … | … | … | … | … | … |
| 501 | 0.06263 | 0.0 | 11.93 | 0.0 | 0.573 | 6.593 | 69.1 | 2.4786 | 1.0 | 273.0 | 21.0 |
| 502 | 0.04527 | 0.0 | 11.93 | 0.0 | 0.573 | 6.120 | 76.7 | 2.2875 | 1.0 | 273.0 | 21.0 |
| 503 | 0.06076 | 0.0 | 11.93 | 0.0 | 0.573 | 6.976 | 91.0 | 2.1675 | 1.0 | 273.0 | 21.0 |
| 504 | 0.10959 | 0.0 | 11.93 | 0.0 | 0.573 | 6.794 | 89.3 | 2.3889 | 1.0 | 273.0 | 21.0 |
| 505 | 0.04741 | 0.0 | 11.93 | 0.0 | 0.573 | 6.030 | 80.8 | 2.5050 | 1.0 | 273.0 | 21.0 |

| | 11 | 12 |
|---|---|---|
| 0 | 396.90 | 4.98 |
| 1 | 396.90 | 9.14 |
| 2 | 392.83 | 4.03 |
| 3 | 394.63 | 2.94 |

```
4     396.90  5.33
..       …    …
501   391.99  9.67
502   396.90  9.08
503   396.90  5.64
504   393.45  6.48
505   396.90  7.88

[506 rows x 13 columns]
```

[10]: `df.columns = dataset.feature_names`

[11]: `df`

[11]:
```
          CRIM    ZN  INDUS  CHAS    NOX     RM   AGE     DIS  RAD    TAX  \
0      0.00632  18.0   2.31   0.0  0.538  6.575  65.2  4.0900  1.0  296.0
1      0.02731   0.0   7.07   0.0  0.469  6.421  78.9  4.9671  2.0  242.0
2      0.02729   0.0   7.07   0.0  0.469  7.185  61.1  4.9671  2.0  242.0
3      0.03237   0.0   2.18   0.0  0.458  6.998  45.8  6.0622  3.0  222.0
4      0.06905   0.0   2.18   0.0  0.458  7.147  54.2  6.0622  3.0  222.0
..         …     …      …     …      …      …     …       …    …      …
501    0.06263   0.0  11.93   0.0  0.573  6.593  69.1  2.4786  1.0  273.0
502    0.04527   0.0  11.93   0.0  0.573  6.120  76.7  2.2875  1.0  273.0
503    0.06076   0.0  11.93   0.0  0.573  6.976  91.0  2.1675  1.0  273.0
504    0.10959   0.0  11.93   0.0  0.573  6.794  89.3  2.3889  1.0  273.0
505    0.04741   0.0  11.93   0.0  0.573  6.030  80.8  2.5050  1.0  273.0

     PTRATIO       B  LSTAT
0       15.3  396.90   4.98
1       17.8  396.90   9.14
2       17.8  392.83   4.03
3       18.7  394.63   2.94
4       18.7  396.90   5.33
..        …      …      …
501     21.0  391.99   9.67
502     21.0  396.90   9.08
503     21.0  396.90   5.64
504     21.0  393.45   6.48
505     21.0  396.90   7.88

[506 rows x 13 columns]
```

[12]: `df["prices"]=dataset.target`

[13]: `df`

[13]:
```
         CRIM    ZN  INDUS  CHAS    NOX     RM   AGE     DIS  RAD    TAX  \
0     0.00632  18.0   2.31   0.0  0.538  6.575  65.2  4.0900  1.0  296.0
1     0.02731   0.0   7.07   0.0  0.469  6.421  78.9  4.9671  2.0  242.0
2     0.02729   0.0   7.07   0.0  0.469  7.185  61.1  4.9671  2.0  242.0
3     0.03237   0.0   2.18   0.0  0.458  6.998  45.8  6.0622  3.0  222.0
4     0.06905   0.0   2.18   0.0  0.458  7.147  54.2  6.0622  3.0  222.0
..        ...   ...    ...   ...    ...    ...   ...     ...  ...    ...
501   0.06263   0.0  11.93   0.0  0.573  6.593  69.1  2.4786  1.0  273.0
502   0.04527   0.0  11.93   0.0  0.573  6.120  76.7  2.2875  1.0  273.0
503   0.06076   0.0  11.93   0.0  0.573  6.976  91.0  2.1675  1.0  273.0
504   0.10959   0.0  11.93   0.0  0.573  6.794  89.3  2.3889  1.0  273.0
505   0.04741   0.0  11.93   0.0  0.573  6.030  80.8  2.5050  1.0  273.0

     PTRATIO       B  LSTAT  prices
0       15.3  396.90   4.98    24.0
1       17.8  396.90   9.14    21.6
2       17.8  392.83   4.03    34.7
3       18.7  394.63   2.94    33.4
4       18.7  396.90   5.33    36.2
..       ...     ...    ...     ...
501     21.0  391.99   9.67    22.4
502     21.0  396.90   9.08    20.6
503     21.0  396.90   5.64    23.9
504     21.0  393.45   6.48    22.0
505     21.0  396.90   7.88    11.9

[506 rows x 14 columns]
```

[14]: `df.describe()`

[14]:
```
             CRIM          ZN       INDUS        CHAS         NOX          RM  \
count  506.000000  506.000000  506.000000  506.000000  506.000000  506.000000
mean     3.613524   11.363636   11.136779    0.069170    0.554695    6.284634
std      8.601545   23.322453    6.860353    0.253994    0.115878    0.702617
min      0.006320    0.000000    0.460000    0.000000    0.385000    3.561000
25%      0.082045    0.000000    5.190000    0.000000    0.449000    5.885500
50%      0.256510    0.000000    9.690000    0.000000    0.538000    6.208500
75%      3.677083   12.500000   18.100000    0.000000    0.624000    6.623500
max     88.976200  100.000000   27.740000    1.000000    0.871000    8.780000

              AGE         DIS         RAD         TAX     PTRATIO           B  \
count  506.000000  506.000000  506.000000  506.000000  506.000000  506.000000
mean    68.574901    3.795043    9.549407  408.237154   18.455534  356.674032
std     28.148861    2.105710    8.707259  168.537116    2.164946   91.294864
min      2.900000    1.129600    1.000000  187.000000   12.600000    0.320000
25%     45.025000    2.100175    4.000000  279.000000   17.400000  375.377500
50%     77.500000    3.207450    5.000000  330.000000   19.050000  391.440000
```

```
75%      94.075000      5.188425    24.000000  666.000000    20.200000  396.225000
max     100.000000     12.126500    24.000000  711.000000    22.000000  396.900000

             LSTAT       prices
count   506.000000  506.000000
mean     12.653063   22.532806
std       7.141062    9.197104
min       1.730000    5.000000
25%       6.950000   17.025000
50%      11.360000   21.200000
75%      16.955000   25.000000
max      37.970000   50.000000
```

[15]: `df.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  prices   506 non-null    float64
dtypes: float64(14)
memory usage: 55.5 KB
```

[16]: `df.prices`

```
[16]: 0      24.0
      1      21.6
      2      34.7
      3      33.4
      4      36.2
             …
      501    22.4
      502    20.6
```
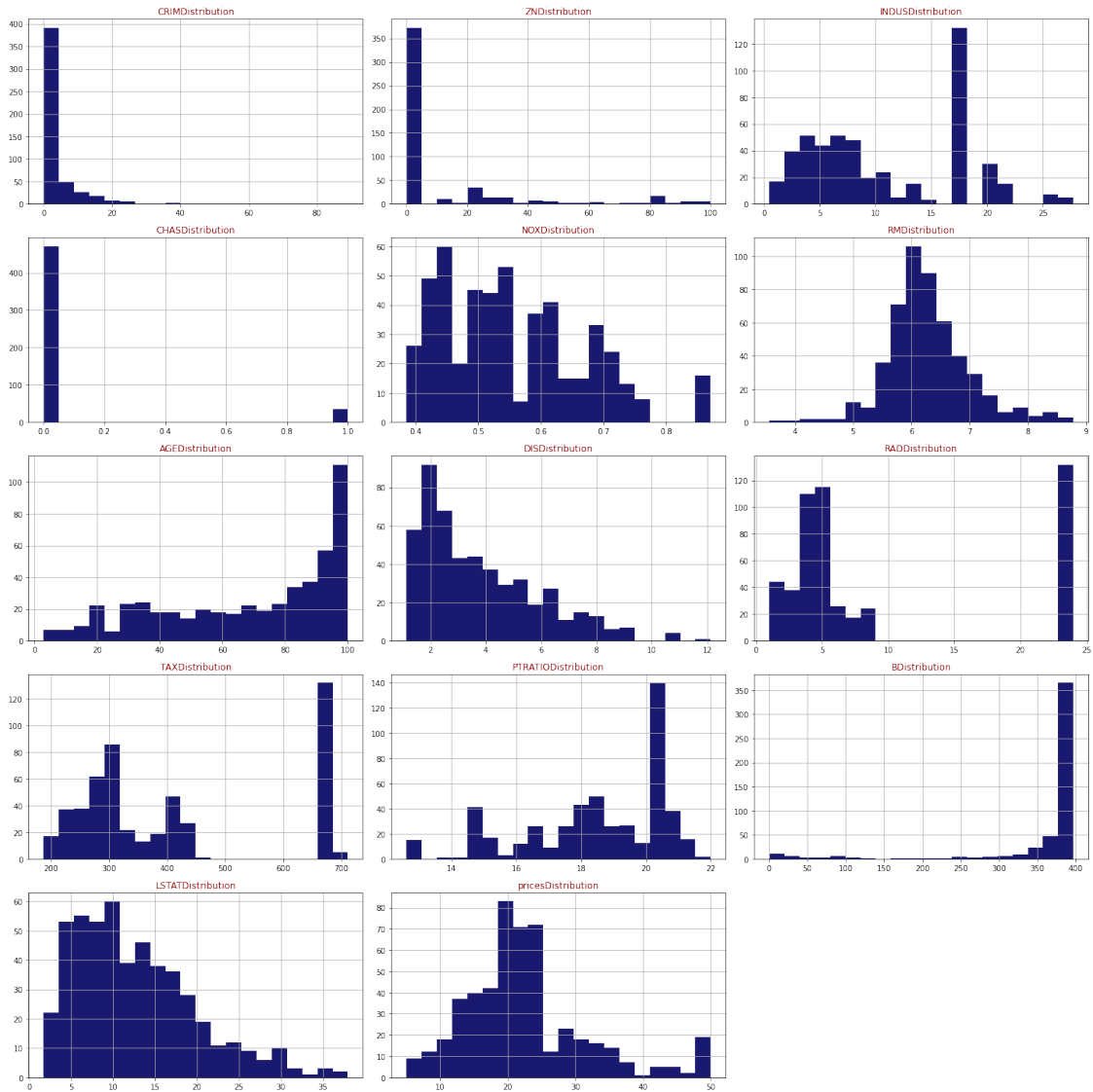
```
503    23.9
504    22.0
505    11.9
Name: prices, Length: 506, dtype: float64
```

[17]: 
```python
import matplotlib.pyplot as plt
from sklearn.linear_model import LinearRegression
```

[18]: 
```python
def draw_plots(df, var, rows, cols):
    fig=plt.figure(figsize=(20,20))
    for i, f in enumerate(var):
        ax=fig.add_subplot(rows,cols,i+1)
        df[f].hist(bins=20,ax=ax, facecolor='midnightblue')
        ax.set_title(f+'Distribution',color='DarkRed')

    fig.tight_layout()
```
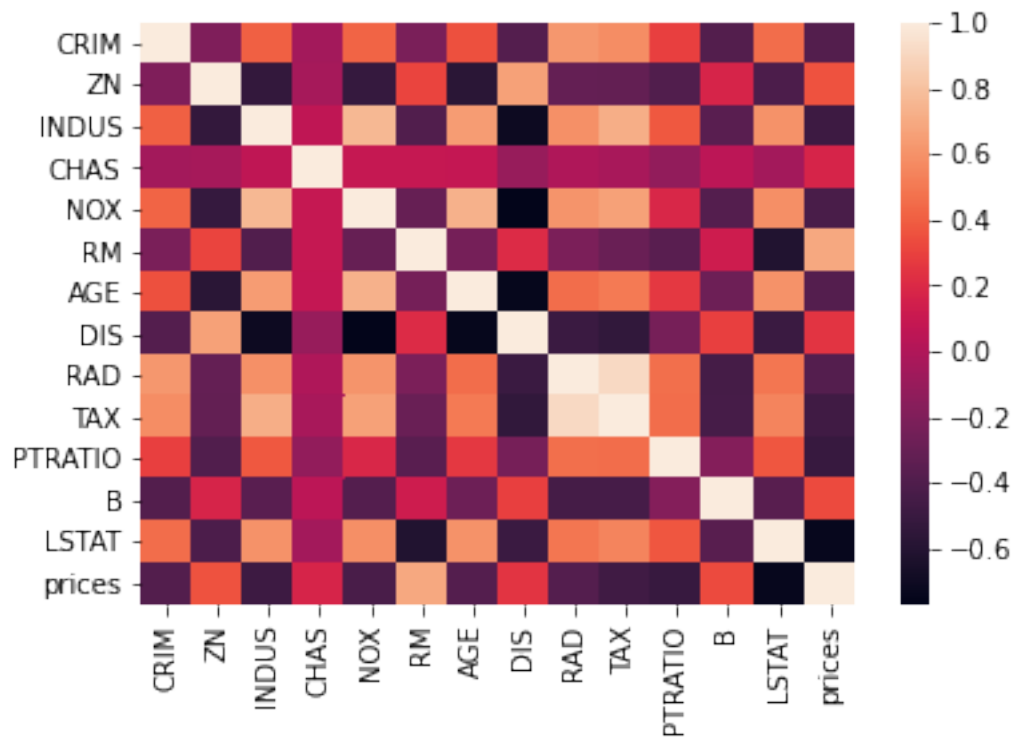
[19]: 
```python
plt.show()
draw_plots(df,df.columns,5,3)
```
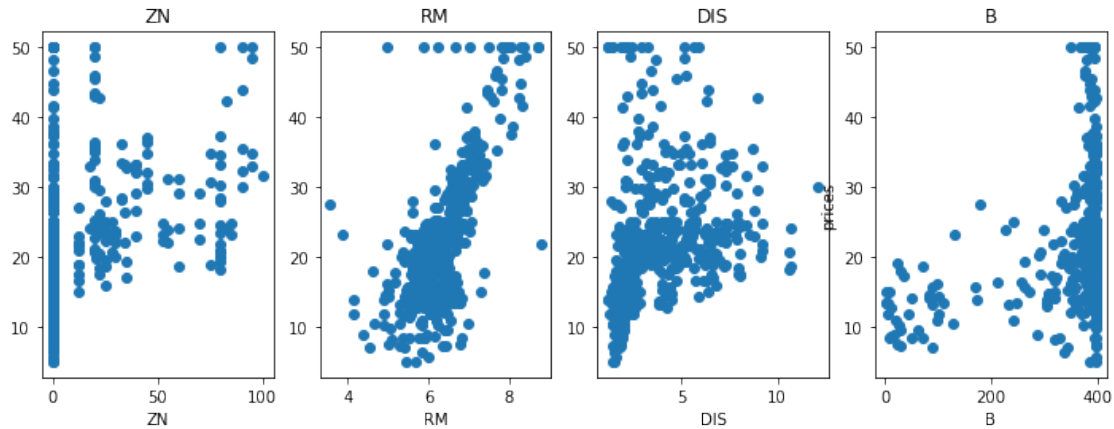
```
[26]: sns.heatmap(df.corr())
```

```
[26]: <AxesSubplot:>
```

```
[33]: X= df[['ZN', 'RM','DIS','B']]
      Y= df['prices']
      plt.figure(figsize=(12, 4))
      predictors = ['ZN', 'RM','DIS','B']
      target = df['prices']
      for i, col in enumerate(predictors):
          plt.subplot(1, len(predictors) , i+1)
          x = df[col]
          y = target
          plt.scatter(x, y, marker='o')
          plt.title(col)
          plt.xlabel(col)
      plt.ylabel('prices')
```

[33]: Text(0, 0.5, 'prices')

```
[40]: #accuracy
      from sklearn.model_selection import train_test_split
      from sklearn.metrics import r2_score
      x_train,x_test,y_train,y_test = train_test_split(X,Y,test_size = 0.3)
      reg = LinearRegression()
      reg.fit(x_train,y_train)
      y_pred = reg.predict(x_test)
      r2_score(y_test,y_pred)
```

[40]: 0.5393900127100043

[41]: df

[41]:
```
        CRIM    ZN  INDUS  CHAS    NOX     RM   AGE     DIS  RAD    TAX  \
0    0.00632  18.0   2.31   0.0  0.538  6.575  65.2  4.0900  1.0  296.0
1    0.02731   0.0   7.07   0.0  0.469  6.421  78.9  4.9671  2.0  242.0
2    0.02729   0.0   7.07   0.0  0.469  7.185  61.1  4.9671  2.0  242.0
3    0.03237   0.0   2.18   0.0  0.458  6.998  45.8  6.0622  3.0  222.0
4    0.06905   0.0   2.18   0.0  0.458  7.147  54.2  6.0622  3.0  222.0
..       ...   ...    ...   ...    ...    ...   ...     ...  ...    ...
501  0.06263   0.0  11.93   0.0  0.573  6.593  69.1  2.4786  1.0  273.0
502  0.04527   0.0  11.93   0.0  0.573  6.120  76.7  2.2875  1.0  273.0
503  0.06076   0.0  11.93   0.0  0.573  6.976  91.0  2.1675  1.0  273.0
504  0.10959   0.0  11.93   0.0  0.573  6.794  89.3  2.3889  1.0  273.0
505  0.04741   0.0  11.93   0.0  0.573  6.030  80.8  2.5050  1.0  273.0

     PTRATIO       B  LSTAT  prices
0       15.3  396.90   4.98    24.0
1       17.8  396.90   9.14    21.6
2       17.8  392.83   4.03    34.7
3       18.7  394.63   2.94    33.4
4       18.7  396.90   5.33    36.2
```

```
..       …      …     …      …
501     21.0  391.99   9.67    22.4
502     21.0  396.90   9.08    20.6
503     21.0  396.90   5.64    23.9
504     21.0  393.45   6.48    22.0
505     21.0  396.90   7.88    11.9

[506 rows x 14 columns]
```

[45]: `print(df['prices'].max())`

```
50.0
```

[ ]: