```python
import pandas as pd
import seaborn as sns
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```python
from sklearn.datasets import load_boston
boston = load_boston()
print(boston.DESCR)
```

```
Boston House Prices dataset
===========================

Notes
------
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.
http://archive.ics.uci.edu/ml/datasets/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.

**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.
   - many more! (see http://archive.ics.uci.edu/ml/datasets/Housing)
```

```python
boston_df = boston.data
```

```python
boston.keys()
```

```
dict_keys(['data', 'target', 'feature_names', 'DESCR'])
```

```python
boston.feature_names
```

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

```python
df= pd.DataFrame(data= boston.data , columns=boston.feature_names )
```

```python
df.head()
```

|   | 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 [ ]:    df['prices']=boston.target
```
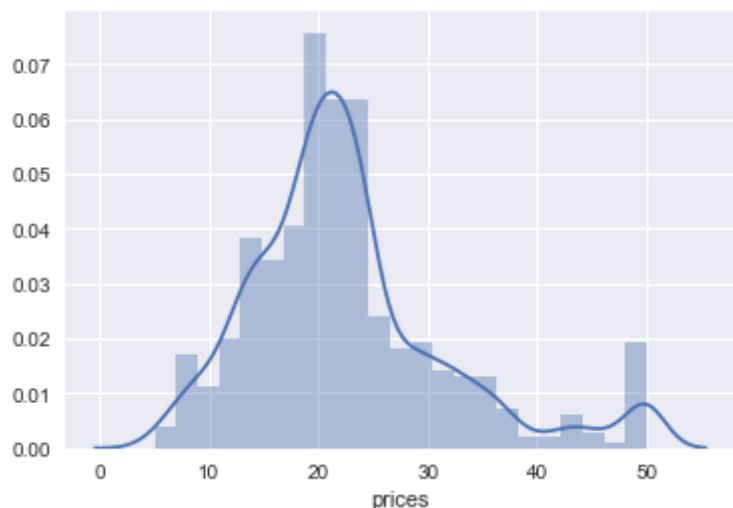
```
In [ ]:    df.head()
```

Out[ ]:

|   | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LSTAT | prices |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|-------|--------|
| 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 |

DATA VISUALIZATION

```
In [ ]:    sns.distplot(df['prices'])
```

Out[ ]:    <matplotlib.axes._subplots.AxesSubplot at 0x162d7920588>



TRAINING THE MODEL

```
In [ ]:    df.columns
```

Out[ ]:    Index(['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
                  'PTRATIO', 'B', 'LSTAT', 'prices'],
                 dtype='object')

```
In [ ]:    X=df[['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX',
                  'PTRATIO', 'B', 'LSTAT']]
```

```
In [ ]:    y=df['prices']
```

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

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

```
In [ ]:    from sklearn.linear_model import LinearRegression
```
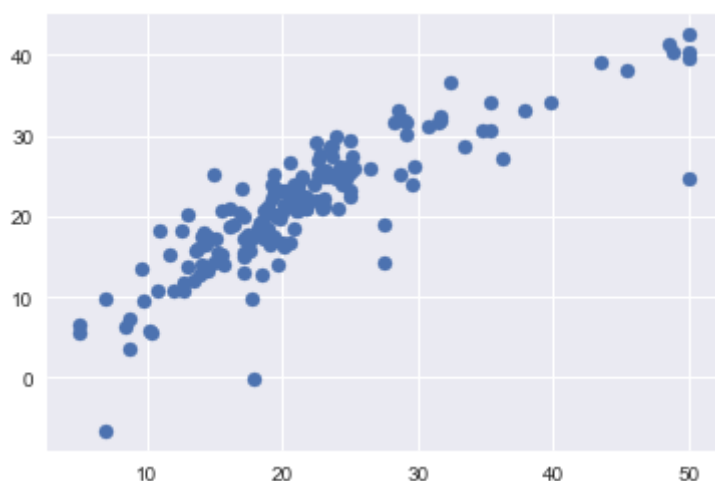
```
In [ ]:    lm=LinearRegression()
```

```
In [ ]:    lm.fit(X_train,y_train)
```

Out[ ]:    LinearRegression(copy_X=True, fit_intercept=True, n_jobs=1, normalize=False)

```
In [ ]:    predictions= lm.predict(X_test)
```

```
In [ ]:    plt.scatter(y_test,predictions)
```

Out[ ]:    <matplotlib.collections.PathCollection at 0x162df77d7f0>



Regression Evaluation Metrics

```
In [ ]:    from sklearn import metrics
```

```python
print('Mean_Absolute_Error : ' , metrics.mean_absolute_error(y_test,predictions))
```

Mean_Absolute_Error :   3.15128783659

```python
print('Mean_Squared_Error : ' , metrics.mean_squared_error(y_test,predictions) )
```

Mean_Squared_Error :   20.7471433603

```python
print('Root_Mean_Squared_Error : ', np.sqrt(metrics.mean_squared_error(y_test,predictions)))
```

Root_Mean_Squared_Error :   4.55490322184

```python
print('Mean_Absolute_Error : ' , metrics.mean_absolute_error(y_test,predictions))
```

Mean_Absolute_Error :   3.15128783659

```python
print('Mean_Squared_Error : ' , metrics.mean_squared_error(y_test,predictions) )
```

Mean_Squared_Error :   20.7471433603

```python
print('Root_Mean_Squared_Error : ', np.sqrt(metrics.mean_squared_error(y_test,predictions)))
```

Root_Mean_Squared_Error :   4.55490322184