```
1 import numpy as np
2 import pandas as pd
3 import sklearn
```

```
1 from sklearn.datasets import load_boston
2 df = load_boston()
```

```
1 df.keys() # returns all the keys of the dataset dictionary
```

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

```
1 print(df.DESCR)  # Info about the dataset
```

.. _boston_dataset:

Boston house prices dataset
---------------------------

**Data Set Characteristics:**

    :Number of Instances: 506

    :Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14

    :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 U
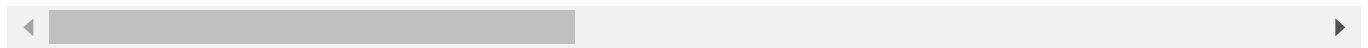
The Boston house-price data of Harrison, D. and Rubinfeld, D.L. 'Hedonic
prices and the demand for clean air', J. Environ. Economics & Management,

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 r
problems.

.. topic:: References

   - Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sc
   - Quinlan,R. (1993). Combining Instance-Based and Model-Based Learning. In Proceeding

```
1 boston=pd.DataFrame(df.data,columns=df.feature_names)
2 boston.head()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LST/ |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|------|
| 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 | 4.9 |
| 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 | 9.1 |
| 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 | 4.0 |
| 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 | 2.9 |
| 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 | 5.3 |

```
1 boston['MEDV']=df.target
2 boston.head()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LST/ |
|---|------|-----|-------|------|-------|-------|------|--------|-----|-------|---------|--------|------|
| 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 | 4.9 |
| 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 | 9.1 |
| 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 | 4.0 |
| 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 | 2.9 |
| 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 | 5.3 |

```
1 boston.isnull()
```

| | CRIM | ZN | INDUS | CHAS | NOX | RM | AGE | DIS | RAD | TAX | PTRATIO | B | LS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| 1 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| 2 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| 3 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| 4 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | |
| 501 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |
| 502 | False | False | False | False | False | False | False | False | False | False | False | False | Fa |

```
1 boston.isnull().sum()
```

```
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
MEDV        0
dtype: int64
```

```
1 from sklearn.model_selection import train_test_split
2 X=boston.drop('MEDV',axis=1)
3 Y=boston['MEDV']
4 X_train,X_test,Y_train,Y_test=train_test_split(X,Y,test_size=0.15, random_state=5)
5 print(X_train.shape)
6 print(X_test.shape)
7 print(Y_train.shape)
8 print(Y_test.shape)
```

```
(430, 13)
(76, 13)
(430,)
(76,)
```

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error
3
4
5
```

```
1 ## FITTING MODEL ON THE TRAINING DATASET
2 lin_model = LinearRegression()
3 lin_model.fit(X_train,Y_train)
```

```
LinearRegression(copy_X=True, fit_intercept=True, n_jobs=None, normalize=False)
```

```
 1 y_train_predict = lin_model.predict(X_train)
 2 rmse = (np.sqrt(mean_squared_error(Y_train, y_train_predict)))
 3
 4 print("The model performance for training set")
 5 print('RMSE is {}'.format(rmse))
 6 print("\n")
 7
 8 # on testing set
 9 y_test_predict = lin_model.predict(X_test)
10 rmse = (np.sqrt(mean_squared_error(Y_test,y_test_predict)))
11
12 print("The model performance for testing set")
13 print('RMSE is {}'.format(rmse))
```

```
The model performance for training set
RMSE is 4.710901797319796


The model performance for testing set
RMSE is 4.687543527902972
```