

```
In [1]: 1 from sklearn.datasets import load_boston
        2 boston = load_boston()
        3 import matplotlib.pyplot as plt
```

```
In [2]: 1 print(boston.data.shape)
```

```
(506, 13)
```

```
In [3]: 1 print(boston.feature_names)
```

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

```
In [4]: 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.  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 [5]: 1 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 (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>) (<http://archive.ics.uci.edu/ml/datasets/Housing>)

In [6]:

```
1 import pandas as pd
2 bos = pd.DataFrame(boston.data)
3 print(bos.head())
```

|   | 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 |   |

|   | 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 |

In [7]:

```
1 bos['PRICE'] = boston.target
2
3 X = bos.drop('PRICE', axis = 1)
4 Y = bos['PRICE']
```

In [ ]:

```
1
```

In [8]:

```
1 from sklearn.model_selection import train_test_split
2 X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.3, ra
3
4 print(X_train.shape)
5 print(X_test.shape)
6 print(Y_train.shape)
7 print(Y_test.shape)
```

(354, 13)

(152, 13)

(354,)

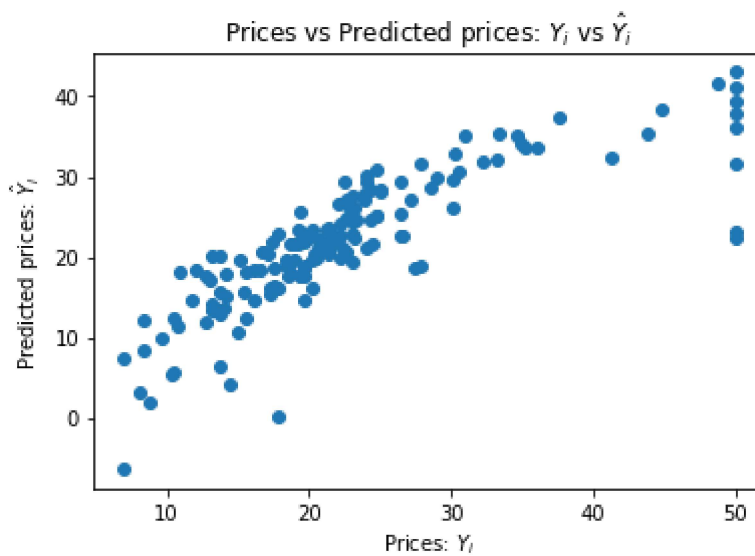
(152,)

In [9]:

```
1 # Standardizing the data
2
3 from sklearn.preprocessing import StandardScaler
4 scaler = StandardScaler()
5
6 X_train = scaler.fit_transform(X_train)
7 X_test = scaler.transform(X_test)
```

In [ ]: 1

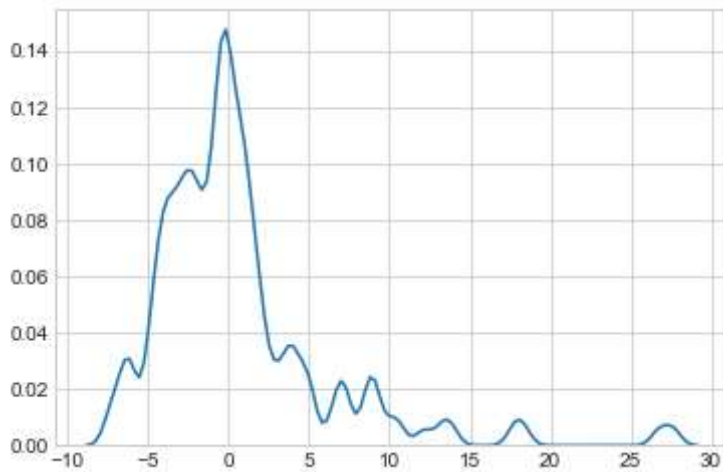
```
In [10]: 1 # code source:https://medium.com/@haydar_ai/Learning-data-science-day-9-LinearRegression
2 from sklearn.linear_model import LinearRegression
3
4 lm = LinearRegression()
5 lm.fit(X_train, Y_train)
6
7 Y_pred = lm.predict(X_test)
8
9 plt.scatter(Y_test, Y_pred)
10 plt.xlabel("Prices: $Y_i$")
11 plt.ylabel("Predicted prices: $\hat{Y}_i$")
12 plt.title("Prices vs Predicted prices: $Y_i$ vs $\hat{Y}_i$")
13 plt.show()
```



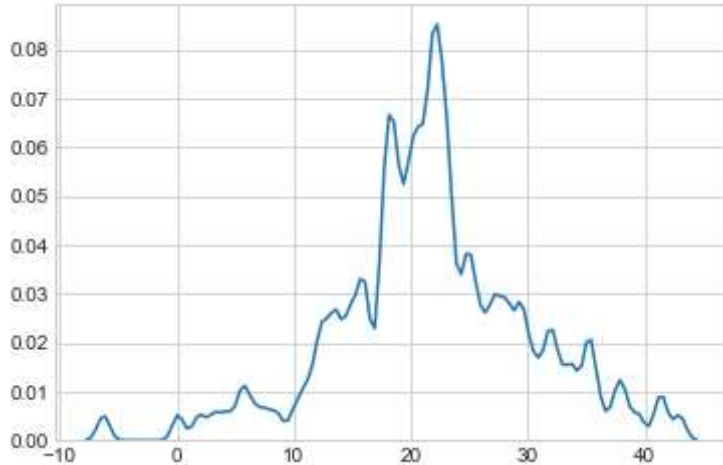
```
In [11]: 1 from sklearn.metrics import mean_squared_error
2 import seaborn as sns;
3 import numpy as np;
4
5
6 mse = (mean_squared_error(Y_test, Y_pred))
7 print('ROOT Mean Squared Error of linear regression',np.sqrt(mse))
```

ROOT Mean Squared Error of linear regression 5.541049738742553

```
In [12]: 1 delta_y = Y_test - Y_pred;
2
3 import seaborn as sns;
4 import numpy as np;
5 sns.set_style('whitegrid')
6 sns.kdeplot(np.array(delta_y), bw=0.5)
7 plt.show()
8
```



```
In [13]: 1 sns.set_style('whitegrid')
2 sns.kdeplot(np.array(Y_pred), bw=0.5)
3 plt.show()
```



## Defining the SGD Linear Algorithm

```
In [15]: 1 X.shape
```

```
Out[15]: (506, 13)
```

```
In [40]: 1 X_train.shape
```

```
Out[40]: (354, 13)
```

```

In [57]: 1 def dL_dW(X,e,N):
2         x11= (-2/N)*(X.T @ (e))
3         # for derivative functn w.r.t to 'w' we are taking the dot product (using
4         return x11
5 def dL_db(e,N):
6     return (-2/N)*np.sum(e)
7     # for derivative functn w.r.t to 'b' we are takinh summation of error e
8 def gradient_descent(learning_rate ,iterations):
9     '''
10    We are defining the SGD fuctn here for each of the training data in X.
11    the error is calculated and funtn is is updated in order to reduce the er
12    This process is repeated for a fixed number of iterations, the learning r
13    '''
14    N = len(X_train)
15    w1 = np.random.rand(1,(X_train.shape[1]))
16    b1 = np.random.rand()
17    optimal_w = []
18    optimal_b = []
19    for j in range (iterations):
20        # getting the error e , which returns me a vector of (354,1)
21        e = (Y_train[:,np.newaxis] - (X_train @ w1.T) - b1)
22        # getting the w1 , which returns me a vector of (1,13)
23        w1 = (w1) - (learning_rate*(dL_dW(X_train,e,N))).T
24        b1 = (b1) - (learning_rate*dL_db(e,N))
25    return w1,b1

```

```

In [81]: 1 w,b = gradient_descent(learning_rate = 0.1,iterations =10000)
2         w,b

```

```

Out[81]: (array([[ -1.26415881,  0.94329906, -0.16687636,  0.18653568, -1.49252028,
                  2.79557313, -0.29648219, -2.72594888,  2.76899352, -2.1378414 ,
                  -2.09193889,  1.16450017, -3.29650834]]), 22.556214689265538)

```

```

In [70]: 1 o= w.T
2         o.shape

```

```

Out[70]: (13, 1)

```

```
In [37]: 1 from sklearn.metrics import mean_squared_error
2
3 Y_pred_sgd = (X_test @ (o))+(b)
4 mse = (mean_squared_error(Y_test, Y_pred_sgd))
5 print('ROOT Mean Squared Error',np.sqrt(mse))
6 plt.scatter(Y_test, Y_pred_sgd)
7 plt.xlabel("Prices: $Y_i$")
8 plt.ylabel("Predicted prices: $\hat{Y}_i$")
9 plt.title("Prices vs Predicted prices: $Y_i$ vs $\hat{Y}_i$")
10 plt.show()
```

ROOT Mean Squared Error 5.541049688511123



```
In [84]: 1 from prettytable import PrettyTable
2
3 x = PrettyTable()
4
5 x.field_names = ["Method", "Metric to compare"]
6
7 x.add_row(["Linaer Regrsson Function", 'ROOT Mean Squared Error is 5.541049738742553'])
8 x.add_row(["SGD Function", 'ROOT Mean Squared Error is 5.541049688511123'])
9
10 print(x)
```

| Method                   | Metric to compare                            |
|--------------------------|--|
| Linaer Regrsson Function | ROOT Mean Squared Error is 5.541049738742553 |
| SGD Function             | ROOT Mean Squared Error is 5.541049688511123 |

```
In [ ]: 1
```