

Assignment 6

Implementing the SGD on Linear Regression for Boston House Price Dataset.

Objective:

To implement custom SGD on boston house price dataset and compare the results of custom SGD with the SKlearn SGD dataset results.

```
In [0]: 1 # Importing the required librarie
        2 import warnings
        3 warnings.filterwarnings("ignore")
        4 import pandas as pd
        5 import numpy as np
        6 import matplotlib.pyplot as plt
        7 from sklearn.datasets import load_boston           #Used to Load the Boston dataset
        8 from sklearn.model_selection import train_test_split
        9 from sklearn.preprocessing import StandardScaler
       10 from sklearn.linear_model import SGDRegressor
       11 from sklearn.metrics import mean_squared_error, r2_score
       12 from sklearn.model_selection import train_test_split
       13 from prettytable import PrettyTable
       14 import seaborn as sns
       15
```

Understanding the Dataset

```
In [0]: 1 Boston_DS=pd.DataFrame(load_boston().data,columns=load_boston().feature_names)
2 print("-"*120)
3 print(f'Shape of the Boston dataset: {(Boston_DS.shape)}')
4 print("-"*120)
5 print(f'No of features present in the Boston dataset: {(load_boston().feature_names)}')
6 print("-"*120)
7 Boston_DS.head(2)
```

Shape of the Boston dataset: (506, 13)

No of features present in the Boston dataset: ['CRIM' 'ZN' 'INDUS' 'CHAS' 'NOX' 'RM' 'AGE' 'DIS' 'RAD' 'TAX' 'PTRATIO' 'B' 'LSTAT']

```
Out[726]:
```

	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	396.9	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	396.9	9.14

Segregating the target from the features

```
In [0]: 1 Boston_DS['Price'] = load_boston().target
2 X_SKLN = Boston_DS.drop('Price', axis = 1)
3 Y_SKLN = Boston_DS['Price']
```

Splitting the data into Train and Test

```
In [0]: 1 X_SKLN_Train,X_SKLN_Test,Y_SKLN_Train,Y_SKLN_Test = train_test_split(X_SKLN,Y_SKLN,test_size =0.33,random_state=42)
2 print(X_SKLN_Train.shape)
3 print(X_SKLN_Test.shape)
4 print(Y_SKLN_Train.shape)
5 print(Y_SKLN_Test.shape)
6 X_SKLN_Train.mean()
```

```
(339, 13)
```

```
(167, 13)
```

```
(339,)
```

```
(167,)
```

```
Out[728]: CRIM      3.351324
ZN      11.716814
INDUS   11.261858
CHAS     0.076696
NOX     0.557498
RM       6.327324
AGE     68.940118
DIS      3.762468
RAD      9.483776
TAX     409.132743
PTRATIO  18.261652
B       358.431475
LSTAT    12.497611
dtype: float64
```

Standardizing the Data

```
In [0]: 1 std = StandardScaler()
2 X_SKLN_Train = std.fit_transform(X_SKLN_Train)
3 X_SKLN_Test = std.transform(X_SKLN_Test)
4 SKLN_SGD_RG = SGDRegressor()
5 SKLN_SGD_RG.fit(X_SKLN_Train, Y_SKLN_Train)
6 Y_SKLN_PRD = SKLN_SGD_RG.predict(X_SKLN_Test)
7 SkLearn_w=SKLN_SGD_RG.coef_
8 print(f"Sklearn's SGD regressor Coefficients:{(SKLN_SGD_RG.coef_)}")
9 print(f"Sklearn's SGD regressor Y Intercept:{(SKLN_SGD_RG.intercept_)}")
```

```
Sklearn's SGD regressor Coefficients:[-0.89378768  0.71539225  0.14751301  0.89272288 -1.67100139  2.85790581
-0.41028437 -2.85475367  1.32968721 -0.67137791 -2.01879448  1.04490366
-3.92669325]
```

```
Sklearn's SGD regressor Y Intercept:[22.97865716]
```

Implementing the own SGD Regressor for Linear Regression

```
In [0]: 1 OWN_Boston_DS = load_boston()
2 OWN_Boston_DS.data.shape
3 OWN_Boston_DS.feature_names
4 OWN_Boston_DS.target.shape
5 OWN_Boston_DS_data = pd.DataFrame(OWN_Boston_DS.data, columns = OWN_Boston_DS.feature_names)
6 print(OWN_Boston_DS_data.head())
```

	CRIM	ZN	INDUS	CHAS	NOX	...	RAD	TAX	PTRATIO	B	LSTAT
0	0.00632	18.0	2.31	0.0	0.538	...	1.0	296.0	15.3	396.90	4.98
1	0.02731	0.0	7.07	0.0	0.469	...	2.0	242.0	17.8	396.90	9.14
2	0.02729	0.0	7.07	0.0	0.469	...	2.0	242.0	17.8	392.83	4.03
3	0.03237	0.0	2.18	0.0	0.458	...	3.0	222.0	18.7	394.63	2.94
4	0.06905	0.0	2.18	0.0	0.458	...	3.0	222.0	18.7	396.90	5.33

[5 rows x 13 columns]

Normalizing the data

```
In [0]: 1 mean= OWN_Boston_DS_data.mean()
2 std= OWN_Boston_DS_data.std()
3 OWN_Boston_DS_data = (OWN_Boston_DS_data - mean)/std
4 print(OWN_Boston_DS_data.head())
```

	CRIM	ZN	INDUS	...	PTRATIO	B	LSTAT
0	-0.419367	0.284548	-1.286636	...	-1.457558	0.440616	-1.074499
1	-0.416927	-0.487240	-0.592794	...	-0.302794	0.440616	-0.491953
2	-0.416929	-0.487240	-0.592794	...	-0.302794	0.396035	-1.207532
3	-0.416338	-0.487240	-1.305586	...	0.112920	0.415751	-1.360171
4	-0.412074	-0.487240	-1.305586	...	0.112920	0.440616	-1.025487

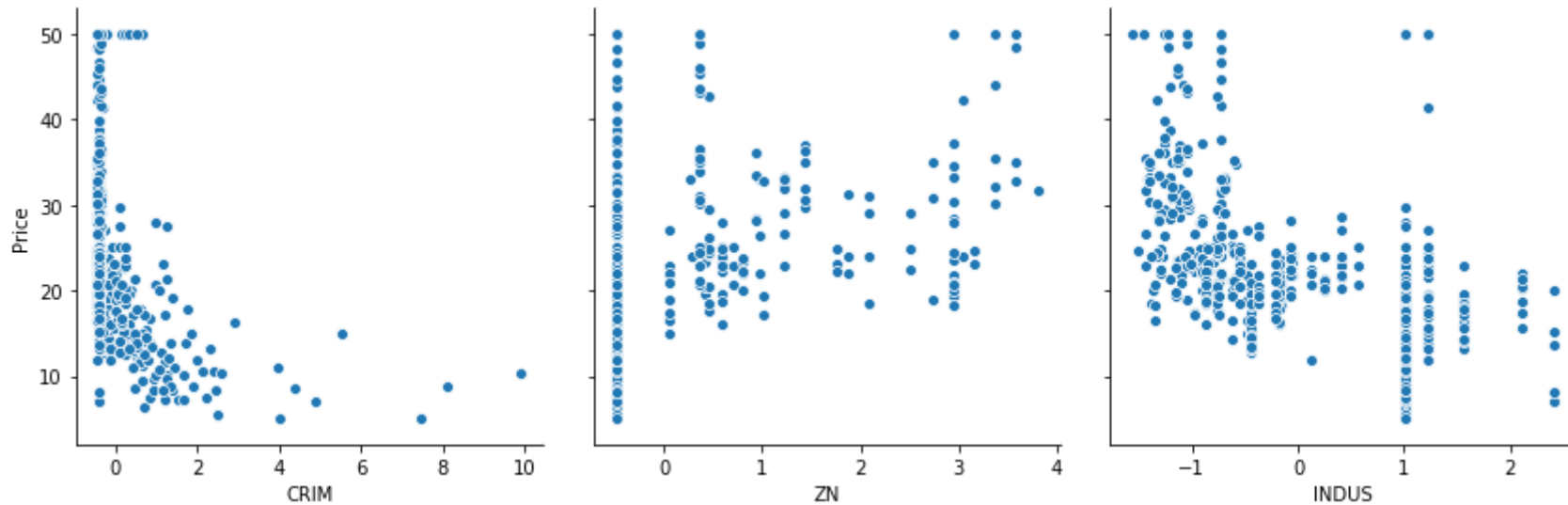
[5 rows x 13 columns]

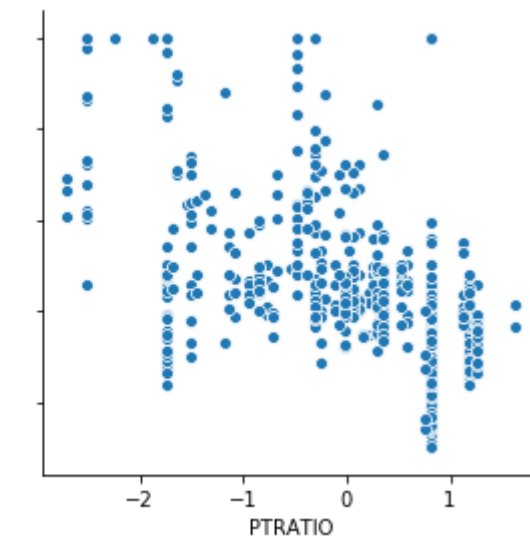
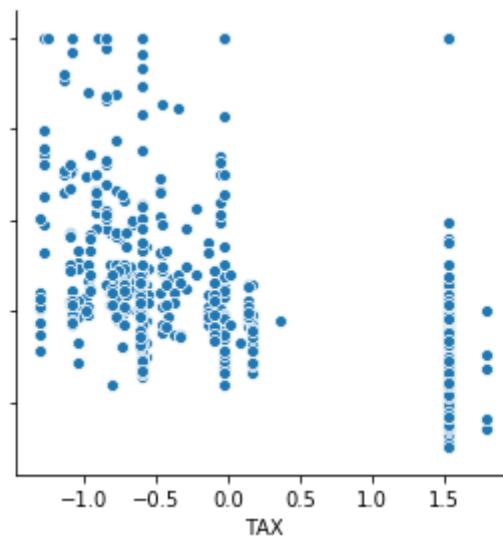
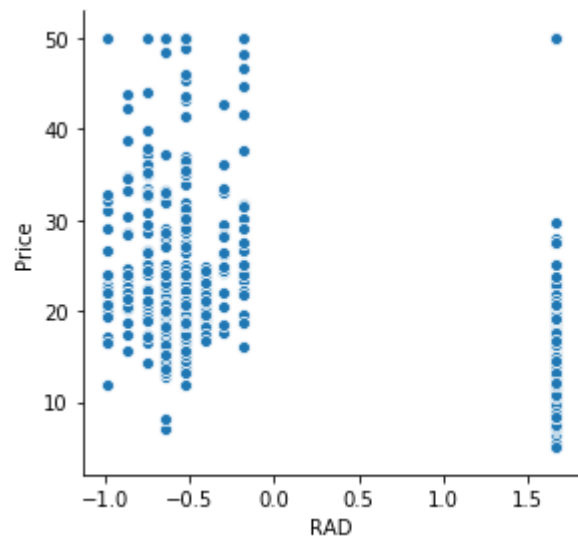
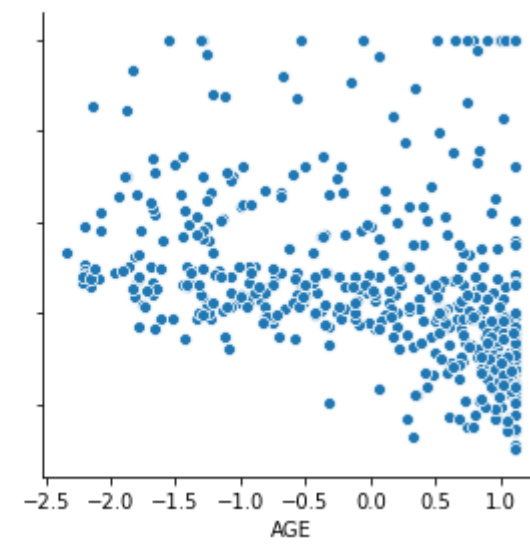
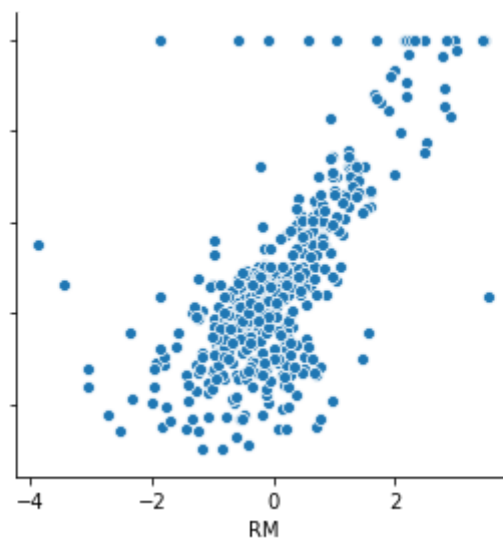
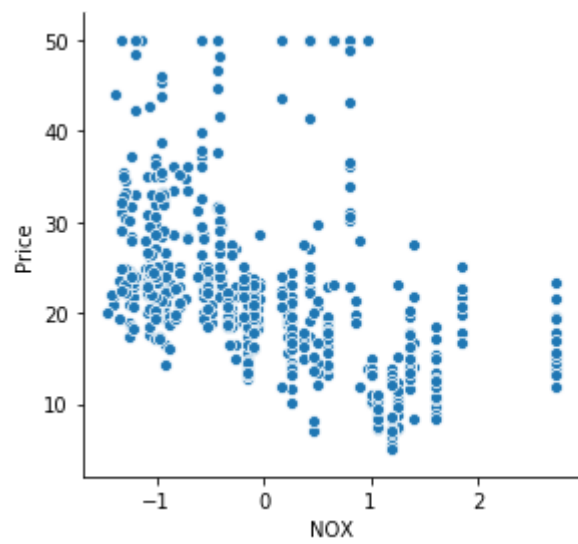
```
In [0]: 1 OWN_Boston_DS_data["Price"] = OWN_Boston_DS.target
2 OWN_Boston_DS_data.head()
3 OWN_Boston_Y = OWN_Boston_DS_data["Price"]
4 OWN_Boston_X = OWN_Boston_DS_data.drop("Price", axis = 1)
5 OWN_Train_X,OWN_Test_X,OWN_Train_Y,OWN_Test_Y=train_test_split(OWN_Boston_X,OWN_Boston_Y,test_size=0.33,random_state=42)
6 print(OWN_Train_X.shape,OWN_Train_Y.shape,OWN_Test_X.shape,OWN_Test_Y.shape)
7 OWN_Train_X["Price"] = OWN_Train_Y
8 x_sh=OWN_Train_X.shape[1]-1
```

(339, 13) (339,) (167, 13) (167,)

Comparing the price with the each features in the data

```
In [0]: 1 #https://stackoverflow.com/questions/31966494/compare-1-independent-vs-many-dependent-variables-using-seaborn-pairplot-in-an-h
2 plt.close();
3 Plot1 = sns.pairplot(data=OWN_Boston_DS_data,height=4,
4                     y_vars=['Price'],
5                     x_vars=['CRIM', 'ZN', 'INDUS'])
6 Plot2 = sns.pairplot(data=OWN_Boston_DS_data,height=4,
7                     y_vars=['Price'],
8                     x_vars=['NOX', 'RM', 'AGE'])
9 Plot3 = sns.pairplot(data=OWN_Boston_DS_data,height=4,
10                    y_vars=['Price'],
11                    x_vars=['RAD', 'TAX', 'PTRATIO'])
12 Plot4 = sns.pairplot(data=OWN_Boston_DS_data,height=4,
13                    y_vars=['Price'],
14                    x_vars=['CHAS', 'DIS', 'B', 'LSTAT'])
```





OBSERVATION:

- 1.The prices increase as the value of RM increases linearly.Also we can see few outliers
- 2.The prices tend to decrease with an increase in LSTAT.
- 3.CRIM rate seems to be high at the place of price ranges from 10 to 30.

```
In [0]: 1 print(OWN_Boston_DS_data.head(5))
```

	CRIM	ZN	INDUS	CHAS	...	PTRATIO	B	LSTAT	Price
0	-0.419367	0.284548	-1.286636	-0.272329	...	-1.457558	0.440616	-1.074499	24.0
1	-0.416927	-0.487240	-0.592794	-0.272329	...	-0.302794	0.440616	-0.491953	21.6
2	-0.416929	-0.487240	-0.592794	-0.272329	...	-0.302794	0.396035	-1.207532	34.7
3	-0.416338	-0.487240	-1.305586	-0.272329	...	0.112920	0.415751	-1.360171	33.4
4	-0.412074	-0.487240	-1.305586	-0.272329	...	0.112920	0.440616	-1.025487	36.2

[5 rows x 14 columns]

```
In [0]: 1 #https://www.geeksforgeeks.org/ml-r-squared-in-regression-analysis/  
2 def Loss_Function(intr_DS,weig_DS,X_DS,Y_DS):  
3     loss = 0  
4     for i in range(0, len(X_DS)):  
5         Exp1=Y_DS[:,i] - (np.dot(X_DS[i] , weig_DS) + intr_DS)  
6         loss += (Exp1) ** 2  
7     return loss/len(X_DS)  
8
```

OBSERVATION:

The loss function is calculated based on the below formula

$$\sum_{i=1}^K [y[i] - x[i].w^T + b]$$

Where W= Weight and b=Intercept


```

In [0]: 1 # https://stackoverflow.com/questions/50328545/stochastic-gradient-descent-for-linear-regression-on-partial-derivatives
2
3 def Own_regressor(weig, intr, Full_DS, X1, Y1, LR):
4
5     """Implementation of own SGD Gradient Descent for Linear regression"""
6
7     Weig_deriv=0
8     Intr_deriv=0
9     Iterations=1000
10    loss_Train=[]
11    loss_Test=[]
12    j=0
13    while j < 1000:
14        Full_DS_Batch=Full_DS.sample(100)
15        x = np.asmatrix(Full_DS_Batch.drop("Price", axis = 1))
16        y = np.asmatrix(Full_DS_Batch["Price"])
17        for i in range(len(x)):
18            tmp=y[:,i]-np.dot(x[i],weig)+intr
19            Weig_deriv+=np.dot(-2*x[i].T,tmp)
20            Intr_deriv+=(-2*tmp)
21        weig_update=weig-(LR*Weig_deriv)
22        b_update=intr-(LR*Intr_deriv)
23        if (weig==weig_update).all():
24            break
25        else:
26            weig=weig_update
27            intr=b_update
28            LR=LR/2
29        # Evaluating the loss in the custom built SGD regressor of train data
30        loss_TR=Loss_Function(intr,weig,x,y)
31        loss_Train.append(loss_TR)
32        # Evaluating the loss in the custom built SGD regressor of test data
33        loss_TS=Loss_Function(intr,weig,np.asmatrix(X1),np.asmatrix(Y1))
34        loss_Test.append(loss_TS)
35        j=j+1
36    return weig,intr,loss_Train,loss_Test

```

OBSERVATION:

1. At initial, keeping the Weight and Intercept values as zero
2. Setting up the iteration value
3. Setting up the batch of datasize
4. Calculating the

$$\frac{\partial L}{\partial \mathbf{w}} \text{ and } \frac{\partial L}{\partial \mathbf{b}}$$

Where W= Weight and b=Intercept

$$\frac{\partial \mathbf{L}}{\partial \mathbf{w}} = \sum_{i=1}^K [(-2.xi)(yi - (xi. w^T) + b)]$$
$$\frac{\partial \mathbf{L}}{\partial \mathbf{b}} = \sum_{i=1}^K [(-2)(yi - (xi. w^T) + b)]$$

5.At each iteration learning rate is reduced by half.

6.And the loop was keep executing till the weight get saturated.

7.At some point weight gets saturated. Finding that point and taking its previous value.

```

In [0]: 1 #https://matplotlib.org/3.1.1/tutorials/introductory/pyplot.html
2 lr = 0.001
3 # Choosing random values for w and b
4 intercept = np.random.rand()
5 wgt = np.random.rand(x_sh)
6 bias = np.asmatrix(wgt).T #Transpose of weight matrix
7 weig_, intr_, loss_Train_, loss_Test_ = Own_regressor(bias, intercept, OWN_Train_X, OWN_Test_X, OWN_Test_Y, lr)
8 print(f"SGD Coefficient: {(weig_)}")
9 print(f"Y_Intercept: {(intr_)}")
10
11 print(f"Train Loss = {(loss_Train_)}")
12 print(f"Test Loss= {(loss_Test_)} ")
13
14 OWN_Prediction=(np.dot(np.asmatrix(OWN_Test_X), weig_) + intr_)
15 OWN_Prediction_List=np.array(OWN_Prediction).T[0]
16 # Error Plot
17 plt.close();
18 plt.figure()
19 plt.plot(range(len(loss_Train_)), np.reshape(loss_Train_,[len(loss_Train_), 1]), label = "Train Loss",c='red', ls='-',lw=2)
20 plt.plot(range(len(loss_Test_)), np.reshape(loss_Test_, [len(loss_Test_), 1]), label = "Test Loss",c='blue', ls='-',lw=2)
21 plt.title("Error Plot")
22 plt.xlabel("No of Iterations")
23 plt.ylabel("Error's")
24 plt.legend()
25 plt.show()

```

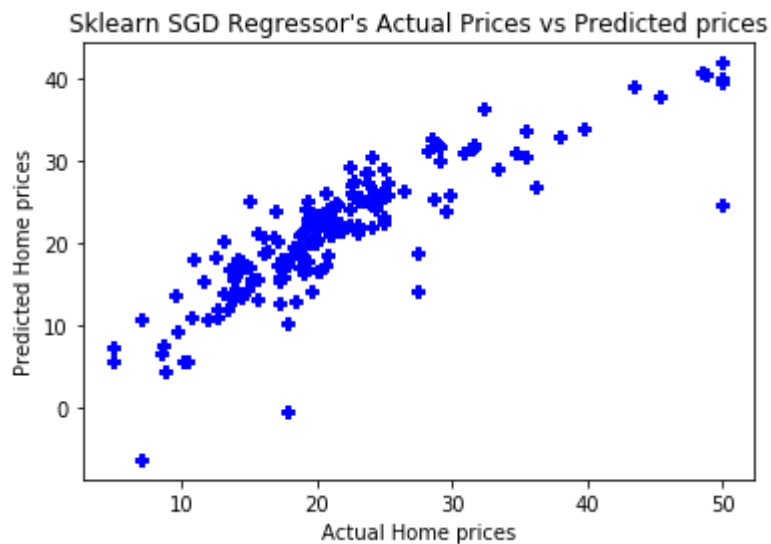
SGD Coefficient: [[2.58298431e-01]

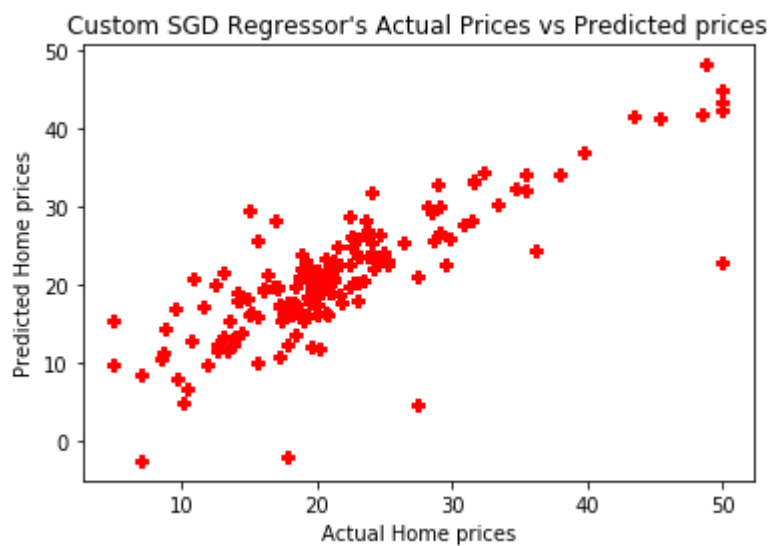
[3.39221477e-01]
 [-7.37036275e-01]
 [1.06994685e+00]
 [-4.42579777e-02]
 [4.87969387e+00]
 [8.78820383e-01]
 [-7.40355353e-01]
 [4.58860351e-01]
 [-7.82629557e-04]
 [-2.52489286e+00]
 [9.83409962e-01]
 [-3.01768581e+00]]

Y_Intercept: [[22.35480459]]

Train Loss = [matrix([[370.20906386]]), matrix([[206.28317897]]), matrix([[107.22916446]]), matrix([[63.09551347]]), matrix([[4
 6.62212318]]), matrix([[44.68280846]]), matrix([[19.99057045]]), matrix([[28.32658378]]), matrix([[23.63833973]]), matrix([[19.
 31282885]]), matrix([[32.69288678]]), matrix([[37.32332292]]), matrix([[45.33978117]]), matrix([[27.19064905]]), matrix([[30.54
 485685]]), matrix([[31.40479638]]), matrix([[22.16166989]]), matrix([[23.80463751]]), matrix([[28.07904338]]), matrix([[27.6488
 0512]]), matrix([[28.02435142]]), matrix([[23.37228194]]), matrix([[26.39519442]]), matrix([[35.32755431]]), matrix([[30.714889
 46]]), matrix([[41.37173085]]), matrix([[36.56274588]]), matrix([[27.86868173]]), matrix([[44.96784203]]), matrix([[19.1549883


```
In [0]: 1 #https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.markers.MarkerStyle.html#matplotlib.markers.MarkerStyle
2 #https://matplotlib.org/3.1.1/api/_as_gen/matplotlib.pyplot.scatter.html
3 # Applying sklearn SGD Regressor.
4 plt.close();
5 plt.figure(1)
6 plt.subplot(111)
7 plt.scatter(Y_SKLN_Test, Y_SKLN_PRD,marker='P',c='blue')
8 plt.xlabel("Actual Home prices")
9 plt.ylabel("Predicted Home prices")
10 plt.title("Sklearn SGD Regressor's Actual Prices vs Predicted prices")
11 plt.show()
12
13 # Applying own SGD Regressor
14 plt.close();
15 plt.figure(2)
16 plt.subplot(111)
17 plt.scatter([OWN_Test_Y], [(np.dot(np.asmatrix(OWN_Test_X), weig_) + intr_)],marker='P',c='red')
18 plt.xlabel("Actual Home prices")
19 plt.ylabel("Predicted Home prices")
20 plt.title("Custom SGD Regressor's Actual Prices vs Predicted prices")
21 plt.show()
```





```
In [0]: 1 # Sklearn SGD Regression
        2 print(f"MSE of Sk learn's prediction:{(mean_squared_error(Y_SKLN_Test, Y_SKLN_PRD))}")
        3 print(f"r2_score of Sk learn's prediction:{(r2_score(Y_SKLN_Test, Y_SKLN_PRD))}")
```

MSE of Sk learn's prediction:21.228706973845743
r2_score of Sk learn's prediction:0.7194883008348725

```
In [0]: 1 # https://stackoverflow.com/questions/40901445/function-to-calculate-r2-r-squared-in-r
        2 #import math
        3 OWN_loss = Loss_Function(intr_, weig_, np.asmatrix(OWN_Test_X), np.asmatrix(OWN_Test_Y))
        4 print(f"MSE of own SGD_reg :{(float(OWN_loss))}")
        5
        6 X1=np.asmatrix(OWN_Test_X)
        7 Y1=np.asmatrix(OWN_Test_Y)
        8 for i in range(0, len(X1)):
        9     Avg_Y = np.mean(Y1)
       10     Exp2=(Y1[:,i] - Avg_Y)
       11     Exp3=(Y1[:,i] - (np.dot(X1[i], weig_) + intr_))
       12     Sum_of_squares = sum((Exp2)**2)
       13     Residual_sum = sum((Exp3)**2)
       14     R_Square = 1-(Residual_sum/Sum_of_squares)
       15
       16 print(f"r2_score of own SGDreg :{(float(R_Square))}")
```

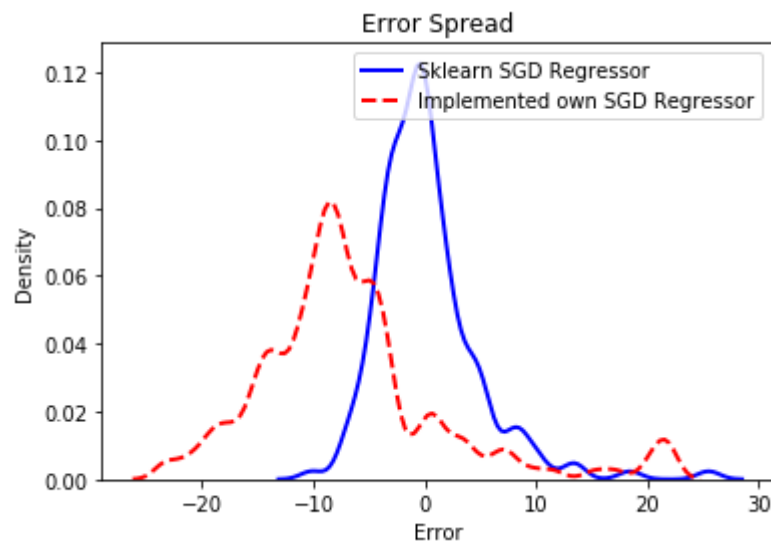
MSE of own SGD_reg :25.476540352979335
r2_score of own SGDreg :0.9992409296653334

OBSERVATION:

Calculating the r square from the below formula
Where

$$r^2 = 1 - \frac{\text{Residual Sum of squares}}{\text{Sum of squares}}$$

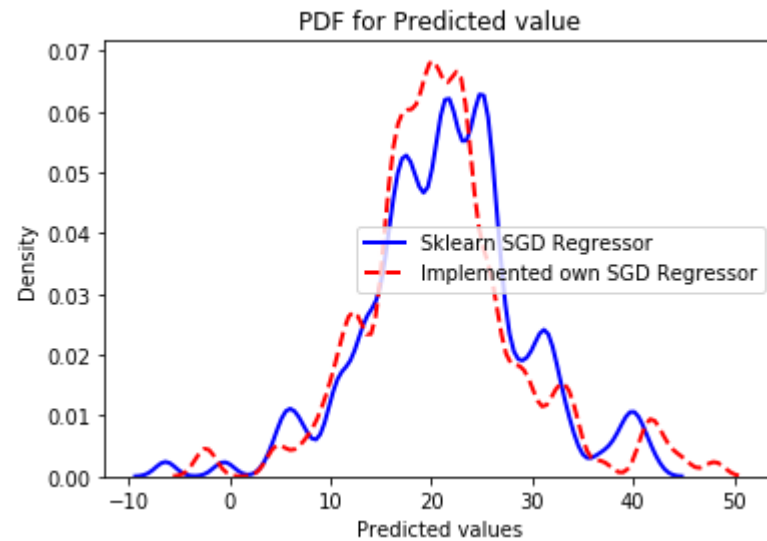
```
In [0]: 1 # Error distribution
2 OWN_Loss_Y = np.asmatrix(OWN_Test_Y) - (OWN_Prediction)
3 SKLN_Loss_Y = Y_SKLN_Test - Y_SKLN_PRD
4
5 sns.kdeplot(np.array(SKLN_Loss_Y), label = "Sklearn SGD Regressor", bw =1,c='blue', ls='-',lw=2)
6 sns.kdeplot(np.asarray(OWN_Loss_Y)[0], label = "Implemented own SGD Regressor", bw =1,c='red', ls='--',lw=2)
7 plt.title("Error Spread")
8 plt.xlabel("Error")
9 plt.ylabel("Density")
10 plt.legend()
11 plt.show()
```



OBSERVATION:

From above plot SKLearn SGD is working better than custom implementation regressor. Because the over all distribution of SKlearn with respectiv error is low.

```
In [0]: 1 # Predicted value distribution
2 sns.kdeplot(Y_SKLN_PRD, label = "Sklearn SGD Regressor", bw =1,c='blue', ls='-',lw=2)
3 sns.kdeplot(OWN_Prediction_List, label = "Implemented own SGD Regressor", bw =1,c='red', ls='--',lw=2)
4 plt.title("PDF for Predicted value")
5 plt.xlabel("Predicted values")
6 plt.ylabel("Density")
7 plt.show()
```



OBSERVATION:

From above plot you can see the most of the predicted values are overlap.

Comparison with Pretty Table


```
In [0]: 1 PT = PrettyTable()
2 PT.field_names=['S.No','Custom SGD Regressor','SkLearn SGD Regressor']
3 for i in range(x_sh):
4     PT.add_row([i+1,float(weig_[i]),SkLearn_w[i]])
5 print(PT)
```

S.No	Custom SGD Regressor	SkLearn SGD Regressor
1	0.2582984312172637	-0.8937876809095562
2	0.3392214765870696	0.7153922469815946
3	-0.7370362747081937	0.14751301403165415
4	1.0699468487113792	0.8927228844675043
5	-0.04425797767794621	-1.6710013907984682
6	4.8796938688994365	2.8579058082668083
7	0.87882038297057	-0.4102843692600765
8	-0.7403553528758017	-2.854753671182647
9	0.45886035113737783	1.3296872106166655
10	-0.0007826295570831534	-0.6713779089268122
11	-2.5248928602744205	-2.0187944797150936
12	0.983409961680586	1.044903661065243
13	-3.0176858146709526	-3.9266932523379334

SUMMARY:

Overall when we compare the custom implemeneted SGD regressor with the SKlearn SGD regressor they are working little similar. But still we can see the Sklearn SGD regressor model is performing good when we consider the MSE and r square values. Hence i conclude the Sklearn SGD regressor model have better performance when ccompare to the custom SGD regressor.

REFERENCE:

To create a mathematical equations in the Jupyter Notebook

<https://www.math.ubc.ca/~pwalls/math-python/jupyter/latex/> (<https://www.math.ubc.ca/~pwalls/math-python/jupyter/latex/>)

<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Working%20With%20Markdown%20Cells.html> (<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Working%20With%20Markdown%20Cells.html>)

<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Typesetting%20Equations.html> (<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Typesetting%20Equations.html>)

<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Typesetting%20Equations.html> (<https://jupyter-notebook.readthedocs.io/en/stable/examples/Notebook/Typesetting%20Equations.html>)