

Programmed By : Rithik Tripathi

[Connect with me on LinkedIn \(https://www.linkedin.com/in/rithik-tripathi-data-scientist/\)](https://www.linkedin.com/in/rithik-tripathi-data-scientist/)

1 This notebook explains how cost function is used to find best value of parameters using LINEAR REGRESSION

Cost Function ¶

In [1]:

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.metrics import mean_squared_error as mse
```

C:\Users\rkt7k\anaconda3\lib\site-packages\scipy__init__.py:146: UserWarning: A NumPy version >=1.16.5 and <1.23.0 is required for this version of SciPy (detected version 1.24.1
warnings.warn(f"A NumPy version >={np_minversion} and <{np_maxversion}")

In [2]:

```
1 # creating the sample dataset
2 experience = [1.2,1.5,1.9,2.2,2.4,2.5,2.8,3.1,3.3,3.7,4.2,4.4]
3 salary     = [1.7,2.4,2.3,3.1,3.7,4.2,4.4,6.1,5.4,5.7,6.4,6.2]
4
5 df = pd.DataFrame({
6     'experience' : experience,
7     'salary' : salary
8 })
9 df.head()
```

Out[2]:

	experience	salary
0	1.2	1.7
1	1.5	2.4
2	1.9	2.3
3	2.2	3.1
4	2.4	3.7

Defining a Linear Regression Equation to Predict Salary

(Variables taken as constant for demonstration)

In [3]:

```
1 pred_salary = np.array(experience)*1.7 - 0.7
2 df['pred_salary'] = pred_salary
```

In [4]:

```
1 df.head()
```

Out[4]:

	experience	salary	pred_salary
0	1.2	1.7	1.34
1	1.5	2.4	1.85
2	1.9	2.3	2.53
3	2.2	3.1	3.04
4	2.4	3.7	3.38

Calculating error

In [5]:

```
1 df['error'] = df['salary'] - df['pred_salary']
2 error = salary - pred_salary # storing in array
```

Calculating MSE

In [6]:

```
1 # method 1
2 np.square(error).sum()/len(error)
```

Out[6]:

0.3383500000000003

In [7]:

```
1 # method 2
2 mse(df.salary, df.pred_salary)
```

Out[7]:

0.3383500000000003

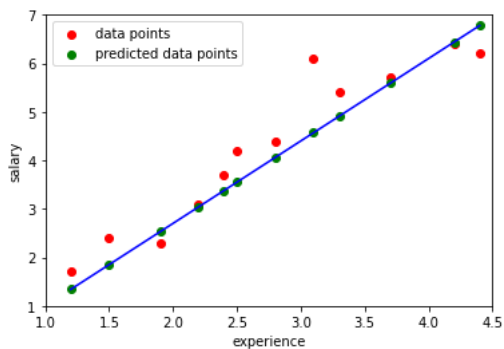
Visualizing Data

In [8]:

```
1 plt.scatter(df.experience, df.salary, color = 'red', label = 'data points')
2 plt.scatter(df.experience, df.pred_salary, color = 'green', label = 'predicted data points')
3 plt.plot(df.experience, df.pred_salary, '-b')
4 plt.xlim(1,4.5)
5 plt.ylim(1,7)
6 plt.xlabel('experience')
7 plt.ylabel('salary')
8 plt.legend()
9
10 # data follows linear curve
11 # blue line demonstrates the best fit line
```

Out[8]:

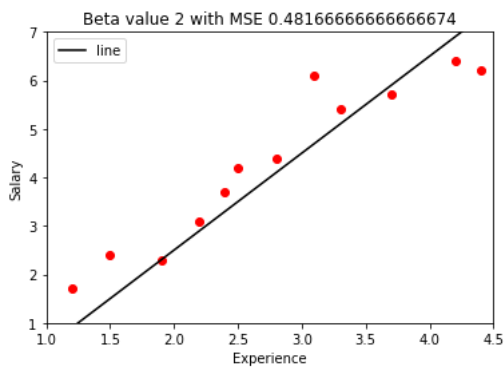
<matplotlib.legend.Legend at 0x15daa054970>



Try playing around with the values of beta and b here to observe the change in fitting line

In [9]:

```
1 beta = 2
2 b = 1.5
3
4 #predicted points
5 line1 = np.array(experience)*beta - b
6
7 # Plotting the line
8 plt.scatter(df.experience, df.salary, color = 'red')
9 plt.plot(df.experience, line1, color = 'black', label = 'line')
10 plt.xlim(1,4.5)
11 plt.ylim(1,7)
12 plt.xlabel('Experience')
13 plt.ylabel('Salary')
14 plt.legend()
15 MSE = mse(df.salary, line1)
16 plt.title("Beta value "+str(beta)+" with MSE "+ str(MSE))
17 MSE = mse(df.salary, line1)
```



Calculating error for a range of values of Beta

In [15]:

```
1 def calc_error(beta, df):
2
3     b = 1.1 # considering b as constant
4
5     pred_salary = df.experience*beta + b # predictions for each data point
6
7     MSE = mse(df.salary, pred_salary) # calculating MSE for current Beta value predicitions
8
9     return MSE
```

In [19]:

```
1 # Calculating Cost (Errors) for raneg of Beta (slope) values
2 slope = [i/100 for i in range(0,300)]
3 costs = []
4
5 for i in slope:
6     cost = calc_error(beta= i, df= df)
7     costs.append(cost)
8
9 # Arranging in a DataFrame
10 cost_df = pd.DataFrame({
11     'Beta' : slope,
12     'Cost' : costs
13 })
14 cost_df.head()
```

Out[19]:

	Beta	Cost
0	0.00	12.791667
1	0.01	12.585876
2	0.02	12.381806
3	0.03	12.179455
4	0.04	11.978824

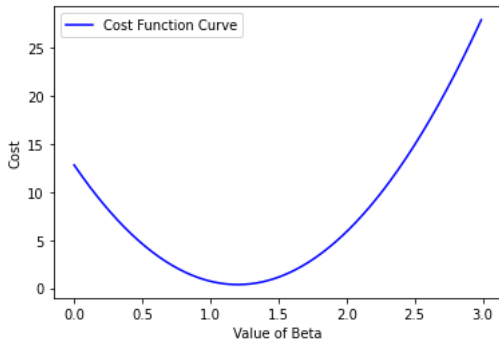
Visualising Cost v/s Beta

In [20]:

```
1 # plotting the cost values corresponding to every value of Beta
2 plt.plot(cost_df.Beta, cost_df.Cost, color = 'blue', label = 'Cost Function Curve')
3 plt.xlabel('Value of Beta')
4 plt.ylabel('Cost')
5 plt.legend()
```

Out[20]:

<matplotlib.legend.Legend at 0x15dac390eb0>



In [35]:

```
1 # Method 1 => Calc. Lowest cost and then Locate respective Beta with index Location of cost
2 min_cost = cost_df[cost_df['Cost'] == cost_df['Cost'].min()].index
3 best_beta = cost_df.iloc[min_cost]['Beta']
4 best_beta
```

Out[35]:

```
120    1.2
Name: Beta, dtype: float64
```

In [40]:

```
1 # Method 2 => Sort Dataframe w.r.t Cost and get the beta value of 0th row
2 cost_df.sort_values(by='Cost')['Beta'].iloc[0]
```

Out[40]:

```
1.2
```