

GRIP - The Sparks Foundation

Task - Prediction using Supervised ML

To Predict the percentage of marks of the students based on the number of hours they studied

In [1]:

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 from sklearn.model_selection import train_test_split
6 from sklearn.linear_model import LinearRegression
7 from sklearn.metrics import mean_absolute_error
```

In [2]:

```
1 data = pd.read_csv('http://bit.ly/w-data')
2 data.head(5)
```

Out[2]:

	Hours	Scores
0	2.5	21
1	5.1	47
2	3.2	27
3	8.5	75
4	3.5	30

In [3]:

```
1 data.isnull == True
```

Out[3]:

False

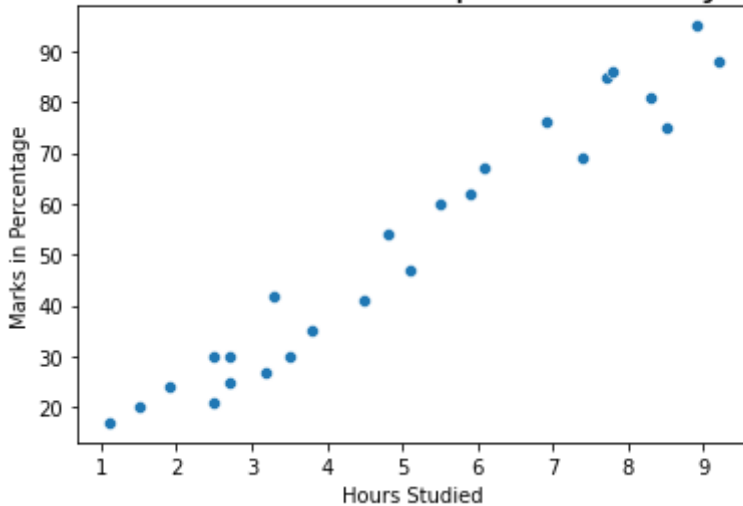
As we have found no null data in the dataset, we can directly move ahead to visualize the data

In [4]:



```
1 sns.scatterplot(y= data['Scores'], x= data['Hours'])
2 plt.title('Marks obtained with respect to Study Hours',size=20)
3 plt.ylabel('Marks in Percentage', size=10)
4 plt.xlabel('Hours Studied', size=10)
5 plt.show()
6
```

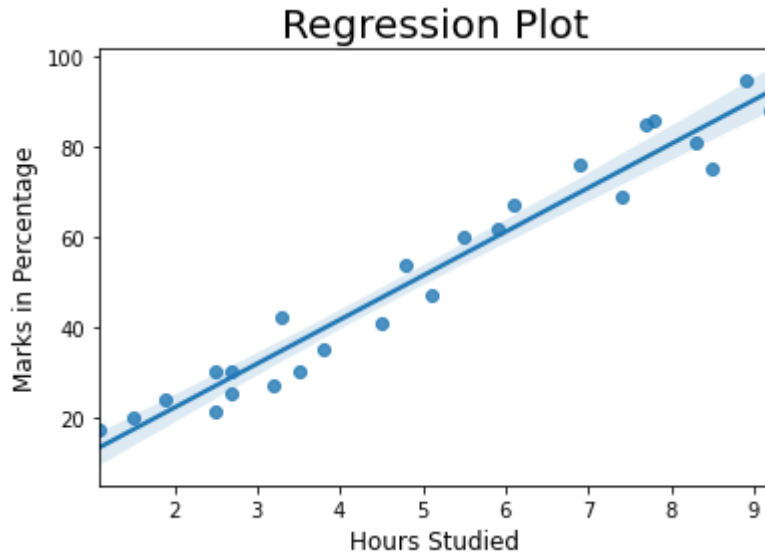
Marks obtained with respect to Study Hours



From the above scatter plot there looks to be correlation between the 'Marks in Percentage' and 'Hours Studied', Lets plot a regression line to confirm the correlation.

In [5]:

```
1 sns.regplot(x= data['Hours'], y= data['Scores'])
2 plt.title('Regression Plot',size=20)
3 plt.ylabel('Marks in Percentage', size=12)
4 plt.xlabel('Hours Studied', size=12)
5 plt.show()
6 print(data.corr())
```



	Hours	Scores
Hours	1.000000	0.976191
Scores	0.976191	1.000000

It is very clear from the plot that these two variables are positively related.

Training the Model

Splitting the Data

In [6]:

```
1 X = data.iloc[:, :-1].values
2 y = data.iloc[:, 1].values
3 train_X, val_X, train_y, val_y = train_test_split(X, y, random_state = 0)
```

Fitting the Data into the model

In [7]:

```
1 regression = LinearRegression()
2 regression.fit(train_X, train_y)
```

Out[7]:

LinearRegression()

Predicting the Percentage of Marks

In [8]:



```
1 pred_y = regression.predict(val_X)
2 prediction = pd.DataFrame({'Hours': [i[0] for i in val_X], 'Predicted Marks': [k for k
3 prediction
```

Out[8]:

	Hours	Predicted Marks
0	1.5	16.844722
1	3.2	33.745575
2	7.4	75.500624
3	2.5	26.786400
4	5.9	60.588106
5	3.8	39.710582
6	1.9	20.821393

Comparing the Predicted Marks with the Actual Marks

In [9]:



```
1 compare_scores = pd.DataFrame({'Actual Marks': val_y, 'Predicted Marks': pred_y})
2 compare_scores
```

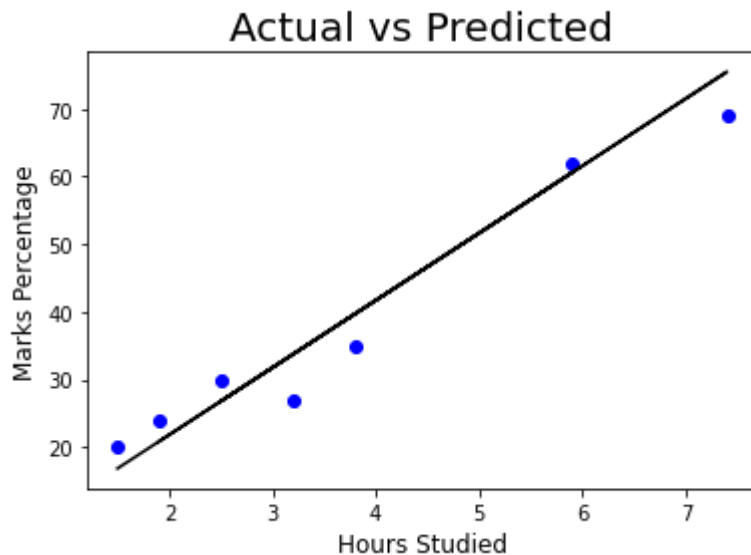
Out[9]:

	Actual Marks	Predicted Marks
0	20	16.844722
1	27	33.745575
2	69	75.500624
3	30	26.786400
4	62	60.588106
5	35	39.710582
6	24	20.821393

Visually Comparing the Predicted Marks with the Actual Marks

In [10]:

```
1 plt.scatter(x=val_X, y=val_y, color='blue')
2 plt.plot(val_X, pred_y, color='Black')
3 plt.title('Actual vs Predicted', size=20)
4 plt.ylabel('Marks Percentage', size=12)
5 plt.xlabel('Hours Studied', size=12)
6 plt.show()
```



Evaluating the Model

In [11]:

```
1 print('Mean absolute error: ',mean_absolute_error(val_y,pred_y))
```

Mean absolute error: 4.130879918502482

Small value of Mean absolute error states that the chances of error or wrong forecasting through the model are very less.

What will be the predicted score of a student if he/she studies for 9.25 hrs/ day?

In [12]:



```
1 hours = [9.25]
2 answer = regression.predict([hours])
3 print("Score = {}".format(round(answer[0],3)))
```

Score = 93.893

According to the regression model if a student studies for 9.25 hours a day he/she is likely to score 93.89 marks.

Author- Vedangi Sharma

In []:



1