

The Spark Foundation

Data Science & Business Analytics Internship jan-2022

Author :- Shamu Vishwakarma

Task 1 :- Predict the percentage of student based on the no of study hours

Objective :- predict the percentage of an student based on the no. of study hours.

Importing all necessary libraries

```
In [3]: 1 import pandas as pd
2 # for data manipulation & working with csv files
3
4 import numpy as np
5 # for numerical manipulation
6
7 import matplotlib.pyplot as plt
8 # for plotting graphs
9
10 import seaborn as sns
11 # for making statistical graphics
12
13 from sklearn.model_selection import train_test_split
14 # for splitting data set
15
16 from sklearn.linear_model import LinearRegression
17 # for linear regression
18 %matplotlib inline
19
```

Reading data from remote link

```
In [4]: 1 url="http://bit.ly/w-data"
2 Student_Data=pd.read_csv(url)
3 print("Data succesfully loaded")
```

Data succesfully loaded

```
In [5]: 1 #Reading Data set  
        2 Student_Data
```

Out[5]:

	Hours	Scores
0	2.5	21
1	5.1	47
2	3.2	27
3	8.5	75
4	3.5	30
5	1.5	20
6	9.2	88
7	5.5	60
8	8.3	81
9	2.7	25
10	7.7	85
11	5.9	62
12	4.5	41
13	3.3	42
14	1.1	17
15	8.9	95
16	2.5	30
17	1.9	24
18	6.1	67
19	7.4	69
20	2.7	30
21	4.8	54
22	3.8	35
23	6.9	76
24	7.8	86

```
In [8]: 1 # Checking shape of data set
        2 print(Student_Data.shape)
        3 Student_Data.head() #for reading top five Rows
```

(25, 2)

Out[8]:

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

```
In [9]: 1 #for reading Bottom five rows
        2 Student_Data.tail()
```

Out[9]:

	Hours	Scores
20	2.7	30
21	4.8	54
22	3.8	35
23	6.9	76
24	7.8	86

```
In [10]: 1 # Checking null values in data set
         2 Student_Data.isnull().sum()
```

Out[10]: Hours 0
Scores 0
dtype: int64

```
In [11]: 1 Student_Data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 25 entries, 0 to 24
Data columns (total 2 columns):
#   Column  Non-Null Count  Dtype
---  -
0   Hours   25 non-null        float64
1   Scores  25 non-null        int64
dtypes: float64(1), int64(1)
memory usage: 528.0 bytes
```

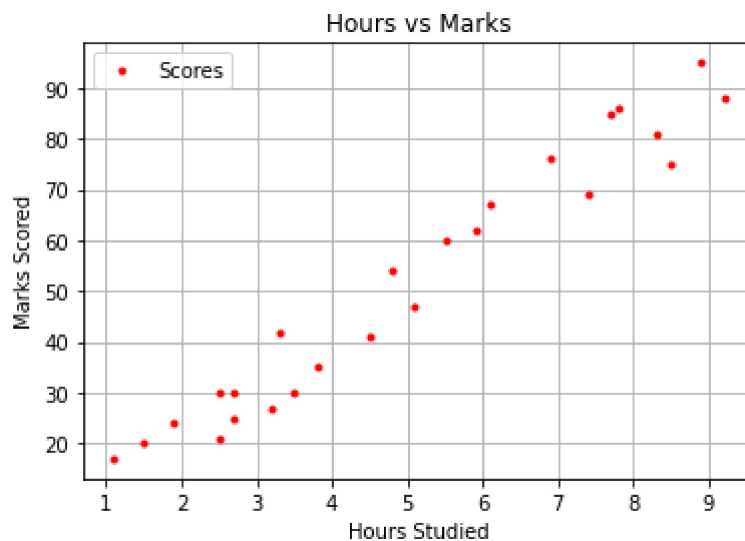
```
In [12]: 1 # Checking numerical data
         2 Student_Data.describe()
```

Out[12]:

	Hours	Scores
count	25.000000	25.000000
mean	5.012000	51.480000
std	2.525094	25.286887
min	1.100000	17.000000
25%	2.700000	30.000000
50%	4.800000	47.000000
75%	7.400000	75.000000
max	9.200000	95.000000

Now we will plot the Graph using matplotlib to understand relation between columns

```
In [13]: 1 # Plotting the distribution of scores
         2 Student_Data.plot(x='Hours', y='Scores', style='.',color='red')
         3 plt.title('Hours vs Marks')
         4 plt.xlabel('Hours Studied')
         5 plt.ylabel('Marks Scored')
         6 plt.grid()
         7 plt.show()
```



Here We observe that there is linear relationship Between the Marks scored by the student & their Respective Study Hours. So, we will use simple linear regression supervised Machine Learning Model to predict the Further values.

```
In [14]: 1 # Correlation coeff is 0.976191
          2 # which is a strong positive correlation.
          3
          4 # Checking correlation between columns
          5 Student_Data.corr()
```

Out[14]:

	Hours	Scores
Hours	1.000000	0.976191
Scores	0.976191	1.000000

Preparing the data

We are going to divide this dataset column (i.e Hours,Scores) into "attribute" (inputs) & "label" (outputs), here Hours is attribute & Scores are label

```
In [15]: 1 X = Student_Data.iloc[:, :-1].values
          2 y = Student_Data.iloc[:, 1:].values
```

```
In [16]: 1 X
```

Out[16]: array([[2.5],
[5.1],
[3.2],
[8.5],
[3.5],
[1.5],
[9.2],
[5.5],
[8.3],
[2.7],
[7.7],
[5.9],
[4.5],
[3.3],
[1.1],
[8.9],
[2.5],
[1.9],
[6.1],
[7.4],
[2.7],
[4.8],
[3.8],
[6.9],
[7.8]])

In [17]: 1 y

Out[17]: array([[21],
[47],
[27],
[75],
[30],
[20],
[88],
[60],
[81],
[25],
[85],
[62],
[41],
[42],
[17],
[95],
[30],
[24],
[67],
[69],
[30],
[54],
[35],
[76],
[86]], dtype=int64)

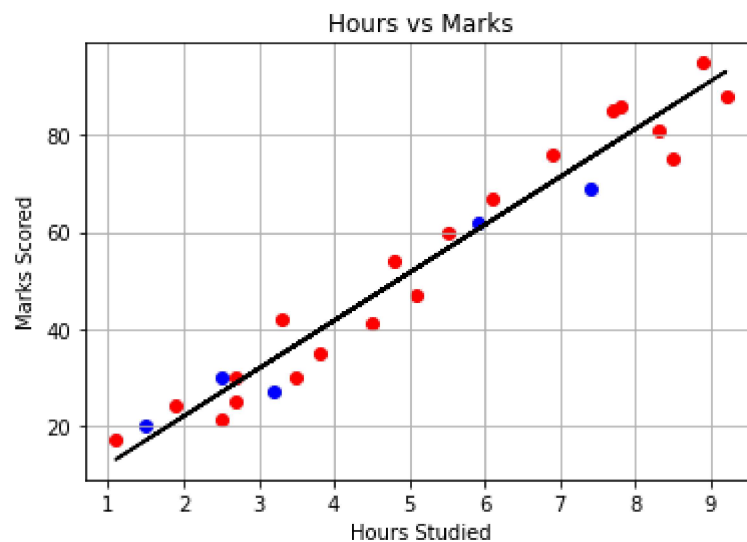
```
In [18]: 1 from sklearn.model_selection import train_test_split
2 # Splitting the data into train & test sets
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, rand
4 #here 80% of our data is training data and 20% is the testing data
5
6 print('rows in the total set: {}'.format(Student_Data.shape[0]))
7 print('rows in the training set: {}'.format(X_train.shape[0]))
8 print('rows in the test set: {}'.format(X_test.shape[0]))
```

rows in the total set: 25
rows in the training set: 20
rows in the test set: 5

```
In [19]: 1 from sklearn.metrics import accuracy_score
2 from sklearn.linear_model import LinearRegression
3 regressor = LinearRegression()
4 regressor.fit(X_train, y_train)
```

Out[19]: LinearRegression()

```
In [20]: 1 # Plotting the regression line
2 line = regressor.coef_*X+regressor.intercept_
3
4 # Plotting for the test data
5 plt.scatter(X_train, y_train,color="red")
6 plt.scatter(X_test, y_test,color="blue")
7 plt.plot(X, line, color="black");
8 plt.title('Hours vs Marks')
9 plt.xlabel('Hours Studied')
10 plt.ylabel('Marks Scored')
11 plt.grid()
12 plt.show()
```



Testing our Linear Regression Model

```
In [21]: 1 print(X_test) # Testing data - In Hours
2 y_pred = regressor.predict(X_test) # Predicting the scores
```

```
[[1.5]
 [3.2]
 [7.4]
 [2.5]
 [5.9]]
```

```
In [22]: 1 print(y_test)
          2 print(y_pred)
```

```
[[20]
 [27]
 [69]
 [30]
 [62]]
[[16.88414476]
 [33.73226078]
 [75.357018 ]
 [26.79480124]
 [60.49103328]]
```

```
In [23]: 1 # Comparing Actual vs Predicted
          2 df = pd.DataFrame({'Actual': [y_test], 'Predicted': [y_pred]})
          3 df
```

Out[23]:

	Actual	Predicted
0	[[20], [27], [69], [30], [62]]	[[16.884144762398037], [33.73226077948984], [7...

What will be the predicted score if a student study for 9.10 hrs/day?

```
In [24]: 1 # now we are ready to test with your own data
          2 hours = 9.10
          3 own_pred = regressor.predict([[hours]])
          4 print("No of Hours = {}".format(hours)+" hr.")
          5 print("Predicted Score = {}".format(own_pred[0]))
```

No of Hours = 9.1 hr.
Predicted Score = [92.20513402]

Hence we can conclude that if a student is involved in 9.10 hours per day , then there is a possibility that the percentage comes out to be 92.20513402


```

In [27]: 1 from tkinter import *
2 def alert_popup(title, message):
3     """Generate a pop-up window ."""
4     root = Tk()
5     root.title(title)
6     w = 300      # popup window width
7     h = 200      # popup window height
8     sw = root.winfo_screenwidth()
9     sh = root.winfo_screenheight()
10    x = (sw - w)/2
11    y = (sh - h)/2
12    root.geometry('%dx%d+%d+%d' % (w, h, x, y))
13    m = message
14    w = Label(root, text=m, width=50, height=10)
15    w.pack()
16    b = Button(root, text="OK", command=root.destroy, width=10)
17    b.pack()
18    mainloop()
19
20 alert_popup("Predictions", own_pred[0])

```

Evaluation of linear regression model

The final step is to evaluate the performance of algorithms. This step is particularly important to compare how well different algorithms perform on particular data set for simplicity here, we have chosen the means square error. There are many such metrics

```

In [26]: 1 from sklearn import metrics
2 print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))

```

Mean Absolute Error: 4.183859899002975

In []: 1

In []: 1

In []: 1

In []: 1

In []: 1

In []: 1

In []: 1

In []: 1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1

In []:

1