

GRIP -Graduate Rotational Internship Program

Data Science & Business analytics Internship

Linear Regression with Python Scikit Learn

In this section we will see how the Python Scikit-Learn library for machine learning can be used to implement regression functions. We will start with simple linear regression involving two variables

Simple Linear Regression

In this regression task we will predict the percentage of marks that a student is expected to score based upon the number of hours they studied. This is a simple linear regression task as it involves just two variables.

Author : Dharani R

```
In [1]: # Importing all libraries required in this notebook
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
```

```
In [2]: # Reading data from remote link
url = "https://raw.githubusercontent.com/AdiPersonalWorks/Random/master/student_scores%20-%20student_scores.csv"
s_data = pd.read_csv(url)
print("Data imported successfully")

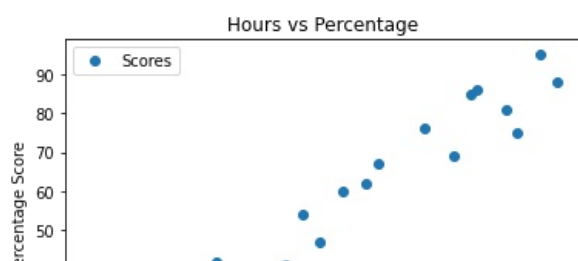
s_data.head(10)
```

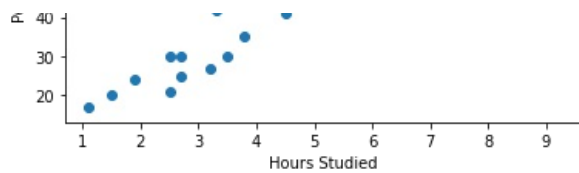
Data imported successfully

```
Out[2]:
```

	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

```
In [3]: # Plotting the distribution of scores
s_data.plot(x='Hours', y='Scores', style='o')
plt.title('Hours vs Percentage')
plt.xlabel('Hours Studied')
plt.ylabel('Percentage Score')
plt.show()
```





From the graph above, we can clearly see that there is a positive linear relation between the number of hours studied and percentage of score.

Preparing the data

The next step is to divide the data into "attributes" (inputs) and "labels" (outputs).

```
In [7]: X = s_data.iloc[:, :-1].values
        y = s_data.iloc[:, 1].values
```

Now that we have our attributes and labels, the next step is to split this data into training and test sets. We'll do this by using Scikit-Learn's built-in `train_test_split()` method:

```
In [8]: from sklearn.model_selection import train_test_split
        X_train, X_test, y_train, y_test = train_test_split(X, y,
                                                            test_size=0.2, random_state=0)
```

Training the Algorithm

We have split our data into training and testing sets, and now is finally the time to train our algorithm.

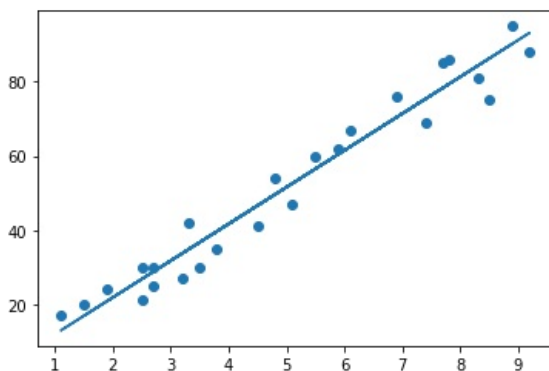
```
In [9]: from sklearn.linear_model import LinearRegression
        regressor = LinearRegression()
        regressor.fit(X_train, y_train)

        print("Training complete.")
```

Training complete.

```
In [10]: # Plotting the regression line
        line = regressor.coef_*X+regressor.intercept_

        # Plotting for the test data
        plt.scatter(X,y)
        plt.plot(X, line);
        plt.show()
```



Making Predictions

Now that we have trained our algorithm, it's time to make some predictions.

```
In [11]: print(X_test) # Testing data - In Hours
```

```
y_pred = regressor.predict(X_test) # Predicting the scores
```

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

```
In [12]: # Comparing Actual vs Predicted  
df = pd.DataFrame({'Actual': y_test, 'Predicted': y_pred})  
df
```

```
Out[12]:
```

	Actual	Predicted
0	20	16.884145
1	27	33.732261
2	69	75.357018
3	30	26.794801
4	62	60.491033

```
In [13]: # You can also test with your own data  
hours = 9.25  
own_pred = regressor.predict([[hours]])  
print("No of Hours = {}".format(hours))  
print("Predicted Score = {}".format(own_pred[0]))
```

```
No of Hours = 9.25  
Predicted Score = 93.69173248737535
```

```
In [14]: from sklearn import metrics  
print('Mean Absolute Error:',  
      metrics.mean_absolute_error(y_test, y_pred))
```

```
Mean Absolute Error: 4.183859899002975
```

```
In [ ]:
```

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js