

TASK 1 - Prediction using Supervised Machine Learning

In the given task, we need to predict the percentage of a student on the basis of number of hours studied using the Linear Regression supervised machine learning algorithm.

Steps Involved

- Step 1 - Importing the dataset
- Step 2 - Visualising the dataset
- Step 3 - Data preparation
- Step 4 - Training the algorithm
- Step 5 - Model Visualisation
- Step 6 - Making predictions
- Step 7 - Evaluating the model

DONE BY: NOTAM KEDARI

STEP 1: Importing The Dataset

In [1]:

Importing all the required libraries

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns

# To ignore the warnings
import warnings as wg
wg.filterwarnings("ignore")
```

In [7]:

Reading data from remote link

```
url = "https://raw.githubusercontent.com/AdiPersonalWorks/Random/master/student_scores%20-%20student_scores.csv"
df = pd.read_csv(url)
```

In [8]:

now let's observe the dataset

```
df.head()
```

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]:

df.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]:

To find the number of columns and rows

```
df.shape
```

Out[10]:

(25, 2)

In [11]:

To find more information about our given dataset

```
df.info()
```

Out[11]:

```
<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]:

df.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

In [13]:

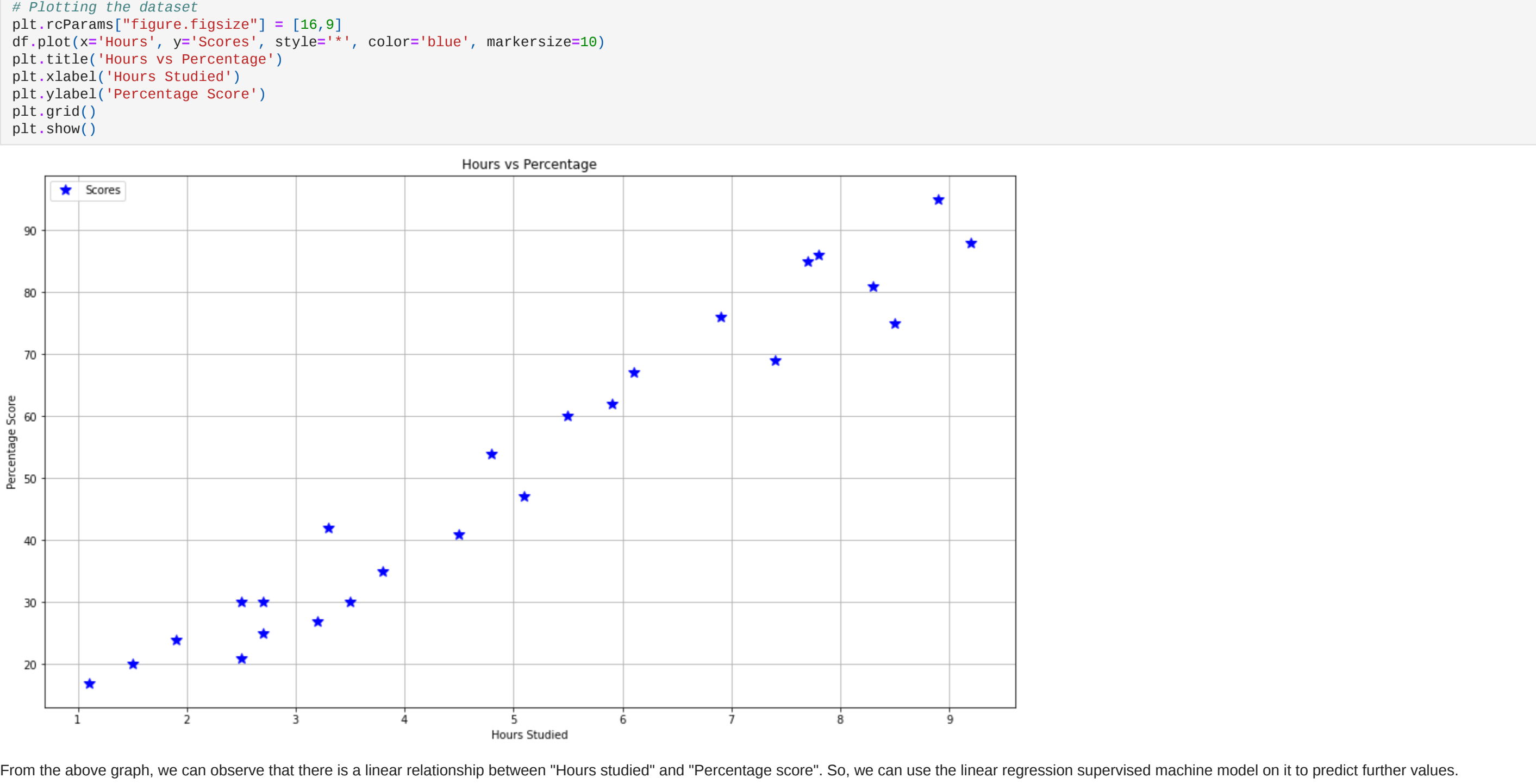
now we will check if our dataset contains null or missings values

```
df.isnull().sum()
```

Out[13]:

```
Hours      0
Scores     0
dtype: int64
```

Step 2: Visualising The Dataset



From the above graph, we can observe that there is a linear relationship between "Hours studied" and "Percentage score". So, we can use the linear regression supervised machine model on it to predict further values.

In [15]:

we can also use .corr to determine the corelation between the variables

```
df.corr()
```

Out[15]:

	Hours	Scores
Hours	1.000000	0.976191
Scores	0.976191	1.000000

STEP 3 - Data Preparation

In [16]:

df.head()

Out[16]:

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

In [17]:

using iloc function we will divide the data

```
X = df.iloc[:, :1].values
y = df.iloc[:, 1:].values
```

In [18]:

X

Out[18]:

```
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],
       [5.9],
       [7.8]])
```

In [19]:

y

Out[19]:

```
array([[21],
       [47],
       [27],
       [75],
       [30],
       [28],
       [88],
       [60],
       [81],
       [25],
       [85],
       [62],
       [41],
       [42],
       [17],
       [95],
       [30],
       [24],
       [67],
       [69],
       [38],
       [54],
       [35],
       [76],
       [86]], dtype=int64)
```

In [20]:

Splitting data into training and testing data

```
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)
```

STEP 4: TRAINING THE ALGORITHM

In [21]:

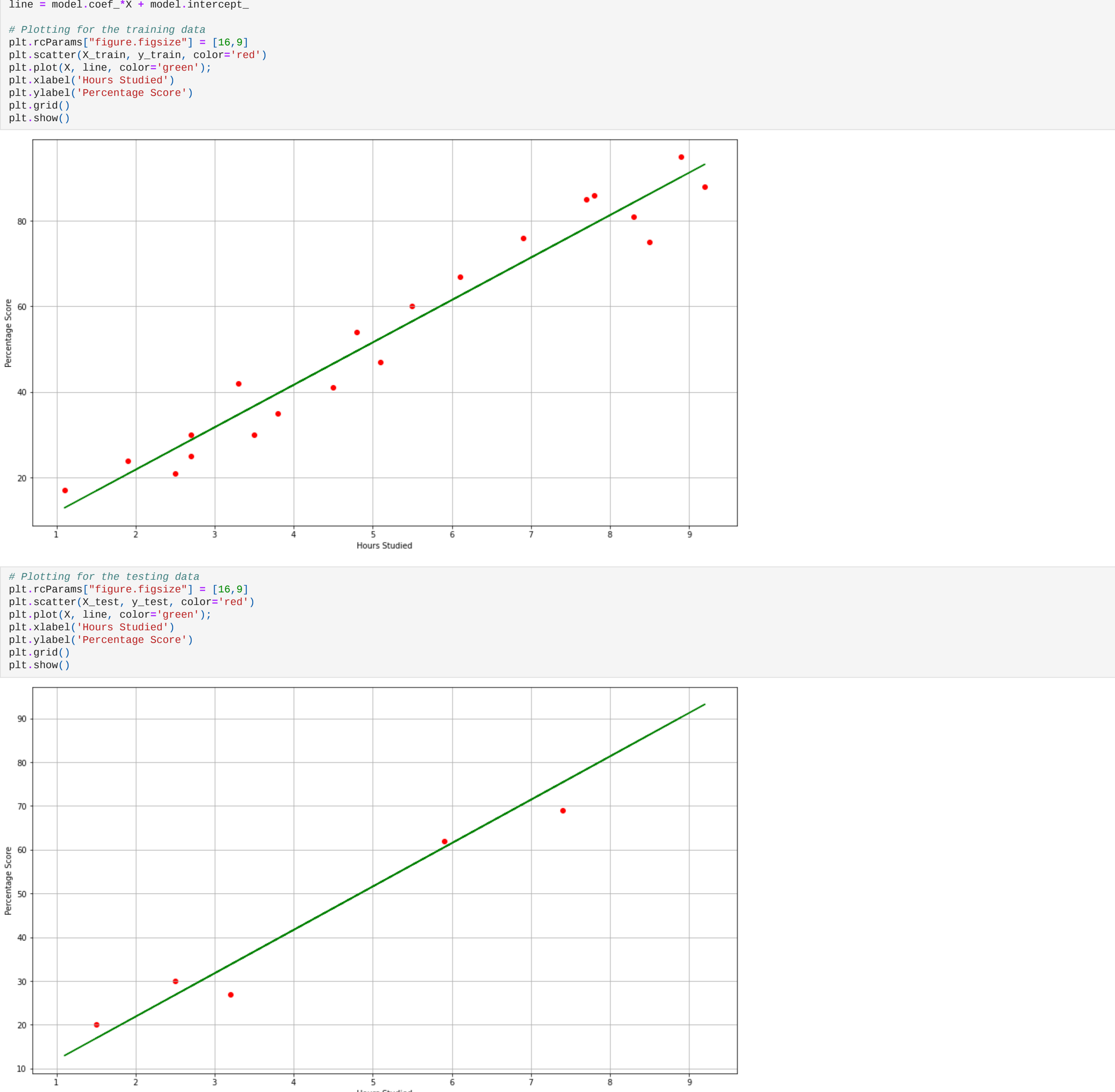
from sklearn.linear_model import LinearRegression

```
model = LinearRegression()
model.fit(X_train, y_train)
```

Out[21]:

LinearRegression()

STEP 5: MODEL VISUALISATION



STEP 6: PREDICTION OF DATA

In [24]:

print(X_test) # Testing data - In Hours

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

Out[24]:

```
[[1.5]
 [3.2]
 [7.4]
 [2.5]
 [5.8]]
```

In [25]:

Comparing Actual vs Predicted

```
y_test
```

Out[25]:

```
array([[28],
       [27],
       [69],
       [38],
       [62]], dtype=int64)
```

In [26]:

y_pred

Out[26]:

```
array([[16.88414476],
       [33.73226878],
       [75.357618 ],
       [26.79480124],
       [68.49183328]])
```

In [27]:

Comparing Actual vs Predicted

```
comp = pd.DataFrame({ 'Actual': [y_test], 'Predicted': [y_pred] })
comp
```

Out[27]:

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

In [28]:

Testing with your own data

```
hours = 9.25
own_pred = model.predict([[hours]])
print("The predicted score if a person studies for", hours, "hours is", own_pred[0])
```

The predicted score if a person studies for 9.25 hours is [93.69173249]

Hence, it can be concluded that the predicted score if a person studies for 9.25 hours is 93.69173248737538

STEP 7: EVALUATING THE MODEL

In [29]:

from sklearn import metrics

```
print('Mean Absolute Error:', metrics.mean_absolute_error(y_test, y_pred))
```

Out[29]:

```
Mean Absolute Error: 4.18385989902975
```

DONE BY: NOTAM KEDARI

In []:

Out[]: