

# GRIP : THE SPARKS FOUNDATION

Data Science and Business Analytics

Task 1 : prediction using supervised ML

In this task i will predict the percentage score of a student based on the number of hours studied. In this task, i have used two variables where the feature is the no. of hours studied and target value is the percentage score. This can be achieved with the help of Linear Regression.

## importing required libraries

```
In [46]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
```

```
In [47]: # reading the data
data_link = "stud.csv"
data = pd.read_csv(data_link)
data.head(27)
```

Out[47]:

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

In [48]: `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
```

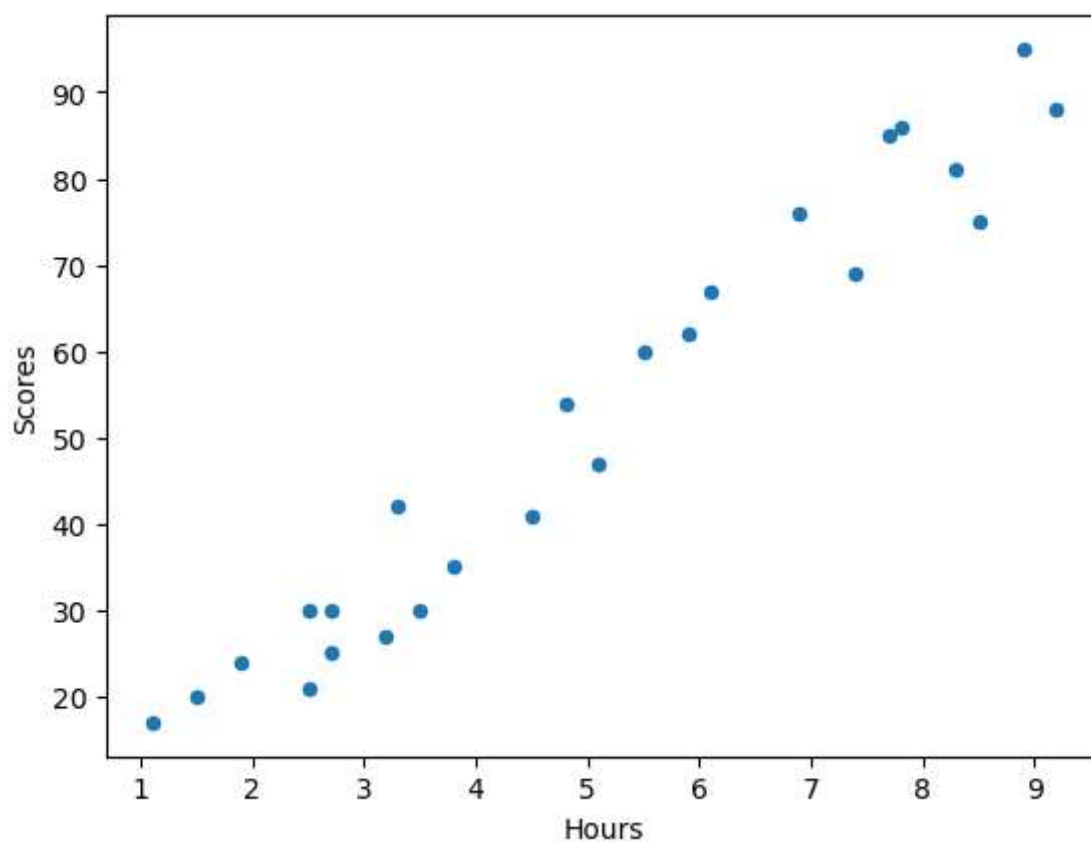
In [49]: `data.describe()`

Out[49]:

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

```
In [12]: import matplotlib.pyplot as plt  
%matplotlib inline
```

```
In [50]: data.plot(kind= 'scatter', x="Hours", y="Scores");  
plt.show()
```



```
In [15]: data.corr(method= "pearson")
```

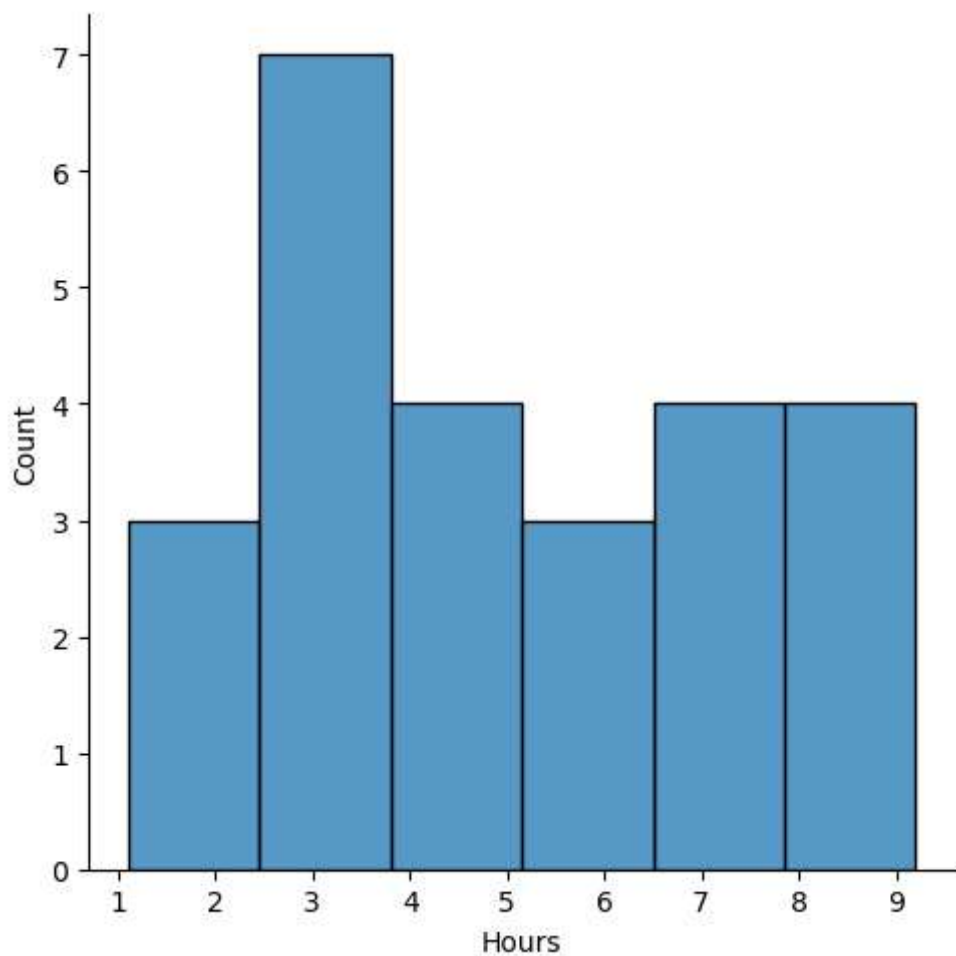
Out[15]:

	Hours	Scores
<b>Hours</b>	1.000000	0.976191
<b>Scores</b>	0.976191	1.000000

```
In [51]: hours = data["Hours"]  
scores = data["Scores"]
```

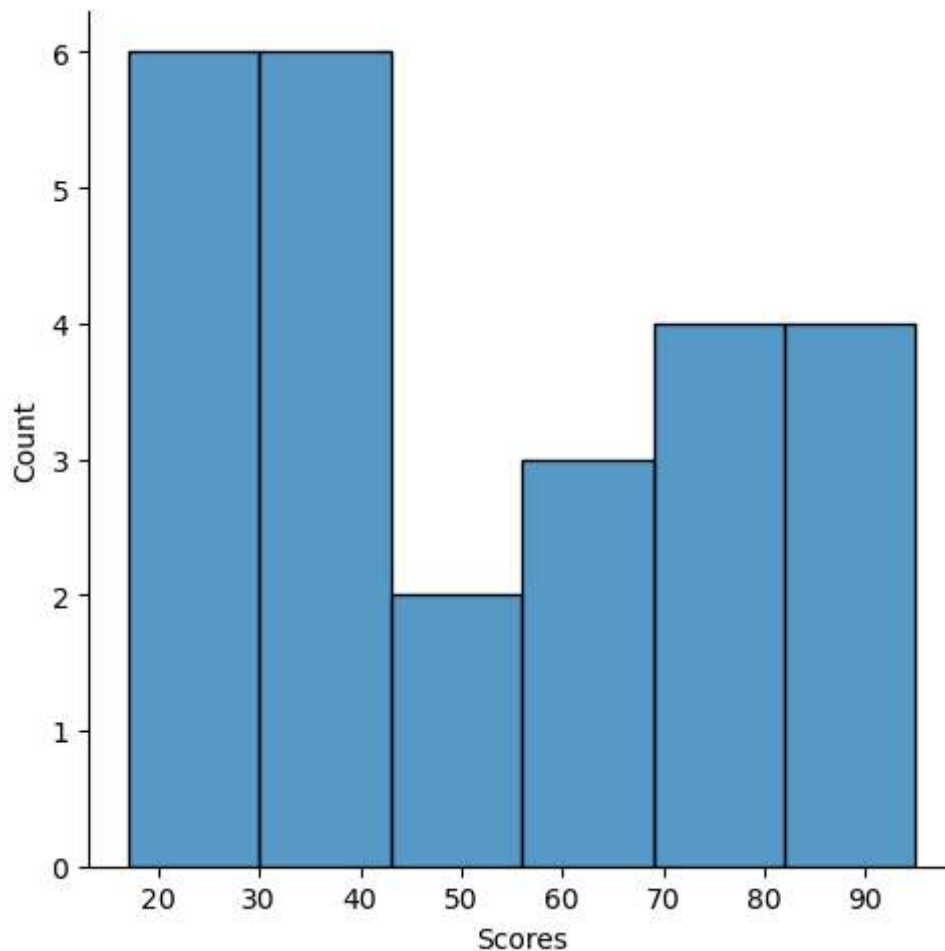
```
In [52]: sns.displot(hours)
```

```
Out[52]: <seaborn.axisgrid.FacetGrid at 0x1d90b97ad90>
```



```
In [53]: sns.displot(scores)
```

```
Out[53]: <seaborn.axisgrid.FacetGrid at 0x1d90cfe88e0>
```



## Linear Regression

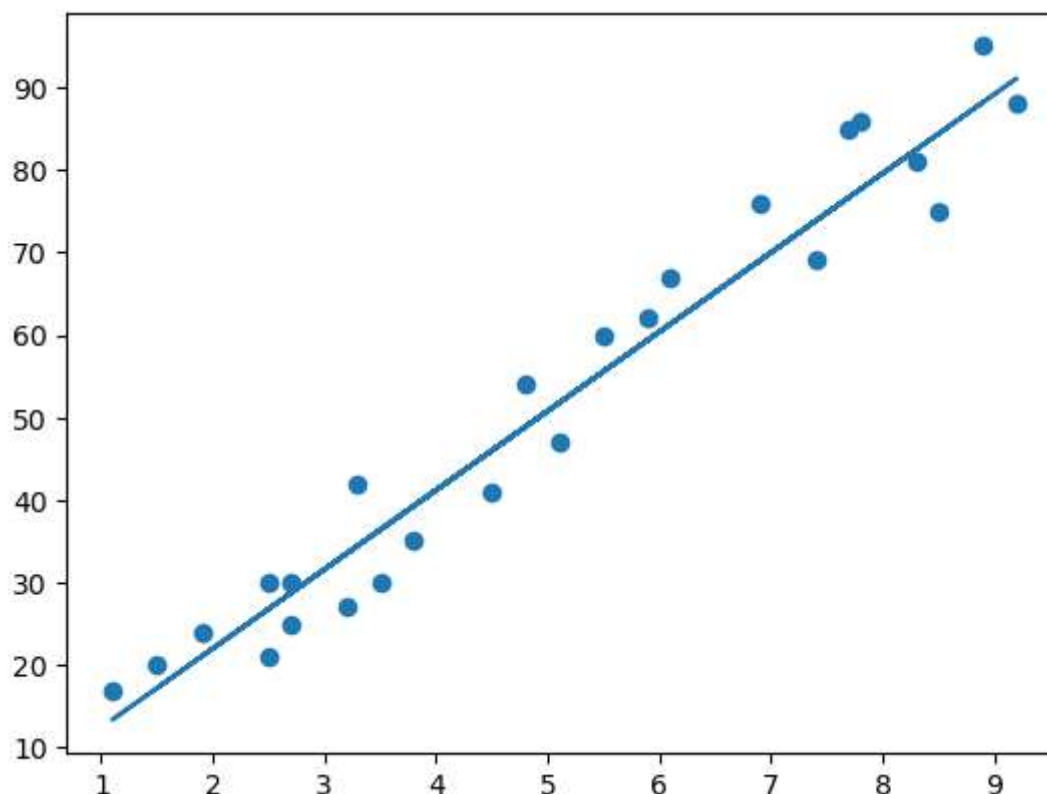
```
In [58]: a = data.iloc[:, :-1].values  
        b = data.iloc[:, 1].values
```

```
In [76]: from sklearn.model_selection import train_test_split  
        a_train, a_test, b_train, b_test = train_test_split(a, b, test_size=0.2, random_state=50)
```

```
In [60]: from sklearn.linear_model import LinearRegression  
        reg = LinearRegression()  
        reg.fit(a_train, b_train)
```

```
Out[60]: LinearRegression()
```

```
In [67]: m = reg.coef_           #coef_ and intercept_ are attributes of LinearRegression  
        c = reg.intercept_  
        line = m*a + c  
  
        plt.scatter(a, b)  
        plt.plot(a, line)  
        plt.show()
```



```
In [68]: b_pred = reg.predict(a_test)
```

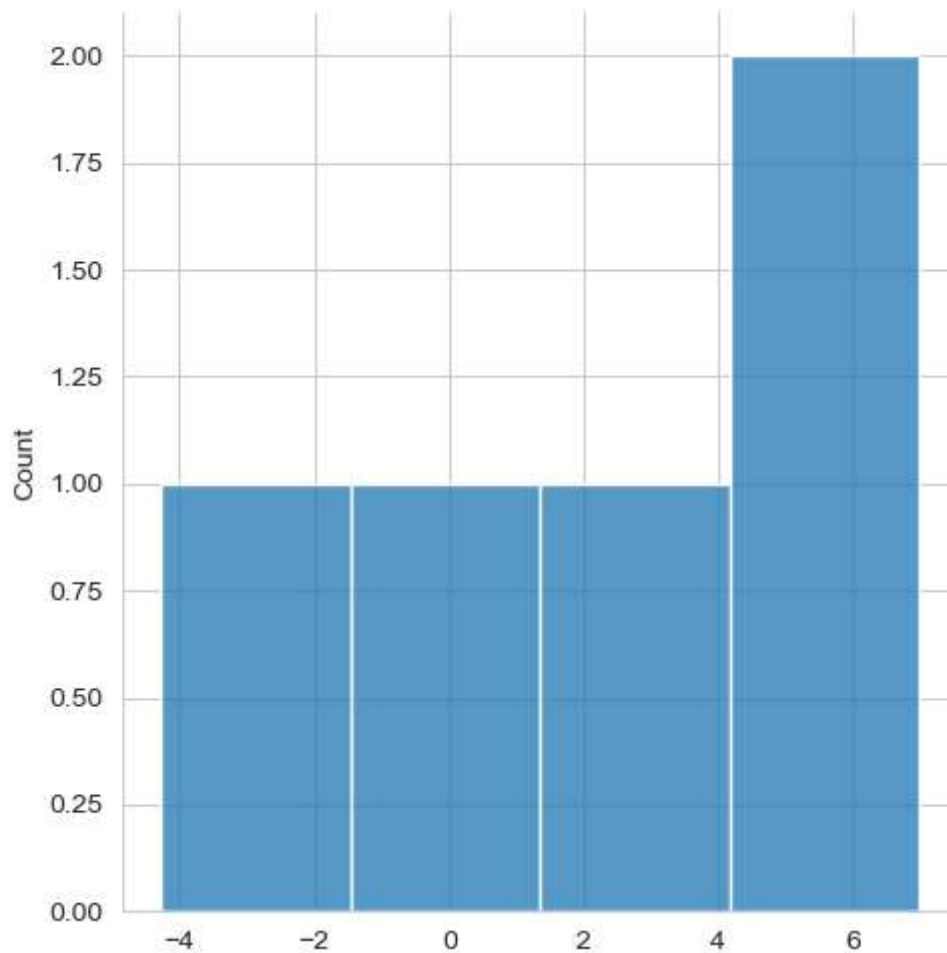
```
In [84]: actual_pred = pd.DataFrame({"target": b_test, "predicted":b_pred})  
actual_pred
```

```
Out[84]:
```

	target	predicted
0	95	88.211394
1	30	28.718453
2	76	69.020122
3	35	39.273652
4	17	13.365436

```
In [85]: sns.set_style("whitegrid")  
sns.displot(np.array(b_test - b_pred))
```

```
Out[85]: <seaborn.axisgrid.FacetGrid at 0x1d90e0ce9d0>
```



**what would be the predicted score of a student if he/she studies for 9.25 hours/day?**

```
In [90]: h = 9.25
s = reg.predict([[h]])
print("if a student studies {} hours a day then he/she will score {} % in exam".format(h, s))

if a student studies 9.25 hours a day then he/she will score [91.56986604] % in exam
```