

▼ Name: Gore Aniket Machhindra

Email: aniket.m.gore.1901@gmail.com

Role: Data Science and Business Analytics Intern

Task 1: Prediction using Supervised ML

```
1 # Importing all libraries required in this notebook
2 import pandas as pd
3 import numpy as np
4 import numpy as np
5 import seaborn as sns
6 import matplotlib.pyplot as plt
7 %matplotlib inline
```

```
1 url = "http://bit.ly/w-data"
2 data = pd.read_csv(url)
3 print("Data imported successfully")
4 data.head()
```

☞ Data imported successfully

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

```
1 data.tail()
```

Hours Scores



```
1 data.shape
```

```
(25, 2)
```

```
22      3.8      35
```

```
1 data.describe()
```

	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



```
1 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
```

```
1 data.corr()
```

	Hours	Scores
Hours	1.000000	0.976191
Scores	0.976191	1.000000



```
1 #Collecting X and Y
2 X=data['Hours'].values
3 Y=data['Scores'].values
```

```

1 #Mean X and Y
2 mean_x=np.mean(X)
3 mean_y=np.mean(Y)
4 #Total number of values
5 n=len(X)
6 #Using the formula to calculate b1 and b2
7 numer=0
8 denom=0
9 for i in range(n):
10     numer+=(X[i]-mean_x)*(Y[i]-mean_y)
11     denom+=(X[i]-mean_x)**2
12 b1=numer/denom
13 b0=mean_y-(b1*mean_x)
14 #Print coefficients
15 print(b1,b0)

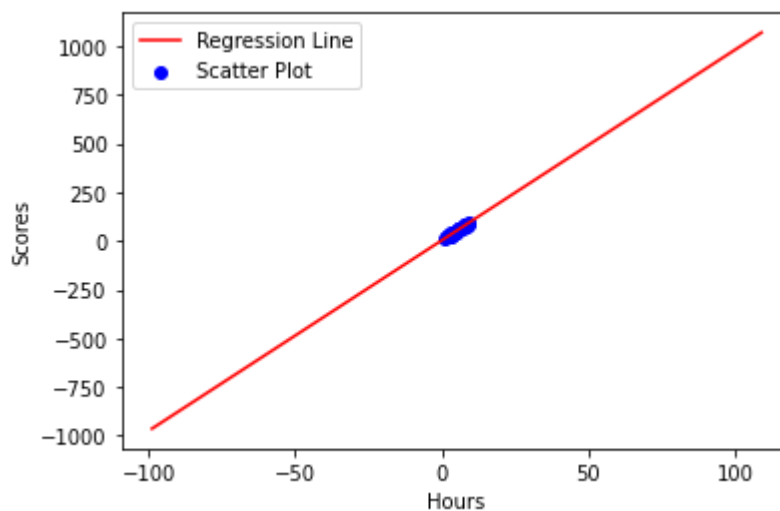
```

9.775803390787475 2.4836734053731746

```

1 #Plotting Values and Regression Line
2 max_x=np.max(X)+100
3 min_x=np.min(X)-100
4 #Calculating line values x and y
5 x=np.linspace(min_x,max_x,1000)
6 y=b0+b1*x
7 #Plotting Line
8 plt.plot(x,y,color='red',label='Regression Line')
9 #Plotting Scatter Points
10 plt.scatter(X,Y,c='blue',label='Scatter Plot')
11 plt.xlabel('Hours')
12 plt.ylabel('Scores')
13 plt.legend()
14 plt.show()

```



```

1 ss_t=0
2 ss_r=0
3 for i in range(n):

```

```

4 y_pred=b0+b1*X[i]
5 ss_t+=(Y[i]-mean_y)**2
6 ss_r+=(Y[i]-y_pred)**2
7 r2=1-(ss_r/ss_t)
8 print(r2)

```

0.9529481969048356

```

1 #Creating Model using Scikit Learn Library
2 from sklearn.linear_model import LinearRegression
3 from sklearn.metrics import mean_squared_error
4 #Cannot use Rank 1 matrix in scikit Learn
5 X=X.reshape((n,1))
6 #Creating Model
7 reg=LinearRegression()
8 #Fitting training data
9 reg=reg.fit(X,Y)
10 #Y Prediction
11 Y_pred=reg.predict(X)
12 #Calculating R2 Score
13 r2_score=reg.score(X,Y)
14 print(r2_score)

```

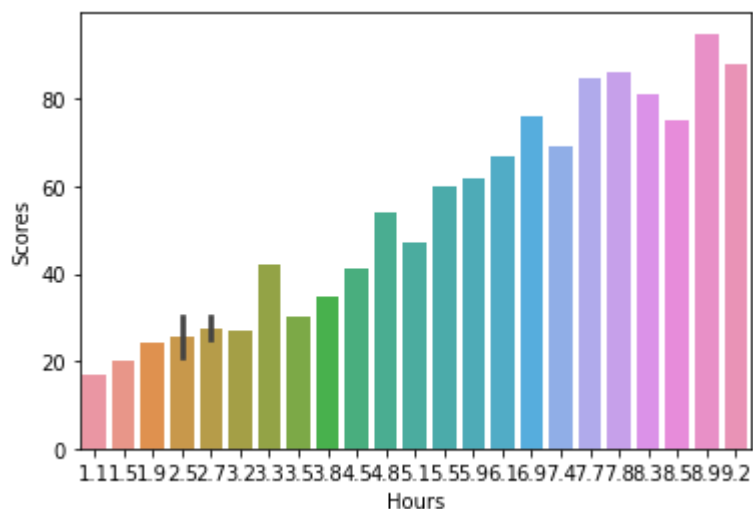
0.9529481969048356

```
1 sns.barplot(data['Hours'], data['Scores'])
```

```

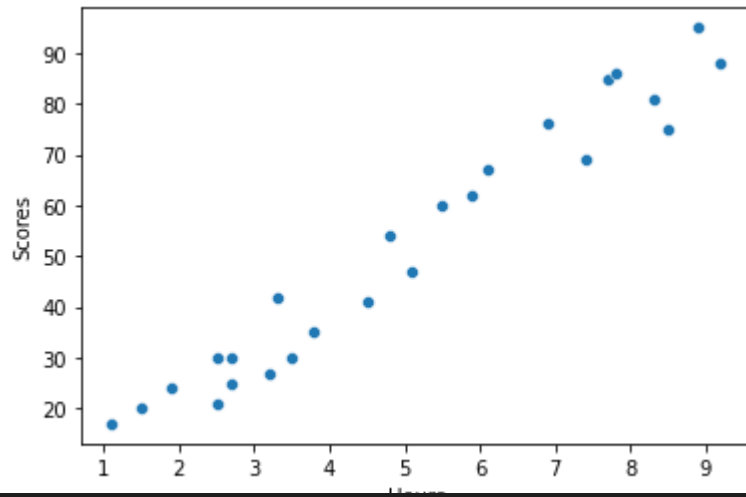
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7ffa51262a50>

```

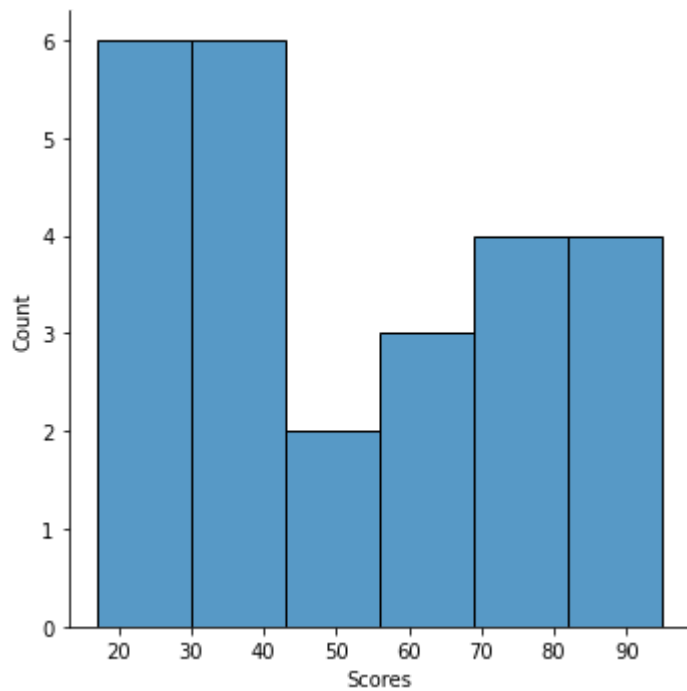


```
1 sns.scatterplot(data['Hours'], data['Scores'])
```

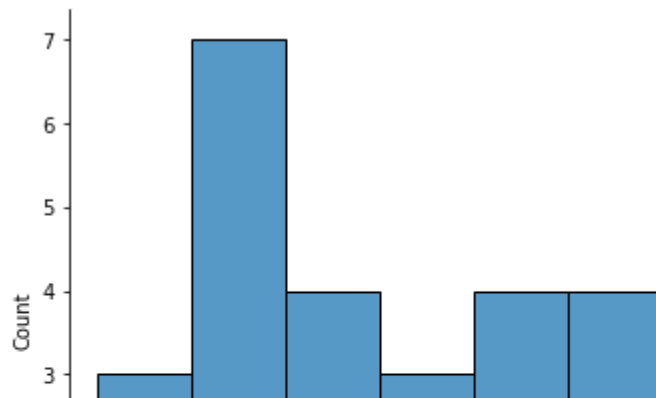
```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7ffa516c2990>
```



```
1 sns.displot(data['Scores'])
2 plt.show()
```

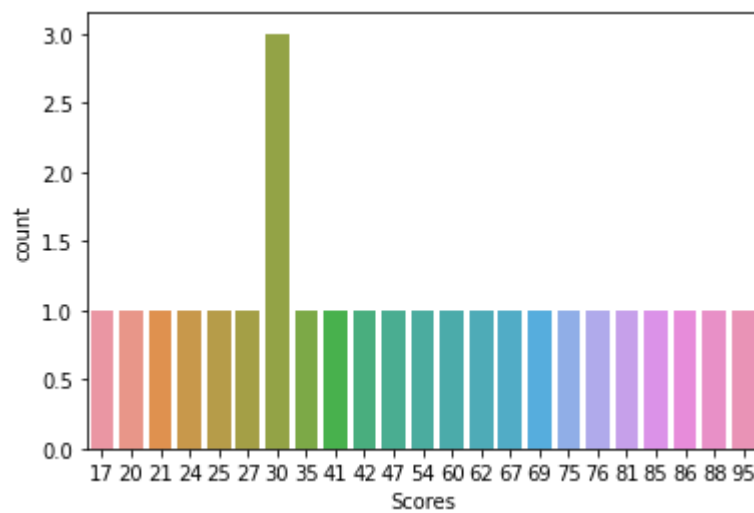


```
1 sns.displot(data['Hours'])
2 plt.show()
```



```
1 sns.countplot(data['Scores'])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword arguments: {'x': 'Scores', 'y': 'count'}. This warning will be removed in a future version of Seaborn.
FutureWarning
<matplotlib.axes._subplots.AxesSubplot at 0x7ffa4f35ca90>
```



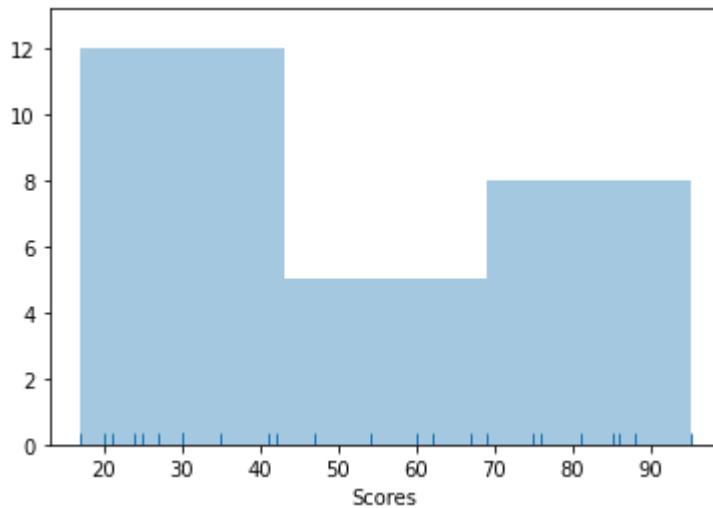
```
1 sns.countplot(data['Hours'])
```

```
2
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword arguments: {\"x\": 0, \"y\": 1} instead of positional arguments. This will result in an error in the future.  
FutureWarning  
<matplotlib.axes._subplots.AxesSubplot at 0x7ffa469f7250>
```

```
1 sns.distplot(data['Scores'],kde=False,rug=True)  
2 plt.show()
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2619: FutureWarning: `distplot` is a deprecated alias for `displot`. Please adjust your code to use `displot` instead.  
warnings.warn(msg, FutureWarning)  
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2103: FutureWarning: The ` rug` parameter is deprecated in favor of the `rug` parameter of the `FacetGrid` class.  
warnings.warn(msg, FutureWarning)
```



```
1 sns.jointplot(data['Hours'],data['Scores'],kind="reg")  
2 plt.show()
```

100

20 1

20 1



92.90985477015732

```
_pred=list(data['predicted_Scores'].values)
```

```
mse=(np.sqrt(s/len(data)))/mean_y
```



```
3 rmse
```

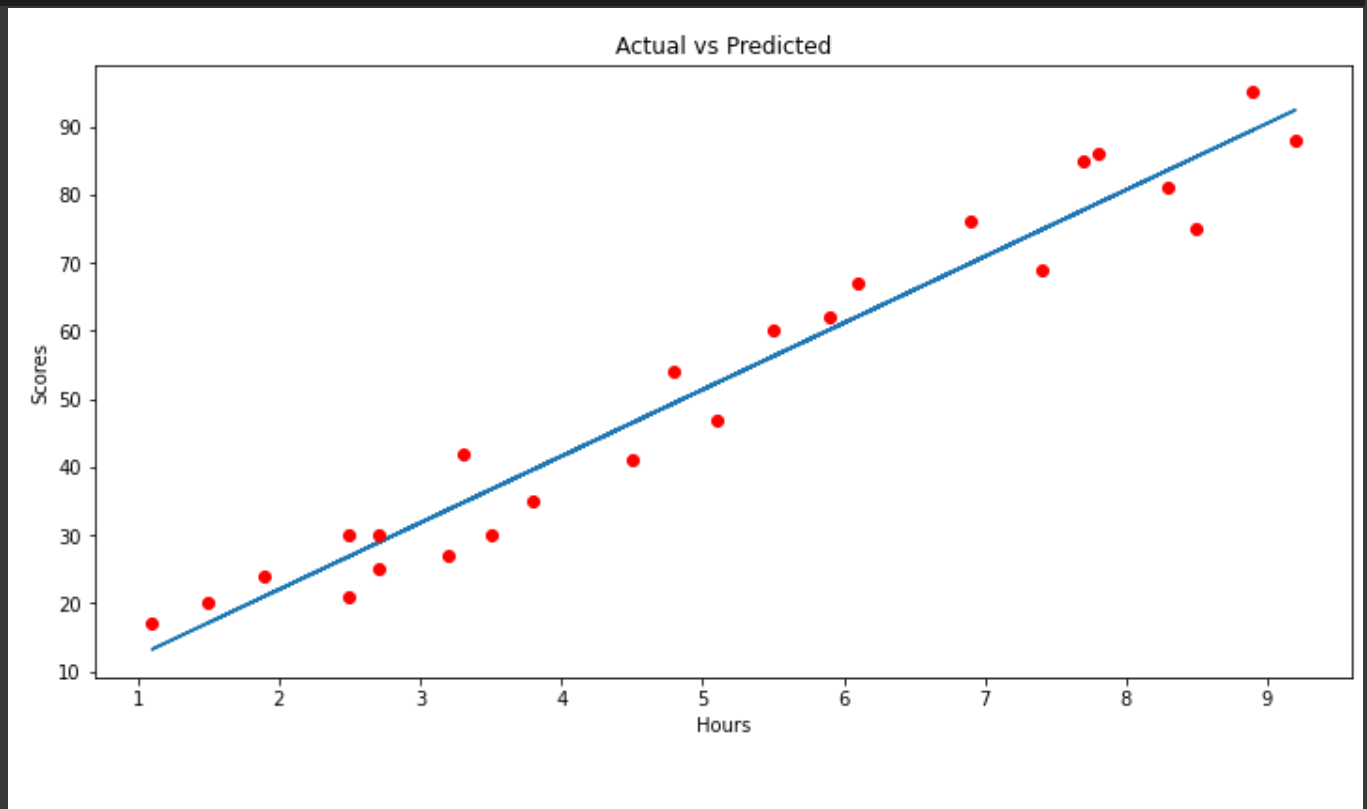
```
0.10439521325937494
```

```
1 import statsmodels.formula.api as smf
2 model=smf.ols('Scores ~ Hours',data=data)
3 model=model.fit()
```

```
/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning:
import pandas.util.testing as tm
```

```
1 data['pred_ols']=model.predict(data['Hours'])
```

```
1 plt.figure(figsize=(12,6))
2 plt.plot(data['Hours'],data['pred_ols']) #Regression line
3 plt.plot(data['Hours'],data['Scores'],'ro') #Scatter plot showing actual data
4 plt.title('Actual vs Predicted')
5 plt.xlabel('Hours')
6 plt.ylabel('Scores')
7 plt.show()
```



We can observe that the predicted value for 9.25 hours is around 92

✓ 0s completed at 5:37 PM

