

▼ To find the best fit line for given data using Linear Regression

▼ Assignment on Linear Regression:

The following table shows the results of a recently conducted study on the correlation of the number of hours spent driving with the risk of developing acute backache. Find the equation of the best fit line for this data.

Number of hours spent driving (x)	Risk Score on a scale of 0-100 (y)
10	95
9	80
2	10
15	50
10	45
16	98
11	38
16	93

```
1 # Importing libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import seaborn as sns
6 from sklearn.linear_model import LinearRegression
```

```
1 # Function to plot the given data
2 def plotData(X, y, x_label='driving_hours', y_label='risk_score'):
3     plt.figure(figsize=(6,6))
4     plt.xlabel(x_label)
5     plt.ylabel(y_label)
6     plt.scatter(X,y)
7     return plt
```

```
1 # Function to plot the regression line along with the given data
2 def plotRegressionLine(X, y, y_pred, x_label='driving_hours', y_label='risk_score'):
3     plt = plotData(X, y, x_label, y_label)
4     plt.plot(X, y_pred, color='red', linewidth=3)
5     plt.show()
```

▼ Evaluating the model - Coefficient of Determination (R2 Score)

$$\text{Coefficient of Determination} \rightarrow R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

$$\text{Sum of Squares Total} \rightarrow SST = \sum (y - \bar{y})^2$$

$$\text{Sum of Squares Regression} \rightarrow SSR = \sum (y' - \bar{y}')^2$$

$$\text{Sum of Squares Error} \rightarrow SSE = \sum (y - y')^2$$

```

1 # Function to calculate R2 score
2 def calcR2Score(X, y, y_pred):
3     y_mean = y.mean()
4     SStot = sum((y-y_mean)**2)
5     SSres = sum((y-y_pred)**2)
6     r2_score = 1-(SSres/SStot)
7     return r2_score

```

Pearson's Correlation Coefficient is a linear correlation coefficient that returns a value of between -1 and +1. A -1 means there is a strong negative correlation and +1 means that there is a strong positive correlation. A 0 means that there is no correlation (this is also called zero correlation)

```

1 # Function to calculate Pearson's Correlation Coefficient
2 def correlationCoef(X,y):
3     X_mean = X.mean()
4     y_mean = y.mean()
5     num = sum((X-X_mean)*(y-y_mean))
6     den = (sum((X-X_mean)**2)*sum((y-y_mean)**2))**0.5
7     coef = num/den
8     return coef

```

```

1 data = {
2     'driving_hours': [10, 9, 2, 15, 10, 16, 11, 16],
3     'risk_score':[95, 80, 10, 50, 45, 98, 38, 93]
4 }
5 df = pd.DataFrame.from_dict(data)

```

```
1 df.head()
```

	driving_hours	risk_score
0	10	95
1	9	80
2	2	10
3	15	50
4	10	45

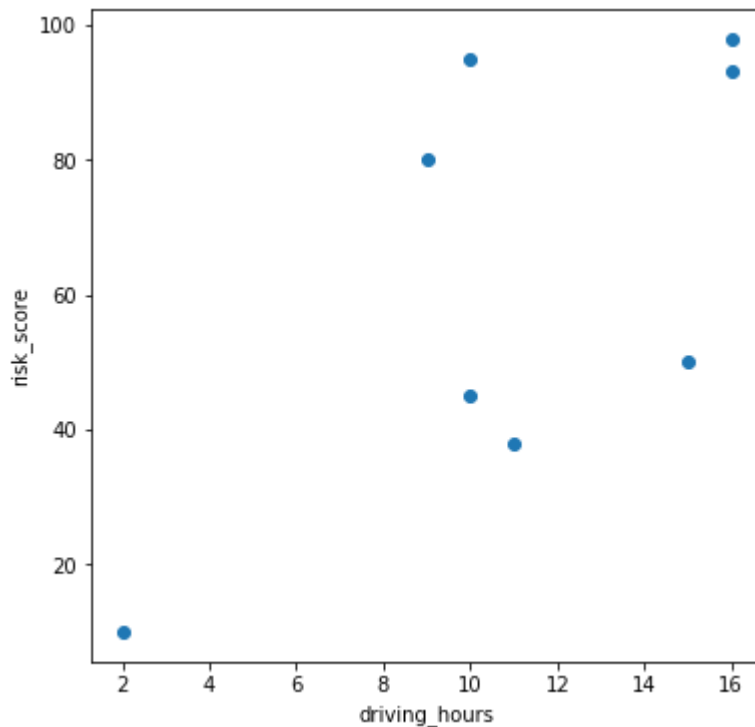
```
1 X = np.array(df['driving_hours'])
2 y = np.array(df['risk_score'])
```

```
1 correlationCoef(X,y)
```

```
0.6611314653759117
```

```
1 plotData(X, y)
```

```
<module 'matplotlib.pyplot' from 'C:\\Users\\Megh Khaire\\Anaconda3\\envs\\Pytho
```



▼ Fitting the Regression Line

$$\hat{Y} = b_0 + b_1 x$$

$$\beta_1 = \frac{\sum_{i=1}^m (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^m (x_i - \bar{x})^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

```
1 # Function to calculate coefficients of the line
2 def calcCoefficient(X,y):
3     X_mean = X.mean()
4     y_mean = y.mean()
5     # Σ [ (xi - x)(yi - y) ] / Σ [ (xi - x)2]
6     coef = (sum((X-X_mean)*(y-y_mean))/sum((X-X_mean)**2))
7     intercept = y_mean - coef*X_mean
8     return coef, intercept
```

```
1 W,W0 = calcCoefficient(X,y)
2 y_pred = W*X + W0
3 print("Coefficient: ",W)
4 print("Intercept: ",W0)
5 print("R2 score: ",calcR2Score(X, y, y_pred))
6 print(f"Equation: Y = {W}*X+{W0}")
7 plotRegressionLine(X, y, y_pred)
```

Coefficient: 4.58789860997547
 Intercept: 12.584627964022893
 R2 score: 0.43709481451010035

▼ Sklearn implementation

Using Sklearn's LinearRegression api for evaluating the given dataset

```

1 X = np.array(X).reshape(-1,1)
2
3 model = LinearRegression()
4 model.fit(X, y)
5 y_pred = model.predict(X)
6 print("Coefficient: ",model.coef_[0])
7 print("Intercept: ",model.intercept_)
8 print("R2 score: ",calcR2Score(X, y, y_pred))
9 print(f"Equation: Y = {model.coef_[0]}*X+{model.intercept_}")
10 plotRegressionLine(X, y, y_pred)

```

Coefficient: 4.587898609975469
 Intercept: 12.584627964022907
 R2 score: 0.43709481451010035
 Equation: Y = 4.587898609975469*X+12.584627964022907

