```python
In [3]:  import pandas as pd
```

That line of code is to import Pandas for me to be able to use his functions to explore, manipulate, analyze and clean the dataset.

```python
In [4]:  df = pd.read_csv('Student_Performance.csv')
```

Pandas function read_csv() is to load the dataset in Jupyter Notebook.

```python
In [5]:  duplicates = df.duplicated()
         print(duplicates)

0        False
1        False
2        False
3        False
4        False
         ...
9995     False
9996     False
9997     False
9998     False
9999     False
Length: 10000, dtype: bool
```

As we can see, there are no duplicates in the Dataset.

```python
In [6]:  miss_val = df.isnull().any(axis=1)
         print(miss_val)

0        False
1        False
2        False
3        False
4        False
         ...
9995     False
9996     False
9997     False
9998     False
9999     False
Length: 10000, dtype: bool
```

```python
In [11]: df['Extracurricular Activities'] = df['Extracurricular Activities'].replace({'Yes':1, 'No':0})
```

As we can see, there is no row with a field column missing value.

```python
In [12]: import numpy as np
         import sklearn
         from sklearn.linear_model import LinearRegression
```

#MULTIPLE LINEAR REGRESSION

I'm importing numpy library because, in the implementation of the multiple linear regression, I'll need it to perform mathematical operations such as vector operations and matrix manipulations. Moreover, the dependent variable and independent variables are represented as NumPy arrays in scikit-learn and Pandas, so we need to use NumPy to operate on these arrays within our model.

Now I'll be selecting the independent variables also called features and the dependent variable(target varaiable).

```python
In [13]: X = df[['Hours Studied','Previous Scores','Extracurricular Activities','Sleep Hours']]
         Y = df['Performance Index']
```

```python
In [ ]:
```

Create and fit the linear regression model

```python
In [14]: model = LinearRegression()
         model.fit(X, Y)

Out[14]: LinearRegression()
```

Make predictions

```python
In [15]: predictions = model.predict(X)
```

Print model summary (coefficients, R-squared, etc.)

```python
In [17]: print(model.coef_)   # Coefficients
         print(model.intercept_)   # Intercept
         print(model.score(X, Y))  # R-squared value

[2.85673907 1.01868958 0.62741988 0.48194429]
-33.24005596528127
0.9879162180086278
```

```python
In [ ]: The coefficients* represent the slope/tilt of the linear relationship between the dependent variable and each one of the
        independent variables.
        The intercept represent the average expected value for the dependent variable when all independent variables are equal to Zero.
        The R-values associated with each coefficient indicate whether the relationship between that predictor variable and the response is statistically significant.
        A low p-value (typically < 0.05) suggests a significant relationship, while a high p-value suggests no significant effect.
```

```python
In [ ]:                              #SIMPLE LINEAR REGRESSION
```

Create the model

```python
In [19]: model = LinearRegression()
```

Fit the model (X: explanatory variable; y: response variable)

```python
In [20]: X = df[['Hours Studied']]
         y = df['Sleep Hours']
         model.fit(X, y)

Out[20]: LinearRegression()
```

Get coefficients (intercept and slope)

```python
In [21]: intercept = model.intercept_
         slope = model.coef_[0]

         print(f"Intercept (b0): {intercept:.2f}")
         print(f"Slope (b1): {slope:.2f}")

Intercept (b0): 6.53
Slope (b1): 0.00
```

```python
In [ ]:
```