<pre>from sklearn from sklearn import warni warnings.fil</pre>	as as pd y as np lotlib.pyplot as plt orn as sns n.model_selection import train_test_split n.preprocessing import LabelEncoder n.linear_model import LinearRegression n.metrics import mean_squared_error, r2_score ings lterwarnings('ignore') mdas as pd:
This adds the import num This adds the a	name pd to the pandas library. It is used to manipulate and analyze data, especially when working with tabular data (DataFrames). mpy as np: alias np to the numpy library, which is used for numerical computation. It is employed for mathematical computations and array manipulation. atplotlib.pyplot as plt: natplotlib.pyplot, a tool for making plots and visualizations, including scatter plots, graphs, and histograms.
Imports seabor from sklea This function d from the p	aborn as sns: orn, a Matplotlib-based statistical data visualization package that makes it easier to create visually appealing and educational plots. arn.model_selection import train_test_split: divides your dataset into two sets—one for model training and one for testing—is imported from Scikit-learn. oreprocessing sklearn import LabelEncoder: learn's LabelEncoder, a tool for encoding categorical variables (such strings or labels) as numerical values.
To build and tra	inear model in sklearn import LinearRegression: rain a linear regression model, import the Scikit-learn LinearRegression class. arn.metrics import r2_score, mean_squared_error: the r2_score and mean_squared_error functions. These measures are employed to assess a regression model's performance. rnings:
Reading df = pd.read print("\nThe df.shape	the csv file to Dataframe d_csv('Student_Performance.csv') e number of rows and columns in the dataset are:") E rows and columns in the dataset are:
Printig the df.head() Hours Studie	e has 10000 rows and six columns e first 5 data of the dataframe fied Previous Scores Extracurricular Activities Sleep Hours Sample Question Papers Practiced Performance Index 7 99 Yes 9 1 1 91.0 4 82 No 4 2 65.0
df.info() <class #="" 'panda="" 1="" column<="" columns="" data="" rangeindex:="" td=""><td>8 51 Yes 7 2 45.0 5 52 Yes 5 2 36.0 7 75 No 8 5 5 66.0 score.frame DataFrame State Frame State State</td></class>	8 51 Yes 7 2 45.0 5 52 Yes 5 2 36.0 7 75 No 8 5 5 66.0 score.frame DataFrame State Frame State
3 Sleep Ho 4 Sample Q 5 Performa dtypes: float memory usage: Perspectiv	S Scores 10000 non-null int64 rricular Activities 10000 non-null object ours 10000 non-null int64 Question Papers Practiced 10000 non-null int64 ance Index 10000 non-null float64 c64(1), int64(4), object(1)
and the continue of df.describe (Hours S count 10000.00 mean 4.9	Studied Previous Scores Sleep Hours Sample Question Papers Practiced Performance Index
25% 3.0 50% 5.0 75% 7.0 max 9.0	40.00000 40.00000 4.000000 0.00000 10.00000 10.00000 54.00000 54.00000 5.00000 2.00000 40.00000 69.00000 85.00000 8.00000 7.00000 7.00000 71.00000 000000 99.00000 9.00000 9.00000 9.00000 100.00000 ws interesting trends, such as a moderate correlation between the hours studied, the previous scores, and the number of sample papers practiced. For instance, students who studied longer or practiced more papers tend to have higher performance in and deviation indicates significant variation in performance despite similar efforts. Sleep seems to be somewhat consistent, with most students averaging between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night. It might be worth exploring whether there's any correlation between 5 and 8 hours per night.
df.isna().su Hours Studie Previous Sco Extracurricu Sleep Hours Sample Quest Performance	g the unwanted and missing values um() ed 0 ores 0 ular Activities 0 tion Papers Practiced 0 Index 0
df.duplicate 127 The dataframe construction.	s: There is no need to impute (fill in) any missing data or remove any rows or columns because of missing values because there are none. Ready for Modeling: By removing one of the most frequent preprocessing procedures, the lack of missing data data for model training. Data Integrity: The dataset is clean if there are no missing values, which is crucial for building strong machine learning models.
print("\nThe df.shape The number of (9873, 6) It appears that	e number of rows after removing the duplicate values is:") E rows after removing the duplicate values is: t your dataset has shrunk from 10,000 rows to 9,873 rows after the duplicate rows were eliminated. This indicates that 127 duplicate rows were effectively eliminated from the original dataset. 0,000 rows were initially present. 9,873 rows after duplicates are eliminated There are 10,000 - 9,873 = 127 duplicate rows.
<pre>import matpl import seabo # List of nu numeric_colu # Create box plt.figure(f for i, colum plt.subp</pre>	umeric columns umns = ['Hours Studied', 'Previous Scores', 'Sleep Hours', 'Sample Question Papers Practiced', 'Performance Index'] xxplots for each numeric column figsize=(15, 10)) mn in enumerate(numeric_columns, 1): plot(2, 3, i)
	Boxplot of Hours Studied Boxplot of Previous Scores Boxplot of Sleep Hours
1 2 Boxpl	3 4 5 6 7 8 9 40 50 60 70 80 90 100 4 5 6 7 8 9 Hours Studied Previous Scores Boxplot of Performance Index
<pre>def detect_o outliers for colu # Ca</pre>	umn in columns: alculate Q1 (25th percentile) and Q3 (75th percentile)
Q1 = Q3 = IQR # Ca lowe uppe # Fi outl return of	<pre>= df[column].quantile(0.25) = df[column].quantile(0.75) = Q3 - Q1 alculate lower and upper bounds for outliers er_bound = Q1 - 1.5 * IQR er_bound = Q3 + 1.5 * IQR ind the outliers liers[column] = df[(df[column] < lower_bound) (df[column] > upper_bound)] outliers tliers for numeric columns</pre>
<pre>outliers = d # Print the for column, print(f" Number of out Number of out Number of out Number of out Number of out</pre>	<pre>detect_outliers_iqr(df, numeric_columns) number of outliers detected for each feature outlier_data in outliers.items(): "Number of outliers in {column}: {len(outlier_data)}") cliers in Hours Studied: 0 cliers in Previous Scores: 0 cliers in Sleep Hours: 0 cliers in Sample Question Papers Practiced: 0 cliers in Performance Index: 0</pre>
from sklearn # Initialize encoder = La # Encode the	the box plot and IQR approach, no outliers were found in the dataset, which is fantastic because it indicates that the data is already well-behaved and free of extreme or aberrant values that could distort the research. This suggests that your dataset is already solves, both of which you have already verified. g the categorical value n.preprocessing import LabelEncoder s the LabelEncoder abelEncoder abelEncoder() s 'Extracurricular Activities' column rricular Activities') = encoder.fit_transform(df('Extracurricular Activities'))
0 1 2 3	sied Previous Scores Extracurricular Activities Sleep Hours Sample Question Papers Practiced Performance Index 7 99 1 9 1 91.0 4 82 0 4 2 65.0 8 51 1 7 2 45.0 5 52 1 5 2 36.0 7 75 0 8 5 66.0
# Save the udf.to_csv('C	e cleaned csv file updated DataFrame to a new CSV file (optional) cleaned_csv_file.csv', index=False) ing the Correlation Matrix the correlation matrix for all features and the target variable matrix = df.corr()
print ("Correprint (correlation we hours Studied Previous Score Extracurricul Sleep Hours Sample Questi Performance I	0.915135 Lar Activities 0.026075 0.050352 ion Papers Practiced 0.043436
Prior Scores The Performant Scores. The performant Hours Studies The Performant	Takeaways for correlation between the x variables and the target y variable is (0.915135): Ince Index, the target variable, and this attribute have the strongest positive association. A very strong positive link is shown by a correlation of 0.92. This indicates that one of the most crucial features for forecasting the performance index is probably erformance index typically rises in tandem with an increase in Previous Scores. (0.375332): Ince Index (0.38), and this attribute have a moderately positive link. This implies that although study hours have a beneficial impact on performance, it is not as significant as Previous Scores. Although the association is weaker, increasing the number improve the performance index.
The Performan might not be a Sleep Hours The association not.	llar Activities (0.026075): noce Index and this attribute have a very weak positive association (0.03). This suggests that the performance index is not significantly impacted by extracurricular activities. Unless there is an indirect effect or non-linear link that we haven't yet identifie a particularly helpful characteristic for the model. is (0.050352): on between sleep hours and the performance measure is likewise quite poor (0.05), indicating that sleep hours hardly affect anything. Since this feature doesn't seem to add much to the prediction, it's worth thinking about whether you want to preserve ample Question Papers (0.043436):
The correction_	shows a very weak positive connection (0.04) with the Performance Index, just like Sleep Hours and Extracurricular Activities. This implies that the quantity of rehearsed question papers is not a reliable indicator and does not have a substantial correlation. Telation Matrix Table the correlation matrix for all feature columns matrix = df[['Hours Studied', 'Previous Scores', 'Extracurricular Activities',
Print (correl Hours Studied Previous Scor Extracurricul Sleep Hours Sample Questi Hours Studied Previous Scor	Hours Studied Previous Scores \ 1
Hours Studied Previous Scor Extracurricul Sleep Hours Sample Questi Visualizii import seabo	ng the Correlation using heat map
<pre>plt.figure(f # Plot the h sns.heatmap(# Add title</pre>	the matplotlib figure figsize=(10, 8)) the atmap for the correlation matrix (correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f', linewidths=0.5) and show the plot Correlation Matrix of X Variables') Correlation Matrix of X Variables 1.0
	Previous Scores0.01 1.00 0.01 0.01 0.01 -0.01 -0.6
Extra	Sleep Hours - 0.00 0.01 -0.02 1.00 0.00 -0.02 -0.02 -0.02 -0.04 -0.02
Sample Questi	Honra Strong Papers Practiced - 0.02 0.01 0.01 0.00 1.00 1.00 1.00 Papers Practiced - 0.02 Scores Strong Previous Scores Strong Previous Stron
	orn as sns lotlib.pyplot as plt
<pre>y = df['Perf # Define the independent_ # Create sca plt.figure(f # Loop throu for i, var i plt.subp sns.scat</pre>	<pre>target variable (y) formance Index'] e independent variables (X) _vars = ['Hours Studied', 'Previous Scores', 'Extracurricular Activities',</pre>
plt.xlab plt.ylab	<pre>pel('Performance Index') yout to prevent overlap of subplots ayout()</pre>
Serformance Index	2 3 4 5 6 7 8 9 40 50 60 70 80 90 100 Hours Studied Previous Scores
Scat	ster Plot: Extracurricular Activities vs Performance Index Scatter Plot: Sleep Hours vs Performance Index 100
100 - x 80 -	0.2 0.4 0.6 0.8 1.0 4 5 6 7 8 9 Extracurricular Activities Plot: Sample Question Papers Practiced vs Performance Index
import scipy	2 4 6 8 Sample Question Papers Practiced y.stats as stats
<pre>X = df[['Hou</pre>	efficient, p_value = stats.pearsonr(X[var], y) 'Correlation between {var} and Performance Index: {corr_coefficient:.4f}, p-value: {p_value:.4f}')
Correlation b Correlation b Correlation b Correlation b Interpretat Correlation k Correlation Correlation Correlation	between Hours Studied and Performance Index: 0.3753, p-value: 0.0000 between Previous Scores and Performance Index: 0.9151, p-value: 0.0000 between Extracurricular Activities and Performance Index: 0.0261, p-value: 0.0096 between Sleep Hours and Performance Index: 0.0504, p-value: 0.0000 between Sample Question Papers Practiced and Performance Index: 0.0434, p-value: 0.0000 tion of Correlation Test Results: between Hours Studied and Performance Index: between Hours Studied a
Correlation Correlation & Correlation Correlation Covalue: 0.0096 I	between Previous Scores and Performance Index: perficient: 0.9151 Interpretation: This indicates a strong positive relationship between previous scores and the performance index. The performance index appears to be highly influenced by previous scores, suggesting that students with higher prior so the performance index is highly reliable and unlikely to be due to chance. between Extracurricular Activities and Performance Index: perficient: 0.0261 Interpretation: This indicates a very weak positive relationship between extracurricular activities and performance index. The correlation is so weak that it suggests that extracurricular activities have little to no impact on performance in Interpretation: Despite the very weak correlation, the p-value is less than 0.05, suggesting that the relationship is statistically significant. However, the weak correlation means this variable has minimal influence on performance. between Sleep Hours and Performance Index: perficient: 0.0504 Interpretation: This indicates a very weak positive relationship between sleep hours and performance index. While there is a slight positive trend, it is not a strong or meaningful relationship. p-value: 0.0000 Interpretation: The p-value
than 0.05, indication to Correlation Corre	between Sample Question Papers Practiced and Performance Index: befficient: 0.0434 Interpretation: This indicates a very weak positive relationship between the number of sample question papers practiced and the performance index. p-value: 0.0000 Interpretation: The p-value is less than 0.05, so the relationship is indicated. However, the very weak correlation suggests that practicing sample papers has a negligible effect on performance. a predictive model.
statistically sign	dependent variables (X) and the target variable (y) ('Performance Index', axis=1) # All columns except 'Performance Index' formance Index'] data into training and testing sets (80% train, 20% test) n.model_selection import train_test_split test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42) the linear regression model n.linear_model import LinearRegression earRegression() model on the training data _train, y_train)
# Define ind X = df.drop(y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the	n the test set del.predict(X_test) performance metrics n.metrics import mean_squared_error, r2_score
# Define ind X = df.drop(Y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on Y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th	re(y_test, y_pred) he results
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print(f'Mean print(f'Mean print(f'Root print(f'Root) print(f'Root print(f'Root print(f'Root) print(f'Root print(f'Root)	rety_test, y_pred) the results in Squared Error (MSE): {mse:.4f}') t Mean Squared Error (MSE): {mse:.4f}') t Mean Squared Error (RMSE): {rmse:.4f}') Error (MSE): 4.3059 Jared Error (RMSE): 2.0751 **): 0.9884 If the Model Evaluation: The mean squared error (MSE). The average squared difference between the expected and actual values is known as the MSE. A better model fit is indicated by a reduced MS Description of the model is indicated by a reduced management of the model is indicated by a reduced management of the model is indicated by a reduc
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print(f'Mean print(f'Root print(f'R-sq Mean Squared Root Mean Squ R-squared (R² Results of 4.3059 is the This is a comp The root RM In the same un average predict R2 (R-square The model exp variable.	part (mse) re(y_test, y_pred) he results n Squared Error (MSE): {mse:.4f}') t Mean Squared Error (MSE): {rmse:.4f}') quared (RF): {r2:.4f}') Error (MSE): 4.3059 pared Error (RMSE): 2.0751)): 0.3984 f the Model Evaluation: we mean squared error (MSE). paratively tiny value, suggesting that, on average, the model's predictions closely match the actual values. The average squared difference between the expected and actual values is known as the MSE. A better model fit is indicated by a reduced MSE (mean squared error): 2.0751 nits as the goal variable (Performance Index), the RMSE provides us with an estimate of the average deviation between the projected and actual values. Given that the RMSE and Performance Index are on the same scale, this figure indicates that the cition error of the model is roughly 2.0751 units.
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print(f'Rean print	The Model Evaluation: In mean squared error (MSE): (mean, 42)1) The Model Evaluation: In mean squared error (MSE): The Model Evaluation (MSE): The Model Evalu
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rnse = np.sq r2 = r2_scor # Display th print(f'Mean print(f'Rean print(f'Re	In common of the
Building # Define ind X = df.drop(y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rx = np.sq r2 = r2_scor # Display th print(f'Mean print(f'Rean	The Model Evaluation The de page of the Confinence Squared and the State Confirence Squared and the State Confinence Squared and the State Confirence Squared and the State Confinence Squared and the State Confirence Squared and th
Building # Define ind X = df.drop(y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sg r2 = r2_scor # Display th print(f'Nean print(f'Root print(f'Ro	As a control of the c
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print (f'Mean print (f'Reot print (coeffic }) # Calculate from sklearn from sklear	The Model Evolution of the common to the com
Building # Define ind X = df.drop(Y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on Y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print(f'Read print(f'Read print(f'Read print(f'Read Results of 4.3059 is the This is a comp The root RM In the same un average predict R2 (R-square Results of 4.3059 is the This is a comp The root RM In the same un average predict R2 (R-square The model exp variable. Table of # Get the fe features = X # Create a D coefficients "Coeffic }) # Add the in intercept_df coefficients # Display th print(coeffic }) # Add the in intercept_df coefficients # Create a D coefficients	The Model Evaluation:
Building # Define ind X = df.drop() y = df['Perf # Split the from sklearn X_train, X_t # Initialize from sklearn model = Line # Train the model.fit(X_ # Predict on y_pred = mod # Calculate from sklearn import numpy mse = mean_s rmse = np.sq r2 = r2_scor # Display th print (f'Mean print (f'Reot print (fore) # Get the fe feature = x # Create a D coefficients 'Feature 'Coeffic }) # Add the in interpretat Hours Studie For every addif Sleep Hours Extracurricul For every addif Sleep Hours Every extra ho Intercept: When all other We can ok With a value of that, although in an effect on the # Plot Actua pt.figure (f spl. scatter pt. pl. scatter pt. pl. scatter pt. scatter pt	A STATE OF THE CONTROL OF THE CONTRO

Course Title: GCIS-523-0B: Statistical Computing