# Student Performance Analyzer Report

*(Abdullah, AI-ASSIGNMENT#1, SU92-BSAIM-F24-002)*

## Introduction

This report provides a detailed explanation of the Student Performance Analyzer program. The project follows a complete data science workflow—from data loading and preprocessing to analysis, model training, evaluation, and saving the trained model for future use. Each step has been designed to ensure clarity, efficiency, and scalability in understanding student performance data and predicting outcomes using machine learning techniques.

## Step 1: Data Loading

The process begins with importing the dataset using the pandas library. The 'load_data()' method reads the CSV file ('StudentsPerformance.csv') and stores it in a pandas DataFrame. To verify that the data has been successfully imported, the first few rows are displayed using df.head(). This step establishes the foundation for all subsequent analysis.

## Step 2: Data Understanding

In this stage, the 'understand_data()' method provides a comprehensive overview of the dataset. It displays the data types, memory usage, and the count of non-null values using df.info(). Statistical summaries of numerical columns are obtained through df.describe(), while df.isnull().sum() identifies any missing values. Together, these commands help in understanding the dataset's structure, integrity, and quality.

## Step 3: Data Preprocessing

Preprocessing ensures the dataset is clean and ready for analysis. The 'preprocess_data()' method uses one-hot encoding (via pd.get_dummies()) to convert categorical variables into numerical format, as most machine learning algorithms require numeric input. Missing values, if any, are removed using df.dropna(). This step guarantees consistency and prevents model training issues caused by incomplete data.

## Step 4: Univariate Analysis

Univariate analysis focuses on understanding individual features of the dataset. The 'univariate_analysis()' function generates histograms for all numerical columns, allowing visualization of data distribution. This helps in detecting skewness, understanding variability, and identifying possible outliers in each feature.

## Step 5: Bivariate Analysis

The 'bivariate_analysis()' function examines relationships between variables using a correlation heatmap generated with seaborn. This visualization highlights the strength and direction of relationships between features and the target variable ('math score'). Identifying correlated variables assists in feature selection and better model interpretation.

## Step 6: Exploratory Data Analysis (EDA)

The 'exploratory_data_analysis()' method provides deeper insights into the dataset. It first visualizes the distribution of the target variable ('math score') using a histogram with a

density curve. Boxplots are then created to detect outliers in numerical features, followed by countplots for categorical features to observe their frequency distribution. This comprehensive approach helps understand how features interact with one another.

### Step 7: Data Splitting

Once the dataset is preprocessed, it is divided into training and testing subsets using the train_test_split() function. Eighty percent of the data is allocated for training the model, while the remaining twenty percent is reserved for testing. The target variable ('math score') is separated from the feature set, ensuring proper supervised learning structure.

### Step 8: Model Training (Support Vector Machine)

In this stage, a Support Vector Regressor (SVR) model from the scikit-learn library is utilized. The 'train_model()' function fits the training data into the SVR model to establish patterns and relationships. Support Vector Machines are effective in handling high-dimensional data and generating accurate predictions.

### Step 9: Model Evaluation

The 'evaluate_model()' method tests the trained SVM model on the test dataset. It calculates the Mean Squared Error (MSE), a key performance metric for regression tasks. A lower MSE value indicates that the model has achieved higher accuracy in predicting the 'math score'. This evaluation confirms the model's predictive efficiency.

### Step 10: Model Saving

The final step involves preserving the trained model using the pickle library. The 'save_model()' function serializes the model object and stores it as 'svm.model.pk'. This allows future reuse of the model without retraining, ensuring time efficiency and reproducibility.

### Conclusion

The Student Performance Analyzer demonstrates a complete machine learning workflow: starting from raw data import, preprocessing, and visualization to model training and evaluation. The use of Support Vector Regression provides a reliable method for predicting academic performance. This project effectively illustrates how data science principles can be applied to real-world educational analytics.