# Student Performance Prediction using Linear Regression

**Created by: Dulhara Lakshan (2025.03.14)**

## 1. Introduction

This project aims to predict student performance based on various factors such as study habits, participation, and stress levels using **Linear Regression**. The model is trained on a dataset containing student records and predicts final grades based on multiple independent variables.

## 2. Dataset Overview

The dataset contains student-related attributes, including:

- **Demographic Information**: Age, Gender

- **Study Habits**: Study Hours per Week, Online Courses Completed

- **Engagement**: Participation in Discussions, Use of Educational Technology

- **Performance Metrics**: Assignment Completion Rate, Exam Score, Attendance Rate

- **Lifestyle Factors**: Time Spent on Social Media, Sleep Hours per Night

- **Psychological Factors**: Preferred Learning Style, Self-Reported Stress Level

- **Target Variable**: Final_Grade (Encoded as a categorical variable)

## 3. Data Preprocessing

### 3.1 Handling Missing Values

The dataset was checked for missing values using df.isnull().sum(). No missing values were found.

### 3.2 Encoding Categorical Features

Categorical variables were encoded as follows:

- **Binary Encoding**:

    - Gender: Female → 0, Male → 1

    - Participation_in_Discussions: No → 0, Yes → 1

    - Use_of_Educational_Tech: No → 0, Yes → 1

- **One-Hot Encoding**:

    - Preferred_Learning_Style: Converted into three columns (Kinesthetic, Reading/Writing, Visual)

o   Self_Reported_Stress_Level: Converted into three columns (Low, Medium, High)

- **Encoding Target Variable**:

  o   Final_Grade: Categorical values were converted into numerical values using LabelEncoder.

### 3.3 Feature Selection

The independent variables (**features**) selected for training the model are:

```
FEATURES = [

  "Age", "Gender", "Study_Hours_per_Week", "Online_Courses_Completed",

  "Participation_in_Discussions", "Assignment_Completion_Rate (%)", "Exam_Score (%)",

  "Attendance_Rate (%)", "Use_of_Educational_Tech", "Time_Spent_on_Social_Media
(hours/week)",

  "Sleep_Hours_per_Night", "Preferred_Learning_Style_Kinesthetic",

  "Preferred_Learning_Style_Reading/Writing", "Preferred_Learning_Style_Visual",

  "Self_Reported_Stress_Level_Low", "Self_Reported_Stress_Level_Medium"

]
```

### 4. Model Training

The dataset was split into training and testing sets (80%-20%) using train_test_split(). The model used for training was **Linear Regression**, implemented as follows:

```
from sklearn.linear_model import LinearRegression

model = LinearRegression()

model.fit(X_train, y_train)

The trained model produced the following coefficients:

print(f"Intercept (c): {model.intercept_}")

print(f"Coefficients (m): {model.coef_}")
```

## 5. Feature Importance

THE IMPORTANCE OF EACH FEATURE WAS VISUALIZED USING A BAR PLOT OF THE REGRESSION COEFFICIENTS.

```
PLT.FIGURE(FIGSIZE=(10, 6))

SNS.BARPLOT(X="COEFFICIENT", Y="FEATURE", DATA=COEF_DF, PALETTE="COOLWARM")

PLT.AXVLINE(0, COLOR="BLACK", LINEWIDTH=1.2)

PLT.TITLE("FEATURE IMPORTANCE (LINEAR REGRESSION COEFFICIENTS)")

PLT.XLABEL("COEFFICIENT VALUE")

PLT.YLABEL("FEATURES")

PLT.SHOW()
```

## 6. Model Evaluation

The model was evaluated using **Mean Squared Error (MSE)** and **R² Score**.

```
FROM SKLEARN.METRICS IMPORT MEAN_SQUARED_ERROR, R2_SCORE

Y_PRED = MODEL.PREDICT(X_TEST)

MSE = MEAN_SQUARED_ERROR(Y_TEST, Y_PRED)

R2 = R2_SCORE(Y_TEST, Y_PRED)

PRINT(F"MEAN SQUARED ERROR: {MSE:.2F}")

PRINT(F"R² SCORE: {R2:.4F}")
```

A scatter plot was used to compare actual vs. predicted values:

```
PLT.SCATTER(Y_TEST, Y_PRED, COLOR='BLUE')

PLT.XLABEL("ACTUAL SCORES")

PLT.YLABEL("PREDICTED SCORES")

PLT.TITLE("ACTUAL VS PREDICTED STUDENT SCORES")

PLT.SHOW()
```

## 7. Individual Feature Analysis

Scatter plots were generated to visualize how each independent variable impacts the final grade.

for feature in features:

```
PLT.FIGURE(FIGSIZE=(6, 4))

SNS.SCATTERPLOT(X=DF[FEATURE], Y=DF["FINAL_GRADE"], ALPHA=0.5, COLOR="BLUE")

PLT.XLABEL(FEATURE)

PLT.YLABEL("FINAL_GRADE")

PLT.TITLE(F"FINAL GRADE VS {FEATURE}")

PLT.SHOW()
```

## 8. Making Predictions

A new student's data was used to predict their final grade.

```
NEW_STUDENT = NP.ARRAY([[18, 1, 40, 10, 1, 85, 75, 90, 1, 15, 7, 0, 1, 0, 0, 1]])

PREDICTED_GRADE = MODEL.PREDICT(NEW_STUDENT)

PRINT(F"PREDICTED FINAL GRADE: {PREDICTED_GRADE[0]}")
```
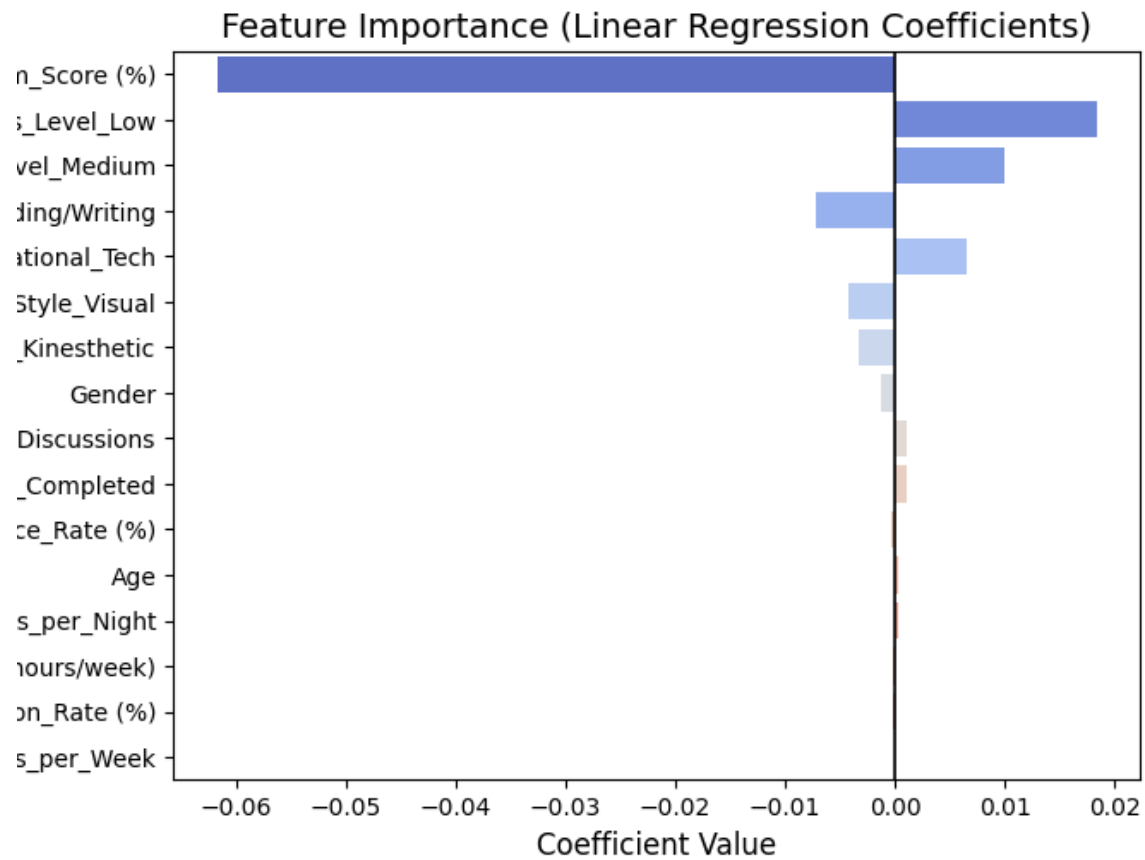
## 9. Conclusion

This project successfully applied **Linear Regression** to predict student performance based on multiple academic and lifestyle factors. The model demonstrated meaningful relationships between study habits, participation, and final grades. Future improvements could include:

- Testing alternative models (e.g., Decision Trees, Random Forest)

- Expanding feature selection with additional psychological and social factors

- Implementing feature scaling and polynomial regression for better performance.

This project provides a valuable tool for educators and students to identify key factors influencing academic success.

**Feature Coefficient**



Feature Importance (Linear Regression Coefficients)

Based on these Linear Regression coefficients, we can interpret the impact of each feature on student grades:

| Feature | Coefficient |
| --- | --- |
| Age | 3.20149716e-04 |
| Gender | -1.33543458e-03 |
| Study_Hours_per_Week | -7.44875421e-05 |
| Online_Courses_Completed | 9.93713199e-04 |
| Participation_in_Discussions | 1.02266361e-03 |
| Assignment_Completion_Rate (%) | -2.33892195e-04 |
| Exam_Score (%) | -6.17185666e-02 |
| Attendance_Rate (%) | -4.05044953e-04 |
| Use_of_Educational_Tech | 6.47230358e-03 |
| Time_Spent_on_Social_Media (hours/week) | -2.50046174e-04 |
| Sleep_Hours_per_Night | 2.57071811e-04 |
| Preferred_Learning_Style_Kinesthetic | -3.37131296e-03 |
| Preferred_Learning_Style_Reading/Writing | -7.27618642e-03 |
| Preferred_Learning_Style_Visual | -4.23511487e-03 |
| Self_Reported_Stress_Level_Low | 1.83073644e-02 |
| Self_Reported_Stress_Level_Medium | 9.88368918e-03 |

Most Influential Factors:

- Exam Score (%) (-0.0617) – The strongest predictor of final grades (negative coefficient means lower scores significantly decrease final grades).
- Self-Reported Stress Level (Medium) (+0.0099) – Moderate stress slightly improves performance, possibly due to motivation.
- Self-Reported Stress Level (Low) (+0.0183) – Lower stress levels positively influence grades more than medium stress.
- Use of Educational Tech (+0.0064) – Using tech tools has a small but positive effect on student success.

Negligible or Negative Impact:

- Study Hours per Week (-0.00007) – Almost no direct impact, suggesting quality > quantity in study sessions.

- Attendance Rate (-0.0004) – Surprisingly low effect, indicating attendance alone isn't a strong predictor.
- Learning Styles (Kinesthetic: -0.0033, Reading/Writing: -0.0072, Visual: -0.0042) – No learning style has a major influence, meaning adaptability may be key.
- Time Spent on social media (-0.00025) – Slight negative impact, but not as harmful as expected.

Key Takeaways:

- Grades are most impacted by exam performance, stress levels, and educational tech usage.

- Study hours alone don't guarantee higher grades—effective learning strategies matter more.

- Stress management plays a crucial role, as moderate stress may be beneficial.

- Social media usage has minimal impact, contradicting common assumptions.

This analysis can help educators and students focus on what truly improves academic performance.