# Student Performance Analysis Project

## Objective:

To analyze student data and identify key factors affecting academic performance. This analysis helps provide insights for educators and students to improve outcomes.

**Dataset:** Student Performance Dataset (Math Subject)
**Tools Used:** Python, Pandas, Matplotlib, Seaborn

## Load and View the Dataset

We load the student-mat.csv dataset using Pandas and view the first few rows.

```
In [19]:  import pandas as pd

          df = pd.read_csv("C:\\Users\\PRAKASH ROUT\\Downloads\\student\\student-mat.csv",
          df.head()
```

Out[19]:

| | school | sex | age | address | famsize | Pstatus | Medu | Fedu | Mjob | Fjob | ... | fan |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | GP | F | 18 | U | GT3 | A | 4 | 4 | at_home | teacher | ... | |
| **1** | GP | F | 17 | U | GT3 | T | 1 | 1 | at_home | other | ... | |
| **2** | GP | F | 15 | U | LE3 | T | 1 | 1 | at_home | other | ... | |
| **3** | GP | F | 15 | U | GT3 | T | 4 | 2 | health | services | ... | |
| **4** | GP | F | 16 | U | GT3 | T | 3 | 3 | other | other | ... | |

5 rows × 33 columns

## Explore the Dataset

Let's check the dataset's shape, columns, data types, and null values.

```
In [20]:  #Check Basic Info
          df.info()
          df.shape
          df.isnull().sum()
          df.describe()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 395 entries, 0 to 394
Data columns (total 33 columns):
 #   Column      Non-Null Count  Dtype
---  ------      --------------  -----
 0   school      395 non-null    object
 1   sex         395 non-null    object
 2   age         395 non-null    int64
 3   address     395 non-null    object
 4   famsize     395 non-null    object
 5   Pstatus     395 non-null    object
 6   Medu        395 non-null    int64
 7   Fedu        395 non-null    int64
 8   Mjob        395 non-null    object
 9   Fjob        395 non-null    object
 10  reason      395 non-null    object
 11  guardian    395 non-null    object
 12  traveltime  395 non-null    int64
 13  studytime   395 non-null    int64
 14  failures    395 non-null    int64
 15  schoolsup   395 non-null    object
 16  famsup      395 non-null    object
 17  paid        395 non-null    object
 18  activities  395 non-null    object
 19  nursery     395 non-null    object
 20  higher      395 non-null    object
 21  internet    395 non-null    object
 22  romantic    395 non-null    object
 23  famrel      395 non-null    int64
 24  freetime    395 non-null    int64
 25  goout       395 non-null    int64
 26  Dalc        395 non-null    int64
 27  Walc        395 non-null    int64
 28  health      395 non-null    int64
 29  absences    395 non-null    int64
 30  G1          395 non-null    int64
 31  G2          395 non-null    int64
 32  G3          395 non-null    int64
dtypes: int64(16), object(17)
memory usage: 102.0+ KB
```

Out[20]:

| | age | Medu | Fedu | traveltime | studytime | failures | fam |
|---|---|---|---|---|---|---|---|
| count | 395.000000 | 395.000000 | 395.000000 | 395.000000 | 395.000000 | 395.000000 | 395.0000 |
| mean | 16.696203 | 2.749367 | 2.521519 | 1.448101 | 2.035443 | 0.334177 | 3.9443 |
| std | 1.276043 | 1.094735 | 1.088201 | 0.697505 | 0.839240 | 0.743651 | 0.8966 |
| min | 15.000000 | 0.000000 | 0.000000 | 1.000000 | 1.000000 | 0.000000 | 1.0000 |
| 25% | 16.000000 | 2.000000 | 2.000000 | 1.000000 | 1.000000 | 0.000000 | 4.0000 |
| 50% | 17.000000 | 3.000000 | 2.000000 | 1.000000 | 2.000000 | 0.000000 | 4.0000 |
| 75% | 18.000000 | 4.000000 | 3.000000 | 2.000000 | 2.000000 | 0.000000 | 5.0000 |
| max | 22.000000 | 4.000000 | 4.000000 | 4.000000 | 4.000000 | 3.000000 | 5.0000 |

# Understand the Features

Explore categorical and numerical columns to understand their distributions.

```
In [21]:  #Check data types and unique values:

          df['sex'].value_counts()
          df['studytime'].value_counts()
          df['failures'].value_counts()
```

```
Out[21]:  failures
          0    312
          1     50
          2     17
          3     16
          Name: count, dtype: int64
```

```
In [22]:  # Understand target variable

          df['average_score'] = df[['G1', 'G2', 'G3']].mean(axis=1)
```

# Data Cleaning

```
In [24]:  # Handle categorical data if needed
          # Convert binary columns (yes/no) to 1/0
          df['schoolsup'] = df['schoolsup'].map({'yes': 1, 'no': 0})
```

```
In [25]:  # Check and remove duplicates
          df.drop_duplicates(inplace=True)
```

# Exploratory Data Analysis (EDA)

In this step, we explore the dataset visually to uncover patterns, trends, and relationships between different features and student performance.

We use **Matplotlib** and **Seaborn** libraries to create informative visualizations.

```
In [26]:  import matplotlib.pyplot as plt
          import seaborn as sns

          sns.set(style="whitegrid")
```

## 1. Distribution of Average Scores

We analyze how the average scores of students are distributed across the dataset.

```
In [28]:  sns.histplot(df['average_score'], kde=True)
          plt.title('Distribution of Average Student Scores')
          plt.show()
```

Distribution of Average Student Scores

## 2. Study Time vs Average Score

We visualize how the amount of time students dedicate to studying impacts their average scores.

```
In [29]: sns.boxplot(x='studytime', y='average_score', data=df)
         plt.title('Study Time vs Average Score')
         plt.show()
```

Study Time vs Average Score

### 3. Gender vs Performance

We compare average scores between male and female students to see if there's a performance gap

```
In [12]:  # gender comparison
          sns.boxplot(x='sex', y='average_score', data=df)
          plt.title('Gender vs Average Score')
          plt.show()
```

Gender vs Average Score

## 4. Failures vs Average Score

We explore how the number of past class failures affects student performance.

In [13]:
```python
# Failure VS Performance
sns.barplot(x='failures', y='average_score', data=df)
plt.title('Failures vs Average Score')
plt.show()
```

Failures vs Average Score

## Correlation Analysis

In this step, we analyze how numeric features are correlated with each other, especially with the final grade ( G3 ) and the `average_score` .

A correlation matrix helps us identify:

- Strong positive or negative relationships
- Multicollinearity
- Key influencing factors for student performance

```
In [30]: # co-relations heatmap
         plt.figure(figsize=(12,8))
         sns.heatmap(df.corr(numeric_only=True), annot=True, cmap='coolwarm')
         plt.title('Correlation Heatmap')
         plt.show()
```

## Correlation Heatmap

| | age | Medu | Fedu | traveltime | studytime | failures | schoolsup | famrel | freetime | goout | Dalc | Walc | health | absences | G1 | G2 | G3 | average_score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| age | 1 | -0.16 | -0.16 | 0.071 | -0.0041 | 0.24 | | 0.054 | 0.016 | 0.13 | 0.13 | 0.12 | -0.062 | 0.18 | -0.064 | -0.14 | -0.16 | -0.13 |
| Medu | -0.16 | 1 | 0.62 | -0.17 | 0.065 | -0.24 | | -0.0039 | 0.031 | 0.064 | 0.02 | -0.047 | 0.047 | 0.1 | 0.21 | 0.22 | 0.22 | 0.22 |
| Fedu | -0.16 | 0.62 | 1 | -0.16 | 0.0092 | -0.25 | | -0.0014 | 0.013 | 0.043 | 0.0024 | 0.013 | 0.015 | 0.024 | 0.19 | 0.16 | 0.15 | 0.18 |
| traveltime | 0.071 | -0.17 | -0.16 | 1 | -0.1 | 0.092 | | -0.017 | 0.017 | 0.029 | 0.14 | 0.13 | 0.0075 | 0.013 | -0.093 | -0.15 | -0.12 | -0.13 |
| studytime | -0.0041 | 0.065 | 0.0092 | -0.1 | 1 | -0.17 | | 0.04 | -0.14 | -0.064 | -0.2 | -0.25 | -0.076 | 0.063 | 0.16 | 0.14 | 0.098 | 0.13 |
| failures | 0.24 | -0.24 | -0.25 | 0.092 | -0.17 | 1 | | -0.044 | 0.092 | 0.12 | 0.14 | 0.14 | 0.066 | 0.064 | -0.35 | -0.36 | -0.36 | -0.38 |
| schoolsup | | | | | | | | | | | | | | | | | | |
| famrel | 0.054 | -0.0039 | -0.0014 | -0.017 | 0.04 | -0.044 | | 1 | 0.15 | 0.065 | -0.078 | -0.11 | 0.094 | -0.044 | 0.022 | -0.018 | 0.051 | 0.022 |
| freetime | 0.016 | 0.031 | -0.013 | 0.017 | -0.14 | 0.092 | | 0.15 | 1 | 0.29 | 0.21 | 0.15 | 0.076 | -0.058 | 0.013 | -0.014 | 0.011 | 0.0038 |
| goout | 0.13 | 0.064 | 0.043 | 0.029 | -0.064 | 0.12 | | 0.065 | 0.29 | 1 | 0.27 | 0.42 | -0.0096 | 0.044 | -0.15 | -0.16 | -0.13 | -0.15 |
| Dalc | 0.13 | 0.02 | 0.0024 | 0.14 | -0.2 | 0.14 | | -0.078 | 0.21 | 0.27 | 1 | 0.65 | 0.077 | 0.11 | -0.094 | -0.064 | -0.055 | -0.073 |
| Walc | 0.12 | -0.047 | 0.013 | 0.13 | -0.25 | 0.14 | | -0.11 | 0.15 | 0.42 | 0.65 | 1 | 0.092 | 0.14 | -0.13 | -0.085 | -0.052 | -0.088 |
| health | -0.062 | -0.047 | 0.015 | 0.0075 | -0.076 | 0.066 | | 0.094 | 0.076 | -0.0096 | 0.077 | 0.092 | 1 | -0.03 | -0.073 | 0.098 | 0.061 | -0.08 |
| absences | 0.18 | 0.1 | 0.024 | -0.013 | 0.063 | 0.064 | | -0.044 | -0.058 | 0.044 | 0.11 | 0.14 | -0.03 | 1 | -0.031 | 0.032 | 0.034 | 0.0059 |
| G1 | -0.064 | 0.21 | 0.19 | -0.093 | 0.16 | -0.35 | | 0.022 | 0.013 | -0.15 | -0.094 | -0.13 | -0.073 | -0.031 | 1 | 0.85 | 0.8 | 0.92 |
| G2 | -0.14 | 0.22 | 0.16 | -0.15 | 0.14 | -0.36 | | -0.018 | 0.014 | -0.16 | -0.064 | -0.085 | 0.098 | 0.032 | 0.85 | 1 | 0.9 | 0.97 |
| G3 | -0.16 | 0.22 | 0.15 | -0.12 | 0.098 | -0.36 | | 0.051 | 0.011 | -0.13 | -0.055 | -0.052 | 0.061 | 0.034 | 0.8 | 0.9 | 1 | 0.96 |
| average_score | -0.13 | 0.22 | 0.18 | -0.13 | 0.13 | -0.38 | | 0.022 | 0.0038 | -0.15 | -0.073 | -0.088 | -0.08 | 0.0059 | 0.92 | 0.97 | 0.96 | 1 |

# Key Insights & Recommendations

## 🔍 Key Insights:

1. **Previous Grades (G1, G2):**

   - Strongly correlated with the final grade (G3).
   - Students with higher scores in G1 and G2 tend to perform well in G3.

2. **Study Time:**

   - More study time is generally associated with better performance.
   - Students who study more than 2 hours show higher average scores.

3. **Failures:**

   - Number of past class failures negatively affects the final grade.
   - Students with 0 past failures perform significantly better.

4. **Parental Education:**

   - A slight positive impact on student performance, especially from mother's education level.

5. **Gender:**

   - No major difference in performance between male and female students.

## ✅ Recommendations:

- **Early Intervention:** Track student grades from G1 and G2 to identify those who may need extra help before final exams.
- **Encourage Study Time:** Promote study habits of more than 2 hours per week to improve overall performance.
- **Support Struggling Students:** Provide extra tutoring for students with a history of failures.
- **Continue School Support Programs:** These help improve performance and should be maintained or expanded.
- **Parental Engagement:** Educating parents about their influence can positively impact student outcomes.

---

## ▓ Conclusion:

This analysis provides a clear picture of the key factors that influence student academic performance. With the right support systems and timely interventions, educators and parents can work together to help students succeed.

In [ ]: