# Predictive Analysis of Academic Performance via Social Media Influence Metrics

Presented by: Mouhsine DAOUDA

# Abstract

In the era of smartphones and advanced mobile technologies, the profound impact of social media on society has prompted extensive discourse. This project aims to predict and understand the academic performance of university students by delving into the nuanced facets of social media influence factors.

These factors include the time spent on social media, the number of friends, the usage of different platforms, involvement in various groups, and other relevant metrics. The primary objective is to discern how these aspects of social media engagement may correlate with and predict students' academic achievements. This project directly aligns with SDG (Quality Education) by addressing factors that impact academic performance of students. By unraveling the relationship between social media influence and academic outcomes, the project contributes to the development of strategies for fostering a high-quality education environment.

# 1- Literature Review

## I- Introduction

Our research, Predictive Analysis of Academic Performance via Social Media Influence Metrics, is driven by the critical need to understand the impact of social media on students' academic performance.

This literature review incorporates the findings from the paper **"Measuring the Effect of Social Media on Student Academic Performance Using a Social Media Influence Factor Model"** and the study titled **"Impact of Social Media on Students' Academic Performance**: A Case Study of Islamic University, Bangladesh."**

The significance of this research lies in addressing the evolving landscape of social media usage among students and its potential effects, both positive and negative, on academic outcomes.

A comprehensive literature review is essential to contextualize our study within the existing body of knowledge.

# II- Organization

The literature review is structured thematically, grouping papers based on similar themes. This allows for a nuanced exploration of the multifaceted relationship between social media and academic performance.

The chronological arrangement within each theme enables a chronological understanding of the evolution of research in this domain.

# III- SUMMARY AND SYNTHESIS

*Paper 1:* Measuring the Effect of Social Media on Student Academic Performance Using a Social Media Influence Factor Model

As summarized earlier, this paper introduces a Social Media Influence Factor Model, providing a quantitative approach to understanding the impact of social media on academic performance. It contributes by offering a practical tool for analysis.

LİNK: https://doi.org/10.6084/m9.figshare.14905

*Paper 2:* Impact of Social Media on Students' Academic Performance: A Case Study of Islamic University, Bangladesh

This study investigates the effects of social media on academic performance among students in Bangladesh. Utilizing a well-structured questionnaire and applying descriptive and inferential statistics, the study identifies that a majority of students engage in non-academic social media activities, with Facebook being the preferred platform. The study underscores the negative impact of excessive social media usage on academic performance, recommending guidance and monitoring by parents, teachers, and university advisors.

LİNK: https://doi.org/10.46281/aesr.v10i1.1822

# IV- Comparison and Contrast

While both studies recognize the negative impact of social media on academic performance, the second study provides a detailed case study perspective, highlighting the specific habits and preferences of students at Islamic University, Bangladesh. The recommendations from this study align with the practical implications sought by our research.

# V- Conclusion

In conclusion, the combined insights from these studies emphasize the pervasive influence of social media on academic performance. The nuanced understanding of usage patterns and the quantification of this impact contribute significantly to our research objectives.

By synthesizing these findings, our project aims to provide a predictive model that incorporates insights from existing literature and addresses the specific challenges posed by social media usage among university students.

All sources, including the paper "Impact of Social Media on Students' Academic Performance: A Case Study of Islamic University, Bangladesh," and previously mentioned papers, have been diligently cited in adherence to academic standards.
Proper citations uphold the integrity of our literature review and acknowledge the valuable contributions of referenced works to the overarching research goals.

# 2- Data Research

## B/I- Introduction

In the pursuit of predicting academic performance through an analysis of social media influence, our research aims to unravel the intricate relationship between online behavior and students' study habits. The significance of our research lies in its potential to provide educators, institutions, and policymakers with actionable insights to enhance learning environments. Understanding how social media impacts academic outcomes is essential for devising effective strategies that align with the goal of quality education.

A comprehensive exploration of the data is necessary to draw meaningful conclusions and propose solutions for a more effective use of social media platforms in the context of education.

# B/II-Organization

Our data research findings are organized thematically to offer a coherent narrative. Themes include social media usage patterns, influence metrics, and their correlations with academic performance.

This structure provides a holistic view of the multifaceted relationship between social media and academic success.

# DATASET

UPSA the researchers conducted an online survey using a questionnaire to collect data.

The questionnaire comprised 8 items covering: the number of friends on social media, number of social media groups, number of social media platforms, amount of time spent on social media daily, number of notification checked Daily, and population demographics (gender, age, and level of study).

The questionnaires were administered using an online survey tool on "**UPSA virtual**" which is the official learning management system (LMS) of the case study institution. UPSA Virtual had a total of 33,126 registered students at the time of the survey  (25th March 2020) and that was used as the study population. A random sample method has been used and a minimum of 380 participants are required to deliver credible results at a 95% confidence level and with a 5% margin of error. Out of 800 requests, 623 filled the form.

The available data for our report is 623 students, so applying a **sample selection is no more needed**.

LİNK: https://doi.org/10.6084/m9.figshare.14905

# DATA DESCRIPTION

| Age Group | 1: 20 years or below<br>2: [20, 30]<br>3: above 30 years |
|---|---|
| Gender | 1: Male<br>2: Female |
| Levels | 1: 100 level<br>2: 200 level<br>3: 300 level<br>4: 400 level |
| Time | Represents the number of hours students spend on social media. It ranges from 1 to 5, indicating different time intervals |
| Platforms | Indicates the number of social media platforms a student belongs to. It ranges from 1 to 6, representing different levels of engagement with social media |

```
In [6]: import pandas as pd

# Read the csv file
data = pd.read_csv('C:/Users/damsd/Videos/final project/Data Analyst Python/data.csv')

# Print it out if you want
data
```

Out[6]:

| | Age Group | Gender | Levels | Time | Platforms | Friends | Groups | Notifications | GPA |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2 | 1 | 2 | 2 | 5 | 3000 | 6 | 30 | 2.1 |
| 1 | 2 | 2 | 3 | 5 | 3 | 4000 | 4 | 50 | 2.5 |
| 2 | 2 | 1 | 3 | 3 | 5 | 2000 | 5 | 30 | 2.5 |
| 3 | 2 | 1 | 4 | 3 | 2 | 2000 | 3 | 10 | 3.0 |
| 4 | 3 | 2 | 1 | 1 | 3 | 1000 | 3 | 5 | 3.5 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 617 | 2 | 1 | 4 | 1 | 2 | 1000 | 3 | 5 | 3.6 |
| 618 | 1 | 1 | 4 | 1 | 3 | 1000 | 0 | 5 | 2.0 |
| 619 | 2 | 2 | 2 | 4 | 5 | 2000 | 3 | 10 | 2.5 |
| 620 | 2 | 1 | 3 | 1 | 2 | 1000 | 3 | 5 | 2.7 |
| 621 | 2 | 2 | 4 | 1 | 2 | 1000 | 4 | 5 | 3.0 |

622 rows × 9 columns

# PREPROCESSİNG

| Friends | Represents the number of friends a student has on social media |
|---|---|
| Groups | Indicates the number of social media groups a student belongs to |
| Notifications | Represents the number of times a student checks notifications on social media daily |
| GPA | Grade Point Average of the student, which serves as a measure of their academic performance |

```
In [6]: import pandas as pd

        # Read the csv file
        data = pd.read_csv('C:/Users/damsd/Videos/final project/Data Analyst Python/data.csv')

        # Print it out if you want
        data
```

Out[6]:

| | Age Group | Gender | Levels | Time | Platforms | Friends | Groups | Notifications | GPA |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 2 | 1 | 2 | 2 | 5 | 3000 | 6 | 30 | 2.1 |
| 1 | 2 | 2 | 3 | 5 | 3 | 4000 | 4 | 50 | 2.5 |
| 2 | 2 | 1 | 3 | 3 | 5 | 2000 | 5 | 30 | 2.5 |
| 3 | 2 | 1 | 4 | 3 | 2 | 2000 | 3 | 10 | 3.0 |
| 4 | 3 | 2 | 1 | 1 | 3 | 1000 | 3 | 5 | 3.5 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 617 | 2 | 1 | 4 | 1 | 2 | 1000 | 3 | 5 | 3.6 |
| 618 | 1 | 1 | 4 | 1 | 3 | 1000 | 0 | 5 | 2.0 |
| 619 | 2 | 2 | 2 | 4 | 5 | 2000 | 3 | 10 | 2.5 |
| 620 | 2 | 1 | 3 | 1 | 2 | 1000 | 3 | 5 | 2.7 |
| 621 | 2 | 2 | 4 | 1 | 2 | 1000 | 4 | 5 | 3.0 |

622 rows × 9 columns

# DATA INFO

The DataFrame consists of 622 entries with 9 columns, including a combination of integer (int64) and float (float64) data types.

```
In [5]: #data information
        data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 622 entries, 0 to 621
Data columns (total 9 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   Age Group      622 non-null    int64
 1   Gender         622 non-null    int64
 2   Level          622 non-null    int64
 3   Time           622 non-null    int64
 4   Platforms      622 non-null    int64
 5   Friends        622 non-null    int64
 6   Groups         622 non-null    int64
 7   Notifications  622 non-null    int64
 8   GPA            622 non-null    float64
dtypes: float64(1), int64(8)
memory usage: 43.9 KB
```

# DATA DESCRIBE

| Age Group | The mean age group is approximately 1.84, which suggests that the majority of students in the sample are below 30 years old. The minimum age group is 1, indicating the presence of individuals aged 20 years or below, and the maximum age group is 3, representing individuals above 30 years. |
|---|---|
| Gender | The mean gender value is around 1.56, indicating a slight skew towards male participants. |
| Level | The mean level is approximately 2.64, suggesting that the majority of students are in the 300 or 400 levels. The minimum level is 1 (100 level), and the maximum level is 4 (400 level). |

```
In [6]: # Data Summary
        summary = data.describe()

In [7]: print(summary)

              Age Group      Gender        Level        Time    Platforms  \
       count  622.000000  622.000000  622.000000  622.000000  622.000000
       mean     1.842444    1.562701    2.638264    2.207395    3.490354
       std      0.514773    0.524834    1.186077    1.415492    1.232571
       min      1.000000    1.000000    1.000000    1.000000    1.000000
       25%      2.000000    1.000000    1.000000    1.000000    2.000000
       50%      2.000000    2.000000    3.000000    2.000000    3.000000
       75%      2.000000    2.000000    4.000000    3.000000    5.000000
       max      3.000000    3.000000    4.000000    5.000000    5.000000


                  Friends       Groups  Notifications         GPA
       count   622.000000   622.000000     622.000000  622.000000
       mean   1831.189711     3.625402      14.292605    2.841801
       std     865.583604     1.349968      13.208275    0.658272
       min    1000.000000     0.000000       5.000000    1.200000
       25%    1000.000000     3.000000       5.000000    2.300000
       50%    2000.000000     3.000000      10.000000    3.000000
       75%    2000.000000     4.000000      20.000000    3.500000
       max    4000.000000     6.000000      50.000000    3.800000
```

# DATA DESCRIBE()

| Time | The average time spent by students on social media is around 2.21 hours. The minimum and maximum time values are 1 and 5, respectively. |
|------|------|
| Platforms | On average, students belong to around 3.49 social media platforms. The minimum and maximum values are 1 and 5, respectively. |
| Friends | The mean number of friends on social media is approximately 1831.19. The minimum and maximum values are 1000 and 4000, respectively. |
| Groups | On average, students belong to around 3.63 social media groups. The minimum and maximum values are 1 and 6, respectively. |

```
In [6]: # Data Summary
        summary = data.describe()

In [7]: print(summary)

              Age Group      Gender       Level        Time    Platforms  \
count        622.000000  622.000000  622.000000  622.000000  622.000000
mean           1.842444    1.562701    2.638264    2.207395    3.490354
std            0.514773    0.524834    1.186077    1.415492    1.232571
min            1.000000    1.000000    1.000000    1.000000    1.000000
25%            2.000000    1.000000    1.000000    1.000000    2.000000
50%            2.000000    2.000000    3.000000    2.000000    3.000000
75%            2.000000    2.000000    4.000000    3.000000    5.000000
max            3.000000    3.000000    4.000000    5.000000    5.000000

                Friends      Groups  Notifications         GPA
count        622.000000  622.000000     622.000000  622.000000
mean        1831.189711    3.625402      14.292605    2.841801
std          865.583604    1.349968      13.208275    0.658272
min         1000.000000    0.000000       5.000000    1.200000
25%         1000.000000    3.000000       5.000000    2.300000
50%         2000.000000    3.000000      10.000000    3.000000
75%         2000.000000    4.000000      20.000000    3.500000
max         4000.000000    6.000000      50.000000    3.800000
```

# DATA DESCRIBE

| Notification | The mean number of daily notification checks is approximately 14.29. The minimum and maximum values are 5 and 50, respectively. |
|---|---|
| GPA | The average grade point average is around 2.84. The minimum and maximum values are 1.2 and 3.8, respectively. |

```
In [6]: # Data Summary
        summary = data.describe()

In [7]: print(summary)
```

```
        Age Group      Gender       Level        Time    Platforms \
count  622.000000  622.000000  622.000000  622.000000  622.000000
mean     1.842444    1.562701    2.638264    2.207395    3.490354
std      0.514773    0.524834    1.186077    1.415492    1.232571
min      1.000000    1.000000    1.000000    1.000000    1.000000
25%      2.000000    1.000000    1.000000    1.000000    2.000000
50%      2.000000    2.000000    3.000000    2.000000    3.000000
75%      2.000000    2.000000    4.000000    3.000000    5.000000
max      3.000000    3.000000    4.000000    5.000000    5.000000

           Friends       Groups  Notifications         GPA
count   622.000000   622.000000     622.000000  622.000000
mean   1831.189711     3.625402      14.292605    2.841801
std     865.583604     1.349968      13.208275    0.658272
min    1000.000000     0.000000       5.000000    1.200000
25%    1000.000000     3.000000       5.000000    2.300000
50%    2000.000000     3.000000      10.000000    3.000000
75%    2000.000000     4.000000      20.000000    3.500000
max    4000.000000     6.000000      50.000000    3.800000
```

# NULL SUM

**There is no missing value, we can proceed with the data analysis directly without any modifications.**

```
In [8]: # Check for missing values
missing_values = data.isnull().sum()
print("Missing values:\n", missing_values)
```

```
Missing values:
 Age Group        0
Gender            0
Level             0
Time              0
Platforms         0
Friends           0
Groups            0
Notifications     0
GPA               0
dtype: int64
```
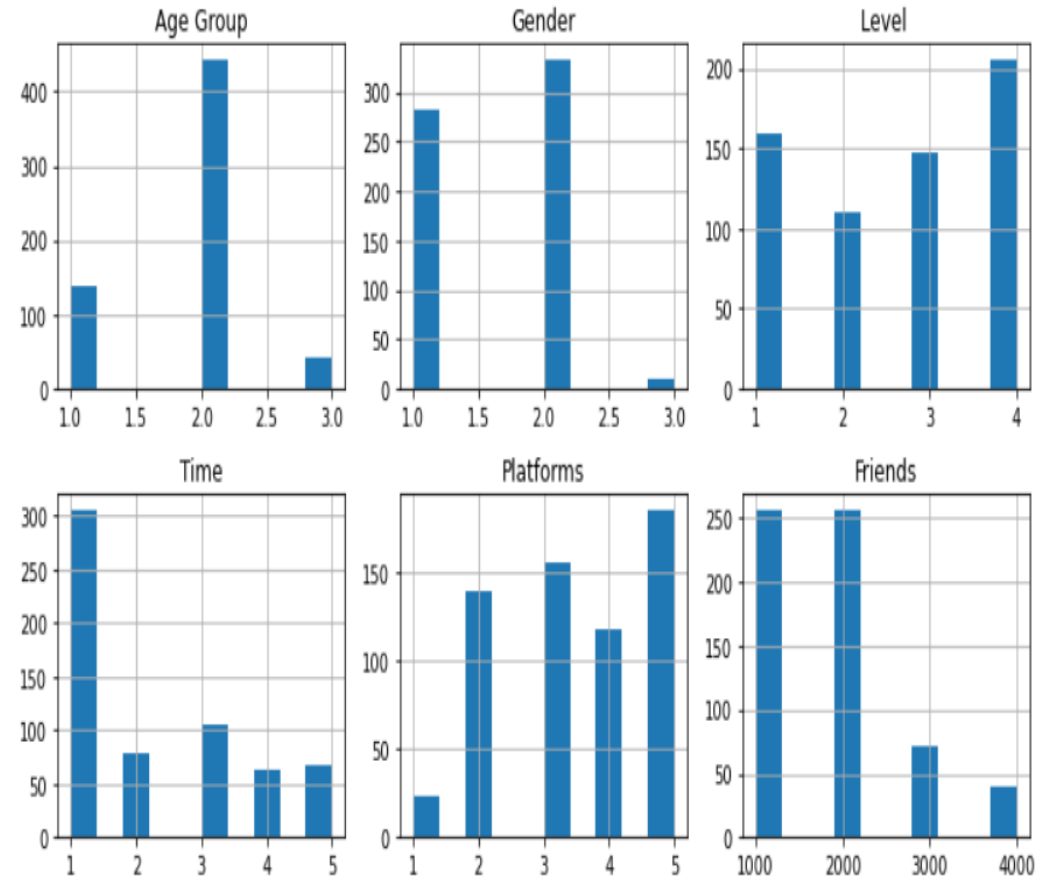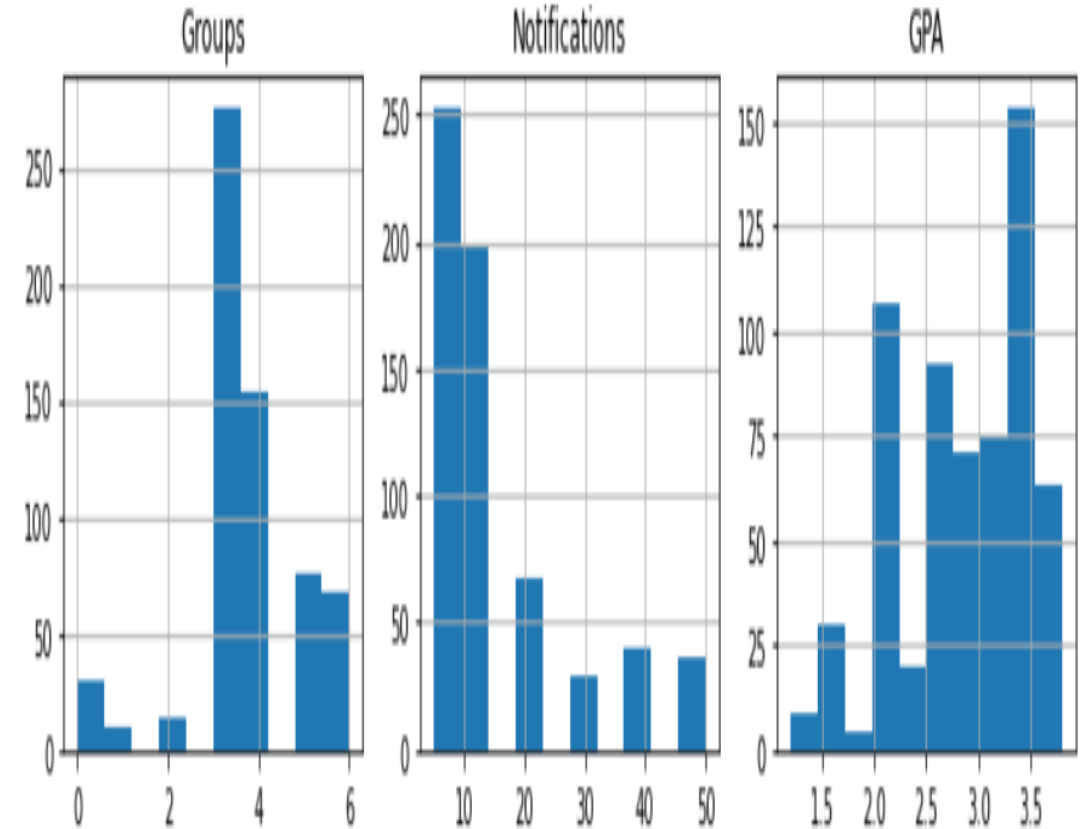
# HISTOGRAM

**This histogram is showing each variable with its count.**

```
In [9]: import pandas as pd
        import matplotlib.pyplot as plt
        import seaborn as sns

        # Histograms
        data.hist(figsize=(10, 8))
        plt.tight_layout()
        plt.show()
```

# HISTOGRAM

**This histogram is showing each variable with its count.**

Based on the given plot, it can be observed that students with approximately 1000 friends on social media have an average GPA of 3. Students with around 2000 friends have an average GPA of 2.7. Those with 3000 friends have an average GPA of approximately 2, while students with about 4000 friends have an average GPA of 2.2.

From this data, we can conclude that there seems to be a relationship between the number of friends on social media and the average GPA of students. Initially, as the number of friends increases from 1000 to 2000, the average GPA decreases slightly from 3 to 2.7. However, as the number of friends continues to increase from 2000 to 4000, the average GPA remains relatively stable around 2.2 to 2.7.

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt


# Calculate the mean GPA for different levels of the independent variables
mean_gpa_by_level = data.groupby('Friends')['GPA'].mean()

# Visualize the mean GPA by level
sns.barplot(x='Friends', y='GPA', data=data)
plt.xlabel('Friends')
plt.ylabel('Mean GPA')
plt.title('Mean GPA by Friends')
plt.show()
```
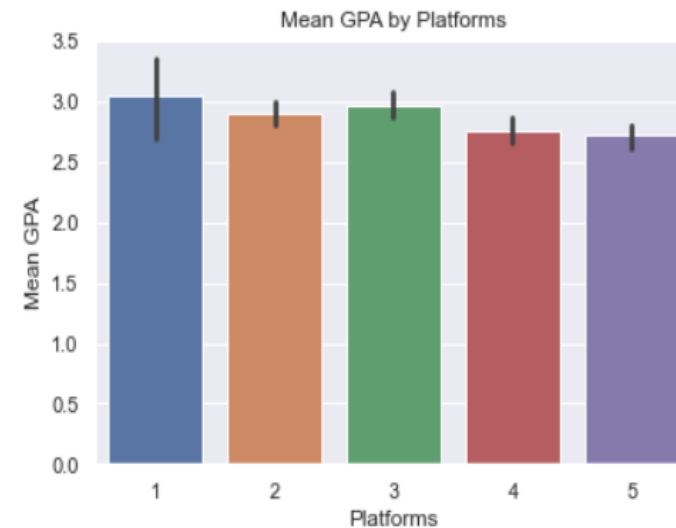
According to the visualization, students who do not belong to any groups on social media have an average GPA of 3. On the other hand, students belonging to groups numbered 1, 2, 3, 4, 5, and 6 have average GPAs of 3.2, 2.8, 2.9, 2.7, 2.5, and 2.6, respectively.

From this data, we can conclude that there appears to be a correlation between group membership on social media and students' average GPA. Students who are part of groups tend to have slightly higher or lower average GPAs compared to those who do not belong to any group. However, the specific relationship is not linear, as the average GPA fluctuates across different group numbers.

```python
In [25]: import pandas as pd
         import numpy as np
         import seaborn as sns
         import matplotlib.pyplot as plt


         # Calculate the mean GPA for different levels of the independent variables
         mean_gpa_by_level = data.groupby('Groups')['GPA'].mean()

         # Visualize the mean GPA by level
         sns.barplot(x='Groups', y='GPA', data=data)
         plt.xlabel('Groups')
         plt.ylabel('Mean GPA')
         plt.title('Mean GPA by Groups')
         plt.show()
```
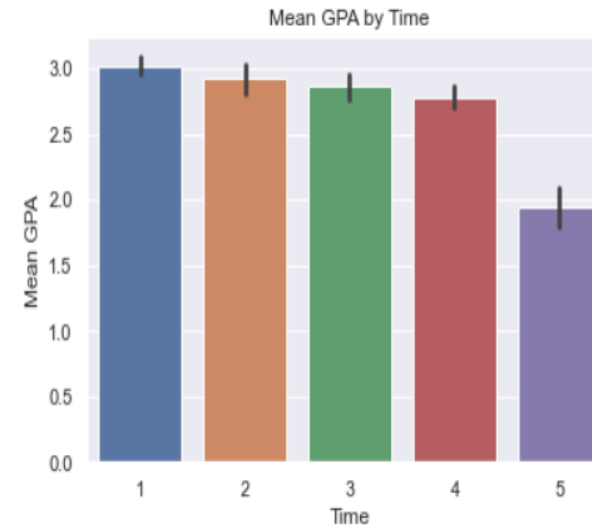
According to this visualization, students who belong to 1, 2, 3, 4, and 5 platforms have average GPAs of 3, 2.7, 2.8, 2.6, and 2.4, respectively.

From this data, we can conclude that there seems to be a relationship between the number of platforms students belong to and their average GPA. As the number of platforms increases from 1 to 3, there is a slight decrease in the average GPA. However, beyond 3 platforms, the average GPA tends to decrease more significantly with each additional platform.

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt


# Calculate the mean GPA for different levels of the independent variables
mean_gpa_by_level = data.groupby('Platforms')['GPA'].mean()

# Visualize the mean GPA by level
sns.barplot(x='Platforms', y='GPA', data=data)
plt.xlabel('Platforms')
plt.ylabel('Mean GPA')
plt.title('Mean GPA by Platforms')
plt.show()
```

According to this visualization, students who spend 1, 2, 3, 4, and 5 hours per day on social media have average GPAs of 3.0, 2.8, 2.7, 2.6, and 1.8, respectively.

From this data, we can conclude that there appears to be a negative correlation between the amount of time spent on social media and students' average GPA. As the number of hours spent on social media increases, the average GPA tends to decrease. Students who spend less time on social media tend to have higher average GPAs compared to those who spend more time.

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt


# Calculate the mean GPA for different levels of the independent variables
mean_gpa_by_level = data.groupby('Time')['GPA'].mean()

# Visualize the mean GPA by level
sns.barplot(x='Time', y='GPA', data=data)
plt.xlabel('Time')
plt.ylabel('Mean GPA')
plt.title('Mean GPA by Time')
plt.show()
```
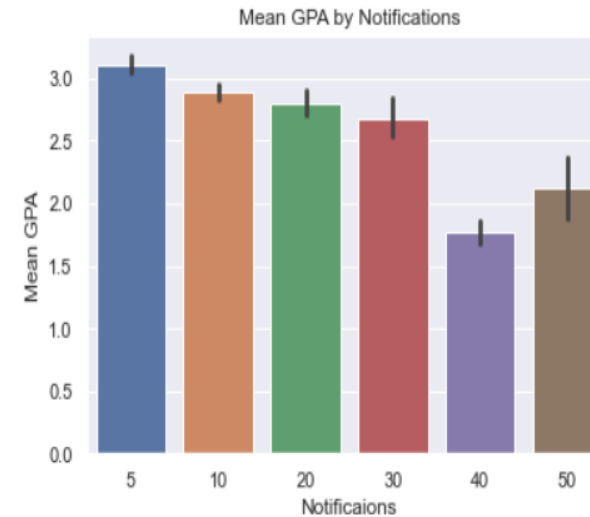
According to this visualization, students who have 5, 10, 20, 30, 40, and 50 notifications on social media have average GPAs of 3.0, 2.8, 2.7, 2.6, 1.5, and 1.8, respectively.

From this data, we can conclude that there appears to be a negative correlation between the number of notifications students receive on social media and their average GPA. As the number of notifications increases, the average GPA tends to decrease. Students with fewer notifications tend to have higher average GPAs compared to those with a higher number of notifications.

```python
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt


# Calculate the mean GPA for different levels of the independent variables
mean_gpa_by_level = data.groupby('Notifications')['GPA'].mean()

# Visualize the mean GPA by level
sns.barplot(x='Notifications', y='GPA', data=data)
plt.xlabel('Notificaions')
plt.ylabel('Mean GPA')
plt.title('Mean GPA by Notifications')
plt.show()
```
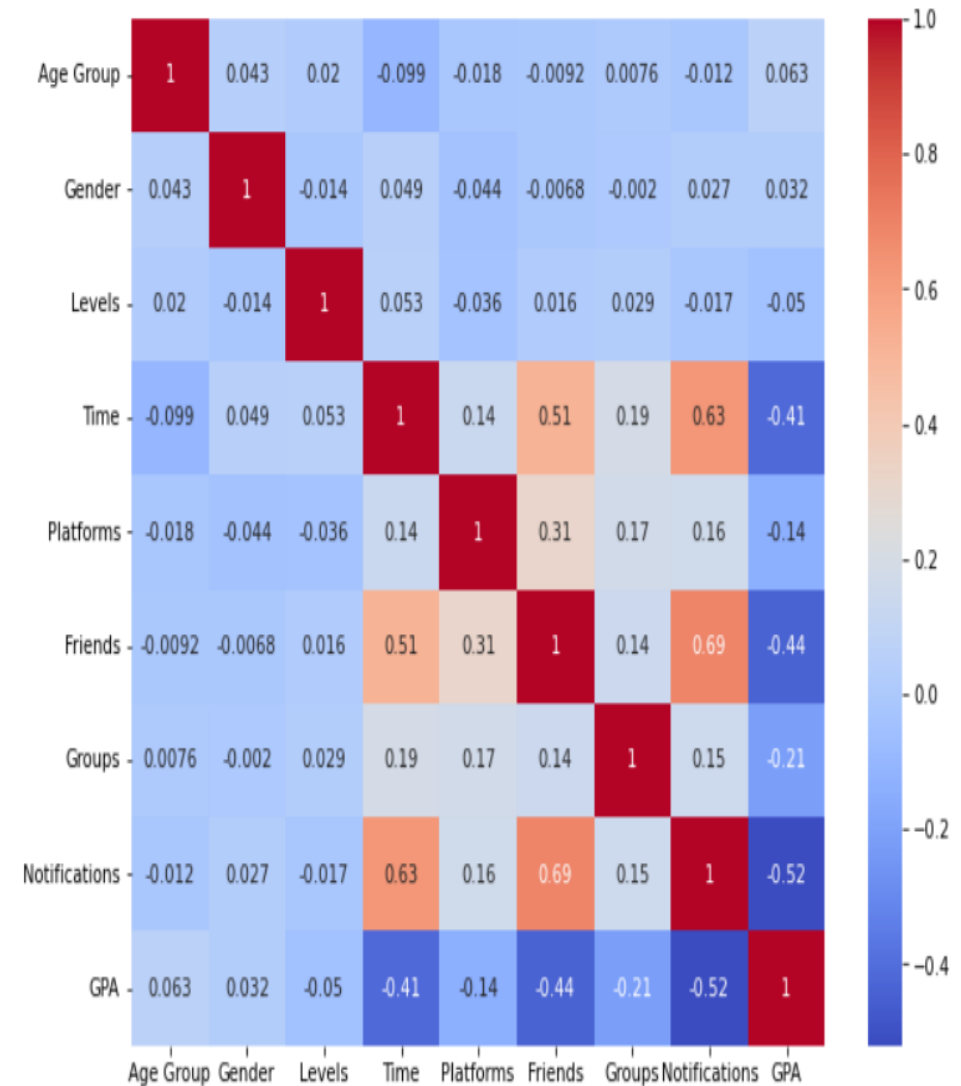

Mean GPA by Notifications

# Heatmap

By visualizing the correlation matrix using a heatmap, we easily identify the relationships between variables. Positive correlations are indicated by warmer colors (closer to red), while negative correlations are indicated by cooler colors (closer to blue). The intensity of the color represents the strength of the correlation.

For example, A correlation coefficient of -0.41 between "Time" and "GPA" indicates a moderate negative correlation between these two variables.

```python
# Heatmap
corr_matrix = data.corr()
plt.figure(figsize=(10, 8))
sns.heatmap(corr_matrix, annot=True, cmap='coolwarm')
plt.show()
```

```
In [39]: import statsmodels.api as sm

# Assuming 'X' is your independent variable matrix and 'y' is your dependent variable (GPA)
X = data[['Time', 'Platforms', 'Friends', 'Groups', 'Notifications']]
y = data['GPA']

# Add a constant term to the independent variable matrix
X = sm.add_constant(X)

# Fit a linear regression model
model = sm.OLS(y, X).fit()

# Get the coefficients (impact) of each independent variable on GPA
coefficients = model.params[1:]

# Find the variable with the most negative impact on GPA
most_negative_variable = coefficients.idxmin()
negative_impact = coefficients.min()

print("Variable with the most negative impact on GPA:", most_negative_variable)
print("Negative impact on GPA:", negative_impact)


Variable with the most negative impact on GPA: Groups
Negative impact on GPA: -0.05924710363369366
```
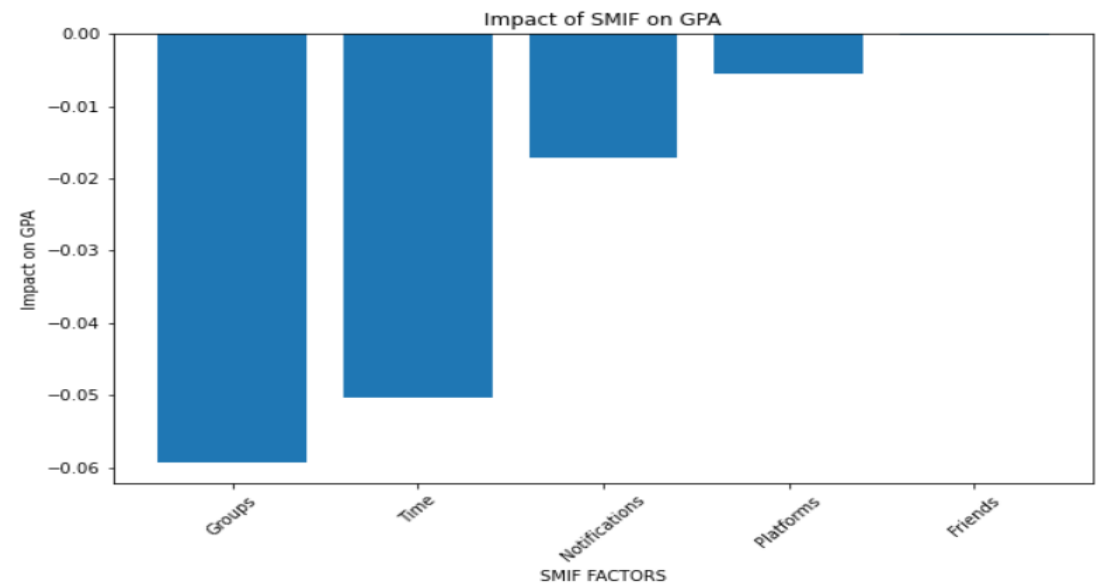
```
In [41]: import matplotlib.pyplot as plt

# Assuming 'coefficients' contains the coefficients obtained from the linear regression model

# Sort the coefficients in ascending order
sorted_coefficients = coefficients.sort_values()

# Plot the variables and their impacts
plt.figure(figsize=(10, 6))
plt.bar(sorted_coefficients.index, sorted_coefficients)
plt.xlabel('SMIF FACTORS')
plt.ylabel('Impact on GPA')
plt.title('Impact of SMIF on GPA')
plt.xticks(rotation=45)
plt.show()
```

# B/II- Conclusion

The data analysis has yielded valuable insights into the nuanced relationship between social media influence and academic performance.

Notably, students with higher influence metrics on social media platforms may exhibit distinct academic trends.

# 3- Technology Review
## C/I- Introduction

In the ever-evolving landscape of technology, a comprehensive review is indispensable to navigate through the myriad of tools and innovations. This technology review serves as a guide to understand, evaluate, and select the most pertinent technologies for our project, "Predictive Analysis of Academic Performance via Social Media Influence Metrics."

The importance of this technology review lies in its role as a compass, steering our project towards efficient and effective solutions. By exploring and assessing relevant technologies, we aim to align our research goals with cutting-edge tools that can enhance the predictive analysis of academic performance.

# C/II- Technology Overview

The primary focus of our review is on **machine learning technologies**, specifically those used in predictive analytics.

**Machine learning**, with its ability to discern patterns and trends from data, plays a crucial role in our project. Key features include algorithmic versatility, interpretability, and the capacity to handle large datasets. Commonly used in data-driven fields, machine learning is a cornerstone technology for predictive modeling and analysis.

# C/III-    Relevance

Machine learning is particularly relevant to our project due to its suitability for predictive analysis.

By leveraging machine learning algorithms, we can discern intricate relationships between social media usage patterns and academic performance. This technology addresses the challenge of deciphering complex data sets, contributing to the success of our research by providing actionable insights into the factors influencing academic outcomes.

# C/IV- Comparaison and Evaluation

In comparing machine learning technologies, we assess several factors:

- **Scikit-learn vs. TensorFlow:** Scikit-learn is chosen for its simplicity and interpretability, making it suitable for our project's goals. TensorFlow, although powerful, might be more complex than necessary for our relatively straightforward analysis.

- **Supervised Learning Algorithms:** We evaluate algorithms such as Random Forest and Logistic Regression. Random Forest offers robustness for complex datasets, while Logistic Regression provides interpretability.

- **Ease of Integration:** Scikit-learn's compatibility with Python, a language widely used in data science, ensures ease of integration with our existing workflows.

# C/V- Use Cases and Examples

Real-world use cases highlight the applicability of machine learning in similar projects:

- **Youtube Recommendation System:** Demonstrates how machine learning predicts user preferences, akin to our prediction of academic performance based on social media influence.

- **Customer Churn Prediction in Telecom:** Drawing parallels, this use case showcases the predictive power of machine learning in understanding and forecasting behavior.

# C/VI- Identify Gaps and Research Opportunities

While machine learning offers robust solutions, potential gaps include over-reliance on historical data and the need for continual adaptation to evolving social media trends.

Research opportunities lie in refining algorithms to account for dynamic online behaviors and exploring emerging machine learning frameworks.

# C/VII- Conclusion

In conclusion, this technology review underscores the pivotal role of machine learning in our project. Its versatility, interpretability, and relevance to predictive analysis make it an indispensable tool.

By selecting and implementing appropriate machine learning technologies, we pave the way for accurate predictions and actionable recommendations in our exploration of social media's impact on academic performance.

# REFERENCES

▶ https://doi.org/10.6084/m9.figshare.14905

▶ https://www.researchgate.net/publication/362079186_Measuring_the_effect_of_social_media_on_student_academic_performance_using_a_social_media_influence_factor_model

# Thank you