# STUDENT PERFORMANCE FACTORS

**Visual Analytics Techniques Homework**

Merve Pakcan Tufenk

This project seeks to analyze and understand the factors that influence student performance.
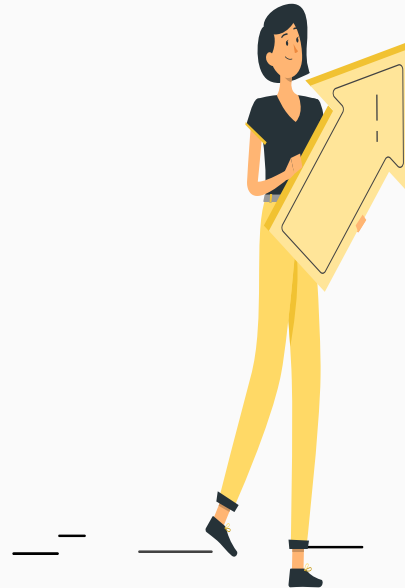
**USER:** Parents, administrators, educators, psychologists, counsellors

**DATA:** A comprehensive dataset of student performance metrics

**TASK:** Uncovering patterns, relationships, and key drivers of academic success to inform effective interventions

**SDG 4**
Quality Education

**SDG 10**
Reduced Inequalities

Ssas
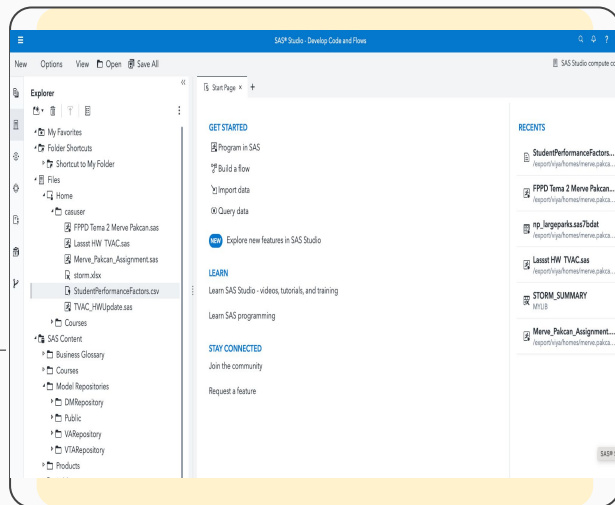
# 6600 ROWS

## 20 COLUMNS

Student Performance Factors dataset from Kaggle

Uploaded CAS using *SAS Studio - Develop Code and Flows* menu.

Loaded file into the CAS environment, within the CASUSER library, using *Manage Data* menu, **indicated by the green 'In-memory data' status.**
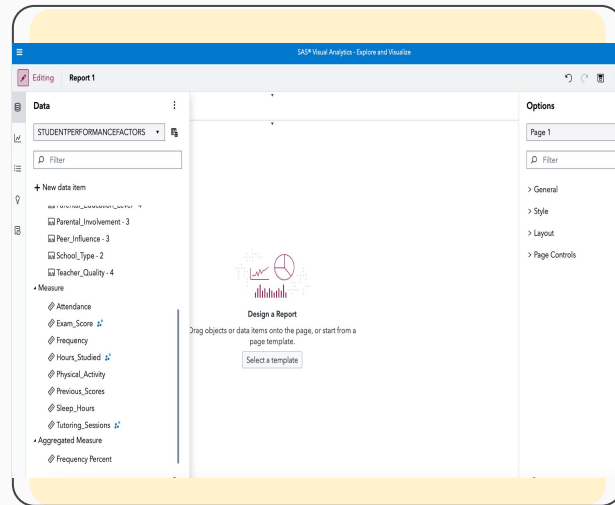
13 categorical variables, 7 numerical variables in *Explore and Visualize* menu.

# DATA WRANGLING



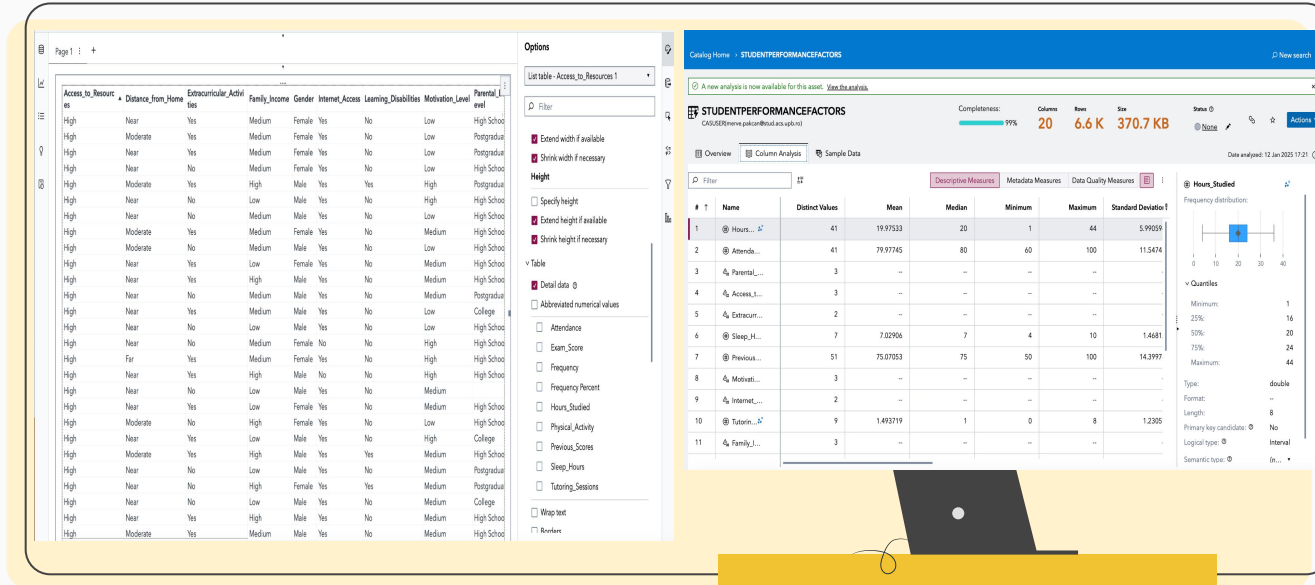List table created, **no new data items** are needed. Data analyzed in Discover Information Assets part. Cross Tabulation is not proper due to high number of categorical variables

# 5836 ROWS

Uploaded clean data to discover information assets part, shows that data is clean to analyze

## DATA PREPARATION



```
1   libname mylib clear;
2
3   /* Defining the "libmerv1" library to access the directory containing the data files */
4   libname libmerv1 '/export/viya/homes/merve.pakcan@stud.acs.upb.ro/casuser/';
5
6   /* Importing the dataset */
7 ⊖ proc import datafile='/export/viya/homes/merve.pakcan@stud.acs.upb.ro/casuser/StudentPerformanceFactors.csv'
8       dbms=csv out=libmerv1.studperformance replace;
9       guessingrows=12767; /* Ensures SAS scans all rows to correctly determine column types */
10      getnames=yes;
11  run;
12
13  /* Defining a macro variable for the dataset */
14  %let dataset = libmerv1.studperformance
15
16  /* Displaying dataset structure and variables */
17 ⊖ proc contents data=&dataset;
18  run;
19
20  /* Checking a sample of the dataset to ensure values loaded correctly */
21 ⊖ proc print data=libmerv1.studperformance(obs=10);
22  run;
23
24  /* Identifying Missing Values */
25  /* Using proc means with nmiss to identify missing values in numeric columns. */
26 ⊖ proc means data=libmerv1.studperformance n nmiss;
27      var Hours_Studied Attendance Sleep_Hours Previous_Scores Exam_Score;
28  run;
29
30  /* Checking for Missing Values in all categorical variables */
31 ⊖ proc freq data=libmerv1.studperformance;
32      tables _all_ / missing;
33  run;
34
35  /* Discovered missing values in the categorical variables */
36  /* Missing values: Teacher_Quality, Parental_Education_Level, and Distance_from_Home */
37 ⊖ proc freq data=&dataset;
38      tables Teacher_Quality Parental_Education_Level Distance_from_Home / missing;
39  run;
40
```

The FREQ Procedure

| Teacher_Quality | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
|  | 78 | 1.18 | 78 | 1.18 |
| High | 1947 | 29.47 | 2025 | 30.65 |
| Low | 657 | 9.94 | 2682 | 40.59 |
| Medium | 3925 | 59.41 | 6607 | 100.00 |

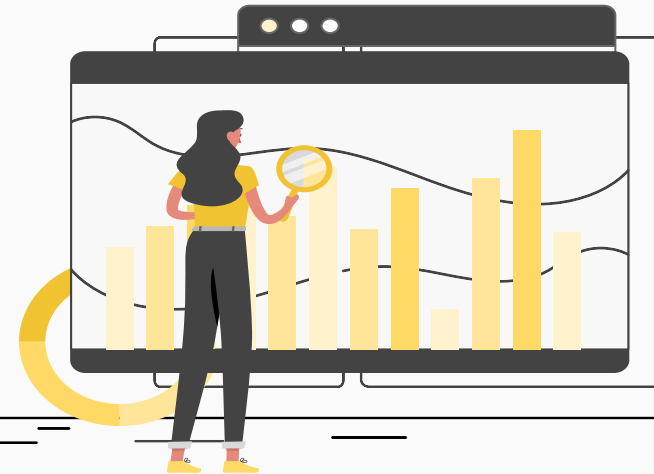| Parental_Education_Level | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
|  | 90 | 1.36 | 90 | 1.36 |
| College | 1989 | 30.10 | 2079 | 31.47 |
| High School | 3223 | 48.78 | 5302 | 80.25 |
| Postgraduate | 1305 | 19.75 | 6607 | 100.00 |

| Distance_from_Home | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
|  | 67 | 1.01 | 67 | 1.01 |
| Far | 658 | 9.96 | 725 | 10.97 |
| Moderate | 1998 | 30.24 | 2723 | 41.21 |
| Near | 3884 | 58.79 | 6607 | 100.00 |

The MEANS Procedure

| Variable | N | Mean | Std Dev | Minimum | Maximum |
|---|---|---|---|---|---|
| Hours_Studied | 5836 | 20.0080535 | 5.7822916 | 4.0000000 | 36.0000000 |
| Attendance | 5836 | 80.0219328 | 11.4995264 | 60.0000000 | 100.0000000 |
| Sleep_Hours | 5836 | 7.0412954 | 1.4696616 | 4.0000000 | 10.0000000 |
| Previous_Scores | 5836 | 75.1077793 | 14.3481094 | 50.0000000 | 100.0000000 |
| Exam_Score | 5836 | 66.9883482 | 3.2301389 | 59.0000000 | 75.0000000 |

---

ⓘ A new analysis is now available for this asset. View the analysis.

### CLEANED_NOP_DATA
CASUSER(merve.pakcan@stud.acs.upb.ro)

Completeness: 100%  
Columns **20**  
Rows **5.8 K**  
Size **335.5 KB**  
Status ⓘ None

Actions

Overview | Column Analysis | Sample Data

Date analyzed: 14 Jan 2025 20:10

**Sensitive**
Information Privacy

(none found)
Time Period Covered

(none found)
Top Areas Covered

(none found)
Top Languages

∨ Summary

This dataset describes information about the following entities: **gender, family name**. The most important column is **Exam_Score**. The storage format is **CAS**. The data has values that could be considered **private**.

∨ Description

Describe the purpose of this asset and what it tracks or measures.

|  | varchar | | | double | |
|---|---|---|---|---|---|
| Filter | | | | | |
| Name/Label | Length | Semantic Type | Informatio... | Terms | |
| ⊞ Hours_Studied | 8 | (none) ▾ | None | (none) ▾ | |
| ⊞ Attendance | 8 | (none) ▾ | None | (none) ▾ | |
| ⊞ Parental_Involvement | 6 | (none) ▾ | None | (none) ▾ | |
| ⊞ Access_to_Resources | 6 | (none) ▾ | None | (none) ▾ | |
| ⊞ Extracurricular_Activ... | 3 | (none) ▾ | None | (none) ▾ | |
| ⊞ Sleep_Hours | 8 | (none) ▾ | None | (none) ▾ | |
| ⊞ Previous_Scores | 8 | (none) ▾ | None | (none) ▾ | |

∨ Contacts (0)
No contacts are assigned.

∨ Tags (0)
No tags are assigned.

∨ Properties

Asset type: In-memory data

Date modified: 14 Jan 2025 20:07
Modified by: --

Date created: 14 Jan 2025 20:07
Created by: merve.pakcan@stud.acs...

Last accessed: 14 Jan 2025 20:09
Accessed by: merve.pakcan@stud.acs...

Library: CASUSER(merve.pakca...

Eliminated missing values and outliers(small amount compared to dataset size), with using SAS Develop and Code part. No duplicate found.

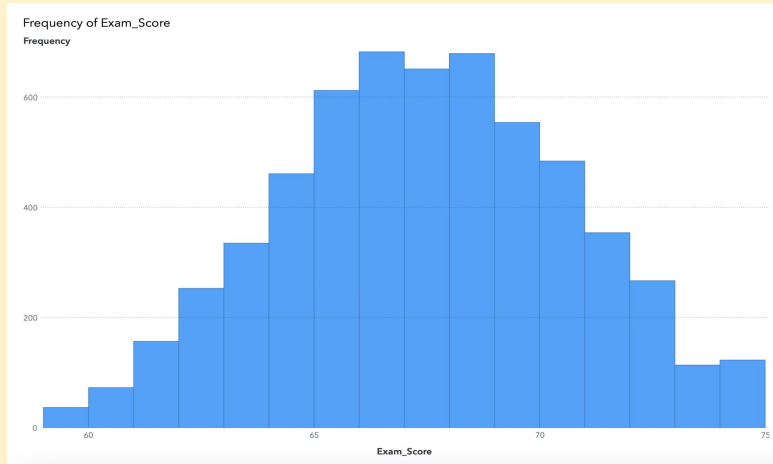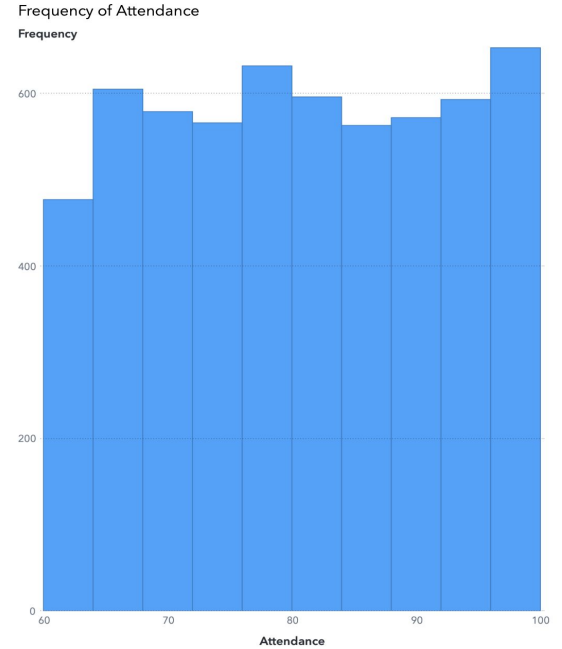| QUESTION | SUB-QUESTION | VARIABLES | NOTES |
|---|---|---|---|
| What factors most influence student performance? | Which variables have the strongest correlation with Exam_Score? | Dependent Variable: Exam_Score<br><br>Independent Variable: Attendance, Hours_Studied, Previous_Scores, Parental_Involvement, Access_to_Resources, Extracurricular_Activities, Sleep_Hours, Previous_Scores, Motivation_Level, Internet_Access, Tutoring_Sessions, Family_Income, Teacher_Quality, School_Type, Peer_Influence, Physical_Activity, Learning_Disabilities, Parental_Education_Level, Distance_from_Home, Gender | Needed to identify key variables of success. |
| What socioeconomic factors most impact exam results? | Which socioeconomic variables have the strongest correlation with Exam_Score? | Dependent Variable: Exam_Score<br><br>Independent Variable: Parental_Education_Level, Family_Income, Access_to_Resources ,Internet_Access, Distance_from_Home, Peer_Influence, Parental_Involvement | Exploring educational equity and disparities. |

## Target Variable



Frequency of Exam_Score
Frequency

**Shape:** Bell-shaped distribution, **Center:** Around 67, **Spread:** Between 59 and 75



Frequency of Attendance
Frequency

Uniform distribution

Scatter Plot of Selected Measures

Hours_Studied

Exam_Score

Positive correlation between hours studied and exam scores: as study hours increase, exam scores tend to increase



Exam_Score by Parental_Involvement
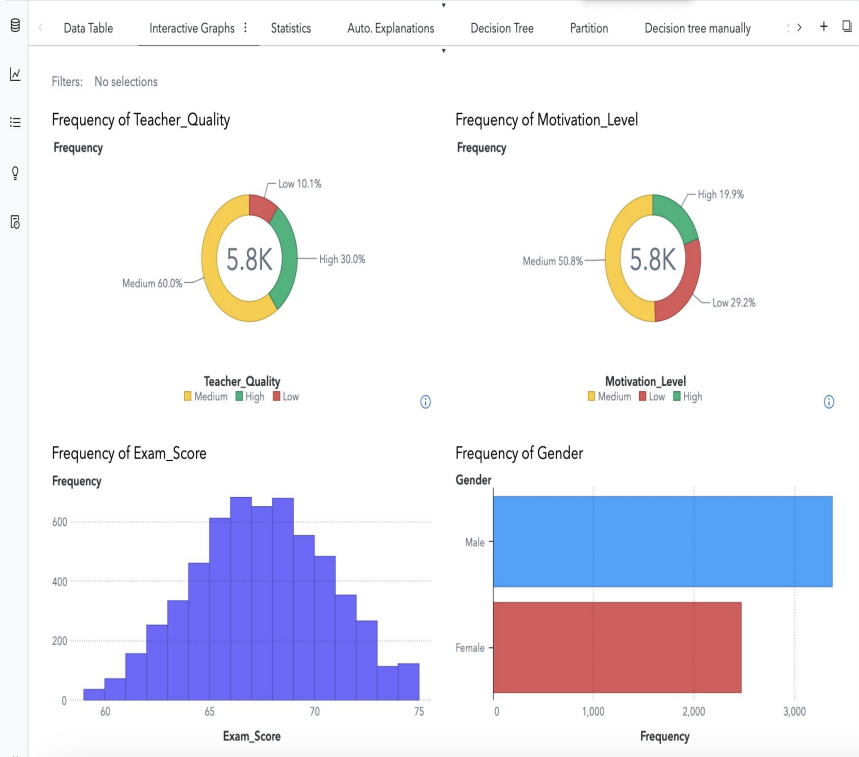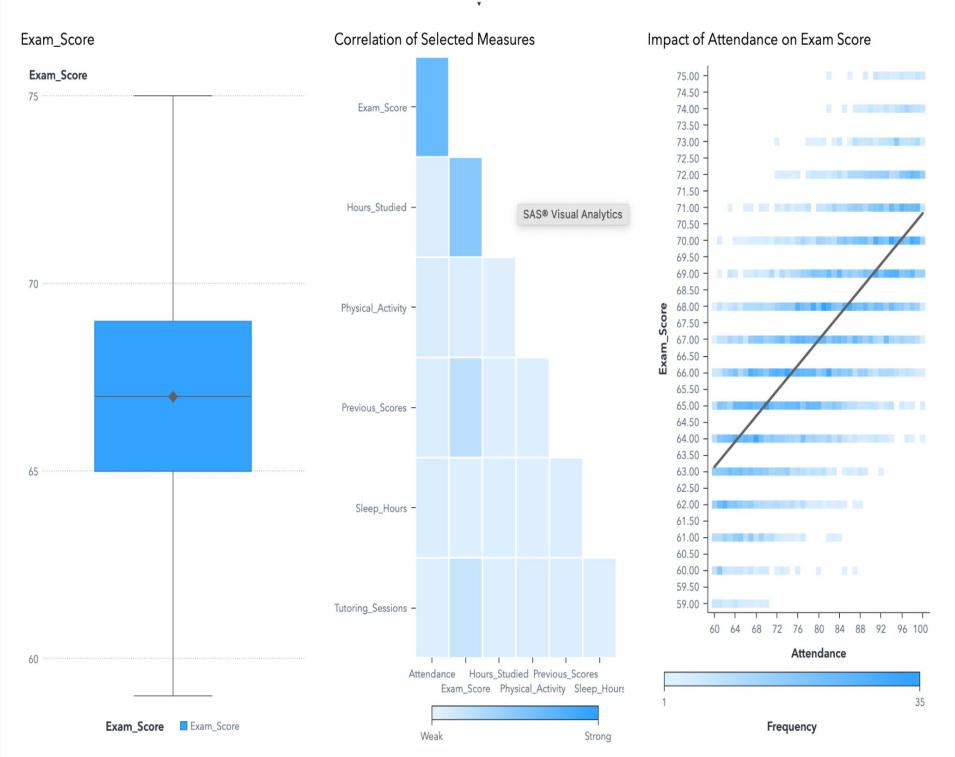
Exam_Score

Parental_Involvement

Exam_Score ■ Exam_Score

Positive correlation between parental engagement and academic success

The distribution of **Exam_Score** shows a median around 67. The range is relatively narrow. In correlation matrix and heatmap, a strong linear relationship exists between **Attendance** and **Exam_Score**
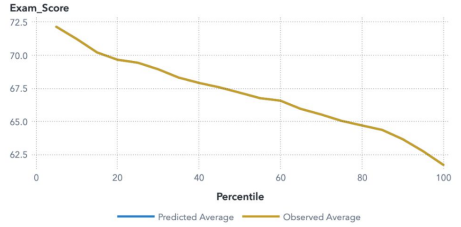
Additional interactive graphs created to understand variables more.

As it seen from automatic explanations and decision tree, the three most related factors are **Attendance**, Hours_Studied, Previous_Score
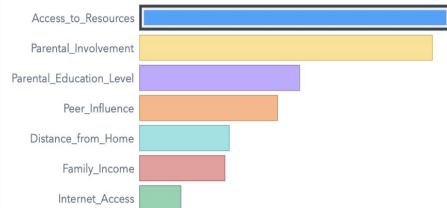
# FOR QUESTION 2

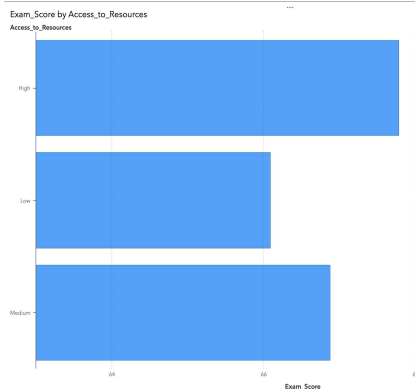What socioeconomic factors most impact exam results?



## What are the characteristics of Exam_Score?

Exam_Score ranges from 59 to 75. Average Exam_Score is 67. Most cases (the middle 80%) have an Exam_Score between 63 and 71. Access_to_Resources best differentiates the highest (top 10%) and the lowest (bottom 10%) Exam_Score cases. The three most related factors are Access_to_Resources, Parental_Involvement, and Parental_Education_Level.

### What factors are most related to Exam_Score?

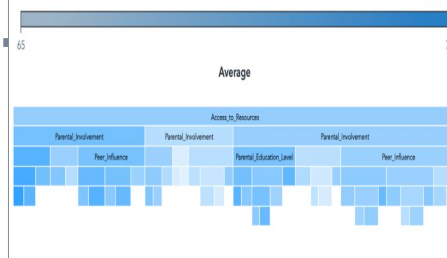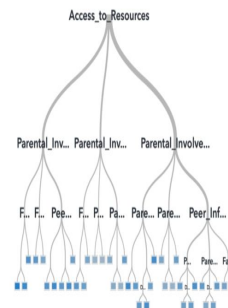### What is the relationship between Exam_Score and Access_to_Resources?

When Access_to_Resources is High, the average of Exam_Score is a high value. When Access_to_Resources is Low, the average of Exam_Score is a low value. The most common Access_to_Resources value is Medium.

As it seen from automatic explanations and decision tree, the three most related socioeconomic factors are **Access_to_Resources**, Parental_Involvement, Parental_Education_Level.
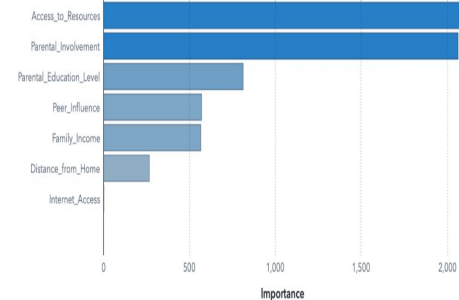
# THANK YOU!

What was easy: Exploring patterns with visualization part
What was challenging:Some technical issues at Sas

| QUESTION | FINDINGS | DECISONS |
|---|---|---|
| What factors most influence student performance? | **Attendance** and **Hours_Studied** are most related factors | The findings can support policymakers, educators, and institutions in designing targeted programs to improve student performance and reduce educational inequalities. |
| What socioeconomic factors most impact exam results? | **Access_to_Resources**, **Parental_Involvement** are most related factors | • Empower families through community outreach to enhance parental engagement.<br>• Promoting attendance and study habits while offering support for students that has low score |