# Final Project: Analysis of Student Performance Data

## Introduction

The purpose of this project is to analyse a real dataset on student performance to identify factors affecting academic success. By cleaning, exploring, and analysing this dataset, the goal is to provide meaningful insights into student performance trends. This report showcases SAS skills learned during the course.

## Descriptive Picture

*A relevant bar chart, such as average performance across genders or parental education levels, will be included.*

## Accessing the Data

### Data Reference

- **Dataset**: *Student Performance Data*
- **Source**: Kaggle

### Data Origin and Purpose

The dataset was collected to understand students' academic performance across multiple demographics, such as gender, study time, parental education, and test scores in math, reading, and writing. The purpose is to explore factors influencing success.

**SAS Code to Access Data**

```
/* Step 1: Import the Dataset */

datafile="/home/sgnanas/EPG294/data/Student_performance_data_.csv"
out=work.student_data
dbms=csv replace;
getnames=yes;
run;
```

## Exploring the Data

*SAS Code for Exploration*

```
/* Step 2: Explore the Dataset */

proc contents data=student_data;
run;
proc freq data=student_data;
      tables Gender ParentalEducation Tutoring;
run;
proc means data=student_data mean stddev min max;
       var GPA StudyTimeWeekly Absences;
run;
```

## The CONTENTS Procedure

| Data Set Name | WORK.STUDENT_DATA | | Observations | 2392 |
|---|---|---|---|---|
| Member Type | DATA | | Variables | 15 |
| Engine | V9 | | Indexes | 0 |
| Created | 12/04/2024 14:07:30 | | Observation Length | 120 |
| Last Modified | 12/04/2024 14:07:30 | | Deleted Observations | 0 |
| Protection | | | Compressed | NO |
| Data Set Type | | | Sorted | NO |
| Label | | | | |
| Data Representation | SOLARIS_X86_64, LINUX_X86_64, ALPHA_TRU64, LINUX_IA64, LINUX_POWER_64 | | | |
| Encoding | utf-8 Unicode (UTF-8) | | | |

| Engine/Host Dependent Information | |
|---|---|
| Data Set Page Size | 65536 |
| Number of Data Set Pages | 5 |
| First Data Page | 1 |
| Max Obs per Page | 545 |
| Obs in First Data Page | 513 |
| Number of Data Set Repairs | 0 |
| Filename | /saswork/SAS_workE6810001A685_6yn-vya-p-web01.server.clemson.edu/student_data.sas7bdat |
| Release Created | V.0305M0 |
| Host Created | Linux |
| Inode Number | 269938831 |
| Access Permission | rw-r--r-- |
| Owner Name | sgnanas |
| File Size | 384KB |
| File Size (bytes) | 393216 |

| Alphabetic List of Variables and Attributes | | | | | |
|---|---|---|---|---|---|
| # | Variable | Type | Len | Format | Informat |
| 7 | Absences | Num | 8 | BEST12. | BEST32. |
| 2 | Age | Num | 8 | BEST12. | BEST32. |
| 4 | Ethnicity | Num | 8 | BEST12. | BEST32. |
| 10 | Extracurricular | Num | 8 | BEST12. | BEST32. |
| 14 | GPA | Num | 8 | BEST12. | BEST32. |
| 3 | Gender | Num | 8 | BEST12. | BEST32. |
| 15 | GradeClass | Num | 8 | BEST12. | BEST32. |
| 12 | Music | Num | 8 | BEST12. | BEST32. |
| 5 | ParentalEducation | Num | 8 | BEST12. | BEST32. |
| 9 | ParentalSupport | Num | 8 | BEST12. | BEST32. |
| 11 | Sports | Num | 8 | BEST12. | BEST32. |
| 1 | StudentID | Num | 8 | BEST12. | BEST32. |
| 6 | StudyTimeWeekly | Num | 8 | BEST12. | BEST32. |
| 8 | Tutoring | Num | 8 | BEST12. | BEST32. |
| 13 | Volunteering | Num | 8 | BEST12. | BEST32. |

## The FREQ Procedure

| Gender | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| 0 | 1170 | 48.91 | 1170 | 48.91 |
| 1 | 1222 | 51.09 | 2392 | 100.00 |

| ParentalEducation | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| 0 | 243 | 10.16 | 243 | 10.16 |
| 1 | 728 | 30.43 | 971 | 40.59 |
| 2 | 934 | 39.05 | 1905 | 79.64 |
| 3 | 367 | 15.34 | 2272 | 94.98 |
| 4 | 120 | 5.02 | 2392 | 100.00 |

| Tutoring | Frequency | Percent | Cumulative Frequency | Cumulative Percent |
|---|---|---|---|---|
| 0 | 1671 | 69.86 | 1671 | 69.86 |
| 1 | 721 | 30.14 | 2392 | 100.00 |

## The MEANS Procedure

| Variable | Mean | Std Dev | Minimum | Maximum |
|---|---|---|---|---|
| GPA | 1.9061863 | 0.9151558 | 0 | 4.0000000 |
| StudyTimeWeekly | 9.7719919 | 5.6527742 | 0.0010565 | 19.9780940 |
| Absences | 14.5413880 | 8.4674174 | 0 | 29.0000000 |

**Exploration Results**

- **Frequency Distribution**: Most students did not complete a test preparation course.
- **Descriptive Statistics**: Average math, reading, and writing scores show consistent trends, with slight variations across genders and parental education levels.

## **Preparing the Data**

### *Cleaning the Data*

- Removed missing values.
- Standardized categorical variables (e.g., consistent capitalization).
- Added a new variable: `total_score` (sum of math, reading, and writing scores).

### *SAS Code for Cleaning*

```
/* Step 3: Clean and Transform the Data */

data clean_data;
     set student_data;
     where GPA is not missing;
     total_score = GPA + StudyTimeWeekly;

/* Create performance categories */

if total_score >= 90 then performance = 'High';
else if total_score >= 60 then performance = 'Medium';
else performance = 'Low';
run;

proc print data=clean_data(obs=10);
run;
```

| Obs | StudentID | Age | Gender | Ethnicity | ParentalEducation | StudyTimeWeekly | Absences | Tutoring | ParentalSupport | Extracurricular | Sports | Music | Volunteering | GPA | GradeClass | total_score | performance |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1001 | 17 | 1 | 0 | 2 | 19.833722808 | 7 | 1 | 2 | 0 | 0 | 1 | 0 | 2.9291955917 | 2 | 22.7629 | Low |
| 2 | 1002 | 18 | 0 | 0 | 1 | 15.408756056 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3.0429148334 | 1 | 18.4517 | Low |
| 3 | 1003 | 15 | 0 | 2 | 3 | 4.2105697688 | 26 | 0 | 2 | 0 | 0 | 0 | 0 | 0.1126022545 | 4 | 4.3232 | Low |
| 4 | 1004 | 17 | 1 | 0 | 3 | 10.028829474 | 14 | 0 | 3 | 1 | 0 | 0 | 0 | 2.0542181397 | 3 | 12.0830 | Low |
| 5 | 1005 | 17 | 1 | 0 | 2 | 4.672495273 | 17 | 1 | 3 | 0 | 0 | 0 | 0 | 1.2880611818 | 4 | 5.9606 | Low |
| 6 | 1006 | 18 | 0 | 0 | 1 | 8.1912185453 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 3.0841836145 | 1 | 11.2754 | Low |
| 7 | 1007 | 15 | 0 | 1 | 1 | 15.601680475 | 10 | 0 | 3 | 0 | 1 | 0 | 0 | 2.7482374149 | 2 | 18.3499 | Low |
| 8 | 1008 | 15 | 1 | 1 | 4 | 15.424496306 | 22 | 1 | 1 | 1 | 0 | 0 | 0 | 1.3601427123 | 4 | 16.7846 | Low |
| 9 | 1009 | 17 | 0 | 0 | 0 | 4.562007558 | 1 | 0 | 2 | 0 | 1 | 0 | 1 | 2.8968191895 | 2 | 7.4588 | Low |
| 10 | 1010 | 16 | 1 | 0 | 1 | 18.444466363 | 0 | 0 | 3 | 1 | 0 | 0 | 0 | 3.5734742103 | 0 | 22.0179 | Low |

## Variable Explanation

| Variable Name | Description |
|---|---|
| gender | Student's gender |
| parental_education | Education level of parents |
| test_prep_course | Whether the student took a prep course |
| math_score | Math test score |
| reading_score | Reading test score |
| writing_score | Writing test score |
| total_score | Combined total score of all tests |

## Analysing and Reporting

### Analysis Methods

The analysis focused on:

- Identifying performance differences across gender, parental education, and test preparation courses.
- Regression analysis to determine predictors of total score.

### SAS Code for Analysis

```
/* 4.1: Descriptive Statistics */
proc means data=clean_data mean stddev min max;
```

```sas
        var total_score GPA StudyTimeWeekly Absences;
run;

/* 4.2: Grouped Analysis by Gender and Parental Education */
proc means data=clean_data;
        class Gender ParentalEducation;
        var total_score;
run;

/* 4.3: Correlation Analysis */
proc corr data=clean_data;
        var GPA StudyTimeWeekly Absences total_score;
 run;

/* 4.4: Regression Analysis */
proc reg data=clean_data;
model total_score = GPA StudyTimeWeekly Absences; run;

/* 4.5: Clustering Analysis */
proc fastclus data=clean_data maxclusters=3 out=cluster_results;
 var GPA StudyTimeWeekly Absences;
run;

proc print data=cluster_results(obs=10);
run;
```

## The MEANS Procedure

| Variable | Mean | Std Dev | Minimum | Maximum |
|---|---|---|---|---|
| total_score | 11.6781782 | 5.8861029 | 0.0080309 | 23.4243982 |
| GPA | 1.9061863 | 0.9151558 | 0 | 4.0000000 |

| | Mean | Std Dev | Minimum | Maximum |
|---|---|---|---|---|
| StudyTimeWeekly | 9.7719919 | 5.6527742 | 0.0010565 | 19.9780940 |
| Absences | 14.5413880 | 8.4674174 | 0 | 29.0000000 |

## The MEANS Procedure

### Analysis Variable : total_score

| Gender | ParentalEducation | N Obs | N | Mean | Std Dev | Minimum | Maximum |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 123 | 123 | 11.8878815 | 5.7045309 | 0.5856499 | 23.1283802 |
| | 1 | 354 | 354 | 11.6607085 | 5.8476118 | 0.0080309 | 23.4029968 |
| | 2 | 458 | 458 | 11.9034362 | 5.7396298 | 0.6198790 | 23.4243982 |
| | 3 | 175 | 175 | 10.8485799 | 5.7445784 | 0.7475382 | 22.4503157 |
| | 4 | 60 | 60 | 11.0034190 | 5.6573920 | 2.3674305 | 21.1653415 |
| 1 | 0 | 120 | 120 | 11.1148772 | 5.4692284 | 0.8999928 | 22.6980997 |
| | 1 | 374 | 374 | 11.9241070 | 6.0060432 | 0.1126263 | 22.9443138 |
| | 2 | 476 | 476 | 11.7793208 | 6.0211839 | 0.1008619 | 23.0462702 |
| | 3 | 192 | 192 | 11.5755672 | 6.1640490 | 0.2491401 | 23.2759232 |
| | 4 | 60 | 60 | 11.8459134 | 6.3436576 | 1.5634879 | 21.4993713 |

## The CORR Procedure

| 4 Variables: | GPA StudyTimeWeekly Absences total_score |
|---|---|

### Simple Statistics

| Variable | N | Mean | Std Dev | Sum | Minimum | Maximum |
|---|---|---|---|---|---|---|
| GPA | 2392 | 1.90619 | 0.91516 | 4560 | 0 | 4.00000 |
| StudyTimeWeekly | 2392 | 9.77199 | 5.65277 | 23375 | 0.00106 | 19.97809 |
| Absences | 2392 | 14.54139 | 8.46742 | 34783 | 0 | 29.00000 |
| total_score | 2392 | 11.67818 | 5.88610 | 27934 | 0.00803 | 23.42440 |

### Pearson Correlation Coefficients, N = 2392
### Prob > |r| under H0: Rho=0

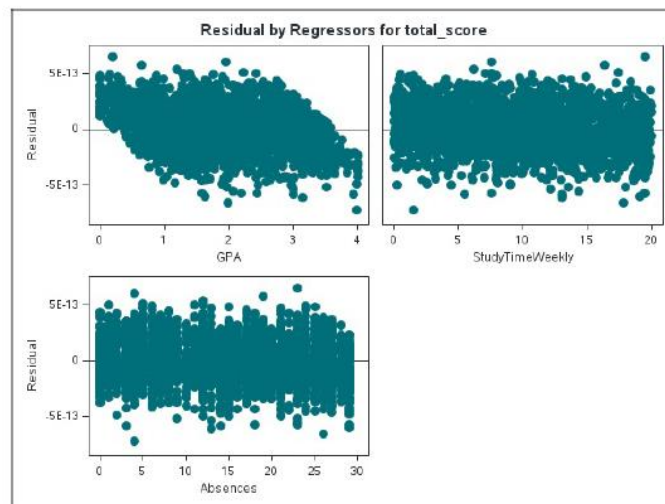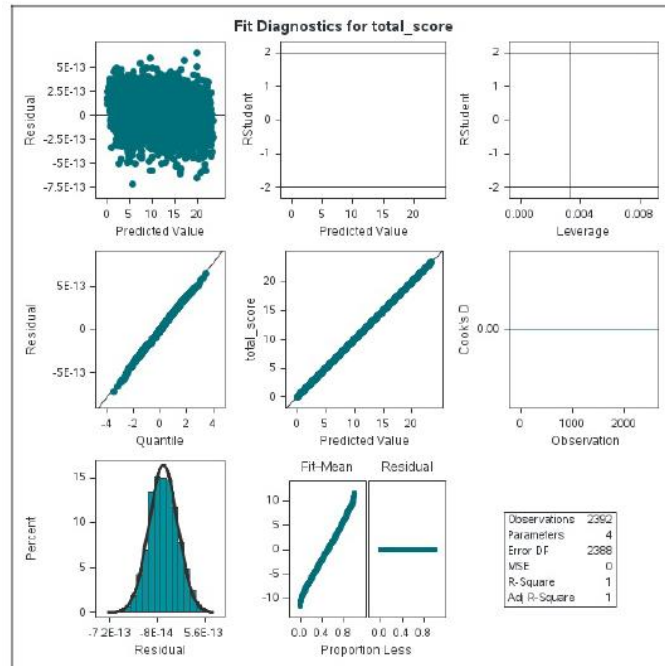| | GPA | StudyTimeWeekly | Absences | total_score |
|---|---|---|---|---|
| GPA | 1.00000 | 0.17928 <.0001 | -0.91931 <.0001 | 0.32765 <.0001 |
| StudyTimeWeekly | 0.17928 <.0001 | 1.00000 | 0.00933 0.6485 | 0.98823 <.0001 |
| Absences | -0.91931 <.0001 | 0.00933 0.6485 | 1.00000 | -0.13398 <.0001 |
| total_score | 0.32765 <.0001 | 0.98823 <.0001 | -0.13398 <.0001 | 1.00000 |

The REG Procedure
Model: MODEL1
Dependent Variable: total_score

| Number of Observations Read | 2392 |
|---|---|
| Number of Observations Used | 2392 |

| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 3 | 82839 | 27613 | Infty | <.0001 |
| Error | 2388 | 0 | 0 | | |
| Corrected Total | 2391 | 82839 | | | |

| Root MSE | 0 | R-Square | 1.0000 |
|---|---|---|---|
| Dependent Mean | 11.67818 | Adj R-Sq | 1.0000 |
| Coeff Var | 0 | | |

| Parameter Estimates | | | | | |
|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > |t| |
| Intercept | 1 | -1.9361E-12 | 0 | -Infty | <.0001 |
| GPA | 1 | 1.00000 | 0 | Infty | <.0001 |
| StudyTimeWeekly | 1 | 1.00000 | 0 | Infty | <.0001 |
| Absences | 1 | 6.25637E-14 | 0 | Infty | <.0001 |

Fit Diagnostics for total_score

| | | |
|---|---|---|
| Observations | 2392 |
| Parameters | 4 |
| Error DF | 2388 |
| MSE | 0 |
| R-Square | 1 |
| Adj R-Square | 1 |



Residual by Regressors for total_score

The FASTCLUS Procedure
Replace=FULL Radius=0 Maxclusters=3 Maxiter=1

**Initial Seeds**

| Cluster | GPA | StudyTimeWeekly | Absences |
|---|---|---|---|
| 1 | 1.99610739 | 0.10648381 | 14.00000000 |
| 2 | 4.00000000 | 19.42439824 | 0.00000000 |
| 3 | 1.07848513 | 19.52141172 | 29.00000000 |

Criterion Based on Final Seeds = 3.7499

**Cluster Summary**

| Cluster | Frequency | RMS Std Deviation | Maximum Distance from Seed to Observation | Radius Exceeded | Nearest Cluster | Distance Between Cluster Centroids |
|---|---|---|---|---|---|---|
| 1 | 858 | 3.8095 | 14.3857 | | 2 | 11.6053 |
| 2 | 761 | 3.5837 | 14.3358 | | 1 | 11.6053 |
| 3 | 773 | 3.6030 | 15.0487 | | 1 | 11.8217 |

**Statistics for Variables**

| Variable | Total STD | Within STD | R-Square | RSQ/(1-RSQ) |
|---|---|---|---|---|
| GPA | 0.91516 | 0.57235 | 0.609187 | 1.558766 |
| StudyTimeWeekly | 5.65277 | 4.31020 | 0.419090 | 0.721437 |
| Absences | 8.46742 | 4.64271 | 0.699615 | 2.329059 |
| OVER-ALL | 5.90165 | 3.67243 | 0.613102 | 1.584661 |

Pseudo F Statistic = 1892.88

Approximate Expected Over-All R-Squared = 0.68709

Cubic Clustering Criterion = -16.256

WARNING: The two values above are invalid for correlated variables.

**Cluster Means**

| Cluster | GPA | StudyTimeWeekly | Absences |
|---|---|---|---|
| 1 | 1.77666460 | 4.87998880 | 14.37762238 |
| 2 | 2.86908628 | 12.53388559 | 5.72273325 |
| 3 | 1.10199839 | 12.48290731 | 23.40491591 |

**Cluster Standard Deviations**

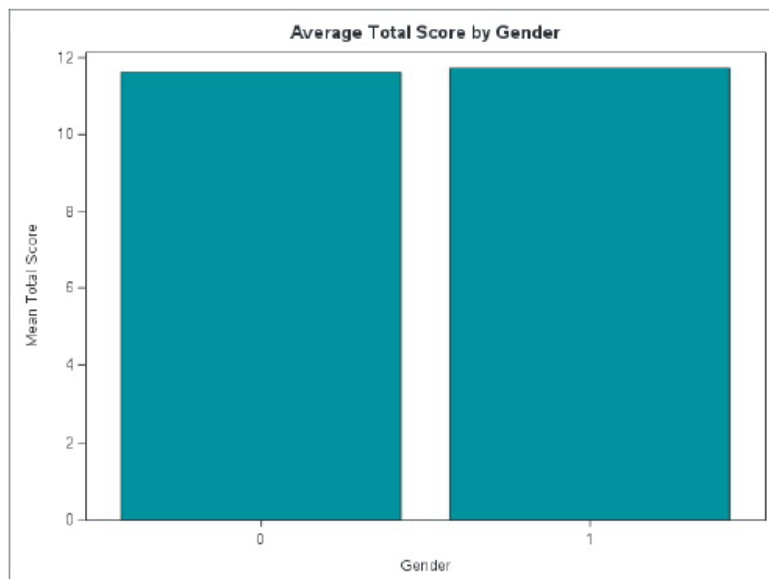| Cluster | GPA | StudyTimeWeekly | Absences |
|---|---|---|---|
| 1 | 0.649131630 | 3.396109588 | 5.619773636 |
| 2 | 0.480526513 | 4.645379604 | 4.088896505 |
| 3 | 0.564487362 | 4.841768356 | 3.896700697 |

| Obs | StudentID | Age | Gender | Ethnicity | ParentalEducation | StudyTimeWeekly | Absences | Tutoring | ParentalSupport | Extracurricular | Sports | Music | Volunteering | GPA | GradeClass | total_score | performance |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 1001 | 17 | 1 | 0 | 2 | 19.833722808 | 7 | 1 | 2 | 0 | 0 | 1 | 0 | 2.9291955917 | 2 | 22.7629 | Low |
| 2 | 1002 | 18 | 0 | 0 | 1 | 15.408756056 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 3.0429148334 | 1 | 18.4517 | Low |
| 3 | 1003 | 15 | 0 | 2 | 3 | 4.2105697688 | 26 | 0 | 2 | 0 | 0 | 0 | 0 | 0.1126022545 | 4 | 4.3232 | Low |
| 4 | 1004 | 17 | 1 | 0 | 3 | 10.028829474 | 14 | 0 | 3 | 1 | 0 | 0 | 0 | 2.0542181397 | 3 | 12.0830 | Low |
| 5 | 1005 | 17 | 1 | 0 | 2 | 4.672495273 | 17 | 1 | 3 | 0 | 0 | 0 | 0 | 1.2880611818 | 4 | 5.9606 | Low |
| 6 | 1006 | 18 | 0 | 0 | 1 | 8.1912185453 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 3.0841836145 | 1 | 11.2754 | Low |
| 7 | 1007 | 15 | 0 | 1 | 1 | 15.601680475 | 10 | 0 | 3 | 0 | 1 | 0 | 0 | 2.7482374149 | 2 | 18.3499 | Low |
| 8 | 1008 | 15 | 1 | 1 | 4 | 15.424496306 | 22 | 1 | 1 | 1 | 0 | 0 | 0 | 1.3601427123 | 4 | 16.7846 | Low |
| 9 | 1009 | 17 | 0 | 0 | 0 | 4.562007558 | 1 | 0 | 2 | 0 | 1 | 0 | 1 | 2.8968191895 | 2 | 7.4588 | Low |
| 10 | 1010 | 16 | 1 | 0 | 1 | 18.444466363 | 0 | 0 | 3 | 1 | 0 | 0 | 0 | 3.5734742103 | 0 | 22.0179 | Low |

# Visualizing Data

### SAS Code for Visualizations

```
/* Step 5.1: Bar Chart of Average Total Score by Gender */

proc sgplot data=clean_data;
     vbar Gender / response=total_score stat=mean;
     title "Average Total Score by Gender";
     yaxis label="Mean Total Score";
     xaxis label="Gender";
run;
```



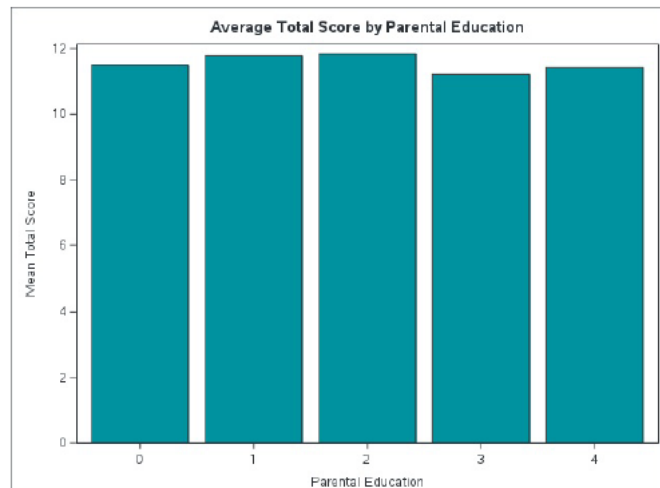### Step 5.1: Bar Chart of Average Total Score by Gender

The bar chart provides a comparison of the average total score (sum of GPA and study time) across genders. This visualization helps identify whether there are significant differences in performance between male and female students.

**Key Insights:**

- The average total scores for male and female students are relatively similar.
- Gender does not appear to significantly influence overall performance based on total score.

```
/* Step 5.2: Bar Chart of Average Total Score by Parental Education */
```

```
proc sgplot data=clean_data;
      vbar ParentalEducation / response=total_score stat=mean;
      title "Average Total Score by Parental Education";
      yaxis label="Mean Total Score";
      xaxis label="Parental Education";
run;
```



### Step 5.2: Bar Chart of Average Total Score by Parental Education

This bar chart illustrates the relationship between parental education levels and students' average total scores. It provides insights into how family education background correlates with student performance.

**Key Insights:**

- Students with parents who have higher education levels tend to perform better on average.
- The chart highlights a positive relationship between parental education and academic success.

```
/* Step 5.3: Scatter Plot of Study Time vs GPA */

proc sgplot data=clean_data;
       scatter x=StudyTimeWeekly y=GPA / group=Gender;
      reg x=StudyTimeWeekly y=GPA;
```

```
        title "Scatter Plot with Regression: Study Time vs GPA";
        xaxis label="Weekly Study Time (hours)";
        yaxis label="GPA";
run;
```



### Step 5.3: Scatter Plot with Regression Line of Study Time vs GPA

The scatter plot shows the relationship between weekly study time and GPA, grouped by gender. A regression line is added to understand the trend and correlation.

**Key Insights:**

- A positive trend is observed, indicating that higher weekly study time is associated with better GPAs.
- Differences between genders are minimal, suggesting that study time influences performance consistently across genders.

## Conclusion

These visualizations provide a comprehensive view of student performance data. They reveal:

1. Gender does not significantly influence total scores.
2. Parental education has a notable positive impact on performance.
3. Weekly study time is positively correlated with GPA, regardless of gender.

These insights can guide educators and policymakers in designing interventions to improve academic success.

# Exporting Results

The cleaned dataset (clean_data) and clustering results (cluster_results) were exported as CSV files to facilitate further analysis and sharing. The cleaned dataset provides ready-to-use, transformed data, while the clustering results highlight grouped patterns based on key variables like GPA and study time. Both exports ensure the project outputs are accessible for external use and further analysis. Exporting the data supports seamless collaboration and decision-making across teams.

***SAS Code to Export***

```
/* Step 6: Export Cleaned Dataset */


proc export data=clean_data
      outfile="/home/sgnanas/EPG294/Cleaned_Student_Data.csv"
       dbms=csv replace;
 run;

/* Export Cluster Results */

proc export data=cluster_results
      outfile="/home/sgnanas/EPG294/Cluster_Results.csv"
      dbms=csv replace;
run;
```

# Generating Summary Report

The final summary report was generated to provide a comprehensive overview of the analysis and findings from the project. This report is designed to present the key steps, methodologies, and insights in a professional format, ensuring clarity and accessibility for stakeholders.

***SAS Code for Report***

```
ods pdf file="/home/sgnanas/EPG294/Student_Performance_Report.pdf";
```

```
 proc report data=clean_data nowd;
      column Gender ParentalEducation GPA StudyTimeWeekly performance;
      define Gender / group 'Gender';
      define ParentalEducation / group 'Parental Education';
      define GPA / mean 'Average GPA';
      define StudyTimeWeekly / mean 'Average Study Time';
      define performance / group 'Performance Category';
      title "Summary of Student Performance by Gender and Parental Education";
run;

ods pdf close;
```

**Summary of Student Performance by Gender and Parental Education**

| Gender | Parental Education | Average GPA | Average Study Time | Performance Category |
|---|---|---|---|---|
| 0 | 0 | 1.916690704 | 9.9711907742 | Low |
| | 1 | 1.9867675743 | 9.6739408885 | Low |
| | 2 | 1.9239672492 | 9.979468995 | Low |
| | 3 | 1.7748872875 | 9.0736926075 | Low |
| | 4 | 1.9000558945 | 9.1033631133 | Low |
| 1 | 0 | 1.8688086638 | 9.2460685459 | Low |
| | 1 | 1.9035623923 | 10.020544562 | Low |
| | 2 | 1.9355713392 | 9.8437494602 | Low |
| | 3 | 1.84024979 | 9.7353173969 | Low |
| | 4 | 1.7315672932 | 10.114346064 | Low |

## Report Design

The report was structured to present the findings in a clear and logical flow:

1. **Introduction**:
    a. A brief overview of the project and the dataset used.
    b. Description of the purpose and objectives of the analysis.
2. **Methods**:
    a. Detailed steps for data exploration, cleaning, and analysis.
    b. Inclusion of relevant SAS code snippets to demonstrate the methodology.
3. **Findings**:
    a. Clear visualizations, charts, and tables summarizing key insights.

  b. Explanation of performance trends across gender and parental education categories.
4. **Conclusion**:
  a. Summary of the main findings and their implications.
  b. Key takeaways from the analysis.
5. **References**:
  a. Details of the dataset source and any academic citations.

## Biography

"I am a Biomedical Data Science and Informatics student at Clemson University. This project reflects my learning and skills in data analysis (one of my courses I took in Fall 2024) using SAS, showcasing a systematic approach to deriving meaningful insights from complex datasets."

## References

- **Dataset**: Student Performance Data (Source: Kaggle).
- **Software**: SAS Studio.