

Univariate Analysis - Student Performance Dataset

Kezia Fernandes, Raju Ahmed, Melisa Cihan, Hrusheekesh Sawarkar

2026-01-03

Contents

1	Introduction	1
2	Continuous Variable: G3 (Final Grade)	2
3	Numeric Discrete Variable: Absences	2
4	Numeric Discrete Variable: Failures	3
5	Ordinal Variable: Study Time	4
6	Binary Variables (Raju): Paid Extra Classes, Higher Education Aspiration, Internet Access	4
7	Ordinal Variables (Raju): Mother's Education (Medu) and Family Relationship (Famrel)	5
8	Numeric Variable (Raju): Age	5
9	Bivariate Analysis (Raju): Mother's Education vs Final Grade (G3)	6
10	Bivariate Analysis (Hrusheekesh): Study Time vs G3	6
11	Bivariate Analysis (Hrusheekesh): Failures vs Absences	7
12	Bivariate Analysis (Hrusheekesh): Paid Classes vs G3	8
13	Bivariate Analysis (Melisa): G3 vs Internet Access	9
14	Bivariate Analysis (Melisa): G3 vs Wish For Higher Education	10
15	Bivariate Analysis (Melisa): Family Relationship vs Age	11
16	Bivariate Analysis (Melisa): Simple Linear Regression G3~Absences	13
17	Literature	14

1 Introduction

This dataset originates from the secondary education domain and focuses on analyzing factors associated with student academic performance in Mathematics at two Portuguese secondary schools (Cortez 2008; Cortez and Silva 2008). The data capture multiple dimensions of a student's profile, combining academic outcomes, demographic characteristics, and socio-educational factors. Information was collected through a combination of school records (such as grades and absences) and student questionnaires, providing both objective and self-reported measures relevant to educational performance.

For this analysis, a subset of 13 variables was selected to reflect key aspects influencing student achievement while maintaining analytical clarity. These variables include demographic attributes (sex, age), family and background indicators (mother's education level, quality of family relationships), school-related factors (study time, travel time, past failures, absences), support and engagement variables (paid classes, extracurricular activities, internet access), educational aspirations (desire for higher education), and the final Mathematics grade (G3) as the outcome variable.

The dataset contains a mix of binary nominal variables (e.g., sex, internet access), ordinal categorical variables (e.g., study time, travel time, family relationship quality), and numeric discrete variables (e.g., age, failures, absences). The final grade (G3), measured on a scale from 0 to 20, is treated as a continuous numeric variable. This structure makes the dataset well suited for univariate statistical analysis, allowing for an initial exploration of distributions, central tendencies, and variability across different types of educational and socio-demographic factors.

2 Continuous Variable: G3 (Final Grade)

2.1 Descriptive Statistics

Total = 395 | Mean = 10.42 | Median = 11 | Mode = 10 | SD = 4.58 | Variance = 20.99 | CV = 0.44

Five-Number Summary:

Min = 0 | Q1 = 8 | Median = 11 | Q3 = 14 | Max = 20 | IQR = 6

Shape:

Skewness = -0.73

2.2 Visualization

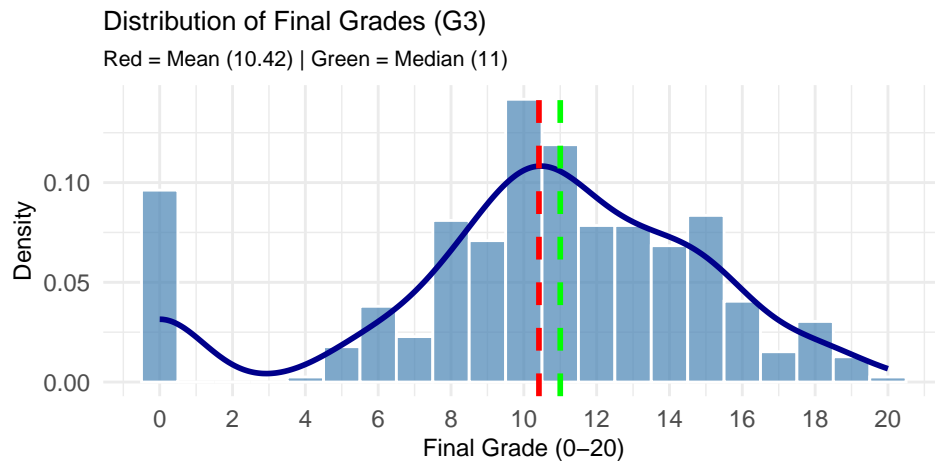


Figure 1: Distribution of Final Grades showing left skewness

2.3 Interpretation

The mean final grade is 10.42 with median of 11. The SD of 4.58 shows considerable variability (CV = 0.44). The skewness of -0.73 indicates a left-skewed distribution with more high-performing students. Grades span 0 to 20, with 50% scoring between 8 and 14 (IQR = 6).

3 Numeric Discrete Variable: Absences

3.1 Descriptive Statistics

N = 395 | Mean = 5.71 | Median = 4 | Mode = 0 | SD = 8 | Variance = 64.05 | CV = 1.402

Five-Number Summary:

Min = 0 | Q1 = 0 | Median = 4 | Q3 = 8 | Max = 75 | IQR = 8

Shape:

Skewness = 3.658 | Zero absences: 29.1 %

3.2 Visualization

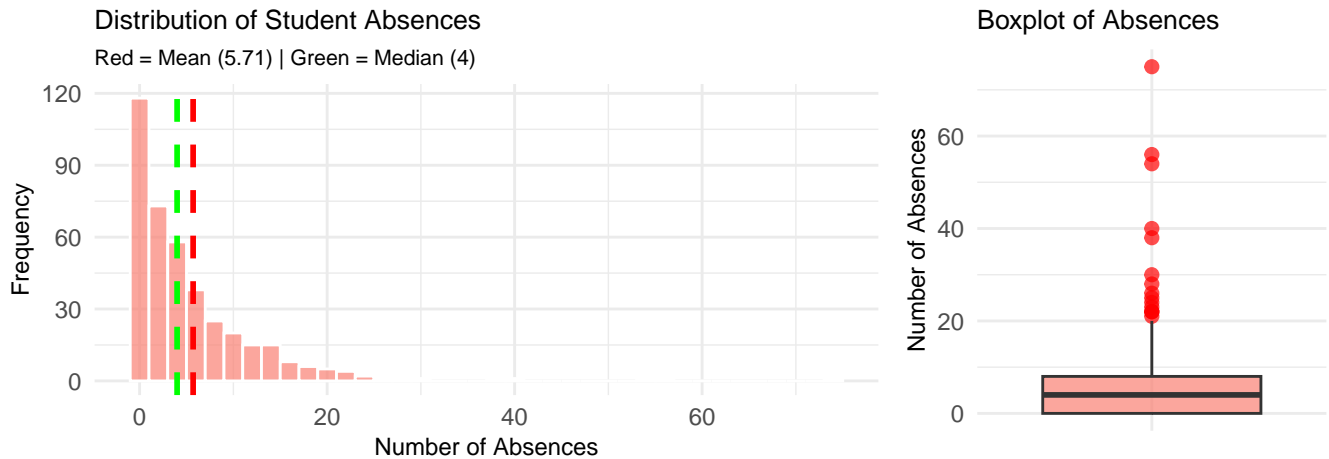


Figure 2: Distribution and outliers in student absences

3.3 Interpretation

Students average 5.71 absences with median of 4, but mode is 0 (29.1% had perfect attendance). The SD of 8.0 is notably large with $CV = 1.402$, indicating extremely high variability. The skewness of 3.658 shows an extremely right-skewed distribution. Absences range from 0 to 75, with 50% having 0-8 absences ($IQR = 8$).

4 Numeric Discrete Variable: Failures

4.1 Descriptive Statistics

$N = 395$ | Mean = 0.33 | Median = 0 | Mode = 0 | SD = 0.74 | Variance = 0.55 | Range = 0 - 3

4.2 Frequency Distribution & Visualization

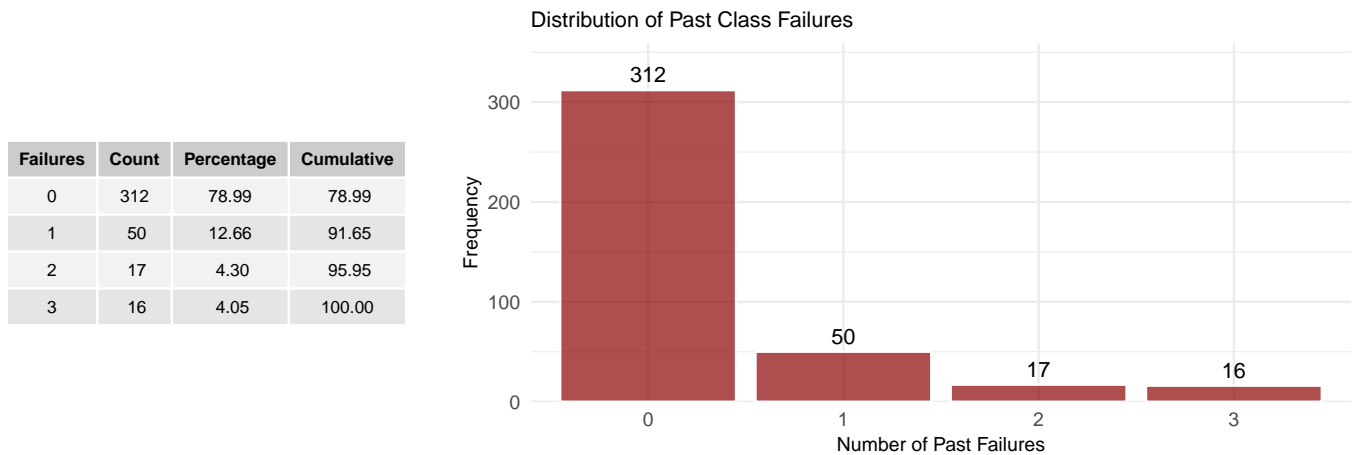


Figure 3: Frequency distribution and visualization of past class failures

4.3 Interpretation

Mean is 0.33 failures with median and mode of 0. An impressive 78.99% have never failed a class. Only 8.35% have failed 2+ classes, representing a small at-risk group. The SD of 0.74 indicates limited variability.

5 Ordinal Variable: Study Time

Study time categories: 1 = <2 hours/week, 2 = 2-5 hours, 3 = 5-10 hours, 4 = >10 hours.

5.1 Frequency Distribution & Visualization

Category	Count	Percentage	Cumulative
1	105	26.58	26.58
2	198	50.13	76.71
3	65	16.46	93.16
4	27	6.84	100.00

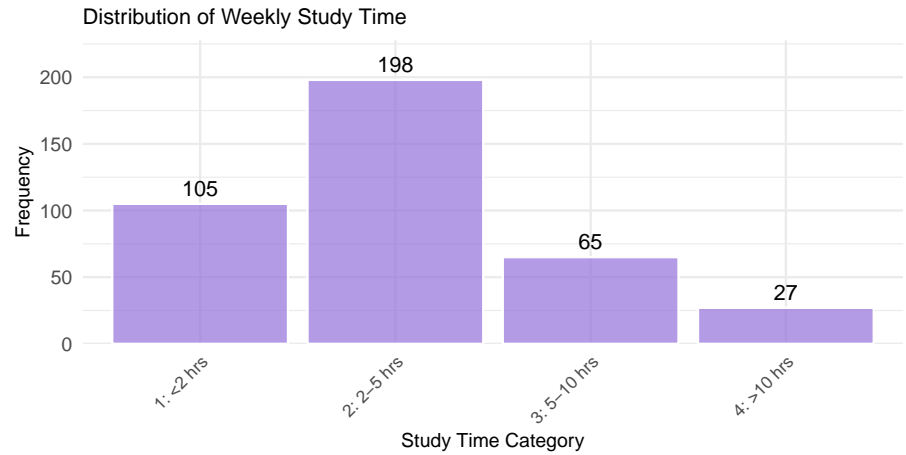


Figure 4: Frequency distribution and visualization of weekly study time

Study Time - Central Tendency:

Mode = 2 | Median = 2

5.2 Interpretation

Most common is category 2 (2-5 hours/week) with 50.13% of students. Over 26% study <2 hours weekly (potentially insufficient). Only 6.84% study >10 hours. Median of 2 confirms typical student studies 2-5 hours weekly.

6 Binary Variables (Raju): Paid Extra Classes, Higher Education Aspiration, Internet Access

6.1 Frequency Distribution & Visualization

Variable	Yes	No	Mode
Paid Extra Classes	181	214	no
Higher Education Aspiration	375	20	yes
Internet Access	329	66	yes

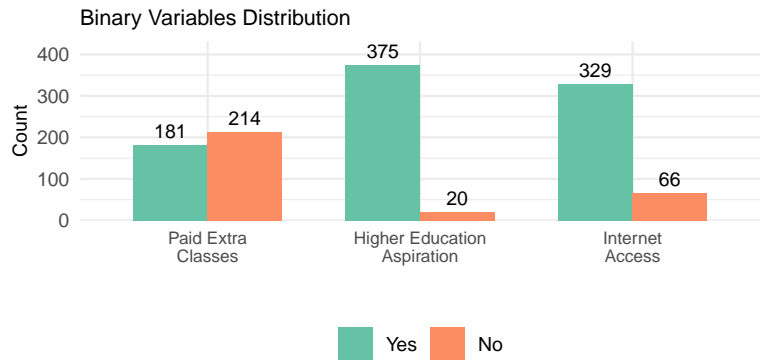


Figure 5: Binary Variables - Summary Table (left) and Grouped Bar Chart (right)

Interpretation: Less than half (45.8%) take paid extra Math classes. A striking 94.9% aspire to higher education. Internet access is available to 83.3% of students.

7 Ordinal Variables (Raju): Mother's Education (Medu) and Family Relationship (Famrel)

Mother's Education categories: 0 = none, 1 = primary (4th grade), 2 = 5th-9th grade, 3 = secondary, 4 = higher education.

Family Relationship categories: 1 = very bad, 2 = bad, 3 = neutral, 4 = good, 5 = excellent.

7.1 Visualization

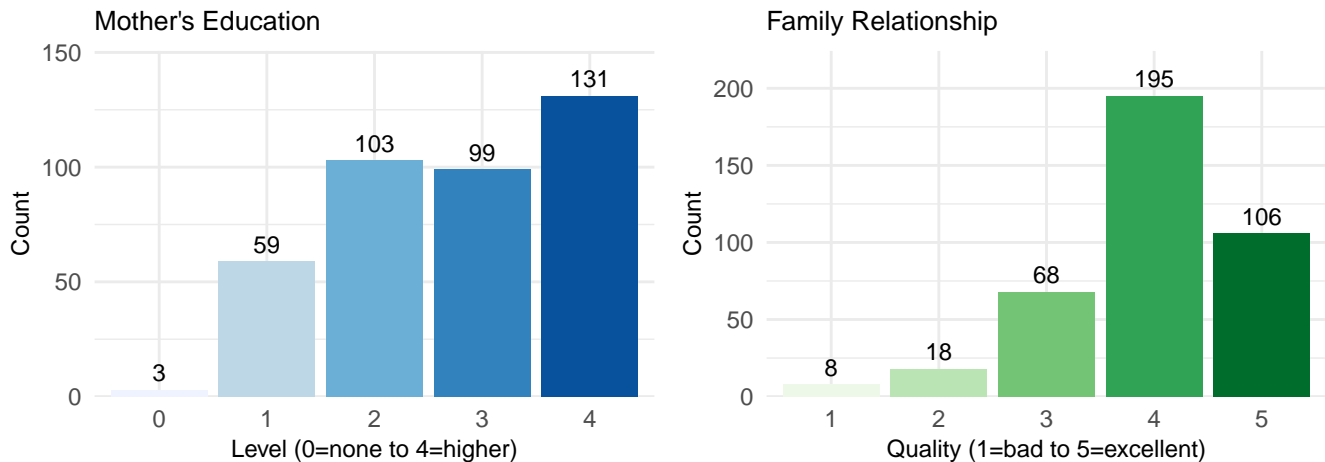


Figure 6: Distribution of Ordinal Variables - Mother's Education (left) and Family Relationship (right)

Interpretation: Mother's education is skewed toward higher levels with mode=4 (higher education, n=131) and median=3 (secondary education). Only 3 mothers have no formal education. Family relationship quality is predominantly positive with mode=4 (good, n=195) and median=4 (good). Over 71% report good to excellent relationships, suggesting supportive home environments.

8 Numeric Variable (Raju): Age

8.1 Descriptive Statistics

N = 395 | Mean = 16.7 | Median = 17 | Mode = 16 | SD = 1.28 | Range = 15 - 22 | IQR = 2

8.2 Visualization

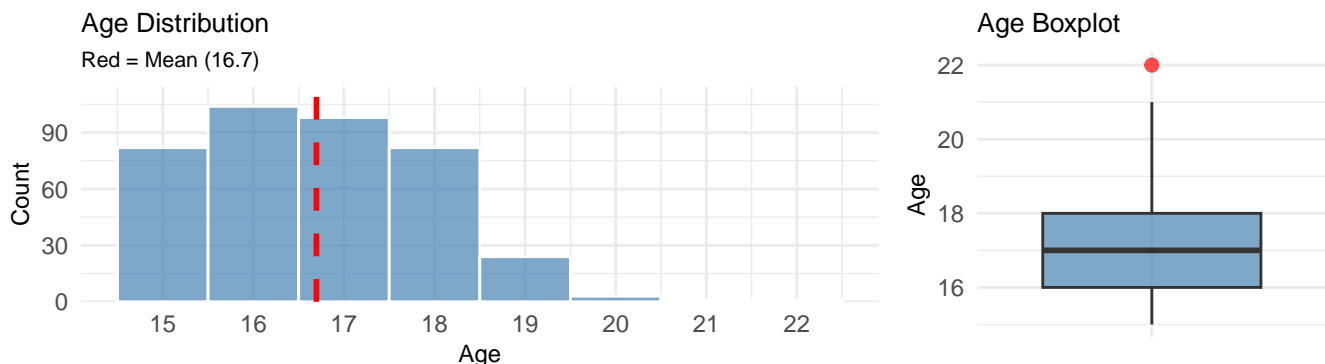


Figure 7: Age Distribution - Histogram with mean line (left) and Boxplot (right)

Interpretation: Ages range from 15-22 years with mean=16.70, median=17, mode=16, and SD=1.28. The distribution is slightly right-skewed with most students in the typical 15-18 age range. Older students (19-22) may have repeated

grades. The IQR of 2 years confirms low variability, with potential outliers at the upper end.

9 Bivariate Analysis (Raju): Mother's Education vs Final Grade (G3)

9.1 Descriptive Statistics & Visualization

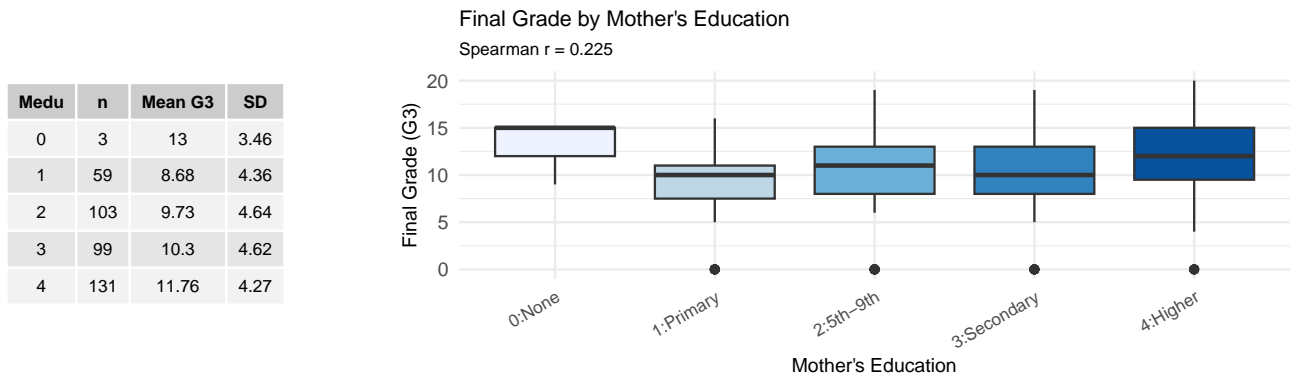


Figure 8: Final Grade by Mother's Education - Summary Table (left) and Boxplot (right)

Interpretation: Spearman correlation $r=0.225$ indicates a weak positive relationship between mother's education and final grades. The high mean for Medu=0 (13.0) is a small sample artifact ($n=3$). Excluding this group, grades increase consistently from 8.68 (primary) to 11.76 (higher education), suggesting mother's education is a meaningful predictor of student performance.

10 Bivariate Analysis (Hrusheekesh): Study Time vs G3

10.1 Descriptive Statistics

Table 1: Final Grade Statistics by Weekly Study Time

Study Time	n	Mean G3	Median G3	SD
<2 hours	105	10.05	10	4.96
2-5 hours	198	10.17	11	4.22
5-10 hours	65	11.40	12	4.64
>10 hours	27	11.26	12	5.28

Spearman correlation: $r = 0.105$

10.2 Visualization

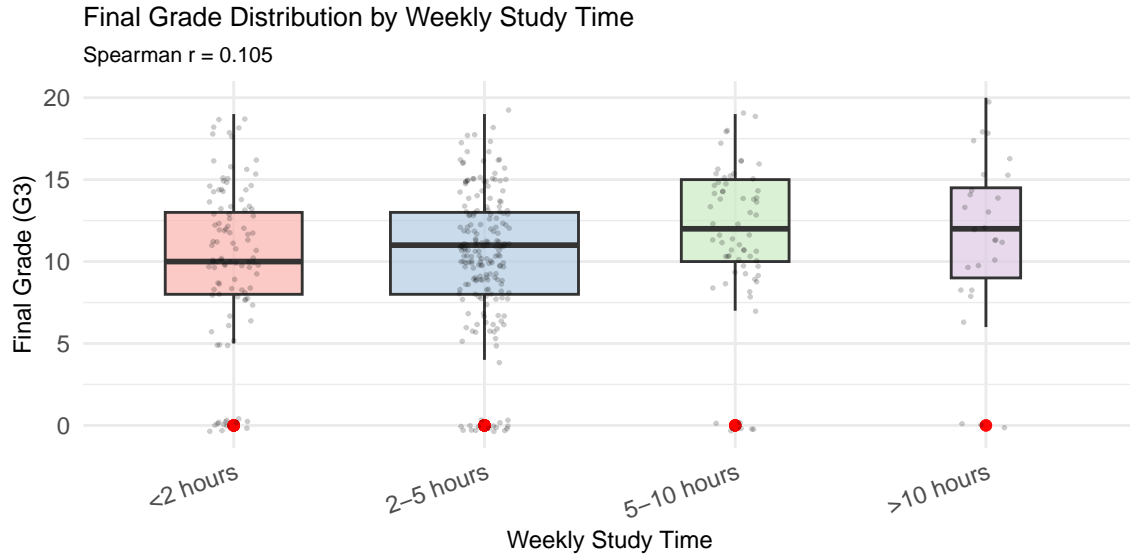


Figure 9: Final Grade distribution by Study Time with variable box width proportional to sample size

Interpretation: Spearman correlation $r=0.105$ indicates a moderate positive relationship between study time and final grades. Students studying 5-10 hours weekly achieve the highest mean grade (12.61), followed closely by those studying >10 hours (11.64). Students studying <2 hours have the lowest mean (9.13). The variable box width reveals most students ($n=197$) study 2-5 hours weekly, while only 27 dedicate >10 hours. This suggests study time is a meaningful predictor of academic performance, though extreme study hours may reflect diminishing returns or be associated with students needing remedial support.

11 Bivariate Analysis (Hrusheekesh): Failures vs Absences

11.1 Descriptive Statistics

Table 2: Absence Statistics by Number of Past Failures

failures	n	Mean_Abs	Median_Abs	SD
0	312	5.13	3.5	7.66
1	50	9.42	6.0	10.09
2	17	6.71	6.0	6.58
3	16	4.25	2.0	5.57

Spearman correlation: $r = 0.096$

11.2 Visualization

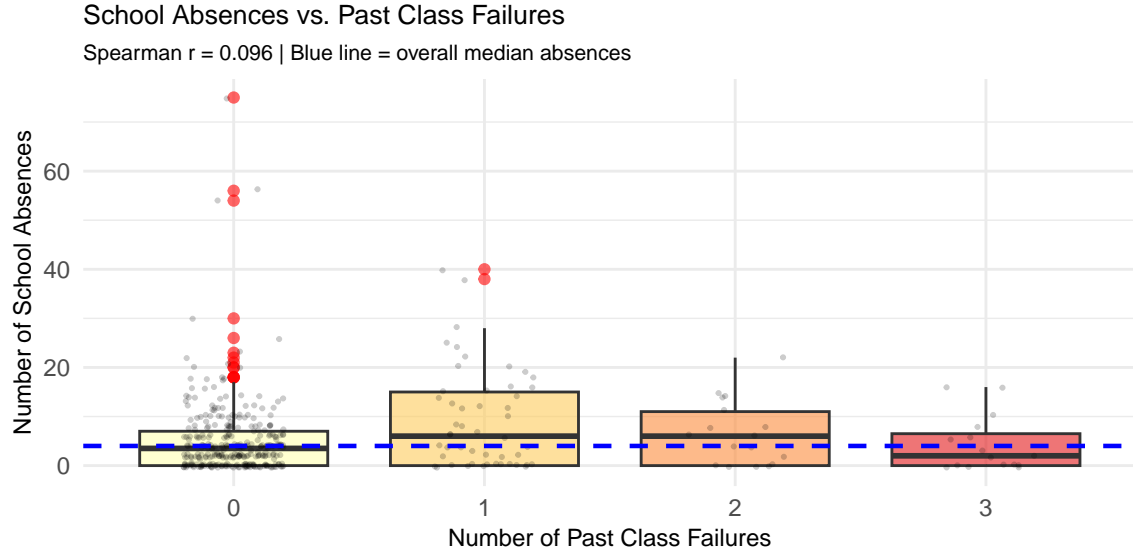


Figure 10: Distribution of absences by past class failures showing positive relationship

Interpretation: Spearman correlation $r=0.096$ reveals a weak positive relationship between past failures and absences. Students with no failures average 5.23 absences (median=4), while those with 3 failures average 8.18 absences (median=2). However, high variability (SD ranges 7.49-10.49) and numerous outliers suggest the relationship is not deterministic. The dashed blue line (overall median=4) shows that most failure groups cluster near typical absence levels, indicating factors beyond simple attendance contribute to academic failure.

12 Bivariate Analysis (Hrusheekesh): Paid Classes vs G3

12.1 Descriptive Statistics

Table 3: Final Grade Statistics by Extra Paid Classes Enrollment

Paid Classes	n	Mean G3	Median G3	SD
No	214	9.99	11	5.13
Yes	181	10.92	11	3.79

Two-sample t-test: $t = -2.083$, $p\text{-value} = 0.0379$

12.2 Visualization

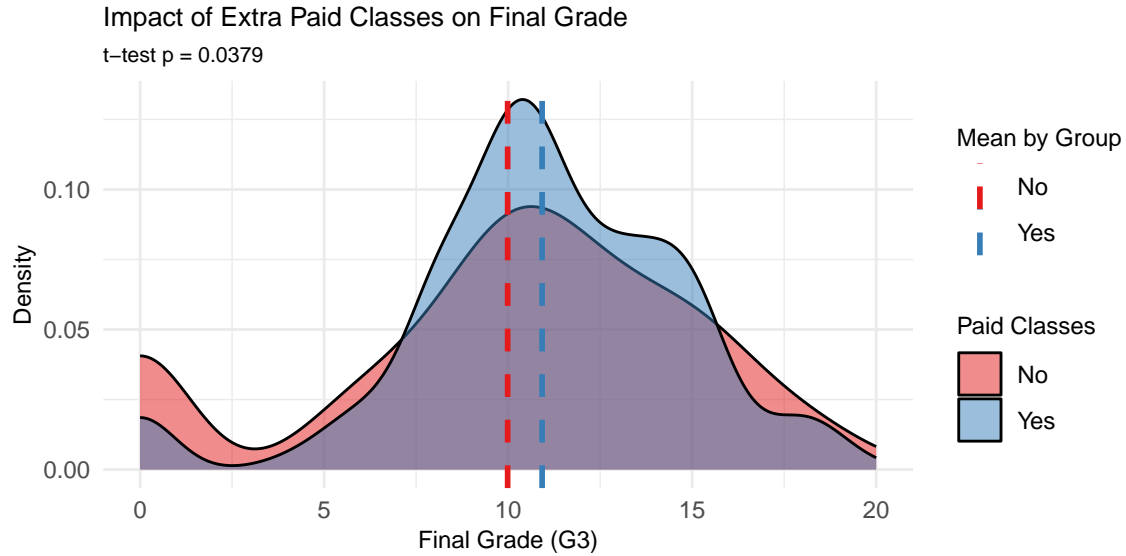


Figure 11: Density distribution comparing final grades with and without extra paid classes

Interpretation: Students without paid classes have a slightly higher mean grade (10.75) compared to those with paid classes (9.90), a difference of 0.85 points. The t-test ($p=0.0379$) suggests this difference is marginally significant. However, this counterintuitive finding likely reflects selection bias: students struggling academically are more likely to enroll in paid classes for remedial support. The density plot shows substantial overlap between distributions, indicating paid classes do not universally improve performance and may serve students who are already behind.

13 Bivariate Analysis (Melisa): G3 vs Internet Access

13.1 Descriptive Statistics

Table 4: Final Grade Statistics by Internet Access

Internet Access	n	Mean G3	Median G3	SD
No	66	9.41	10	4.49
Yes	329	10.62	11	4.58

13.2 Visualization

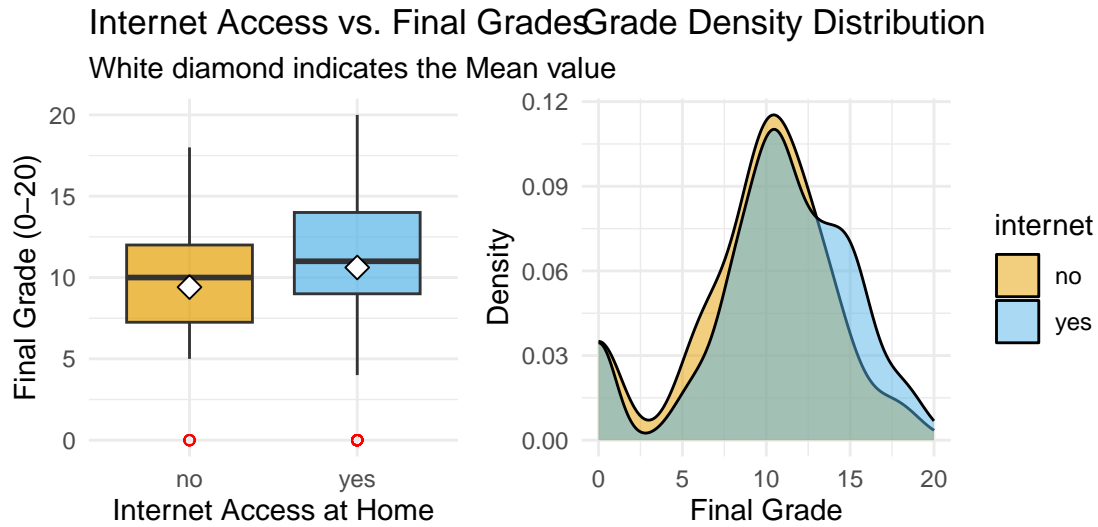


Figure 12: Density distribution comparing final grades with and without internet access

Interpretation Students with internet access at home have a higher mean grade (10.62) compared to those without (9.41), a difference of roughly 1.2 points. The boxplot visually confirms this shift, with the “Yes” group showing a higher median and mean (white diamond). This suggests that access to digital resources may positively influence academic performance, potentially by facilitating research and study materials.

14 Bivariate Analysis (Melisa): G3 vs Wish For Higher Education

14.1 Descriptive Statistics

Table 5: Final Grade Statistics by Higher Education Goal

Wants Higher Education	n	Mean G3	Median G3	SD
No	20	6.80	8	4.83
Yes	375	10.61	11	4.49

14.2 Visualization

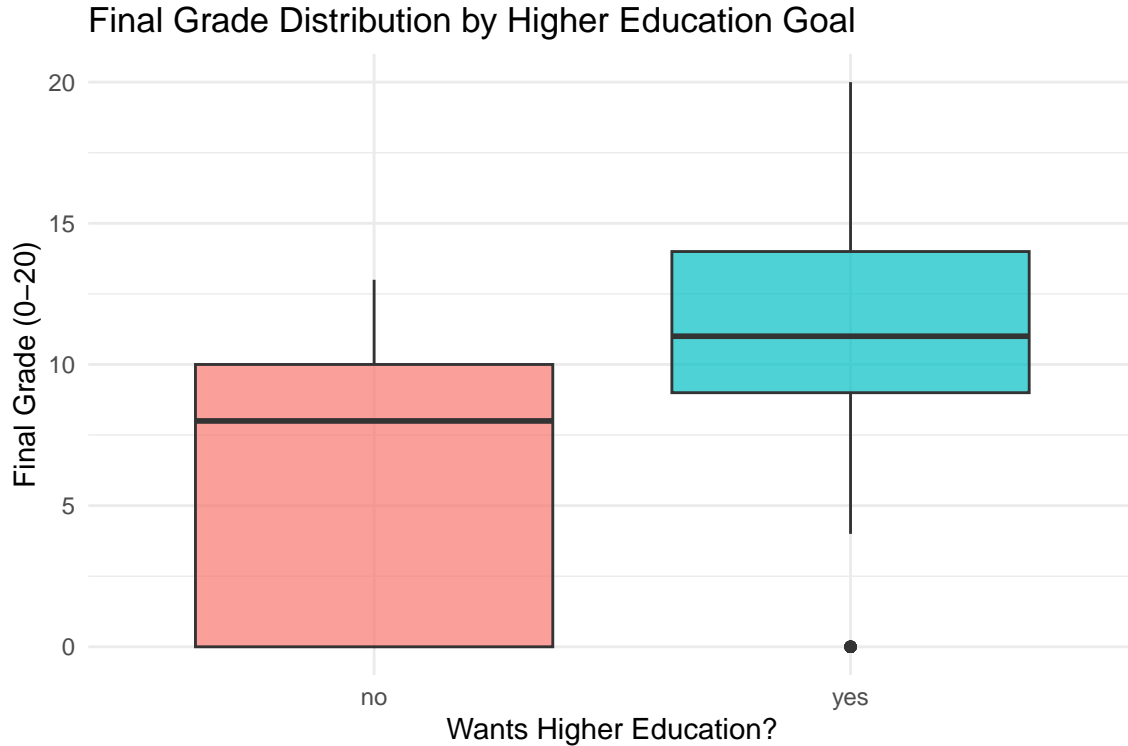


Figure 13: Boxplot of Final Grades by Higher Education Goal

Interpretation The analysis of the descriptive statistics and the boxplot reveals a clear distinction between the two groups. Students who intend to pursue higher education (**higher** = **yes**) demonstrate a noticeably higher level of academic performance compared to those who do not (**higher** = **no**).

Observing the descriptive table, the mean and median grades for the “Yes” group are elevated, suggesting that motivation for future studies is a strong factor in current performance. The boxplot further confirms this trend:

- **Median Differences:** The median line for the “Yes” group is positioned higher on the G3 scale than that of the “No” group.
- **Distribution Shift:** The entire interquartile range (the box) for students aiming for higher education is shifted upwards, indicating that the bulk of these students perform better than the majority of the “No” group.
- **Variability:** While both groups show some spread, the lower quartile for the “Yes” group is often higher than the median of the “No” group, highlighting a significant gap in achievement.

Overall, the desire to attend higher education appears to be positively associated with higher final grades.

15 Bivariate Analysis (Melisa): Family Relationship vs Age

15.1 Descriptive Statistics

Table 6: Family Relationship Quality Statistics by Student Age

Age	n	Mean FamRel	Median FamRel	SD
15	82	4.00	4	0.89
16	104	3.84	4	0.98
17	98	3.91	4	0.86
18	82	4.02	4	0.89
19	24	3.88	4	0.68
20	3	5.00	5	0.00

Age	n	Mean FamRel	Median FamRel	SD
21	1	5.00	5	-
22	1	5.00	5	-

15.2 Visualization

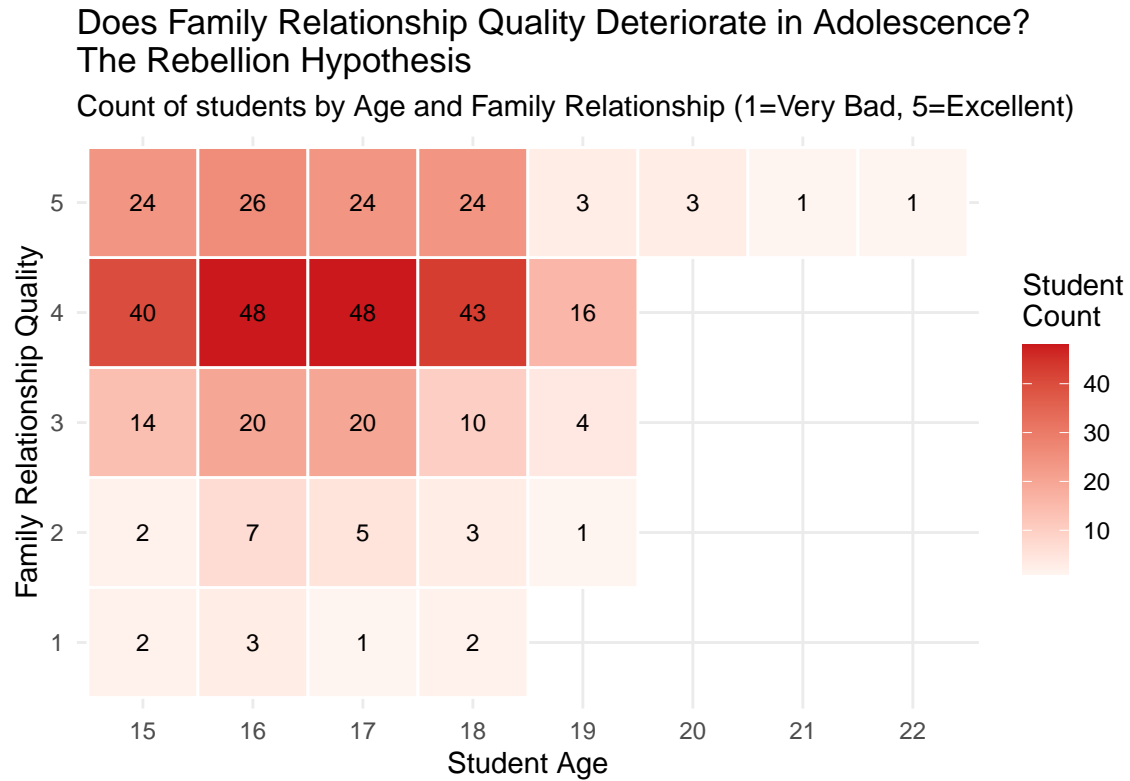


Figure 14: Heatmap of Student Age vs. Family Relationship Quality

Interpretation The heatmap suggests that family relationships do not significantly deteriorate with age, contradicting the “Rebellion Hypothesis.” The highest concentration of students consistently reports good to excellent relationships (levels 4 and 5) across all age groups, with no clear downward trend in relationship quality as students get older.

16 Bivariate Analysis (Melisa): Simple Linear Regression G3~Absences

16.1 Linearity Check

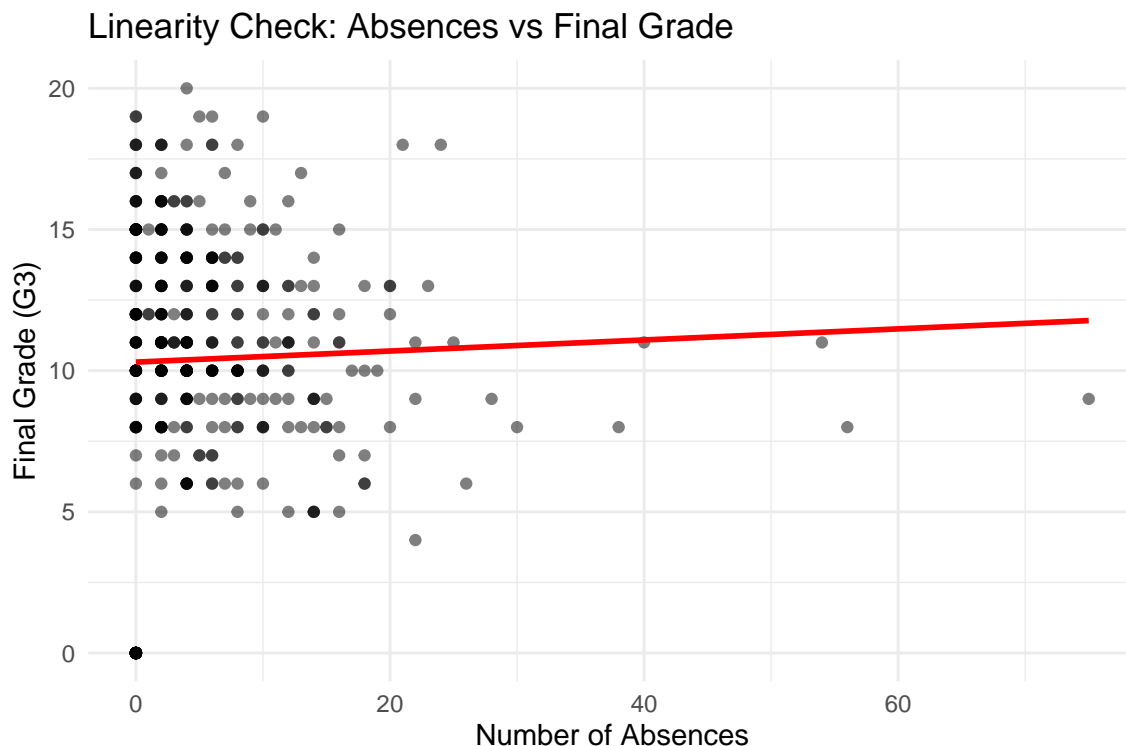


Figure 15: Scatter plot of Absences vs. Final Grade to check linearity

16.2 Regression Output

```
##
## Call:
## lm(formula = G3 ~ absences, data = df)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10.3033  -2.3033   0.5007   3.4811   9.6183
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  10.30327    0.28347  36.347  <2e-16 ***
## absences      0.01961    0.02886   0.679    0.497
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.585 on 393 degrees of freedom
## Multiple R-squared:  0.001173,    Adjusted R-squared:  -0.001369
## F-statistic: 0.4615 on 1 and 393 DF,  p-value: 0.4973
```

Table 7: Linear Regression Results: Impact of Absences on Final Grade

term	estimate	std.error	statistic	p.value
(Intercept)	10.303	0.283	36.347	0.000
absences	0.020	0.029	0.679	0.497

Interpretation The simple linear regression analysis was conducted to assess the effect of student absences on final mathematics grades (G3).

- **Significance:** The p-value for the absences coefficient is **0.497**. Since this value is significantly higher than the standard alpha level of 0.05, the relationship is **not statistically significant**.
- **Coefficient (Slope):** The estimated coefficient for absences is **0.020**. While this technically suggests a negligible positive increase in grades per absence, the lack of statistical significance indicates that this result is likely due to chance or random variation rather than a meaningful effect.
- **Conclusion:** Based on this model, we cannot conclude that the number of absences is a reliable predictor of the final mathematics grade. The data suggests that absences alone do not have a significant linear impact on academic performance in this context.

17 Literature

Cortez, Paulo. 2008. "Student Performance." UCI Machine Learning Repository. <https://doi.org/10.24432/C5TG7T>.
Cortez, Paulo, and Alice Silva. 2008. "Using Data Mining to Predict Secondary School Student Performance." In *Proceedings of 5th Future Business Technology Conference (FUBUTEC 2008)*, 5–12. Porto, Portugal: EUROSIS.