# Students Perfomance in written tests
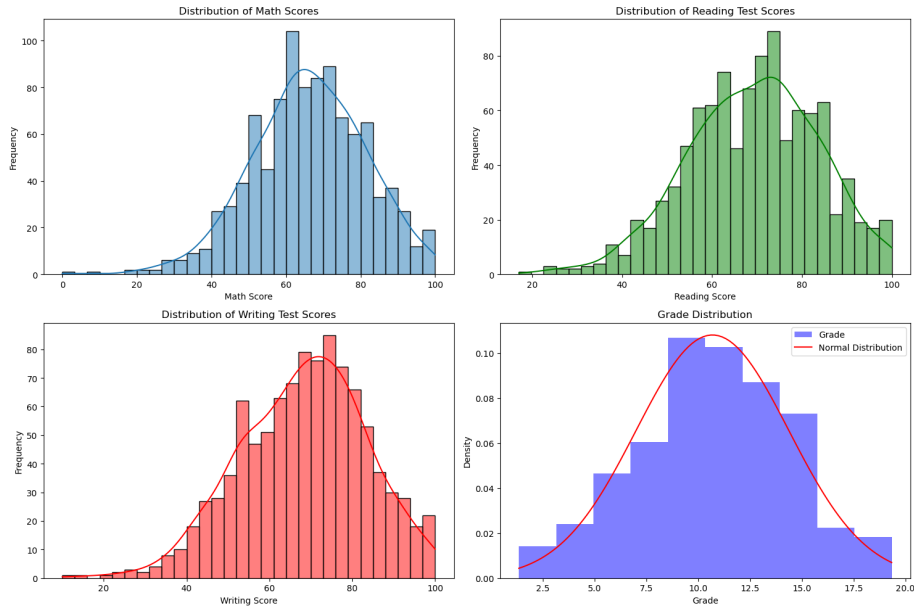
Mykyta Panchenko, Alen Sahinpasic

February 2024

## 1 Introduction

This part of the project advances our exploration of the Student's performance in tests dataset, which was introduced in our earlier work. Moreover, we will introduce a new dataset in order to make the regression more credible. Here, we focus on the correlation between a different variables that describe student's status and their performance. By analyzing this relationship, our objective is to ultimately develop a predictive model capable of estimating a student's performance based on some of the variables we will introduce later into the report.
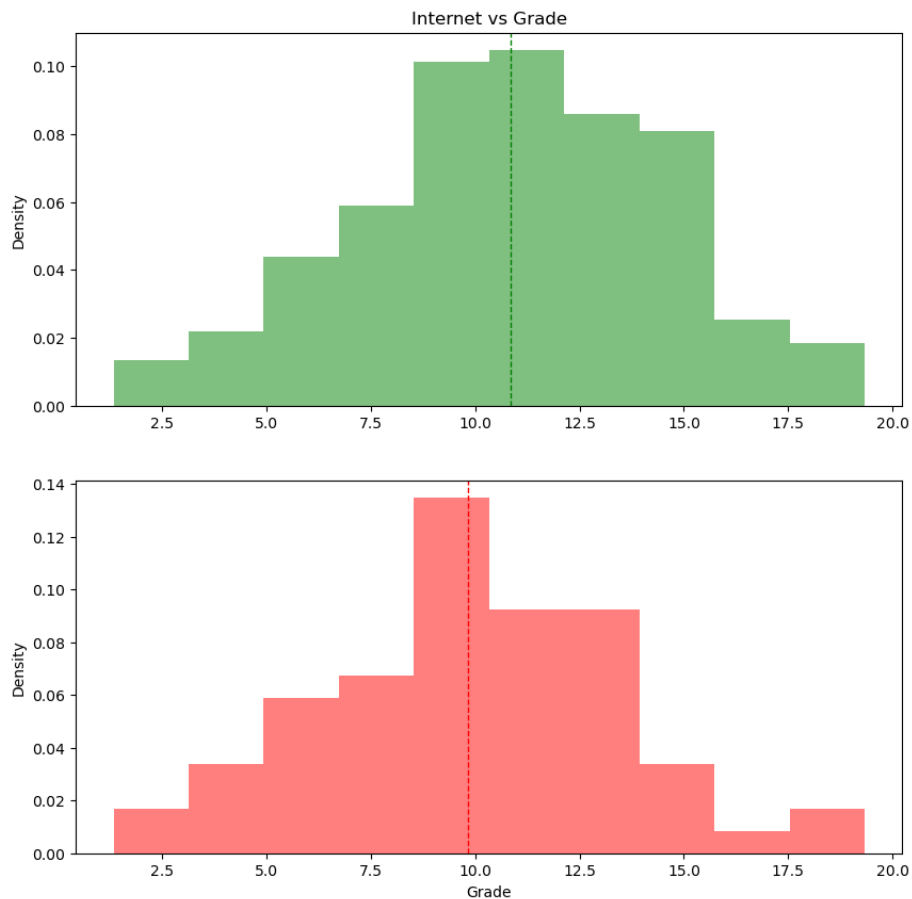
## 2 Distribution and Fitting

In the last report, we assumed that the distribution of scores for each test is approximately normal. In this part, we fitted normal distributions and performed Shapiro-Wilk tests. For the first 3 distributions which are a part of the first dataset, we got a 0.99 statistic with p≈0. On the other hand, with a second dataset, which is not artificial, we got p=0.0505, which means we can't reject the null hypothesis at 5 percent significance level.

# 3 Hypothesis testing

Hypothesis one: Performance of people who have access to the internet is the same as the performance of those who don't on 95 percent confidence level.
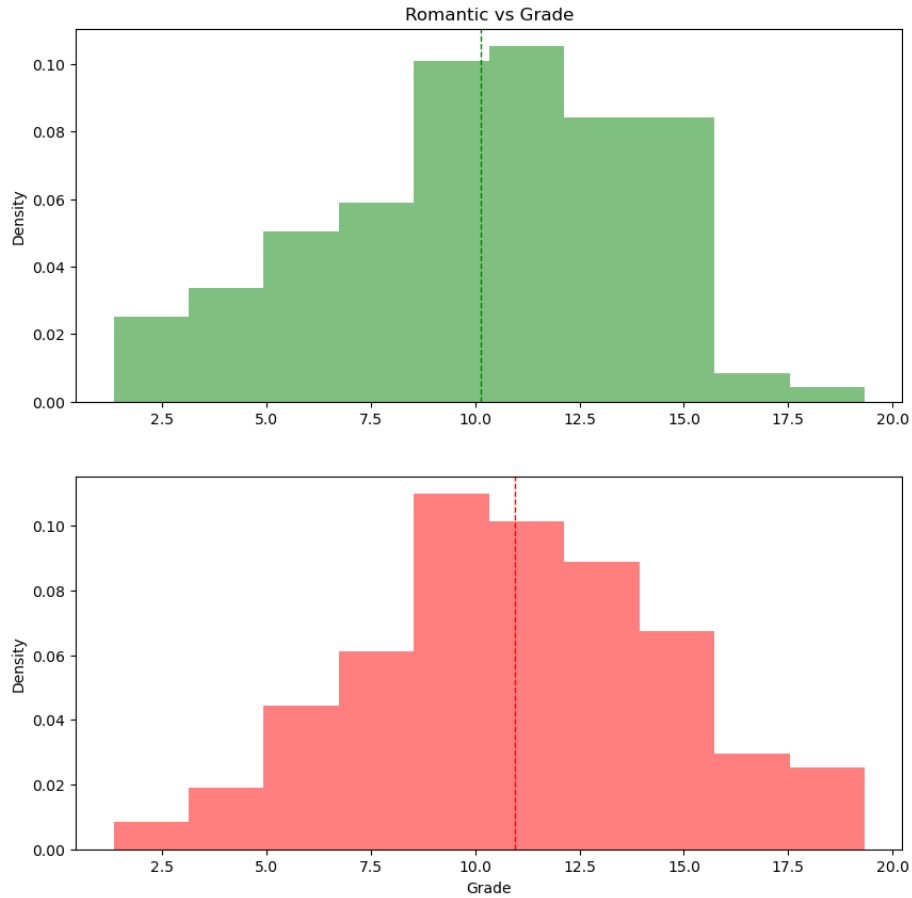


```python
stats.ttest_ind(df2[df2.internet == 'yes'].Grade, df2[df2.internet == 'no'].Grade, equal_var=False)
```
```
TtestResult(statistic=2.1027435949194566, pvalue=0.03811743647022728, df=95.52440308821478)
```

The P-value is 0.04 with equal variance and 0.038 with Welch's t-test, therefore we can reject the null hypothesis on a 5 percent confidence level and say that the performance is not the same.

Hypothesis two: Performance of people who are in a romantic relationship is the same as the performance of those who aren't on 95 percent confidence level.
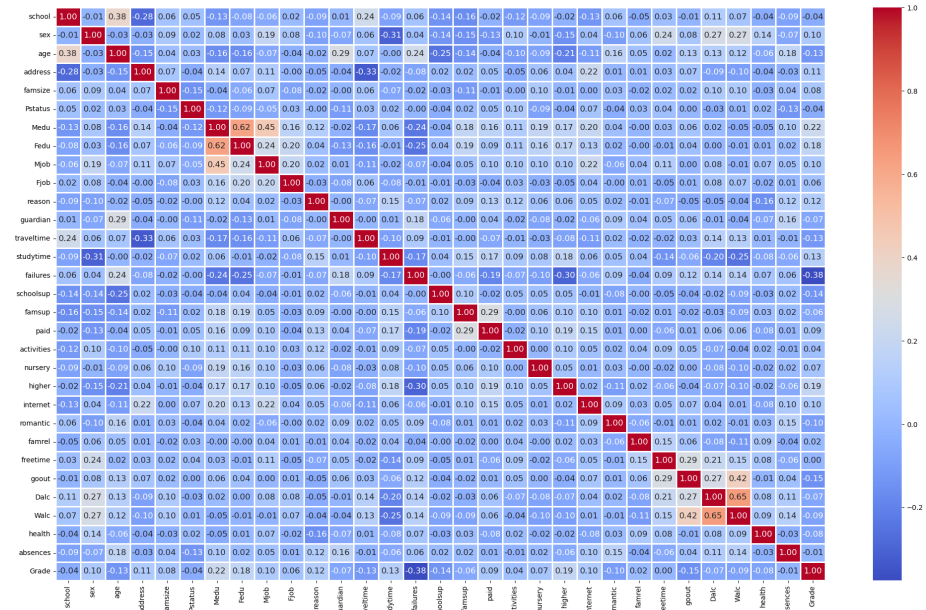


Romantic vs Grade

```
stats.ttest_ind(df2[df2.romantic == 'yes'].Grade, df2[df2.romantic == 'no'].Grade, equal_var = False)

TtestResult(statistic=-2.0439114822017004, pvalue=0.041965111150106205, df=261.249873225931)
```

The P-value is 0.042 both with the Student's and Welch's t-test, therefore we can reject the null hypothesis on a 5 percent confidence level and say that the performance is not the same.
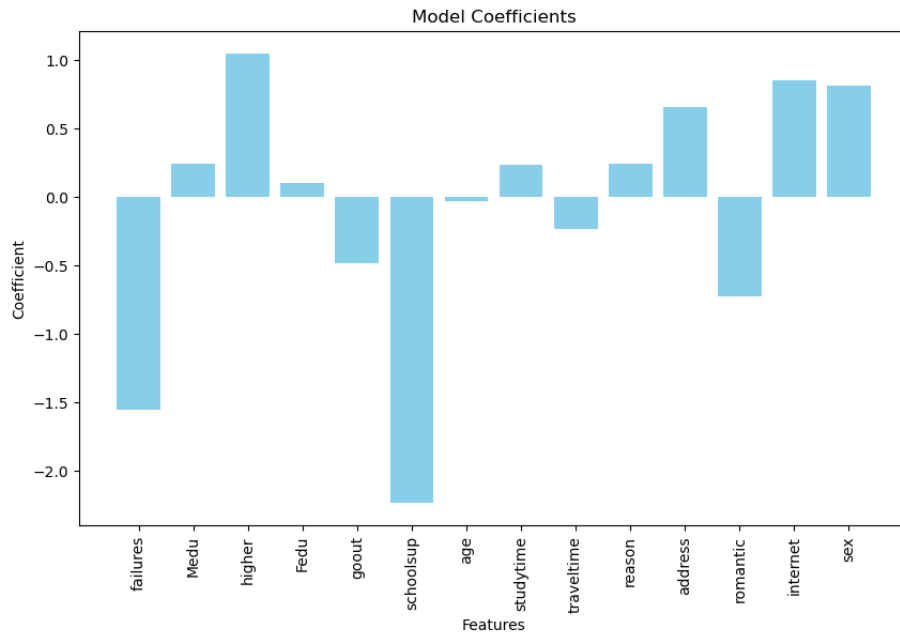
# 4 Linear Regression

First of all, we plotted a correlation matrix to see which variables are mostly correlated to the Grade. We selected features that had the correlation coefficient of $> 0.1$.



```python
target_column_name = 'Grade'
correlations = df2.corr()[target_column_name]
sorted_features = correlations.abs().sort_values(ascending=False)
correlation_threshold = 0.1
selected_features = sorted_features[sorted_features >= correlation_threshold].index.tolist()
print(selected_features[1:])
```

```
['failures', 'Medu', 'higher', 'Fedu', 'goout', 'schoolsup', 'age', 'studytime', 'traveltime', 'reason', 'address', 'romantic', 'internet', 'sex']
```

We proceeded to build a linear regression model by dividing the dataset into train and test sets and obtained the coefficients of the regression model with an Intercept of 9.885774365061723 Then, we proceeded to test each feature if it's



Model Coefficients

significant or not and build a model.We got an $R^2$ of 23.44 with Mean Absolute error of 2.74.

We would also like to present a comparison between predicted and actual values in the dataset. As we can see, $R^2$ of 23.44 clearly displays the perfomance of the model.

| | Actual value | Predicted value |
|---|---|---|
| 343 | 5.666667 | 10.277512 |
| 97 | 9.000000 | 8.724014 |
| 183 | 8.666667 | 11.684757 |
| 288 | 14.333333 | 12.505546 |
| 315 | 11.666667 | 9.887602 |

# 5    Conclusions

In conclusion, our exploration of the dataset has provided valuable insights. Through a combination of distribution fitting, statistical tests, and regression analysis, we've established a nuanced understanding of how variables interact within the confines of the dataset.

According to the linear regression, we have found out that the best performance would show someone who has parents with higher education, no failed classes in the past, who is aiming to continue studying in the university after school, who lives less than 15 minutes to school, has access to the internet and studies more than 10 hours a week. Variables that contributed the most were: Number of classes student failed in the past, Whether or not student receives extra educational support, and if the student wants to resume their studies.