# Analysis of the Influence of the SocioEconomic Environment in the Results of Mathematics in the Colombian University Admission Exams

**Jorge Arias Martí & Andrés Mejía Rodríguez**

**Master in Statistics for Data Science**

**October 2021**

# Introduction

Various studies have shown the correlation between socioeconomic status and results in several academic tests around the world. As an illustration, Dixon et. al [1] found that in the United States, a lower SES (Socioeconomic Status) is related to a lower entry rate at universities. Fernández and Martínez [2] found that in Spain between the years 1970-2012 "students from higher socio-economic backgrounds show fewer fluctuations in school failure than students from lower socio-economic background."

The main objective of this project is to evaluate how the socioeconomic conditions of the students affect their academic performance. To do so, we will analyze the grades obtained in the Colombian set of exams *Saber 11*, an exam that Colombian students must face in order to graduate at high school and be admitted in a higher education institution .In this exam, students are tested in Mathematics, Science, Reading, English and Education for Citizenship [3].It is expected then, that students in a well-off situation would obtain a higher score than their peers.

We will analyze data provided by the ICFES (*Instituto Colombiano para la Evaluación de la Educación),* the institution which oversees the exam. The original dataset is *SB11_20202.csv* and contains the recorded results of all the 504872 students who took the test in the second 2020 call. This call is taken mostly by students whose schools use the "A" calendar (school year starting in January) who are the majority of the schools in the country. From this dataset we have extracted the following variables:

ESTU_CONSECUTIVO: It's simply a public ID of the student [4].Not useful for analysis.

FAMI_ESTRATOVIVIENDA: Stratum of the dwelling. Ranging from 1 to 6 and based on the conditions of the students' living place, Colombians who are in a similar socioeconomic condition, are assigned by the same number. The higher this number is, the better their socioeconomic status is. It is a Categorical numerical ordinal and discrete variable [4, 5]. Figure 1 shows the distribution of the students by the stratum to which they belong.

PUNT_INGLES: Score obtained in the English exam. It is a numerical quantitative discrete variable [4]. Its range is from 0 to 100.

PUNT_MATEMATICAS: Score obtained in the Mathematics exam. It is a numerical quantitative discrete variable [4]. Its range is from 0 to 100.

PUNT_LECTURA-CRITICA: Score obtained in the Critical Reading exam. Numerical quantitative discrete variable [4]. Its range s from 0 to 100.

The results of this work can be used to compare the yield of recent exam takers with future candidates' yield as long as the scales and methodology of the exam do not change. We also expect this work to reflect in the future the closure of gaps between the academic performance of students belonging to different socioeconomic backgrounds.

It would be also interesting to make statistical inference to examinees from other nationalities that also face one exam in order to access to higher education. Ecuador is a candidate for this, as it has a similar GPD per capita (US $5600.4) to that of Colombia (US $5332.8) [6].
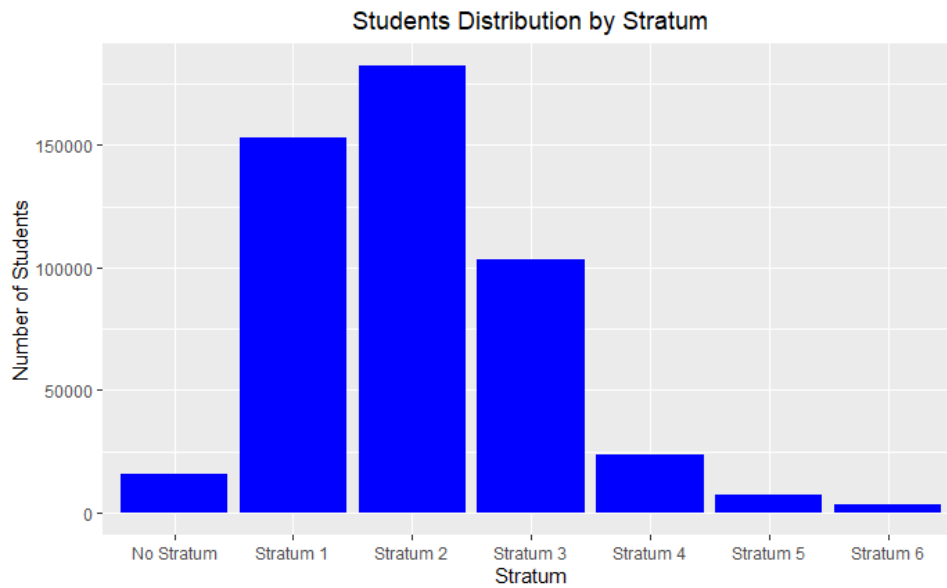
**Figure 1**. It shows the distribution of the examinees by their stratum. Note that the distribution is quite asymmetric and that it is quite centered to the left. We can see that there is a higher number of students in the lower stratum than in the upper ones.

## Model Selection

As an exploratory analysis we will analyze the distribution of the students by their score in each test and their summary statistics. These plots are shown in figure 2.
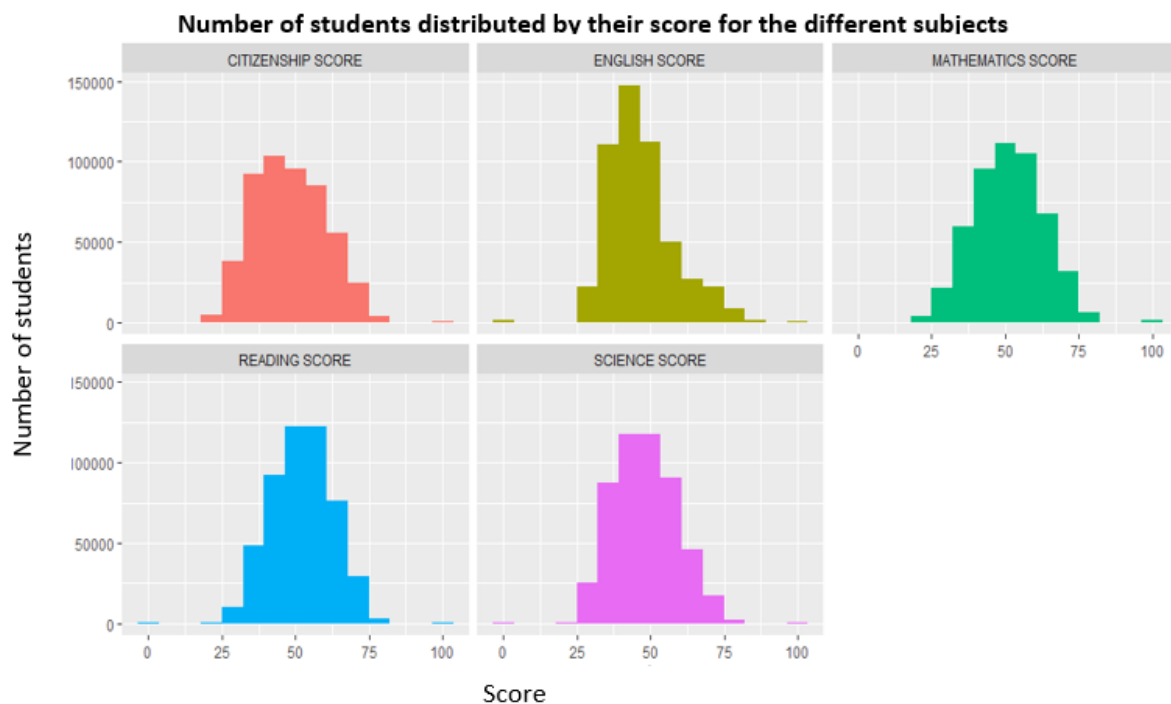


**Figure 2**. It shows the distribution of the examinees by their score in the different exams. At first sight, it seems that it is quite possible that Mathematics, Reading and Science could be approximated by a Normal distribution.

As we see in figure 2, Mathematics, Reading and Science have a more symmetrical and centered distribution than English and Citizenship. Table 1 shows some statistics of all these distributions:

| Task | MEAN | STANDARD DEVIATION | KURTOSIS | SKEWNESS |
|---|---|---|---|---|
| SCIENCE SCORE | 48.197 | 10.500 | 3.130 | 0.257 |
| ENGLISH SCORE | 46.884 | 11.313 | 4.855 | 0.862 |
| READING SCORE | 52.157 | 10.158 | 3.042 | 0.012 |
| MATHEMATICS SCORE | 51.020 | 11.648 | 3.264 | 0.146 |
| CITIZENSHIP SCORE | 48.234 | 11.971 | 2.906 | 0.309 |

**Table 1.** Some statistics about the distribution of the students by exam. The kurtosis includes the excess. As it can be observed, all subjects except English have a Kurtosis close to 3, and a skewness close to zero. Apparently, they are good candidates to be approached by a Normal. The mean is the sample mean, taking our whole population as a sample.

Normal distribution has a kurtosis of 3 and skewness of 0 [7], and if we compare these two parameters with the ones in Table 1, our distributions seem to have a high degree of normality. To check whether this is true or not, we show the qq-plots, which are shown below in figure 3:
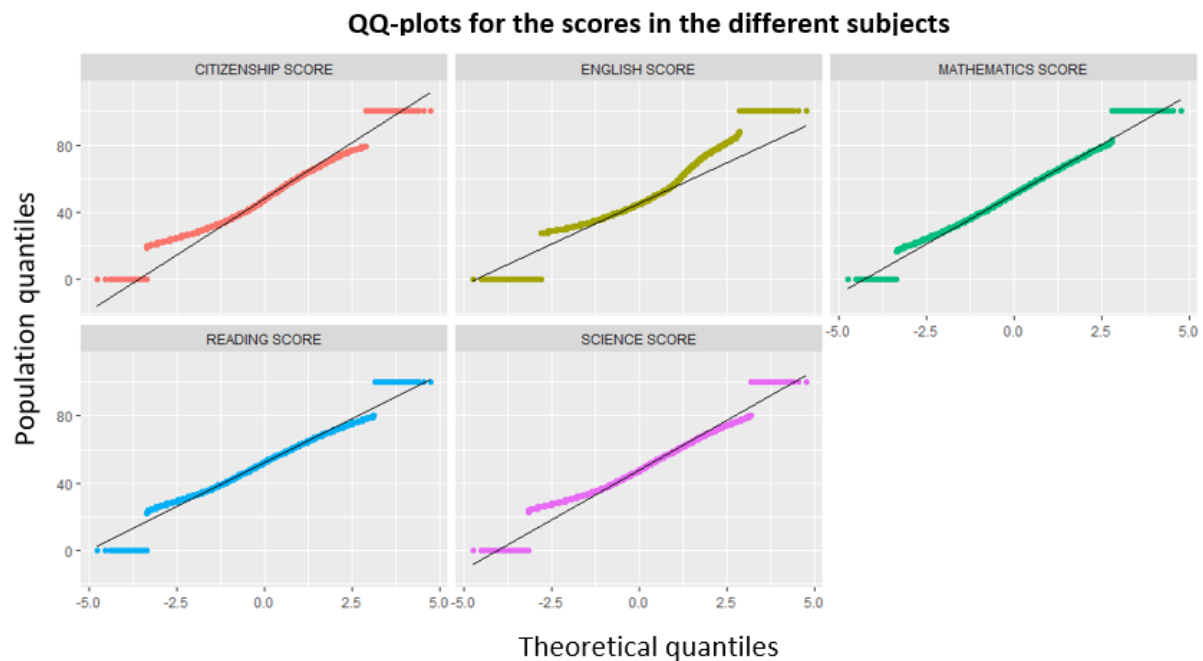


**Figure 3.** It Shows the qq-plots for the five exams. As it can be seen, Mathematics and Reading are the ones who are closer to the Normal theorical quantiles. Flat tails are due to the reduced number of students located at the extremes in figure 2.

Figure 3 shows that Mathematics and Reading are the subjects whose quantiles are closer to the normal ones. That is the reason because we have considered that they are more treatable than the others, so, from now on, we will work only with the Mathematics score, leaving aside the others.

Having used the normal distribution as a model for the results of the Mathematics test we would like to estimate now both its mean and its variance. We will do this with the method of moments. By definition, the k-moment in the discrete case has the following expression [9]:

$$a_k = \frac{\sum_{i=1}^n x_i^k}{n}$$

So, applying the definition $a_1 = E[X] = \frac{\sum_{i=1}^n x_i}{n}$ which means that $a_1$ is the sample mean. On the other hand, $a_2 = E[X^2]$

For the continuous case the definition is the following one:

$$a_k = \int_{-\infty}^{\infty} x f(x)$$

where $f(x)$ is the probability density. Replacing $f(x)$ by the normal distribution formula, its two first moments have the next analytical expressions:

$$a_1 = \mu$$

$$a_2 = \mu^2 + \sigma^2$$

Finally, we get the next equation system:

$$\mu = \overline{x_n}$$

$$\mu^2 + \sigma^2 = \frac{\sum_{i=1}^n x_i^2}{n}$$

The solution of the system is

$$\hat{\mu} = \overline{X_n} \quad \text{and}$$

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n x_i^2}{n} - \hat{\mu}^2 = \frac{\sum_{i=1}^n x_i^2}{n} - \overline{X_n}^2 = S_n^2$$

Where $S_n^2$ denotes the sample variance. We can extract two conclusions from here: for normal distributions, the sample mean is a good estimator of the population mean, and that the sample variance is a good estimator of the population variance.

Having solved the system, and replacing values, we obtain that $\hat{\mu} = 51.02$ and $\hat{\sigma}^2 = 135.67$.


## One-sample inference

We will consider two estimators for the population mean in the Mathematics test. The first one that will be used is the sample mean, which is an unbiased and consistent estimator of the mean. However, as it uses all available information, several factors like the possibility of leaving in blank some answers or a possible fraud committed by some examinees, make influence in the value of the mean. We should also take into account that the sample mean is quite sensitive to outliers, as all scores have the same contribution, so it is not a robust estimator. Its formula is given by the next expression:

$$\overline{x_n} = \frac{\sum_{i=1}^{n} x_i}{n}$$

Considering this we will also use a robust estimator, in this case the trimmed mean. The trimmed mean consists of removing a given percentage of the highest and lowest observations and calculating the mean in those that remain. The trimmed mean is not as sensitive to outliers as the sample mean, but we instead lose efficiency as we don't use all the available information. This estimator also depends on how symmetric the underlying distribution is in order to guarantee the unbiasedness, but in figure 2 we see that the Mathematics score distribution seems to be quite symmetric, so this assumption is supposed to be correct.

The expression of the trimmed mean is the following one:

$$T_\alpha = \frac{1}{[(1-\alpha)n]} \sum_{i=[n\alpha/2]}^{[(1-\alpha/2)n]} x_i$$

Where $\alpha$ is half of the proportion of data not taken into account. In our case, we take $\alpha = 0.1$, so the 10% of the data are not included (the highest 5% and the lowest 5%).

The following table includes a summary of the estimators in the Mathematics test of the sample mean and the trimmed mean and the coefficient of variation in each case.

| | Sample mean | Trimmed Mean |
|---|---|---|
| Value | 51.02 | 50.93 |
| Standard Error | 0.0163 | 0.0170 |
| Coefficient of Variation | 0.000321 | 0.000336 |

Table 2. Value, Standard error and coefficient of variation for the Sample mean and for the Trimmed mean.

We see in table 2 that both means values are close one to each other. In this table it is included also the Standard Error and the Coefficient of variation.

For the sample mean the standard error is given by:

$$\sigma_{\bar{x}} = \frac{S}{\sqrt{n}}$$

While for the trimmed mean, we need first to define the Windsorized sum of squared deviations, which has the next expression:

$$se^2 = [n\alpha/2](X_{n\alpha} - \bar{X})^2 + \sum_{i=[n\alpha/2]}^{[n(1-\alpha/2)]} (X_{n\alpha} - \bar{X})^2 + [n\alpha/2](X_{[n(1-\alpha/2)]} - \bar{X})^2$$

And the standard error is:

$$\sigma_{\bar{x}} = \frac{se_{wk}}{\sqrt{(n - 2n\alpha)(n - 2n\alpha - 1)}}$$

Finally, the CV is defined as

$$CV = \frac{\sigma_{\bar{x}}}{\hat{\mu}}$$

Next step is to find the confidence interval for the population mean, and to do so, we will use the sample mean (not the trimmed). Since we have a sample with $n > 500000$ ,we can use the asymptotic approach for the mean confidence interval.

Under this approximation, the confidence interval for the population mean has the following expression:

$$CI_{1-\alpha}(\mu) = \overline{X_n} \pm Z_{\alpha/2}\frac{S'}{\sqrt{n}}$$

Where S' is the sample quasivariance (although our sample covers the whole population) which value is $S' = 11.648$ , $Z_{\alpha/2}$ is the pivotal quantity of the normal distribution for $\alpha/2$ and $n = 504872$. In addition, finding a 95% confidence interval means that $\alpha = 1 - 0.95 = 0.05$ , so $Z_{\alpha/2} = 1.96$  Replacing these values in the equation we obtain the following values:

$$CI_{95}(\mu) = [50.968, 51.052]$$

As we see, the CI has a length of 0.084, so it is quite a narrow interval for a high confidence level as it is the 95% level.

We will also find the proportion of students who had a score higher than 80 in the Mathematics exam. To achieve this goal, we will create the binary variable *Greatt80*, whose value is 1 if the student had a score greater or equal than 80 and zero otherwise.

The proportion estimator will be the sample mean of this variable, because if we sum all the $x_i$,the obtained result is the number of students with a score higher than 80, and if we divide by $n$ (the total examinees number), the proportion of the sample will be obtained. This proportion is a consistent estimator of the population proportion, by algebra of consistency

The results obtained can be seen in the following table.

| #Total examinees at the mathematics task | # Students with a score higher than 80 | $\widehat{p}$ |
|---|---|---|
| 504872 | 2399 | 0.00475 |

**Table 3.** It shows the number of total examinees, and the number of examinees with a score equal or higher to 80. p̂ is the proportion estimator.

In table 3 we can observe that the proportion of students with a score higher than 80 in the Mathematics exam is really low, only 0.475 % of them had this or a superior mark.

Now, as our variable *Greatt80* is binary and take as a value 1 in case of success (reaching a score of 80 or higher) and zero otherwise, it follows a Bernoulli distribution, so the sum of this

variable will follow a Binomial distribution. p̂ is defined as the division between this sum and n,so as a consequence, p follows a Binomial distribution [9]. The next expression shows It:

$$Greatt80_i \sim \text{Bern(p)}; \sum_{i=1}^{n} Greatt80_i \sim \text{Binomial(n, p)}$$

The variance for a binomial distribution is given by $S^2 = p(1-p)$. As $\hat{p}$ is a consistent estimator for $p$ [9], we can try estimating the population variance taking $\hat{S}^2 = \hat{p}(1-\hat{p})$. Then, the estimator for the population variance is $\hat{S}^2 = 0.00473$.

We can now try whether the condition of the normal approximation is true or not. If we take the product $np$ we obtain that $n\hat{p} = 2398 > 5$, which is one of the conditions. The other one is $n\hat{p}(1-\hat{p}) > 5$. In our case $n\hat{p}(1-\hat{p}) = 2386$, so the normal approximation is correct [9].

In addition, with $n$ large, $\hat{p}$ satisfies that [9]:

$$\frac{(\hat{p}-p)}{\sqrt{\hat{p}(1-\hat{p})/n}} \xrightarrow{d} N(0,1)$$

So, the confidence interval has the following expression:

$$CI_{1-\alpha}(p) = \hat{p} \pm Z_{\alpha/2}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Again, our objective is to find a confidence interval with $\alpha = 1 - 0.95 = 0.05$. p̂ has the value shown in table 3, and $n$ is the total number of students, which is also shown in the table. Replacing all these values in the formula, we have the following confidence interval:

$$CI_{95}(\text{p}) = [0.00456, 0.00494]$$

As we see, the interval has a length of $3.8 \cdot 10^{-4}$, so we have obtained a quite narrow CI for the population proportion for a level of confidence of 95%.


## Inference with more than one sample

We will now explore how the stratum relates to the results of the Mathematics test. We will first visualize the scores distribution in every stratum using a boxplot, shown in figure 4

Looking at the graph it seems that the mean score and its variability increases for stratum 1 to 4. In the table 4 we summarize the results of the analysis of the means using the same methods used in one-sample inference.
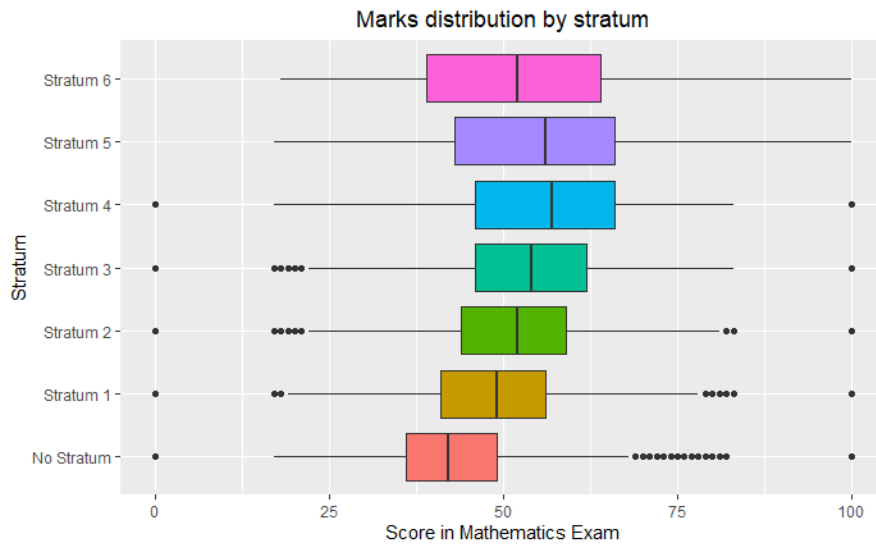
**Figure 4.** Boxplot for the score in the mathematics exam. We see that for stratum 1 to stratum 4, the mean increases.

| Stratum | Maths Score Average | # of Students | Standard Deviation | Standard Error | Coefficient of Variation |
|---|---|---|---|---|---|
| Stratum 1 | 48.747 | 152852 | 10.820 | 0.028 | 0.003 |
| Stratum 2 | 51.503 | 182322 | 11.083 | 0.026 | 0.003 |
| Stratum 3 | 53.786 | 103550 | 11.589 | 0.036 | 0.004 |
| Stratum 4 | 55.794 | 23463 | 13.467 | 0.088 | 0.005 |
| Stratum 5 | 54.823 | 7019 | 14.805 | 0.177 | 0.008 |
| Stratum 6 | 51.546 | 3083 | 15.335 | 0.276 | 0.010 |

**Table 4.** It shows the scores mean in the mathematics exam, as well as the number of exam takers by stratum, the Standard Deviation, the Standard Error, and the CV. The mean used hear is the normal mean ,i.e. $\frac{\sum x_i}{n}$ .Last three parameters has been obtained by the same procedure applied in one sample inference.

Note that the mean score for stratum 6 seems to be very close to the mean score for stratum 2. From Figure 4, it seems that the variance is heteroscedastic, as it increases as the strata also increases. We will check this hypothesis doing a pairwise test of the equality of variances. The test used will be a F-test for variance given by:

$$\frac{S_x}{S_y} = F_{n_1-1, n_2-1}$$

Where $F$ is the Fisher-Snedecor's distribution with $n_1 - 1$ and $n_2 - 1$ degrees of freedom ($n_1$ and $n_2$ are the population in the Stratum x and Stratum y, respectively)

To do this we will use the *var.test* function in R. This results in rejecting the null hypothesis of the variances being equal in all cases with p-values smaller than $2.2 \cdot 10^{-16}$ except in one case. The p-value of the test comparing strata 5 and 6 is 0.0206 that rejects the null hypothesis for $\alpha = 0.05$.

The null hypothesis of the means being equal or not will be tested with the Welsh t-test, which is ideal to compare means in groups with different sizes and variances [12], as it is our case.

The variable is defined as:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_\Delta} ; s_\Delta = \sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

Which follows a t distribution with degrees of freedom v [12].

$$v = \frac{(\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y})^2}{\frac{(s_x/n_x)^2}{n_x - 1} + \frac{(s_y/n_y)^2}{n_y - 1}}$$

This test is implemented in R using the *t.test* function. Doing these pairwise comparisons between strata rejects all hypothesis of the means being equal with p-values lower than $2.2 \cdot 10^{-16}$, with one exception. When comparing strata 2 and 6 we get a p-value of 0.87, so we cannot reject the null hypothesis that the mean score is equal.

We can also find a confidence interval for the difference of means using this distribution, the formula of the CI becomes the following one:

$$CI(\mu_X - \mu_Y) = \bar{X} - \bar{Y} \pm t_{\alpha/2}\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$$

Replacing in the formula we obtain the next CI:

$$CI_{95}(\mu_6 - \mu_2) = [-0.587, 0.500]$$

Note that this interval contains zero, which is equivalent to affirm that there is nothing about the null hypothesis of the means being equal.

For each stratum we will find the proportion of students who obtained a score greater or equal to 80 points in every stratum. To do so we will use the variable *greatt80* defined before. The estimator of proportion is the mean of this variable. The obtained results are shown in the following table:

| Stratum | # Students in the Stratum | # Students with a mark great or equal to 80 | $\widehat{p}$ | MSE |
|---|---|---|---|---|
| Stratum 1 | 152852 | 361 | 0.00236 | 0.00187 |
| Stratum 2 | 182322 | 700 | 0.00383 | 0.00297 |
| Stratum 3 | 103550 | 691 | 0.00661 | 0.00509 |
| Stratum 4 | 23463 | 397 | 0.0169 | 0.0130 |
| Stratum 5 | 7019 | 142 | 0.0201 | 0.0146 |
| Stratum 6 | 3083 | 44 | 0.0143 | 0.0119 |

**Table 5.** It Shows the scores mean in the mathematics exam, as well as the number of exam takers by stratum. The mean used here is the normal mean, i.e., $\frac{\sum x_i}{n}$ .We can also see the population proportion estimator for each stratum,$\hat{p}$ and the Mean Squared error.

The MSE values that appear in table 5 are the different variances for each $\hat{p}_i$. As $\hat{p}$ is an unbiased estimator, the MSE is equal to the variance, because the bias is zero.

Besides, in the table we can observe that for the stratum between 1 and 5, $\hat{p}$ increases its value as the strata is higher while stratum 6 has a lower $\hat{p}$ even than Stratum 4. If this case is ignored, a higher stratum seems to mean a higher proportion of students with a score higher than 80.

We will compare the proportion of students who obtained a score greater than 80 in Mathematics between strata 2 and 6. Previously it has been mentioned that in general, $\hat{p}$ has a higher value for high stratum. To verify if this statement is essentially true or is due to chance, firstly we will find the confidence interval for the difference between $\hat{p}_2$ and $\hat{p}_6$ (whose values are, respectively 0.00383 and 0.0143). The confidence interval for the difference between two proportions is given by the next expression:

$$CI_{1-\alpha}(\hat{p}_6 - \hat{p}_2) = (\hat{p}_6 - \hat{p}_2) \pm Z_{\alpha/2}\sqrt{\frac{\hat{p}_6(1-\hat{p}_2)}{n_6} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

where $n_2$ and $n_5$ are the number of students of students in the stratum 6 and 2, and $Z_{\alpha/2}$ is the pivotal quantity for the normal. Again, α=0.05 because we want to find the CI again for a confidence level of 95%. Applying this formula gives as the following interval:

$$CI_{95}(\hat{p}_6 - \hat{p}_2) = [0.0062, 0.0146]$$

As the confidence interval does not contain the zero, it is contrary to the fact that the proportions are equal.

Previously it has been mentioned the fact that $\hat{p}$ follows a binomial distribution that can be approximated to a normal one. We check that $n_6\hat{p}_6 = 44 > 5$ and $n_6\hat{p}_6(1-\hat{p}_6) = 43 > 5$. On the other hand $n_2\hat{p}_2(1-\hat{p}_2) = 697 > 5$, so normal approximation is valid for $\hat{p}_2$ and $\hat{p}_6$.

We can now define our test statistic:

$$Z = \frac{\hat{p}_6 - \hat{p}_2}{\sqrt{\frac{\hat{p}_6(1-\hat{p}_6)}{n_6} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$$

Our null hypothesis H$_o$ is the equality of proportions; Besides, we fix an α value of 0.05. The value of Z is $Z = 4.88$. Furthermore $Z_{\alpha/2} = 1.64$, so Z >$Z_{\alpha/2}$, which means that the null hypothesis of the proportions being equal should be rejected. This is consistent with the confidence interval finding, which did not contain the zero value.

## Conclusions

In first place, the distribution of the scores in Mathematics have been studied considering the normal approximation. Although the statistical analysis carried out supports this statement, the approximation may be not exact.

Besides, although apparently it seems to be some kind of correlation between the obtained average score and the stratum, comparing the means between the stratum 2 and 6 by confidence intervals and hypothesis testing, the conclusion at we arrive is that the difference between both is not statistically significant to be able to affirm that they are different.

On the other hand, the proportion of students with higher score seems to be larger for high stratum than for the low ones. This statement is supported not only by the confidence interval of the difference between them, but also by the hypothesis testing, so both results are again coherent between them.

However, we should consider the fact that the analyzed data is about the students who took the exam in the second 2020 call. This can have led us to a possible bias. As an illustration, more than half of the students in stratum 6 took their exam in the first call (3083 vs 2134). Further analysis can be made including students in both callings in a year.

In addition, it would have been interesting to analyze the score obtained in other subjects such English or Citizenship and comparing if SES have a stronger influence in the score obtained in this subject than in Mathematics.

# References

[1] Dixon-Roman, Ezekiel & Everson, Howard & Mcardle, John. (2013). *Race, Poverty and SAT Scores: Modeling the Influences of Family Income on Black and White High School Students' SAT Performance*. Teachers College Record. 115.

[2] María Fernández-Mellizo & José Saturnino Martínez-García (2017) *Inequality of educational opportunities: School failure trends in Spain (1977–2012)*, International Studies in Sociology of Education, 26:3, 267-287, DOI: 10.1080/09620214.2016.1192954

[3] Ministerio de Educación Nacional - República de Colombia. (2021). *Pruebas saber*. Retrieved from https://www.mineducacion.gov.co/1621/w3-article-244735.html

[4]ICFES. *Diccionario de variables periodo 20191 a 20202*. Retrieved from https://icfesgovcomy.sharepoint.com/personal/dataicfes_icfes_gov_co/_layouts/15/onedrive.aspx?ct=1589296771489&or=OWA%2DNT&cid=5cc96871%2D447f%2D0e87%2D9de5%2D3893e123b5ba&originalPath=aHR0cHM6Ly9pY2Zlc2dvdmNvLW15LnNoYXJlcG9pbnQuY29tLzpmOi9nL3BlcnNvbmFsL2RhdGFpY2Zlc19pY2Zlc19nb3ZfY28vRWtMWWVpLWRxRGxGdVJiOWhsZjIOElCaFpIbXdraFJKdFFkFkzd05OanN0TkNvQT9ydGltZT1SWllqbFqWUN0OTVJVZw&id=%2Fpersonal%2Fdataicfes%5Ficfes%5Fgov%5Fco%2FDocuments%2FDataIcfes%2F4%2E%20Saber11%2F2%2E%20Documentaci%C3%B3n%2F1%2E%20Saber11%2F2%2E%20Diccionarios%20Saber11%2FDiccionario%20Saber11%202019%2D1%20a%202020%2D2%2Epdf&parent=%2Fpersonal%2Fdataicfes%5Ficfes%5Fgov%5Fco%2FDocuments%2FDataIcfes%2F4%2E%20Saber11%2F2%2E%20Documentaci%C3%B3n%2F1%2E%20Saber11%2F2%2E%20Diccionarios%20Saber11

[5] Secretaría distrital de planeación. *Estratificación económica-Generalidades.* Retrieved from http://www.sdp.gov.co/gestion-estudios-estrategicos/estratificacion/generalidades

[6] *GDP per capita (current US$). World Development Indicators*. World Bank. Retrieved from https://data.worldbank.org/indicator/NY.GDP.PCAP.CD?most_recent_value_desc=true

[7] Stan Brown, Oak Road Systems *Measures of Shape: Skewness and Kurtosis*, 2008–2016.

[8] Taboga, Marco (2017). "Normal distribution - Maximum Likelihood Estimation", Lectures on probability theory and mathematical statistics, Third edition. Kindle Direct Publishing. Online appendix. https://www.statlect.com/fundamentals-of-statistics/normal-distribution-maximum-likelihood.

[9] Isabel Molina Peralta and Eduardo García Portugués (2021)*, A First Course on Statistical Inference*, *v0.8*

[10] Tukey JM, McLaughlin DH (1963) *Less Vulnerable Confidence and Significance Procedures for Location Based on a Single Sample: Trimming/Winsorization 1*. Sankhya A, 25:331–352.

[11] Altman, D. G., Machin, D., Bryant, T. M., and Gardner, M. J., eds. (2000). *Statistics with Confidence: Confidence Intervals and Statistical Guidelines.* BMJ Books, London.

[12] Lu, Zhenqiu & Yuan, Ke-Hai. (2010). Welch's t test. 10.13140/RG.2.1.3057.9607.