

A Comprehensive Study on US College Admissions

Tsirmpas Dimitris

Athens University of Economics and Business

Department of Informatics

May 22, 2023



Athens University of Economics and Business

Department of Statistics

Greece

Contents

1	Abstract	2
2	Introduction	2
3	Exploratory Analysis	2
4	Variable correlations	4
4.1	Writing scores influenced by gender	4
4.2	Writing scores influenced by previous program	5
5	Predictive / interpretative models	5
5.1	Building the base model	6
5.2	Identifying test score causation	10
6	Conclusions & Discussion	14
7	References	14
8	Addendum	16
8.1	Exploratory Analysis	16
8.2	OLS model preconditions	16

1 Abstract

US College Admissions have and continue to be a subject of great debate among scholars and analysts. Such educational institutions have an interest in selecting the most qualified applicants using limited data, while the applicants themselves often protest admission requirements, especially those deemed discriminatory in nature. In this study we discover relationships between the candidate's gender and previous program and their overall test scores. We also identify a positive correlation between test scores, which we attribute to a common, unknown variable called "Competence".

2 Introduction

The aim of this study is to investigate possible links and relationships between a candidate's characteristics and their performance in multiple standardized tests. We employ a random sample of 200 students who applied to continue their studies in their respective universities. Their application consisted of three standardized tests testing their skills and knowledge in mathematics, social studies and creative writing. An overview of the data contained can be found in Table 1.

The data and replication code can be found in our GitHub repository ¹. We make the assumption that the dataset has been acquired through random, unbiased sampling. We also make the assumption that the records using different IDs represent different students (and as such are considered independent samples).

The study is structured as follows: In Section 3 we make general observations about our data, identify key relationships and form our first hypotheses. We follow up these hypotheses in Section 4 with robust analyses and in Section 5 by employing various regression models. Section 6 features an overview and discussion about our findings. Finally, we include graphs, tables and supporting documents in the report's Addendum (Section 8).

3 Exploratory Analysis

The numerical variables contained in the dataset are described in Table 2. We observe that they are all almost symmetrical ($-0.5 \leq skew \leq 0.5$), and feature moderate negative (right) skewness with a mean/median hovering just above a score of 50. This indicates most students score around the baseline, most of which pass the exams with a mediocre grade.

We next study relationships between the candidates' characteristics (not including test scores). We run χ^2 tests on *Gender*, *Race*, *Program* and *School Type*. The only

¹<https://github.com/dimits-exe/collegeanalysis>

Name	Type	Description	Range
Id	Nominal	The student's ID	[1-200]
Gender	Binary	The student's gender	{male, female}
Race	Nominal	The student's race	{white, latin-american, asian, african-american}
Schtype	Binary	The type of the student's secondary education institution	{public, private}
Prog	Nominal	The student's previous study cycle	{general, vocation, academic }
Write	Numeric	The grade on the writing test	[0-100]
Math	Numeric	The grade on the mathematics test	[0-100]
Socst	Numeric	The grade on the social studies test	[0-100]

Table 1: An overview of the data used in this study.

Var.	Obs.	Mean	Std	Median	Trim	Min	Max	Skew	Kurt	SE
Write	200	52.77	9.48	54	53.36	31	67	-0.4	-0.7	0.67
Math	200	52.65	9.37	52	52.33	33	75	0.28	-0.6	0.66
Socst	200	52.41	10.74	52	52.99	26	71	-0.3	-0.5	0.76

Table 2: Summary statistics on the numerical data used in the study.

statistically significant relationship in our dataset is between *School Type* and *Program* ($p_value = 0.015$), which can be seen in Table 3.

Finally, we study the relationships between the three subjects. As shown in Table 4, there is a very statistically significant ($p_value = 0.0000$), positive (Pearson's $r > 0.5$) relationship between all three subjects. This could be indicative of either one of the variables influencing the other, or an unknown, interfering variable which positively affects the three test scores. We hypothesize the latter, as the existence of such a variable indicating the student's general competence in tests makes intuitive sense. We will refer to this potential, unknown variable as "*Competence*" in this report.

Since the three subjects are strongly correlated we can explore correlations between

School Type	Previous Program			Total
	general	academic	vocation	
public	39	81	48	168
	23.2%	48.2%	28.6%	100%
private	6	24	2	32
	18.8%	75%	6.2%	100%
Total	45	105	50	200
	22.5%	52.5%	25%	100%

Table 3: χ^2 test between *School Type* and *Program*. Notice the overwhelming majority of candidates who attended private schools having an academic background prior to applying.

	Writing	Math	Social Studies
Writing			
Math	0.62 0.000 ****		
Social Studies	0.60 0.000 ****	0.54 0.000****	

Table 4: Pearson’s correlation coefficient (Holm’s correction) between the tests and their p-values. Stars indicate significance scores: > 1:”, 0.1:’*’, 0.01: ’**’, 0.001: ’***’, < 0.0001: ’****’.

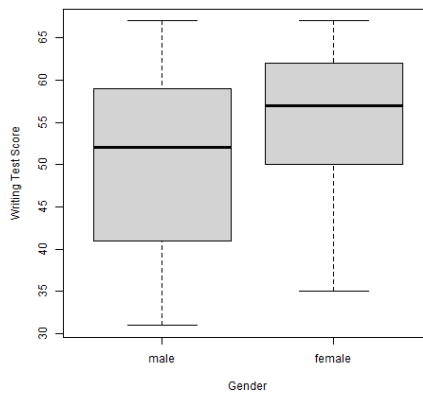


Figure 1: Boxplot displaying the writing score by gender.

the rest of the factors and any of the tests, assuming that a correlation with one strongly indicates with the other two as well.

We notice a probable correlation between gender and writing scores, as shown in Figure 1, as well as between the student’s program and writing scores, as shown in Figure 2.

4 Variable correlations

4.1 Writing scores influenced by gender

Our exploratory analysis indicated a possible discrepancy between the results of writing tests between men and women, as well as between different programs. We thus investigate whether gender and the candidate’s current program play a role in writing test score.

We begin by verifying the preconditions necessary for the standard t-test in order to compare the genders’ scores. We compute the mean differences of the samples by subtracting the global mean by the women’s scores [2], and conclude they are not normal

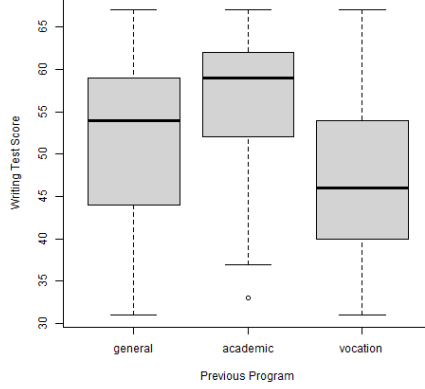


Figure 2: Boxplot displaying the writing score by current program. Notice the significant differences in means.

(Shapiro-Wilk normality test, $p_value = 0.0024$). The variances are not homogeneous (Levene's test $p_value = 0.0022$), although this should not dissuade us from using a parametric t-test since the relatively large sample size ($N = 200$) and balanced groups ($N_{women} = 109, N_{men} = 91$) means the precondition's violation is not significant [3]. We will use a non-parametric test since the distribution of the differences do not have a normal kurtosis (Jarque-Bera Normality Test [1], $p_value = 0.026$), and thus their mean is not suitable.

We conclude there is a statistically significant difference between the writing scores of men and women (Wilcoxon Rank Sum Test, $p_value = 0.0009$) with women having on average 5 more score than men (Wilcoxon Rank Sum Test with $H_a = less$, $p_value = 0.0004$).

4.2 Writing scores influenced by previous program

We will now verify the preconditions for the parametric ANOVA test in order to test which programs are correlated with the writing tests' scores. The variances of the residuals are homogeneous (Levene's test $p_value = 0.0022$), but not normal (Shapiro-Wilk normality test, $p_value = 0.002$). We will thus use a non-parametric test.

We discover there is a statistically significant difference between the groups (Kruskal-Wallis rank sum test, $p_value = 4e - 08$). We can consult the boxplots in Figure 2, where we observe significant differences between all three groups.

5 Predictive / interpretative models

5.1 Building the base model

Having confirmed our hypotheses regarding the relationships between the various variables and the writing test scores, we attempt to build a model which will estimate a candidate's math and social study scores, testing our hypothesis that the three scores are influenced by the same external factors.

Since there seem to be strong correlations between most independent variables and the writing scores, and since we already established a strong correlation between the scores themselves (attributed to the candidate's "*Competence*"), we will be using a simple, Ordinary Least Squares (OLS) model.

We initially build an OLS model estimating the writing scores which involves all the available variables. This model exhibits a good fit ($R^2_{adj} = 0.4753, F = 37.06, p_value = 2.2e - 16, df = 194$). We also verify all the preconditions for the regression model; the residuals appear to be normal (Shapiro-Wilk normality test $p_value = 0.2099$), homogeneous (Levene's Test, $p_value = 0.1057$) and non-correlated (Durbin-Watson test, $p_value = 0.236$) and there are no significant outliers (see Addendum).

Our current model exhibits high confidence about the gender ($p_value_{female} = 0.0056$), previous program ($p_value_{progacademic} = 0.0025$), writing scores ($p_value_{write} = 3.57e - 08$), social study scores ($p_value_{socst} = 0.0053$) and the constant ($p_value_{intercept} = 8.29e - 09$). However, it exhibits low confidence about the candidate's race ($p_value_{raceasian} = 0.05$) and no confidence about his school type ($p_value = 0.6502$). In order to build an optimal model, we consider dropping these two variables. Dropping the candidate's race along with the school type leads to a slight decrease of $R^2_{adj} = 0.4753$. Dropping only the school type on the other hand leads to an improved $R^2_{adj} = 0.4882$. Since the model maintains its confidence about the other variables, we keep this model as the optimal one. A summary of the optimal model can be found in Table 5.

We now follow the same procedure for the social studies test. We fit a model involving all available variables, which explains our sample to a decent degree ($R^2_{adj} = 0.4535, F = 17.7, p_value = 2.2e - 16, df = 188$). We detect an outlier in our precondition verification, a student with an above-average grade in writing but an abysmal one in social studies. Since this individual exerts a significant influence in our model (as it is largely dependent on the writing scores for its estimations) we remove them from the sample. The other preconditions are met, the residuals are sufficiently normal (Shapiro-Wilk normality test $p_value = 0.0125$), homogeneous (Levene's Test, $p_value = 0.3394$) and non-correlated (Durbin-Watson test, $p_value = 0.678$).

As mentioned above, this model is reliant on the math and especially on the writing scores, while the rest of the variables are mostly non-statistically significant. We thus employ a stepwise procedure in order to eliminate variables deemed statistically insignificant while retaining, or increasing our models goodness of fit. A summary of the resulting model can be found in Table 6.

The results seem to verify our main hypothesis, that the test scores are predominately caused by the unknown "Competence" value. While other variables such as the candidate's past program, have a statistically significant influence in our model, we can see that the model's estimations are consistently based on the writing and other-lesson's test scores. We can additionally rule out this relationship being a result of multi-co-linearity between the test variables, as seen in the previous auto-correlations tests.

Table 5: Linear regression model predicting math test scores, taking into account other test scores.

	<i>Dependent variable:</i>
	math
genrefemale	−2.828*** (−4.799, −0.858)
progacademic	3.788*** (1.346, 6.229)
progvocation	−0.375 (−3.161, 2.410)
write	0.404*** (0.266, 0.542)
socst	0.167*** (0.052, 0.283)
raceasian	4.978* (−0.017, 9.972)
raceafrican-amer	−1.103 (−5.104, 2.898)
racewhite	2.324 (−0.686, 5.334)
Constant	20.334*** (13.730, 26.938)
Observations	200
R ²	0.509
Adjusted R ²	0.488
Residual Std. Error	6.702 (df = 191)
F Statistic	24.729*** (df = 8; 191)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Table 6: Linear regression model predicting social study test scores, taking into account other test scores.

	<i>Dependent variable:</i>
	socst
raceasian	−5.642* p = 0.060
raceafrican-amer	0.778 p = 0.745
racewhite	0.308 p = 0.865
progacademic	2.264 p = 0.131
progvocation	−2.939* p = 0.077
write	0.479*** p = 0.000
math	0.233*** p = 0.005
Constant	14.585*** p = 0.001
Observations	199
R ²	0.477
Adjusted R ²	0.458
Residual Std. Error	7.844 (df = 191)
F Statistic	24.899*** (df = 7; 191)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

5.2 Identifying test score causation

The results above, although encouraging, do not rule out the alternative hypothesis we posed in the Introduction of this report, that the definite correlation between the test scores is caused by one of the test scores themselves influencing the others. We thus repeat the experiments of the previous subsection while omitting the writing scores. If our alternative hypothesis was correct we would expect our models to no longer have a good performance, while having a much smaller statistical significance on the other lesson's model.

We begin by constructing an OLS regression model which tries to estimate the math score by considering the candidate's characteristics and his scores in the social studies test. Our initial model including all the variables scores a $R^2_{adj} = 0.3999$ score. We next verify all the necessary preconditions; the residuals appear to be normal (Shapiro-Wilk normality test $p_value = 0.8403$), homogeneous (Levene's Test, $p_value = 0.1641$) and non-correlated (Durbin-Watson test, $p_value = 0.534$) and there are no significant outliers.

This model again displays high confidence in the following variables: Race ($p_value_{raceasian} = 0.002$, $p_value_{racewhite} = 0.012$), Program ($p_value_{progacademic} = 0.0005$), Social Study Test Scores ($p_value_{socst} = 4.25e - 09$) and the Constant ($p_value_{socst} = 2e - 16$). By employing a backwards procedure we end up with our final model with a total $R^2_{adj} = 0.4022$, found in Table 7.

We repeat the procedure for estimating the social study test scores without relying on the writing tests. Our base model, which considers all the available variables, displays a $R^2_{adj} = 0.3393$, which is considerably worse than the respective math model. This may be because of the previously mentioned reliance on the writing scores, which are no longer available to the model. Additionally, similarly to the previous models, the model lacks confidence in almost all other variables; the only statistical significant variable other than the math scores ($p_value_{math} = 4.25e - 09$) and the Constant ($p_value_{Intercept} = 5.27e - 08$) is the candidate's Program ($p_value_{progvocation} = 0.0283$).

We again verify the preconditions necessary for the linear regression model. The residuals appear to be normal (Shapiro-Wilk normality test $p_value = 0.7449$), homogeneous (Levene's Test, $p_value = 0.2314$) and non-correlated (Durbin-Watson test, $p_value = 0.952$) and there are no significant outliers.

Because of the many possible variables that are candidates for removal we can again employ a stepwise model selection algorithm. The best model by AIC keeps only the math and program variables (as expected) but explains less of the data ($R^2_{adj} = 0.3383$). Since the reduction in R^2_{adj} is minimal, and since the new model can reach this score by discarding almost all other variables (which as described above are essentially considered as noise by our model), we will be using this model as optimal (Table 8).

These findings seem to contradict our alternative hypothesis. While our models' performance certainly degraded, they still achieve comparable results, with their performance loss being explained by the degree to which "*Competence*" can be measured. In other words, this unknown variable can be approximated more accurately by considering both other tests, instead of just one. This is further proof that our initial hypothesis appears to be correct.

Table 7: Linear regression model predicting math test scores, without relying on the writing tests.

	<i>Dependent variable:</i>
	math
raceasian	8.074*** p = 0.003
raceafrican-amer	-1.261 p = 0.567
racewhite	4.026** p = 0.015
progacademic	4.726*** p = 0.0005
progvocation	-1.039 p = 0.499
socst	0.335*** p = 0.000
Constant	29.617*** p = 0.000
Observations	200
R ²	0.420
Adjusted R ²	0.402
Residual Std. Error	7.243 (df = 193)
F Statistic	23.318*** (df = 6; 193)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

Table 8: Linear regression model predicting social study test scores, without relying on the writing tests.

	<i>Dependent variable:</i>
	socst
progacademic	2.833* p = 0.085
progvocation	-3.829** p = 0.037
math	0.486*** p = 0.000
Constant	26.285*** p = 0.000
Observations	200
R ²	0.348
Adjusted R ²	0.338
Residual Std. Error	8.733 (df = 196)
F Statistic	34.908*** (df = 3; 196)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

6 Conclusions & Discussion

In this report we studied extensively the relationships between different characteristics of US university applicants. We observed that characteristics such as race, gender, schooling and previous programs are uncorrelated, with the exception of a positive relationship between private schooling and academic background. We verified that female candidates score on average higher than males in writing tests, as well as that the candidate's previous program influences their writing test scores.

Additionally, we observed a definite positive correlation between candidate test scores. We hypothesized this was the product of an unknown variable, called "*Competence*" which similarly influences all test scores and posed an alternative hypothesis stating that one of the test scores was the cause of the correlation. We show evidence of this variable's existence by constructing linear regression models and disprove the alternative hypothesis by constructing such models without access to the common variable *Writing Score*.

This study highlights that other test scores are consistently robust and important variables when attempting to assess future test scores, even if these test scores are on different disciplines (such as social studies and mathematics). Further research is warranted to check whether past test scores can be used to consistently predict candidate performance, as well as from how far in the past, and within which disciplines these scores would be useful.

We note that these findings are only representative for our sample. Our statistical models were used exclusively as explanatory models and should not be used for prediction. We also note that the relationships presented in this report are only representative of our sample, and should not be generalized for the general population. Finally, we warn against extrapolating any relationships from our tests, since they only imply a correlation, not necessarily causation.

7 References

- [1] Carlos M. Jarque and Anil K. Bera. "Efficient tests for normality, homoscedasticity and serial independence of regression residuals". In: *Economics Letters* 6.3 (1980), pp. 255–259. ISSN: 0165-1765. DOI: [https://doi.org/10.1016/0165-1765\(80\)90024-5](https://doi.org/10.1016/0165-1765(80)90024-5). URL: <https://www.sciencedirect.com/science/article/pii/0165176580900245>.
- [2] Geoffrey R Loftus and Michael EJ Masson. "Using confidence intervals in within-subject designs". In: *Psychonomic bulletin & review* 1.4 (1994), pp. 476–490.

- [3] Donald W Zimmerman. "A note on preliminary tests of equality of variances". In: *British Journal of Mathematical and Statistical Psychology* 57.1 (2004), pp. 173–181.

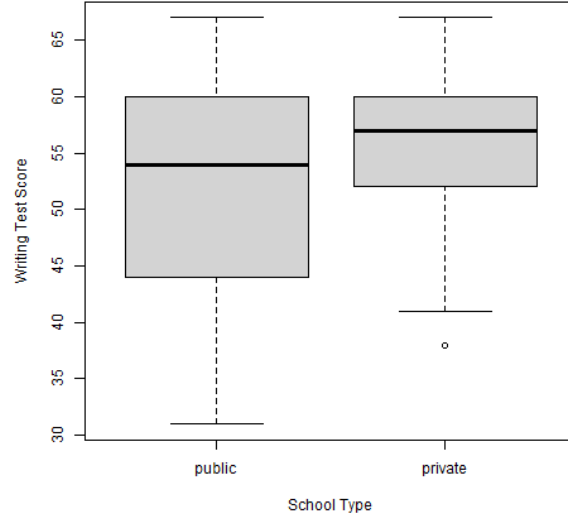


Figure 3: Boxplot displaying the writing score by school type.

8 Addendum

8.1 Exploratory Analysis

In this section we include tables and figures which were used in the exploratory analysis of the data in the Introduction.

8.2 OLS model preconditions

In order to visually confirm the normal distribution of the model's residuals, we plot their boxplots for each of the 4 quantiles. We expect these boxplots to resemble those of the normal distribution, centered on $y = 0$ and with 95% of their values not going above/below the $y = 1.95$ and $y = -1.95$ respectively. Figures 5, 6 show models including the writing scores. Figures 7, 8 show models including only the other lesson's scores.

We also check for outliers by plotting the normalized residuals against the model's estimations. The red lines denote the $y = 1.95$ and $y = -1.95$ values respectively and we expect 95% of the points to be within them. Any value outside of $[-3, 3]$ indicates a strong outlier which must be investigated. Figures 9, 10 show models including the writing scores. Figures 11, 12 show models including only the other lesson's scores.

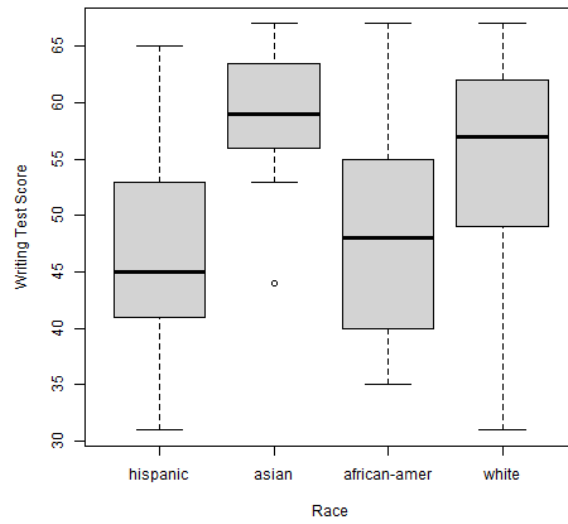


Figure 4: Boxplot displaying the writing score by candidate race.

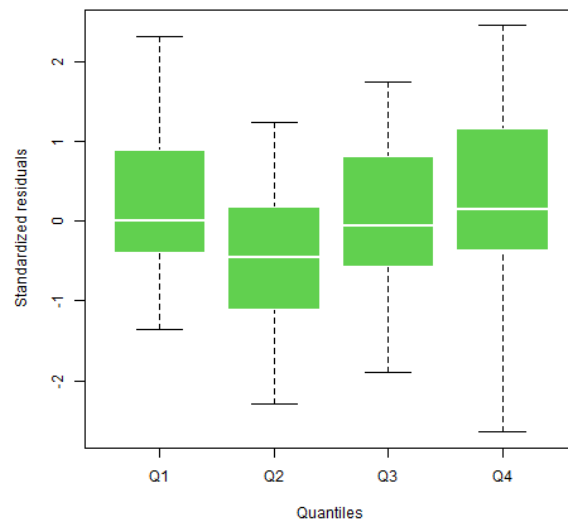


Figure 5: The normalized residuals of the math model plotted for each of the 4 quantiles.

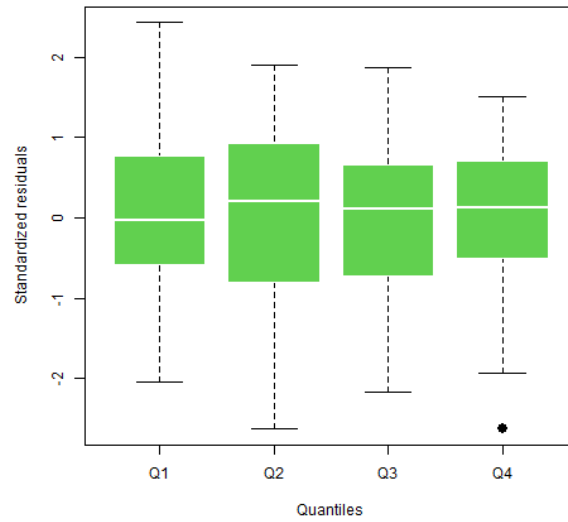


Figure 6: The normalized residuals of the social studies model plotted for each of the 4 quantiles.

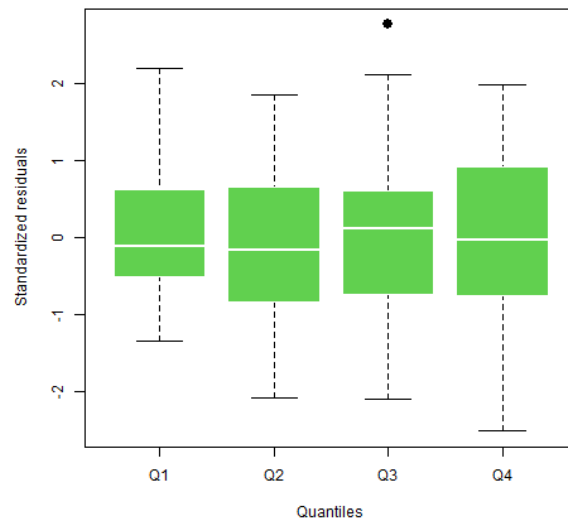


Figure 7: The normalized residuals of the math model plotted for each of the 4 quantiles.

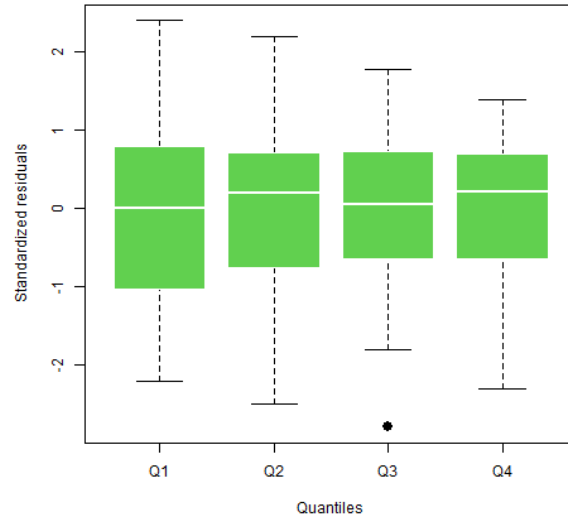


Figure 8: The normalized residuals of the math model plotted for each of the 4 quantiles.

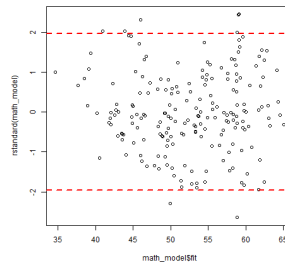


Figure 9: The normalized residuals of the math model plotted against the model's estimations.

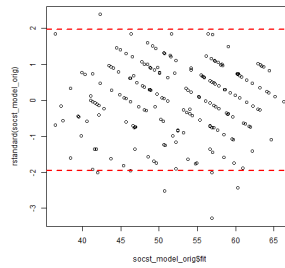


Figure 10: The normalized residuals of the social studies model plotted against the model's estimations. Notice the one outlier above the $y = 3$ line, which represents the data point that was investigated and removed.

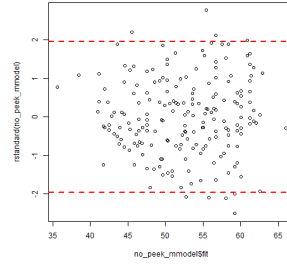


Figure 11: The normalized residuals of the math model plotted against the model's estimations. We notice three potential outliers. These values are considered non-anomalous, as they stray sufficiently away from the $y = 3$ and $y = -3$ brackets, and are to be expected in a sample of 200 values.

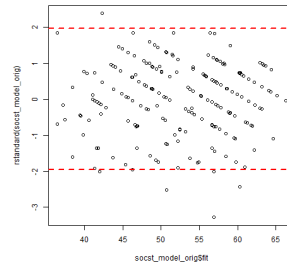


Figure 12: The normalized residuals of the social studies model (without the writing scores variable) plotted against the model's estimations. We don't consider the values above and below the $y = 3$ and $y = -3$ brackets respectively for the same reasons as in Figure 11.