

How school type, language, special education services, family income, and parental education level influence OSSLT first attempt results in Ontario schools

STA302 Final Project Part 1

Xuanle Zhou Luhan Wang Junyi Hou

April 10, 2025

1 Introduction

The Ontario Secondary School English Literacy Test (OSSLT) is mandatory for high school graduation in Ontario, therefore English language learning is a significant focus for both parents and students. This paper aims to investigate how school type, language, special education services, family income, and parental education level influence OSSLT first attempt results in Ontario schools. Zhang et al. (2020) found that family income and parental education level significantly contribute to a student's academic success. Their study was conducted in China, and revealed that higher family income and more advanced parental education are correlated with better student performance. This supports and shapes our hypothesis that students with higher family income and parental education level will perform better on the OSSLT. Bernhofer and Tonin (2022) showed that students perform better when taught in their first language, which questions if English language school students will perform better on the OSSLT compared to those in non-English language schools. Lastly, Aseery (2024) explored how technology and multimedia elements in religious education classes could enhance English language learning. Aseery's findings suggest that multimedia tools in religious education classes improve student engagement and motivation, which enhances learning outcomes. We would expect schools supplying these technologies in 2025. Therefore, we hypothesize that religious schools will have higher OSSLT pass rates.

While Zhang et al. (2020) concluded that higher income and parental education level lead to higher achievement, there are exceptions, as many successful individuals come from lower-income backgrounds. We also expect that students receiving special education services may

Table 1: OSSLT First Attempt Pass Rate Descriptive Statistics

	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
OSSLT_First_Attempt_	737	82.45	11.93	85	83.92	10.38	0	100	100	-1.93	6.93	0.44

perform worse on the OSSLT due to specific learning disabilities, despite receiving accommodations. This research question fits well with the concept of multiple linear regression, which examines how multiple predictor variables collaboratively influence a response variable. Therefore, we have selected multiple linear regression as our analysis method. Since the main goal is to observe patterns between variables, this model will focus on interpretability.

This research will benefit those seeking an accurate analysis of the factors that influence English learning outcomes, particularly in the context of the OSSLT. The response variable, OSSLT results, serves as an effective measure of students' English proficiency, as it is both a pass/fail test and provides continuous data.

2 Data description

The dataset, available on the Ontario Data Catalogue (Ontario 2024b), provides insights into schools in Ontario, supporting policy-making, and educational research. This study repurposes it to investigate and predict the OSSLT first-attempt pass rate. Data were collected from schools, school boards, EQAO, and Statistics Canada through online forms, surveys, phone interviews, and in-person visits, then compiled by Ontario Data Catalogue (Ontario 2024a).

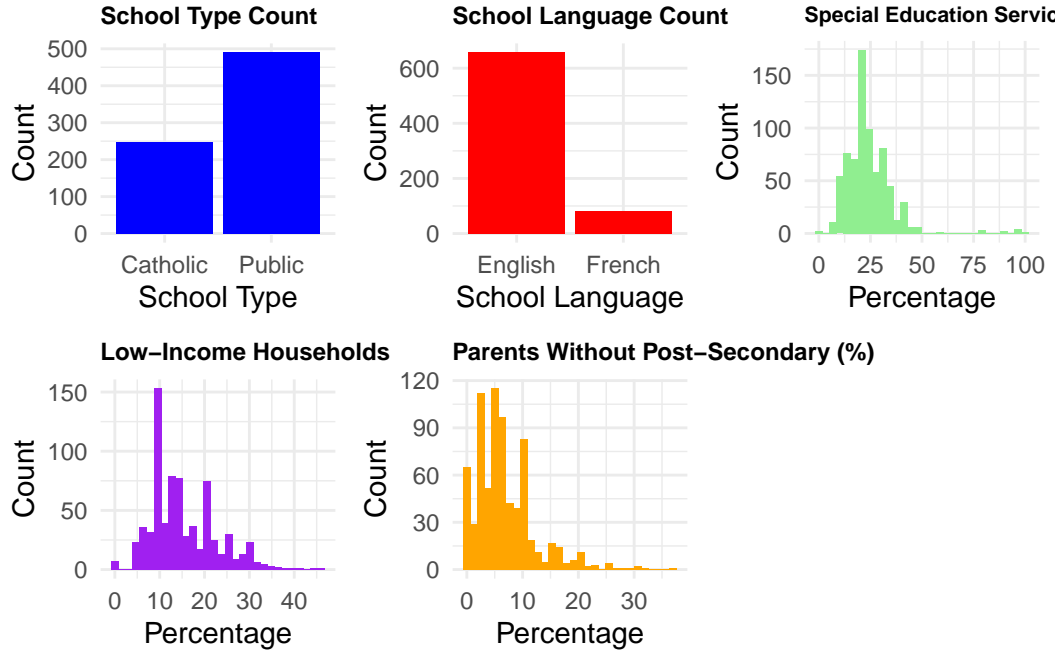
The **OSSLT_First_Attempt_PassRate**, the response variable, measures the percentage of students passing Ontario Secondary School Literacy Test on their first attempt, ranging from 0 to 100. The mean of 82.45 and median of 85 indicate high pass rates. The dataset originally had 4,926 observations, reduced to 737 after cleaning, ensuring statistical reliability. Despite being bounded, the pass rate is continuous, suitable for linear regression.

School Type is categorical, with two types: Catholic and Public. Most schools are public. Cheema (2024) noted, private schools generally outperform public schools in literacy. We expect Catholic schools to have higher OSSLT pass rates due to structured curriculum and discipline.

School Language is binary, English or French. Most schools operate in English, which is expected to correlate with higher OSSLT pass rates.

Students receiving special education services often exhibit lower literacy achievement and slower progress, as noted by Vaughn and Wanzek (2014). Our model aims to capture this pattern. The mean of this predictor variable is 24.07%, with a median of 22%, includes outliers where 100% of students receive special education services.

Table 2: Histograms for Selected Predictors



The percentage of school-aged children in low-income households has a mean of 15.27% and skewness of 0.88, indicating some schools have significantly higher concentrations. As Nadeem, Akhtar, and Ahmad (2021) found, lower-income students often have lower literacy skills, which we expect to correlate with lower OSSLT pass rates.

The percentage of students whose parents lack post-secondary credentials averages 6.76%, with skewness of 1.56 and kurtosis of 3.73, suggesting a slight right skew. As Davis-Kean, Tighe, and Waters (2021) states, parental education influences children's academic success, making this a relevant predictor.

3 Preliminary results

Table 3: Regression Preliminary Results

	Coefficient	Standard_Error	t_Statistic	p_Value
(Intercept)	105.176	1.005	104.633	0.000
School_TypePublic	-1.069	0.652	-1.640	0.101
School_LanguageFrench	5.459	0.984	5.547	0.000
Special_Ed_Pct	-0.621	0.026	-23.760	0.000

	Coefficient	Standard_Error	t_Statistic	p_Value
Low_Income_Pct	-0.296	0.048	-6.117	0.000
No_Parent_Degree_Pct	-0.465	0.065	-7.104	0.000

3.1 Residual Analysis

3.1.1 Linear Models Assumptions:

1. Linearity

$$E(Y_i|X = \mathbf{x}_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}$$

2. Constant Error Variance (Homoscedasticity)

$$Var(Y_i|X = \mathbf{x}_i) = \sigma^2$$

3. Uncorrelated and Normal Errors

$$Cov(e_i, e_j) = 0 \text{ for } i \neq j \text{ and } e_i \sim N(0, \sigma^2)$$

3.1.2 Assumption Check

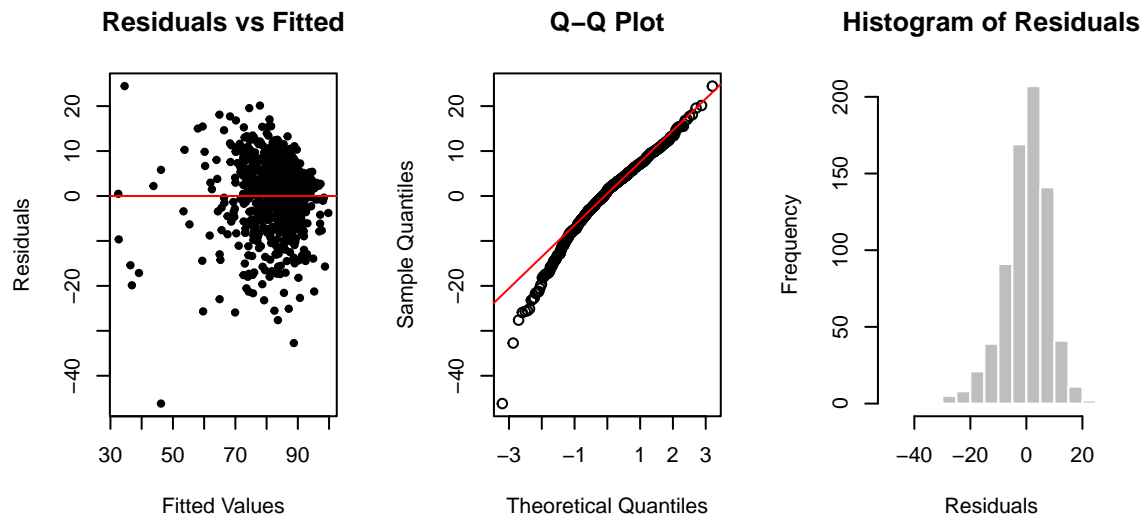


Figure 1: Residual Plots

1. **Linearity & Homoscedasticity:** The residuals vs. fitted plot shows no clear pattern, suggesting linearity. Slight heteroscedasticity is observed.
2. **Normality:** The Q-Q plot and the histogram of residuals suggest residuals are approximately normal, though slight deviations exist at the left tail.
3. **Independence:** No evident pattern in the residual plot suggests residuals are independent.

3.2 Model Interpretation & Discussion

3.2.1 Key Findings and interpretation

- The **intercept (105.18)** represents the estimated pass rate for a **Catholic, English-language school with 0% special education, 0% low-income students, and 0% students whose parents have no degree**. This provides a reference point for understanding the model's predictions.
- **School Language (French vs. English)** and the three numeric variables (**Special_Ed_Pct**, **Low_Income_Pct**, **No_Parent_Degree_Pct**) are strongly associated with the **OSSLT pass rate**.
- **School Type (Public vs. Catholic)** does not show a statistically significant difference in pass rate in this model.
- Higher proportions of **special education students, low-income students, and students whose parents have no degree** are each associated with a **lower pass rate**.
- Conversely, being a **French-language school** is associated with a **higher pass rate** relative to the English.
- The model explains **54% of the variation in pass rates**, which is reasonable for educational data, suggesting these variables collectively have a substantial but not complete ability to predict pass rates.

3.2.2 Comparison to Literature

Our findings align with prior research while offering insights specific to Ontario:

- **Family Income & Parental Education:** Consistent with Zhang et al. (2020), our results confirm that higher family income and parental education correlate with better OSSLT pass rates.
- **School Language:** Contrary to Bernhofer and Tonin (2022), our study shows French-language schools had higher OSSLT pass rates than English-language schools, indicating other factors like curriculum or funding may play a role. Further investigation is needed.

Table 4: Descriptive Statistics for Selected Predictors

	vars	n	mean	sd	median	trimmed	mad	min	max	range	skew	kurtosis	se
School_Type*	1	737	1.66	0.47	2	1.71	0.00	1	2	1	-0.70	-1.52	0.02
School_Language*	2	737	1.11	0.31	1	1.01	0.00	1	2	1	2.51	4.31	0.01
Special_Ed_Pct	3	737	24.07	11.53	22	22.91	8.90	0	100	100	2.72	13.46	0.42
Low_Income_Pct	4	737	15.27	7.30	13	14.58	5.93	0	46	46	0.88	0.66	0.27
No_Parent_Degree_Pct	5	737	6.76	5.42	5	6.06	4.45	0	37	37	1.56	3.73	0.20

- **Special Education:** Higher proportions of special education students negatively impact OSSLT success, aligning with expectations.
- **School Type:** No significant difference was found between public and Catholic schools, despite Aseery (2024) suggesting that religious schools may benefit from enhanced multimedia learning tools.

4 Model Selection

4.1 Response Variable Transformation and Assumption Comparison

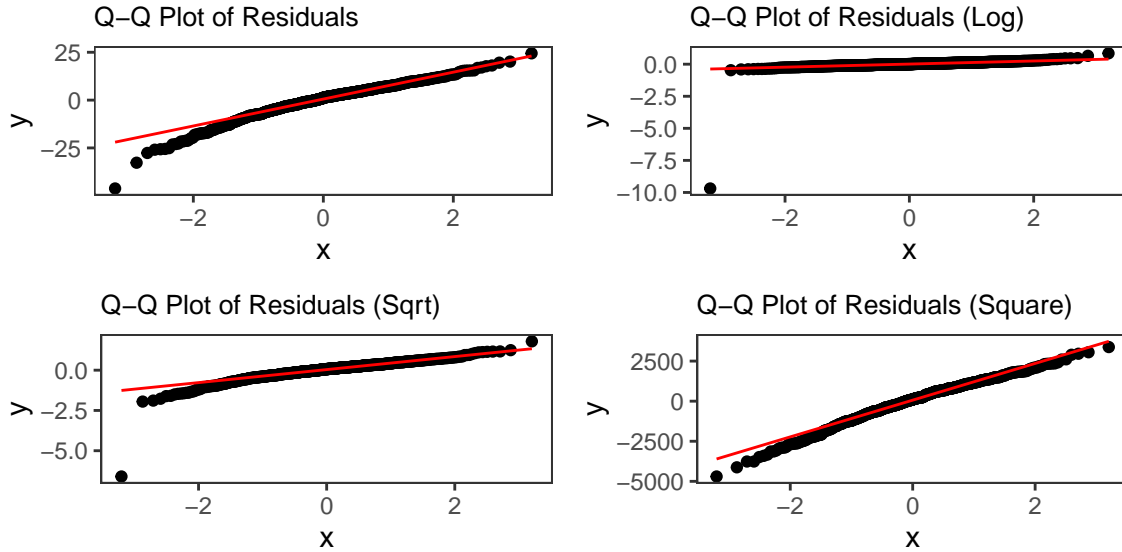


Figure 2: Residual Plots

After fitting our original model, we consider some transformations. The first method is Box Cox, which uses maximum likelihood to choose transformation so the residuals are approximately normally distributed. We found 2 as Box Cox lambda, suggesting squared transforma-

tion. Therefore we Try Y^2 Transformation, we also tried many other transformations such as $1/Y$, and decided to display log and square root as they improved from original model.

Normality: The Q-Q plot of residuals suggest residuals are approximately normal for square root and log, though slight deviations exist at the left tail. Note that the scale for the Square root graph is a lot bigger, minimal deviation indicates violation.

4.2 Response Variable Transformation Preview with Residuals

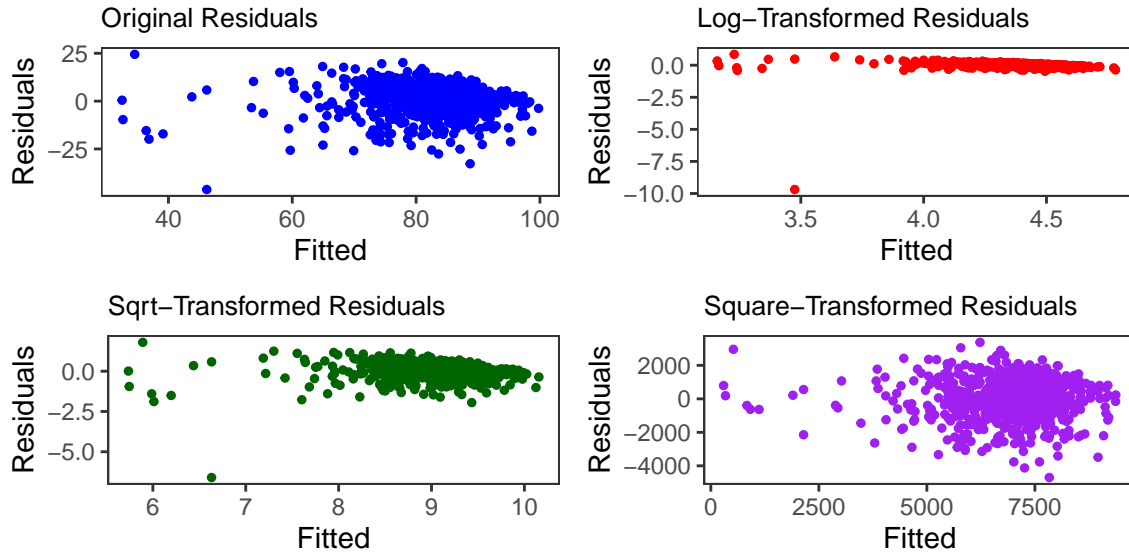


Figure 3: Residual Plots

Linearity & Homoscedasticity: The Square Root residuals vs. fitted plot is ideal, as the Y scale range is a lot smaller than other transformations and there is no clear patterns. This suggest linearity. Slight heteroscedasticity is observed around $x=10$.

4.3 Response Variable Transformation Metrics and Decision

We eliminate Square since the model performs worse in every aspect. Square Root and Log Transformation both rapidly improve from the original model based on AIC and BIC. However, R^2 dropped by approximately 34% compared to original for Log, indicating log model's ability to explain variance significantly decreased. In contrast, Square Root Model R^2 only decreased by less than 2% for better Normality, linearity and Homoscedasticity. The difference in AIC and between Log and Square Root is not significant enough to replace R^2 . Therefore, we apply square root transformation on our response variable.

Table 5: R^2 , AIC, BIC

Model	R_Squared	AIC	BIC
Original	0.5415	5179.0167	5211.2348
Log-Transformed	0.1999	687.2028	719.4209
Sqrt-Transformed	0.5265	1168.6607	1200.8788
Square	0.5057	12565.1800	12597.3981

4.4 Y transformed model summary, VIF and Confidence Interval

The confidence interval contains 0 and p value > 0.05 for School_Type, suggest dropping this predictor. Variance Inflation Factors > 5 for all, indicating no predictor have issues with multicollinearity.

Table 6: VIF, Confidence Interval and Summary

Predictor	Estimate	Std. Error	t value	p-value	Lower 95% CI	Upper 95% CI	VIF
(Intercept)	10.5120	0.0662	158.87	$< 2e-16$	10.3821	10.6419	
School Type - Public	-0.0530	0.0429	-1.23	0.217	-0.1372	0.0313	1.07
School Language - French	0.2982	0.0648	4.60	4.91e-06	0.1710	0.4253	1.06
% in Special Education	-0.0414	0.0017	-24.05	$< 2e-16$	-0.0447	-0.0380	1.02
% in Low Income	-0.0198	0.0032	-6.24	7.59e-10	-0.0261	-0.0136	1.41
% in No Education - Parent	-0.0241	0.0043	-5.59	3.16e-08	-0.0325	-0.0156	1.42

4.5 X Transformation Assumption Preview

Other than dropping School Type, we also consider X transformation. After testing different combinations of X transformations, we found that adding a squared special education term slightly increase model performance. Therefore we compare four models: Original with $\text{Sqrt}(Y)$, remove school type, add squared special education, and both remove and add. All X transformations models shows similar residual plots.

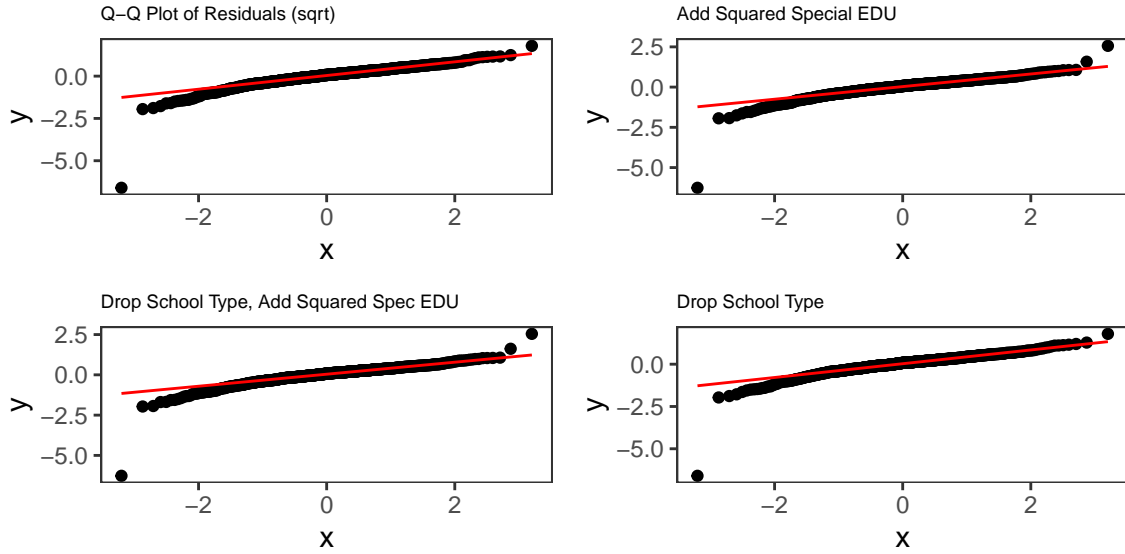


Figure 4: QQ of all possible X Transformation

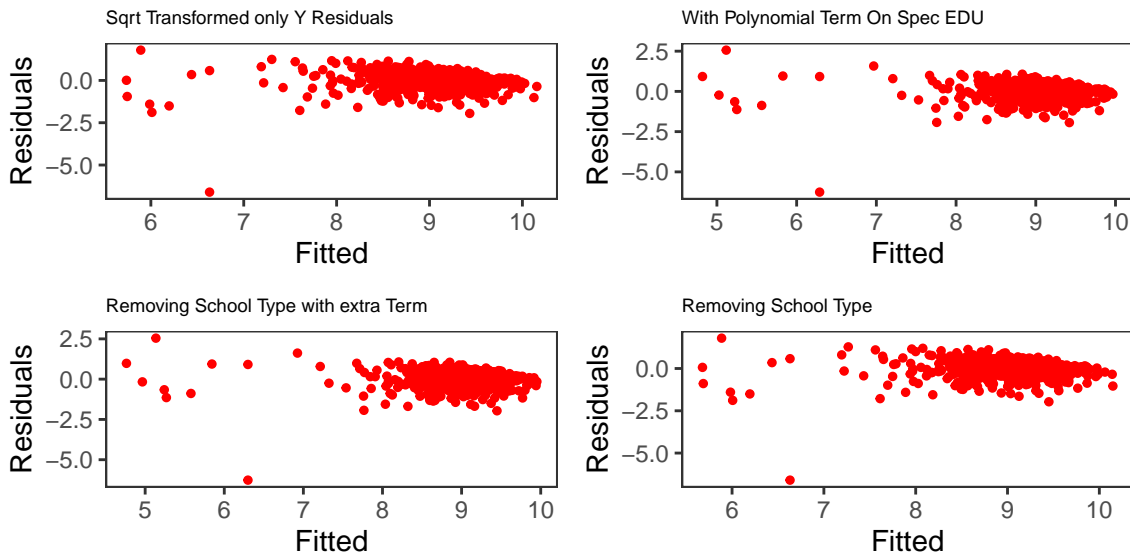


Figure 5: Residual Plots

4.6 Predictor Transformation Metrics and Decision

R^2 , AIC, BIC, RSS, F and P value all improved from models that has the Squared Special Education predictor. We decide to drop school type, since all metrics does not differ by much,

indicating that this predictor has minimal impact on the model, or no relationships. Finally, we can conclude the final model with 5 Predictors: School Language, Special Education Proportion, Special Education Squared, Low Income Proportion, Parents with minimal education proportion.

Table 7: Predictor Transformation Decision

Model	R_Squared	AIC	BIC
Original (Y) Transformed	0.5265	1168.661	1200.879
X Transform	0.5458	1138.996	1175.817
X Transform with One Less Predictor	0.5444	1140.411	1172.629
One Less Predictor	0.5262	1168.196	1195.812

Table 8: Anova F score for X Transformation

	Model	F_statistic	df1	df2	p_value	RSS
value	Transformwith_x	148.4284	6	730	0	198.0364
value1	OneLessPredictor	176.8579	5	731	0	198.9560
value2	drop_No_Add	205.3582	4	732	0	207.1613
value3	Originalsqrt	164.7091	5	731	0	206.7302

4.7 Outlier Detection and Removal

To detect outliers, influential and leverage points. We tested observation for standardized and studentized residual for outliers; hat for leverage points; DFFITS, DFBETAS and Cook's Distance for influential points.

These columns (211, 385, 102, 118, 225, 484, 486, 488, 533) appear under several different tests, After observing their plot, we confirm to fit the model after removing these.

Table 9: Outlier Detection and Removal

Observation	Std. Residuals ($ t_{\text{oni}} > 4$)	Studentized Residuals ($ t_{\text{student}} > 3$)	High Leverage (Hat)	High Cook's D	High DF-FITS	High DF-BETAS	Priority
211	Yes (Extreme residual)	Yes	Yes	Yes	Yes	Yes	Highest (Extreme residual + all influence metrics)

Observation	Std. Residuals ($ t_{\text{oni}} > 4$)	Studentized Residuals ($ r_{\text{student}} > 3$)	High Lever- age (Hat)	High Cook's D	High DF- FITS	High DF- BE- TAS	Priority
385	Yes (Extreme residual)	Yes	Yes	Yes	Yes	Yes	Highest (Extreme residual + all influence metrics)
102	-	Yes	Yes	Yes	Yes	Yes	Very High (Studentized residual + all influence metrics)
118	-	Yes	Yes	Yes	Yes	Yes	Very High (Studentized residual + all influence metrics)
225	-	Yes	Yes	Yes	Yes	Yes	Very High (Studentized residual + all influence metrics)
484	-	Yes	-	Yes	Yes	Yes	High (Studentized residual + influential)
486	-	Yes	-	Yes	Yes	Yes	High (Studentized residual + influential)
488	-	Yes	-	Yes	Yes	Yes	High (Studentized residual + influential)
533	-	Yes	-	Yes	Yes	Yes	High (Studentized residual + influential)

4.8 Assess Model Performance after removing outlier

AIC, BIC, R^2 , F score, RSS all improved significantly after removing outliers, so we conclude the removing outlier process.

Table 10: Assess Model Performance after removing outlier

Model	R_Squared	AIC	BIC	F_statistic	RSS
All-Transformed	0.5444	1140.41	1172.63	176.8579	198.9560
Outliers Removed	0.5947	831.53	863.66	214.3851	131.0277

4.9 QQ and Residual final comparsion

Both graphs improved with no outliers near the bound of scales.

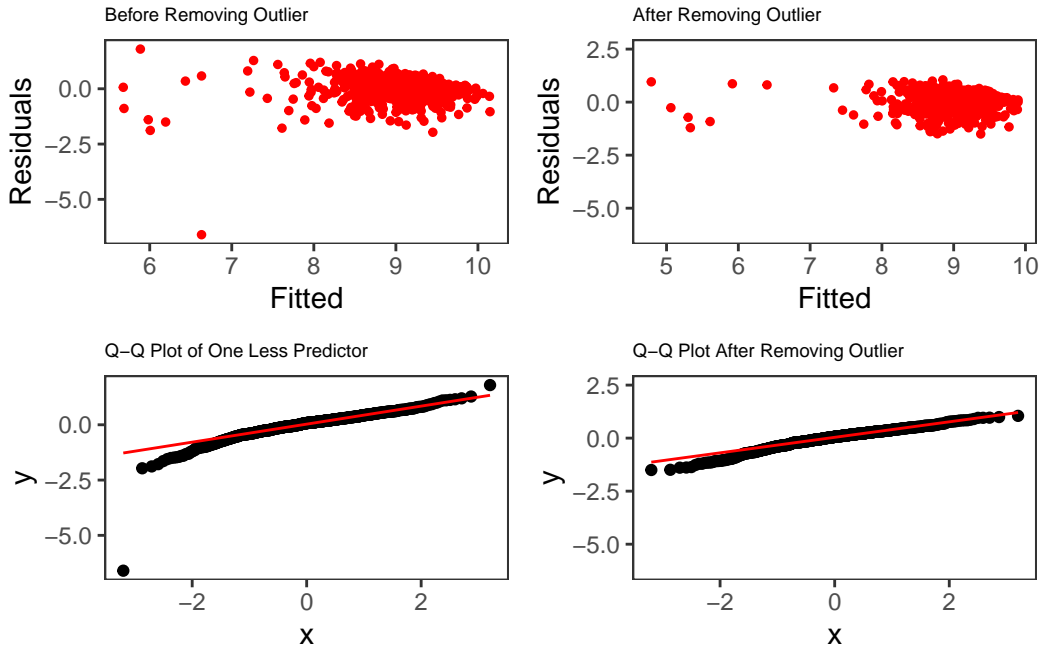


Figure 6: Residual Plots

5 Final model inference and results

Table 11: Regression Coefficients with 95% Confidence Intervals

Predictor	Estimate	Std. Error	t value	p-value	Lower 95% CI	Upper 95% CI
(Intercept)	9.9797	0.0757	131.80	< 2e-16	9.8310	10.1283
School Language - French	0.3108	0.0514	6.04	2.42e-09	0.2099	0.4118
% in Special Education	-0.0136	0.0037	-3.64	0.000292	-0.0209	-0.0062

% in Low Income	-0.0116	0.0026	-4.41	1.21e-05	-0.0168	-0.0064
% in No Educated Parent	-0.0310	0.0035	-8.89	< 2e-16	-0.0379	-0.0242
% in Special Education (Squared)	-0.0003	0.0000	-7.33	6.17e-13	-0.0004	-0.0002

5.1 Model Interpretation

The regression results presented in Table 11 provide meaningful insight into how school and family level factors influence OSSLT first attempt success rates in Ontario. Among the school characteristics, the language of instruction stands out as a significant predictor: schools offering instruction in French are associated with higher OSSLT performance, holding all other variables constant. Specifically, French-language schools are predicted to have a 0.31 unit increase in the square root of the OSSLT first attempt pass rate compared to English-language schools, with a 95% confidence interval ranging from 0.2099 to 0.4118, indicating a consistently positive and statistically significant effect.

A particularly important finding concerns the impact of special education. Rather than following a simple linear pattern, the percentage of students receiving special education services shows a nonlinear effect on OSSLT outcomes. The model includes both a linear and a squared term for this predictor. The linear term has a negative coefficient of -0.0136, while the squared term is also negative and statistically significant, with a coefficient of -0.0003. This suggests that the negative effect of special education becomes increasingly severe as the proportion of students in special education rises, indicating a compounding disadvantage in schools with especially high concentrations of these students.

In addition, two other indicators of socioeconomic disadvantage, the percentage of students from low-income households and the percentage whose parents have no post-secondary education, are both significantly associated with lower OSSLT performance. Each one-percentage-point increase in students from low-income households corresponds to a 0.0116 unit decrease in the square root of the OSSLT pass rate (95% CI: -0.0168 to -0.0064), while each additional percentage point of students whose parents lack formal post-secondary credentials is associated with a 0.0310 unit decrease (95% CI: -0.0379 to -0.0242). The narrow confidence intervals across all predictors indicate that the estimated effects are both precise and robust.

5.2 Comparing with Literature

The findings from our final regression model align closely with much of the existing literature on factors influencing student literacy outcomes. Consistent with Zhang et al. (2020), we found that both family income and parental education level are significant predictors of OSSLT success: schools with higher percentages of students from low-income households and students whose parents lack post-secondary education showed notably lower OSSLT pass rates. This supports the broader claim that socioeconomic status plays a critical role in shaping educational achievement. Similarly, our results reinforce the observations of Vaughn and Wanzek

(2014), as schools with higher proportions of students receiving special education services were significantly associated with lower OSSLT outcomes, likely due to the academic challenges these students face, even with accommodations.

However, our results diverge from the expectation presented by Bernhofer and Tonin (2022), who suggest students perform better when taught in their first language. In our analysis, schools offering instruction in French had significantly higher OSSLT performance, even though the OSSLT is administered in English. This suggests that French language instruction may be associated with school environments or educational practices that contribute positively to student literacy, despite the language difference. It may also reflect broader institutional or cultural differences between French and English schools in Ontario that warrant further exploration, such as school funding models, curriculum focus, or community engagement.

Our hypothesis regarding school type was not supported in the final model. School type, which identifies whether a school is Catholic or Public, was excluded due to a lack of statistical significance as discussed above. This outcome contrasts with the findings of Cheema (2024), who reported that private schools tend to outperform public schools in literacy achievement. While we expected Catholic schools to demonstrate higher OSSLT performance due to structured curricula or the potential influence of religious education resources, our model did not find a meaningful difference once other variables were accounted for. This suggests that variation in OSSLT performance across schools is more strongly explained by socioeconomic and instructional factors than by school type alone.

Table 12: Model Fit Statistics

Metric	Value
R-squared	0.5975
Adjusted R-squared	0.5947
AIC	831.5327
BIC	863.6648
Residual Std. Error	0.4260

5.3 Model Performance Assessment

The performance of the final multiple linear regression model, as shown in Table 12, can be evaluated using several statistical metrics that assess both goodness of fit and model parsimony. One of the most interpretable metrics is the R-squared value, which in this model is 0.5975. This indicates that approximately 60% of the variation in the square root of the OSSLT first-attempt pass rate is explained by the predictors included in the model. In the context of educational research, where student performance can be influenced by many unmeasured social, psychological, and institutional factors, an R-squared value above 0.5 is considered

relatively strong. It suggests that the model captures a substantial portion of the meaningful variance across schools in Ontario.

The Adjusted R-squared value, which accounts for the number of predictors in the model, is 0.5947. While slightly lower than the unadjusted R-squared, this is expected and confirms that the included predictors contribute meaningfully to explaining the outcome without overfitting the data. The minimal difference between the two values suggests that the model achieves a good balance between explanatory power and complexity. This strengthens confidence that the model's performance is not artificially inflated by the number of predictors used.

Beyond explanatory power, model selection criteria such as the Akaike Information Criterion (AIC), and the Bayesian Information Criterion (BIC) provide important insights into model efficiency and generalizability. AIC are measures of model fit that penalize complexity, with lower values indicating better-fitting models. In our model, the AIC is 831.5327.

The distinction between AIC and AICc lies in their intended use: AICc is a bias-corrected version of AIC that is particularly useful when the sample size is small or when the number of estimated parameters is a moderate to large fraction of the sample size. According to the rule of thumb provided by Burnham and Anderson (2004), AICc is preferred over AIC when the sample size $n \leq 40(p + 2)$ where p is the number of predictors. Based on this criterion, our dataset includes around eight hundred observations and only four predictors, which indicates that AIC is a more suitable measure for evaluating our model. This explains how only AIC is used for analysis throughout this report.

6 Discussion and conclusion

7 Author contributions

- **Junyi Hou** : Contributed to Introduction section and Model Selection section
- **Luhan Wang**: Contributed to Primary Model Results and Diagnostics section and Discussion and Conclusion section
- **Xuanle Zhou** : Contributed to the Data Description section and Final Model Inference and Results section

All team members contributed to the overall analysis, editing, and refinement of the final report.

References

- Aseery, Ahmad. 2024. "Enhancing Learners' Motivation and Engagement in Religious Education Classes at Elementary Levels." *British Journal of Religious Education* 46 (1): 43–58.
- Bernhofer, Juliana, and Mirco Tonin. 2022. "The Effect of the Language of Instruction on Academic Performance." *Labour Economics* 78: 102218.
- Burnham, Kenneth P, and David R Anderson. 2004. "Multimodel Inference: Understanding AIC and BIC in Model Selection." *Sociological Methods & Research* 33 (2): 261–304.
- Cheema, Jehanzeb Rashid. 2024. "Difference in Literacy Between Private and Public Schools: Evidence from a Survey of 61 Economies." *International Journal of Research in Education and Science* 10 (2): 218–40.
- Davis-Kean, Pamela E, Lauren A Tighe, and Nicholas E Waters. 2021. "The Role of Parent Educational Attainment in Parenting and Children's Development." *Current Directions in Psychological Science* 30 (2): 186–92.
- Nadeem, Tahir, Nasreen Akhtar, and Masood Ahmad. 2021. "A Study of the Relationship Between Family Income and Literacy Level." *STATISTICS, COMPUTING AND INTERDISCIPLINARY RESEARCH* 3 (2): 59–69.
- Ontario, Government of. 2024a. "Find Your School." <https://www.ontario.ca/page/find-your-school>.
- . 2024b. "School Information and Student Demographics." <https://data.ontario.ca/dataset/school-information-and-student-demographics/resource/e0e90bd5-d662-401a-a6d2-60d69ac89d14>.
- Vaughn, Sharon, and Jeanne Wanzek. 2014. "Intensive Interventions in Reading for Students with Reading Disabilities: Meaningful Impacts." *Learning Disabilities Research & Practice* 29 (2): 46–53.
- Zhang, Feng, Ying Jiang, Hua Ming, Yi Ren, Lei Wang, and Silin Huang. 2020. "Family Socio-Economic Status and Children's Academic Achievement: The Different Roles of Parental Academic Involvement and Subjective Social Mobility." *British Journal of Educational Psychology* 90 (3): 561–79.