

## Introduction

This project aims to investigate the relationship between the overall academic performance of high school students and their institutional setting. In this study, we analyze data from 372 high schools in New York City. The data consisted of the average SAT scores of 2014-2015 school's cohorts, along with various school and cohort's attributes, such as the school's borough and ethnicity proportion. Our intention is to identify variables that potentially affect the overall academic outcomes, and to quantify the extent that such variables have.

## Data

Data	Source	Description
score.csv	<a href="https://www.kaggle.com/nycopendata/high-schools">https://www.kaggle.com/nycopendata/high-schools</a>	Average SAT scores(Math, Reading, Writing), along with various attributes of 435 schools in NYC. The data pertained to 2014-2015 cohorts
demographics.csv	<a href="http://schools.nyc.gov/NR/rdonlyres/46093164-D8AA-40DD-A400-8F80CEBC8DD5/0/DemographicSnapshot201112to201516Public_FINAL.xlsx">http://schools.nyc.gov/NR/rdonlyres/46093164-D8AA-40DD-A400-8F80CEBC8DD5/0/DemographicSnapshot201112to201516Public_FINAL.xlsx</a>	Contains information about the gender proportion of each school
survey_2014.csv	<a href="http://schools.nyc.gov/documents/misc/2014%20Public%20Data%20File%20SUPPRESSED.xlsx">http://schools.nyc.gov/documents/misc/2014%20Public%20Data%20File%20SUPPRESSED.xlsx</a>	2014 survey result collected from parents and teachers

## Associated Variables

In this analysis, we combine data from the mentioned 3 files. We narrow down the associated features from the original sources to just 26 variables. We use R-script(mungdata.R) to perform the data munging and collect the processed result in processed\_score.csv. Noted that the total number of observation that we analyze is reduced to 372 instances because of the missing SAT score in some of the data in score.csv. The explanation of each variable in processed\_score.csv is shown in the following table.

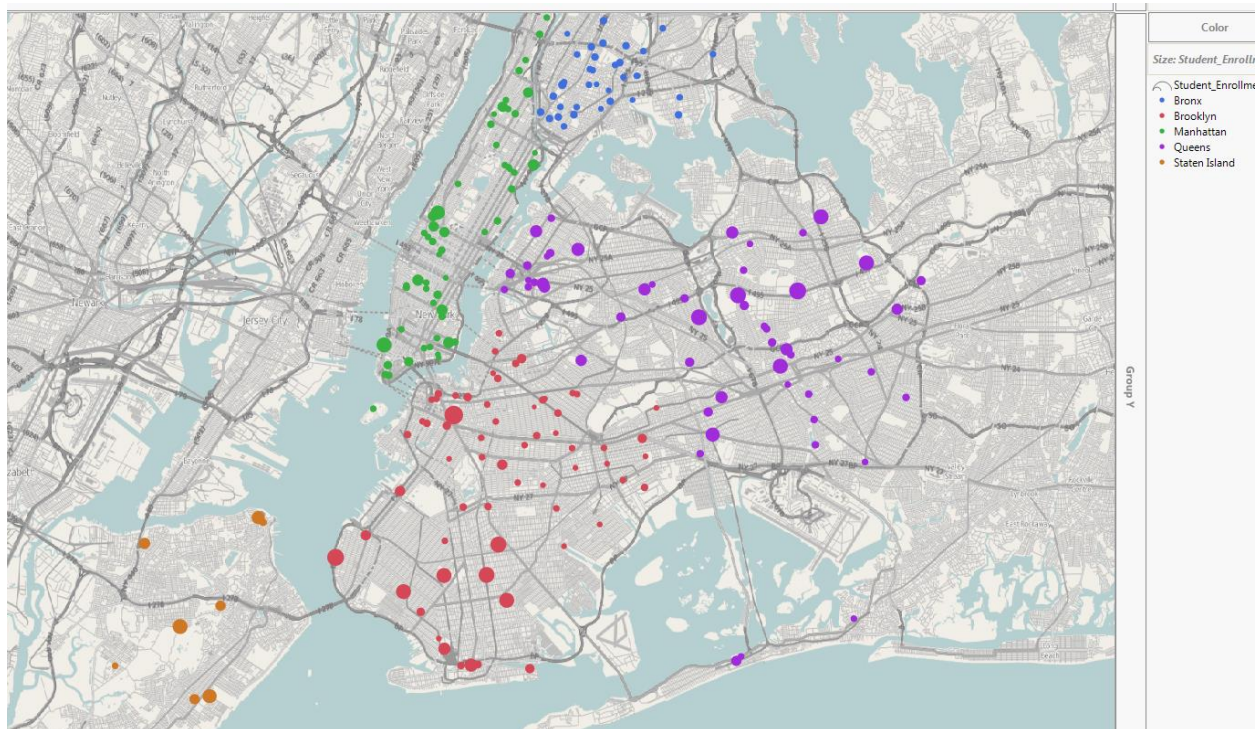
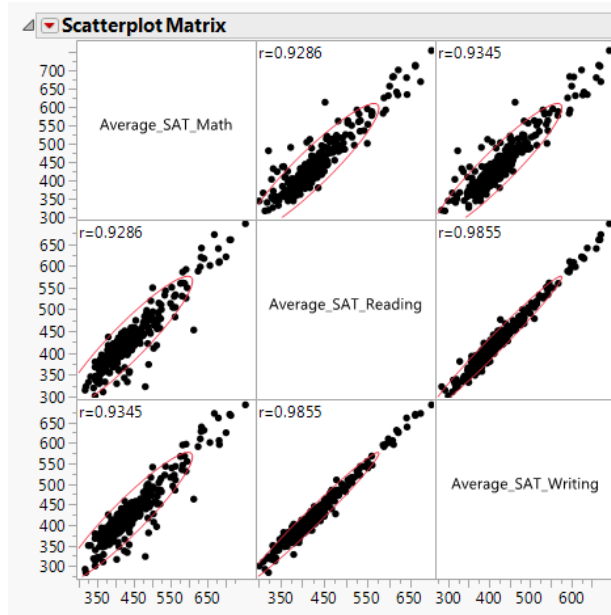
Variable Name	Description	Type
DBN	School's unique identifier	Character
Borough	School's Borough. Comprised of 5 area: Staten Island, Queens, Manhattan, Brooklyn, Bronx	Character
City	City where the school is located	Character
Latitude	School's Latitude	Numeric

Longitude	School's Longitude	Numeric
Start_Time	School's Opening hour	Numeric(e.g.: convert from 8:15 AM to 8.15)
End_Time	School's Ending hour	Numeric(e.g.: convert from 4:00 PM to 16.00)
Student_Enrollment	Number of school's enrollment	Numeric
Percent_White	%White students in 2014-2015 cohort	Numeric
Percent_Black	%Black students in 2014-2015 cohort	Numeric
Percent_Hispanic	%Hispanic students in 2014-2015 cohort	Numeric
Persent_Asian	%Asian students in 2014-2015 cohort	Numeric
*Average_SAT_Math	Average SAT Math score of 2014-2015 cohort	Numeric
*Average_SAT_Reading	Average SAT Reading score of 2014-2015 cohort	Numeric
*Averate_SAT_Writing	Average SAT Writing score of 2014-2015 cohort	Numeric
Female_Percent	%Female students in 2014-2015 cohort	Numeric
Male_Percent	%Male students in 2014-2015 cohort	Numeric
Disabilities_Percent	%Disability students in 2014-2015 cohort	Numeric
EngLearner_Percent	%English learner students in 2014-2015 cohort	Numeric
Poverty_Percent	%Poverty students in 2014-2015 cohort	Numeric
Parent_Response_Rate	Parent response rate on 2014 school's survey	Numeric
Teacher_Response_Rate	Teacher response rate on 2014 school's survey	Numeric
Instructional_Core_Satisfaction	%Response regarding instructional satisfaction	Numeric
Systems_for_Improvement_Satisfaction	%Response regarding system satisfaction	Numeric
School_Culture_Satisfaction	%Response regarding culture satisfaction	Numeric
Class_Hours	School's operating duration	Numeric(difference in hour: end_time – open_time)

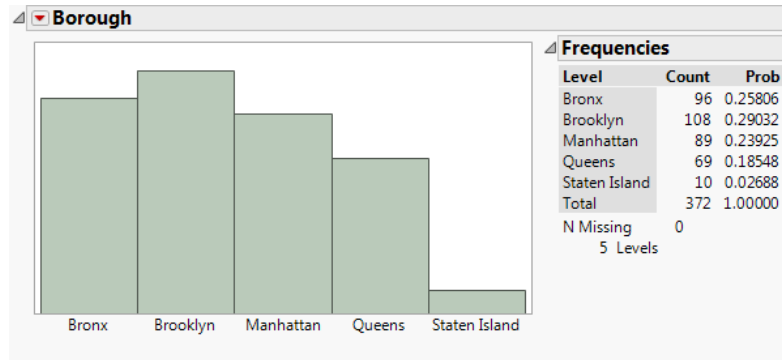
\* - Dependent variable

### Exploratory Analysis

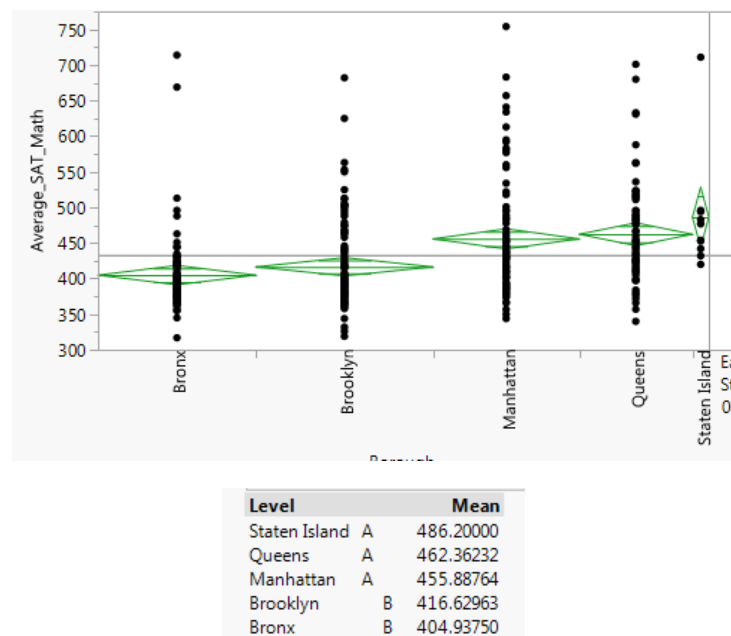
Firstly, we look at the relationship among all SAT scores. They are(unsurprisingly) highly correlated. So, in our analysis, we will put more emphasis on the SAT-Math score and later apply our findings to SAT-Reading and SAT-Writing scores.



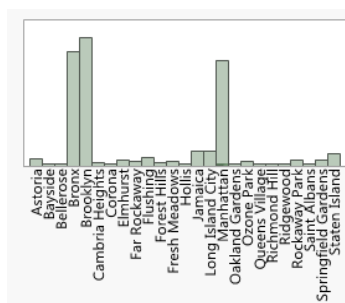
We indicate the location of each school in the above graph. The circle's size corresponds to the size of enrollment. The school's distribution in each borough can be summarized as follows:



We investigate the effect of spatial information (Borough) on SAT Math score by performing ANOVA. We find that this variable maybe helpful in predicting the SAT score.



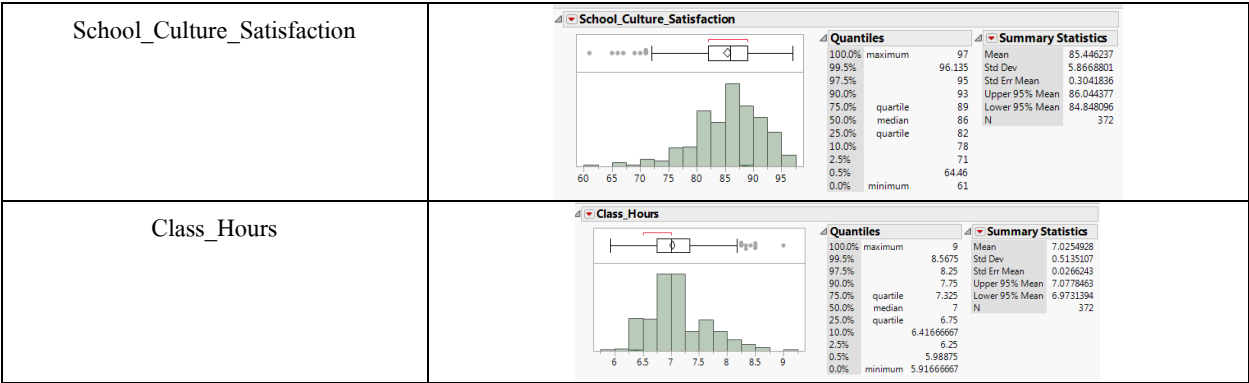
We also investigate the distribution of schools in each City but find that this variable is too fined-grained and decide to drop it as fear of running into overfitting.



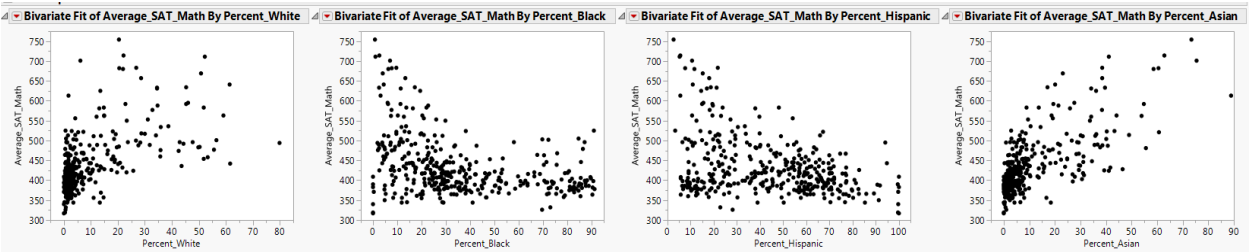
The statistics of other variables are shown in the following table.

Variable	Statistics
Start_Time	<div><div><div><div><div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div><div></div></div></div></div></div></div>





For preliminary analysis, the scatter plot between ethnicity proportion and SAT-Math score shows some predictive power and indicates that these variables should be included in the model.

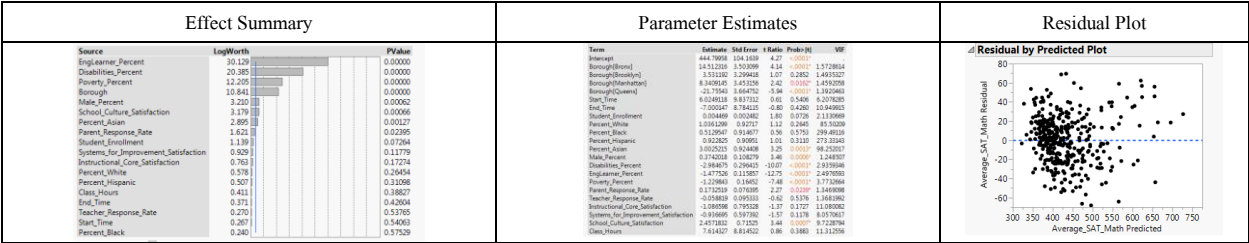


Project Objective

Uncover the relationship between the academic setting and cohort’s academic outcomes(as measured by the average SAT scores). We also quantify the variable’s effect on 3 different sections of SAT exams(Math, Reading, Writing) to determine their predictive power on each section. Linear Regression is chosen as our base model to fit the data on because of its simplicity and interpretability.

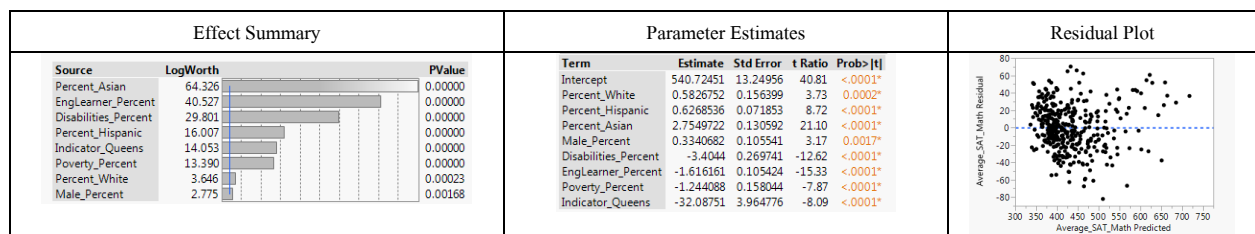
Model Development

As stated earlier, our model will be developed based on using Average SAT-Math score as a dependent variable. The initial model building consisted of all independent variables (exclude DBN, City, Latitudes, Longitudes, and Female\_Percent (as this variable is reflected in Male\_Percent)). We obtain the following model:

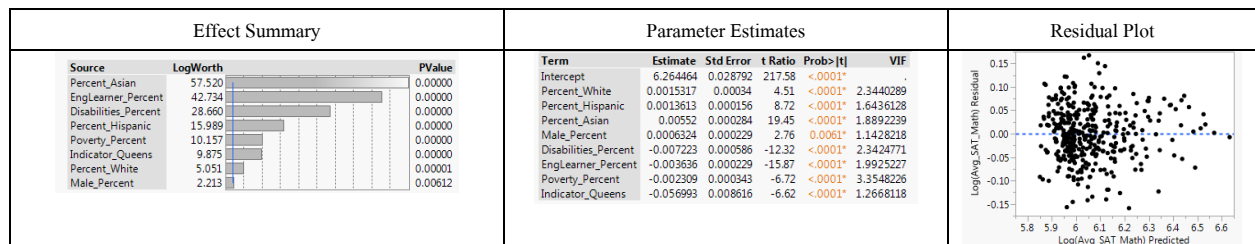


Although the base model shows a strong predictive power (RSquare = 0.87 and RMSE = 26), it has many undesirable properties; the model contains many variables that are not statistically significant, some independent variables are highly correlated (as shown by VIF), and the Residual Plot shows Heteroscedasticity problem (verified by Park-Test). To attenuate these effects, we perform a series of model development, which can be summarized as follows:

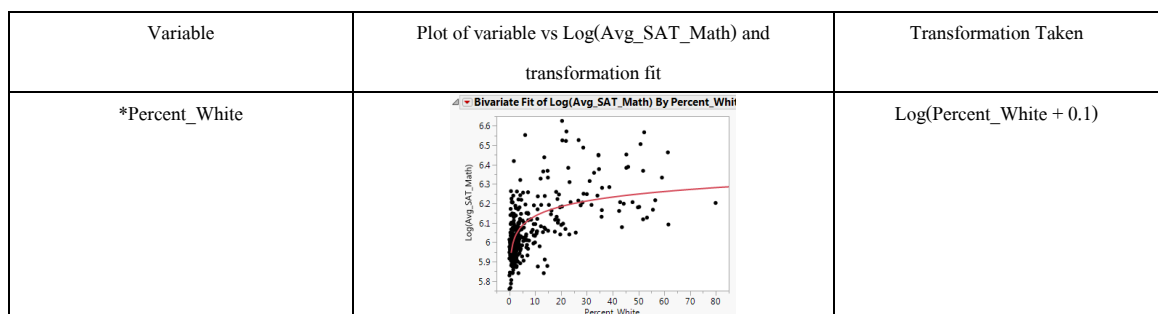
1. Manually create a dummy variable based on Borough. As opposed to the one generated by JMP, this will allow us to remove an individual borough that we found not significant. Bronx is treated at the base level since it has the lowest Average SAT-Math score means.
2. Re-fit the model. Iteratively remove variables with high VIF and P-Value that exceeds 0.01 significant threshold.



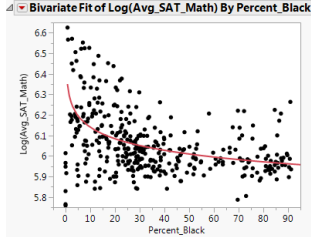
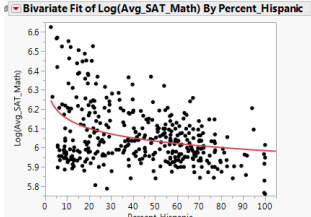
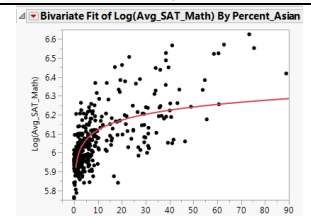
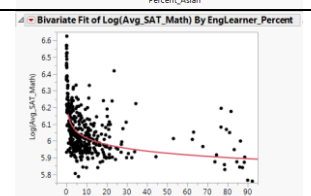
3. As Heteroscedasticity is still presented, we apply Log transformation to Average SAT-Math score and re-fit the model. Drop any unnecessary variable as stated in 2)



4. We try to eliminate Heteroscedasticity by plotting every independent variable against Log(Avg SAT Math) and transform them appropriately if we think that leads to a more linear relationship.

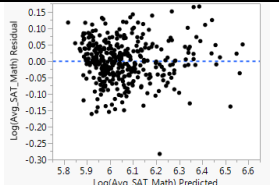




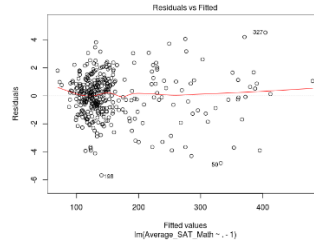
*Percent_Black		$\text{Log}(\text{Percent\_Black} + 0.1)$
*Percent_Hispanic		$\text{Log}(\text{Percent\_Hispanic} + 0.1)$
*Percent_Asian		$\text{Log}(\text{Percent\_Asian} + 0.1)$
*EngLearner_Percent		$\text{Log}(\text{EngLearner\_Percent})$

\* Add 0.1 to the original value before applying the Log transformation because some instances have 0 value

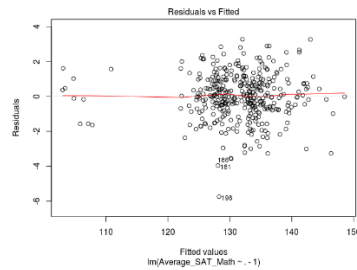
- Re-fit the model using the transformed variables. Drop any unnecessary variable as stated in 2)

Effect Summary			Parameter Estimates						Residual Plot
Source	LogWorth	PValue	Term	Estimate	Std Error	t Ratio	Prob> t	VIF	
Log(Percent_Asian_adj)	26.680	0.00000	Intercept	6.1194874	0.06625	92.37	<.0001*	.	
Log(EngLearner_Percent_adj)	21.873	0.00000	Indicator_Queens	-0.040766	0.009608	-4.22	<.0001*	1.3138748	
Student_Enrollment	12.281	0.00000	Indicator Staten Island	-0.058806	0.022263	-2.64	0.0086*	1.206477	
Log(Percent_Black_adj)	7.769	0.00000	Student_Enrollment	4.0076e-5	5.348e-6	7.49	<.0001*	1.6645134	
School_Culture_Satisfaction	7.332	0.00000	Disabilities_Percent	-0.001553	0.00059	-2.63	0.0086*	1.9542378	
Systems_for_Improvement_Satisfaction	6.633	0.00000	Poverty_Percent	-0.001524	0.000327	-4.66	<.0001*	2.5098501	
Poverty_Percent	5.344	0.00000	Systems_for_Improvement_Satisfaction	-0.005927	0.001124	-5.27	<.0001*	4.7956064	
Indicator_Queens	4.504	0.00003	School_Culture_Satisfaction	0.0075669	0.001355	5.58	<.0001*	5.8679113	
Indicator Staten Island	2.065	0.00862	Log(Percent_Black_adj)	-0.022268	0.003859	-5.77	<.0001*	1.8124027	
Disabilities_Percent	2.054	0.00884	Log(EngLearner_Percent_adj)	-0.037004	0.003532	-10.48	<.0001*	2.5515644	
			Log(Percent_Asian_adj)	0.0377099	0.003195	11.80	<.0001*	1.9398177	

- We have lessened the effect of Heteroscedasticity but the issue is still presented. Now, we turn to Weighted Least Squares Regression approach. We switch to using R to conduct the analysis at this point (as performing the analysis in JMP can be quite tedious). The analysis code can be found in analysis.R . We use the transformed data collected from step 4). As the observations come from aggregated result, we firstly try to weight the data by the enrollment size (that is, multiply every variable by  $\sqrt{\text{Student\_Enrollment}}$ ) and fit the regression model with no intercept. The residual plot indicates that Heteroscedasticity is still presented.



- Now, we try Weighted Least Squares Regression with two-stage approach; firstly, fit the regression model using transformed variables in step 4) and then use the mean square residual of each borough as a weight for WLS. The residual plot indicate that Heteroscedasticity issue is now fixed.



- We import the transformed data back to JMP(weighted\_score.csv) and drop any unnecessary variable as stated in 2).

We obtain the following model:

Effect Summary			Parameter Estimates				Residual Plot
Source	LogWorth	PValue	Term	Estimate	Std Error	t Ratio	Prob >  t
1/Weighted	253.427	0.00000	Student_Enrollment/Weighted	0.0000395	5.377e-6	7.34	<.0001*
Log(Percent_Asian_adj)/Weighted	26.231	0.00000	Log(Percent_Black_adj)/Weighted	-0.021531	0.003834	-5.62	<.0001*
Log(EngLearner_Percent_adj)/Weighted	21.764	0.00000	Log(Percent_Asian_adj)/Weighted	0.0348968	0.003159	11.08	<.0001*
Student_Enrollment/Weighted	11.858	0.00000	Disabilities_Percent/Weighted	-0.00174	0.000577	-3.02	0.0027*
School_Culture_Satisfaction/Weighted	7.881	0.00000	Log(EngLearner_Percent_adj)/Weighted	-0.036992	0.003542	-10.44	<.0001*
Systems_for_ImprovementSatisfaction/Weighted	7.493	0.00000	Poverty_Percent/Weighted	-0.001402	0.000324	-4.33	<.0001*
Log(Percent_Black_adj)/Weighted	7.407	0.00000	Systems_for_ImprovementSatisfaction/Weighted	-0.006209	0.001116	-5.55	<.0001*
Poverty_Percent/Weighted	4.719	0.00002	School_Culture_Satisfaction/Weighted	0.0078705	0.001353	5.82	<.0001*
Queens/Weighted	3.817	0.00015	1/Weighted	6.1170387	0.06605	92.61	<.0001*
Disabilities_Percent/Weighted	2.561	0.00274	Queens/Weighted	-0.037877	0.009895	-3.83	0.0002*

We then perform Park-Test to verify for Heteroscedasticity. The P-Value for  $\hat{Y}$  when regressed on  $r^2$  is 0.75 and when regressed on  $\log(r^2)$  is 0.68. This indicates that the Heteroscedasticity is now fixed.

The model in step 8) is our selected model. The model equation is

$$\frac{\text{Log}(\text{AvgSatMath})}{\text{Weight}} = 6.1 * \frac{1}{\text{Weighted}} + 0.000039 * \frac{\text{StudentEnrollment}}{\text{Weighted}} - 0.02 * \frac{\text{Log}(\text{PercentBlack})}{\text{Weighted}} + 0.03 * \frac{\text{Log}(\text{PercentAsian})}{\text{Weighted}} - 0.001 * \frac{\text{DisabilityPercent}}{\text{Weighted}} - 0.03 * \frac{\text{Log}(\text{EngLearnerPercent})}{\text{Weighted}} - 0.001 * \frac{\text{PovertyPercent}}{\text{Weighted}} - 0.006 * \frac{\text{SystemForImprovementSatisfaction}}{\text{Weighted}} + 0.007 * \frac{\text{SchoolCultureSatisfaction}}{\text{Weighted}} - 0.03 * \frac{\text{Queens}}{\text{Weighted}}$$

Where Weighted is the means of Residual Square per group and has the following values:

Group	Weighted
Bronx	0.0463926118015564

Brooklyn	0.0447919246176774
Manhattan	0.0446127098917998
Queens	0.0484199559569425
Staten Island	0.0584236025734739

The equation can be simplified to

$\begin{aligned} \text{Log(AvgSatMath)} = & 6.1 * 0.000039 * \text{StudentEmrollment} - 0.02 * \text{Log(PercentBlack)} + 0.03 * \text{Log(PercentAsian)} - 0.01 \\ & * \text{DisabilityPercent} - 0.03 * \text{Log(EngLearnerPercent)} - 0.001 * \text{PovertyPercent} - 0.006 \\ & * \text{SystemForImprovementSatisfaction} + 0.007 * \text{SchoolCultureSatisfaction} - 0.03 * \text{Queens} \end{aligned}$
---