

# Exploring Racial Demographics and SAT Performance: A Borough-Level Study in New York City

HANNAH YU

April 11, 2024

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Literature Review . . . . .	2
1.2	Research Contribution . . . . .	3
<b>2</b>	<b>Data Analysis</b>	<b>4</b>
2.1	Summary Statistics . . . . .	5
2.2	Visualization . . . . .	7
<b>3</b>	<b>Geographic Analysis</b>	<b>11</b>
<b>4</b>	<b>Model</b>	<b>17</b>
<b>5</b>	<b>Regression Results</b>	<b>20</b>
<b>6</b>	<b>Machine Learning</b>	<b>28</b>
6.1	Decision Tree . . . . .	28
6.2	Random Forest Model . . . . .	31
<b>7</b>	<b>Conclusion</b>	<b>32</b>

---

# 1 Introduction

The impact of student racial demographics on academic achievements is widely acknowledged in the education industry. Examining factors affecting learning outcomes is crucial in fostering an inclusive learning environment that encourages students' personal growth and development. Understanding the education disparities among students with various racial backgrounds is essential for policymakers in developing regulations supporting education equality. With such an idea in mind, this study examines the relationship between SAT scores and student racial demographics in New York City.

## 1.1 Literature Review

Despite New York City being one of the most diverse cities in the US, the city still faces challenges in stopping racial segregation and providing equal access to quality education (Bonastia, 2023). According to the research by (Card and Rothstein, 2007), highly segregated cities increase the black-and-white test score gap. Additionally, (Penney, 2017), noted that the shortage of minority teachers may contribute to the difference in academic performance as students as students perform better with teachers of the same race.

To avoid confounding factors and offer a more comprehensive view of the education disparities in New York City, this study also examines factors relating to students' racial demographics that may affect their SAT performances. Studies have found public school quality tends to parallel the economic status of their neighbourhood, highlighting the interplay between socioeconomic factors and educational outcomes (Kemple, Farley and Stewart, 2019). Students attending schools in wealthier neighbourhoods are more likely to have wealthier family backgrounds. According to Everson and Millsap (2004), students' family backgrounds directly and indirectly affect their SAT performances. As noted by Battle and Lewis (2002) in their studies, socioeconomic factors are also big contributors to differences in students' academic performance by race. African American students typically are outperformed by their white students counterparts. However, with socioeconomic status being controlled, minorities perform better than their white peers (Battle and Lewis, 2002). With these established relationships in previous literature, my paper also examines the effect of the school's socioeconomic status, measured by the average housing prices of the area.

A school's overall academic and social environment also has a great impact on its students' academic performance. Defining a school's climate by various dimensions of the school environment including academic expectations, communications, engagements, safety and respect, Davis and Warner

---

(2015) found a positive association between school climate and students' academic performance in New York City. In some conditions, the school climate's effect even outweighs the effect of students' backgrounds. Researchers have found consistent associations between higher academic achievements and lower dropout rates (Fetler, 1989). Similarly, giving students more opportunities for challenges and initiatives may improve their grades. Studies have proven that participating in Advanced Placement (AP) classes helps students score higher in the SATs: students with similar initial achievements who took AP classes score slightly higher than those who did not (McKillip and Rawls, 2013). In addition, partaking in extracurricular activities boosts students' SAT scores, especially among minorities and socioeconomically disadvantaged students (Everson and Millsap, 2005). With previous research proving the impact of school environment on academic achievements, this study explores a variety of school characteristics including enrollment statistics, and academic and extracurricular programs offered.

## **1.2 Research Contribution**

Using comprehensive information about each school from the dataset, this research aims to unravel the relationship between racial demographics and academic success across different boroughs in New York City. The analysis reveals that higher test scores are more likely to be associated with white or Asian students. Unlike much literature that focuses on the full SAT score, this research also individually studies the impact of these related factors on individual SAT subsection scores: Math, Reading, and Writing, each worth 800 points. Notably, I find that schools with more Asian populations tend to score higher in SAT Math. On the other hand, schools with more white students prevail in both Reading and Writing.

From the maps plotted in the study, it is evident that education performances also differ by area, and the differences can both be attributed to economic status measured by housing prices, students' enrollment rates, and the students' racial identities. However, from the OLS regression results, individual boroughs do not directly impact the school's average SAT scores. The difference in scores across boroughs is more likely to be attributed to differences in education resources. The analyses regarding school programs, including AP, language, and extracurricular offerings match the findings in previous studies. High schools providing more challenging classes and more diverse extracurriculars are linked with higher student performance on the SATs. This research seeks to shed light on the long-lasting issue of education disparities existing among students with various backgrounds. Given

---

that New York City’s public school system is the largest by far across the United States, the findings may also have broader implications for education systems across the nation.

The remainder of the paper is structured into several sections. Section 2.1 displays the data after cleaning and loading. Section 2.2 presents plots, histograms, and figures related to the research question. Section 3 provides a geographic analysis of the research topic. The regression models are presented in section 4 and their results are in section 5. To bring more evidence to the findings, section 6 includes two machine learning models and their results. Finally, the paper concludes at section 7.

## 2 Data Analysis

This paper utilizes the dataset "Average SAT Scores for NYC Public Schools" published by NYC Open Data and College Board compiled from Kaggle (*Average SAT Scores for NYC Public Schools*, 2014). This rich dataset encompasses details of every accredited public high school in New York City of the 2014-2015 school year, including school information, demographic breakdown, enrollment numbers, and average SAT scores.

In addition to the data retrieved from Kaggle, additional school information from the 2014-2015 DOE High School Directory dataset from NYC Open Data provided by the Department of Education is added (NYC Open Data, 2014). Abundant school characteristics such as course and extracurricular offerings facilitate a more comprehensive examination of the relationship between students’ backgrounds and SAT scores. Lastly information on New York City property prices was retrieved from Kaggle as well to observe the socioeconomic differences of different areas (NYC Property Sales, n.d.).

The raw SAT score dataset contains 435 observations each representing one public high school in New York City. In the dataset, several observations contain missing values for both student racial breakdown and SAT score results. After removing the 61 observations with missing values and performing other basic data-cleaning procedures on the remaining 374 observations, I created several new variables in aid of subsequent analysis. *TotalAverageScore* was generated by aggregating the sum of Average SAT Math, Reading, and Writing into a single score, on a 2400-point scale. Combining these subscores into a single score simplifies the data analysis. Additionally, four demographic variables representing the percentage (scaled out of 100) of the students in each school that belong to specific racial categories (*White*, *Black*, *Hispanic*, and *Asian*) were created.

---

## 2.1 Summary Statistics

In my endeavour to assess disparities that affect academic achievement, I selected the dependent variable *TotalAverageScore* as the metric for student performance. My chosen independent variables in representing student demographics are *White*, *Black*, *Hispanic*, and *Asian*. These four groups are the predominant population in the United States, and together usually comprise more than 95% of the student population in each school. In addition, to investigate the relationship between school size and academic outcomes, I also included the variable *StudentEnrollment* in the analysis.

The *TotalAverageScore* represents the mean SAT score of the students who took the SAT during the 2014-2015 school year at each school in the dataset. This score is based on the old 2400-point scale, where the SAT Reading, Writing, and Math each contribute 800 points. The average SAT score serves as a reliable indicator of academic performance because of its standardized nature and its reflection on the educational outcomes of students.

Each of the four variables represents the percentage of students belonging to one of the four main racial demographics in New York City schools (e.g. *White* at 34 means 34% of students in that school are white). Accurately identifying the racial composition of each school is crucial for my investigation into the impact of student demographics on academic achievement. Various studies have shown correlations between students' race and their academic performance in different contexts. Therefore, analyzing the relationship between these demographic variables and SAT scores is a critical initial step in understanding educational disparities.

The variable *StudentEnrollment* represents the total number of students enrolled in each school during the same time period. I seek to represent school size with this variable and explore how school size correlates with SAT scores. Factors that influence education outcomes such as resource allocation, class sizes, and overall learning environment may all be correlated with the size of the school. Therefore, including this variable in my analysis is essential for identifying disparities in the educational composite.

	Student Enrollment	White (%)	Black (%)	Hispanic (%)	Asian (%)	Average SAT Score (out of 2400)
<b>count</b>	374	374.0	374.0	374.0	374.0	374
<b>mean</b>	756	8.5	35.4	43.9	10.4	1275
<b>std</b>	774	13.4	25.4	24.5	14.4	195
<b>min</b>	142	0.0	0.0	2.6	0.0	924
<b>25%</b>	397	1.3	16.4	20.8	1.6	1157
<b>50%</b>	482	2.6	28.8	45.3	4.2	1226
<b>75%</b>	660	9.4	50.1	63.4	11.1	1327
<b>max</b>	5447	79.9	91.2	100.0	88.9	2144

Table 1: Summary Statistics

Table 1 provides a comprehensive overview of the main variables of interest. There are in total 374 observations, or schools, in the dataset, which is represented by "count" in the table.

With a mean of 756 students per school and a standard deviation of 774, I observe a wide range of enrollment sizes, ranging from 142 to 5447 students. This wide range indicates significant variability in school size across the dataset, which has a potential impact on average SAT scores.

With mean percentages ranging from 8.5% for white students to 24.5% for Hispanic students, the racial composition of schools exhibits considerable diversity. The standard deviations ranging from 13.4% for White students to 25.4% for black students highlight considerable diversity in racial demographics among schools. The large standard deviation of 25.4% for black students suggests more pronounced variability in the proportion of Black students across schools compared to other racial groups.

While New York City is comprised of 30% white individuals, the mean percentage of white students in public schools is the lowest. This scenario can be explained by a high portion of white students choosing to attend private schools as the US Department of Education reports that while only 14% of schools in New York City are private, 40% of its students are white (Di et al., 2021). Similar observations can be made toward Asian students as well: while 16% of public school students are Asian, they represent 62% of the students in specialized private schools (Harris and Hu, 2018). The diverse racial composition of schools underscores the importance of considering demographic factors in analyzing SAT scores. Racial demographic variabilities suggest underlying differences in socioeconomic backgrounds, educational resources, and cultural influences that may affect academic achievement.

The wide range of average SAT scores from 924 to 2144 highlights disparities in academic achieve-

ment across schools. While the highest average score reaches 2144, the mean score of 1275 falls significantly lower. This distinction prompts us to question the factors contributing to such a gap. This paper attempts to address this question by examining various factors, including student racial demographics, school location, and school size.

## 2.2 Visualization

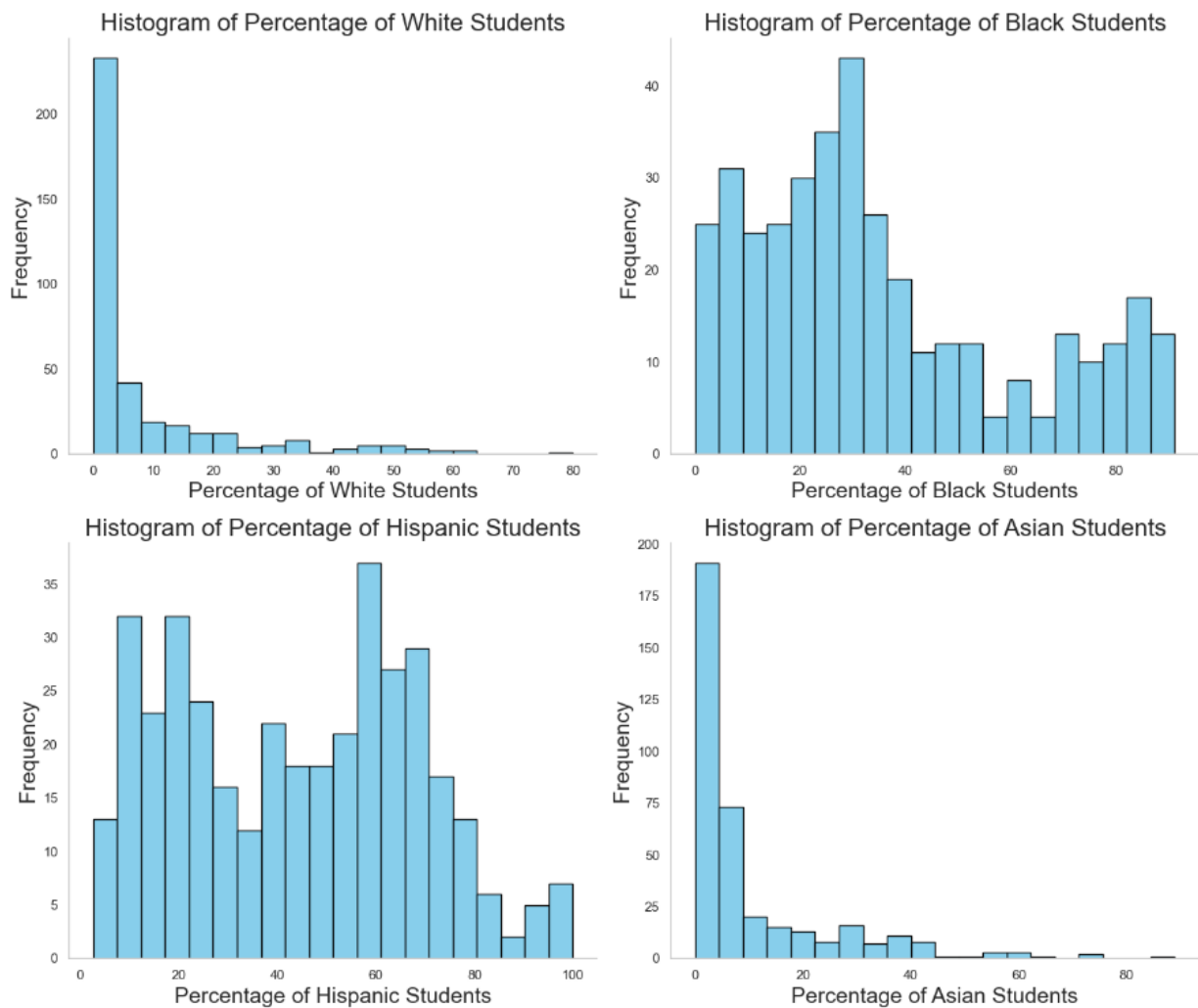


Figure 1: Distribution of Student Racial Demographics

The histograms from Figure 1 provide a visual representation of the distribution of students within each racial category across the dataset. From observing the graphs, I notice that both *White* and *Asian* exhibit similar distribution shapes, with the majority of schools having less than 20% representation of these groups. This trend confirms the preference for private schools over public schools among white and Asian communities.

In contrast, the distribution of black and Hispanic students appears more evenly spread and fol-

---

lows a symmetrical, bell-shaped curve. This observation suggests that black and Hispanic students are more widely represented across schools, with a notable proportion of schools having moderate to high percentages of these groups.

These patterns highlight potential differences in representation across different racial groups as some groups are more prevalent than others. Further exploration of educational access and outcomes among racial groups needs to take the representation differences into account.

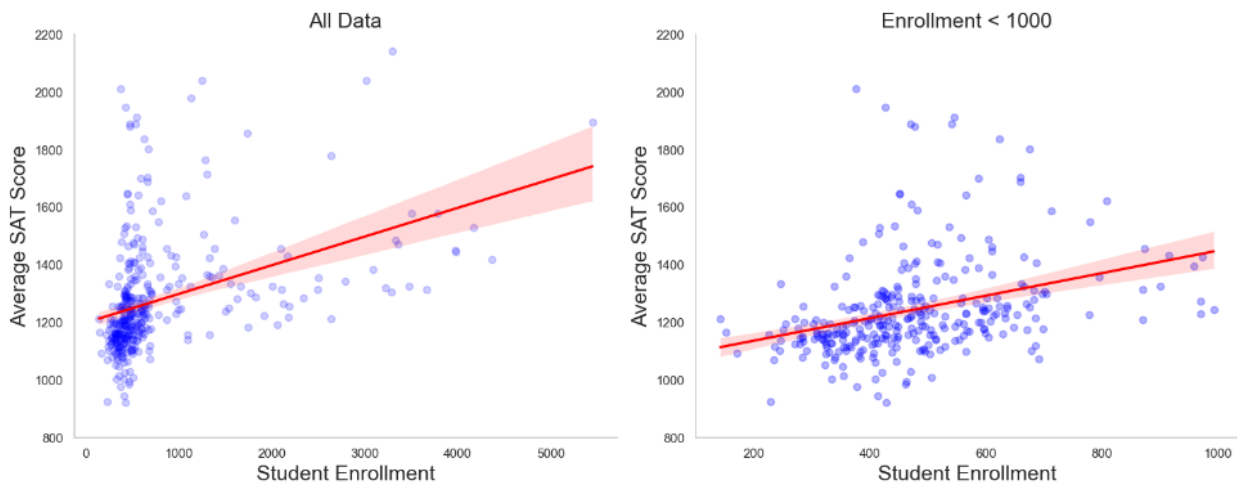


Figure 2: Relationship between Student Enrollment and Average SAT Score

The scatter plots from Figure 2 depicts the relationship between student enrollment and SAT scores. Since most school sizes are less than 1000 students, I created another plot focusing on schools with less than or equal to 1000 students. Both plots exhibit positive relations between student enrollment and average SAT scores. There is a general increase in average SAT scores when the school size is larger. The outcome may be attributed to larger schools having better resources such as well-equipped libraries, advanced laboratories, and specialized teachers. Larger schools also can offer a wider range of courses to allow students to expand their knowledge and exercise their critical thinking abilities. Students in larger schools with exposure to these extra resources and opportunities are more likely to have better academic performances.



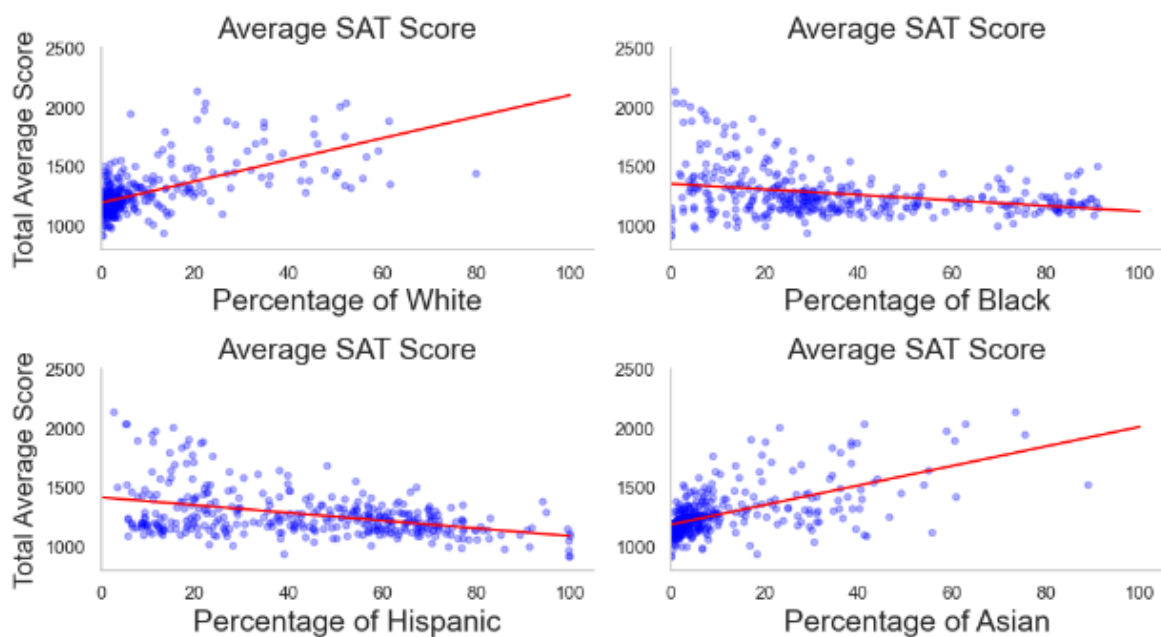


Figure 3: Relationship between Student Race and Average SAT Scores

Figure 3 shows that the trend observed in previous analyses persists: white and Asian students tend to exhibit a similar pattern, as do black and Hispanic students. In these scatter plots, I observe a positive association between these racial groups and academic performance, suggesting schools with a higher proportion of white or Asian students tend to have higher average SAT scores. As indicated by the negative correlation, schools with a higher proportion of black or Hispanic students tend to display lower average SAT scores.

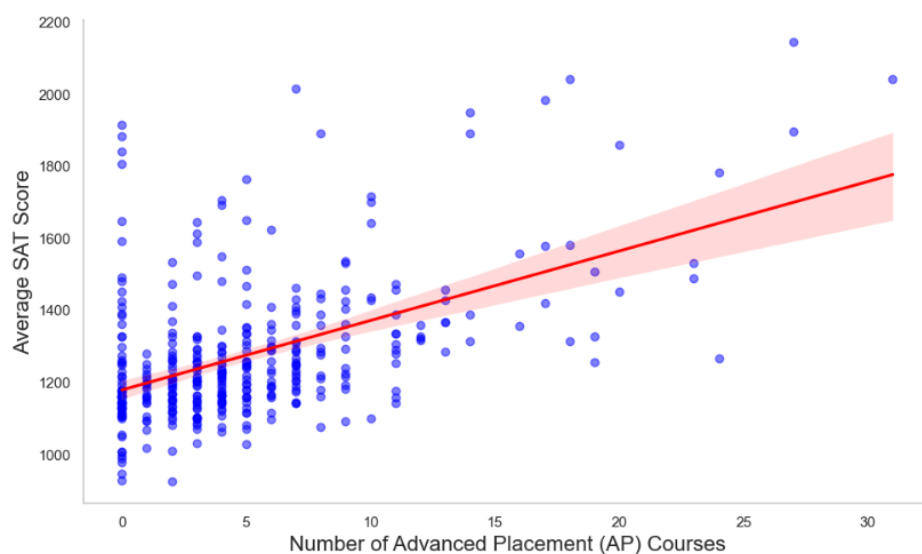


Figure 4: Number of AP Courses vs. Average SAT Score

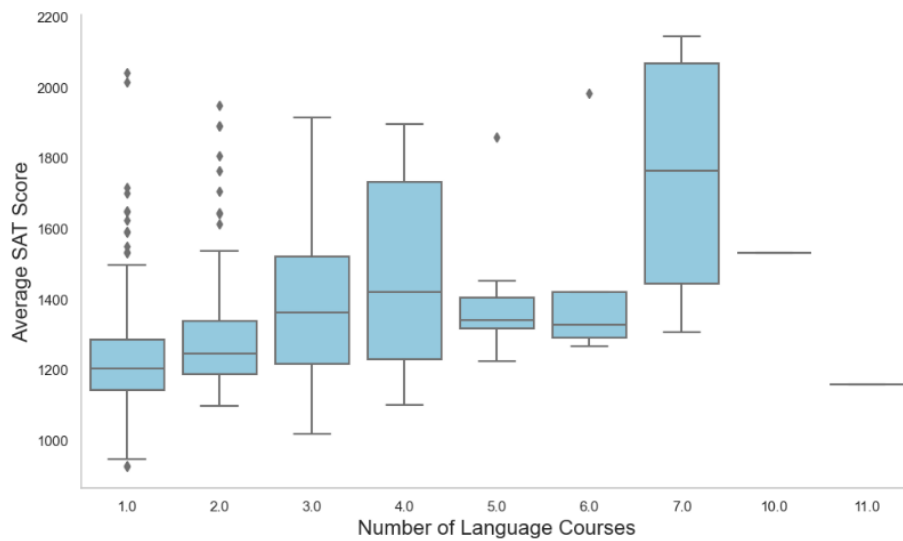


Figure 5: Number of Language Courses vs. Average SAT Score

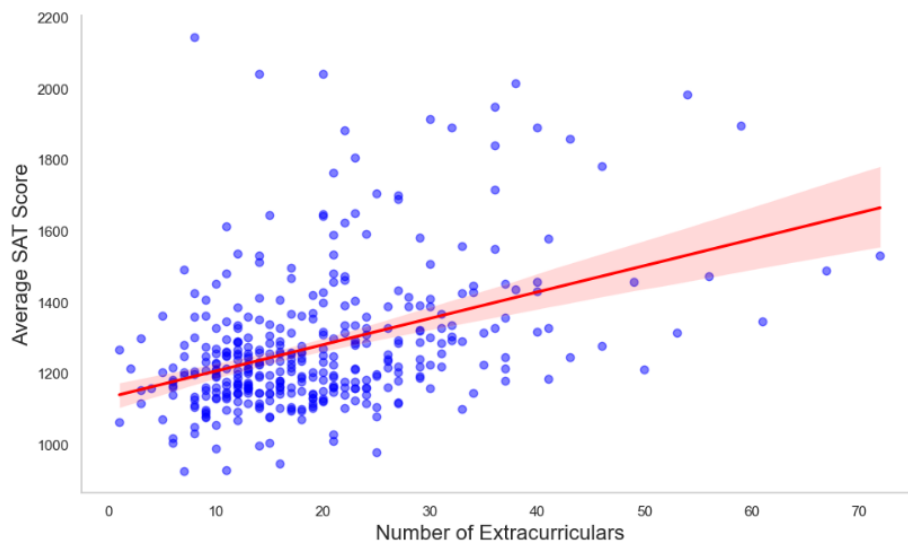


Figure 6: Number of Extracurricular Activities vs. Average SAT Score

Figure 4, 5, 6 showcase the relationship between average SAT scores (by school) and the number of AP courses, language courses, and extracurricular activities offered at school, respectively. I observe an upward positive trend for all three of the graphs, indicating that more student opportunities, both academic and extracurricular are all associated with higher SAT scores.

Advanced Placement (AP) classes are first-year college-level classes offered in high school for students seeking academic challenges and aiming to enhance their college applications. Taking many AP classes suggests the student's readiness for college education. High schools offering many AP classes imply that their students are more prone to take many and prepare for college. Their initiative in taking advanced courses probably would also make them score higher on the SATs to boost their

chance of admission to better universities. Similarly, taking multiple language classes also reflects the student’s willingness to engage in academic rigour and implies their possible commitment to making similar efforts to achieve high SAT scores. Therefore, schools with more AP and language course offerings are associated with higher SAT scores.

Extracurricular participation is another significant part of college applications. Students aiming for top universities dedicate considerable time outside of class to various extracurricular activities to showcase their interests and strengths. Therefore, schools offering extensive extracurricular opportunities attract more students aiming for reputable colleges, who are also motivated to achieve high SAT scores. Moreover, schools capable of providing these diverse programs and opportunities are often well-funded, indicating their ability to hire quality teachers and provide resources conducive to academic success.

### 3 Geographic Analysis

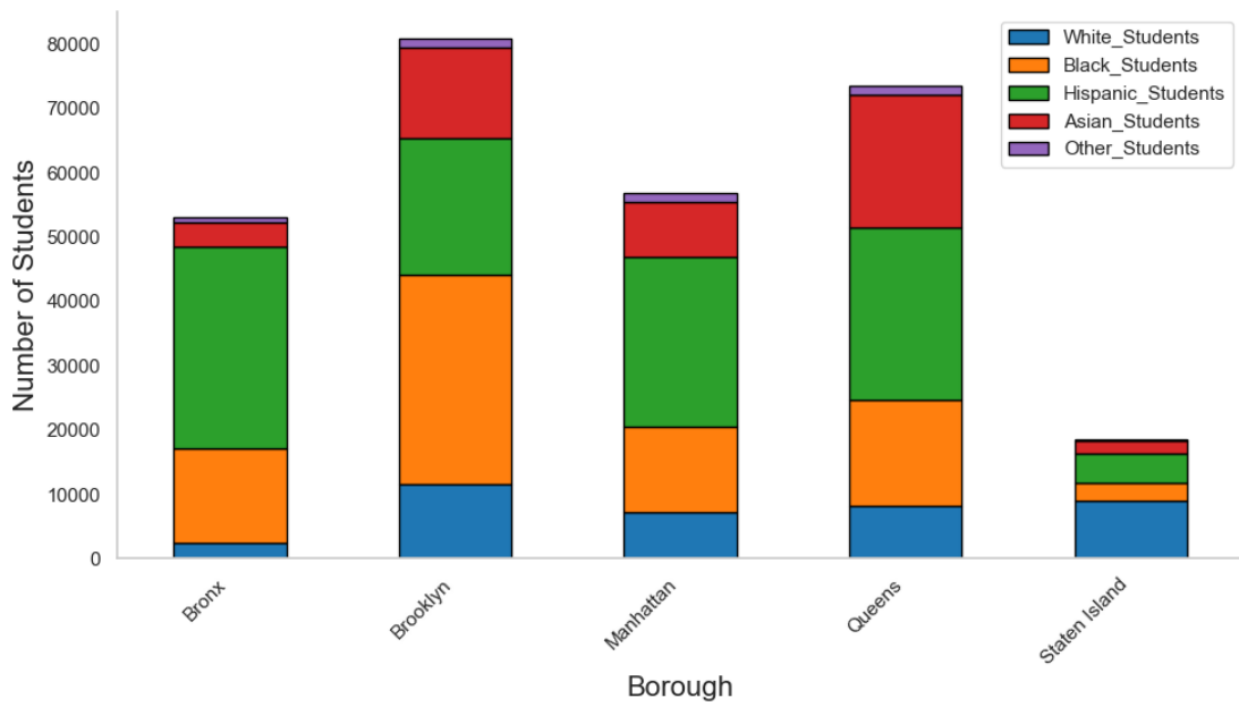


Figure 7: Number of Students by Borough and Race

New York City comprises five governing boroughs: the Bronx, Brooklyn, Manhattan, Queens, and Staten Island (Boroughs of New York City, 2020). Each borough exhibits distinct socioeconomic status, with Manhattan typically leading in income and economic wealth, wealthier areas often possess more resources to invest in public school education. Studies have consistently demonstrated

that students with wealthy backgrounds tend to have a significant advantage in scoring higher on standardized tests such as the SAT (Bonastia, 2023). Therefore, conducting an extensive analysis of education disparities needs to consider the geographical location of schools.

Figure 7 provides a breakdown of the number of students by borough. White students are more prevalent in Brooklyn and less so in the Bronx. Even though Staten Island has the smallest amount of student population, almost half of its students are white. Brooklyn is also the borough with the most amount of black students, while Staten Island has the least number of black students. Hispanic students show a slight preference for the Bronx, while similar black students, dislike Staten Island. Asians on the other hand take over Queens and are the least prevalent in Staten Island.

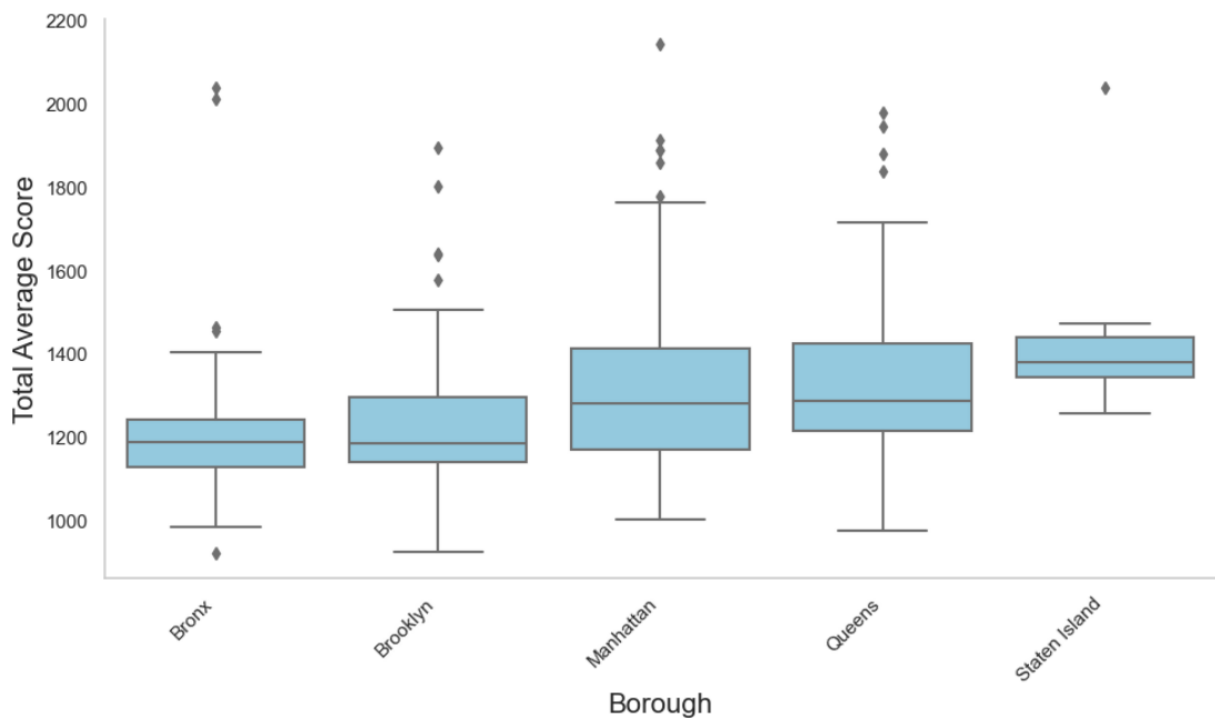


Figure 8: Number of Students by Borough and Race

Figure 8 shows the distribution of SAT scores across different boroughs. Given the variation in education resources across areas, I aim to observe the direct differences in educational outcomes among boroughs. From the plot, it becomes evident that one Manhattan school exhibits the highest average SAT scores, while Bronx school records the lowest. This finding matches the previous analysis where areas with more white and Asian students are prone to higher scores and areas with more black and Hispanic students with lower scores. Notably, Staten Island has the highest median score among all boroughs. This result may be attributed to the high portion of white students in that location, highlighting intriguing disparities in educational achievement across New York City.

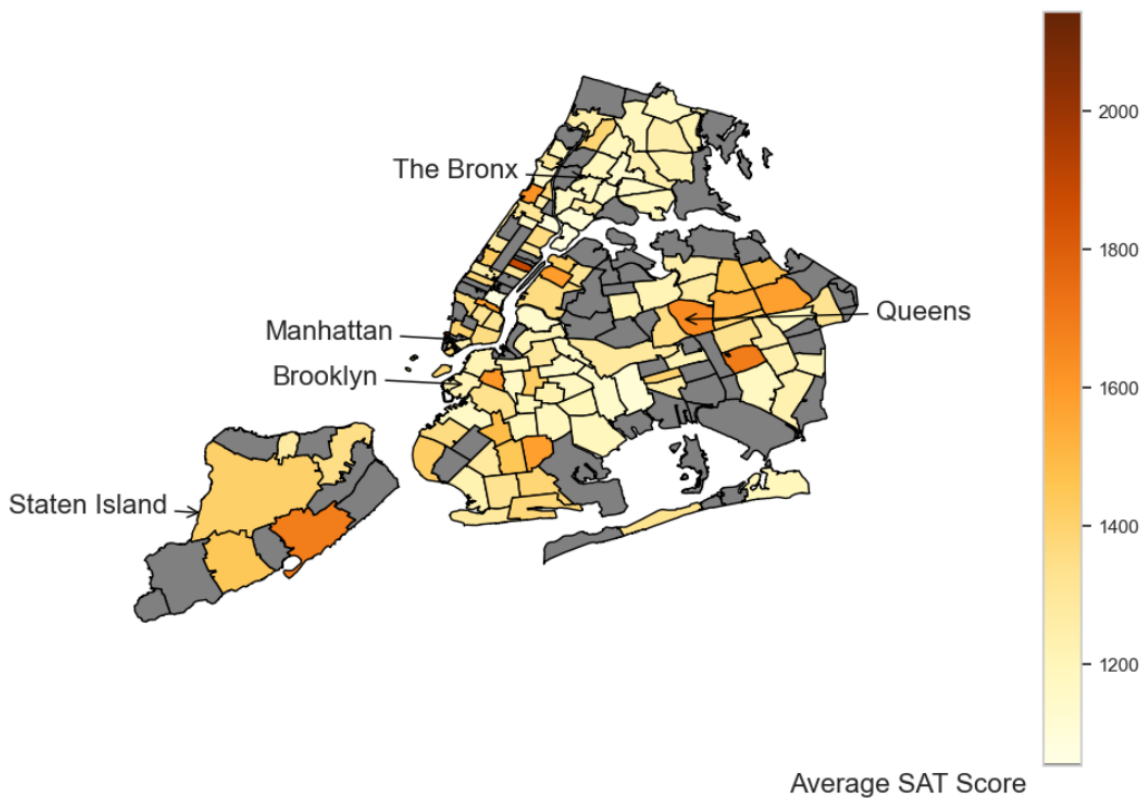


Figure 9: Distribution in Average SAT Scores in New York City

To gain insights into the spatial distribution and potential correlations between the variables and different neighbourhoods of New York City. The following maps provide a more comprehensive understanding of the city's dynamics and characteristics

Figure 9 focuses on the difference in average SAT scores across different zip code areas in New York City. The colour scale spans from the minimum to the maximum SAT score. Lighter colours indicate lower average SAT scores, while darker colours represent higher scores. Areas where there is no available public school average SAT data are depicted in gray. The largely consistent colour distribution shows that in most areas, the average SAT score ranges from 1250 to 1500. Each of the five boroughs has one or more particular zip code areas that stand out with an average SAT score higher than 1750. An area in Manhattan, notably, exhibits average SAT scores of more than 2000. These noticeable outliers are worth additional examination to understand the underlying factors contributing to these disparities.

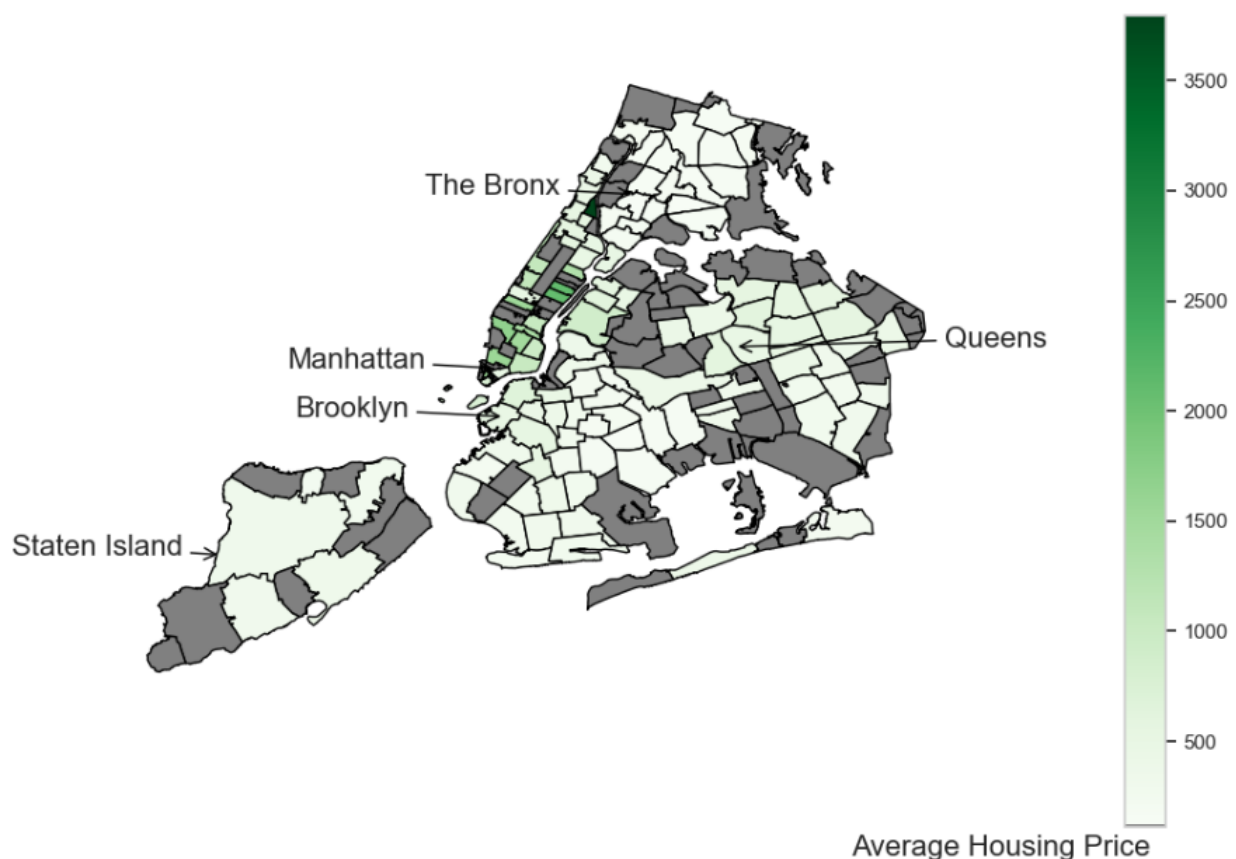


Figure 10: Distribution of Housing Prices in New York City

Notable financial disparities in the city might contribute to the difference in academic performance. Figure 10 categorizes the average housing prices across each zip code. From the map, it becomes evident that there isn't a strong correlation between the average SAT scores and housing prices apart from Manhattan. Most areas' average price per square foot is around \$500. However, Manhattan exhibits an abnormally high housing price compared to the rest of the city. In certain wealthy neighbourhoods in Manhattan, the average prices are as high as \$2500. The most expensive area in Manhattan with housing prices more than \$3000 interestingly is also the area with the highest average score in Figure 9. Along with the high prices in Manhattan, several areas in Brooklyn bordering Manhattan also exhibit high housing prices. Manhattan, being a global hub for wealth and luxury, exhibiting such remarkable housing costs is not surprising. Given that affluent individuals in Manhattan can afford top-notch education resources and teachers, it is logical for Manhattan to have the most number of outliers in high average SAT scores in the cities. Interestingly, Staten Island, with the highest median SAT scores and a majority of white students appears to have low housing prices. This scenario can be attributed to its geographical isolation being an island, housing availability, and

other local economic characteristics.

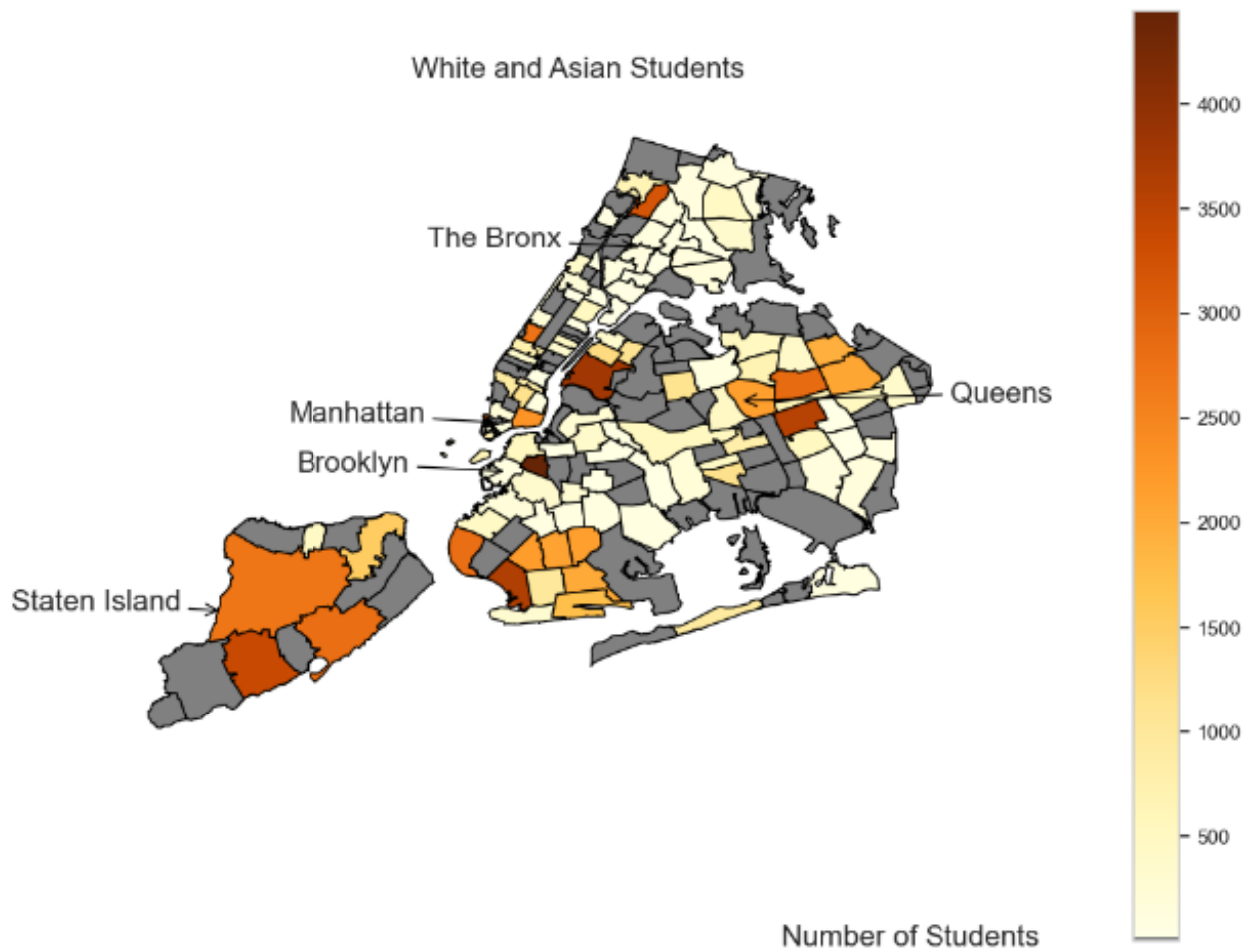


Figure 11: Geographical Distribution of White and Asian Students in New York City

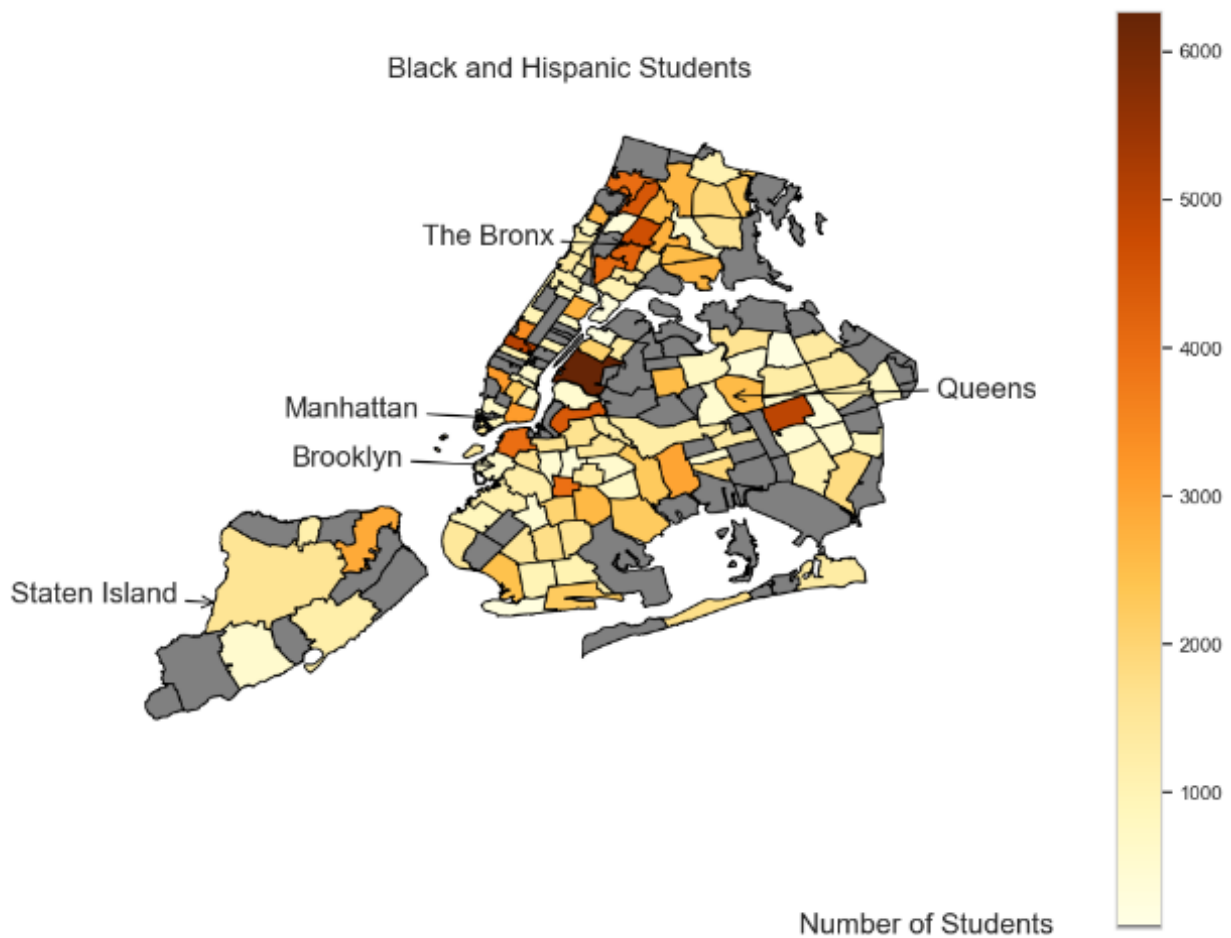


Figure 12: Geographical Distribution of Black and Hispanic Students in New York City

In addition to examining average SAT scores, it is equally important to understand demographic variations among explanatory variables. Building upon previous findings regarding the consistent trends in the distribution of white and Asian students, as well as black and Hispanic students, their demographic patterns were measured collectively. To identify unique areas where students of these racial groups may concentrate, Figure 11 and Figure 12 each illustrate the distribution of students by zip code. The two maps generally follow the trend detected earlier, with white and Asian students gathering in Staten Island and certain areas of each of the boroughs except the Bronx, while black and Hispanic students, representing a larger population, were dispersed across all areas, with a notable concentration in the Bronx.

Combining the three maps, it is generally the case where areas with outlying high SAT scores are associated with more white and Asian students. There is one Manhattan area where more than 4000 white and Asian students are located. This area also hosts similar numbers of students from other racial backgrounds. Interestingly, this corresponds to the area with the highest average SAT score on



---

Figure 9. Despite black and Hispanic students comprising a larger proportion of New York City public school populations, the disproportionate representation of white and Asian students in this particular area suggests a potential correlation between their presence and the high SAT scores observed. This phenomenon might be attributed to Manhattan having several elite public high schools that only admit students based on their test scores. The detailed breakdown by area and link between the two maps further strengthen the notion that racial demographics may significantly influence academic performances in New York City.

## 4 Model

To estimate the relationship between academic performance and student demographic factors, I employ Ordinary Least Squares (OLS) regression models to estimate the parameters. An OLS model is solved by finding the parameters that minimize the sum of squared residuals, i.e.

$$\min_{\hat{\beta}} \sum_{i=1}^N \hat{u}_i^2$$

where  $\hat{u}_i$  is the difference between the observation and the predicted value of the dependent variable.

The 4 main regression frameworks are presented below. I first inspect the relationship between racial identities and related factors and the average SAT score of a school. I then break down SAT scores into their subsections - Math, Reading, and Writing - to examine their relationship separately in Frameworks 2-4.

The variables I choose to inspect are the following: *White*, *Black*, *Hispanic*, *Asian*, *Other*, *StudentEnrollment*, *PercentTested*, *SchoolMinutes*, *HousingPrice(persqft)*, *Bronx*, *Brooklyn*, *Manhattan*, *Queens*, *StatenIsland*, *num\_AP*, *num\_language*, and *num\_extracurricular*. The five main racial variables are the focus of the paper as I attempt to investigate the impact of student racial identities on their test scores. Measuring the number of students in the school, *studentenrollment* controls the role of school size in the regression as bigger schools may have the advantage of having more resources. Recognizing the importance of peer influence, *PercentTested* controls for the percent of students in that school taking the SAT. To find out if longer school time affects students' test performance, *SchoolMinutes* measures the number of minutes a school day has. *HousingPrice(persqft)* measures the average housing price in the area where the school is located

as schools in wealthier areas might be better funded with more equipment and teachers. Measuring the geographical importance of the school on students' test scores, the five borough dummy variables are included. The variables *num\_AP*, *num\_language*, and *num\_extracurricular* are included to measure the impact of having diverse programs and extracurriculars on students' SAT performance.

$$\begin{aligned} \text{Total Average Score}_i = & \beta_0 + \beta_1 \text{Black}_i + \beta_2 \text{Hispanic}_i + \beta_3 \text{Asian}_i + \beta_4 \text{Other}_i \\ & + \beta_5 \text{Student Enrollment}_i + \beta_6 \text{Percent Tested}_i + \beta_7 \text{School Minutes}_i \\ & + \beta_8 \text{Bronx}_i + \beta_9 \text{Brooklyn}_i + \beta_{10} \text{Manhattan}_i + \beta_{11} \text{Queens}_i + \beta_{12} \text{AP}_i \\ & + \beta_{13} \text{Language}_i + \beta_{14} \text{Extracurricular}_i + u_i \end{aligned} \quad (1)$$

Firstly, I regress *TotalAverageScore* on the main explanatory variables: *Black*, *Hispanic*, *Asian*, and *Other*. *White* is excluded due to multicollinearity; therefore, the constant captures the average effect of *White*. The first model measures the simple relationship between racial identities and SAT scores without the interference of controls. To measure the effect of school characteristics on SAT scores, Regression 2 includes the controls *StudentEnrollment*, *PercentTested*, and *SchoolMinutes*. Regression 3 includes the additional variable *HousingPrice(persqft)*. However, housing price does not improve the model based on later evaluation of the results, it is dropped in later regressions. Regression 4 takes the score difference across boroughs into examination by including the dummy variables for four of five of the boroughs in New York City (Staten Island) is excluded for multicollinearity issues. Lastly, Regression 5 further assesses the influence of school opportunities such as the number of AP courses, language courses, and extracurriculars offered more classes and opportunities might have a positive influence on test scores. The results for the five models are presented in Table 1.

$$\begin{aligned} \text{Average Score (SAT Math)}_i = & \beta_0 + \beta_1 \text{Average Score (SAT Reading)}_i \\ & + \beta_2 \text{Average Score (SAT Writing)}_i + \beta_3 \text{Black}_i + \beta_4 \text{Hispanic}_i \\ & + \beta_5 \text{Asian}_i + \beta_6 \text{Other}_i + \beta_7 \text{Student Enrollment}_i \\ & + \beta_8 \text{Percent Tested}_i + \beta_9 \text{School Minutes}_i + \beta_{10} \text{Bronx}_i \\ & + \beta_{11} \text{Brooklyn}_i + \beta_{12} \text{Manhattan}_i + \beta_{13} \text{Queens}_i + \beta_{14} \text{AP}_i \\ & + \beta_{15} \text{Language}_i + \beta_{16} \text{Extracurricular}_i + u_i \end{aligned} \quad (2)$$

---


$$\begin{aligned}
\text{Average Score (SAT Reading)}_i &= \beta_0 + \beta_1 \text{Average Score (SAT Math)}_i \\
&+ \beta_2 \text{Average Score (SAT Writing)}_i + \beta_3 \text{Black}_i + \beta_4 \text{Hispanic}_i \\
&+ \beta_5 \text{Asian}_i + \beta_6 \text{Other}_i + \beta_7 \text{Student Enrollment}_i \\
&+ \beta_8 \text{Percent Tested}_i + \beta_9 \text{School Minutes}_i + \beta_{10} \text{Bronx}_i \\
&+ \beta_{11} \text{Brooklyn}_i + \beta_{12} \text{Manhattan}_i + \beta_{13} \text{Queens}_i + \beta_{14} \text{AP}_i \\
&+ \beta_{15} \text{Language}_i + \beta_{16} \text{Extracurricular}_i + u_i
\end{aligned} \tag{3}$$

$$\begin{aligned}
\text{Average Score (SAT Writing)}_i &= \beta_0 + \beta_1 \text{Average Score (SAT Math)}_i \\
&+ \beta_2 \text{Average Score (SAT Reading)}_i + \beta_3 \text{Black}_i + \beta_4 \text{Hispanic}_i \\
&+ \beta_5 \text{Asian}_i + \beta_6 \text{Other}_i + \beta_7 \text{Student Enrollment}_i \\
&+ \beta_8 \text{Percent Tested}_i + \beta_9 \text{School Minutes}_i + \beta_{10} \text{Bronx}_i \\
&+ \beta_{11} \text{Brooklyn}_i + \beta_{12} \text{Manhattan}_i + \beta_{13} \text{Queens}_i + \beta_{14} \text{AP}_i \\
&+ \beta_{15} \text{Language}_i + \beta_{16} \text{Extracurricular}_i + u_i
\end{aligned} \tag{4}$$

Regressions from Tables 2, 3, and 4 are very similar. These three tables examine the associations between racial identities and SAT Math, Reading, and Writing separately. Each table starts by controlling only the racial identity variables to observe the simple relationship between these variables. Then the second regression controls only controls for the other two SAT subjects that are not being examined to see if performing well in other SAT subjects would also improve the performance of the subject of examination. The third regression combines racial variables and SAT subject scores. The rest of the regression follows the same pattern as the first table. The addition of control variables that also might be associated with SAT scores allows for a more comprehensive analysis of the factors influencing SAT performance and decreases the chance of omitted variable bias.

## 5 Regression Results

Dependent variable: Total Average Score					
	(1)	(2)	(3)	(4)	(5)
Asian	-0.596 (0.812)	-1.173 (0.751)	-1.539** (0.746)	-1.571* (0.835)	-1.341 (0.815)
Black	-6.050*** (0.530)	-4.905*** (0.524)	-4.885*** (0.522)	-5.115*** (0.569)	-5.104*** (0.540)
Bronx				123.318*** (43.456)	123.964*** (41.339)
Brooklyn				48.804 (41.730)	53.838 (39.713)
Hispanic	-6.590*** (0.553)	-5.444*** (0.547)	-5.456*** (0.540)	-6.603*** (0.612)	-6.465*** (0.583)
Housing Price (per sqft)			0.022 (0.014)		
Manhattan				135.667*** (42.669)	136.062*** (40.678)
Other	15.701*** (4.057)	11.472*** (3.810)	10.548*** (3.781)	8.185** (3.767)	8.466** (3.609)
Percent Tested		3.002*** (0.371)	2.928*** (0.366)	2.699*** (0.362)	2.226*** (0.355)
Queens				57.417 (43.412)	58.472 (41.204)
School Minutes		0.377* (0.193)	0.362* (0.191)	0.324* (0.186)	0.115 (0.184)
Student Enrollment		0.022** (0.009)	0.021** (0.009)	0.030*** (0.009)	-0.012 (0.012)
const	1757.729*** (51.583)	1310.201*** (95.928)	1316.009*** (95.679)	1325.113*** (93.523)	1416.500*** (92.678)
num_AP					9.180*** (1.730)
num_extracurricular					1.092* (0.626)
num_language					-10.547* (5.572)
Observations	374	374	373	374	363
R <sup>2</sup>	0.603	0.669	0.663	0.698	0.735
Adjusted R <sup>2</sup>	0.599	0.662	0.656	0.689	0.725
Residual Std. Error	123.419 (df=369)	113.246 (df=366)	111.362 (df=364)	108.695 (df=362)	102.827 (df=348)
F Statistic	140.216*** (df=4; 369)	105.488*** (df=7; 366)	89.638*** (df=8; 364)	76.077*** (df=11; 362)	69.060*** (df=14; 348)
Note:	*p<0.1; **p<0.05; ***p<0.01				

Table 2: Regression Table 1

---

Regression 1 models the effect of racial identities on average SAT scores. The average influence of being a white student on SAT score is the constant 1757.729. On average compared to white students, a one percentage point increase in the number of black students decreases the average SAT score at the school by 6.05 percentage points, holding all other variables constant. Increases in the percentage of Hispanic students have similar outcomes. An increase in the percentage of Asian students also decreases the SAT score slightly compared to white students, but it is not statistically significant. All variables except for *Asian* are highly statistically significant with a p-value smaller than 0.01.

These findings suggest that on average, white students tend to outperform all minority students on the SAT exam with the exception of students identifying as *Other*. The large score gap between white black and Hispanic students can be attributed to many socioeconomic factors. Black and Hispanic communities often face disadvantages like lower household incomes and limited access to quality education resources. Schools in predominantly poorer neighbourhoods where these individuals might live often receive fewer resources and funding. All of these characteristics might contribute to lower scores in schools with more black and Hispanic populations.

The addition of school characteristics does not vary the racial identities association with SAT scores much. Having more students taking the test creates a better study environment at school and might positively influence their peers to study harder. Higher student enrollments and minutes in school are also positively related to SAT scores student enrollment might be associated with more school resources available, and staying longer in school might allow students to have the opportunity to learn more and practice more. These all positively impacted students' SAT performances.

The addition of *HousingPrice(persqft)* in Regression 3 has no statistically significant impact on the dependent variable. This could be due to the complexity of housing market dynamics, influenced by various factors location and amenities. The complexity makes it difficult to isolate the specific effect of housing prices on educational outcomes.

The addition of borough and school opportunities controls in Regressions 4 and 5 does not greatly impact the influence of racial identity variables. The schools located in The Bronx and Manhattan both have a positive statistically significant influence on average SAT scores compared to living in Staten Island. The dummy variable for *StatenIsland* is omitted to avoid multicollinearity. Therefore, the constant serves as a comparison variable for white students living in Staten Island. More AP classes and extracurriculars also positively impact the SAT scores; while surprisingly one more language course offered decreases the average SAT score. This decrease might be due to time and

---

resource allocation. Language courses might require significant time and commitment from students which takes away their time for SAT preparations.

The addition of new variables from Regression 1 to Regression 5 increases the adjusted R squared, except for Regression 3 where *HousingPrice(per sq ft)* is added. Therefore, Regressions 4 and 5 drop the variable. The addition of borough controls and course offering controls in Regressions 4 and 5 increases the adjusted R squared as it becomes 0.689 and 0.725 respectively. In conclusion, the addition of the new variables improves the model more than expected by chance.

The F-statistic, measuring the differences among sample averages, decreases with the addition of new controls. However, it remains large indicating that there are fairly great differences among sample averages in the models.

Dependent variable: Average Score (SAT Math)						
	(1)	(2)	(3)	(4)	(5)	(6)
Asian	0.810*** (0.274)		1.452*** (0.113)	1.381*** (0.113)	1.560*** (0.127)	1.417*** (0.125)
Average Score (SAT Reading)		0.309** (0.126)	0.546*** (0.085)	0.538*** (0.084)	0.527*** (0.082)	0.569*** (0.081)
Average Score (SAT Writing)		0.748*** (0.120)	0.360*** (0.083)	0.337*** (0.082)	0.352*** (0.080)	0.330*** (0.079)
Black	-1.973*** (0.179)		-0.137 (0.085)	-0.125 (0.086)	-0.099 (0.094)	-0.079 (0.092)
Bronx					-4.751 (6.562)	-3.441 (6.335)
Brooklyn					2.205 (6.236)	4.020 (6.017)
Hispanic	-2.011*** (0.187)		0.057 (0.090)	0.058 (0.091)	0.142 (0.106)	0.170 (0.104)
Manhattan					-1.682 (6.457)	-1.303 (6.249)
Other	3.803*** (1.370)		-1.580*** (0.572)	-1.604*** (0.571)	-1.255** (0.567)	-1.034* (0.551)
Percent Tested				0.181*** (0.059)	0.154*** (0.058)	0.147*** (0.056)
Queens					-11.636* (6.500)	-10.083 (6.254)
School Minutes				0.057** (0.029)	0.061** (0.028)	0.057** (0.028)
Student Enrollment				0.002* (0.001)	0.002 (0.001)	0.000 (0.002)
const	575.814*** (17.415)	-11.574 (9.419)	40.313*** (14.323)	16.214 (17.380)	10.463 (17.424)	-2.072 (18.226)
num_AP						0.455* (0.271)
num_extracurricular						0.136 (0.095)
num_language						0.089 (0.847)
Observations	374	374	374	374	374	363
R <sup>2</sup>	0.668	0.874	0.945	0.947	0.951	0.956
Adjusted R <sup>2</sup>	0.664	0.874	0.944	0.946	0.949	0.954
Residual Std. Error	41.667 (df=369)	25.558 (df=371)	16.970 (df=367)	16.700 (df=364)	16.213 (df=360)	15.539 (df=346)
F Statistic	185.540*** (df=4; 369)	1291.174*** (df=2; 371)	1055.297*** (df=6; 367)	728.146*** (df=9; 364)	536.860*** (df=13; 360)	473.009*** (df=16; 346)
Note:					* p<0.1; ** p<0.05; *** p<0.01	

Table 3: Regression Table 2

---

Regression Table 2 models the relationship between racial and related factors on SAT Math scores. Specifically, Regression 1 models the influence of racial identities on average SAT Math scores. The average influence of being a white student on SAT Math score is the constant 575.814. On average compared to white students, a one percentage point increase in the number of black students decreases the average SAT Math score at the school by 1.973 percentage points, holding all other variables constant. Increases in the percentage of Hispanic students have similar outcomes. An increase in the percentage of Asian students increases the SAT Math score by 0.81 points, compared to white students.

The racial identity variables alter slightly throughout the table with new controls added, with only *Asian* and *Others* remaining significant. This finding suggests that schools with more Asian students on average are better in math than other students. Asian students' better performance might be attributed to family emphasis, cultural attitudes toward the importance of math education, and additional support systems within their communities.

Regression 2 proves that having higher scores in other SAT sections is significantly associated with higher SAT math scores. This finding suggests that students who perform well in other SAT sections tend to perform well in math as well. This finding can be explained by the fact that students who excel academically usually perform well in all subjects instead of one. Students who value their test scores would make sure they are prepared for all parts of the SAT exams.

School characteristics such as *PercentTested* continue to have a highly significant positive relationship with SAT scores, suggesting more students taking the test positively influence the average test scores of the whole school. The addition of borough and school opportunities controls in Regressions 4 and 5 does not greatly impact the influence of racial identity variables. More language course offered continues to decrease the average SAT score on average.

The addition of new variables from Regression 1 to Regression 6 increases the adjusted R-squared. The F-statistic, measuring the differences among sample averages, fluctuates with the addition of new controls. The very large F-statistics in Regression 1 and 2 might indicate unusually great differences among sample averages in the models.



Dependent variable: Average Score (SAT Reading)						
	(1)	(2)	(3)	(4)	(5)	(6)
Asian	-0.729*** (0.276)		-0.349*** (0.077)	-0.347*** (0.077)	-0.387*** (0.090)	-0.375*** (0.089)
Average Score (SAT Math)		0.052** (0.021)	0.187*** (0.029)	0.190*** (0.029)	0.198*** (0.031)	0.219*** (0.031)
Average Score (SAT Writing)		0.891*** (0.024)	0.784*** (0.028)	0.787*** (0.028)	0.775*** (0.029)	0.747*** (0.031)
Black	-1.975*** (0.180)		0.041 (0.050)	0.046 (0.051)	0.041 (0.058)	0.021 (0.057)
Bronx					1.508 (4.018)	1.277 (3.929)
Brooklyn					-0.866 (3.817)	-1.512 (3.732)
Hispanic	-2.251*** (0.188)		-0.052 (0.053)	-0.047 (0.054)	-0.085 (0.065)	-0.102 (0.064)
Manhattan					2.249 (3.950)	2.500 (3.873)
Other	5.904*** (1.379)		0.496 (0.337)	0.582* (0.342)	0.498 (0.348)	0.420 (0.343)
Percent Tested				-0.052 (0.035)	-0.052 (0.036)	-0.047 (0.035)
Queens					1.638 (3.995)	1.479 (3.892)
School Minutes				0.002 (0.017)	0.001 (0.017)	-0.006 (0.017)
Student Enrollment				0.000 (0.001)	0.001 (0.001)	0.001 (0.001)
const	590.405*** (17.537)	29.305*** (3.558)	19.412** (8.401)	17.843* (10.293)	21.628** (10.608)	30.617*** (11.182)
num_AP						0.000 (0.169)
num_extracurricular						-0.072 (0.059)
num_language						-0.247 (0.525)
Observations	374	374	374	374	374	363
R <sup>2</sup>	0.545	0.971	0.975	0.975	0.975	0.977
Adjusted R <sup>2</sup>	0.540	0.971	0.974	0.974	0.974	0.976
Residual Std. Error	41.960 (df=369)	10.479 (df=371)	9.919 (df=367)	9.920 (df=364)	9.923 (df=360)	9.636 (df=346)
F Statistic	110.585*** (df=4; 369)	6319.168*** (df=2; 371)	2358.501*** (df=6; 367)	1572.535*** (df=9; 364)	1088.288*** (df=13; 360)	907.206*** (df=16; 346)
Note:					* p<0.1; ** p<0.05; *** p<0.01	

Table 4: Regression Table 3

Dependent variable: Average Score (SAT Writing)						
	(1)	(2)	(3)	(4)	(5)	(6)
Asian	-0.677** (0.282)		-0.156* (0.083)	-0.151* (0.083)	-0.225** (0.096)	-0.213** (0.096)
Average Score (SAT Math)		0.126*** (0.020)	0.136*** (0.031)	0.131*** (0.032)	0.145*** (0.033)	0.144*** (0.035)
Average Score (SAT Reading)		0.892*** (0.024)	0.867*** (0.031)	0.865*** (0.031)	0.851*** (0.032)	0.849*** (0.035)
Black	-2.102*** (0.184)		-0.121** (0.052)	-0.125** (0.053)	-0.142** (0.060)	-0.139** (0.060)
Bronx					3.478 (4.209)	3.627 (4.185)
Brooklyn					1.385 (4.001)	1.577 (3.979)
Hispanic	-2.327*** (0.192)		-0.101* (0.055)	-0.103* (0.057)	-0.139** (0.068)	-0.146** (0.068)
Manhattan					1.824 (4.142)	1.721 (4.131)
Other	5.994*** (1.410)		0.357 (0.355)	0.228 (0.359)	0.163 (0.366)	0.210 (0.366)
Percent Tested				0.065* (0.037)	0.069* (0.037)	0.056 (0.038)
Queens					4.237 (4.183)	4.268 (4.144)
School Minutes				-0.014 (0.018)	-0.016 (0.018)	-0.013 (0.018)
Student Enrollment				-0.001 (0.001)	-0.001 (0.001)	-0.001 (0.001)
const	591.511*** (17.928)	-14.719*** (3.796)	1.083 (8.901)	6.858 (10.829)	7.821 (11.177)	8.002 (12.043)
num_AP						0.077 (0.180)
num_extracurricular						0.055 (0.063)
num_language						-0.069 (0.560)
Observations	374	374	374	374	374	363
R <sup>2</sup>	0.563	0.974	0.974	0.975	0.975	0.976
Adjusted R <sup>2</sup>	0.558	0.974	0.974	0.974	0.974	0.975
Residual Std. Error	42.896 (df=369)	10.486 (df=371)	10.434 (df=367)	10.398 (df=364)	10.401 (df=360)	10.274 (df=346)
F Statistic	118.900*** (df=4; 369)	6880.781*** (df=2; 371)	2317.899*** (df=6; 367)	1556.594*** (df=9; 364)	1077.267*** (df=13; 360)	872.515*** (df=16; 346)
Note:					* p<0.1; ** p<0.05; *** p<0.01	

Table 5: Regression Table 4

---

The design for Regression Tables 3 and 4 are the same as Table 2, with Table 3 exploring SAT Reading scores and Table 4 exploring SAT Writing scores. Having more Asian and Hispanic students are both consistently associated with lower SAT Reading and Writing scores compared to white students. More black student population is also linked with lower SAT Reading scores initially; however, with the addition of controls, this relationship became positive and statistically insignificant; suggesting a weak relationship between the two variables. White students consistently outperforming minority students in reading and writing tests might be due to English proficiencies. Good performance in these tests relies on a good understanding of the English language, culture, and vocabulary. English may not be the first language of many of these students; therefore, they might struggle more to perform well in these exams compared to white students.

Good performance in other sections of the SAT exam continues to have a positive relationship with the subject of the examination. This phenomenon again is due to students tending to balance their strengths across different subjects. For example, students who excel in math have the analytical skill that allows them to read and write critically. This highlights the interconnection of all subjects and underscores the importance of a well-rounded education.

Other than *PercentTested*, school characteristics have no significant influence on both Reading and Writing scores. Notably, the relationship between *studentenrollment* and both scores is almost 0, indicating almost no linear relationship between the two variables. The lack of association might be because larger schools while having more resources also experience challenges like overcrowded classrooms and limited individual attention that may negatively impact students' performance. Borough controls also lack significant relationships with test scores, further validating the fact that while education disparities across areas exist, they are not the main reasons for the issue. Boroughs in New York City usually contain wide ranges of socioeconomic neighbourhoods and even within boroughs, there are significant differences in resources and funding. Therefore, using boroughs as controls might have overlooked the individual nuances across neighbourhoods. School courses and opportunity offerings also have minimum impact on both sections of the SAT. Notably, the number of AP classes the school offers has no linear relationship with SAT Reading scores. This finding suggests that while AP and language and extracurricular offerings have a positive relationship with the total SAT score, they do not have any significant influence on SAT subsection scores. Compared to course offerings, effective teaching skills, students' own language abilities, and access to other resources might contribute more to students well performance in SAT Reading and Writing exams.

The addition of new controls consistently increases the adjusted R squared for SAT Reading and

---

SAT Writing. These increases suggest that the addition of new variables improved the models more than expected by chance. The F-statistics, measuring the differences among sample averages, again fluctuate with the addition of new controls. It grew very large while I only assessed the relationship between SAT subsections (6319.168 and 6880.781 respectively) and gradually decreased to 907.206 for SAT Reading and 872.515 for SAT Writing. The very large F statistics throughout the models indicate great differences among sample averages in the models.

## 6 Machine Learning

### 6.1 Decision Tree

Following the previous regression models and interpretations, this section employs regression tree models in exploring the relationship between racial identity factors on SAT scores. The advantage of using regression trees is that they adapt automatically to feature scales and units. My model aims to provide similar results as my linear regression outputs.

The goal of the regression tree is to minimize MSE by choosing a feature and location to split on. The minimization is based on the minimum objective function:

$$\min_{j,s} \left[ \sum_{i: x_{i,j} \leq s, x_i \in R1} (y_i - \hat{y}_{R1})^2 + \sum_{i: x_{i,j} > s, x_i \in R2} (y_i - \hat{y}_{R2})^2 \right]$$

The decision on the location to split is based on the above minimization problem. To split the tree, the regressor chooses a variable to split on and a threshold from the rectangular region  $R$  containing all values of  $X$ . This process is repeated to form all the branches. We want to stop when  $|R| =$  some chosen minimum size or when depth of tree = some chosen maximum. The full objective function for my regression tree is the following:

$$\min_{tree \in T} \sum (\hat{f}(x) - y)^2 + \alpha |\text{terminal nodes in tree}|$$

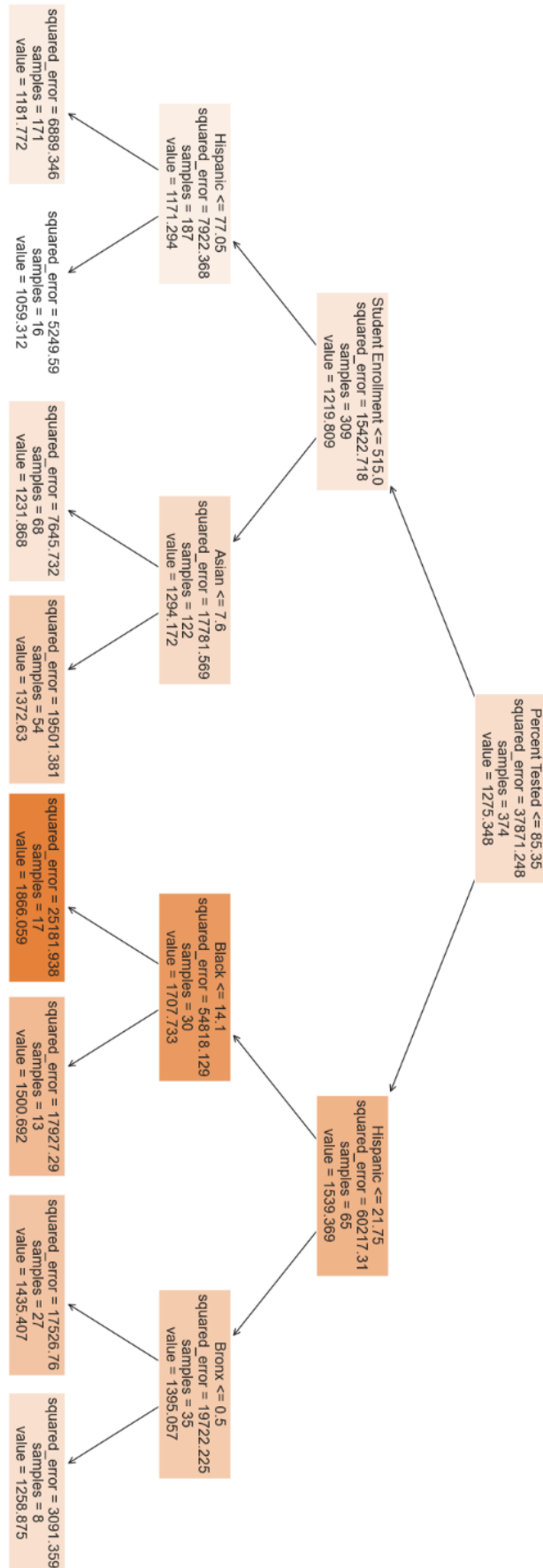


Figure 13: Decision Tree  
29

---

This decision tree is run using the independent variables from my preferred specification Regression 5. The tree highlights the association between all the controls and average SAT scores in school by breaking them down into groups according to the threshold calculated by algorithms. The goal of the tree is to minimize the mean squared error by splitting the variables based on the threshold. The error of prediction (MSE) of the entire decision tree model is 10679.56; it is slightly larger than the MSE from the regression model (10136). Therefore, the regression model's prediction accuracy is slightly larger than that of the tree. The error can be minimized at the expense of having an overly complex tree.

The tree has three levels and started its first split at the variable *PercentTested*. This variable is chosen for the first split due to the importance of a school having many of its students taking the SAT exam. Students who value their education more tend to take the exam, and having many students taking the exam creates a positive learning environment that may potentially linked to higher test scores. There are 8 terminal nodes of the tree. The smallest squared error is 5249.59, which is significantly smaller than the MSE we began with. Using the tree, I can predict that a school with less than 85.35% of students that took the SAT, a student enrollment of less than 515, and less than 77.05% of Hispanic students on average has an average SAT score of 1059. Conversely, a school with more than 85.35% of students that took the SAT, a Hispanic student population of less than 21.75%, black student population of less than 14.1% is predicted to have on average an average SAT score of 1866.

The decision tree provides additional evidence for the regression results. The tree first splits on the variable *PercentTested* which reflects the fact that this variable is highly significant to the 1% level on the regression and a school with a higher percentage of students who take the test is associated with higher average SAT scores. The leaf nodes confirm this assumption that all four nodes with higher average scores stemmed from having a percentage tested higher than 85.35%. \*Hispanic\* being chosen as the next split also matches the regression output which is highly statistically significant and having a higher percentage of Hispanic students is associated with lower scores. From the tree output, I observe that nodes with a Hispanic student population of less than 21.75% are associated with the highest SAT scores (1866 and 1500 respectively). Furthermore, schools with higher student enrollment (more than 515) are associated with higher test scores matching the regression output. From *SchoolEnrollment*, the tree further split into *Hispanic* and *Asian* and their results match the regression output as well. Schools with more than 77.05% Hispanic student population are associated with the lowest average SAT score (1059). This again brings more evidence to our

---

conclusion that schools with a higher number of Hispanic students are related to lower test scores. On the other hand, schools with more than 7.6% Asian student population are linked to higher average scores of 1372. The fact that the score is a lot higher than schools with less than 7.6% (1231) matches the regression output that schools with more Asian populations are associated with higher scores. Moving back to the node *Hispanic* at level 1, this node branches down to *Black* and *Bronx*. Similar to the regression results schools with fewer black students (<14.1%) and less likely to be in the Bronx (<0.5%) both end up with considerably higher average SAT scores.

From the tree, I conclude that along with racial identity differences, many factors contribute to SAT score disparities. However, the fact that the decision tree output comes to the same conclusion as the regression that schools with more racial minorities, in particular, black and Hispanic students are consistently linked to lower scores shows the potential education disparities existing in the city. In addition, schools with fewer students taking the test also is an important factor contributing to lower test scores, which highlights the impact of peer influence and the school environment.

## 6.2 Random Forest Model

A single regression tree's results usually have high variance. Therefore, I introduce a Random Forest Model that reduces the variance by fitting many bootstrapped data. Unlike bootstrapping which has the risk of creating many similar trees, the Random Forest Model selects a random sample of  $X$ s at each node and chooses to split variables from this sample of  $X$ s. Using this model, I am able to further decrease the MSE to 1381.15.

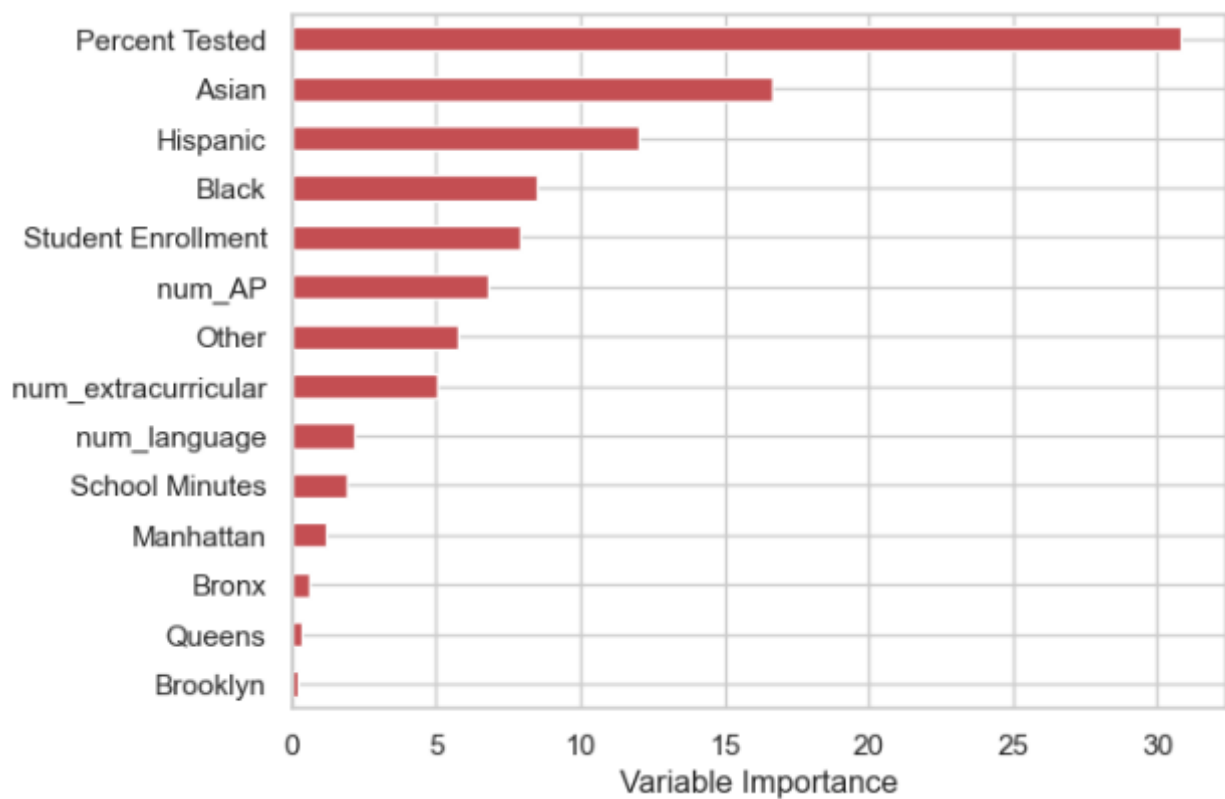


Figure 14: Importance Matrix

Lastly, to measure the influence of each variable in reducing the MSE, I create this importance matrix that ranks the importance of each variable. Under my expectations, *PercentTested* contributes to over 30% of the total contribution of all the variables. The three racial factors *Asian*, *Hispanic*, and *Black* follow *PercentTested*. This suggests that schools with more students identifying as certain races generate considerable impact on students' scores. Following racial identity variables, other school-related factors such as *StudentEnrollment* and *num<sub>AP</sub>* exert some impact on the outcome as well. Borough-related factors have the smallest influence on SAT scores. This finding suggests the difference in SAT scores across the city is not hugely due to area differences.

## 7 Conclusion

This study seeks to unravel the relationship between student racial demographics and academic achievements in New York City. Utilizing data about SAT scores in New York City public schools, I explored various racial demographics and their impact on average SAT scores across areas in New York City by zip code. I notice higher proportions of white and Asian students are consistently associated with higher academic performance, while a higher share of other races tend to demonstrate



---

lower average SAT scores. The trend that schools with more white students outperform their counterparts persists when the SAT exam is broken down into three subject tests: Math, Reading, and Writing. While schools with more Asian students have advantages in gaining higher math scores, they have lower Reading and Writing scores along with the other minority counterparts.

In addition, the investigation of zip code-level performance indicates complex education disparities across different neighbourhoods. The highest average SAT score belongs to a Manhattan school while the lowest goes to the Bronx. It is also observable that affluent areas in Manhattan with remarkably above-average housing prices are associated with above-average SAT scores. This result emphasizes the need to ensure that every student, regardless of their race or location, should have access to quality education.

The study extends beyond racial demographics to related factors such as school statistics, program offerings, and school location. Schools where a higher percentage of their students taking the SAT tests are highly correlated with higher overall SAT scores and subsection scores. Schools with more advanced course offerings such as AP classes and diverse extracurricular opportunities also have higher test scores. More motivated students tend to engage in these opportunities; therefore they would also put equal effort into achieving a high SAT score so that they have a higher chance of getting admission to good universities. On the other hand, schools with more language class offerings lower their average SAT score. This phenomenon can be explained by students allocating more of their time to language classes that would not improve their SAT scores directly.

One caveat of this study is the limit in data I have in hand. Because data for this study was collected almost a decade ago, newer versions of the data might introduce new findings. The analysis would be more comprehensive if cross-sectional data from multiple years could be retrieved instead of one year. With cross-sectional data, I could present the trend of the relationship between racial identities and test scores over the years. Even though there are noticeable SAT score differences across different areas, there is no significant linear relationship between boroughs and test scores because of the wide range of differences within each borough. Therefore, the main goal for potential future stages of this study is to incorporate year-level and smaller neighbourhood-level fixed effects to control for the variation across each year and neighbourhood. Lastly, if possible, quasi-experimental processes such as diff-in-diff or regression discontinuity can be utilized to form a causal inference on the relationship. However, given the currently available data, the analyses performed are highly valuable.

In conclusion, my findings reveal education disparities are still very much present in New York City. The study explains the relationship between demographic factors and education outcomes,

---

along with its covariates. Understanding the importance of this relationship is essential in creating and maintaining a more equitable education system that provides every student with the same chance to succeed.

---

## References

*Average SAT Scores for NYC Public Schools*. 2014.

**URL:** <https://www.kaggle.com/datasets/nycopenda>

Battle, Juanita and Michael Lewis. 2002. "The Increasing Significance of Class: The Relative Effects of Race and Socioeconomic Status on Academic Achievement." *Journal of Poverty* 6(2):21–35.

**URL:** [https://doi.org/10.1300/j134v06n02\\_2](https://doi.org/10.1300/j134v06n02_2)

Bonastia, C. 2023. "Segregation in New York City Schools Continues." Accessed: 2023-11-22.

**URL:** <https://www.thirteen.org/blog-post/segregation-new-york-city-schools-continues/gazetteterrymurphy>

Boroughs of New York City. 2020. "Boroughs of New York City."

**URL:** [https://en.wikipedia.org/wiki/Boroughs\\_of\\_New\\_York\\_City](https://en.wikipedia.org/wiki/Boroughs_of_New_York_City)

Card, D. and J. Rothstein. 2007. "Racial segregation and the black–white test score gap." *Journal of Public Economics* 91(11-12):2158–2184.

Davis, Jane R and Nicholas Warner. 2015. "Schools matter." *Urban Education* 53(8):004208591561354.

**URL:** <https://doi.org/10.1177/0042085915613544>

Di, M., C. Kinga, W. Di, C. Fenelon, K. Flood, E. Milborn and C. Rodriguez. 2021. "Public and private school segregation in New York City."

**URL:** <https://files.eric.ed.gov/fulltext/ED613612.pdf>

Everson, H. T. and R. E. Millsap. 2004. "Beyond individual differences: Exploring school effects on SAT scores." *Educational Psychologist* 39(3):157–172.

Everson, H. T. and R. E. Millsap. 2005. "Everyone Gains: Extracurricular Activities in High School and Higher SAT® Scores." Research Report No. 2005-2.

**URL:** <https://eric.ed.gov/?id=ED562676>

Fetler, M. 1989. "School Dropout Rates, Academic Performance, Size, and Poverty: Correlates of Educational Reform." *Educational Evaluation and Policy Analysis* 11(2):109–116.

---

Harris, E. A. and W. Hu. 2018. “Asian groups see bias in plan to diversify New York’s elite schools.”.

**URL:** <https://www.nytimes.com/2018/06/05/nyregion/carranza-specialized-schools-admission-asians.html>

Kemple, J., C. Farley and K. Stewart. 2019. “Wide gap in SAT/ACT test scores between wealthy, lower-income kids.”.

**URL:** <https://news.harvard.edu/gazette/story/2023/11/new-study-finds-wide-gap-in-sat-act-test-scores-between-wealthy-lower-income-kids/>

McKillip, M. E. M. and A. Rawls. 2013. “A Closer Examination of the Academic Benefits of AP.” *The Journal of Educational Research* 106(4):305–318.

NYC Open Data. 2014. “2014 - 2015 DOE High School Directory | NYC Open Data.” [https://data.cityofnewyork.us/Education/2014-2015-DOE-High-School-Directory/n3p6-zve2/about\\_data](https://data.cityofnewyork.us/Education/2014-2015-DOE-High-School-Directory/n3p6-zve2/about_data). Accessed : [April 10, 2024].

NYC Property Sales. n.d. “NYC Property Sales.” <https://www.kaggle.com/datasets/new-york-city/nyc-property-sales>. Accessed: [April 10, 2024].

Penney, Jeffrey. 2017. “Racial Interaction Effects and Student Achievement.” *Education Finance and Policy* 12(4):447–467.

**URL:** [https://doi.org/10.1162/edfp\\_a00202](https://doi.org/10.1162/edfp_a00202)