An Analysis of the Department of Education Quality Survey and Its Efficacy

Andrew Bowen[1], Glen Dale Davis[1], Josh Forster[1], Shoshana Farber[1], & Charles Ugiagbe[1]

[1] City University of New York

Abstract

Abstract coming soon!

*Keywords:* Educational Outcomes, School Quality, Education

An Analysis of the Department of Education Quality Survey and Its Efficacy

## Introduction

The NYC School Survey seeks to collect data to provide an overview of New York City Schools. Beginning in 2005, the survey looks to collect demographic and achievement data for New York City Public Schools, and provide a standardized rating of various elements of school quality.

The survey has changed over the years. This change has come from recommendations of public policy analysts in order to more accurately define the quality of schools *New York City Schools (2018)*. The 2020-21 academic year report provides a robust dataset defined at the school level with academic and socioeconomic data provided.

**Research Question:** This study aims to determine whether the school ratings within the NYC School Quality Survey accurately reflect educational outcomes, or if other variables related to certain schools can be used as a better proxy.

## Literature Review

One of the main predictors of academic performance is the socioeconomic background of a student. Students from low-income families are nearly four times more likely to drop out of high school than students from wealthy families *Education Statistics (2008)*.

Attempts to use more sophisticated modeling techniques and different sources datasets come from several prior studies. *Bernacki, Chavez, and Uesbeck (2020)* based their modeling off trying to predict based on student digital behavior, rather than social factors. The model in this study reached an accuracy of 75%, and was able to flag early interventions. While this modeling technique attempts to predict the same variable (educational achievement, albeit a different metric where we are predicting college attainment), the base dataset used to train the model and input variables are different.

Similarly, *Musso, Cascallar, Bostani, and Crawford (2020)* attempted to train an artificial neural network (ANN) to identify variable relationships to educational performance data. They modeled educational performance of Vietnamese students in grade 5. They included individual characteristics as well as information related to daily routines in their training data. This method uses a more sophisticated model, and resulted in accuracy in prediction of $95 - 100$ higher than other modeling techniques. However, the training data came in that case from a different country (Vietnam, rather than the United States). Comparing modeling results from this (and other US-centric studies) may not be prudent.

*Yağcı (2022)* predicted final grade exams for Turkish students as well via machine learning models. Their input variables were prior exam grades. These can be a good "vacuum" comparison to compare one set of academci performance to another. However, there is a concern that good exam grades (even in one subject) do not correspond to a higher rate of career success later in life *Afarian and Kleiner (2003)*. Additionally, a parent study also found a correlation of up to 0.3 between academic grades and later job performance *Roth, BeVier, Switzer III, and Schippmann (1996)*.

Measuring the input variables that impact educational outcomes is a difficult task. With so many confounding variables, it can be difficult to determine direct causal relationships that have an outsized impact

## Data Sourcing

The dataset used in this study is published from the NYC School Quality Report for the Academic Year 2020 - 2021. It consists of data from 487 New York City public schools, and 391 variables (in the form of columns). This dataset is defined at the school level, indexed by a school's *district borough number* (DBN).

In addition to the school quality ratings provided from survey responses in the data, there is average and raw academic performance data included. In addition to thesea

academic indicators, there are socioeconomic variables included as well, such as the

percentage of students at a given school in temporary housing services.

## Methodology

We create a 20% holdout set of data to be used later on in order to evaluate the

efficacy of our model's predictive capability. The remaining 80% of the data is to be used for

model training and exploratory data analysis (EDA).

Additionally, we impute both our training and evaluation datasets. Given we are

dealing with continuous numeric (and not categorical variables), we use the *Predictive Mean*

*Matching* imputation method native to the R `mice` package

```
##
## iter imp variable
##  1   1  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  1   2  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  1   3  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  1   4  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  1   5  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  2   1  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  2   2  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  2   3  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  2   4  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  2   5  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  3   1  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  3   2  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  3   3  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
##  3   4  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg
```

```
## 3  5  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 4  1  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 4  2  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 4  3  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 4  4  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 4  5  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 5  1  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 5  2  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 5  3  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 5  4  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg

## 5  5  survey_pp_CT  survey_pp_ES  survey_pp_SE  survey_pp_SF  survey_pp_TR  colleg


##

##  iter imp variable

## 1  1  survey_pp_CT  survey_pp_SE  college_rate

## 1  2  survey_pp_CT  survey_pp_SE  college_rate

## 1  3  survey_pp_CT  survey_pp_SE  college_rate

## 1  4  survey_pp_CT  survey_pp_SE  college_rate

## 1  5  survey_pp_CT  survey_pp_SE  college_rate

## 2  1  survey_pp_CT  survey_pp_SE  college_rate

## 2  2  survey_pp_CT  survey_pp_SE  college_rate

## 2  3  survey_pp_CT  survey_pp_SE  college_rate

## 2  4  survey_pp_CT  survey_pp_SE  college_rate

## 2  5  survey_pp_CT  survey_pp_SE  college_rate

## 3  1  survey_pp_CT  survey_pp_SE  college_rate

## 3  2  survey_pp_CT  survey_pp_SE  college_rate

## 3  3  survey_pp_CT  survey_pp_SE  college_rate

## 3  4  survey_pp_CT  survey_pp_SE  college_rate
```

```
##   3   5  survey_pp_CT  survey_pp_SE  college_rate

##   4   1  survey_pp_CT  survey_pp_SE  college_rate

##   4   2  survey_pp_CT  survey_pp_SE  college_rate

##   4   3  survey_pp_CT  survey_pp_SE  college_rate

##   4   4  survey_pp_CT  survey_pp_SE  college_rate

##   4   5  survey_pp_CT  survey_pp_SE  college_rate

##   5   1  survey_pp_CT  survey_pp_SE  college_rate

##   5   2  survey_pp_CT  survey_pp_SE  college_rate

##   5   3  survey_pp_CT  survey_pp_SE  college_rate

##   5   4  survey_pp_CT  survey_pp_SE  college_rate

##   5   5  survey_pp_CT  survey_pp_SE  college_rate
```
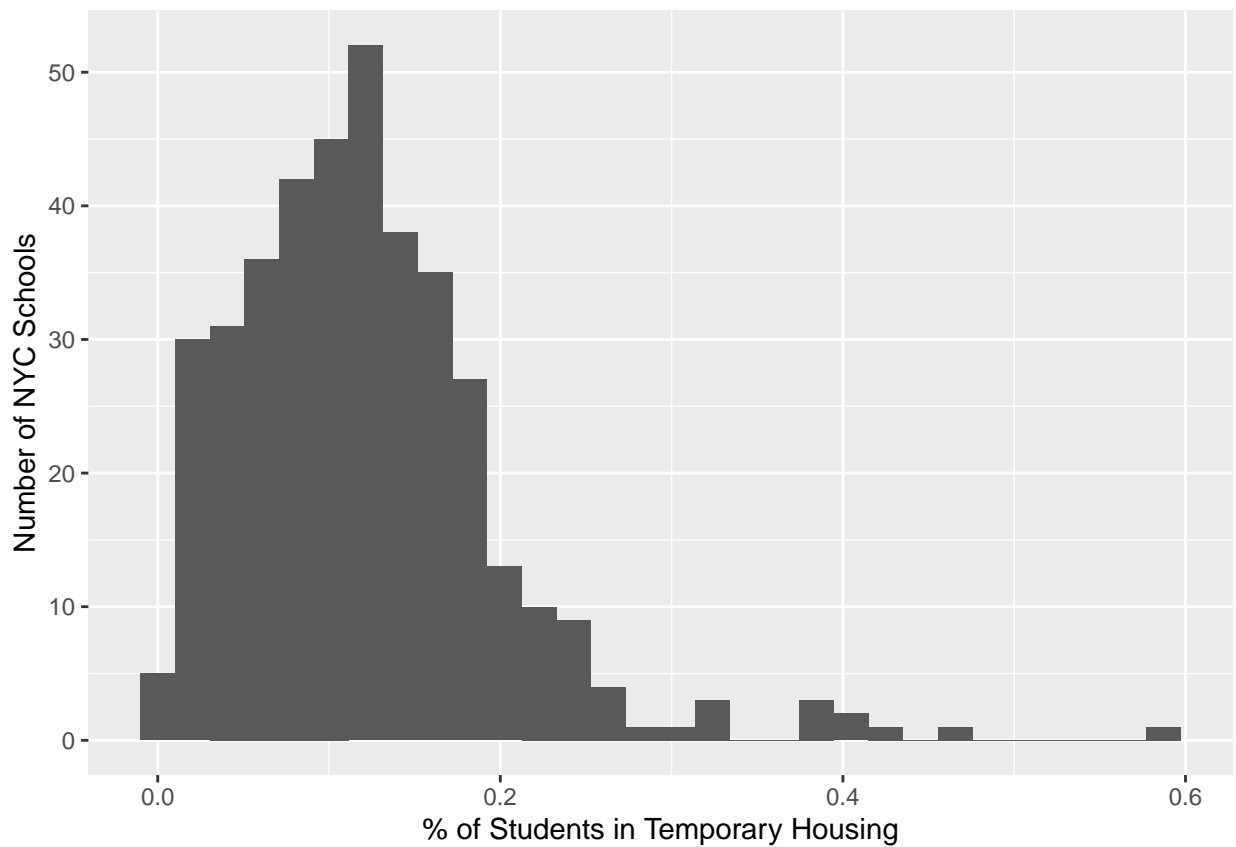
The below plot shows the raw relationship between each survey rating (*Collaborative Teaching*, *Trust*, etc) and the response variables of interest: *Average English/Math SAT scores* per school.

## Experimentation and Results

We can use two variables as a proxy for the school's survey rating in predicting college persistence:

- Percent of Students in Tempoarary Housing (`temp_housing_pct`)
- Economic Need Index (`eni_hs_pct_912`) - this is a measure of the percent of students facing economic hardship at a school *(noauthor_student_2021?)*. This measures the economic hardship faced by students measured along a few criteria:
  - The student is eligible for public assistance from the NYC Human Resources Administration (HRA)
  - The student lived in temporary housing in the past four years
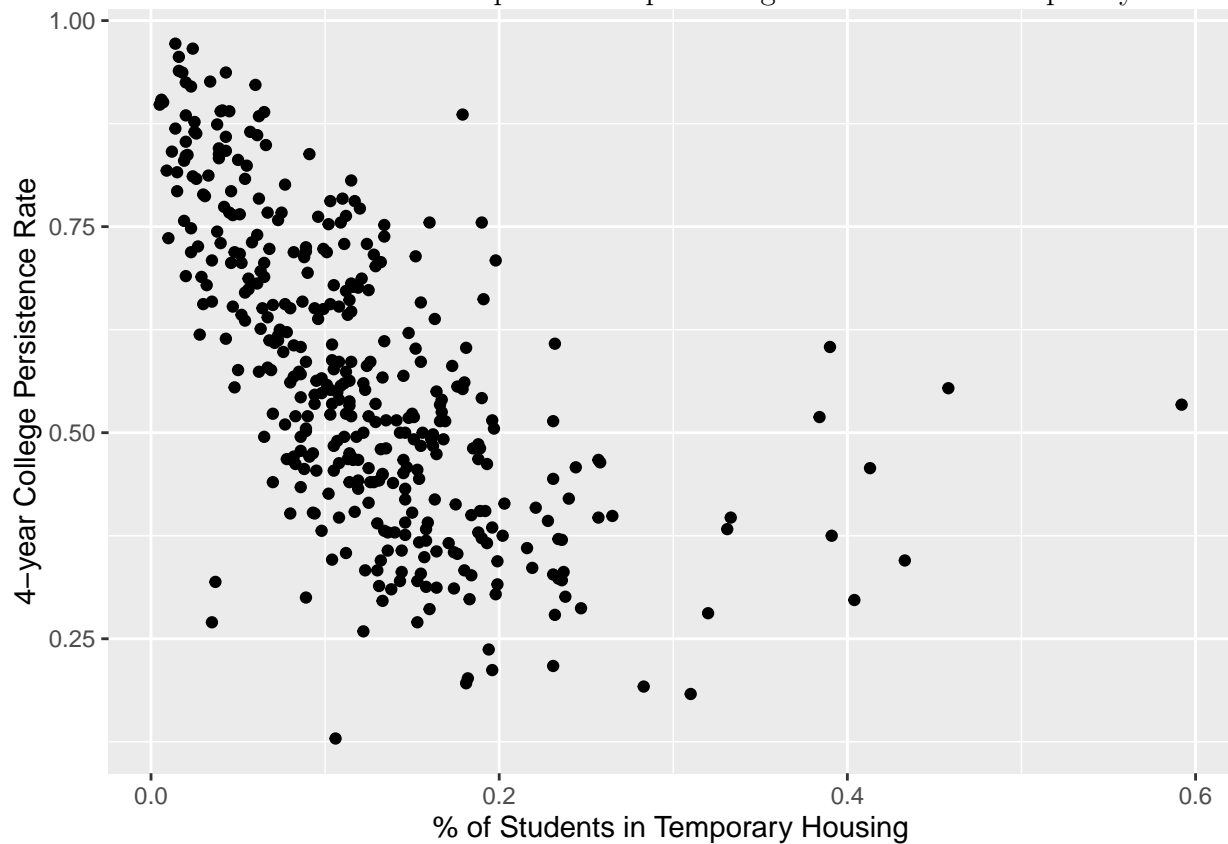  - The student is in high school, has a home language other than English, and

entered the NYC DOE for the first time within the last four years.

## Histogram of train$college_rate



train$college_rate



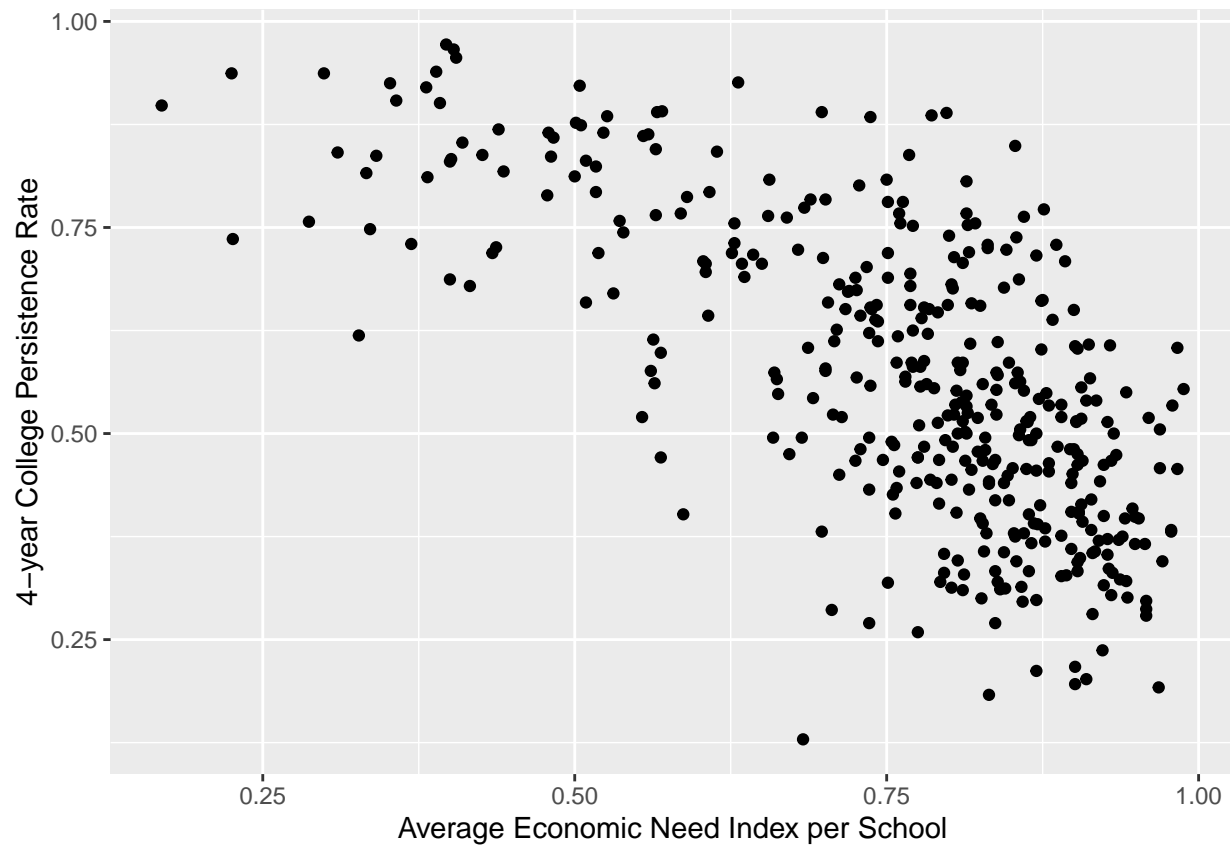% of Students in Temporary Housing

We see this distribution of the percentage of students in temporary housing per school to be skewed left. This will be an important piece of information as we model these relationships later.

First, we should check an assumption of linearity between our predictor and response variables. In this case this a scatter plot of the percentage of students in temporary housing



We see a general linear relationship for schools with lower rates of students in temp housing. However, this linear relationship does **not** visually hold for schools with higher rates of temp housing use.

Plotting the relationship below between a school's economic need index

Again, we see a non-linear relationship between our predictor (*Economic Need Index*) and Outcome Variable (*College Persistence Rate*)
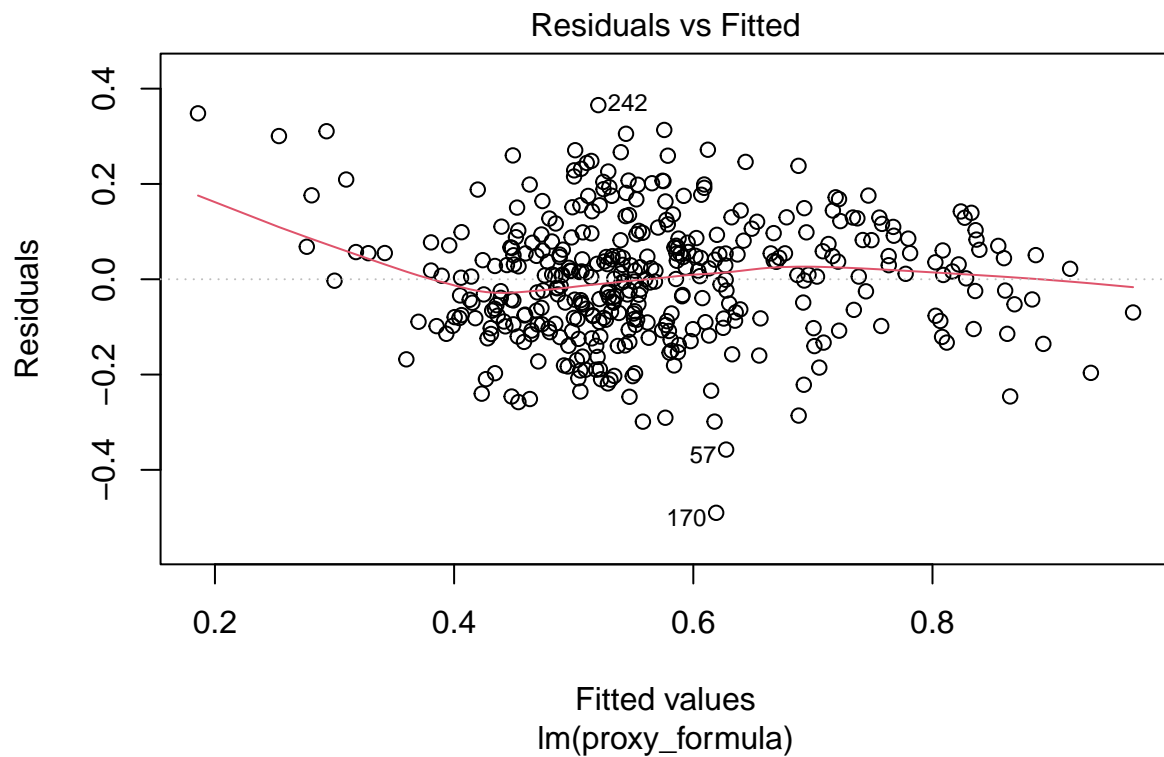
```
##
## Call:
## lm(formula = proxy_formula, data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.49012 -0.08789  0.00611  0.07876  0.36541
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.06702    0.03559  29.982  < 2e-16 ***
```
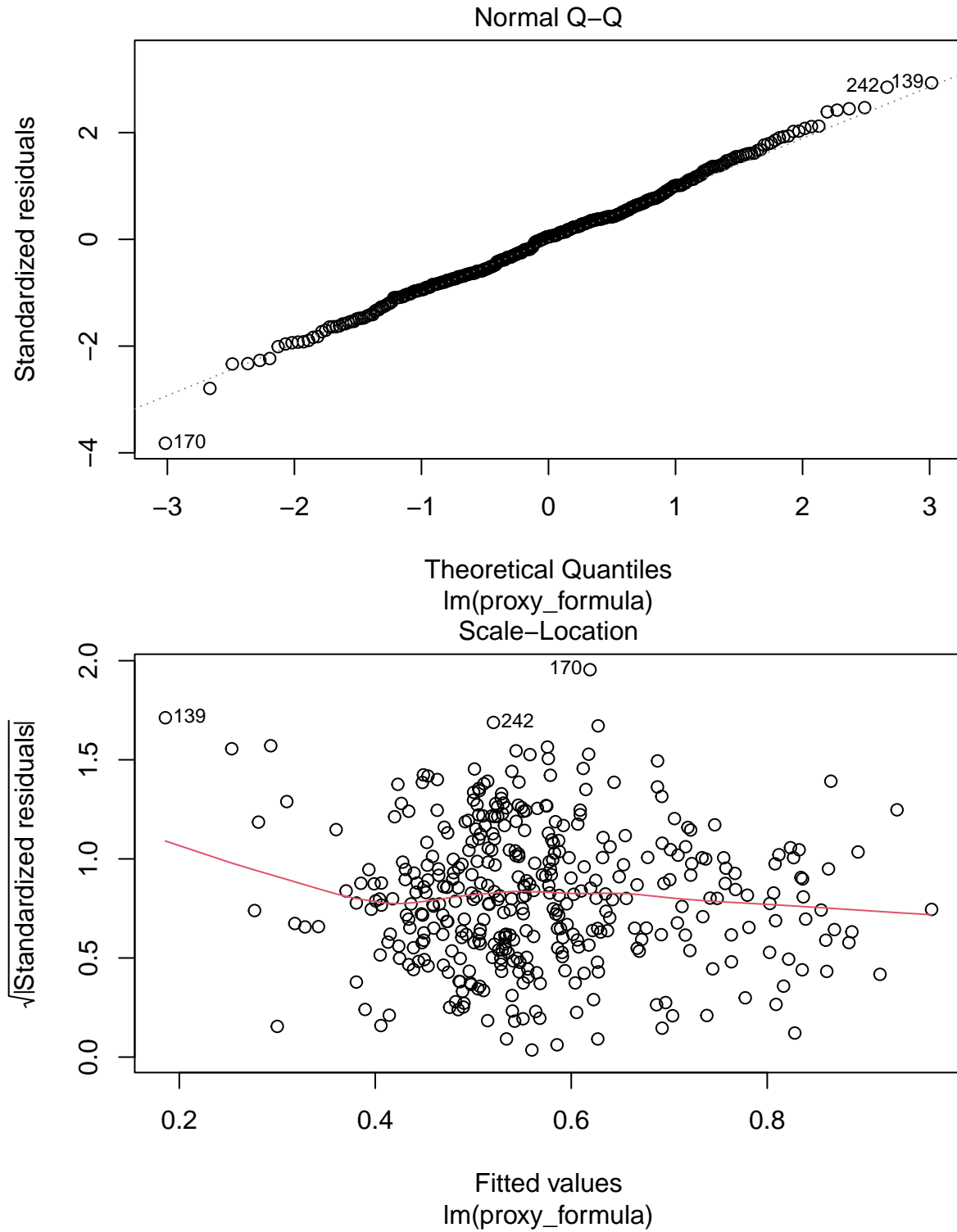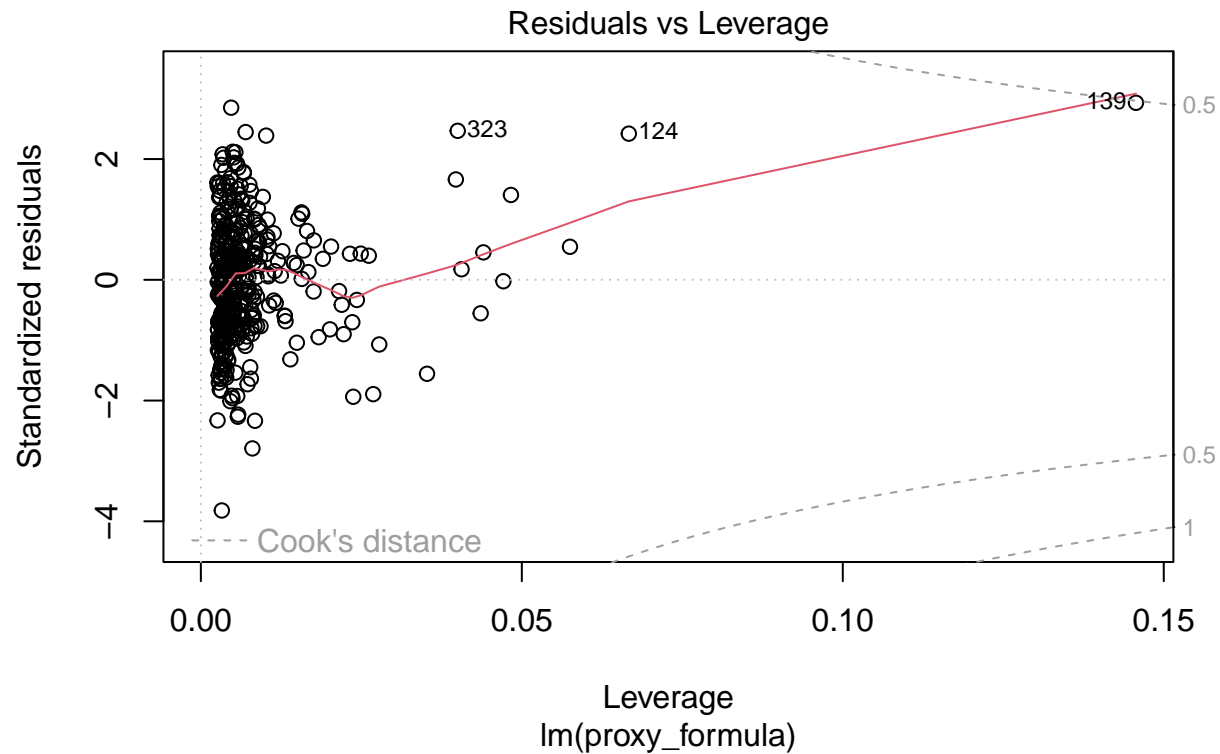
```
## temp_housing_pct -0.54355     0.12098  -4.493  9.3e-06 ***

## economic_need     -0.57142     0.05791  -9.867  < 2e-16 ***

## ---

## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##

## Residual standard error: 0.1285 on 387 degrees of freedom

## Multiple R-squared:  0.493,  Adjusted R-squared:  0.4903

## F-statistic: 188.1 on 2 and 387 DF,  p-value: < 2.2e-16
```



Residuals vs Fitted

Normal Q–Q

lm(proxy_formula)

Scale–Location

Fitted values
lm(proxy_formula)

## Residuals vs Leverage



Given the

```
##
## Call:
## lm(formula = proxy_formula, data = train, weights = weights)
##
## Weighted Residuals:
##     Min      1Q  Median      3Q     Max
## -4.9321 -0.8493  0.0723  0.8035  3.5727
##
## Coefficients:
##                   Estimate Std. Error t value Pr(>|t|)
## (Intercept)        1.05450    0.03252  32.424  < 2e-16 ***
## temp_housing_pct  -0.64311    0.12791  -5.028 7.61e-07 ***
## economic_need     -0.53909    0.05538  -9.734  < 2e-16 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.261 on 387 degrees of freedom
## Multiple R-squared:  0.5205, Adjusted R-squared:  0.518
## F-statistic: 210.1 on 2 and 387 DF,  p-value: < 2.2e-16
```

**Model Evaluation.**

## Conclusion

**TODO**

- Merge/Join in ACT/SAT information by DBN
- Model Selection

# References

Afarian, R., & Kleiner, B. (2003). The relationship between grades and career success. *Management Research News*, *26*, 42–51. https://doi.org/10.1108/01409170310783781

Bernacki, M. L., Chavez, M. M., & Uesbeck, P. M. (2020). Predicting achievement and providing support before STEM majors begin to fail. *Computers & Education*, *158*, 103999. https://doi.org/https://doi.org/10.1016/j.compedu.2020.103999

Education Statistics, N. C. for. (2008). *Percentage of high school dropouts among persons 16 through 24 years old.* Retrieved from https://nces.ed.gov/programs/digest/d08/tables/dt08_110.asp

Musso, M. F., Cascallar, E. C., Bostani, N., & Crawford, M. (2020). Identifying reliable predictors of educational outcomes through machine-learning predictive modeling. *Frontiers in Education*, *5*. https://doi.org/10.3389/feduc.2020.00104

New York City Schools, T. R. A. for. (2018). *Redesigning the Annual NYC School Survey: Lessons from a Research-Practice Partnership.* https://steinhardt.nyu.edu/sites/default/files/2021-01/Lessons_from_a_Research-Practice_Partnership.pdf.

Roth, P. L., BeVier, C. A., Switzer III, F. S., & Schippmann, J. S. (1996). Meta-analyzing the relationship between grades and job performance. *Journal of Applied Psychology*, *81*(5), 548–556. https://doi.org/10.1037/0021-9010.81.5.548

Yağcı, M. (2022). Educational data mining: Prediction of students' academic performance using machine learning algorithms. *Smart Learning Environments*, *9*(1), 11. https://doi.org/10.1186/s40561-022-00192-z

# Appendices

Below is the code used to generate this report. It's also available on GitHub here

```r
knitr::opts_chunk$set(echo = FALSE, warning = FALSE, message = FALSE)
library(tidyverse)
```

```r
library(gridExtra)

library(glue)

library(mice)

# library(autoReg)

library("papaja")

r_refs("r-references.bib")

# Read in our dataset from GitHub

# https://www.opendatanetwork.com/dataset/data.cityofnewyork.us/bm9v-cvch

df <- read.csv("https://data.cityofnewyork.us/api/views/26je-vkp6/rows.csv?date=20231108

label_cols <- c("dbn", "school_name", "school_type")

# Convert needed columns to numeric typing

df <- cbind(df[, label_cols], as.data.frame(lapply(df[,!names(df) %in% label_cols], as.


df$college_rate <- df$val_persist3_4yr_all

df$economic_need <- df$eni_hs_pct_912

set.seed(42)


# Adding a 20% holdout of our input data for model evaluation later

train <- subset(df[sample(1:nrow(df)), ]) %>% sample_frac(0.8)

test  <- dplyr::anti_join(df, train, by = 'dbn')


cols <- c("survey_pp_CT",

          "survey_pp_ES", "survey_pp_SE",

          "survey_pp_SF", "survey_pp_TR",

          "temp_housing_pct", "economic_need",

          "college_rate",

          "val_mean_score_act_math_all",
```

```r
        "val_mean_score_act_engl_all")

train_data <- train[, cols]

imp <- mice(train_data, method="pmm", seed=42)

train <- complete(imp)

test_data <- test[, cols]

imp <- mice(test_data, method="pmm", seed=42)

test <- complete(imp)

hist(train$college_rate)

ggplot(train, aes(x=temp_housing_pct)) + geom_histogram() + labs(x="% of Students in Tem

# Plot temp housing percentage vs college persistence rate

ggplot(train, aes(x=temp_housing_pct, y=college_rate)) + geom_point() + labs(x="% of Stu

# Plot ENI vs college persistence rate

ggplot(train, aes(x=economic_need, y=college_rate)) + geom_point() +

  labs(x="Average Economic Need Index per School", y="4-year College Persistence Rate")

proxy_formula <- college_rate ~ temp_housing_pct + economic_need

proxy_lm <- lm(proxy_formula, train)

summary(proxy_lm)

plot(proxy_lm)

# Calculating weights for WLS

weights <- 1 / lm(abs(proxy_lm$residuals) ~ proxy_lm$fitted.values)$fitted.values^2


#perform weighted least squares regression

wls_model <- lm(proxy_formula, data = train, weights=weights)


summary(wls_model)
```