

# Quality of NYC Schools - Survey Analysis

Sandro Mikautadze

Last compiled on 14/07/2022

# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Introduction</b>                         | <b>3</b> |
| <b>2</b> | <b>Data Cleansing</b>                       | <b>3</b> |
| 2.1      | Initial Remarks on Raw Data . . . . .       | 3        |
| 2.2      | Dataset Loading and Preview . . . . .       | 3        |
| 2.3      | Raw Data Cleaning . . . . .                 | 4        |
| 2.4      | NA Values Inspection . . . . .              | 5        |
| 2.5      | Joining the Datasets . . . . .              | 7        |
| 2.6      | Final Dataset . . . . .                     | 7        |
| <b>3</b> | <b>Data Analysis</b>                        | <b>8</b> |
| 3.1      | Recalling the Goal of the Project . . . . . | 8        |

# 1 Introduction

This is a TLDR. Enjoy!

- Do student, teacher and parent perceptions of NYC school quality appear to be related to demographic and academic success metrics?
- Do students, teachers, and parents have similar perceptions of NYC school quality?

## 2 Data Cleansing

### 2.1 Initial Remarks on Raw Data

In `data\raw-data` 5 files are available: **combined.csv**, **masterfile11\_gened\_final.txt**, **masterfile11\_gened\_final.xlsx**, **masterfile11\_d75\_final.txt** and **masterfile11\_d75\_final.xlsx**.

These files have been downloaded from the following links: - <https://data.cityofnewyork.us/Education/2011-NYC-School-Survey/mnz3-dyi8> [last visited July 7th, 2022] - <https://data.world/dataquest/nyc-schools-data/workspace/file?filename=combined.csv> [last visited July 13th, 2022]

From the *Survey-Data-Dictionary* file in `data\metadata` we can notice that **masterfile11\_gened\_final** and **masterfile11\_d75\_final** differ by a small aspect: **gened** contains information on all community schools, while **d75** from all District 75 schools, that is schools designed to teach and help students with disabilities. As the Dictionary states, “these files display one line of information for each school, by DBN, that includes the response rate for each school, the number of surveys submitted, the size of the eligible survey population at each school, question scores, the percentage of responses selected, and the count of responses selected”.

Both files come with two different formats: *.txt* and *.xlsx*. I decide to work working with *.txt*, because the Excel version requires paid software to be visualized (i.e. Microsoft Excel). Having a look at the *.txt* datasets, we can notice that they are actually saved as tsv (tab separated value) files.

The **combined** dataset has been pre-cleaned as an exercise and contains combined information on different NYC schools based on SAT, AP scores and geographical data.

### 2.2 Dataset Loading and Preview

Importing the **readr** package under **tidyverse**, I will save the datasets as **combined**, **general** and **district**, respectively for **combined.csv**, **masterfile11\_gened\_final.txt** and **masterfile11\_d75\_final.txt**.

```
dim(combined)
```

```
## [1] 479 30
```

```
dim(general)
```

```
## [1] 1646 1942
```

```
dim(district)
```

```
## [1] 56 1773
```

Looking at the Survey Dictionary we can notice that the first columns indicate some characteristics of the school (we'll get into that later). After that, there are some columns that contain aggregate data on the survey. We can identify three groups that responded to the survey: - Students, encoded by **s** - Teachers, encoded by **t** - Parents, encoded by **p**

They were asked questions on 4 main categories: - Safety and Respect, encoded by **saf** - Communication, encoded by **com** - Engagement, encoded by **eng** - Academic expectations, encoded by **aca**

In addition those columns contain at the end a number: 11. We need to be aware of the fact that in the dictionary, that number is 10; so it might represent the year.

**EXAMPLE:** **eng\_p\_11** indicates the engagement score collected in 2011 based on the parent responses.

After the above described columns, we have thousands of columns on the precise survey question and answers.

As far as **combined** goes, we mainly have data on SAT scores with some other info on the different groups of people attending the school, the school's position, the class size, etc. Overall, all these pieces of information might come useful, so I decide to perform no cleaning.

## 2.3 Raw Data Cleaning

Since we don't really care about the specific survey responses that are present in pretty much all columns but the initial ones, I can say that we can exclude them. Moreover, since it would be great to match performance and perception of school quality to the SAT scores, we can exclude Elementary and Middle Schools from the dataset.

```
unique(general$schooltype)
```

```
## [1] "Elementary School"           "Elementary / Middle School"
## [3] "Middle / High School"       "Middle School"
## [5] "High School"               "Elementary / Middle / High School"
## [7] "Early Childhood School"     "YABC"
```

We are going to keep only "High School" rows.

In the **d75** dataset the **schooltype** column has a unique value:

```
unique(district$schooltype)
```

```
## [1] "District 75 Special Education"
```

This value might refer to either elementary school or high school. In this case the **studentsurveyed** column can help us, because, as written in the dictionary, "This field indicates whether or not this school serves any students in grades 6-12". The values that the column takes are the following:

```
unique(district$studentssurveyed)
```

```
## [1] "Yes" "No"
```

Therefore by keeping only the columns with value "Yes" we will only have high schools, which are what we are interested in.

You can find the code of the "reductions" in **src/00-data-processing.r** under the **CLEANING** comment.

```
dim(combined_reduced)
```

```
## [1] 479 26
```

```
dim(general_reduced)
```

```
## [1] 383 23
```

```
dim(district_reduced)
```

```
## [1] 55 23
```

Now we are dealing with a feasible number of variables and they are closer to what we really need. We can combine the data of the survey in a new dataframe, called `survey`.

```
glimpse(survey)
```

```
## Rows: 438
## Columns: 22
## $ dbn      <chr> "01M448", "01M458", "01M509", "01M515", "01M650", "01~
## $ bn       <chr> "M448", "M458", "M509", "M515", "M650", "M696", "M047~
## $ schoolname <chr> "University Neighborhood High School", "Forsyth Satel~
## $ d75      <dbl> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, ~
## $ studentssurveyed <chr> "Yes", "Yes", "Yes", "Yes", "Yes", "Yes", "Yes", "Yes~
## $ schooltype <chr> "High School", "High School", "High School", "High Sc~
## $ saf_p_11  <dbl> 7.9, 8.1, 7.7, 8.3, 9.0, 8.8, 8.9, 7.6, 8.7, 8.0, 7.5~
## $ com_p_11  <dbl> 7.4, 7.0, 7.4, 7.2, 8.4, 8.2, 7.7, 7.0, 8.1, 7.3, 7.1~
## $ eng_p_11  <dbl> 7.2, 6.7, 7.2, 7.4, 8.1, 8.3, 7.9, 6.9, 7.9, 7.1, 6.9~
## $ aca_p_11  <dbl> 7.3, 7.6, 7.3, 7.5, 8.6, 9.1, 8.1, 7.6, 8.3, 7.5, 7.5~
## $ saf_t_11  <dbl> 6.6, 8.5, 6.4, 9.1, 7.6, 8.2, 8.1, 7.3, 8.0, 8.6, 6.6~
## $ com_t_11  <dbl> 5.8, 8.2, 5.3, 7.3, 7.5, 7.4, 6.1, 7.1, 7.7, 8.1, 6.3~
## $ eng_t_11  <dbl> 6.6, 8.9, 6.1, 8.7, 8.3, 7.5, 7.7, 7.8, 7.9, 8.7, 6.8~
## $ aca_t_11  <dbl> 7.3, 8.9, 6.8, 9.1, 8.7, 8.3, 7.2, 7.7, 8.9, 8.9, 7.1~
## $ saf_s_11  <dbl> 6.0, 6.8, 6.4, 8.0, 8.1, 8.3, 7.3, 6.2, 7.4, 7.1, 6.6~
## $ com_s_11  <dbl> 5.7, 6.1, 5.9, 6.3, 6.9, 7.3, 6.3, 5.7, 6.5, 6.5, 6.2~
## $ eng_s_11  <dbl> 6.3, 6.1, 6.4, 7.0, 7.9, 8.0, 7.0, 6.1, 7.3, 7.0, 6.7~
## $ aca_s_11  <dbl> 7.0, 6.8, 7.0, 7.3, 8.4, 8.9, 7.5, 7.2, 7.6, 7.4, 7.5~
## $ saf_tot_11 <dbl> 6.8, 7.8, 6.9, 8.5, 8.3, 8.5, 8.1, 7.0, 7.9, 7.9, 6.9~
## $ com_tot_11 <dbl> 6.3, 7.1, 6.2, 7.0, 7.6, 7.6, 6.7, 6.6, 7.3, 7.3, 6.6~
## $ eng_tot_11 <dbl> 6.7, 7.2, 6.6, 7.7, 8.1, 8.0, 7.5, 6.9, 7.7, 7.6, 6.8~
## $ aca_tot_11 <dbl> 7.2, 7.8, 7.0, 8.0, 8.6, 8.7, 7.6, 7.5, 8.2, 8.0, 7.4~
```

## 2.4 NA Values Inspection

To better clean the data we can have a look at columns with NA values.

```
colSums(is.na(combined_reduced))
```

```
##          dbn          school_name num.of.sat.test.takers
##          0              0              57
##    avg_sat_score    ap.test.takers    total.exams.taken
##          57              0              247
##    exams_per_student  high_score_percent    avg_class_size
##          247              328              44
##          frl_percent    total_enrollment    ell_percent
##          41              41              41
##          sped_percent    selfcontained_num    asian_per
##          41              51              41
##          black_per    hispanic_per    white_per
##          41              41              41
##          male_per    female_per    total.cohort
##          41              41              89
##          grads_percent    dropout_percent    boro
##          111              111              109
##          lat              long
##          109              109
```

```
colSums(is.na(survey))
```

```
##          dbn          bn          schoolname          d75
##          0              0              0              0
## studentssurveyed    schooltype    saf_p_11    com_p_11
##          0              0              0              0
##          eng_p_11    aca_p_11    saf_t_11    com_t_11
##          0              0              0              0
##          eng_t_11    aca_t_11    saf_s_11    com_s_11
##          0              0              3              3
##          eng_s_11    aca_s_11    saf_tot_11    com_tot_11
##          3              3              0              0
##          eng_tot_11    aca_tot_11
##          0              0
```

The first thing that we can notice is that the `highschool` column in the survey dataframe has 424 NA values, out of 438 observations. This means that that column is pretty much unusable, so we will delete it.

In addition, `combined_reduced` has `number.of.exams.with.scores.3.4.or.5` and `high_score_percent` with 328 NA values, which is more than half of the rows in the dataset. So, it is safe to say that those columns are useless and we will delete them.

The final dimensions of the cleaned datasets are the following:

```
dim(combined_reduced)
```

```
## [1] 479 26
```

```
dim(survey)
```

```
## [1] 438 22
```

## 2.5 Joining the Datasets

Now that the necessary cleaning has been done, we can finally join `survey` and `combined_reduced` into one dataset, that we are going to be using for the analysis.

We are going to apply a `left_join` to `combined_reduced` so that we will have all values for schools of which we have SAT data. We will save it as `school_data_raw`. These are the initial dimensions:

```
dim(school_data_raw)
```

```
## [1] 479 47
```

We can eliminate some redundant columns, such `bn` and `schoolname`. In addition, we now know that we are dealing with high schools, so we can drop `schooltype` and `studentssurveyed`.

## 2.6 Final Dataset

Therefore our final cleaned dataset, named `school_data` is the following:

```
glimpse(school_data)
```

```
## Rows: 479
## Columns: 43
## $ dbn                <chr> "01M292", "01M448", "01M450", "01M458", "01M509~
## $ school_name        <chr> "HENRY STREET SCHOOL FOR INTERNATIONAL STUDIES"~
## $ num.of.sat.test.takers <int> 29, 91, 70, 7, 44, 112, 159, 18, 130, 16, 62, 5~
## $ avg_sat_score       <int> 1122, 1172, 1149, 1174, 1207, 1205, 1621, 1246,~
## $ ap.test.takers      <dbl> 2.5, 39.0, 19.0, 2.5, 2.5, 24.0, 255.0, 2.5, 2.~
## $ total.exams.taken   <int> NA, 49, 21, NA, NA, 26, 377, NA, NA, NA, NA, NA~
## $ exams_per_student   <dbl> NA, 1.256410, 1.105263, NA, NA, 1.083333, 1.478~
## $ high_score_percent  <dbl> NA, 20.408163, NA, NA, NA, 92.307692, 50.663130~
## $ avg_class_size      <int> 23, 22, 21, 23, 24, 23, 26, 22, 21, 16, 23, 15,~
## $ frl_percent         <dbl> 88.6, 71.8, 71.8, 72.8, 80.7, NA, 23.0, 69.8, 1~
## $ total_enrollment    <int> 422, 394, 598, 224, 367, NA, 1613, 218, 617, 17~
## $ ell_percent         <dbl> 22.3, 21.1, 5.0, 4.0, 11.2, NA, 0.2, 3.2, 0.2, ~
## $ sped_percent        <dbl> 24.9, 21.8, 26.4, 8.9, 25.9, NA, 2.7, 6.9, 0.8,~
## $ selfcontained_num   <int> 35, 10, 19, 0, 36, NA, 0, 0, 0, 10, 4, 2, 17, 3~
## $ asian_per           <dbl> 14.0, 29.2, 9.7, 2.2, 9.3, NA, 27.8, 0.5, 15.1,~
## $ black_per           <dbl> 29.1, 22.6, 23.9, 34.4, 31.6, NA, 11.7, 45.4, 1~
## $ hispanic_per        <dbl> 53.8, 45.9, 55.4, 59.4, 56.9, NA, 14.2, 49.5, 1~
## $ white_per           <dbl> 1.7, 2.3, 10.4, 3.6, 1.6, NA, 44.9, 4.1, 49.8, ~
## $ male_per            <dbl> 61.4, 57.4, 54.7, 43.3, 46.3, NA, 49.2, 39.9, 3~
## $ female_per          <dbl> 38.6, 42.6, 45.3, 56.7, 53.7, NA, 50.8, 60.1, 6~
## $ total_cohort        <int> 78, 124, 90, NA, 84, 193, 46, 89, 139, 25, 102,~
## $ grads_percent       <dbl> 55.1, 42.7, 77.8, NA, 56.0, 54.4, 100.0, 55.1, ~
## $ dropout_percent     <dbl> 14.1, 16.1, 5.6, NA, 6.0, 18.1, 0.0, 6.7, 0.7, ~
## $ boro                <chr> "Manhattan", "Manhattan", "Manhattan", NA, "Man~
## $ lat                 <dbl> 40.71376, 40.71233, 40.72978, NA, 40.72057, NA,~
## $ long                <dbl> -73.98526, -73.98480, -73.98304, NA, -73.98567,~
## $ d75                 <dbl> NA, 0, NA, 0, 0, 0, NA, 0, 0, 0, 0, 0, 0, 0,~
## $ saf_p_11            <dbl> NA, 7.9, NA, 8.1, 7.7, 8.3, NA, 9.0, 8.8, 8.9, ~
## $ com_p_11            <dbl> NA, 7.4, NA, 7.0, 7.4, 7.2, NA, 8.4, 8.2, 7.7, ~
```

```
## $ eng_p_11          <dbl> NA, 7.2, NA, 6.7, 7.2, 7.4, NA, 8.1, 8.3, 7.9, ~
## $ aca_p_11          <dbl> NA, 7.3, NA, 7.6, 7.3, 7.5, NA, 8.6, 9.1, 8.1, ~
## $ saf_t_11          <dbl> NA, 6.6, NA, 8.5, 6.4, 9.1, NA, 7.6, 8.2, 8.1, ~
## $ com_t_11          <dbl> NA, 5.8, NA, 8.2, 5.3, 7.3, NA, 7.5, 7.4, 6.1, ~
## $ eng_t_11          <dbl> NA, 6.6, NA, 8.9, 6.1, 8.7, NA, 8.3, 7.5, 7.7, ~
## $ aca_t_11          <dbl> NA, 7.3, NA, 8.9, 6.8, 9.1, NA, 8.7, 8.3, 7.2, ~
## $ saf_s_11          <dbl> NA, 6.0, NA, 6.8, 6.4, 8.0, NA, 8.1, 8.3, 7.3, ~
## $ com_s_11          <dbl> NA, 5.7, NA, 6.1, 5.9, 6.3, NA, 6.9, 7.3, 6.3, ~
## $ eng_s_11          <dbl> NA, 6.3, NA, 6.1, 6.4, 7.0, NA, 7.9, 8.0, 7.0, ~
## $ aca_s_11          <dbl> NA, 7.0, NA, 6.8, 7.0, 7.3, NA, 8.4, 8.9, 7.5, ~
## $ saf_tot_11        <dbl> NA, 6.8, NA, 7.8, 6.9, 8.5, NA, 8.3, 8.5, 8.1, ~
## $ com_tot_11        <dbl> NA, 6.3, NA, 7.1, 6.2, 7.0, NA, 7.6, 7.6, 6.7, ~
## $ eng_tot_11        <dbl> NA, 6.7, NA, 7.2, 6.6, 7.7, NA, 8.1, 8.0, 7.5, ~
## $ aca_tot_11        <dbl> NA, 7.2, NA, 7.8, 7.0, 8.0, NA, 8.6, 8.7, 7.6, ~
```

You can find the cleaned dataset in `data/clean-data/school-data.csv`.

## 3 Data Analysis

### 3.1 Recalling the Goal of the Project