# Sunny Futures STEM

Anterra Kennedy
Sasha Prokhorova
Nick Horton

# Table of Contents

# The Goal...

To create long-term improvement in the lives of underprivileged kids by holding fun science demonstrations at outreach events to get them excited about and engaged with STEM.

Sunny Futures STEM

# How can we reach the greatest total number of at-risk kids?

p / 1

# Methodology:

**01.** Identify underperforming schools

**02.** Find nearest subway stations

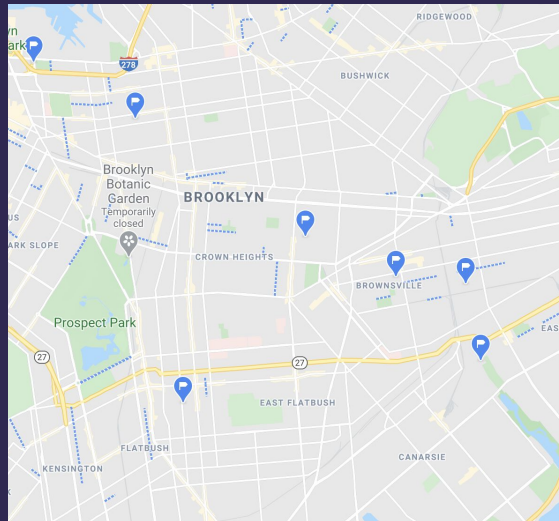**03.** Analyze which stations have the highest weekday afternoon traffic

# Sunny Futures

# Where are at-risk kids?
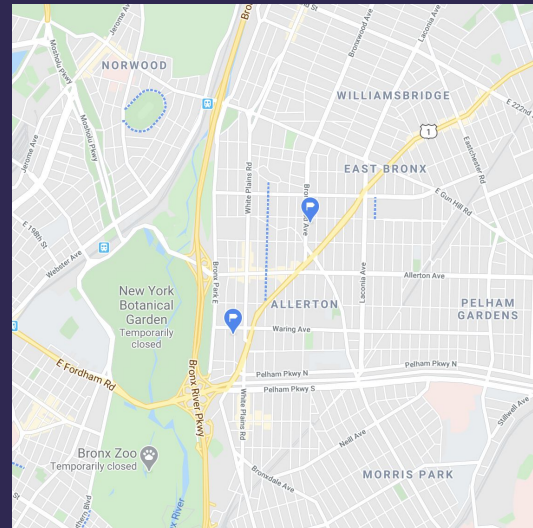
We used NYC Dept of Education's School Performance Data...

... to find the top 20 neediest and worst-performing schools on the metrics of student achievement and economic need index.

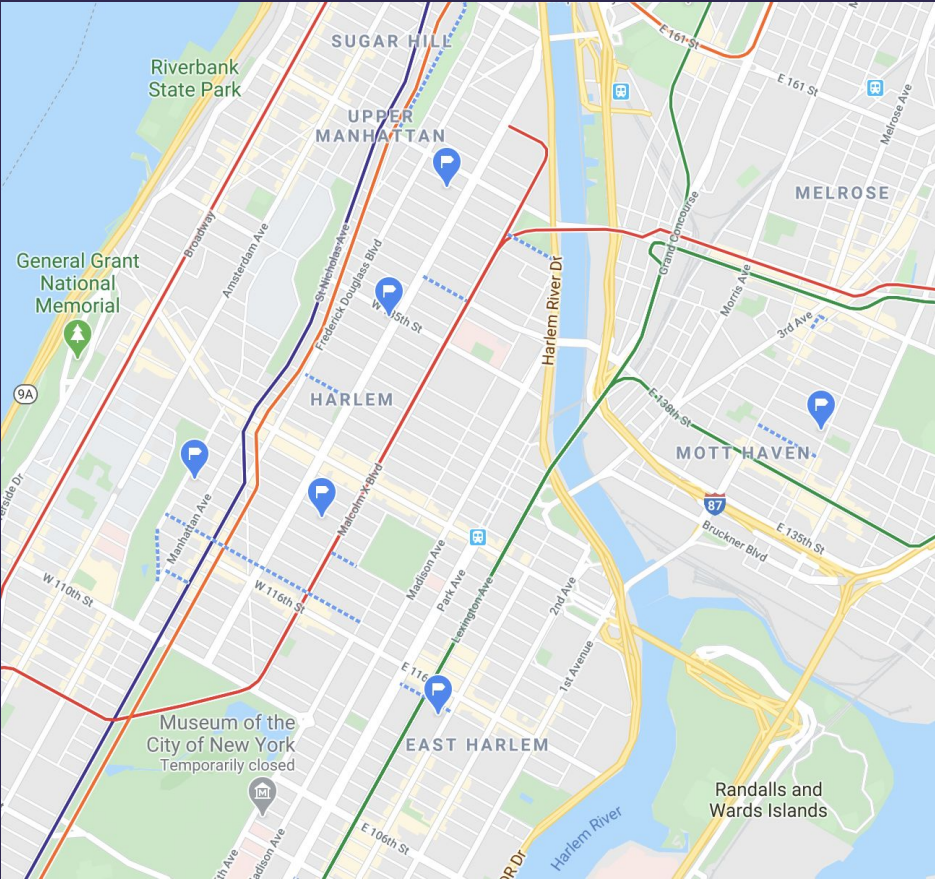| | School Name |
|---|---|
| 296 | New Directions Secondary School |
| 210 | P.S./I.S. 224 |
| 145 | P.S. 194 Countee Cullen |
| 7 | P.S. 034 Franklin D. Roosevelt |
| 641 | School of the Future Brooklyn |
| 114 | James Weldon Johnson |
| 757 | P.S./I.S. 323 |
| 628 | The Fresh Creek School |
| 462 | P.S. 287 Bailey K. Ashford |
| 460 | P.S. 270 Johann DeKalb |
| 577 | M.S. K394 |
| 321 | P.S. 051 Bronx New School |
| 580 | P.S. 399 Stanley Eugene Clark |
| 94 | P.S. 242 - The Young Diplomats Magnet Academy |
| 377 | P.S. 096 Richard Rodgers |
| 154 | Thurgood Marshall Academy for Learning and Social |
| 372 | P.S. 076 The Bennington School |
| 737 | P.S. 251 Paerdegat |
| 965 | P.S. 052 Queens |
| 963 | Cynthia Jenkins School |
| 970 | P.S. 118 Lorraine Hansberry |
| 89 | P.S. 180 Hugo Newman |
| 962 | P.S. 036 Saint Albans School |
| 602 | P.S. 276 Louis Marshall |

Many under-performing schools are nearby...

Sunny Futures

East Brooklyn

East Bronx

Harlem

p / 4

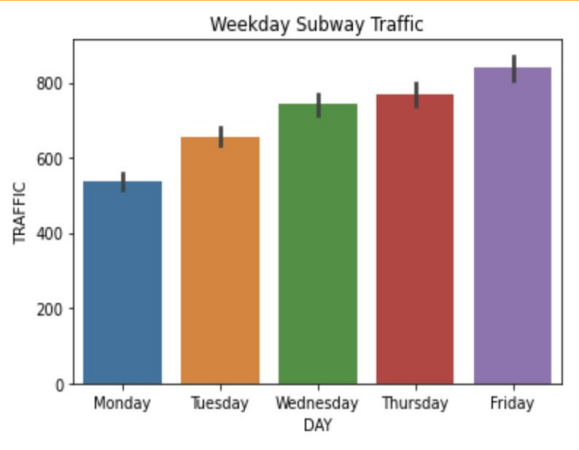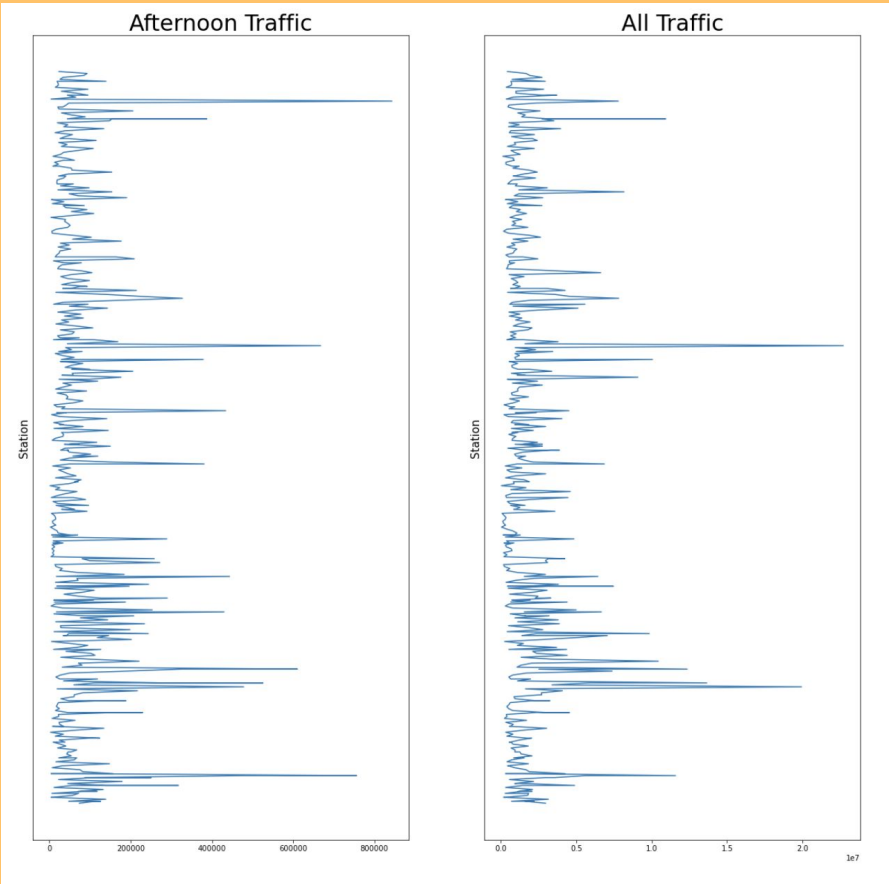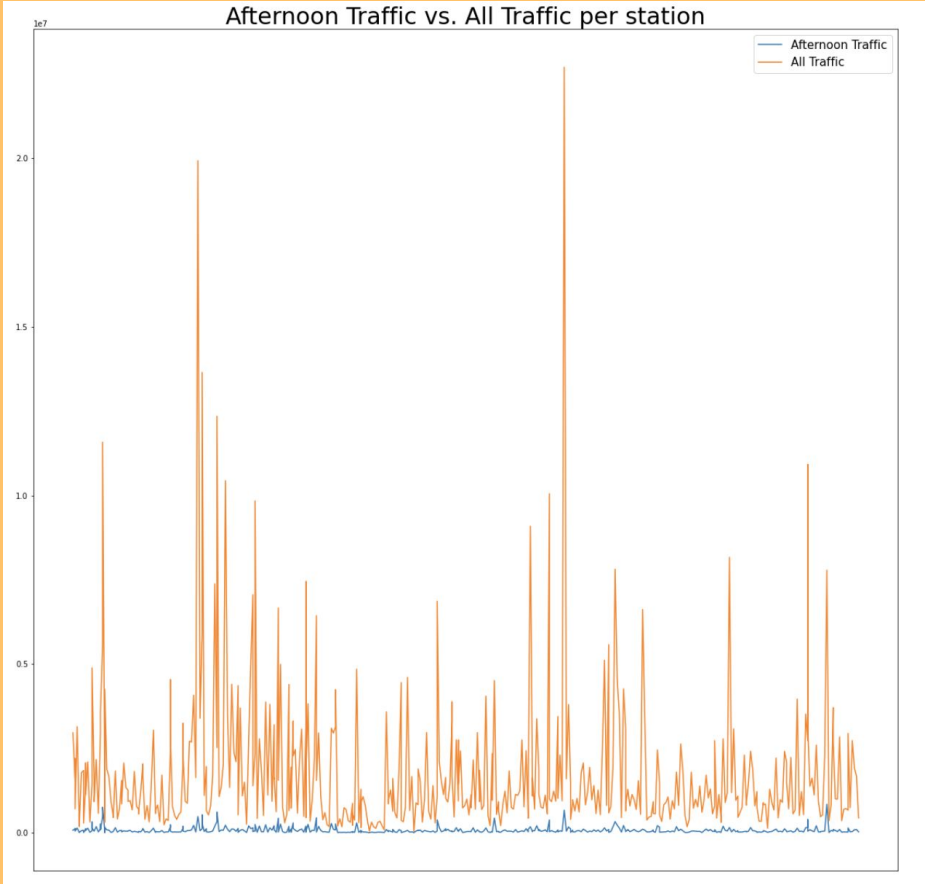| | School Name | Nearest Station | Second Nearest Station |
|---|---|---|---|
| 296 | New Directions Secondary School | 170 ST | N/A |
| 210 | P.S./I.S. 224 | BROOK AV | 3 AV-149 ST |
| 145 | P.S. 194 Countee Cullen | 145 ST | 145 ST |
| 7 | P.S. 034 Franklin D. Roosevelt | 1 AV | N/A |
| 641 | School of the Future Brooklyn | LIVONIA AV | PENNSYLVANIA AV |
| 114 | James Weldon Johnson | 116 ST | N/A |
| 757 | P.S./I.S. 323 | ROCKAWAY AV | N/A |
| 628 | The Fresh Creek School | EAST 105 ST | NEW LOTS AV |
| 462 | P.S. 287 Bailey K. Ashford | YORK ST | N/A |
| 460 | P.S. 270 Johann DeKalb | CLASSON AV | N/A |
| 577 | M.S. K394 | UTICA AV | N/A |
| 321 | P.S. 051 Bronx New School | 42 ST-PORT AUTH | N/A |
| 580 | P.S. 399 Stanley Eugene Clark | CHURCH AV | BEVERLEY ROAD |
| 94 | P.S. 242 - The Young Diplomats Magnet Academy | 125 ST | 125 ST |
| 377 | P.S. 096 Richard Rodgers | PELHAM PKWY | ALLERTON AV |
| 154 | Thurgood Marshall Academy for Learning and Social | 135 ST | 135 ST |
| 372 | P.S. 076 The Bennington School | BURKE AV | N/A |
| 89 | P.S. 180 Hugo Newman | 125 ST | 116 ST |

... and close to subway stations.

p / 5

| | Neighborhood | School Name | Nearest Station | Nearest Station Lines | Second Nearest Station | Second Nearest Station Lines |
|---|---|---|---|---|---|---|
| 0 | East Bronx | P.S. 076 The Bennington School | BURKE AV | 25 | N/A | N/A |
| 1 | East Bronx | P.S. 096 Richard Rodgers | PELHAM PKWY | 25 | ALLERTON AV | 25 |
| 2 | East Brooklyn | M.S. K394 | UTICA AV | AC | N/A | N/A |
| 3 | East Brooklyn | P.S./I.S. 323 | ROCKAWAY AV | 3 | N/A | N/A |
| 4 | East Brooklyn | School of the Future Brooklyn | LIVONIA AV | L | PENNSYLVANIA AV | 3 |
| 5 | East Brooklyn | The Fresh Creek School | EAST 105 ST | L | NEW LOTS AV | 3 |
| 6 | Harlem | James Weldon Johnson | 116 ST | 6 | N/A | N/A |
| 7 | Harlem | P.S. 180 Hugo Newman | 125 ST | ABCD | 116 ST | BC |
| 8 | Harlem | P.S. 194 Countee Cullen | 145 ST | 3 | 145 ST | ABCD |
| 9 | Harlem | P.S. 242 - The Young Diplomats Magnet Academy | 125 ST | 23 | 125 ST | ACBD |
| 10 | Harlem | P.S./I.S. 224 | BROOK AV | 6 | 3 AV-149 ST | 25 |
| 11 | Harlem | Thurgood Marshall Academy for Learning and Social | 135 ST | 23 | 135 ST | BC |
| 12 | other | New Directions Secondary School | 170 ST | BD | N/A | N/A |
| 13 | other | P.S. 034 Franklin D. Roosevelt | 1 AV | L | N/A | N/A |
| 14 | other | P.S. 051 Bronx New School | 42 ST-PORT AUTH | ACENQRS1237W | N/A | N/A |
| 15 | other | P.S. 270 Johann DeKalb | CLASSON AV | G | N/A | N/A |
| 16 | other | P.S. 287 Bailey K. Ashford | YORK ST | F | N/A | N/A |
| 17 | other | P.S. 399 Stanley Eugene Clark | CHURCH AV | 25 | BEVERLEY ROAD | BQ |

# Stations with duplicate names service different subway lines.

# MTA Turnstile Data
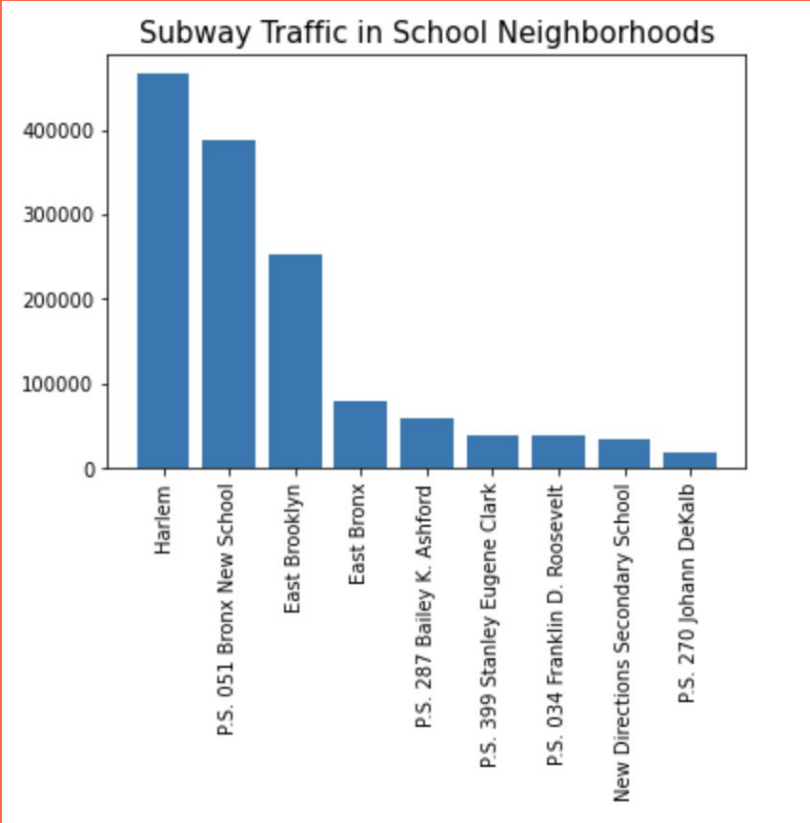
**We filtered the data to isolate weekday afternoons...**



Afternoon Traffic vs. All Traffic per station



Afternoon Traffic | All Traffic



Weekday Subway Traffic

**... which reflects students after-school commute.**

**Subway Traffic Near Schools**

## Result: Harlem
### has the most subway traffic

| | Label | School Traffic |
|---|---|---|
| 2 | Harlem | 465604.0 |
| 2 | P.S. 051 Bronx New School | 388395.0 |
| 1 | East Brooklyn | 251804.0 |
| 0 | East Bronx | 78781.0 |
| 4 | P.S. 287 Bailey K. Ashford | 59548.0 |
| 5 | P.S. 399 Stanley Eugene Clark | 39292.0 |
| 1 | P.S. 034 Franklin D. Roosevelt | 38513.0 |
| 0 | New Directions Secondary School | 34523.0 |
| 3 | P.S. 270 Johann DeKalb | 18718.0 |



**Subway Traffic in School Neighborhoods**

# Where is ideal?

**We want a station with high individual traffic...**

...that can also reach students from more than just the one nearest school.



Traffic at Stations Near Harlem Schools

# All stations near the schools are servicing *the same 3 subway lines.*



Total traffic by Subway Line in Harlem

Students are getting on the same subways at consecutive stops.

# Ideal event placement:



Traffic at Station Near Harlem Schools

**One event at the busiest station on each line:**

- 145th St (ABCD)
- 125 St (23)
- 116 St (6)



p / 11

# Conclusion:

## 3 events, flyers in the week preceding each, and free pizza

This way, we can draw kids from all stops along each line. The pizza will incentivize them to get off their subway 1 or 2 stops up.

# Appendix

future work

# 01

## Future Work
## Areas of Interest

- Expanding into other cities
- Use geo-spatial data to automate

## Advanced Methods

For the purposes of this project, we stuck to simple EDA (exploratory data analysis) and data cleaning. We wanted to analyze clusters of schools based on multiple factors like distance from nearby schools/stations and economic neediness of the school vs its proximity to stations, possibly using regression. There were plenty of unexpected quirks with the data and trial/error with pandas syntax., but these methods could potentially be implemented later on.

# 02

# Navigating Errors...

## Total Station Ridership by Line

| | STATION | LINENAME | TRAFFIC |
|---|---|---|---|
| 310 | GRD CNTRL-42 ST | 4567S | 22706169.0 |
| 88 | 34 ST-HERALD SQ | BDFMNQRW | 19936316.0 |
| 92 | 34 ST-PENN STA | ACE | 13653526.0 |
| 101 | 42 ST-PORT AUTH | ACENQRS1237W | 12353636.0 |
| 28 | 14 ST-UNION SQ | LNQR456W | 11583395.0 |
| 441 | TIMES SQ-42 ST | 1237ACENQRSW | 10921827.0 |
| 105 | 47-50 STS ROCK | BDFM | 10440520.0 |
| 300 | FULTON ST | 2345ACJZ | 10056623.0 |
| 124 | 59 ST COLUMBUS | ABCD1 | 9842010.0 |
| 287 | FLUSHING-MAIN | 7 | 9091166.0 |
| 398 | PATH NEW WTC | 1 | 8170652.0 |
| 91 | 34 ST-PENN STA | 123ACE | 7859477.0 |
| 337 | JKSN HT-ROOSVLT | EFMR7 | 7816060.0 |
| 452 | W 4 ST-WASH SQ | ABCDEFM | 7786708.0 |
| 155 | 86 ST | 456 | 7452296.0 |
| 99 | 42 ST-BRYANT PK | BDFM7 | 7381516.0 |
| 120 | 59 ST | 456NQRW | 7059397.0 |
| 233 | CANAL ST | JNQRZ6W | 6865017.0 |
| 137 | 72 ST | 123 | 6665123.0 |
| 353 | LEXINGTON AV/53 | EM6 | 6621575.0 |

Later on during analysis, something strange became evident... Ridership at Penn Station for lines 123 and ACE together was lower than ridership for line ACE alone, by about 2x.

I went back and checked the data, it turns out someone made an error when compiling the data – I'm guessing the entry for Line 123 + ACE was switched with the entry for Line ACE. If we swapped those, the numbers make sense.

| LINENAME | TRAFFIC |
|---|---|
| 123 | 5899837.0 |
| 123ACE | 7859477.0 |
| ACE | 13653526.0 |

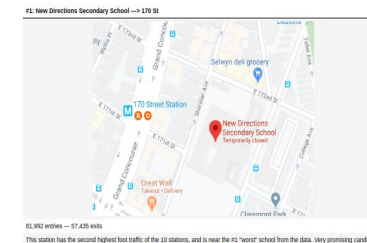5899837.0+7859477.0 = 13759314
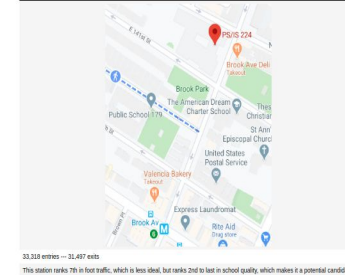
# Sunny Futures

## Earlier Iterations

These were some of our original forays into the data.

We were only using one week of data, compared with the 3 months of data we used in the final project. So we had the right idea with our analysis, but the incomplete data skewed our results.

Ak    Sp    Nh    ;)

# References

- https://infohub.nyced.org/reports/school-quality/school-quality-reports-and-resources

- https://infohub.nyced.org/docs/default-source/default-document-library/201819_ems_sqr_results.xlsx

- http://web.mta.info/developers/turnstile.html

- https://towardsdatascience.com/mta-turstile-data-my-first-taste-of-a-data-science-project-493b03f1708a