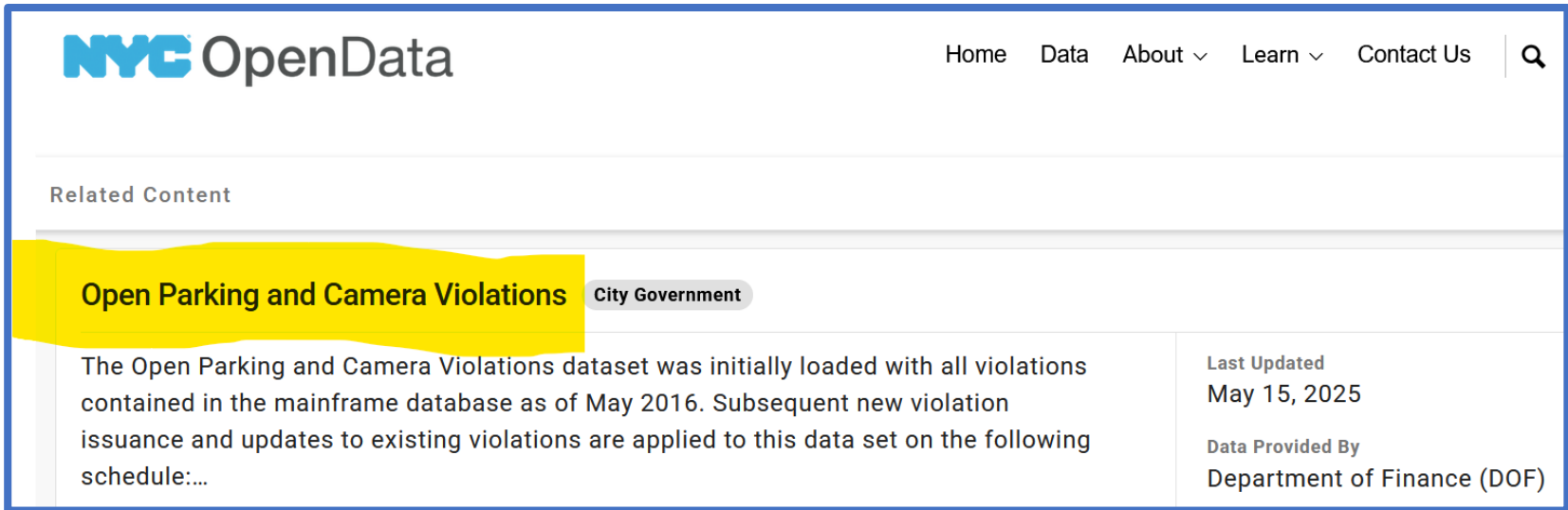# Group 9A – Matt R

Was separated from Group 9 at the ETL process milestone.
Took over KPI #2 which was exploring NYC parking ticket data.
Data was downloaded via manual query from NYC OpenData Open Parking and Camera
Violations Database:
https://data.cityofnewyork.us/City-Government/Open-Parking-and-Camera-
Violations/nc67-uf89/about_data



Data was queried for the months June to December, 2023 in April and May of 2025.
Total amount of rows was approximately 70 million.
Each month contained between 9 million and 11 million rows.

Tools used:

- Extraction: query tool on the NYC Open Data website

- Transform: OpenRefine

OpenRefine is an open-source, freely available tool that allows to cleaning of dirty and inconsistent data on thousands to millions of rows on a local desktop computer.

## OpenRefine

OpenRefine is a powerful free, open source tool for working with messy data: cleaning it; transforming it from one format into another; and extending it with web services and external data.

Our goal is to empower everyone to meaningfully engage with data by providing an accessible open source tool and nurturing a diverse, supportive community.

**Download**

www.OpenRefine.org

### Main features

**Faceting**
Drill through large datasets using facets and apply operations on filtered views of your dataset.

**Clustering**
Fix inconsistencies by merging similar values thanks to powerful heuristics.

**Reconciliation**
Match your dataset to external databases via reconciliation services.

**Infinite undo/redo**
Rewind to any previous state of your dataset and replay your operation history on a new version of it.

**Privacy**
Your data is cleaned on your machine, not in some dubious data laundering cloud.

**Wikibase**
Contribute to Wikidata, the free knowledge base anyone can edit, and other Wikibase instances.

OpenRefine  300 000 July16 July31 2023Open Parking and Camera Violations 20250416 csv  Permalink      Open...  Export ▾  Help

**Facet / Filter**    Undo / Redo 0 / 0    ‹

280,137 rows                                                            Extensions  Wikibase ▾

Show as: **rows** records    Show: 5 **10** 25 50 100 500 1000 rows    « first  ‹ previous  1  - 10  next ›  last »

### Using facets and filters

Use facets and filters to select subsets of your data to act on. Choose facet and filter methods from the menus at the top of each data column.

Not sure how to get started?
**Watch these screencasts**

| All | State | Licer | Issue Date | Viola | Violation | Fine | Pena | Payn | Amo | Prec | County | Ranc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1. | MA | PAS | 2023-07-16T00:00:00Z | 12 | NO STANDING-EXC. AUTH. VEHICLE | 95.0 | 0.00 | 95.00 | 0.00 | 084 | Brooklyn | true | POL |
| 2. | NY | PAS | 2023-07-16T00:00:00Z | 12 | NO STANDING-EXC. AUTH. VEHICLE | 95.0 | 0.00 | 0.00 | 0.00 | 019 | Manhattan | true | TRA |
| 3. | TN | PAS | 2023-07-16T00:00:00Z | 11 | PHTO SCHOOL ZN SPEED VIOLATION | 50.0 | 0.00 | 50.00 | 0.00 | 000 | Brooklyn | true | DEF |
| 4. | NY | PAS | 2023-07-16T00:00:00Z | 01 | REG. STICKER-EXPIRED/MISSING | 65.0 | 0.00 | 65.00 | 0.00 | 043 | Bronx | true | TRA |
| 5. | NY | PAS | 2023-07-16T00:00:00Z | 02 | NO STANDING-DAY/TIME LIMITS | 115.0 | 0.00 | 115.00 | 0.00 | 061 | Brooklyn | true | TRA |
| 6. | NY | PAS | 2023-07-16T00:00:00Z | 08 | FIRE HYDRANT | 115.0 | 0.00 | 115.00 | 0.00 | 122 | Staten Island | true | TRA |
| 7. | NY | PAS | 2023-07-16T00:00:00Z | 03 | NO STANDING-DAY/TIME LIMITS | 115.0 | 0.00 | 115.00 | 0.00 | 001 | Manhattan | true | TRA |
| 8. | NY | PAS | 2023-07-16T00:00:00Z | 06 | NO STANDING-DAY/TIME LIMITS | 115.0 | 0.00 | 115.00 | 0.00 | 077 | Brooklyn | true | TRA |
| 9. | NY | PAS | 2023-07-16T00:00:00Z | 12 | NO PARKING-DAY/TIME LIMITS | 65.0 | 0.00 | 65.00 | 0.00 | 006 | Manhattan | true | TRA |
| 10. | NY | PAS | 2023-07-16T00:00:00Z | 06 | REG STICKER-MUTILATED/C'FEIT | 65.0 | 0.00 | 65.00 | 0.00 | 045 | Bronx | true | TRA |

## - Loading/Visualizations: Orange Data Mining
Orange Data Mining is a powerful open source tool for visualizing various aspects of data without having to use complicated Python programming.



www.OrangeDataMining.com

- Within OpenRefine, used based GREL [General Refine Expression Language]
GREL is a query language that can modify text and perform numerical calculations.

**Custom facet on column Violation**

Expression          Language | General Refine Expression Language (GREL) ▾ |

`row.index % 3 == 0`                                    No syntax error.

| **Preview** | **History** | **Starred** | **Help** |

| row | value | row.index % 3 == 0 |
|-----|-------|--------------------|
| 1. | NO STANDING | true |
| 2. | NO STANDING | false |
| 3. | PHTO SCHOOL ZN SPEED VIOLATION | false |
| 4. | REG. STICKER-EXPIRED/MISSING | true |
| 5. | NO STANDING-DAY/TIME LIMITS | false |
| 6. | FIRE HYDRANT | false |

| OK | Cancel |

## row.index % 3 == 0

# ETL Process Using NYC OpenData, OpenRefine and Orange Data Mining

## Extract Process

1. Go to NYC OpenData website > click on Actions > Query data



2. Under Filters, click on "Select a column to filter"

3. Select the Field > Boolean > search terms > Click "Apply"

**Apply**

Filters | ⊗ Clear all

| Tᴛ Issue Date ▾ | ⋮ | is between ⌄ | ⚠ 7/1/2023 | ⋮ | AND 7/4/2023 |

4. Remove columns that are not needed or have large amounts of data > find column > click on 3 horizontal lines > click on "Exclude column from the query":

☰  🔗 Su...
           su...
                    ☰

▼ Filter

↑ Sort ascending

↓ Sort descending

⌐ Column order

🗒 Group and aggregate

⊖ Exclude column from the query

Description

None provided

5. Click on "Export" button:

Export

6. Click on "Download" to retrieve CSV file:

## Export dataset ✕

Only the data returned by your current query will be exported.

Download file | API endpoint

Export format
CSV

Cancel | Download

# Transform Process

7. Import CSV file into OpenRefine > Browse to file > Next

# 8. Initial upload of data > click on "Create project" button:

## 9. Creation of Project:



## 10. Remove columns that are not needed > click on down arrow of column > Edit column > Remove this column:

11.  Generate text facets to understand the different types of data that are in the column > click on arrow of column data that you are interested in > Facet > Text facet > analysis box on left side:

**12. Click on "Cluster" button to begin clustering process:**



Clustering refers to finding similar data entries that could be combined under one topic.

13. Select "Method" and "Keying function" [in this case, Key Collision and Metaphone3:



7. Click on "Cluster" to begin process

14. Find values that can be consolidated, merge and type in replacement value in "New Cell Value":



15. Click on one of the "Merge Selected" to begin process:

16. Convert text to number: Facet > Edit cells > Common transforms > To Number:

17.  Text is converted into numbers [turns green]:

| Fine | Pena |
|------|------|
| 95 | 0.00 |
| 95 | 0.00 |
| 50 | 0.00 |
| 65 | 0.00 |
| 115 | 0.00 |
| 115 | 0.00 |
| 115 | 0.00 |
| 115 | 0.00 |
| 65 | 0.00 |
| 65 | 0.00 |

18. Extract sample size, eg: 1 out of every 3 rows.  Choose any columne > Facet > Custom text facet:



19. In the "Expression" box, type:

row.index % 3 == 0

20. This will create a new column that will set every 3rd row to "True":

**Custom facet on column Violation**

Expression                    Language [ General Refine Expression Language (G  L) ⌄ ]

`row.index % 3 == 0`                                              No syntax error.

| Preview | History | Starred | Help |

| row | value | r  index % 3 == 0 |
|-----|-------|-------------------|
| 1. | NO STANDING | true |
| 2. | NO STANDING | false |
| 3. | PHTO SCHOOL ZN SPEED VIOLATION | false |
| 4. | REG. STICKER-EXPIRED/MISSING | true |
| 5. | NO STANDING-DAY/TIME LIMITS | false |
| 6. | FIRE HYDRANT | false |

[ OK ]  [ Cancel ]

21. Click "OK" on lower right hand corner:

22. A text facet with "True" and "False" will be created:



23. Click on "True" to extract every 3rd value:

24. Extract sample size: Export > Comma-separated value:

Open... | Export ▾

OpenRefine project archive to file

Tab-separated value

Comma-separated value

HTML table

Excel (.xls)

Excel 2007+ (.xlsx)

ODF spreadsheet

Custom tabular...

SQL...

Templating...

Wikibase edits...

QuickStatements file

Wikibase schema

25. CSV file will be 1/3 the size of the original file:

| | | | |
|---|---|---|---|
| 300-000-July16-July31-2023Open-Parkin... | 4/25/2025 9:44 PM | Excel.CSV | 10,259 KB |

Overall, was able to reduce 70 million rows covering 7 months to 439,000 rows.



## **Load Process**

Load file into Orange Data Miner:

Open Orange Data Mining and drag the CSV File Import widget into the work area:

Double click, find file and load file into Orange Data Miner:

# Main File with 439,000 rows with basic analysis:

# More detailed analysis workflow by breaking out individual topics by CSV file:



County CSV file: →

Correspondence Analysis

Data

Data — Distributions

Data

County

Pivot Table

Time CSV file: →

Time — Data — Data Table (2) — Selected Data → Data — Distributions (1)

Precincts CSV file: →

Precincts — Data — Group by — Data — Data Table — Selected Data → Data — Bar Plot

Selected Data → Data — Scatter Plot (1)

Violations CSV file: →

Violations — Data — Group by (1) — Data — Data Table (1) — Selected Data → Data — Bar Plot (1)

Selected Data → Data — Scatter Plot (2)

# County analysis:







Note that Manhattan accounts for the highest number of tickets.



## County

| Count | Bronx | Brooklyn | Manhattan | No Borough | Queens | Staten Island | Total |
|---|---|---|---|---|---|---|---|
| **Bronx** | 116893.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | **116893.0** |
| **Brooklyn** | 0.0 | 240218.0 | 0.0 | 0.0 | 0.0 | 0.0 | **240218.0** |
| **Manhattan** | 0.0 | 0.0 | 253287.0 | 0.0 | 0.0 | 0.0 | **253287.0** |
| **No Borough** | 0.0 | 0.0 | 0.0 | 67844.0 | 0.0 | 0.0 | **67844.0** |
| **Queens** | 0.0 | 0.0 | 0.0 | 0.0 | 236452.0 | 0.0 | **236452.0** |
| **Staten Island** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 30794.0 | **30794.0** |
| **Total** | **116893.0** | **240218.0** | **253287.0** | **67844.0** | **236452.0** | **30794.0** | **945488.0** |

# Time analysis by the hour:



The time of day with the most tickets is from 10 AM to 12 PM

| | Hour Number |
|---|---|
| 1 | 0 |
| 2 | 16 |
| 3 | 15 |
| 4 | 12 |
| 5 | 14 |
| 6 | 11 |
| 7 | 13 |
| 8 | 15 |
| 9 | 9 |
| 10 | 10 |
| 11 | 7 |
| 12 | 18 |
| 13 | 8 |
| 14 | 11 |
| 15 | 10 |
| 16 | 10 |
| 17 | 9 |
| | 11 |

# Precinct analysis:

Close up of Precinct scatter plot.

Top Precincts:

a. 19th
b. 14th
c. 18th
d. 13th
e. 1st

The 19th Precinct is on the Upper East Side has the highest number of tickets, followed by the 14th [Midtown South] and the 18th [Midtown North]. These areas have some of the highest traffic volumes in NYC.



19th Precinct:

# Violations analysis:



**Selected Top 5 in Data Table**

Top 5 Violations scatter plot :

Top violation is Speeding in a School Zone taken by an automatic speed camera.

# Unsupervised learning refers to algorithms that discover patterns and relationships

# Unsupervised learning refers to algorithms that discover patterns and relationships



**t-SNE:** [t-Stochastic Neighbor Embedding] reduces a large number of dimensions to 2D or 3D visualizations.  Uncovers clusters and patterns in data.

**MDS:** [Mulitdimensional Scaling] also maps high dimensional data to lower dimensions but is less intensive with regards to computer processing.

## t-SNE: [t-Stochastic Neighbor Embedding]



Legend:
- Bronx
- Brooklyn
- Manhattan
- No Borough
- Queens
- Staten Island

## MDS: [Mulitdimensional Scaling]



Legend:
- Bronx
- Brooklyn
- Manhattan
- No Borough
- Queens
- Staten Island

Same dataset but some difference in clustering patterns, where t-SNE is preserving local relationships, while MDS maintains global structures and overall data geometry.

# 2 More Unsupervised Algorithims



**Self-Organizing Map:** neural network based; proximity indicates similarity

**Louvain Clustering:** finds disparate clusters using community detection

# Louvain Clustering Feeding Into Scatter Plot

Brighter color indicates higher Fine Amount Along With Heat Map
Manhattan has the highest concentration of the highest ticket fines.

# Conclusion

<u>Tools Used:</u>
 - <u>Extraction:</u> Query tool on NYC OpenData for Parking Tickets
 - <u>Transform:</u> OpenRefine
 - <u>Load/Visualizations:</u> Orange Data Mining - Additional transformation work done with GREL [General Refine Expression Language]

<u>Project Review:</u>
Having worked previously with data from the GAIA Space Observatory where the data was precise and standardized, it was an intereseting learning process to work with NYC OpenData. By contrast, NYC OpenData was much less precise and chaotic due to humans recording and performing data entry with little quality control.  However, this challenge was an opportunity to learn how to clean and modify the data using OpenRefine.  The Clustering Function in OpenRefine was quite fascinating, as were the results produced by the different algorithimic clustering functions.  Likewise, using Orange Data Mining was a wonderful chance to learn about the different visualizations that data was capable of being presented.  In addition, this application introduced me to the large number of Unsupervised Algorithims that are available to further reveal hidden patterns and trends within large amounts of data.

Both programs required a fairly steep learning curve in order to understand their basic functions.  In doing so, I was able to get a taste of the true power of these programs.

In retrospect, I would have allocated more time towards the download process from NYC OpenData, as I did not anticipate that there would be upwards of 10 million rows of data per month. However, this gave me the chance to learn the true power and usefulness of OpenRefine in handling massive levels of data.

Benefits of using OpenRefine and Orange Data Mining:
 - Do not have to create customized Python code from scratch
 - Both tools are refined and well-tested applications
 - Data does not have to be uploaded to the cloud, saving time and money
 - Since all data is processed locally, the user's privacy is enhanced
 - After using both applications, the skills gained will allow users to take those skills into the field to work on data quickly without having to pay for use of the cloud

After working with both tools in the course of this project, it seems that Baruch could develop a whole course on how to use all aspects of these substantial programs. In doing so, this would allow students to perform sophisticated data manipulations and robust modifications without learning the intracacies of Python programming.

References:
 - <u>Extraction:</u> Query tool on NYC OpenData for Parking Tickets
Website: https://data.cityofnewyork.us/City-Government/Open-Parking-and-Camera-Violations/nc67-uf89/about_data
 - <u>Transform:</u> OpenRefine: https://openrefine.org/
 - <u>Load/Visualizations:</u> Orange Data Mining: https://orangedatamining.com/
 - <u>Additional transformation work:</u> GREL [General Refine Expression Language]: https://openrefine.org/docs/manual/grel