

A Closer Look into NY Parking Violations

Cody Hegwer

Applied CS
CU Boulder
Boulder, CO
cohe7798@colorado.edu

Jeff LaPrade

Applied CS
CU Boulder
Boulder, CO
jela2733@colorado.edu

Adam Grabowski

Applied CS
CU Boulder
Boulder, CO
adgr5757@colorado.edu

PROBLEM STATEMENT

New York is notorious for being a difficult place to find parking in. There are millions of people living in New York City, and many streets are narrow and oftentimes one way. By gathering street parking violations from the 1970s until 2016, the goal is to find trends between common violations in order to avoid them. By using the Vehicle Body Type, Vehicle Color, Violation Time, Street, Intersecting Street, and Plate type attributes specifically, our data will guide drivers in New York to be more aware of which types of cars attract attention when marked for parking violations along with certain streets that are surveyed more often than others.

LITERATURE SURVEY

In 2014, I Quant NY published a study analyzing fire hydrants responsible for large amounts of parking violations within New York City using NYC Open Data. The study found that the most ticketed fire hydrant areas had a bike or protected lane that seemed to leave the hydrant open, but oftentimes violations were left to interpretation. Violations from a single hydrant led to \$25,000 in fines through one year.

Another study found to be interesting to our project was a study published by Angela Ju on

LinkedIn titled "Toronto Parking Ticket Data Analysis." While this study obviously does not analyze parking in New York, its conclusions helped lead to the motivation in our work. The top 5 violation types in Toronto led to 50% of the revenue from all parking violations. This points to confusion by drivers in how certain streets and parking zones lead to many repeat fines. On top of that, it may point to certain city workers who frequent confusing parking zones knowing that people with mistakenly park there without thinking twice.

By analyzing these two studies, we decided parking violations may not always originate from malintent. Oftentimes using a parking space or not that seems open is a split second decision. So, by finding and marking car types that are often flagged, along with streets and street intersections we hope to help parkers in New York to make better decisions and know the risks of parking in certain areas.

PROPOSED WORK

To begin cleaning the data, there are many attributes which are encoded by a number that need to be translated to usable information. Just to list a few, 'Violation Code' is a 2-digit number referring to why the ticket was issued. There are

A Closer Look into NY Parking Violations

3 'Street Code' attributes which are 5-digit numbers, we believe to refer to the location of the violation. There are 'Violation Location' and 'Violation Precinct' which are 3-digit numbers that, from a brief look, tend to be the same.

Another problem is missing values for given attributes. The last 5 columns describe specific violations, such as 'Hydrant Violation' and 'Double Parking Violation'. If we wish to find any interesting information relating to these violations we will have to remove all null rows as the occurrence of said violations are rare, making classification for these pointless. Instead, we would like to classify *given* that the recipient illegally parked in front of a fire hydrant, or *given* that the recipient double parked.

Once the data has been cleaned, we plan to create various predictive models. We will experiment using a subset of our dataset as training data with different attributes as predictors and different classification methods to make the best model. At this stage, it is unclear which attributes will make the best predictors, and examining the accuracy of different models may actually be an answer to some of our initial questions.

Once we believe we have the best model, we may fabricate new data or use another subset of our original dataset to test its predictions. Again, exactly how what attributes will be used in building the model is unclear at this stage: our observations and their relevance to our different inquiries will determine this.

Our inquiries will be very different than of the related work we have identified. We may use classification or predictive modelling in a similar fashion, but ultimately we will be looking answer different questions.

DATA SET

<https://www.kaggle.com/new-york-city/ny-parking-violations-issued>

Attributes:

Summons Number: Ordinal Integer
Plate ID: Ordinal String
Registration State: Ordinal String
Plate Type: Ordinal String
Issue Date: Ordinal Integer
Violation Code: Ordinal Integer
Vehicle Body Type: Ordinal String
Vehicle Make: Ordinal String
Issuing Agency: Ordinal String
Street Code1: Ordinal Integer
Street Code2: Ordinal Integer
Street Code3: Ordinal Integer
Vehicle Expiration Date: Ordinal Integer
Violation Location: Ordinal String
Violation Precinct: Ordinal Integer
Issuer Precinct: Ordinal Integer
Issuer Code: Ordinal Integer
Issuer Command: Ordinal String
Issuer Squad: Ordinal String
Violation Time: Ordinal String
Time First Observed: Ordinal String
Violation County: Ordinal String
Violation In Front Of Or Opposite: Ordinal String
Number: Ordinal String
House Number: Ordinal String
Street: Ordinal String
Street Name: Ordinal String
Intersecting Street: Ordinal String
Date First Observed: Ordinal Integer
Law Section: Ordinal Integer
Sub Division: Ordinal String
Days Parking In Effect: Ordinal String
From Hours In Effect: Ordinal String
To Hours In Effect: Ordinal String
Vehicle Color: Ordinal String

Unregistered Vehicle?: Boolean String

Vehicle Year: Ordinal Integer

Feet From Curb: Quantitative Integer

EVALUATION METHODS

Predominantly, we will be analyzing the accuracy and precision of the predictive model to answer our questions. Due to the nature of these inquiries, the main events we wish to predict are future events that, for the sake of this study, we would not have access to. This means that the quality of our model is the main concern of our evaluation. All of the knowledge we wish to gain is based on the prediction, and so its effectiveness very well may answer our questions as well as being our reference point. The precision of our classification will be more readily evaluated since we may predominantly be testing on data already in our dataset as opposed to new, more recent data.

TOOLS

There are two main platforms we wish to explore our data with. Tableau for visualization and Python for raw data calculations, statistics, and classification. Specifically, we'd like to utilize Numpy, Pandas, and Matplotlib libraries in Python.

1 Tableau

Tableau is a tool for rapid creation for data visualization, allowing the user to input any data, typically a .csv, and generates graphs based on user selected attributes.

A major benefit to using Tableau is its ability to quickly and easily map locations globally, identify borders, and fill in these areas with desired attributes. For example, given US counties voting

results for the 2016 presidential elections, Tableau automatically identifies all the counties locations on a map and, with user instruction, fills in the counties with red and blue based on the counties voting preferences.

This seems interesting to us, as we would like to see if we could discover some trends if the data was mapped to locations, ideally street addresses, in New York City. There are two glaring limitations to this though. One, is if we find that the street addresses are too general and only reduce to the actual street name or we find that majority of tickets issued provided no information for this attribute. Second, we are unfamiliar with if Tableau's location mapping can actually handle street names.

2 Python

The programming language we decided to use for our actual data mining is Python. It is the language we are most familiar with in our college career. It is simple and a good tool for data manipulation. Python also has many well flushed out libraries intended for use by data scientists. The following will be used extensively.

2.1 Numpy/Pandas

Numpy and Pandas are the main libraries being used for the project. Both contain many statistical measures and tools. Additionally, Pandas contains data structures to convert .csv files into easily manipulatable tables.

2.2 Matplotlib

The Matplotlib library contains many options for graph and chart creation. It pairs nicely with Numpy and Pandas as the data we are concerned about are converted to usable arrays.

MILESTONES

There are 5 remaining project turn-ins before the end of the semester. 2 of which allude to the completion of the project. This section is dedicated to setting a timeline for ourselves to abide by as well as giving an outline of the work for following weeks.

1 Project Progress Report

The next major assignment is a Project Progress Report. It should be an update to this document, specifically an update on the milestones we've completed. The following are goals we would like to complete for this assignment:

1.1 Confirmation Meeting(s)

Brief meetings to confirm details about the data set, which attributes we wish to focus on, what hypotheses we hope to prove, etc. Should be regular to keep other milestones on track.

1.2 Data Cleaning

1.3 Model Building and Analysis

2 Final Project Report

At this point, the data mining project is assumed to be completed. The following milestones represent outstanding tasks that should come after our previous Progress report:

2.1 Cleanup/Commenting of Code

2.2 Prepare Presentation