# Exploratory Analysis of Parking Violation

Rishi Bamb     Sanjana Balagar     Shrivatson Ramaratnam Giridharan     Sonali Gupta

*Abstract — Big data is the result of technological advancements that have resulted in the advent of massive amounts of data. Big data refers to datasets that are not only large in size but also contain a lot of different types of data. When properly analyzed, these data can assist industries in making important decisions in a variety of ways. Parking violations are a daily problem in today's fast-paced environment. Parking a vehicle illegally may result in an offense, resulting in a large number of traffic citations being issued. 'Parking Violation' data is one such data set for our exploratory investigation. Every day, millions of automobiles are parked in cities, and New York, as a major metropolis, is no exception, with most residents having parking issues.New York City itself collected approx $957 million in fine revenues In  them more than 59%  that is approx $565 million of the $957 million, come from parking tickets. The analytics and visualizations are performed using various AWS Services.*

## I. INTRODUCTION

[1]      In this project we have used the parking violation dataset from the open nyc and kaggle which contained more than 50 million records and performed various analytics to generate meaningful insights.The NYC Department of Finance collects data on every parking ticket issued in NYC and is responsible for collecting and processing payments of all tickets. Because of the huge number of cars and the limited geography, there are a lot of parking tickets. This prompted us to conduct an exploratory analysis on such data in order to gain insights such as when and where tickets are more likely to be issued, if there is a specific season for it, what types of vehicles are receiving tickets, comparing state data to determine which state has the most tickets issued on a monthly basis, and which vehicle body type receives the most tickets. This analysis is carried out utilizing different AWS services like the S3 for data storage, Redshift for performing queries and analytics , and QuickSight to perform dynamic visualizations with the goal of developing a graphical solution for real-time analytics using the parking dataset.

## II. MOTIVATION

The project's inspiration came from academic coursework.We decided to use the parking violation  dataset for our project which would help in analysing the number of violations issued based on multiple factors.

The project will help us in exploring the big data concepts and tools to understand the parking ticketing system which will moreover help us to gain analytical insights from the available data and help the ticketing department and issuing agencies to analyze the regions where vehicles are getting the highest number of summons and the reasons. We have also seen the rise in the violations according to each fiscal year.By doing this project we will analyse the highest number of tickets raised in each state, the number of summons in each month, and the highest number of violation code committed by the citizens.

## III. LITERATURE SURVEY

### A. Automatic Detection of Parking Violation and Capture of License Plate

Author Z. Liu, W. Chen and C. K. Yeo stated about the automatic detection of parking violations.Illegal parking is a common problem in urban areas, causing traffic congestion and posing a safety concern to other road users. Surveillance video is saved for post-event forensics, and detecting offending vehicles often necessitates manual review. Foreground and background segmentation, which is less resistant to environmental influences, as well as the more robust Single Shot MultiBox Detector (SSD) and its derivatives, are examples of automated detection systems. This study presents a fully automated pipeline for detecting illegal parking from start to finish. You Only Look Once Version 3 is a deep learning-based object detection system that provides reliable and rapid car detection (YOLOv3). The stationary time of the vehicle-in-violation is tracked using template matching and Intersection over Union (IoU) calculations, with built-in error tolerance measures. The license plate is extracted using OpenALPR. Under different conditions, empirical results reveal a high level of vehicle detection and movement tracking accuracy.

### B. Applied research on Real estate price prediction by neural network

Authors N. Lin, E. Liu, F. Tenorio, X. Yang and D. Woodbridge wrote this work in which they highlighted using networked systems and unsupervised machine learning methods on a huge dataset, they investigated the similarities in

parking ticket records in this study. They used an algorithm to cluster existing tickets and dug deeper to identify the distribution of precincts within various clusters using 37 million ticket data (9 GB) gathered by the New York City Department of Finance. This project made use of Amazon Web Services, such as S3, EC2, and EMR, as well as technologies like MongoDB and Apache Spark. The computational time and cost for various EMR settings were studied in this study. They found that when implementing unsupervised learning on a big dataset, as well as storing and managing data, distributed systems offer significant computing benefits. They also discovered that a cluster with more workers is faster than a cluster with fewer workers but more memory space for the data set in use. outcomes based on human behavior, decision-making, and a variety of other elements. Based on the pricing and controlling components, we can analyze and seek to evaluate the property categories that are a great match to the users in our project.

## IV. Proposed work

In this project, we have obtained the dataset from Kaggle and NYC Open data platform. These data are preprocessed using Talend, an etl tool and passed into AWS services like S3 for storage and Redshift for effective analysis of it. And finally, the data is visualized using AWS QuickSight and matplotlib library in python, the resultant showcase the exploratory analysis of the parking violations.
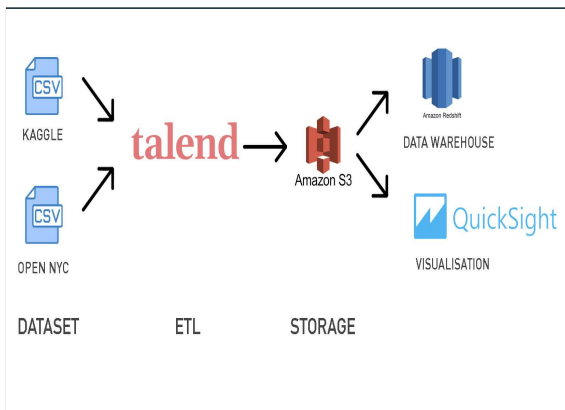
Figure.1

### 1. Data Cleaning - Talend

As we were having four fiscal years data from 2013 to 2017, we cleaned and normalized individual years files using Talend which provides feasible services to process and prepare data. Below diagram is a brief look of one of the normalized files.

Figure.2

### 2. Merging of files

To have data in a single file makes an analysis and working on it easy hence we merged all four files into one using python code.

```python
import os
import glob
import pandas as pd
#set working directory
os.chdir("C:/Users/sonal/Downloads/test")

#find all csv files in the folder
#use glob pattern matching -> extension = 'csv'
#save result in list -> all_filenames
extension = 'csv'
all_filenames = [i for i in glob.glob('*.{}'.format(extension))]

#combine all files in the list
combined_csv = pd.concat([pd.read_csv(f) for f in all_filenames ])
#export to csv
combined_csv.to_csv( "combined_csv.csv", index=False, encoding='utf-8-sig')
```

Figure 3

### 3. Loading transformed data into AWS S3 bucket

Figure. 4

### 3. AWS Redshift

Created cluster on redshift service by providing access to user through an associated IAM role to perform analysis using its query editor.
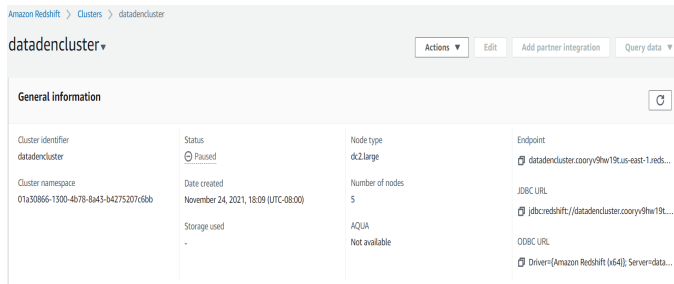
Figure. 5

## 4. Analysis - Redshift Query editor

First we have created a ticket violation table including all required attributes with metadata in a database. Then, loaded data into it from a stored csv file in s3 bucket to perform queries and analysis. The following analysis helps in giving the insights on multiple factors to explore trend and number of summons being issued in four consecutive years.
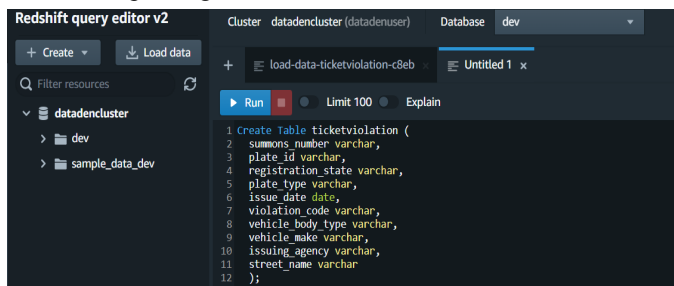
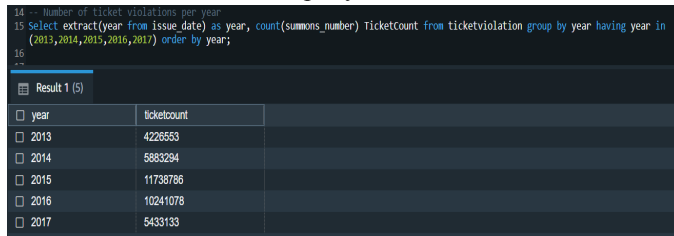Creation of parking violations table:



Figure. 6

Number of ticket violations per year:



Figure. 7

Count of violation in NY:



Figure.8

Max number of violation codes in each registration state or violations most occurred in country:



Figure.9

## V. Data Visualisation

In our project we will be showcasing visualizations on AWS Quicksight and python. Quicksight's analytic platform empowers any skill level target audience to work with data through actionable and insightful visualizations. Below are the screenshots of visual analysis done on Quicksight.We have done visualisations using quicksight where we loaded a JSON manifest file via S3 bucket into quicksight for performing visualisations.

The first visualization is a pie chart which shows the top agencies which are issuing the tickets.· As we can see, the majority of the tickets are issued by T which is the traffic department and it Is obvious since New York is such a populated city and Traffic jams are normal on a day to day basis.The next agency issuing the maximum violation is Transportation department (V).
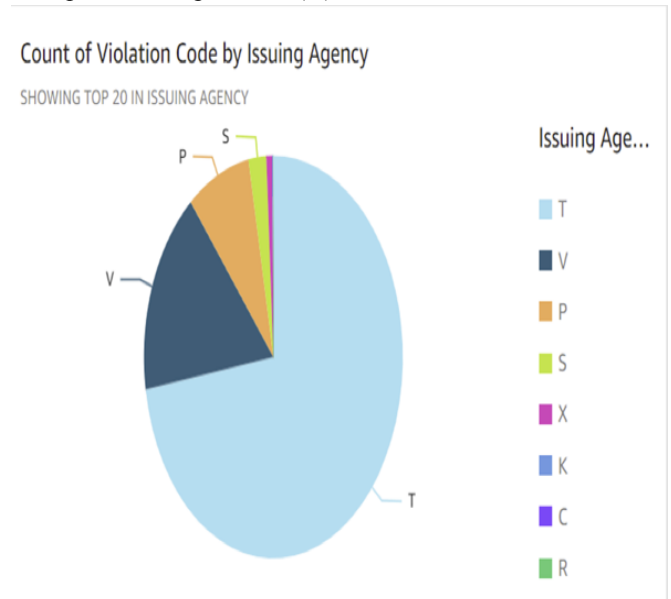


Figure.10

The next bar chart depicts the different violations done by a particular vehicle. From the visualization we can see that ford is the highest and then mg and so on.

Figure.11



Figure.13

The pie chart shows the plate type which has the maximum violations and from this we can see that plate type PAS has the most number of tickets.

The pie chart shows that manufacturer ford has got the maximum number of summons with Toyota being the second.
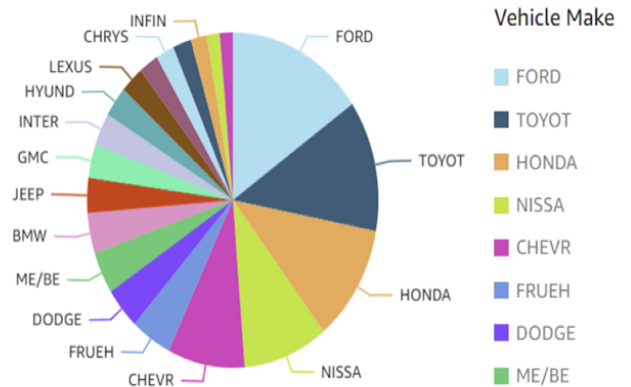


Figure.12



Figure.14

The word chart shows that New York has got the most number of violation code per registration state by looking at the size of NY when compared with other word sizes.

The next geospatial chart shows the count of summons numbers according to the registration state. Here too we can infer that out of North America most of the violations are done in New York and New Jersey being the second.
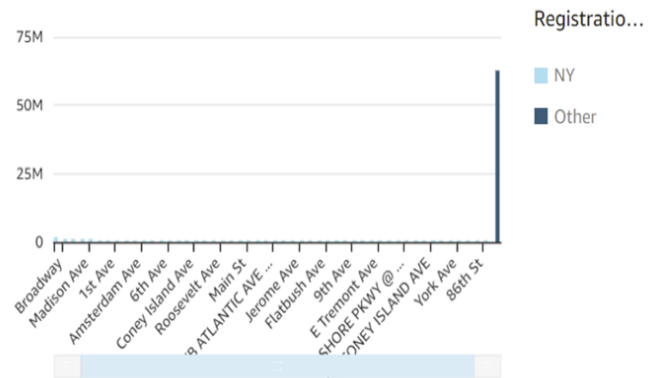
Count of Violation Code by Registration State



Figure.15

The next donor chart shows the number of summons that have been issued according to the vehicle type or the body type . Here we can infer that suburbans contribute to a major part about 35% . We analysed that since suburbans are large in size and they are difficult to park in a street setting it has the highest number of summons issued.

Count of Summons Number by Vehicle Body Type
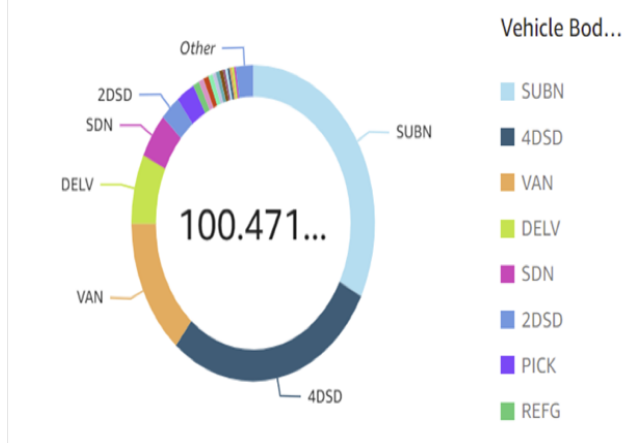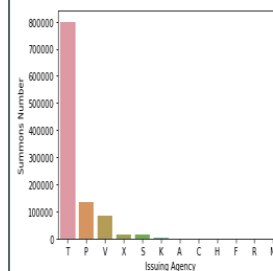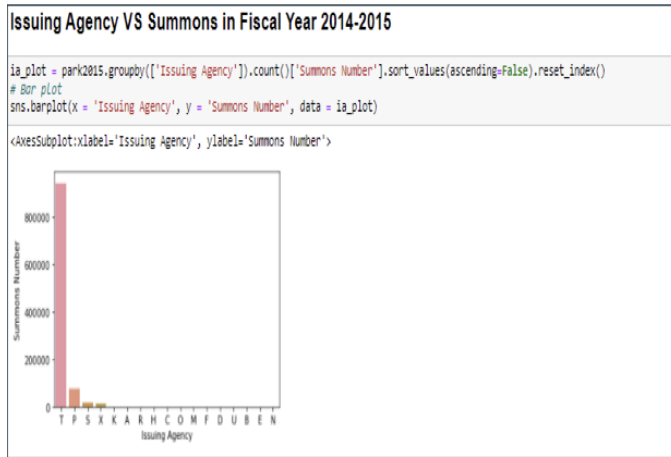
SHOWING TOP 20 IN VEHICLE BODY TYPE



Figure.16

The next bar chart shows the street which has the maximum number of summons by street name . Thus broadway street in New York city has the maximum number of violations due to the fact that it has large number of theatres , pubs , restaurants etc which aid to the violations.

Count of Summons Number by Street Name and Registration State

SHOWING TOP 50 IN STREET NAME AND TOP 2 IN REGISTRATION STATE



Figure.17

Understanding of data is easy when we place it in a visual context. Python offers great libraries and features to represent data graphically. In our project we are generating visualizations using matplotlib. Below are the screenshots of visualizations generated in python.

Issuing Agency VS Summons in Fiscal Year 2013-2014

```
ia_plot = park2013_14.groupby(['Issuing Agency']).count()['Summons Number'].sort_values(ascending=False).reset_index()
# Bar plot
sns.barplot(x = 'Issuing Agency', y = 'Summons Number', data = ia_plot)

<AxesSubplot:xlabel='Issuing Agency', ylabel='Summons Number'>
```



Figure.18

The above bar chart obtained from python shows the number of summons issued by each issuing agency from 2013 - 2015.
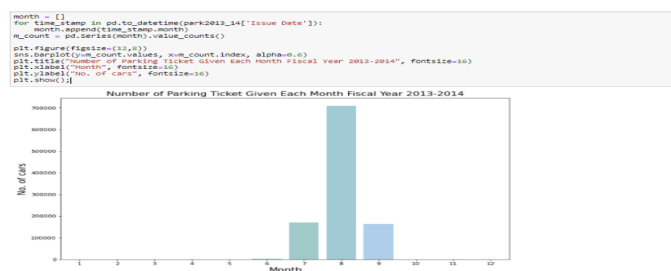


Figure.19

The above bar chart shows the number of tickets raised in each month in the year 2013 - 2014 here we can see that the peak number of tickets starts from the month of june and goes all the way till september.
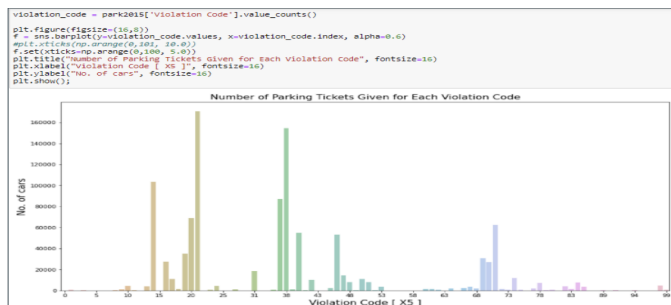


Figure.20

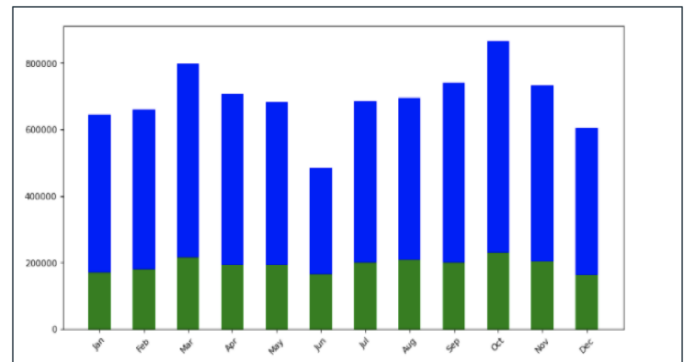The above chart shows the number of tickets raised for each violation code.



Figure.21

The above stacked bar chart shows the number of tickets raised in New York versus the other states in each month for the year 2015 - 2016. Here Blue represents New York and Green represents other states.

## IV. CONCLUSION

We were able to study this large volume of parking dataset released between 2013 and 2017 which consisted of more than 50M records. We performed analytics using quicksight and the visualizations shows that New York gets the highest number of tickets among all the states in United States and in new york on the broadwalk street there are highest number of summons. From our analysis it can also be concluded that in the month of august there are a high number of violations issued as in summers there are many events. It can also be observed that in each month more tickets are raised in New York than in other states.

REFERENCES

[1] Z. Liu, W. Chen and C. K. Yeo, "Automatic Detection of Parking Violation and Capture of License Plate," 2019 IEEE 10th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), 2019, pp. 0495-0500, doi: 10.1109/IEMCON.2019.8936164.

[2] Lin, N., Liu, E., Tenorio, F., Yang, X. and Woodbridge, D., 2019, August. Distributed Data Analytics Framework for Cluster Analysis of Parking Violation. In 2019 IEEE SmartWorld, Ubiquitous Intelligence & Computing, Advanced & Trusted Computing, Scalable Computing & Communications, Cloud & Big Data Computing, Internet of People and SmartCityInnovation(SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP /SCI) (pp. 1958-1963). IEEE.