

A
Project Report
On
“NYC Automated Traffic Volume Counts”

Prepared by

Saurabh Jain | saurabh_jain@csu.fullerton.edu

Pratishtha Soni | pratishthasoni@csu.fullerton.edu

A Report Submitted to
Tseng-Ching James Shen, PhD
Department of Computer Science
College of Engineering and Computer Science (ECS)
California State University, Fullerton (CSUF)
for the Fulfillment of the Requirements for the
CPSC 531-01 Advance Database Management



**California State University, Fullerton
CA - 92831
December, 2022**

TABLE OF CONTENTS

INTRODUCTION.....	3
FUNCTIONALITIES.....	7
ARCHITECTURE & DESGIN.....	8
GITHUB LOCATION OF CODE.....	10
DEPLOYMENT INSTRUCTIONS	11
STEPS TO RUN THE APPLICATION.....	12
TEST RESULT	15

INTRODUCTION

Problem Statement

The vehicular traffic is increasing tremendously these days, simultaneously congestion also increases. To prevent congestion, one option is to increase the capacity by increasing the number of the existing transportation system. A second option is to develop alternatives that increase capacity by improving the efficiency of the existing transportation system. The latter focuses on building fewer lane miles.

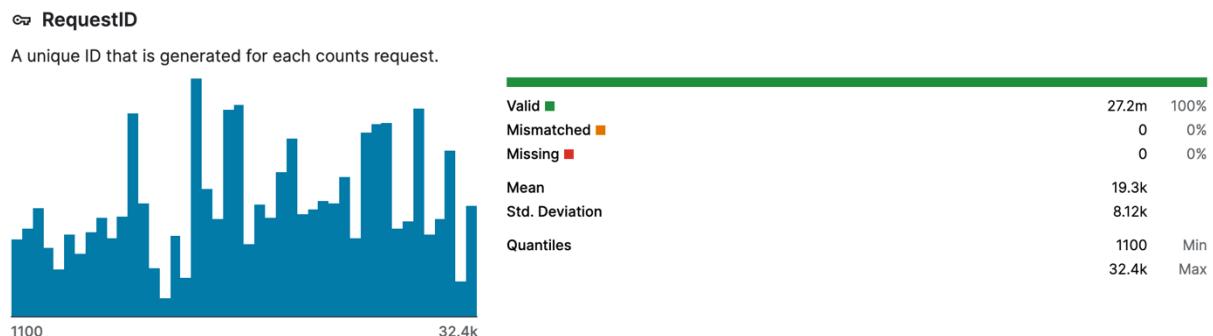
To Overcome these Problems, the New York City Department of Transportation (NYC DOT) uses Automated Traffic Recorders (ATR) to collect traffic sample volume counts at bridge crossings and roadways. These counts do not cover the entire year, and the number of days counted per location may vary from year to year.

Data Description

NYC Automated Traffic Count allows the user to view and access traffic data information of different cities in NYC. The information is displayed in the form of excel. New York City transportation uses automated traffic recorders to collect traffic volume counts at bridge crossings and roadways. The data consists of car volume, averaged every 15 mins for the duration of 2000 to 2020. This traffic data can be used to determine various type of analysis and design patterns in traffic over the year 2000 to 2020.

The information is collected from Kaggle in a form of 1 3GB excel file and consists of 14 rows and 27190511 columns:

RequestID - Unique ID that is generated for each count request.



Boro - Lists which of the five administrative divisions of New York City the location is within, written as a word. Where Brooklyn data is 29% and Queens data is 26%.

A Boro

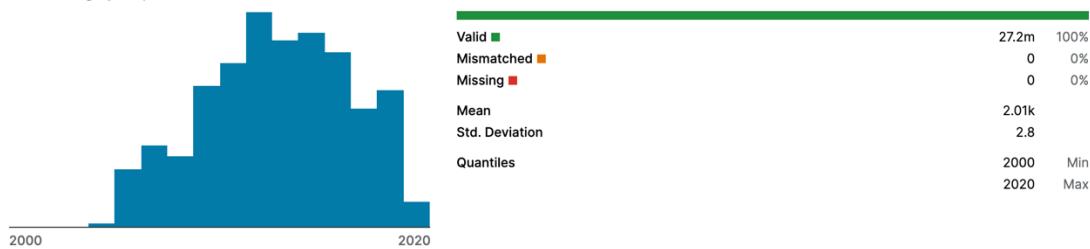
Lists which of the five administrative divisions of New York City the location is within, written as a word



Yr - The two-digit year portion of the date when the count was conducted. (From 2000 to 2020).

Yr

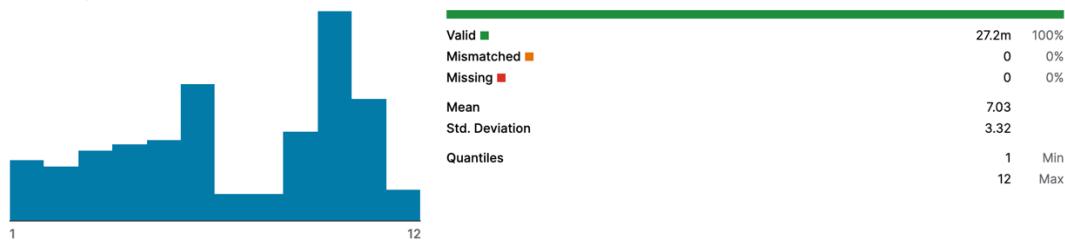
The two digit year portion of the date when the count was conducted.



M - The two-digit month portion of the date when the count was conducted.

M

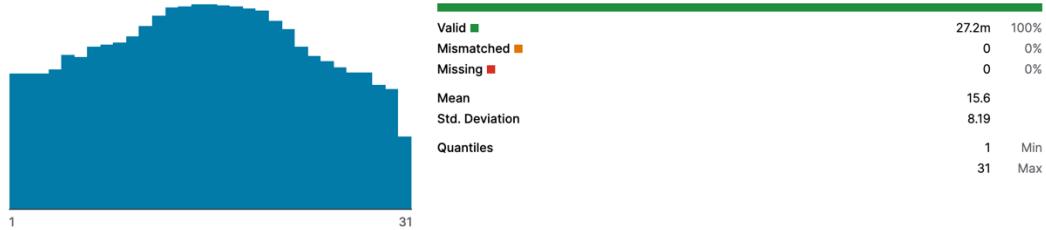
The two digit month portion of the date when the count was conducted.



D - The two-digit day portion of the date when the count was conducted.

D

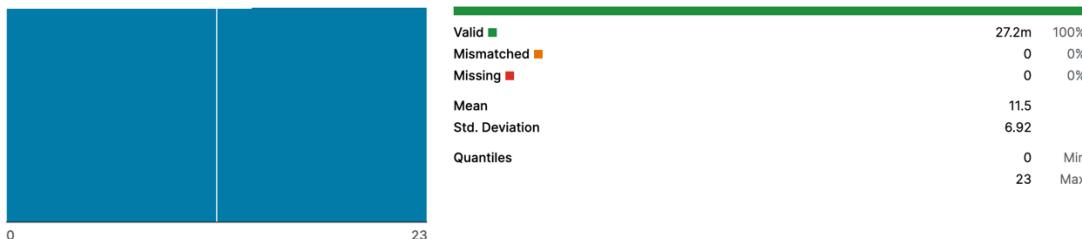
The two digit day portion of the date when the count was conducted.



HH - The two-digit hour portion of the time when the count was conducted.

HH

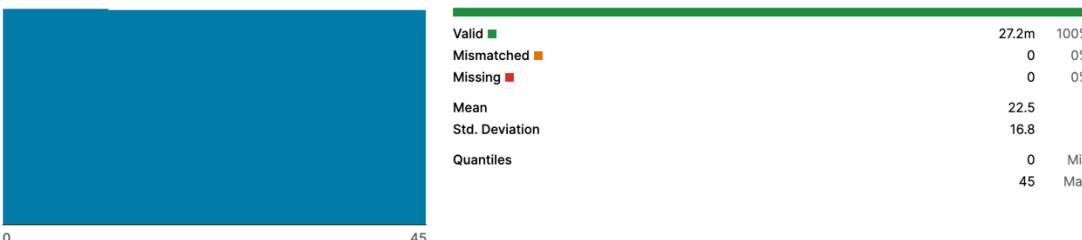
The two digit hour portion of the time when the count was conducted.



MM - The two digits start minute portion of the time when the count was conducted.

MM

The two digit start minute portion of the time when the count was conducted.



Vol - The total sum of counts collected within 15-minute increments.

Vol

The total sum of count collected within a 15 minute increments.



SegmentID - The ID that identifies each segment of a street in the LION street network version 14.

SegmentID

The ID that identifies each segment of a street in the LION street network version 14.



WktGeom - A text markup language for representing vector geometry objects on a map and spatial reference systems of spatial objects.

A WktGeom

A text markup language for representing vector geometry objects on a map and spatial reference systems of spatial objects.



Street - The 'On Street' where the count took place.

fromSt - The 'From Street' where the count took place.

toSt - The 'To Street' where the count took place.

Direction - The text-based direction of traffic where the count took place.

Tool and Technology Used:

- Spark
- Google Collab Notebook
- Google Cloud Platform
- Google Data Proc
- Python Matplotlib

FUNCTIONALITIES

1. Pre-Processing of Dataset:

In this step all the null values from different columns from the dataset is removed using the python language in jupyter notebook. Particularly, pandas library is being used for converting the datatypes (from string to integer and vice versa) and for removing null values and for removing the spaces before, after and in between of each string/int variable.

2. Processing of the Dataset on the Cloud:

The project data is in the form of excel of size **3 GB** and processing this huge data on local is challenging, so processed this bunch of data using Google's cloud dataproc service which is the open-source service for running Apache Hadoop and Apache spark and other tools and frameworks. The dataproc allows a new user 300 free credits to process and play with the data.

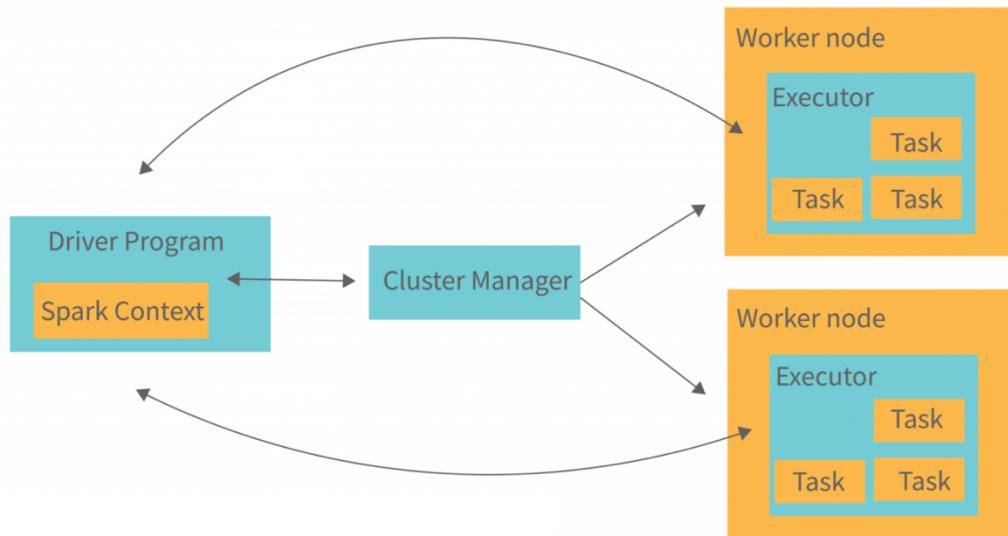
3. Data Analysis/Queries Significance:

- Enhance public safety.
- Reduce congestion.
- Improved access to travel and transit information.
- Generate cost savings to motor carriers, transit operators, toll authorities, and government agencies; and
- Reduce detrimental environmental impacts.

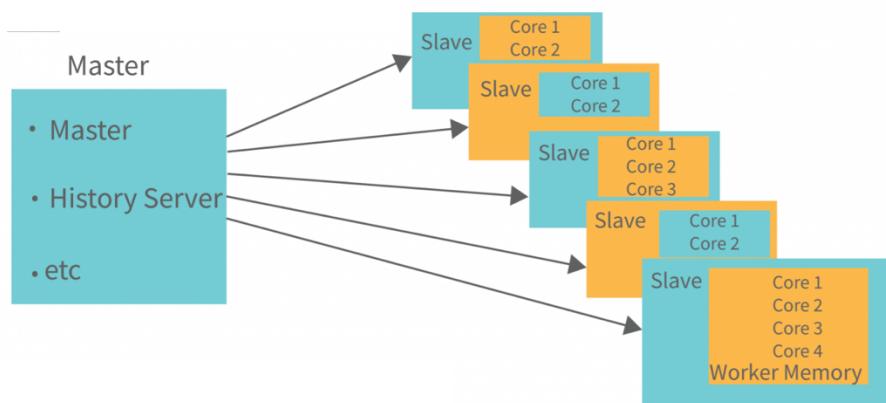
ARCHITECTURE & DESIGN

A high-level view of the Apache Spark application:

In the below diagram, spark uses the master and slave architecture. The driver program connects to a spark cluster, calls the main application, and builds a spark context (which serves as a gateway) to monitor the jobs running in the cluster. The spark context is implemented for all the operations.



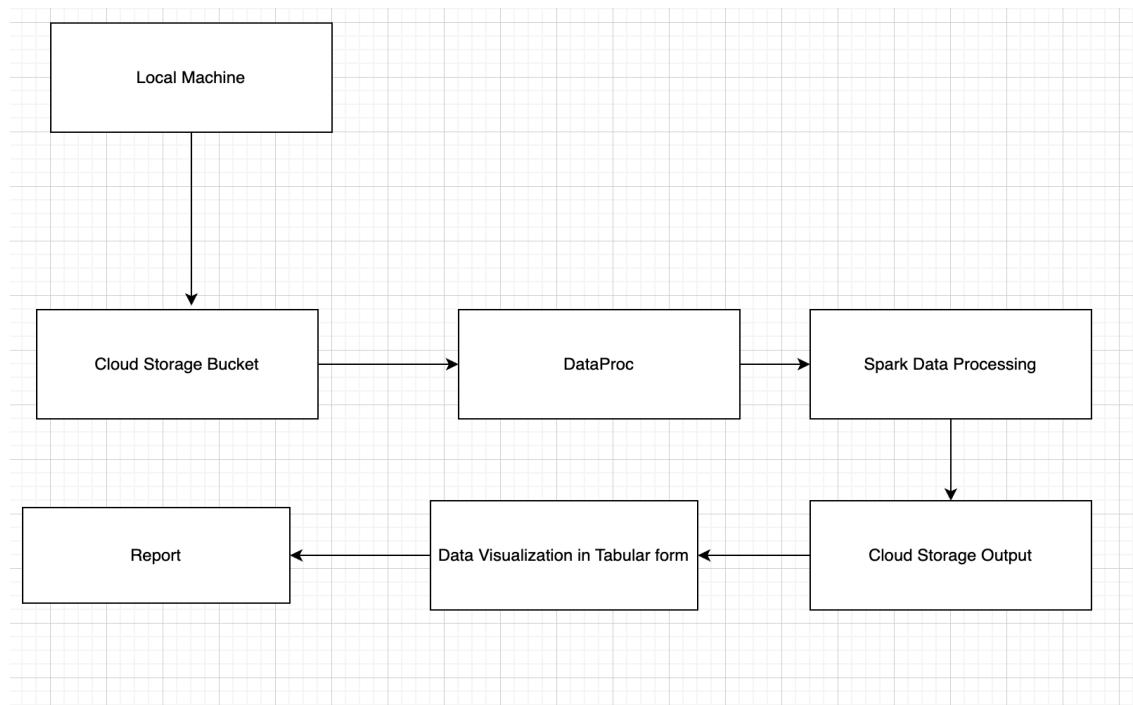
Spark Standalone Architecture



The architecture/implementation diagram of the project:

Methodology:

- **Local Machine:** The computer that you are currently logged on to as a user.
- **Cloud Storage Bucket:** Google Cloud Platform Console to upload files or folders.
- **DataProc:** Dataproc is a managed Spark and Hadoop service that lets you take advantage of open-source data tools for batch processing, querying, streaming, and machine learning.
- **Spark Data Processing:** Apache Spark is an open-source, distributed processing system used for big data workloads. It utilizes in-memory caching and optimized query execution for fast queries against data of any size. Simply put, Spark is a fast and general engine for large-scale data processing.
- **Cloud Storage Output:** Google Cloud Platform Console to upload files or folders.
- **Data Visualization:** Visualization of data using Python Library.

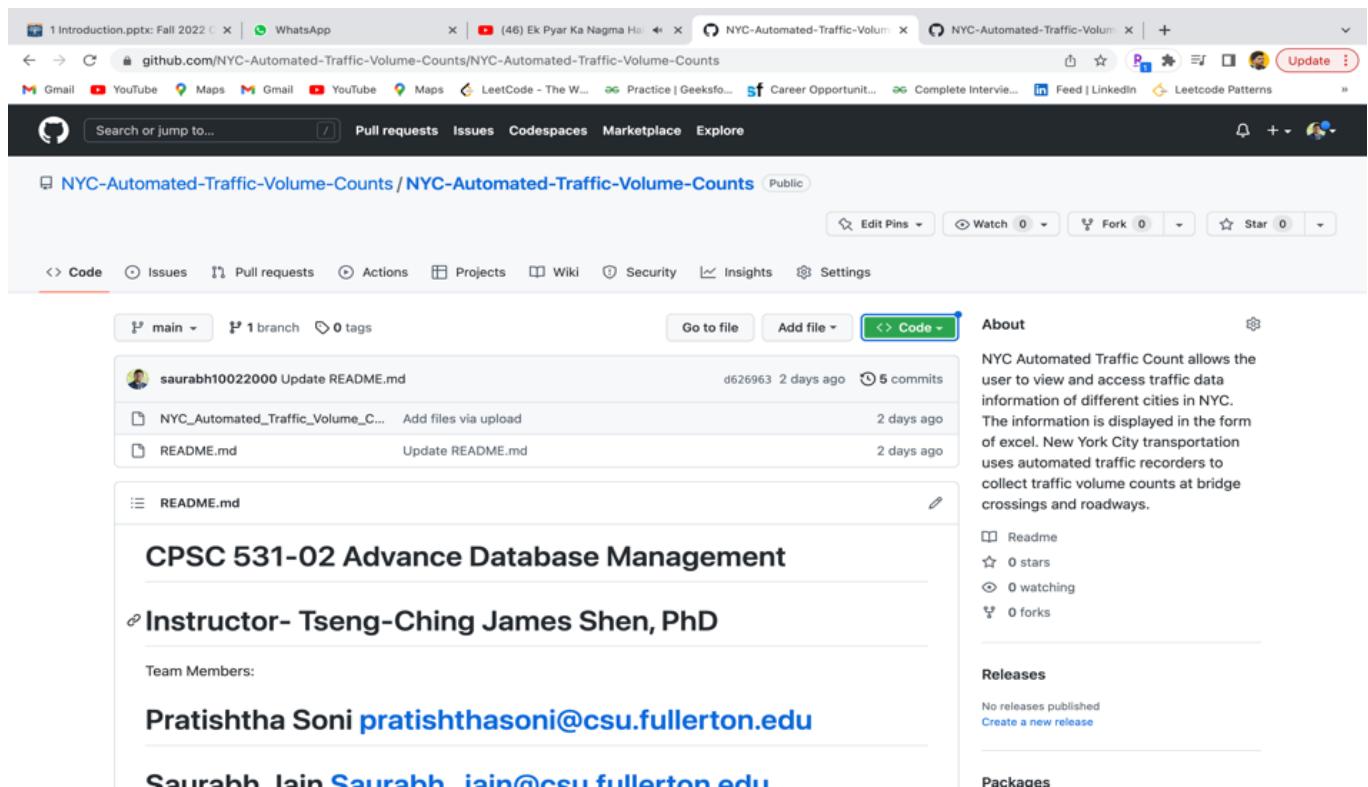


GitHub LOCATION OF CODE

Project Location: <https://github.com/NYC-Automated-Traffic-Volume-Counts/NYC-Automated-Traffic-Volume-Counts.git>.

Steps:

1. Download “NYC_Automated_Traffic_Volume_Counts.ipynb”.
2. Download Data from Kaggle (3GB)
<https://www.kaggle.com/datasets/aadimator/nyc-automated-traffic-volume-counts>
3. Import Data and project files to DataProc or on a Local Machine.
4. Run



The screenshot shows a GitHub repository page for 'NYC-Automated-Traffic-Volume-Counts'. The repository is public and has 5 commits in the main branch. The README.md file contains course information for CPSC 531-02 Advance Database Management, taught by Instructor- Tseng-Ching James Shen, PhD. Team Members listed are Pratishtha Soni (pratishthasoni@csu.fullerton.edu) and Saurabh Jain (Saurabh_jain@csu.fullerton.edu). The repository has 0 stars, 0 forks, and no releases published.

DEPLOYMENT INSTRUCTIONS

Steps to Set-up Spark in Local Machine:

Install JDK

```
!apt-get install openjdk-8-jdk-headless -qq > /dev/null
```

Get Spark installer (Check the path on Spark.org)

```
!wget -q https://archive.apache.org/dist/spark/spark-2.4.3/spark-2.4.3-bin-hadoop2.6.tgz
```

Check if the file is copied

```
!ls
```

Untar the Spark installer

```
!tar -xvf spark-2.4.3-bin-hadoop2.6.tgz
```

Check the spark folder after untar

```
!ls
```

Install findspark - a Python library to find Spark

```
!pip install -q findspark
```

Set environment variables

Set Java and Spark Home based on location where they are stored

```
import os  
os.environ["JAVA_HOME"] = "/usr/lib/jvm/java-8-openjdk-amd64"  
os.environ["SPARK_HOME"] = "/content/spark-2.4.3-bin-hadoop2.6"
```

Create a local Spark Session

```
import findspark  
findspark.init()  
from pyspark.sql import SparkSession  
spark = SparkSession.builder.master("local[*]").getOrCreate()
```

Test Installation

```
df = spark.createDataFrame([{"Google": "Colab", "Spark": "Scala"}, {"Google": "Dataproc", "Spark": "Python"}])  
df.show()
```

STEPS TO RUN THE APPLICATION

1. Create a Bucket and Upload Data.

- In the Google Cloud console, go to the Cloud Storage Buckets page. Go to Buckets.
- Click Create bucket.
- On the Create a bucket page, enter your bucket information. To go to the next step, click Continue. For Name, your bucket, enter a name that meets the bucket name requirements.
- Click Create.
- Upload data and Use the Access link to import data from Bucket to SSH Console.

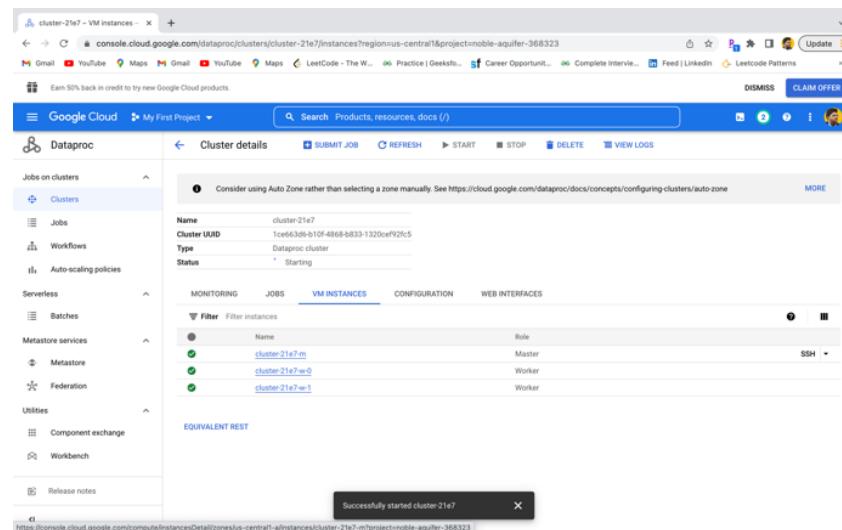
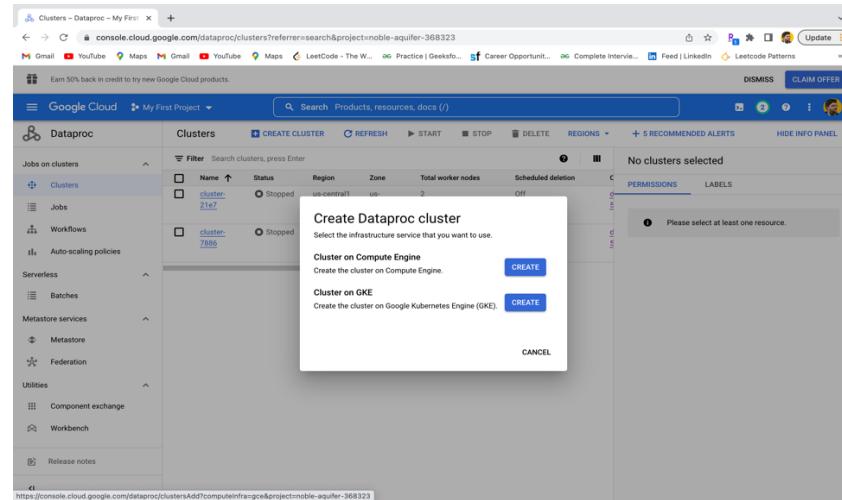
The screenshots illustrate the process of creating a bucket and uploading data to Google Cloud Storage. The top image shows the main Cloud Storage dashboard with several buckets listed. The bottom image provides a detailed view of a specific file's metadata, including download and edit options.

Object details for Automated_Traffic_Volume_Counts.csv

Overview	Value
Type	text/csv
Size	3.1 GB
Created	4 Dec 2022, 21:24:18
Last modified	4 Dec 2022, 21:24:18
Storage class	Standard
Custom time	—
Public URL	Not applicable
Authenticated URL	https://storage.cloud.google.com/nytraffic/Automated_Traffic_Volume_Counts.csv
gsutil URI	gs://nytraffic/Automated_Traffic_Volume_Counts.csv
Permission	Public access
Protection	Not public
Hold status	None
Version history	None
Retention policy	None
Encryption type	Google-managed key

2. Open DataProc Services and create clusters.

- In the Google Cloud console, go to the Dataproc Clusters page. GO TO CLUSTERS
- Click Create cluster.
- In the Create Dataproc cluster dialog, click Create in the Cluster on Compute engine row.
- In the Cluster Name field, enter example-cluster.
- In the Region and Zone lists, select a region and zone
- For all the other options, use the default settings.
- To create the cluster, click Create. Your new cluster appears in a list on the Clusters page. The status is Provisioning until the cluster is ready to use, and then the status changes to Running. Provisioning the cluster might take a couple of minutes.



3. SSH Terminal / Jupyter Notebook to Execution

- Start pyspark.
- customerDF = spark.read.csv("write address of bucket",header=True)
customerDF.show()
- write spark sql queries.
- Use this command to save output in bucket
db.write.save("write address of bucket ",format='csv',header=True)

```

Inbox (1,503) - saurabhjain | WhatsApp | cluster-7886 - Monitoring | https://ssh.cloud.google.com/ | (25) Using PySpark on Dat... | Cloud Monitoring | Data... | + 
← C ssh.cloud.google.com/x/ssh/projects/noble-aquifer-368323/zones/us-central1-f/instances/cluster-7886-m?authuser=0&hl=en_GB&projectNumber=574868...
Gmail YouTube Maps Gmail YouTube Maps LeetCode - The W... Practice | Geeksfo... Career Opportunit... Complete Interview Feed | LinkedIn Leetcode Patterns

SSH-in-browser
Linux cluster-7886-m 5.10.0-0.deb10.16-amd64 #1 SMP Debian 5.10.127-2-bpo10+1 (2022-07-28) x86_64
The programs included with the Debian GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*copyright.

Debian GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Mon Mar 25 04:10:50 2022 from 35.235.244.34
saurabhjain@cluster-7886-m:~$ spark-submit gs://nyctraffic/Automated_Traffic_Volume_Counts.csv gs://nyctraffic/Automated_Traffic_Volume_Counts.csv: Unsupported scheme 'gs'.
saurabhjain@cluster-7886-m:~$ pyspark
Python 3.13.0 | SparkSession created by code-qafoe | (default, Mar 25 2022, 06:04:10)
[Py4J] 10.3.0 on linux
Type "help", "copyright", "credits" or "license" for more information.
Setting default log level to "WARN".
To adjust log level, use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
22/11/14 05:17:06 INFO org.apache.spark.SparkEnv: Registering MemoryTracker
22/11/14 05:17:06 INFO org.apache.spark.SparkEnv: Registering BlockManagerMaster
22/11/14 05:17:06 INFO org.apache.spark.SparkEnv: Registering BlockManagerMasterHeartbeat
22/11/14 05:17:06 INFO org.apache.spark.SparkEnv: Registering OutputCommitCoordinator
Welcome to
      ____          _          _/_/    _/_/_/    _/_/_/_/_/_/_/_/_/      version 3.1.3
      / \     / \ / \ / \ / \ / \     / \ / \ / \ / \ / \
      \ /   / \ / \ / \ / \ / \ / \   / \ / \ / \ / \ / \ / \ / \
      / \   / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \
      \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \
      / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \ / \
Using Python version 3.8.13 (default, Mar 25 2022 06:04:10)
Spark context Web UI available at http://cluster-7886-m.us-central1-f.c.noble-aquifer-368323.internal:46159
Spark context available as 'sc' (master = yarn, app id = application_166840173590_0002).
SparkSession available as 'spark'.
>>> df = spark.read.csv(gs://nyctraffic/Automated_Traffic_Volume_Counts.csv",header=True,inferSchema=True)
>>> trips_df.show(10)
+-----+-----+
|RequestID|Boro|Yr|M|MM|Vol|SegmentID|WktGeom|street|fromSt|toSt|Direction|
+-----+-----+
| 20851|Queens (2015)|6|23|1|30|9| 9896|POINT (105229...| 94 AVENUE| 207 Street|Francis Lewis Bou...| WB|  
| 21231|State Island (2015)|9|14|4|15|6| 9896|POINT (942668....| RICHMOND TERRACE| Wright Avenue| Emeric Court| WB|  
| 22928|Bronx (2017)|6|24|1|30|10| 7283|POINT (1013633....| HUNTS POINT AV| Whittier Street| Rainford Avenue| NB|  
| 27019|Brooklyn (2017)|7|18|1|30|16| 10168|POINT (989255....| FLATBUSH AVENUE| Bright St| Brighton Line| NB|  
| 26734|Manhattan (2017)|11|3|22|0|3555| 137516|POINT (1004175....| WASHINGTON BRIDGE|Harlem River Shore...|Harlem River Shore...| EB|  
| 26015|Bronx (2017)|6|17|1|30|11| 86053|POINT (1021709....| WALLACE AVENUE| Rhinelander Avenue| Bronxdale Avenue| NB|  
| 26323|Manhattan (2017)|11|3|21|0|3599| 100130|POINT (1021709....| S/B AMSTERDAM ST| 20th Street| Cullinan Avenue| SB|  
| 23133|Queens (2016)|3|21|9|15|322| 101100|POINT (1050277....| NORTHERN BOULEVARD| 20 Place| 220 Street| WB|  
| 32417|Queens (2020)|11|14|2|15|18| 147877|POINT (1044172....| MIDLAND PARKWAY| Dalyn Road| Connector| SB|  
| 26198|Bronx (2017)|6|22|4|30|21| 85935|POINT (1021747....| THIERIOT AVENUE| Gieson Avenue| Pelham Line| NB|
+-----+-----+
only showing top 10 rows

```

```

Inbox (1,503) - saurabhjain | WhatsApp | cluster-7886 - Monitoring | https://ssh.cloud.google.com/ | (25) Using PySpark on Dat... | Cloud Monitoring | Data... | + 
← C ssh.cloud.google.com/x/ssh/projects/noble-aquifer-368323/zones/us-central1-f/instances/cluster-7886-m?authuser=0&hl=en_GB&projectNumber=574868...
Gmail YouTube Maps Gmail YouTube Maps LeetCode - The W... Practice | Geeksfo... Career Opportunit... Complete Interview Feed | LinkedIn Leetcode Patterns

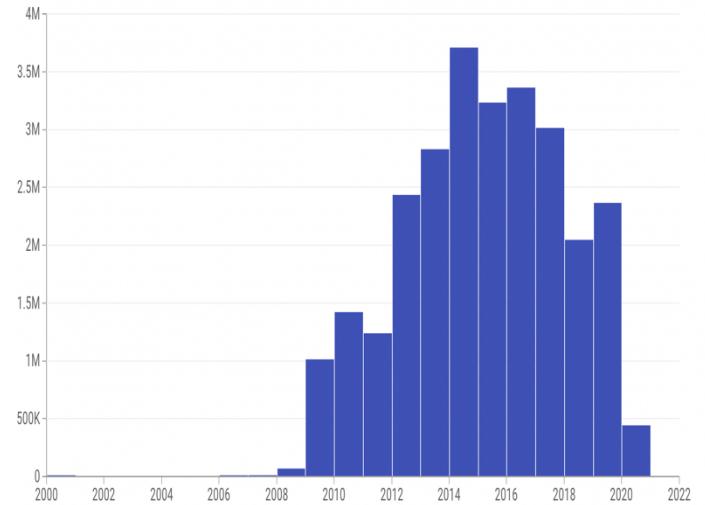
SSH-in-browser
22761|Manhattan (2000)|1|5|7|30|216| 32956|POINT (988684.7 ...|WB UNION SQUARE ...| E 14 ST/BROADWAY| 5 AV| WB|  
7625|Queens (2008)|4|1|5|30|148| 158808|POINT (898240.431...| 141367|POINT (1013833....| JACKIE ROBINSON W...| MILK AV| VERMONT ST| ED|  
7471|Brooklyn (2007)|4|24|1|30|11| 158808|POINT (1013833....| CORNELIA ST| WYCROFF AV| CYPRESS AV| NB|  
7618|Brooklyn (2007)|4|12|22|45|11| 45393|POINT (1011005....| LINDEN BLVD| 78 ST| AMBER ST| WB|  
7624|Queens (2007)|4|18|1|4|0| 190387|POINT (1023782....| 8 AVENUE| DEAD END/BELT PKWY EB EXIT...| ED|  
3442|Bronx (2008)|4|13|8|1|164| 88874|POINT (1020685....| S/B VAN CORTLANDT| E 242 ST/KATONAH AV| KATONAH AVENUE| SB|  
3473|Queens (2008)|6|5|11|0|3291| 138592|POINT (1060197....| E/B LAURELTON C...| DEAD END/BELT PKWY EB EXIT...| ED|  
7489|Queens (2008)|5|26|15|45|79| 153405|POINT (1004881....| E/B KOCSIUSZKO BR.| 54 RD| BQE EB EXIT 35 E-W| EB|  
7499|Queens (2008)|3|17|10|45|11| 136895|POINT (1066765....| E/B UNION TURNPIKE| LANGDALE ST| LANGDALE ST| ED|  
3479|Manhattan (2008)|1|5|1|30|200| 158808|POINT (1013833....| 141367|POINT (1013833....| DEADERICK ST| 5TH AV| ED|  
7521|Queens (2008)|3|17|11|30|128| 156095|POINT (1059170....| E/B MERRICK BLVD| 244 ST| 243 ST| EB|  
7463|Bronx (2008)|5|13|8|45|694| 190807|POINT (1031782....| N/B HUTCHINSON RILL| NEW ENGLAND THRWY| NEW ENGLAND THRWY| NB|  
7464|Bronx (2008)|5|1|8|17|30|0| 154595|POINT (1031782....| HUTCHINSON RIVER...| HUTCHINSON RIVER...| HUTCHINSON RIVER...| NB|  
3475|Manhattan (2008)|7|23|14|15|91| 23255|POINT (981449.3 ...| E/B CHAMBER ST @...| HUDSON ST/W BROADWAY| GREENWICH ST| EB|  
7600|Bronx (2008)|5|1|11|1|30|9| 135125|POINT (1015681....| S/B HENRY HUDSON ...| DEAD END| MOSHOLU PKWY| SB|  
3434|Bronx (2008)|6|23|22|3|0| 88679|POINT (1013140....| NEW YORK AVENUE| CONNECTOR| DEAD END| SB|  
3430|Brooklyn (2008)|3|13|15|30|106| 158808|POINT (1050277....| E/B AVE BRIDGE| 141367|POINT (1013833....| 3 ST| NB|  
7464|Queens (2008)|5|1|7|30|1006| 154955|POINT (1033736....| N/B HUTCHINSON RILL| HUTCHINSON RIVER...| HUTCHINSON RIVER...| NB|
+-----+-----+
only showing top 20 rows
>>> new_results = spark.sql("select * from Traffic order by Yr DESC").show()
+-----+-----+
|RequestID|Boro|Yr|M|Di|HH|MM|Vol|SegmentID|WktGeom|street|fromSt|toSt|Direction|
+-----+-----+
|\F0\fs24\cf0 Req...|Boro|Yr|M|Di|HH|MM|Vol|SegmentID|WktGeom|street|fromSt|toSt|Direction|
| 32384|Manhattan (2020)|10|20|11|30|8| 158808|POINT (898240.431...| EAST 14 STREET| Park Avenue| Morris Avenue| WB\|  
| 32407|Bronx (2020)|10|19|11|0|140| 98936|POINT (1005248....| 141367|POINT (1013833....| QUEENSBOROUGH RD| Dead End| Dead end| WB\|  
| 32384|Manhattan (2020)|11|1|10|0|42| 34257|POINT (898423.209...| CENTRAL PARK WEST| 8 Avenue Line| 8 Avenue Line| NB\|  
| 32407|Manhattan (2020)|10|19|11|30|10| 164620|POINT (1008856....| GRAND CONCOURSE| Concord Line| Concord Line| NB\|  
| 32417|Queens (2020)|11|1|6|8|1| 156673|POINT (1041559....| UTOPIA PARKWAY| 65 Avenue| 67 Avenue| SB\|  
| 32417|Queens (2020)|10|13|15|1|30|194| 150549|POINT (1031410....| CROSS BAY BOULEVARD| Dead End| Connector| NB\|  
| 32384|Manhattan (2020)|10|1|28|11|30|195| 34338|POINT (987874.166...| AMSTERDAM AVENUE| West 60 Street| West 61 Street| NB\|  
| 32407|Bronx (2020)|10|1|3|8|0| 2200| 78877|POINT (1015489....| WESTCHESTER AVENUE| East 167 Street| Home Street| WB\|  
| 32417|Queens (2020)|11|1|8|1|15|167| 150346|POINT (1044907....| MERRICK BOULEVARD| 111 Avenue| Dead end| SB\|  
| 3196|Queens (2020)|3|26|0|0|0| 171277|POINT (1054176....| BEACH 14 STREET| Heyson Road| Seagirt Boulevard| SB\|  
| 32417|Queens (2020)|3|26|1|15|53| 150549|POINT (1031410....| KISSENA BOULEVARD| 41st Street| Barretto Street| ED|  
| 32417|Queens (2020)|11|22|19|15|61| 150539|POINT (103832.08...| 164 STREET| 65 Avenue| 67 Avenue| SB\|  
| 32417|Queens (2020)|10|17|22|30|75| 57192|POINT (1038292....| LIBERTY AVENUE| Waltham Street| 146 Street| EB\|  
| 32407|Bronx (2020)|10|18|1|4|34| 72026|POINT (1007749....| JEROME AVENUE| Goble Place| East Mt Eden Avenue| NB\|  
| 32417|Queens (2020)|10|18|1|0|57| 146374|POINT (1068638....| UNION TURNPIKE| Hewlett Street| Hewlett Street| WB\|  
| 32407|Bronx (2020)|10|20|4|45|22| 70376|POINT (100811.29...| 3 AVENUE|East 154 Street| East 155 Street| NB\|  
| 32384|Manhattan (2020)|10|27|2|0|27| 173249|POINT (1002839....| 145 STREET BRIDGE| Harlen River Shor...| EBB|  
| 32407|Bronx (2020)|10|15|13|0|75| 212069|POINT (1013077....| CROTONA AVENUE| Dead End| Crotona Park North| SB\|
+-----+-----+
only showing top 20 rows
>>>

```

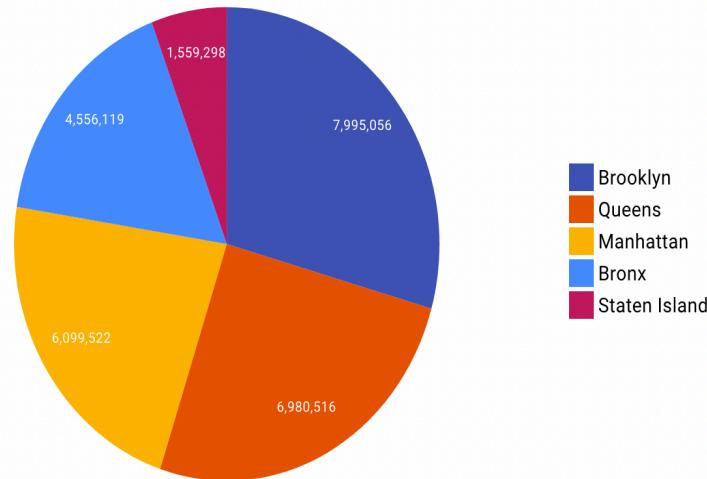
TEST RESULT

Data visualization Using NYC-Automated-Traffic-Volume-Counts

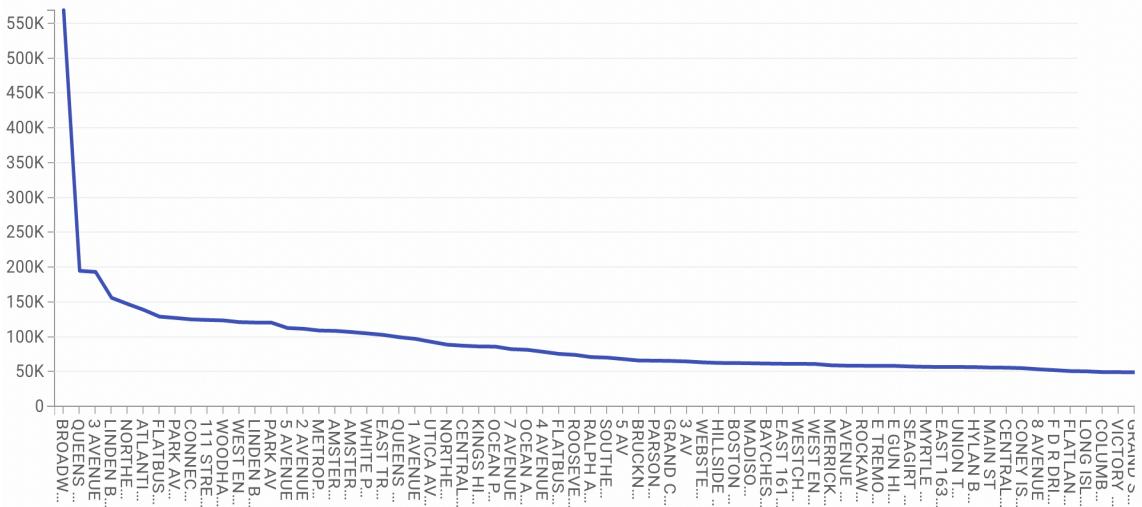
- Variation in Volume of data/vehicle: After pre-processing, the main analysis of this traffic dataset shows the number of increases in the volume of cars from the year 2000 to 2020. (Queries being displayed in the implementation section)



- Busiest Cities of NYC: After pre-processing and running queries and we extracted the top 3 busiest cities and, it found that Brooklyn is the most crowded and busiest city with the count of vehicles (7,995,056, 29%), Queens being the 2nd busiest (6,980,516, 26%) and Manhattan being the 3rd busiest (6,099,052, 22%).

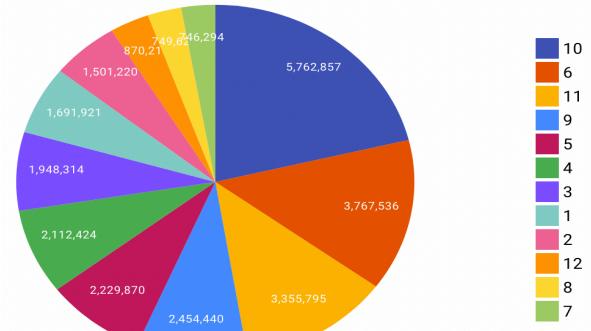


- Most visited streets: From the dataset, we found the most visited and busiest street in different cities in NYC.



- The volume of car increased every 15 mins in New York: This query displays how many numbers of cars increased city of Boro.
- Grouped according to Busiest Month in the year of different Cities: By applying group by, sorting, and count query we found that October was the busiest month in Manhattan in the entire year.

- 1st pie chart shows the number distribution according to month.



- 2nd pie chart shows the busiest city according to the month.

