# Analysis of Los-Angeles Parking citation Data In Hadoop Environment Using Map-Reduce Techniques

Sindhujan Dhayalan,X17170265
Programming for Data Analytics project
M.Sc Data Analytics - National College of Ireland

*Abstract*—**Parking in a busy city like Los-Angeles is troublesome. Chances are that if you drive in Los-Angeles, you are most likely to end up with a parking ticket. Which is a problem when the parking ticket issued in LA averages to $70 in the year 2017 and it doubles when paid late. Statistics show that the city generates large revenues from parking tickets as much that some percentage of the citys budget plan depends on revenues from parking tickets. This project is focused on analyzing big data available on parking tickets in LA open data portal. To analyze the available Big data, Big data processing environment with the combination of Hadoop distributed environment, Map-Reduce programming using Pig and Java, SQL: MySQL and NoSQL: HBase databases for data storage is used. The analysis is focused on finding hidden patterns which can be used to derive knowledge to lessen the burden on people driving caused by parking citations in much populated city like LA. Key findings were even though the parking tickets are expensive, Average ticket pricing has reduced in recent years. The Brand, Body-Style, Colour, Registered state plate of vehicles most involved in parking citations were found. The Top Route and Violation on which parking citations were issued are found.**

*Index Terms*—**Hadoop, Map-Reduce, Pig, HBase, MySQL, JAVA, R Programming.**

## I. INTRODUCTION

Parking tickets are issued for violating parking laws by parking enforcement official set by the State or the city. In cities, highly populated and busy 24/7, parking your vehicle becomes an issue. It becomes hard to find parking space and probably you will end up with a parking ticket. Cities like Los-Angeles, A single parking ticket can cause you more than $65 and late fees double the charges, there are complaints by people for reducing the high rate at parking ticket. It is very hard for everyone driving to pay for expensive parking tickets and people living in the city are unhappy and frustrated about it. But the city depends on high revenues from parking tickets. Published report[1] by City controller Ron Galperin states that LA generated $148 million in gross ticket revenues in the financial year 2015-2016 and Data from Kaggle[2] suggests in the year 2017, People paid a sum of $156 million on parking citations. Most of the revenues go to salaries and administrative costs. The remaining quarter of the revenues is spent on city services through general funds. Ron Galperin has states solutions such as Digital parking signs, Re-evaluation of Street sweeping schedules, Smartphone application on maintaining parking tickets rather than reducing the parking charges. This project is focused on analyzing the data available

on Parking citations in LA city over the period of 2010 to 2018 using Hadoop distributed environment combined with MySQL and HBase databases and using Java and Pig Map-reduce programming. The data sourced from Kaggle is downloaded and stored in MySQL database then loaded to HDFS, where it is processed by Map-reduce programs using Java and Pig, later stored in HBase. The processed data is visualized using Rscript to identify hidden patterns, That can provide insights to reduce the burden of expensive parking tickets on people driving in the city of Los-Angeles. Data analysis is focused on identifying insights on following queries:

1) Average parking fine paid on Each day of a month, Each month of a year and Each year in the period of 2010-2018.
2) Total amount of fines paid by vehicles registered with different states apart from CA.
3) Brand of vehicles and their year of incident involved in parking citations.
4) Top 5 Routes, Violations and Month of the incident involved in parking citations.
5) The Body style, Colour, Brand and Routes of vehicles involved in most parking tickets in the year 2018.

The paper is organized as the following in the upcoming sections, Section 2: Related work, Section 3: Methodology, Section 4: Results followed by, Section 5, Conclusion and future work.

## II. RELATED WORK

Research works using Hadoop MapReduce is mostly focused on Social media analysis, Crime analysis etc. Using Hadoop MapReduce for data analytics of parking citation proposed in this project is a novel approach. But there are many research using the same approach of MapReduce technology applied for this project in other research domains.

Research paper[1] discusses applying Big data analytics on crime data. The framework designed uses Apache sqoop to move data from MySQL database into HDFS.Apache sqoop connects with Relational database management systems and transfers large amount of data from RDBMS to HDFS.Pig Latin Script is used to run MapReduce program on the data stored in HDFS and the output from Pig is stored back to HDFS. The advantages of using Pig Script instead of JAVA script is the less complexity of Pig code. Hence reducing the overall time of development and testing. Pig scripts consume about 5% time compared to writing MapReduce Java programs. But execution time of Pig Scripts are 50% longer

---
[1]http://parking.controlpanel.la/
[2]https://www.kaggle.com/cityofLA/los-angeles-parking-citations

than that of Java MapReduce programs. The output from Pig MapReduce is later used for Visualization to find hidden patterns.

Paper[2] successfully implemented Big data analytics of E-Commerce consumer data on Hadoop platform to understand consumer consumption behaviour. The unstructured data is stored in HDFS as it provides Scalability and data security. Data stored in HDFS is processed using Hadoop's MapReduce framework. The processed data is then transferred to MySQL using Sqoop with help of shell command. Data stored in MySQL was visualized using technologies such as Eclipse platform, jsp, servlet, jdbc, jFreechart.

Hadoop MapReduce Framework on AWS cloud platform is used for data analysis of YouTube data in paper[3]. YouTube video data is collected using YouTube API and stored as CSV file format. The data collected is then fed to HDFS. MapReduce program coded in JAVA programming language is run on the data to perform summary operations to find the top 5 video categories, top 5 video uploaders etc. This output is stored back in HDFS and used for visualization to provide knowledge to end users.

In research paper[4], Author proposed a model to analyze weather logs quickly using Hadoop parallel processing. the pre-processed weather data is loaded into HDFS and MapReduce framework using Java programming is applied to analyze yearly maximum and minimum temperature. The output file from MapReduce is stored in HBase. The data in HBase is visualized in line charts to display the maximum and minimum temperature. The results were that Hadoop did not handle the integration of large number of small files into Large file as the size of the file was bigger than the HDFS block size.

## III. METHODOLOGY

### A. Description of Dataset

The Parking citation Dataset for this project was sourced from Los Angeles open data platform[3], Also available on kaggle[4], the dataset has 9.26 million rows and 19 columns. The 19 columns contain data about the ticket number, Issued date, Issue time, Meter ID, Marked time, RP State plate, Plate expiry date, VIN, Make, Body style, Colour, Location, Route, Agency, Violation code, Violation description, Fine amount, Latitude and Longitude recorded from year 2010 to 2018. The dataset was created on 2018-08-16 and Last updated on 2019-04-23. The Challenges of using dataset is its size as Ms-excel cannot open 9.26 million rows and many records in the dataset contained null and duplicate values.

### B. Data Processing

The Dataset was first pre-processed using RScript, Out of the 19 columns, 10 columns was considered for analysis on queries addressed in this project, they are ticket number, Issued date, RP State plate, Make, Body style, Colour, Route, Agency, Violation description, Fine amount. The tasks carried out using

RScript were, Records with null values were removed, Column names were modified and Issued date column was divided into 3 separate columns: Year, Month and Day, Which will be useful input for Map-Reduce programs. Further, the pre-processed data is imported into MySQL database. For this purpose, first a database 'Parking' and table 'bike' was created in MySQL. Figure1 shows the created schema of 'bike' in MySQL database 'Parking'.



Fig. 1. Schema created in MySQL database

On querying the table, Duplicate Ticket ID records were found, which were removed using MySQL command and finally, A table with 336376 records and 13 fields were produced for use in this project. The Count and Sample of the table created is displayed in figure2 and figure3.



Fig. 2. Table 'bike' count



Fig. 3. MySQL table records - Sample

The pre-processed data in MySQL database is moved and stored in HDFS file system using Sqoop as shown in figure4. Sqoop is a tool designed for transferring big data between Hadoop and RDBMS. Sqoop uses MapReduce to import and export the data, providing fault tolerance and parallel processing[5].



Fig. 4. MySQL to HDFS using SQOOP

---

[3]https://data.lacity.org/A-Well-Run-City/Parking-Citations/wjz9-h9np
[4]https://www.kaggle.com/cityofLA/los-angeles-parking-citations

The dataset in HDFS was processed using 3 Java Map-Reduce programs using eclipse platform and 2 Pig Map-Reduce programs on the basis of the queries to be addressed in the project. Output from the Map-Reduce programs were stored in HDFS and then transferred to HBase. The files are then exported to perform visualization using RScript. Finally the visualization produced was observed and knowledge were derived to answer the queries mentioned in section I. The process flow of the project is displayed in Figure5.
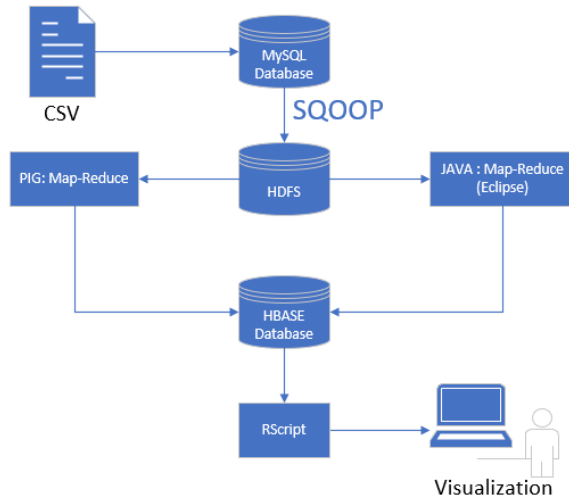


Fig. 5. Process flow

MapReduce tasks performed using JAVA program and Pig Scripts are explained in the following sections.

*C. JAVA Map-Reduce*

MapReduce is a programming model used to process big data. The Map function process key/value pair as input and generate a set of intermediate key/value pairs. Then the reduce function merges all intermediate values associated with same intermediate key. Google uses MapReduce algorithms for processing graphs, texts, Machine learning and performing statistical operations. MapReduce is fault tolerant and is a very useful tool in big data analysis[6]. For this project, JAVA Map-Reduce programs were performed using Eclipse-Hadoop integration. The input file path is from the HDFS and the output file of Eclipse is stored back in HDFS for the 3 tasks discussed below:

Three Classes: Driver, Mapper and Reducer were used to perform the Map-Reduce task. Figure6 displays the classes created for performing the tasks.

**Task1 - State & Fine**
Objective of the task is to find the total amount of fine payed by vehicles with different registration state plate. For this purpose, Words in the fields 'State' as Key and the total amount of fine paid by that state is produced as value output.
Input : State and Fine & Output: Key: State & Value: Fine
**Task2 - Count(Make,Year)**
Objective of the task is to count different Make of vehicle reported in each year on parking citation. For this purpose,
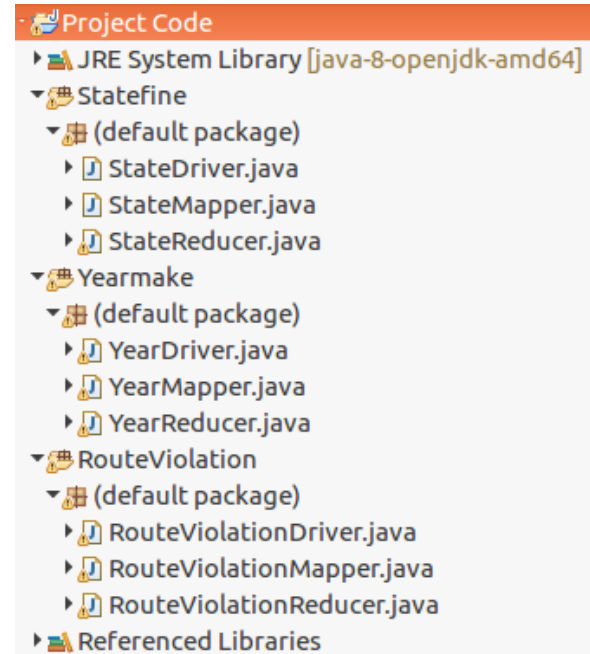


Fig. 6. Java Map-Reduce Classes in Eclipse

Words in the fields 'Make' and 'Year' are concatenated as Key and their occurrence together is counted and produced as value output.
Input : Make and Year
Output: Key: (Make,Year) & Value: Count(Make,Year)
**Task3 - Count(Month,Violence,Route)**
Objective of the task is to count number of parking violation reported in every route in each month. For this purpose, Words in the fields 'Month','Violation' and 'Route' are concatenated as Key and their occurrence together is counted and produced as value output.
Input : Month, Violation and Route
Output:Key:(Month,Violation,Route)
Value: Count(Month,Violation,Route)



Fig. 7. Output of Java Map-Reduce stored in HDFS

*D. Pig Map-Reduce*

Apache Pig is a high level language platform used for analyzing large data sets. Pig Latin can execute Hadoop jobs using MapReduce[7]. Pig-MapReduce process was carried out by creating two Pig Latin scripts for Task 4(pig_cmd.pig) and Task5(pig_cmdcount.pig) and the scripts were run on Apache Pig Grunt shell. The input path was main source dataset stored in HDFS and Output was stored back in HDFS system.

Fig. 8. Output of Task1, Task2 and Task3

**Task4 - Average Fine on Days, Months and Year**
Objective of Task 4 is to calculate the average fine paid on daily, Monthly and yearly basis over a period from 2010 to 2018. Three output files were stored in HDFS in this task for each Days, Months and Years respectively. Figure9 displays the path of output files stored in HDFS. Figure10 and Figure11 displays the Output of Task 4 Map-Reduce program.



Fig. 9. Pig-Task4 Output files stored in HDFS



Fig. 10. Pig-Task4 Average fine on each Month & Year

**Task5 - Count of Body-style, Make, Colour & Route**
Objective of Task5 is to calculate the Count of Body-style, Make, Colour and Route of vehicles involved in parking violation in the year 2018. Four output files were stored in HDFS in this task for Body-style, Make, Colour and Route respectively. Figure12 displays Collage screen-shots of Output of MapReduce performed in Task5.



Fig. 11. Pig-Task4 Average fine on each day of the month



Fig. 12. Pig-Task5 Output

### E. Apache HBase

HBase is an open source, NoSQL column based database written in Java. It runs on top of HDFS (Hadoop distributed file system). It is mainly used to store large data with sparsity in a fault tolerant way. Data is stored in column based model in the form key/Value pair and MapReduce can be used to query the data. All the columns made are grouped into one column family[8].

In this project, The Output from the Map-Reduce programs stored in HDFS is moved to HBase Database. Initially, HBase tables are created with table name and column family. Later, using the ImportTsv command, The data in HDFS file is stored into the HBase table. Figure13 displays the Commands used for importing the data into the table and imported data present in the HBase table.
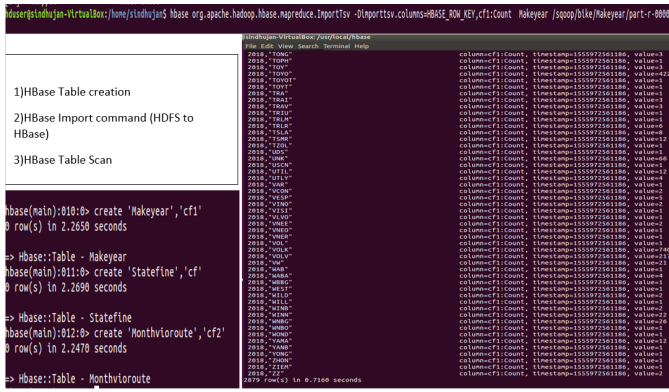
Fig. 13. HBase Table creation & Imported data in HBase table

## IV. RESULTS

The Output files stored in HBase Database were visualized using RScript to gain knowledge, that was used to address the queries mentioned in section I. This section discusses results of visualization of the outputs from MapReduce tasks.

### A. Task 1 - State & Fine



Fig. 14. Word Cloud of State plates with most fines

From the output of Task1, It was observed that vehicles registered with 'CA' state plate had 316910 records (10 times) more parking violations than other states. This is not surprising as the dataset is based on the Parking citations in Los-Angeles which is a city in CA. Being an outlier this value was ignored while visualizing and Aim was to find the vehicles registered outside CA having more parking violation in Los Angeles. Observing Word cloud in figure14, It is inferred that vehicles with state plate registered in AZ, TX and NV are vehicles involved in parking violations apart from CA.

### B. Task 2 - Make & Year

Task 2 output file consists of the count of Year and Make of vehicle with most parking violations. Word cloud was again used to find the most occurring terms and it was observed



Fig. 15. Word Cloud of Make,Year with most fines

from figure15 that vehicles of brands 'FORD', 'CHEV' and 'NIS' were mostly involved in parking violations and Parking violations were more frequent in year 2015 compared to other years.

### C. Task 3 - Month, Violation & Route



Fig. 16. Horizontal bar-plot of top 5 routes with most parking violation with violation description and Month

Task 3 objective was to find the frequently occurring parking violation incident, its location and the month at which these incident happen. It is inferred from the figure16 that 'NO PARK/STREET CLEAN' occurs frequently in the routes '00141' & '00146' around the months between March to May over the period from 2010 to 2018.

### D. Task 4 - Average fine on each Day, Month and Year

Task 4 is carried out to find the Day of the Month, Month of the Year and Year among other Years with most average fine. From the line graphs visualized, It can be inferred that, As of the Day of the Month in figure17, the Average fine value
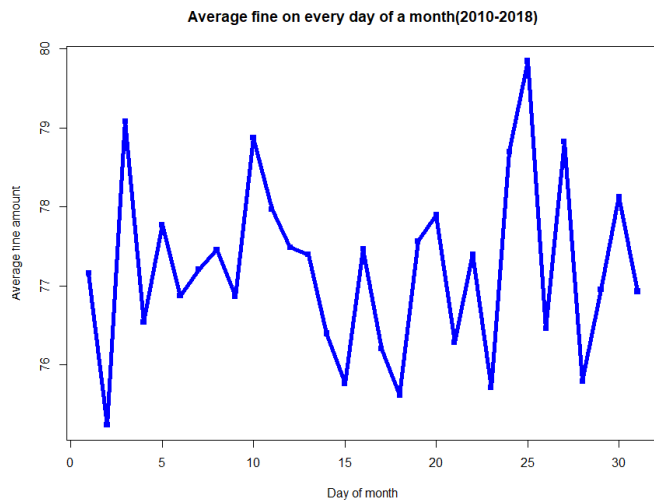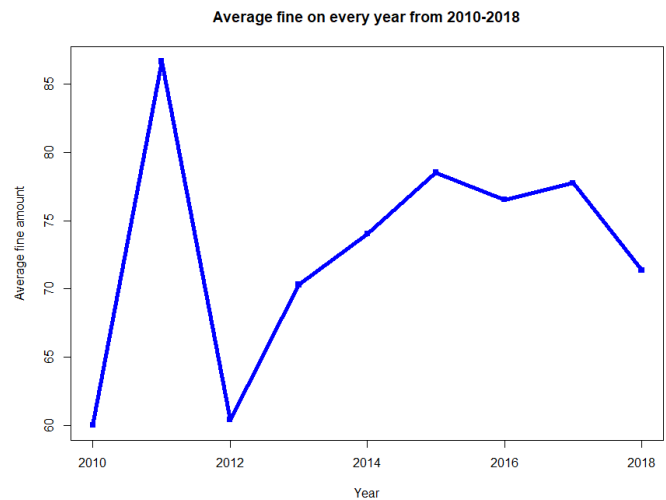
Fig. 17. Average fine every day of Month

is unstable, But the most fine is paid on 25th date of every month.
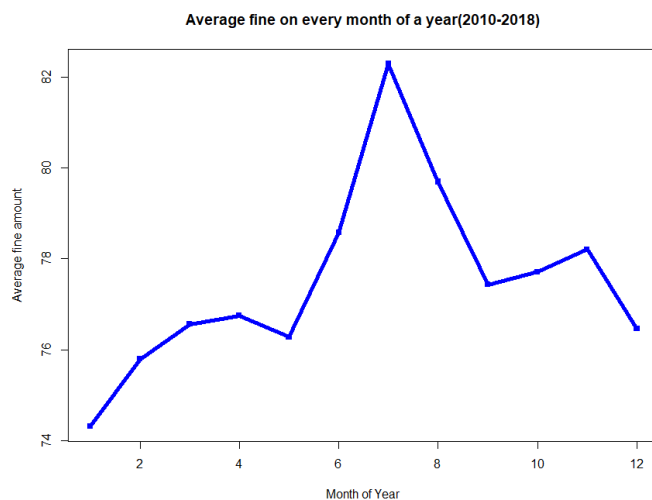


Fig. 18. Average fine on every month of year

Figure18 shows that Most of the fines are paid on months from July to August and July being the month fines are paid the most. Figure19 shows that the average amount of fines paid were highest in 2011 and are gradually decreasing in the recent years from 2016 to 2018.

### E. Task 5 - Body-style, Make, Colour & Route

Task 5 output file contains the count of Body-style, Make and colour of vehicles involved in parking violation incidents and the routes in the violation incidents have happened in year 2018. Figure20 shows top 10 body-style of vehicles involved in parking violation in year 2018 and it is inferred that most of vehicles with body style 'PA'(20420 records) commit about 10 times more parking violations than other body-style of vehicles.



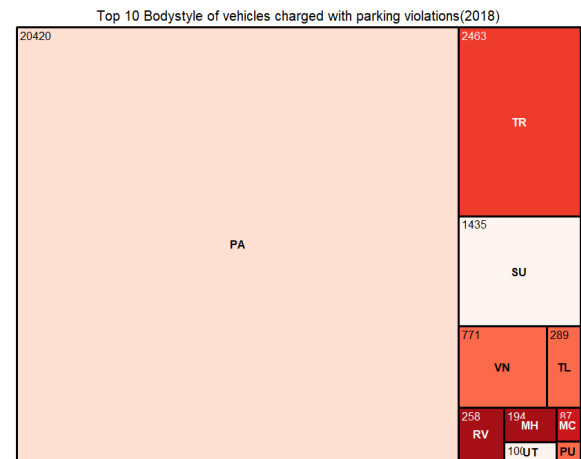Fig. 19. Average fine on every year(2010-2018)



Fig. 20. Top 10 Body-style of vehicles with most parking violations(2018)

Figure21 shows the top 5 colours of vehicles involved in most parking violations and it is inferred from the circle packing that vehicles with colour black, white and grey are involved the most in year 2018.

From Bar chart in figure22, it is inferred that people using 'TOYOTA(TOYO)' brand of vehicles are having trouble with parking the most in the year 2018. the gap between TOYOTA and other brands are significantly large(¿1000). Figure23 displays the Top 10 routes with most parking violations, it is inferred from the graphical representation that Route '00500' has 2719 parking citations in year 2018, The value is 1500 incident more than the second route '00001'.
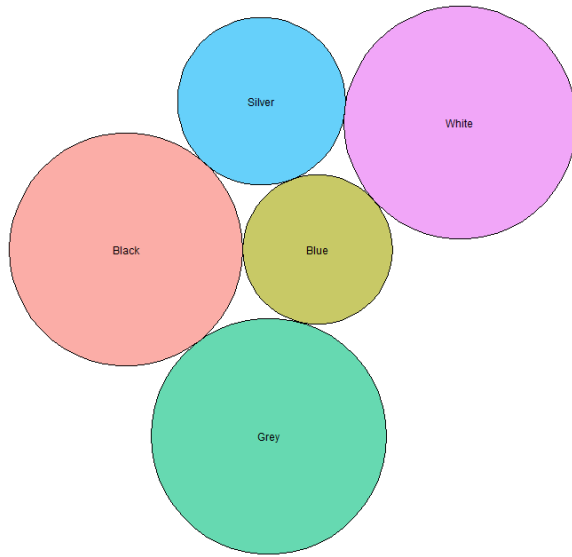
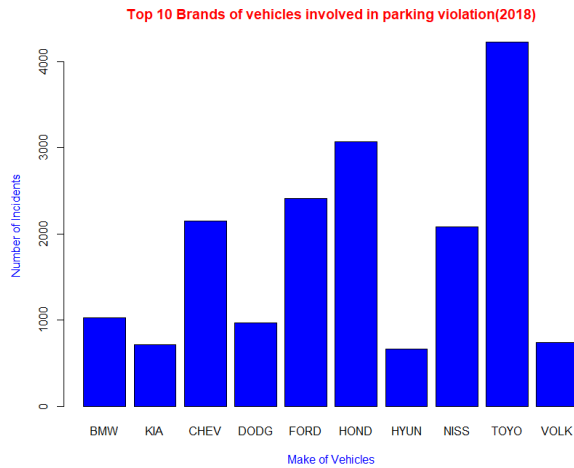Fig. 21. Top 5 Colour of vehicles involved in Parking violations(2018)



Fig. 22. Top 10 Brands of vehicles involved in Parking violations(2018)
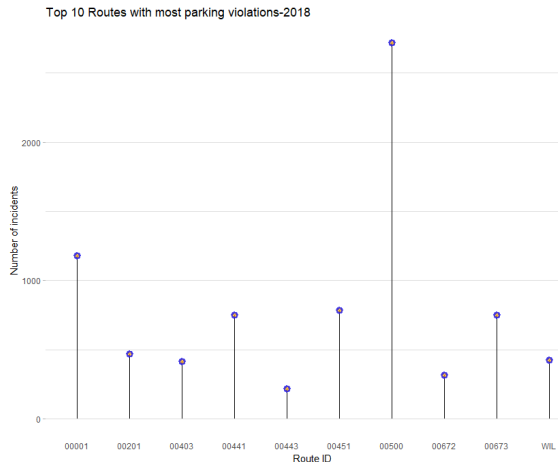


Fig. 23. Top 10 Routes with most Parking violations(2018)

## V. CONCLUSION AND FUTURE WORK

From the previous section, a lot of insights were gathered, Visualizing Data showed many interesting results. Apart from 'CA' registered vehicles, vehicles registered under 'AZ','TX' and 'NV' commit the most parking violation in Los-Angeles. The most violations are made by 'CHEV','FORD' and 'NISS' and frequently in the year 2015. Most violations made were 'NO PARK/STREET CLEAN' in the routes '00141' and '00146' around the months March to May in year range(2010 to 2018). It was found that average fines paid were highest on 25th of each month, Months of July had highest average paid fine and In year 2011, people were paying on average 86 American dollars per parking violation ticket, which has gradually decreased in the recent years. Visualizing data on recent year of 2018, it was found that people using vehicles designed with 'PA' body commit the most parking violations, vehicles colour with either Black, Grey or White are involved the most in parking violations. 'TOYOTA(TOYO)' vehicle users were involved the most in Parking violations and the most violations occurred in Route '00500'. The insights from recent years is that parking violations have decreased compared to previous years. Precautionary measures have been taken on routes '00141' and '00146' which are overall most parking violated routes, but they do not appear in the top 10 parking violated routes in 2018. Using the Hadoop MapReduce environment, we were successfully able to answer all the queries mentioned in the section I. The advantage here is most of the contributors to parking violations at the 1st place have a large gap from the other contributors, hence the government can focus on reducing the causes of parking violation by targeting the top contributors in each attributes as priority. thereby acting on the top most contributors can significantly reduce the overall parking violation problems. Challenges faced were the processing of 9 million rows. the system performance was low while processing of huge amount of data. Future work can be done on analyzing other attributes such as Agency, Issued time, latitude, longitude which were not used in this project to get more interesting knowledge and discover hidden patterns causing parking violations. In future, an approach to use the entire 9 million rows can be attempted.

## REFERENCES

[1] A. Jain and V. Bhatnagar, "Crime data analysis using pig with hadoop," *Procedia Computer Science*, vol. 78, pp. 571 – 578, 2016, 1st International Conference on Information Security Privacy 2015. [Online]. Available: http://www.sciencedirect.com/science/article/pii/S187705091600106X

[2] F. Shao and J. Yao, "The establishment of data analysis model about e-commerces behavior based on hadoop platform," in *2018 International Conference on Intelligent Transportation, Big Data Smart City (ICITBS)*, Jan 2018, pp. 436–439.

[3] P. Merla and Y. Liang, "Data analysis using hadoop mapreduce environment," in *2017 IEEE International Conference on Big Data (Big Data)*, Dec 2017, pp. 4783–4785.

[4] H. Wu, "Decision analysis of the weather log by hadoop," in *International Conference on Communication and Electronic Information Engineering (CEIE 2016)*. Atlantis Press, 2016/10. [Online]. Available: https://doi.org/10.2991/ceie-16.2017.3

[5] A. Scoop, "Apache sqoop," 2019. [Online]. Available: https://sqoop.apache.org/

[6] J. Dean and S. Ghemawat, "Mapreduce: a flexible data processing tool," *Communications of the ACM*, vol. 53, no. 1, pp. 72–77, 2010.

[7] A. Pig, "Welcome to apache pig!" [Online]. Available: https://pig.apache.org/

[8] L. George, *HBase: The Definitive Guide: Random Access to Your Planet-Size Data.* O'Reilly, Aug 2011. [Online]. Available: https://www.dawsonera.com:443/abstract/9781449315771