# Big Data Project - NYC Parking Tickets Analysis

| | |
|---|---|
| PES2201800504 | Aswin A Nair |
| PES2201800107 | R.Darsini |
| PES2201800036 | Ishaan Samant |
| PES2201800505 | Nikita Ganvkar |

## ABSTRACT

Context

The NYC Department of Finance collects data on every parking ticket issued in NYC (~10M per year). This data is made publicly available to aid in ticket resolution and to guide policymakers.

Content of Dataset
The file is roughly organized by fiscal year (July 1 - June 30 of 2015-2016) with the exception of the initial dataset. The column attributes include information such as the vehicle ticketed, the ticket issued, location, and time.
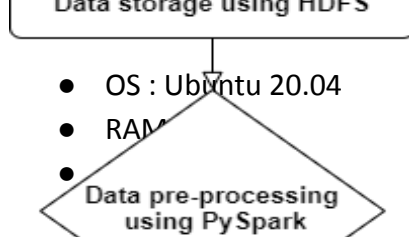
## AIM

The purpose of this project is to conduct exploratory data analysis and queries that will help us understand the data.

## OBJECTIVES

- Using Spark to perform EDA
- Executing queries with Hive and MapReduce
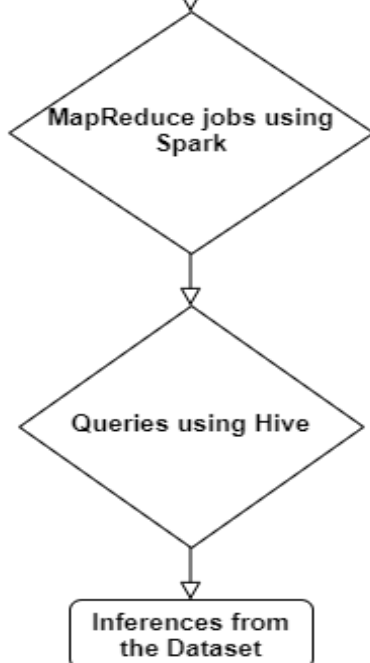- Get basic inferences about the dataset

# SYSTEM CONFIGURATION:



- OS : Ubuntu 20.04
- RAM
-

## DATASET

Parking Tickets: The NYC Department of Finance collects data on every parking ticket issued. EDA and Visualization made publicly available to aid in ticket resolution and to guide policy.
(https://www.kaggle.com/new-york-city/nyc-parking-tickets/data)
Parking Violation Codes: This dataset defines the parking violation codes in NYC and lists the fines.
(https://data.cityofnewyork.us/Transportation/DOF-Parking-Violation-Codes/ncbg-6agr)

PIPELINE

SPARK

Preprocessing

The flowchart nodes read:
- Data storage using HDFS
- Data pre-processing using PySpark
- EDA and Visualization using pyspark
- MapReduce jobs using Spark
- Queries using Hive
- Inferences from the Dataset

```
df = df.drop('Violation Post Code',
'Violation Description',
'No Standing or Stopping Violation',
'Hydrant Violation',
'Double Parking Violation',
'Latitude',
'Longitude',
'Community Board',
'Community Council',
'Census Tract',
'BIN',
'BBL',
'NTA')
df = df.dropna()
```

# EDA - Visualizations
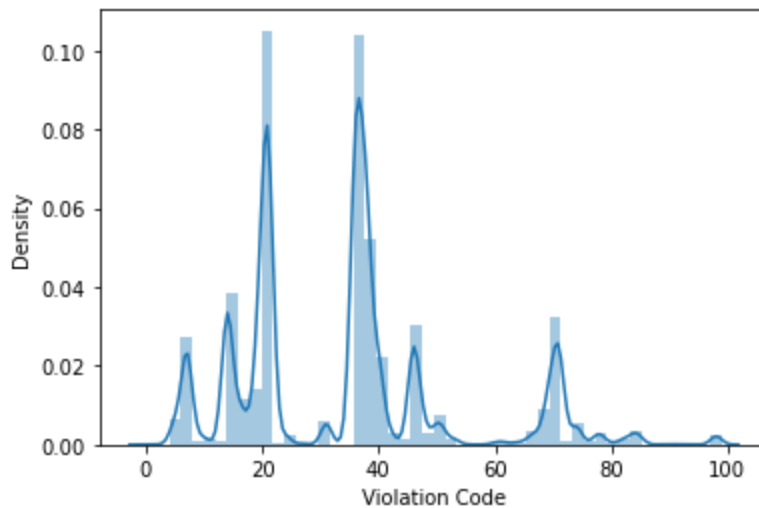
## Parking Tickets given each Month

```
: month = []
  for time_stamp in pd.to_datetime(mini['Issue Date']):
      month.append(time_stamp.month)
  m_count = pd.Series(month).value_counts()

  plt.figure(figsize=(12,8))
  sns.barplot(y=m_count.values, x=m_count.index, alpha=0.6)
  plt.title("Number of Parking Ticket Given Each Month", fontsize=16)
  plt.xlabel("Month", fontsize=16)
  plt.ylabel("No. of cars", fontsize=16)
  plt.show();
```

# Violation Code Distribution

```
sns.distplot(df['Violation Code'])
```



# Parking Ticket per County

```python
violation_county = mini['Violation County'].value_counts()

plt.figure(figsize=(12,8))
f = sns.barplot(y=violation_county.values, x=violation_county.index, alpha=0.6)
# remove labels
plt.tick_params(labelbottom='on')
plt.ylabel('No. of cars', fontsize=16);
plt.xlabel('County', fontsize=16);
plt.title('Parking ticket given in different counties', fontsize=16);
```

Parking ticket given in different counties

Average Amount of Fine for the top Plate Types

```python
import matplotlib.pyplot as plt
plt.xlabel("Plate Type")
plt.ylabel("Avg Fine Amount")
plt.bar(X,Y)
plt.show()
```

# MAP REDUCE

Query: Number of Violations per Registration State

```python
import sys
from pyspark import SparkConf, SparkContext
from csv import reader
conf = SparkConf().setAppName("MR_1")
sc = SparkContext(conf=conf)
line1 = sc.textFile("hdfs://localhost:9000/spark/Parking_Violations.csv")
line1 = line1.mapPartitions(lambda x: reader(x))
violationcodes = line1.map(lambda x: (x[2],1)).reduceByKey(lambda x, y: x + y)
for xs in violationcodes.take(100):
    for x in xs:
        print(x)
```

Output

```
darsini@darsini-VirtualBox:~/hadoop/hadoop$ $SPARK_HOME/bin/spark-submit /home/
darsini/Desktop/spark.py
NE
1626
RI
13296
SD
691
NT
6
YT
14
Registration State
1
AL
5828
GV
1317
AB
243
```

This shows that the state RI (Rhode Island) has the maximum number of violation with about 13,296 violations

Query: Top 20 vehicles in terms of total violations

**MapReduce**

```python
import sys
from pyspark import SparkConf, SparkContext
from csv import reader
conf = SparkConf().setAppName("MR_2")
sc = SparkContext(conf=conf)
line1 = sc.textFile('/map/input/Parking_2016.csv')
line1 = line1.mapPartitions(lambda x: reader(x))
id = line1.map(lambda x: ((x[1],x[2]),1)).reduceByKey(lambda x, y: x + y).sortBy(lambda x: x[1], False)
top20 = sc.parallelize(id.take(20)).map(lambda x: (x[0][0], x[0][1], x[1]))
print(top20.collect())
```

**Output**

# HIVE

Loading two tables:

- Parking Ticket Details (from Kaggle)

```
CREATE EXTERNAL TABLE parking(
SummonsNo INT,
PlateID STRING,
RegistrationState STRING,
PlateType STRING,
IssueDate STRING,
ViolationCode INT,
VehicleBodyType STRING,
VehicleMake STRING,
IssuingAgency STRING,
VehicleExpDate INT,
ViolationLocation STRING,
ViolationPrecinct INT,
IssuerPrecinct INT,
IssuerCode INT,
IssuerCommand STRING,
ViolationTime INT,
StreetName STRING,
LawSection INT,
SubDivision STRING,
VehicleColor STRING,
VehicleYear INT )
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;

load data local inpath '/home/ng/parking2016.csv' into table parking;
select * from parking LIMIT 5;
```

- Details of Violation Codes

```
CREATE EXTERNAL TABLE vi_codes(
ViolationCode INT,
ViolationDef STRING,
FineM INT,
FineA INT )
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE;

load data local inpath '/home/ng/violation_codes.csv' into table vi_codes;
select * from vi_codes LIMIT 5;
```

- Combining two datasets on violation code

```
CREATE TABLE parkingfine AS
SELECT p.SummonsNo, p.PlateID, p.RegistrationState, p.PlateType, p.IssueDate, p.IssuingAgency,
p.ViolationPrecinct, p.ViolationCode, v.ViolationDef, v.FineM, v.FineA
FROM parking p
LEFT OUTER JOIN vi_codes v
ON (p.ViolationCode = v.ViolationCode);
```

- Exporting dataset to be used

```
hive -e 'select * from parkingfine' | sed 's/[\t]/,/g'  > /home/ng/file1.csv
```

## Queries:

**Violation code, the number of violations that have this code**

```
CREATE TABLE codes as
SELECT ViolationCode, count(*)
FROM parking
GROUP BY ViolationCode;
SELECT * FROM codes sort BY `_c1` DESC;
```

Output

```
code_c.violationcode    code_c._c1
21        1531575
36        1253512
38        1143684
14        875607
37        686607
20        611011
46        580517
7         492478
71        488920
40        462517
```

Violation Code 21 which states that *Street Cleaning: No **parking** where **parking** is not allowed by sign, street marking or traffic control device* (description according to the nyc government site) is the most violated code of all the rules.

**Which vehicle body type is most likely to get a parking ticket**

```
CREATE TABLE vehicletype as
SELECT VehicleBodyType, count(*)
FROM parking
GROUP BY VehicleBodyType;
SELECT * FROM vehicletype sort BY `_c1` DESC LIMIT 1;
```

```
vt.vehiclebodytype          vt._c1
SUBN      3466020
```

SUBN - Suburban Vehicles have the highest count of parking tickets

# HIVE vs MapReduce

## Chosen Query: Top 20 vehicles in terms of total violations

### MapReduce

```python
import sys
from pyspark import SparkConf, SparkContext
from csv import reader
conf = SparkConf().setAppName("MR_2")
sc = SparkContext(conf=conf)
line1 = sc.textFile('/map/input/Parking_2016.csv')
line1 = line1.mapPartitions(lambda x: reader(x))
id = line1.map(lambda x: ((x[1],x[2]),1)).reduceByKey(lambda x, y: x + y).sortBy(lambda x: x[1], False)
top20 = sc.parallelize(id.take(20)).map(lambda x: (x[0][0], x[0][1], x[1]))
print(top20.collect())
```

### Output



### Time Taken:

**HIVE**

```
CREATE TABLE codec AS
SELECT PlateID, RegistrationState, count(*)
FROM parkingfine
GROUP BY PlateID, RegistrationState;
CREATE TABLE top_v AS SELECT * FROM codec sort BY `_c2`
SELECT * FROM top_v LIMIT 20;
```

**Output**

```
v_rank.plateid   v_rank.registrationstate            v_rank._c2
56207MG NY       1288
12359MG NY       1084
AP501F   NJ      1030
85989MD NY       1010
14483JY NY       1008
AP300F   NJ      996
47603MD NY       984
12817KA NY       975
```

**Time Taken**

```
Time taken: 58.672 seconds
```