New York Parking Ticket Analytics

Data Processed with an AWS EMR Cluster By Roger Ren, Zhiyi Ren, Karsten Cao, Xinyue Wang



The NYC Department of Finance collects data on every parking ticket issued in NYC. This dataset contains all the information of issued parking tickets in New York City from Aug 2013 to June 2017 with over 20 columns including street name, issue date and violation code.

Data Size: 4.66GB

Format: CSV

Time Window: August 2013 - September 2015

Source: https://www.kaggle.com/new-york-city/nyc-parking-tickets

Data Analytics Goal

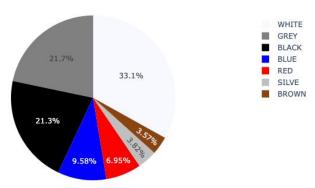
Post Hoc Analysis of Police Behavior

When, where, and which car the police are likely to issue tickets? What are the major reasons for the cars to be stickered?

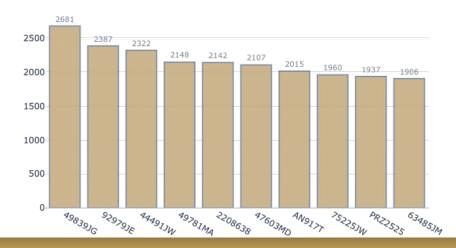
Outcomes

Top7 Car Colors that Get Tickets

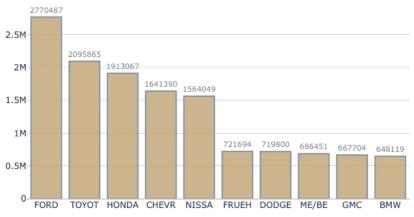
Which Car?



Top10 Car Plate ID that Get Tickets

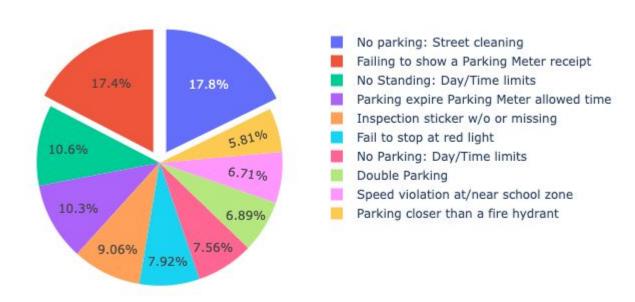


Top10 Brands that Get Tickets



Why?

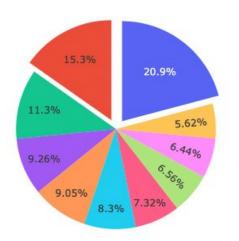
Top10 Violation Reason in NY County



Why?

Top10 Violation Reasons in NY County for June

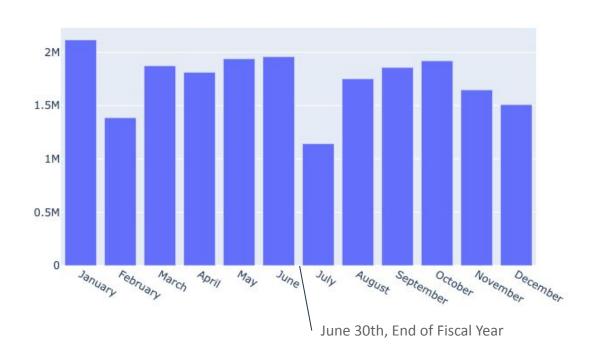
Top10 Violation Reasons in NY County for July





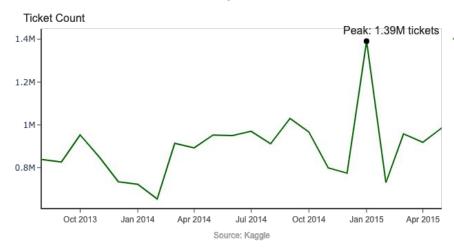
Time?

Count of Issued Parking Tickets by Month

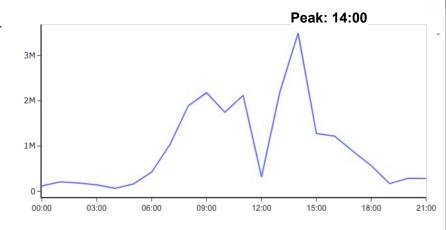


What Time?

The Number of Car Tickers Issued by Month



The Tickets Issued Timeline



Working Time

What Time?

Top20 Date that Get Tickets

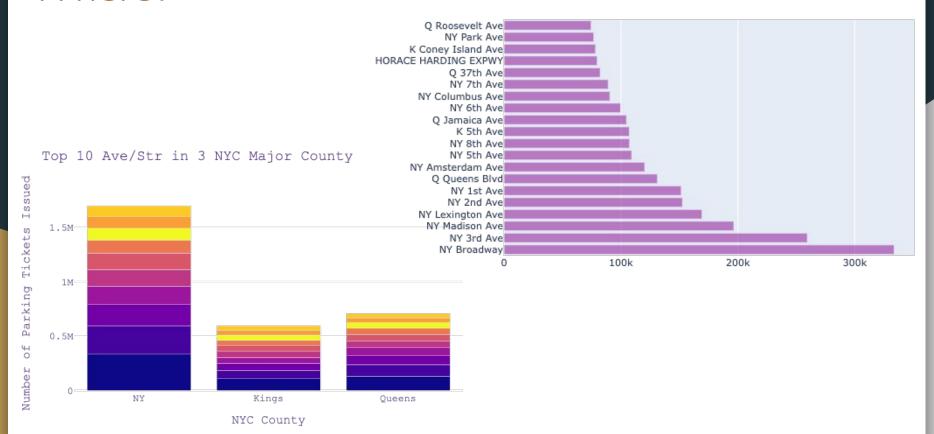




Subway Shutdown

Where?

Top20 Streets with Car Tickets



Where?



Performance Results

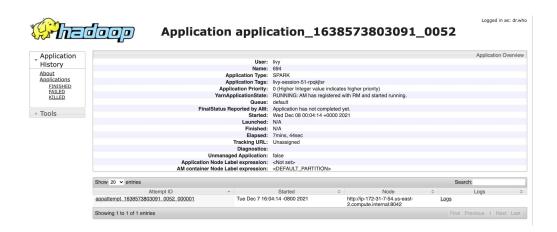
Cluster Settings:

m5.xlarge

4 vCore, 16 GiB memory, EBS only storage EBS Storage: 64 GiB

```
ticonfigure -f
{
    "conf": {
        "spark.pyspark.python": "python3",
        "spark.pyspark.virtualenv.enabled": "true",
        "spark.pyspark.virtualenv.type": "native",
        "spark.pyspark.virtualenv.bin.path": "/usr/bin/virtualenv",
        "spark.app.name":"694"
    }
}
```

Execution Time: 7 Min 44S



Efficiency

Cache the whole data's rdd

PythonRDD[12] at RDD at PythonRDD.scala:53

]: rdd.cache()

- For a single chunk of data operation
- Can imagine the actual power of how many minutes it saved for us

```
Without cache
: #without cache
  start = time.time(
  def trans(data):
         if 'A' in data:
             data = data.replace('A','')
             data = data[:2]+':'+ '00'
          elif 'P' in data:
             data = data.replace('P','')
             d = int(data[1:2])+12
             data = str(d)+':'+ data[2:]
             data = data[:2]+':'+'00'
         data - 'K
      return data
  Q = rdd.map(lambda x: (x[5],x[1])).map(lambda x: (trans(x[0]),x[1])).countByKey()
  Q = sorted(Q.items(), key=lambda x: x[0], reverse=False)
  0 rdd - sc.parallelize(0)
  Q_rdd_df = Q_rdd.filter(lambda x: x[1] > 21).filter(lambda x: x[0] != '').toDF(['Issue_Time','Ticket_Count'])
 Q_rdd_df.createOrReplaceTempView("Q_rdd_view")
  Q_rdd_df.show()
  print("second", time.time() - start)
                    159085
        05:00
       06.00
                    422918
       07:00
                   1030195
       08:00
                   1886260
                   2176025
                   1743116
                   2117534
                   316201
       12:00
                   2182580
       13:00
       14:00
                   3492323
       15:00
                   1276283
       16.00
                  1216548
       17:00
                   881265
       18 - 00
                   561913
       19:00
                   168937
  only showing top 20 rows
  second 66.63361525535583
```

```
]: #with cache
   start = time.time()
   def trans(data):
          if 'A' in data:
               data = data.replace('A','')
               data = data[:2]+':'+ '00'
           elif 'P' in data:
              data = data.replace('P','')
               d = int(data[1:2])+12
               data = str(d)+':'+ data[2:]
               data = data[:2]+':'+'00'
       except:
          data = 'K
       return data
   Q = rdd.map(lambda x: (x[5],x[1])).map(lambda x: (trans(x[0]),x[1])).countByKey()
  0 = sorted(0.items(),key=lambda x: x[0],reverse=False)
  Q_rdd = sc.parallelize(Q)
   Q_rdd_df = Q_rdd.filter(lambda x: x[1] > 21).filter(lambda x: x[0] != '').toDF(['Issue_Time', 'Ticket_Count'])
   Q rdd df.createOrReplaceTempView("Q rdd view")
  print("second", time.time() - start)
                     159085
         05:00
         06:00
                     422918
         07:00
                    1030195
         08:00
         09:00
                    2176025
         10.00
                   1743116
         11:00
                    2117534
         12:00
                    316201
         13:00
                    2182580
         14:00
                    3492323
         15:00
                    1276283
         16:00
                    1216548
         17:00
                    881265
                    561913
         18.00
         19:00
                    168937
   only showing top 20 rows
   second 16.256643295288086
```

Before:66 seconds

After: 16 seconds

Conclusion

Lesson Learned

- Large Dataset
 - a. Samples do not represent the whole
 - b. Data may not be as clean
- EMR Usage
- Coding Together Logistics

Thank you!