



UNIVERSIDAD
DE MÁLAGA

Master en Advanced Analytics on Big Data

Knowledge Area: Big Data

Task 6: Data search and visualization

Author: Juan Morales Conde

Juan Rafael Caro Moreno

Martín Blázquez Moreno

Module: Data Analytics

Málaga, 12th January 2020

Two data sets belonging to the New York public data have been selected.

It has been assumed that the analysis will be performed for an insurance company.

The data sets that have been chosen deal with vehicle collisions along with their characteristics.

The first data set called Motor Vehicle *Collisions – Crashes* provides a lot of information (345 MB formed by 1.63M of rows and 29 columns), the most relevant information for the analysis will be commented:

- CRASH DATE.
- BOROUGH.
- ZIPCODE.
- LOCATION.
- STREET INFORMATION.
- SUMARY PEOPLE INJURED.
- SUMARY PEOPLE KILLED.
- REASON OF THE COLLISION.
- VEHICLE TYPE.

Each row in the database represents the information of an accident.

Dataset URL:

<https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Crashes/h9gi-nx95>

The second set of data is called Motor Vehicle *Collisions – Vehicles* also provides a lot of information (521MB consisting of 3.27M of rows and 25 columns), the most important will be cited and explained:

- VEHICLE INFORMATION
- DRIVER SEX

- DRIVER INFORMATION
- VEHICLE DAMAGE
- REASON OF THE COLLISION
- PUBLIC PROPERTY DAMAGE

Each row in the database represents the vehicle information involved in an accident.

Dataset URL:

<https://data.cityofnewyork.us/Public-Safety/Motor-Vehicle-Collisions-Vehicles/bm4k-52h4>

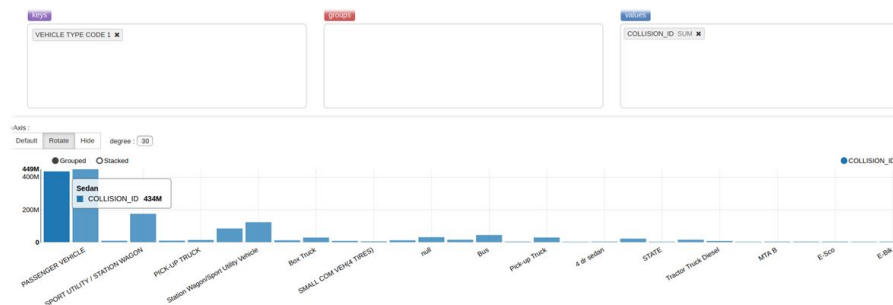
You can relate both sets of data using the *COLLISION_ID* column, since this column is an attribute that uniquely identifies the accident with the vehicle.

The analysis of these data sets could provide value to insurance companies, general traffic management, car brands ...

It is very interesting to carry out this analysis to improve the signaling of the streets, to ensure both pedestrians and cyclists or motorists since it provides a considerable volume of data, in addition to other sets also related to these.

A quick visualization of the *CRASHES* database is carried out.

The following image shows the number of accidents depending on the type of vehicle.



An aggregation function has been used to visualize the number of injuries per neighborhood and the average for cyclists:

```
%pyspark
BD_Vehicles.registerTempTable("VEHICULOS")

Took 0 sec. Last updated by anonymous at January 10 2020, 2:20:27 PM.

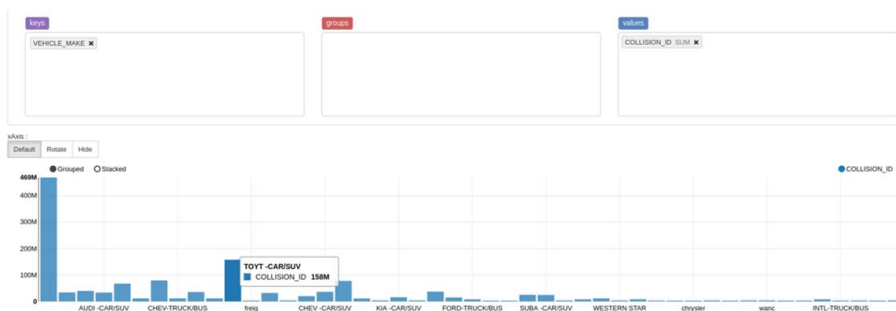
%sql
SELECT COUNT(*) as Number_of_Unlicensed_Driver from VEHICULOS WHERE DRIVER_LICENSE_STATUS = "Unlicensed"

[Table icon] [Bar chart icon] [Pie chart icon] [Line chart icon] [Scatter plot icon] [Download icon] [Settings icon]

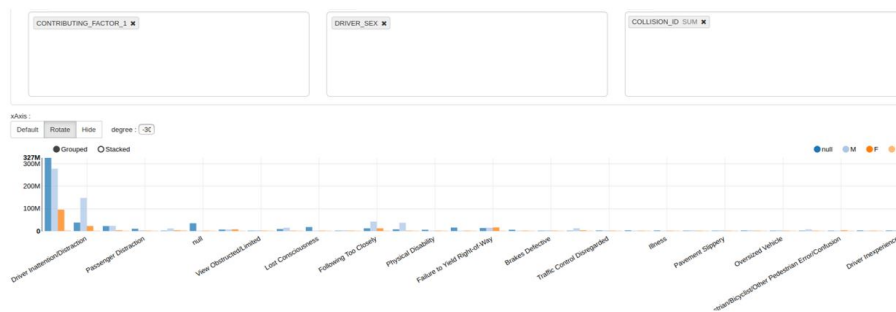
Number_of_Unlicensed_Driver

19563
```

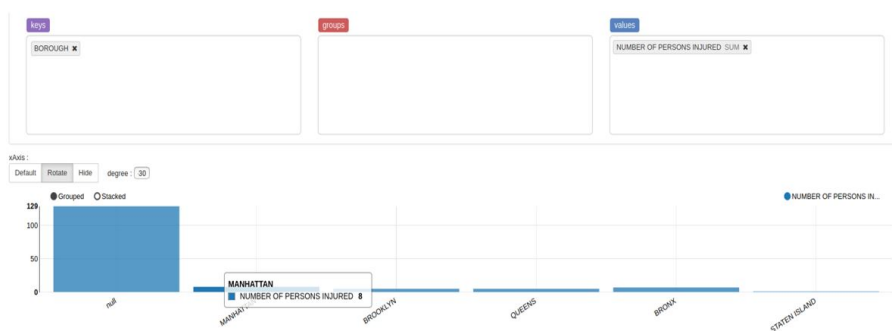
1. What are the brands of vehicles that have suffered the most accidents?



2. Does the sex of people influence the number and cause of accidents?



3. What is the neighborhood where more people have been injured in an accident?



4. How has the number of accidents per year evolved?

