# Motor Vehicle Collisions

Team 5 - Jiaxi Liu, Alice Pascalev, Abhimanyu Barun, Pranav Sai

# Contents

- Topic Introduction
- Data Description
  - Data Source
  - Data Cleaning
  - Dataset Information
- Exploratory Data Analysis
- SMART Questions
  1. How does other variables affect the number of motorist injured?
  2. Effect of time of the day and time of year on the number of vehicles involved
  3. Can we predict number of vehicles involved in the accident?
  4. Modeling Type of Car - 1 involved in a two car accident
  5. Modeling Type of Car - 2 involved in a two car accident
- Issues Faced

# Topic Introduction

# Topic Introduction

**More than 46,000** people die in car crashes each year, according to Annual United States Road Crash Statistics (ASIRT). Furthermore, the latest evidence has found that post-pandemic, reckless driving in the US has surged, with more cases of speeding, unbuckled seatbelts, and impaired driving[1]. Being struck by a vehicle is the leading cause of injury-related death of children in the United States, according to the CDC.[1]

Therefore, the topic of our project is Motor Vehicle Collisions, focusing on motor vehicle collisions, particularly those that require police reports.

1.   Baumgaertner, E. (2021, December 8). *Car crash deaths have surged during COVID-19 pandemic. Here's why*. Los Angeles Times. Retrieved December 13, 2022, from https://www.latimes.com/world-nation/story/2021-12-08/traffic-deaths-surged-during-covid-19-pandemic-heres-why
2.   https://www.cdc.gov/injury/features/child-injury/index.html

# Data Description

# Data Source

We are utilizing the data set **'Motor Vehicle Collisions - Crashes'**, which we accessed at Data.gov. This data is gathered from NYPD and includes data from all motor vehicle collisions in New York City that require a police report. This data set has 1,000,000 rows and 18 variables.

This data was made public in 2014, the year NYC's Vision Zero program was implemented, with the goal of safer streets and zero fatalities. The data is from 2014 until Nov 18, 2022.

*"Police report (MV104-AN) is required to be filled out for collisions where someone is injured or killed, or where there is at least $1000 worth of damage."* [1]

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

# Motor Vehicle Collision - Crashes

CRASH TIME
NUMBER OF PERSONS INJURED
NUMBER OF PERSONS KILLED
NUMBER OF PEDESTRIANS INJURED
NUMBER OF PEDESTRIANS KILLED
NUMBER OF CYCLIST INJURED
NUMBER OF CYCLIST KILLED
NUMBER OF MOTORIST INJURED
NUMBER OF MOTORIST KILLED

CONTRIBUTING FACTOR VEHICLE 1, CONTRIBUTING FACTOR VEHICLE 2, CONTRIBUTING FACTOR VEHICLE 3, CONTRIBUTING FACTOR VEHICLE 4, CONTRIBUTING FACTOR VEHICLE 5

VEHICLE TYPE CODE 1, VEHICLE TYPE CODE 2, VEHICLE TYPE CODE 3, VEHICLE TYPE CODE 4, VEHICLE TYPE CODE 5

# Data Cleaning

- Dropped NA values
- Convert Crash Date and Crash Time columns, to datetime format, split into Year, Month, Day and Hour, Minute columns
- Convert Zip Code from object -> integer
- Convert month number to name
- Create Vehicle Count column
- Due to seven-figure number of observations, only kept top 20 Vehicle Types for each Vehicle Type Code
  - Consolidated categories with misspellings, extra words, case differences (ie. SEDAN, Sedan, 3-door sedan )

# Dataset info

**Time Indicator:**

Year,  Month, Day,  Hour, Minute

**Location Indicator:**

Borough, Zip code, Latitude, Longitude

**# of injured/killed:**

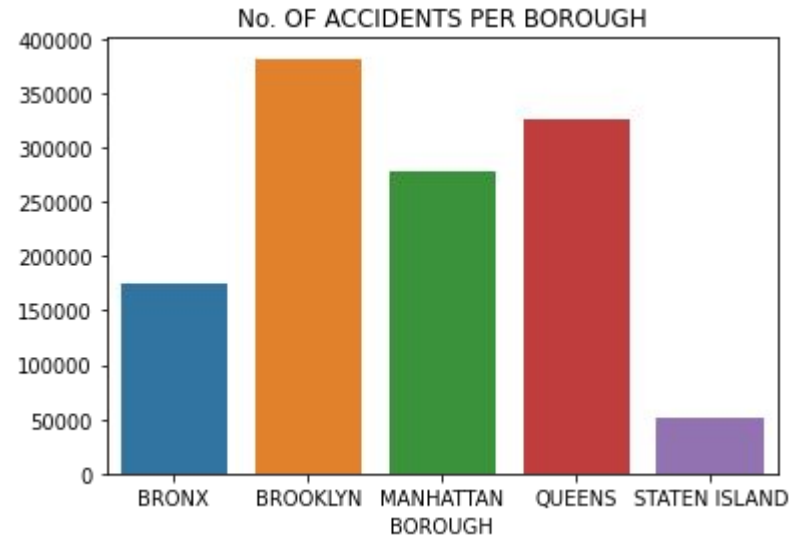Persons, Pedestrians, Cyclist, Motorist

**Vehicle info:**

Contributing factor vehicle 1-5, Vehicle type code 1-5, Vehicle count
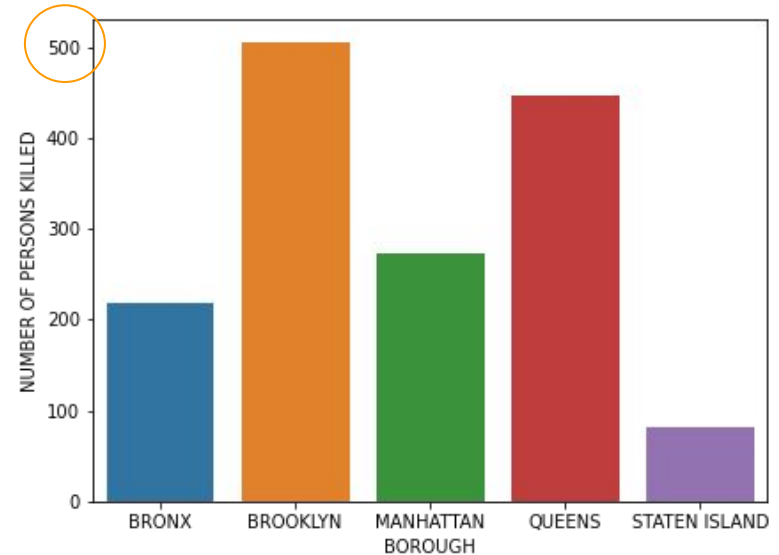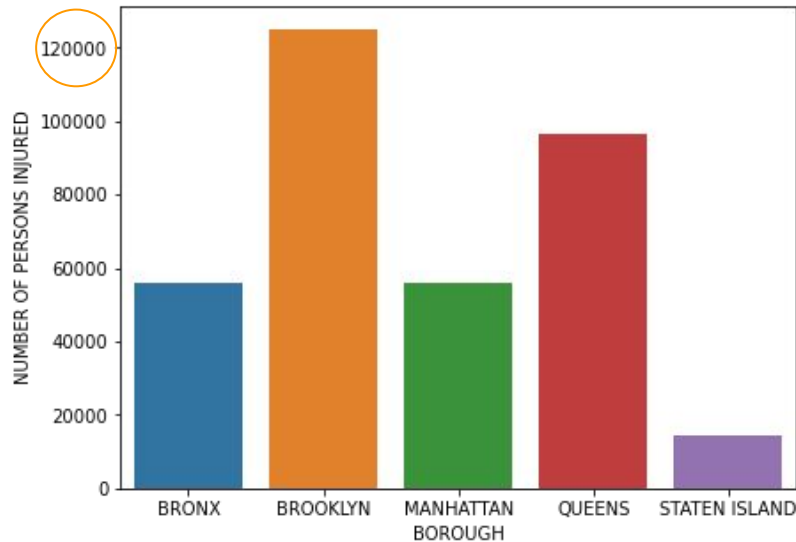
# Exploratory Data Analysis

# Accidents per Borough

- Road accidents are more frequent in Brooklyn, Queens and Manhattan.

- Staten Island has the lowest number of accidents
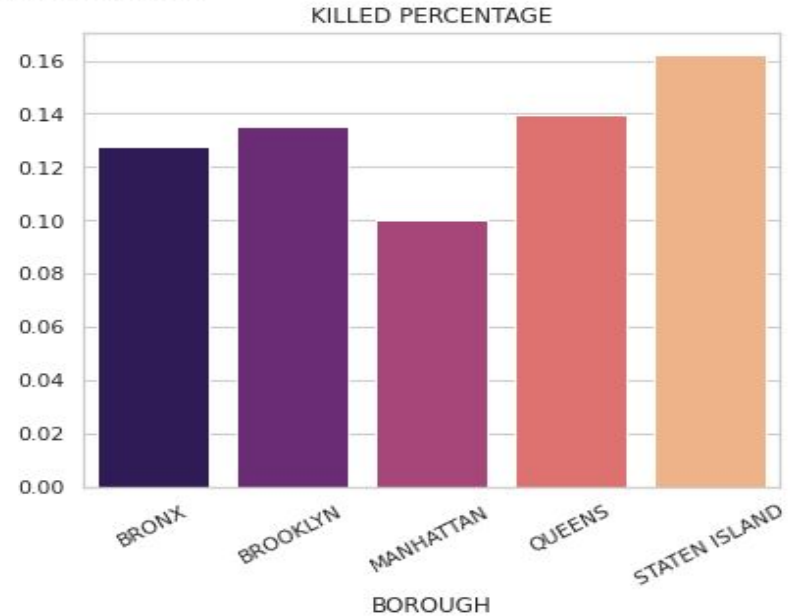
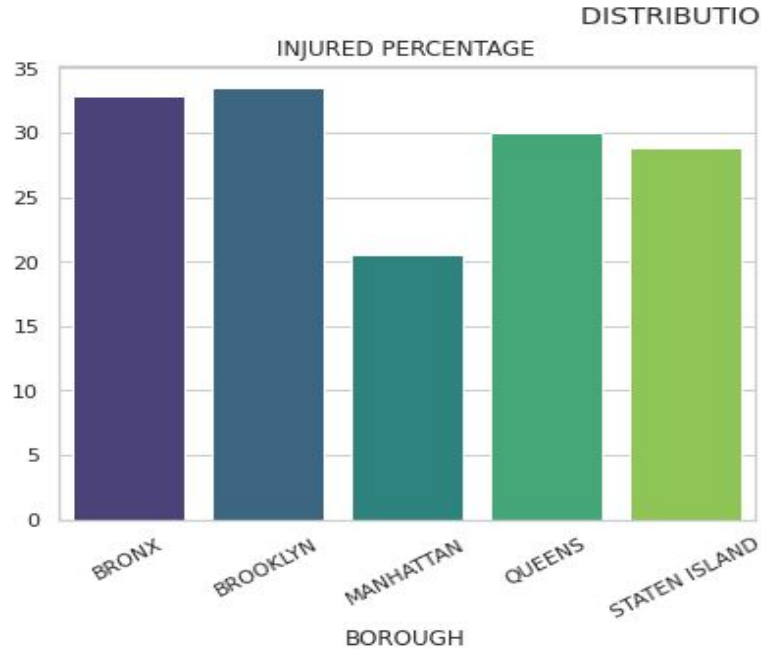- Brooklyn has the highest number of accidents

# Distribution of People Injured and Killed in NYC
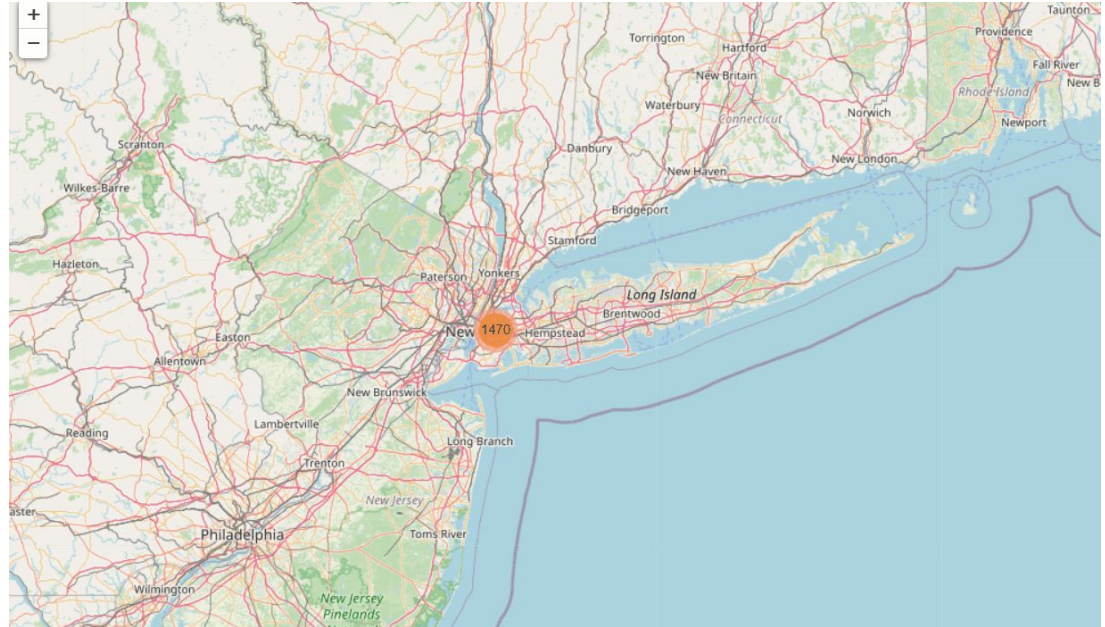


Injured vs Killed in NYC

# Distribution by Percentage



DISTRIBUTION BY PERCENTAGE
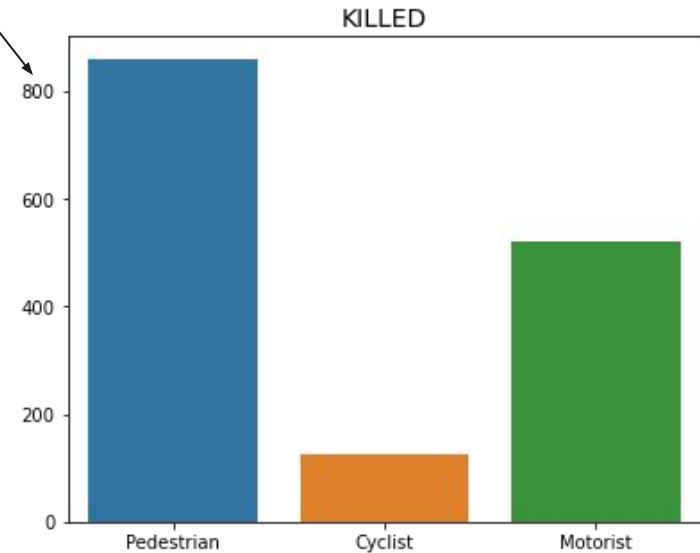
# Impact of location on number of fatal accidents
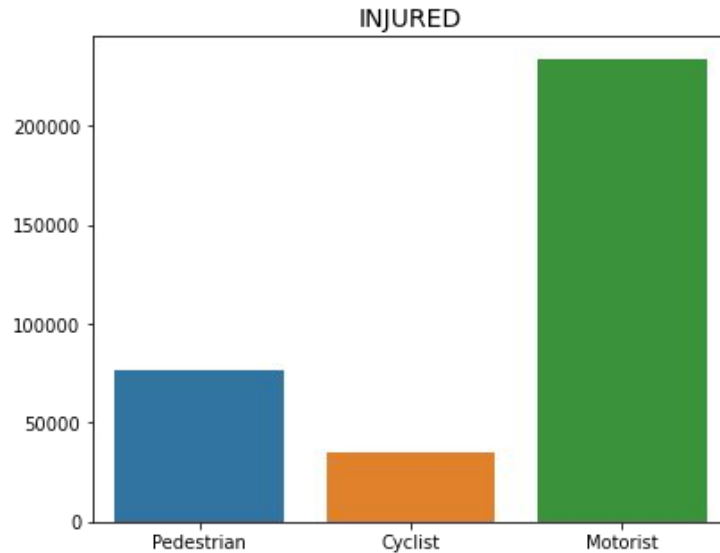
Location wise lethal accidentes

# Killed and Injured per Accident



KILLED vs INJURED PER ACCIDENT

# Vehicles Involved

We see that the vast majority of collisions involved two vehicles



NUMBER OF VEHICLES INVOLVED

# Impact of time of day on location of fatal crashes

Location and time wise lethal accidents

# Accidents and Deaths based on Month



Highest Frequency
(Accidents)

1. **July** 112,401
2. **October** 111,716

Lowest Frequency
(Accidents)

1. **April** 81,082
2. **February** 87,912

# deaths increase by 88%
from month of march to
July

*Please note: Does not include December 2022 statistics*

# Accidents and Deaths based on Hours

Most accidents are observed during the day

And the fatalities are highest at the end of the day from 6PM to 8PM in the evening



Accidents & deaths based on Hour

Rush Hour

# Contributing Factors of Vehicle 1

Over 20000 accidents were due to Driver Inattention or some other distraction.

Failure to Yield Right-of-Way and Driving Unsafely are the two next most common contributing factors to an accident.

Similar amount of Accidents were caused due to errors of another vehicle on road and Following too closely.

Other common contributing factors include - Turning Improperly, Improper lane usage, fatigue and passing too closely.



CONTRIBUTING FACTORS - VEHICLE 1

# Contributing Factors of Vehicle 2

Most accidents were due to Driver Inattention or some other distraction.

Errors of another vehicle on road and Failure to Yield Right-of-Way are the two next most common contributing factors to an accident.

Similar amount of Accidents were caused Improper Lane usage, following too closely, Fatigue, backing unsafely and Turning improperly.



CONTRIBUTING FACTORS - VEHICLE 2

# Data Visualization

# Top 10 Vehicle types involved in accidents



Top 10 vehicle types causing accidents

**Passenger Vehicle** - A motor vehicle designed to carry **10 persons or less** which is constructed either on a truck chassis or with special features for occasional off-road operation.
*15 CRR-NY 55.1*

# Top 10 Vehicle type involvement leading to Death

# Number of Collisions per year

Number of Collisions per Year

# Number of Collisions per month



Number of Collisions per Month

# Number of killed and injured per accident type

# Data Visualization

# SMART QUESTIONS

# 1. How does other variables affect the number of motorist injured?

```
                          OLS Regression Results
==============================================================================
Dep. Variable:     NUMBER OF MOTORIST INJURED   R-squared (uncentered):              0.805
Model:                                    OLS   Adj. R-squared (uncentered):         0.805
Method:                         Least Squares   F-statistic:                     2.212e+06
Date:                        Tue, 13 Dec 2022   Prob (F-statistic):                   0.00
Time:                                17:55:10   Log-Likelihood:                 -1.7912e+05
No. Observations:                     1068904   AIC:                             3.582e+05
Df Residuals:                         1068902   BIC:                             3.583e+05
Df Model:                                   2
Covariance Type:                    nonrobust
==============================================================================
                             coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
NUMBER OF PERSONS INJURED   0.8067      0.000   1959.905      0.000       0.806       0.807
VEHICLE COUNT              -0.0053      0.000    -35.690      0.000      -0.006      -0.005
==============================================================================
Omnibus:                   728443.358   Durbin-Watson:                   1.979
Prob(Omnibus):                  0.000   Jarque-Bera (JB):         98378132.079
Skew:                          -2.385   Prob(JB):                         0.00
Kurtosis:                      49.756   Cond. No.                         3.04
==============================================================================
```
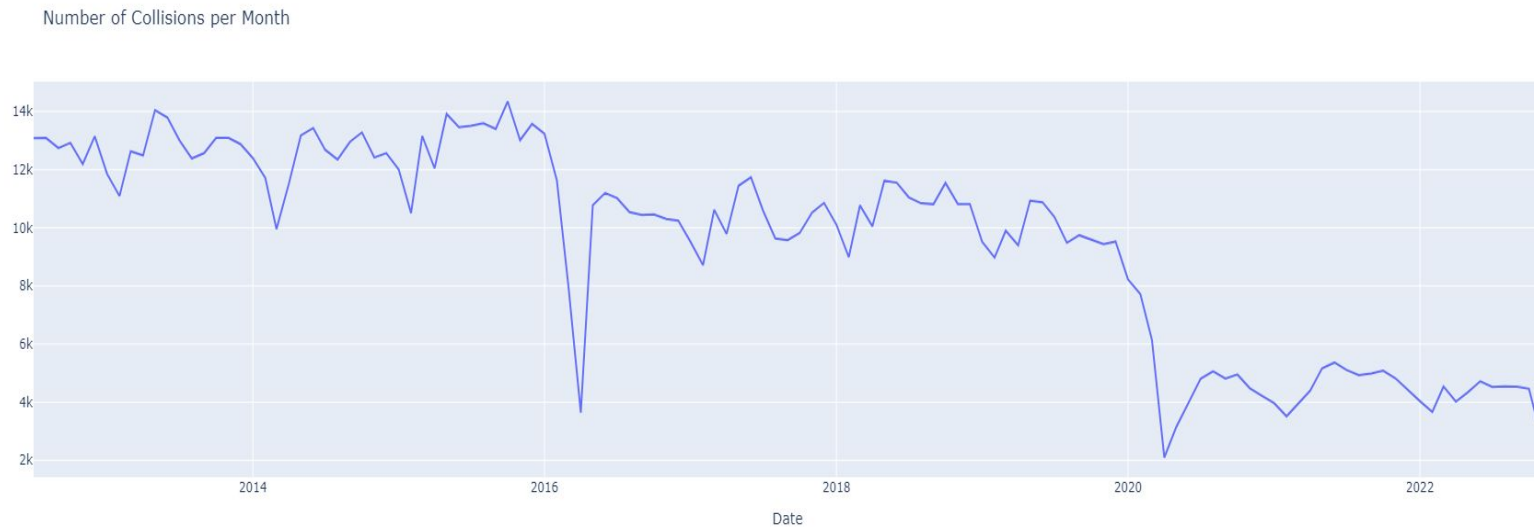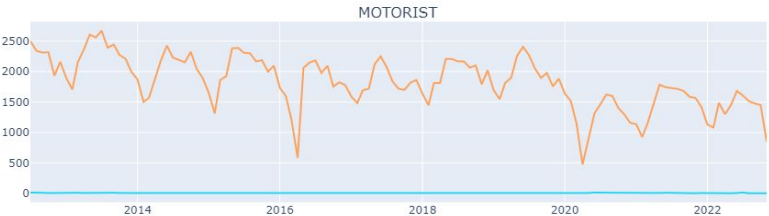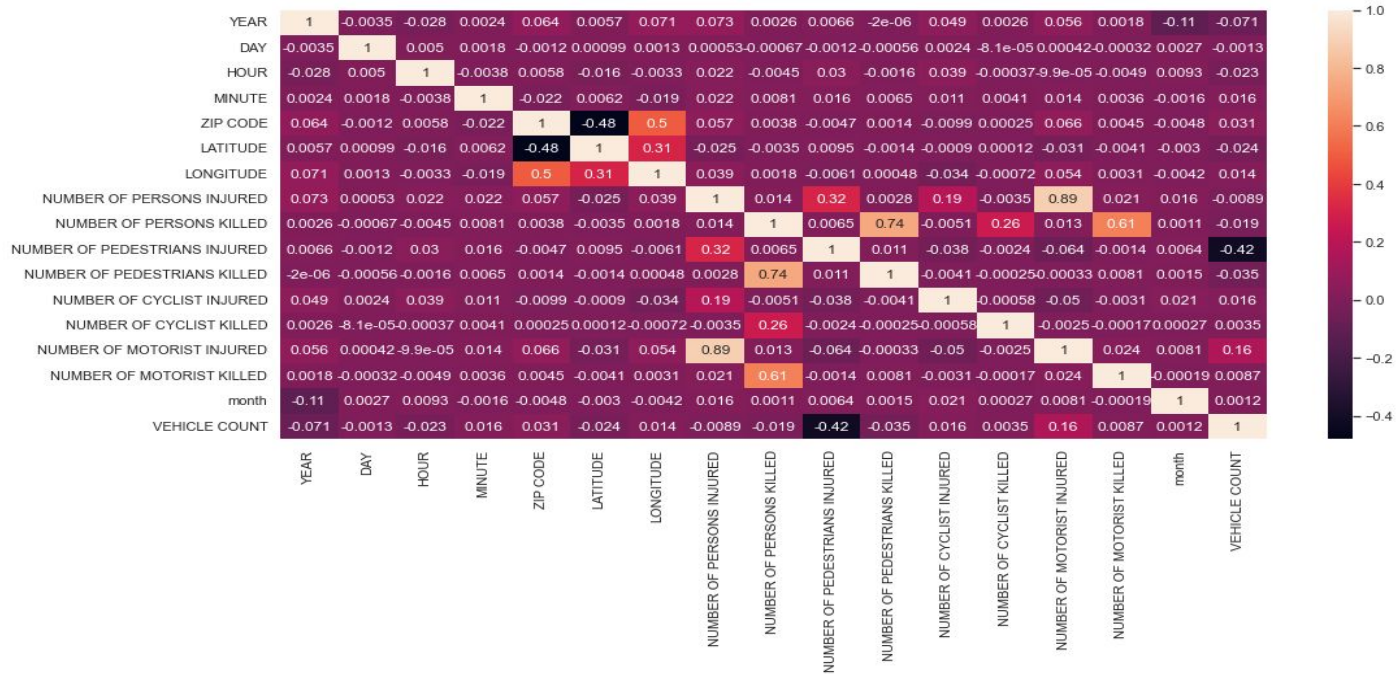
## 2. How does the time of the day and time of year affect the number of vehicles involved?

```
                              OLS Regression Results
==============================================================================
Dep. Variable:          VEHICLE COUNT   R-squared (uncentered):              0.853
Model:                            OLS   Adj. R-squared (uncentered):         0.853
Method:                 Least Squares   F-statistic:                     1.555e+06
Date:                Tue, 13 Dec 2022   Prob (F-statistic):                   0.00
Time:                        17:57:36   Log-Likelihood:                -1.2311e+06
No. Observations:             1068904   AIC:                             2.462e+06
Df Residuals:                 1068900   BIC:                             2.462e+06
Df Model:                           4
Covariance Type:            nonrobust
==============================================================================
                                  coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
month                           0.1069      0.000    615.719      0.000       0.107       0.107
HOUR                            0.0771    9.1e-05    847.793      0.000       0.077       0.077
NUMBER OF PEDESTRIANS INJURED  -0.6828      0.003   -238.208      0.000      -0.688      -0.677
NUMBER OF MOTORIST INJURED      0.2042      0.001    171.231      0.000       0.202       0.207
==============================================================================
Omnibus:                   117351.019   Durbin-Watson:                       1.881
Prob(Omnibus):                  0.000   Jarque-Bera (JB):               327253.178
Skew:                           0.610   Prob(JB):                             0.00
Kurtosis:                       5.421   Cond. No.                             61.4
==============================================================================
```

# 3. Can we predict number of vehicles involved in the accident?



|     | precision | recall | f1-score | support |
|-----|-----------|--------|----------|---------|
| 1.0 | 0.84 | 0.44 | 0.58 | 43523 |
| 2.0 | 0.84 | 0.98 | 0.91 | 229298 |
| 3.0 | 0.16 | 0.01 | 0.01 | 13707 |
| 4.0 | 0.09 | 0.00 | 0.00 | 2914 |
| 5.0 | 0.18 | 0.00 | 0.01 | 1186 |
| | | | | |
| accuracy | | | 0.84 | 290628 |
| macro avg | 0.42 | 0.29 | 0.30 | 290628 |
| weighted avg | 0.80 | 0.84 | 0.80 | 290628 |

# 4. What type of Car - 1 was involved in a two car accident?

```
CLASSIFICATION REPORT:
              precision    recall  f1-score   support

          1       0.60      0.51      0.55       206
          4       0.00      0.00      0.00       192
          5       0.00      0.00      0.00       278
         15       0.00      0.00      0.00       149
         16       0.50      0.85      0.63      2736
         17       0.00      0.00      0.00       324
         19       0.49      0.70      0.58      3317
         20       0.00      0.00      0.00       116
         21       0.31      0.09      0.14      1218
         22       0.40      0.28      0.33      2425
         23       0.30      0.21      0.25      1004
         25       0.00      0.00      0.00        84
         26       0.00      0.00      0.00       294

   accuracy                           0.47     12343
  macro avg       0.20      0.20      0.19     12343
weighted avg       0.39      0.47      0.40     12343
```

# 5. What type of Car - 2 was involved in a two car accident?

```
CLASSIFICATION REPORT:
              precision    recall  f1-score   support

           0       0.51      0.54      0.53       179
           2       0.55      0.10      0.16       221
           3       0.44      0.21      0.28       395
           4       0.00      0.00      0.00       200
           5       0.00      0.00      0.00       286
          16       0.00      0.00      0.00       184
          17       0.46      0.84      0.60      2517
          18       0.00      0.00      0.00       318
          20       0.46      0.69      0.55      3062
          22       0.31      0.09      0.14      1125
          23       0.37      0.26      0.31      2296
          24       0.32      0.23      0.27      1009
          26       0.00      0.00      0.00       222
          27       0.00      0.00      0.00       329

    accuracy                           0.44     12343
   macro avg       0.24      0.21      0.20     12343
weighted avg       0.36      0.44      0.37     12343
```

# Issues Faced with the Dataset

- The **large size** of the dataset **restricted** the type of models we could use, or how we used them:
  - **KNN -** long run time, could not get highly accurate results
  - **Random Forest -** had to adjust n_jobs to -1, and n_estimators to 100, in order to get an output in a more efficient, timely manner
    - **n_jobs=-1 :** use all available CPUs to handle the 1,000,000+ observations
    - **n_estimators=100:** large number of trees for large number of variables, observations
- The data in the survey is added manually i.e. there is not a fixed format.
  - This forced us to remove many rows from the dataset.
  - Same type of vehicles were labelled differently in many rows for ex. Ambulance was mentioned as AMBULANCE, Ambulance and AMBU or Taxi was mentioned as TAXI and Taxi.
  - Unclear categories with only 1 occurrence (ie. 994 , UHUAL, TCN )