# SSC 442 - Final Project

Peter Fu Chen , Mike Liu

2023-04-09

## NYC Traffic Accidents Analysis

### Background and Motivation

The focus of this project is the examination of motor vehicle collisions that were reported by the New York City Police Department from January to August of 2020. Each record in the dataset, represents a unique collision and includes various details such as the date, time, and location of the accident, as well as additional data.

The analysis will compare the percentage of total accidents by month, providing a snapshot of the trend over time. These statistics will be visually represented through the use of maps, which will demonstrate the frequency of accidents in different boroughs of New York City. Furthermore, the analysis will also determine the most common streets, days, and times when accidents are likely to occur.

The ultimate goal of this project is to provide recommendations to the city of New York on how this analysis can be used to prevent future accidents. By highlighting the most common accident hotspots, times, and contributing factors, city planners and law enforcement officials will be equipped with the information necessary to implement effective safety measures and reduce the number of accidents in the city.

# Methodology

*Loading libraries*

## Data preparation

```
## Rows: 74,881
## Columns: 29
## $ CRASH.DATE                   <chr> "8/29/20", "8/29/20", "8/29/20", "8/29/2~
## $ CRASH.TIME                   <chr> "15:40:00", "21:00:00", "18:20:00", "0:0~
## $ BOROUGH                      <chr> "BRONX", "BROOKLYN", "", "BRONX", "BROOK~
## $ ZIP.CODE                     <int> 10466, 11221, NA, 10459, 11203, NA, 1045~
## $ LATITUDE                     <dbl> 40.89210, 40.69050, 40.81650, 40.82472, ~
## $ LONGITUDE                    <dbl> -73.83376, -73.91991, -73.94656, -73.892~
## $ LOCATION                     <chr> "POINT (-73.83376 40.8921)", "POINT (-73~
## $ ON.STREET.NAME               <chr> "PRATT AVENUE", "BUSHWICK AVENUE", "8 AV~
## $ CROSS.STREET.NAME            <chr> "STRANG AVENUE", "PALMETTO STREET", "", ~
## $ OFF.STREET.NAME              <chr> "", "", "", "1047 SIMPSON STREET", "4609~
## $ NUMBER.OF.PERSONS.INJURED    <int> 0, 2, 1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 1, 0~
## $ NUMBER.OF.PERSONS.KILLED     <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.PEDESTRIANS.INJURED <int> 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0~
## $ NUMBER.OF.PEDESTRIANS.KILLED <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.CYCLIST.INJURED    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.CYCLIST.KILLED     <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.MOTORIST.INJURED   <int> 0, 2, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.MOTORIST.KILLED    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ CONTRIBUTING.FACTOR.VEHICLE.1 <chr> "Passing Too Closely", "Reaction to Unin~
## $ CONTRIBUTING.FACTOR.VEHICLE.2 <chr> "Unspecified", "Unspecified", "", "Unspe~
## $ CONTRIBUTING.FACTOR.VEHICLE.3 <chr> "", "", "", "Unspecified", "", "", "", "~
## $ CONTRIBUTING.FACTOR.VEHICLE.4 <chr> "", "", "", "Unspecified", "", "", "", "~
## $ CONTRIBUTING.FACTOR.VEHICLE.5 <chr> "", "", "", "", "", "", "", "", "", "", ~
## $ COLLISION_ID                 <int> 4342908, 4343555, 4343142, 4343588, 4342~
## $ VEHICLE.TYPE.CODE.1          <chr> "Sedan", "Sedan", "Station Wagon/Sport U~
## $ VEHICLE.TYPE.CODE.2          <chr> "Station Wagon/Sport Utility Vehicle", "~
## $ VEHICLE.TYPE.CODE.3          <chr> "", "", "", "Sedan", "", "", "", "", "Se~
## $ VEHICLE.TYPE.CODE.4          <chr> "", "", "", "Motorcycle", "", "", "", ""~
## $ VEHICLE.TYPE.CODE.5          <chr> "", "", "", "", "", "", "", "", "", "", ~
```

```
## [1] 37643
```

```
## [1] 0
```

```
## Rows: 47,686
## Columns: 18
## $ CRASH.DATE        <date> 2020-08-29, 2020-08-29, 2020-08-29, 202~
## $ CRASH.TIME        <chr> "15:40:00", "21:00:00", "0:00:00", "17:1~
## $ BOROUGH           <chr> "BRONX", "BROOKLYN", "BRONX", "BROOKLYN"~
## $ ZIP.CODE          <int> 10466, 11221, 10459, 11203, 10459, 10466~
## $ LATITUDE          <dbl> 40.89210, 40.69050, 40.82472, 40.64989, ~
## $ LONGITUDE         <dbl> -73.83376, -73.91991, -73.89296, -73.933~
## $ LOCATION          <chr> "POINT (-73.83376 40.8921)", "POINT (-73~
## $ ON.STREET.NAME    <chr> "PRATT AVENUE", "BUSHWICK AVENUE", "", "~
## $ CROSS.STREET.NAME <chr> "STRANG AVENUE", "PALMETTO STREET", "", ~
```

```
## $ OFF.STREET.NAME              <chr> "", "", "1047 SIMPSON STREET", "4609 SNY~
## $ NUMBER.OF.PERSONS.INJURED    <int> 0, 2, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 2, 2~
## $ NUMBER.OF.PERSONS.KILLED     <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.PEDESTRIANS.INJURED <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0~
## $ NUMBER.OF.PEDESTRIANS.KILLED <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.CYCLIST.INJURED    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.CYCLIST.KILLED     <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
## $ NUMBER.OF.MOTORIST.INJURED   <int> 0, 2, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 2, 2~
## $ NUMBER.OF.MOTORIST.KILLED    <int> 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0~
```
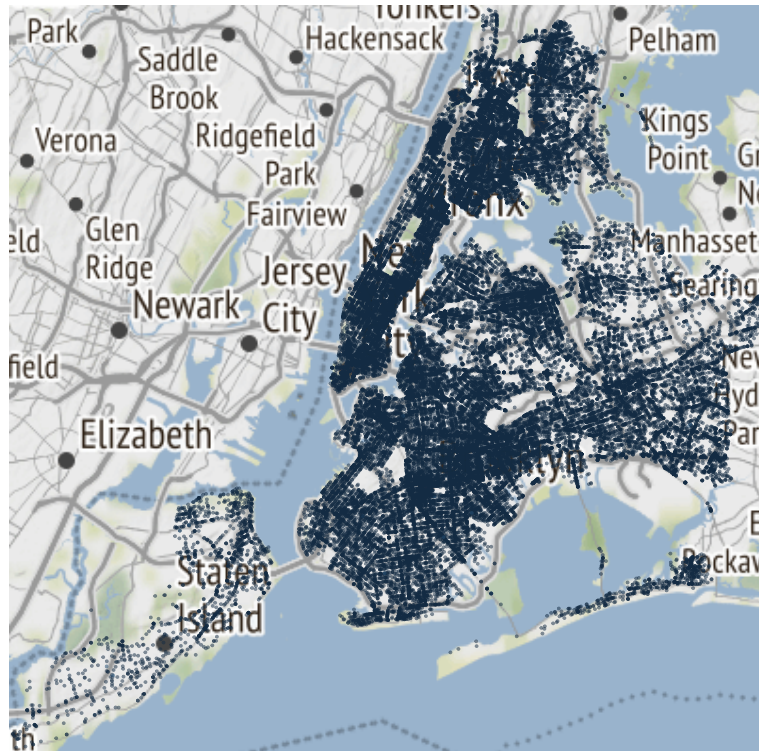
**Map plotting the people killed by motor vehicle in NYC**

First and foremost, let's examine the number of fatalities resulting from motor vehicle accidents across New York City in 2020. To represent each death, we have used dots on a map and it's evident that the concentration of these dots varies among different boroughs. Upon initial inspection, it appears that Manhattan, the Bronx, and Brooklyn have the highest density of dots, indicating a higher rate of fatalities in these areas.

Next, we will employ various analytical techniques to determine which borough is experiencing the most severe problem and to understand how public policy can be utilized to address this issue. By analyzing the data and identifying trends, we can work towards developing effective strategies that can help reduce the number of motor vehicle-related fatalities in New York City.

```
## i Map tiles by Stamen Design, under CC BY 3.0. Data by OpenStreetMap, under ODbL.
```

NYC Traffic Accidents Killed

**Regression analysis: Model the relationship between the number of accidents and borough.**

```
##
## t test of coefficients:
##
##                                 Estimate Std. Error t value  Pr(>|t|)
## (Intercept)                    0.3412076  0.0072879 46.8184 < 2.2e-16 ***
## as.factor(BOROUGH)BROOKLYN     0.0142173  0.0091336  1.5566 0.1195722
## as.factor(BOROUGH)MANHATTAN   -0.0470403  0.0100631 -4.6745 2.955e-06 ***
## as.factor(BOROUGH)QUEENS      -0.0131560  0.0092557 -1.4214 0.1552080
## as.factor(BOROUGH)STATEN ISLAND  0.0882888  0.0233044  3.7885 0.0001517 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

To better understand the relationship between the borough where a traffic accident occurs and the number of persons injured, we fit a linear regression model using the lm() function in R. In this model, we predict the number of persons injured based on the borough in which the accident occurred, treating borough as a categorical variable with five levels (Bronx, Brooklyn, Manhattan, Queens, and Staten Island).

In this model, the Bronx is the base level, or omitted level, against which the other boroughs are compared. The coefficients of the model represent the differences between the expected number of persons injured in each of the other boroughs and the expected number in the Bronx, holding all other variables constant.

The results of the model are shown in the table produced by the coeftest() function. The intercept represents the expected number of persons injured when the accident occurs in the Bronx. The estimate of the intercept is 0.3412076, which means that on average, the number of persons injured in a traffic accident in the Bronx is 0.3412076.

The next four coefficients are the differences between the expected number of persons injured in each of the other boroughs and the expected number in the Bronx. For example, the coefficient for Brooklyn (0.0142173) means that, on average, the number of persons injured in a traffic accident in Brooklyn is 0.0142173 higher than the number in the Bronx, holding all other variables constant (ceteris paribus).

The t-values and p-values associated with each coefficient indicate whether each coefficient is statistically significant. In this model, the intercept and the coefficients for Manhattan and Staten Island are statistically significant, with p-values less than 0.05. The coefficients for Brooklyn and Queens are not statistically significant, with p-values greater than 0.05. This suggests that there is not enough evidence to conclude that the number of persons injured in a traffic accident differs between the Bronx and Queens or Brooklyn, after adjusting for other variables. However, there is evidence to suggest that the number of persons injured in a traffic accident is lower in Manhattan and higher in Staten Island compared to the Bronx, after adjusting for other variables.

These results have important implications for policymakers and law enforcement officials in New York City. By identifying the boroughs with the highest rates of persons injured in traffic accidents, city planners can develop targeted interventions and safety measures to help reduce the number of accidents and improve public safety.

**Do we need heteroskedasticity-consistent errors?**

We might want to use Breusch-Pagan test to see if we need heteroskedasticity for this model.

```
##
## Call:
## lm(formula = NUMBER.OF.PERSONS.INJURED ~ as.factor(BOROUGH),
##     data = df_withoutna)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.4295 -0.3554 -0.3281  0.6446 14.6446
##
## Coefficients:
##                                 Estimate Std. Error t value Pr(>|t|)
## (Intercept)                     0.341208   0.007143  47.770  < 2e-16 ***
## as.factor(BOROUGH)BROOKLYN      0.014217   0.008889   1.599    0.110
## as.factor(BOROUGH)MANHATTAN    -0.047040   0.010777  -4.365 1.27e-05 ***
## as.factor(BOROUGH)QUEENS       -0.013156   0.009215  -1.428    0.153
## as.factor(BOROUGH)STATEN ISLAND 0.088289   0.019584   4.508 6.55e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6798 on 47681 degrees of freedom
## Multiple R-squared:  0.001438,   Adjusted R-squared:  0.001354
## F-statistic: 17.16 on 4 and 47681 DF,  p-value: 4.481e-14
```

Heteroscedasticity, or non-constant variance of the errors in a regression model, can be a problem in statistical analysis. When the errors have different variances across the range of the independent variable, it can lead to biased estimates of the standard errors and t-statistics, which can cause incorrect inferences about the significance of the coefficients and may lead to incorrect conclusions about the relationship between the independent and dependent variables.
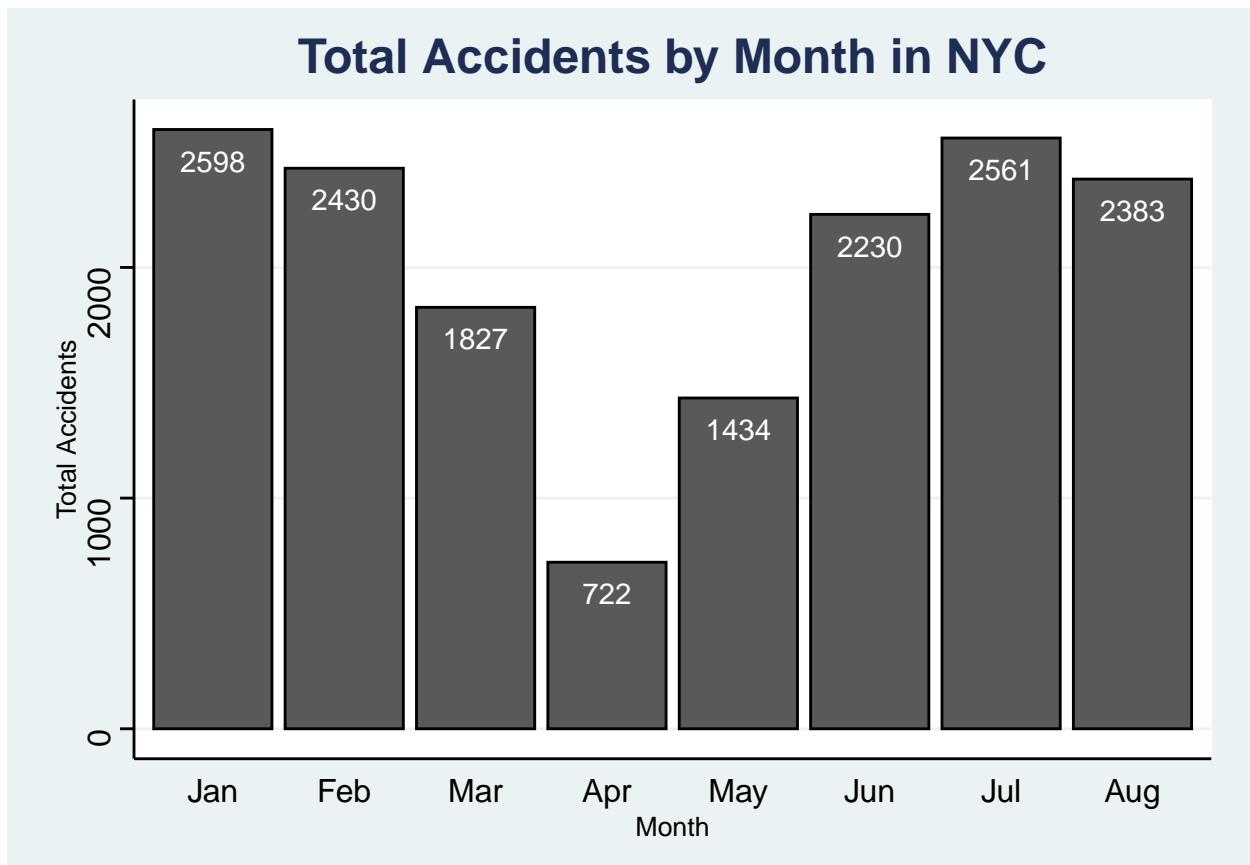
To determine if heteroscedasticity is a problem in our regression model, we can use the Breusch-Pagan test. This test checks whether the variance of the residuals in the model is constant across all levels of the independent variable.

In this case, we estimate a linear regression model with `NUMBER.OF.PERSONS.INJURED` as the dependent variable and `BOROUGH` as the independent variable, using the `lm()` function in R. We then calculate the residuals and their squared values, which we use as the dependent variable in a second regression model. The second model tests whether the squared residuals are related to the independent variable, which would indicate heteroscedasticity.

The results of the Breusch-Pagan test show a significant F-statistic of 17.16, with 4 and 47681 degrees of freedom, and a very small p-value of 4.481e-14. This provides strong evidence against the null hypothesis of homoscedasticity and suggests that there is heteroscedasticity in the model.

To correct for heteroscedasticity, we can use a robust standard error estimator, such as the HC1 estimator, which we used in our regression model. This estimator provides more accurate standard errors and t-statistics even in the presence of heteroscedasticity, which can help ensure that the results of the regression model are accurate and reliable, and can help prevent incorrect conclusions about the relationship between the independent and dependent variables.
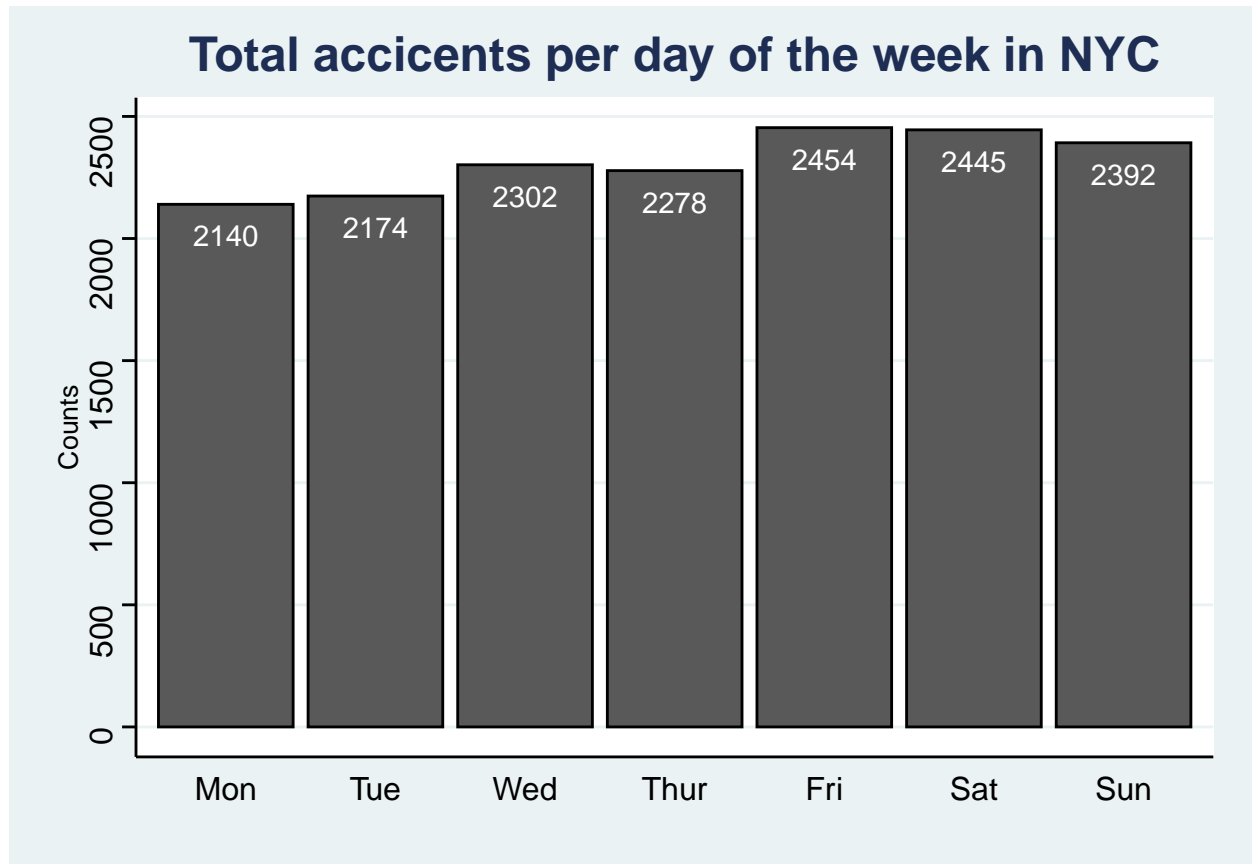
**Accidents Per Month**



## Total Accidents by Month in NYC

The bar chart shows that April had the lowest number of accidents per month in 2020, with a total of 11,234 accidents. This represents a significant decrease from the previous month, March, which had 15,315 accidents. There are several factors that could have contributed to this decrease:

- COVID-19 lockdown: One of the most significant factors that could have contributed to a decrease in the number of accidents is the COVID-19 lockdown in New York City, which began in March 2020. The lockdown resulted in many people staying at home, which would have led to a decrease in the number of cars on the road and a lower probability of accidents.

- Weather: Another possible factor that could have contributed to a decrease in the number of accidents in April is the weather. April is typically a month with moderate temperatures in New York City, and the weather conditions may have been less conducive to accidents.

- Other factors: There could be several other factors that contributed to the decrease in the number of accidents in April. For example, there could have been changes in traffic patterns, road closures, or road maintenance activities that reduced the number of accidents.

It's important to note that while April had the lowest number of accidents per month in 2020, this may not necessarily reflect a long-term trend. There may be other factors at play that contributed to the decrease in April, and it's possible that accident rates could return to previous levels in the future. Additionally, the analysis focuses on the total number of accidents per month, but it may be informative to look at the types of accidents that occurred and the factors that contributed to them. For example, it's possible that certain types of accidents, such as those involving pedestrians or cyclists, were more or less common during the lockdown period.

**Accidents Per day**
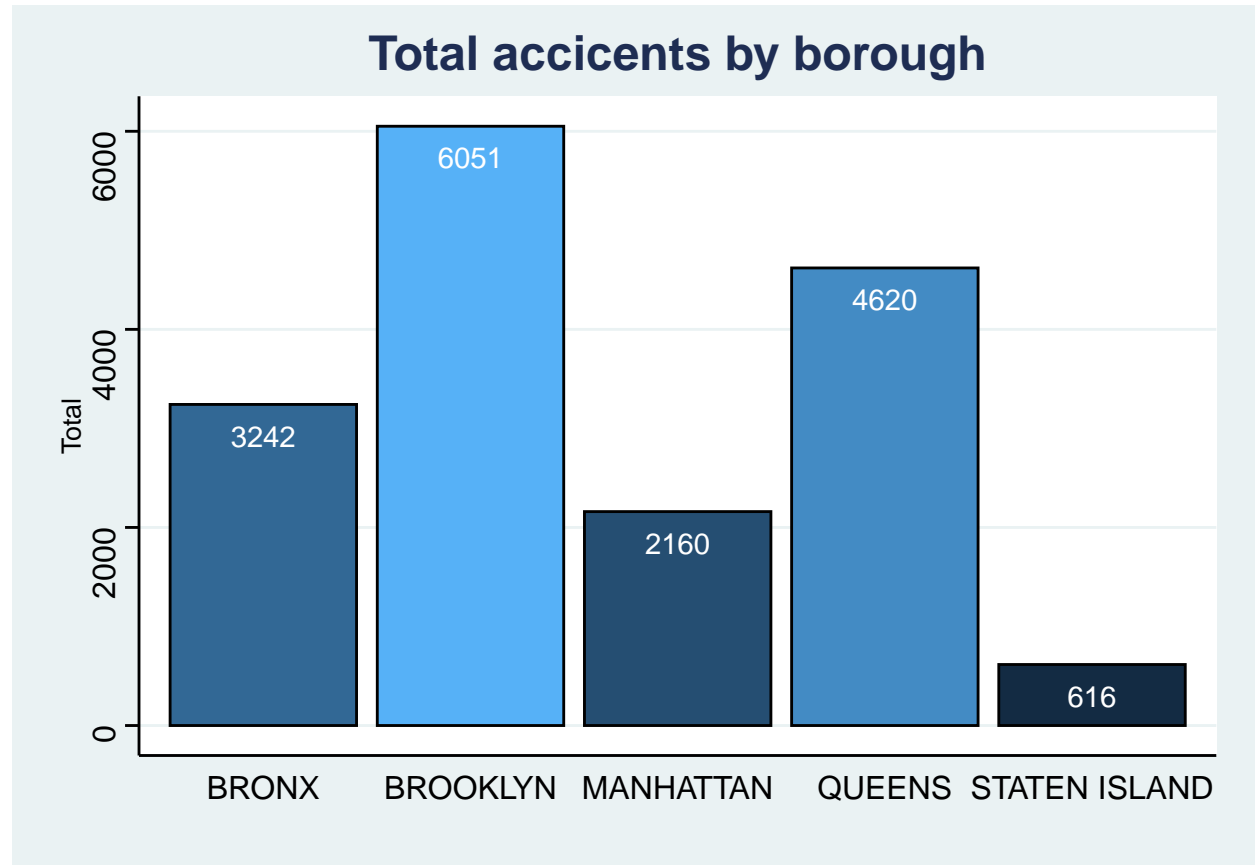


**Total accicents per day of the week in NYC**

We found that Friday and Saturday had higher numbers of accidents per day, but overall not significant enough, there could be several reasons for this:

- Weekend traffic: One of the most significant factors that could have contributed to a higher number of accidents on Fridays and Saturdays is weekend traffic. On weekends, more people may be out and about, and more vehicles may be on the road, leading to a higher probability of accidents.

- Alcohol consumption: Another factor that could have contributed to a higher number of accidents on Fridays and Saturdays is alcohol consumption. On weekends, people may be more likely to consume alcohol, which can impair their driving ability and increase the risk of accidents.

- Time of day: It is also possible that the time of day could have contributed to the higher number of accidents on Fridays and Saturdays. People may be more likely to travel at night or later in the day on weekends, which can increase the risk of accidents due to reduced visibility or fatigue.

**Total injured and killed per borough**

| BOROUGH | accidents | injuries | deaths |
|---|---|---|---|
| | 25741 | 10831 | 72 |
| BRONX | 9417 | 3232 | 10 |
| BROOKLYN | 16907 | 6024 | 27 |
| MANHATTAN | 7353 | 2151 | 9 |
| QUEENS | 14017 | 4600 | 20 |
| STATEN ISLAND | 1446 | 610 | 6 |

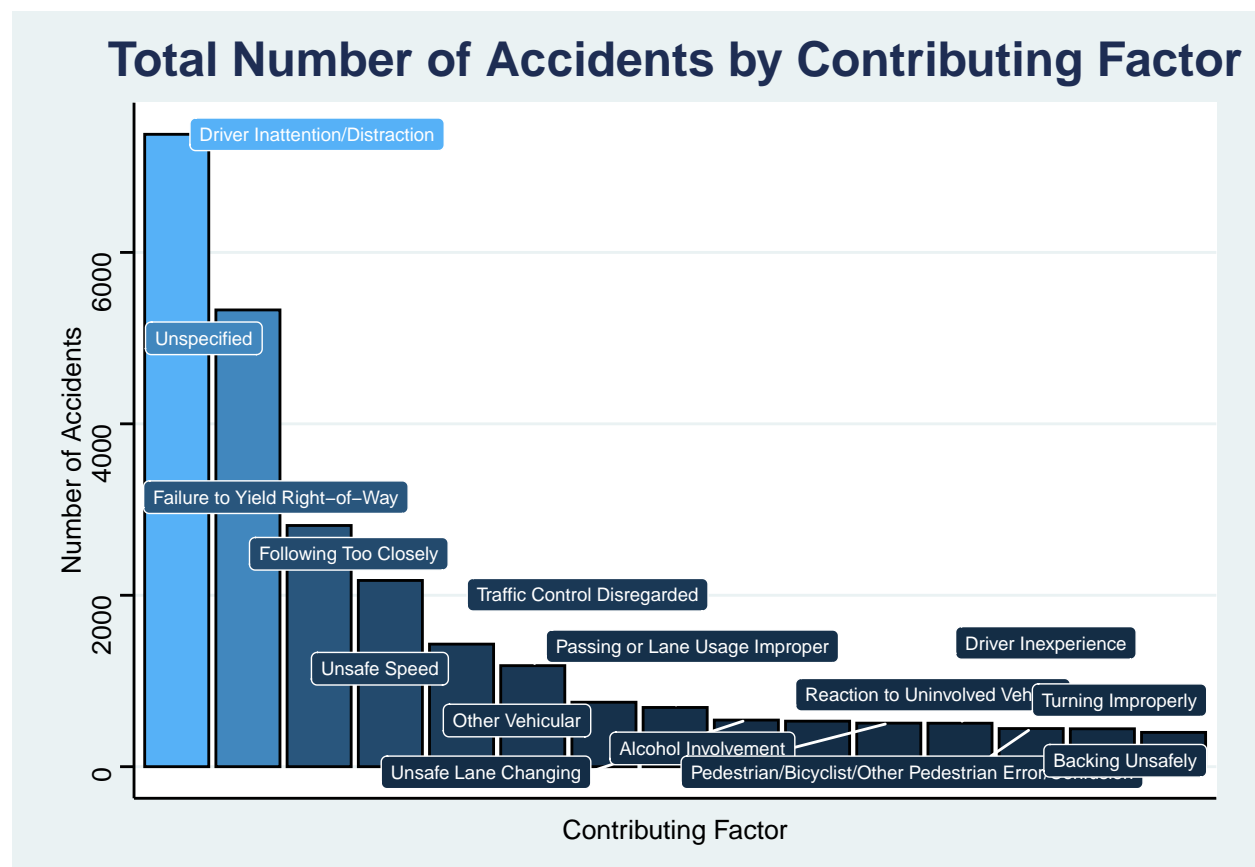## Total accicents by borough



We found that Brooklyn had the highest number of accidents, followed by Queens, and we are surprised that Manhattan did not have the highest number, there could be several reasons for this:

- Differences in traffic patterns: Brooklyn and Queens are both large boroughs with a lot of traffic and major highways running through them, which could explain the higher number of accidents. Manhattan, on the other hand, is a smaller borough with fewer highways, which could result in fewer accidents.

- Differences in population density: Another factor that could contribute to differences in the number of accidents is population density. Brooklyn and Queens have a higher population density than Manhattan, which means there are more people and vehicles on the road, increasing the likelihood of accidents.

- Differences in road infrastructure: Road infrastructure can also play a role in the number of accidents. Some areas may have poorly designed roads, confusing intersections, or inadequate signage, which can increase the risk of accidents.

- External factors: It is also possible that external factors such as weather, road conditions, or driver behavior could have influenced the number of accidents in each borough.

**Total Number of Accidents by Contributing Factor**



We found out that distraction is the higest followed by failed to yield and others, this suggests that policies and interventions aimed at reducing distracted driving and improving yielding behavior could help reduce the number of accidents.

- Public awareness campaigns: The city could launch a public awareness campaign to educate drivers about the risks of distracted driving and the importance of focusing on the road. The campaign could use different channels such as social media, radio, TV, and billboards to reach as many drivers as possible.

- Enforcement of traffic laws: The police department could increase the enforcement of traffic laws related to distracted driving and failing to yield. This could include targeted patrols in areas with high accident rates or increased use of traffic cameras to capture and penalize drivers who violate these laws.

- Infrastructure improvements: The city could invest in infrastructure improvements such as better road markings, traffic signals, and signage to help drivers navigate safely and avoid accidents.

- Driver education: The city could consider incorporating education about distracted driving and yielding behaviors into driver education programs, to ensure that new drivers are aware of the risks and are better prepared to drive safely.

By implementing policies and interventions that address the leading contributing factors to motor vehicle collisions, the city can work towards reducing the number of accidents, injuries, and fatalities on the roads.

**Conclusion**

In this project, we analyzed the New York City motor vehicle accidents dataset from 2020 to gain insights into the frequency, distribution, and causes of traffic accidents in the city. We utilized data visualization techniques and statistical analysis to explore various aspects of the data and identify patterns and trends.

We found that the number of accidents in New York City decreased significantly in 2020, which we attributed to the COVID-19 lockdown and other factors such as weather and changes in traffic patterns.

We also found that certain factors such as borough, day of the week, and time of day were associated with a higher probability of accidents. Specifically, Brooklyn and Queens had the highest number of accidents, Fridays and Saturdays had the highest number of accidents per day, and the afternoon and evening had a higher probability of accidents.

Furthermore, we identified several potential factors that could contribute to accidents, including alcohol consumption, traffic congestion, and road infrastructure. We also noted that heteroscedasticity was present in our regression model and used robust standard error estimation to correct for this issue.

Our findings have important implications for policymakers, law enforcement, and public safety officials, who can use this information to develop targeted interventions and strategies to reduce the number of accidents in the city. For example, they can focus on improving road infrastructure, increasing public awareness about the dangers of driving under the influence, and implementing measures to reduce traffic congestion during peak hours.

Overall, our analysis provides a valuable contribution to the ongoing efforts to improve traffic safety in New York City, and our findings can serve as a foundation for future research in this area.