

## Causality & Mitigations - Analysis of NYC accidents

SIADS 591 Milestone I Project

Amit Jha & Sahil Pujari

September 18, 2021

**Note:** For best experience, please view this report in ArcGIS StoryMap -  
<https://storymaps.arcgis.com/stories/bd7e45b6ae874c45a3642ca2632a8a19>

### Motivation:

Approximately 1.3 million people die each year on the world's roads, and between 20 and 50 million sustain non-fatal injuries [1]. Road traffic crashes cost most countries 3% of their gross domestic product. Road traffic injuries are the leading cause of death for children and young adults aged 5-29 years [2].

We come from India and have first hand experience of frequent road accidents and related fatalities. In 2019 alone there were 449 thousand crashes causing 151 thousand deaths and 451 thousand injuries [3]. That means, a crash is happening every one minute causing a death every three minute in India. This is huge.

Geo tagged data, needed for analysis is not yet available for India and hence we decided to analyze the traffic crashes for New York.

With the power of data science, we want to be able to study factors that contribute towards accidents in the city of New York and suggest strategies that can be taken to reduce them. The concepts taught in the MADS program so far in our curriculum can be applied well towards the datasets we have found for this task.

### Data Sources:

We have used public datasets for this analysis. A brief overview of the datasets and the variables contained is provided below:

1. **NYC Open Data** | Motor Vehicle Collision: The data source contains of 3 datasets - [crashes](#) (crash info detail and factors contributing to the collision), [vehicles](#) (fields related to vehicles involved in the collision), and [persons](#) (fields related to people involved in the collision). The data is obtained by records maintained by New York State for accidents reported between the period of 2015 - 2019.
2. **Traffic Speed Data:** When we started we selected dataset from the common source - [NYC Open Data](#), but provided dataset through this was of 22 GB and contained 52 mn records as of 15-Aug-21, hence was very difficult to process. Fortunately we found another source - [NYC REAL TIME TRAFFIC SPEED DATA FEED\(ARCHIVED\)](#) for the same data provided by [data.BetaNYC](#).

This dataset provides monthly compiled data through different links on the site. For reproducibility, we decided to do webscrap the site to download data for five years(2015-2019) programmatically.

This dataset contains 'real-time' traffic information on five minute intervals from locations where [DOT](#)(Department of Transportation) picks up sensor feeds within the five boroughs, mostly on major arterials and highways. There is one meta location dataset([linkinfo dataset](#)) provided with traffic speed dataset to be able to combine the traffic speed data with the linkinfo(polyline road segment) location data.

3. **Hospitals Locations:** This is a shapefile dataset, provided by Homeland Infrastructure Foundation-Level Data([HIFLD](#)) which provides the location details for the [hospitals](#) available in the USA. This dataset is used to find if a crash incident was close enough to a hospital to get the required treatment in time.
4. **Traffic Volume:** This [dataset](#) provides traffic volume counts collected by DOT for different roadways at different times of the day.
5. **Population:** Have used two population datasets - first provides the [USA population](#) for different years. And the second dataset shapefile, provides US population by ZIP code provided by 2010 US decennial Census accessed via [Kaggle](#).
6. **Driving Cost:** This is the most interesting dataset used in the project and was not available anywhere in the compiled form. So we created it ourselves by compiling the data from a global news and lifecycle publication [INSIDER](#).
7. **Traffic Related Death Rate by Countries:** This [dataset](#) provides traffic related death rate for various countries sourced from the WHO - [Global Status Report on Road Safety 2018](#).

## Data Manipulation Methods:

**Traffic Collision - Crashes:** This is the most important dataset for the project. We started by printing dataset info by using pandas info() method to see data type and missing values data. And we saw pandas were not able to assign efficient data types to fields. Also interestingly coz of data size being huge 1.8 mn records it doesn't show missing values info. Hence we wrote our custom utility to print missing values percentage for the dataset.

As pandas is not able to assign efficient data types, we assigned appropriate data types for the columns through the preprocessing utility we wrote for this dataset.

Further, we wanted to see which crash point is on which road segment or link(linkId is provided through another meta dataset of Traffic Speed dataset), to find the approx speed of the vehicle involved in the crash. Linkinfo dataset contained multiple geolocation points(lat/long) as linkPoints(polyline).

For this we wrote a utility, using below logic:

- Find the max\_lat/long, min\_lat/long for all the polylines(linkId)
- Check if the crash point lies between max/min lat/long points of a polyline

However, we found the linkPoints field to be often corrupted and having junk numeric value. Leading to the crash point linked to any polyline through above utility actually being faraway from the road segment on plotting crash point and corresponding polyline.

Google encoded polyline fields available on linkinfo dataset came to our rescue. Here too, we faced the challenge of escape characters being present into encoded polylines. So we decided to use try/catch and used a regular expression package - "re" to escape character by using re.escape(text) if polyline is not able to decode in try block.

We could decode this field into link points(pairs of lat/long points) using a polyline package and get corresponding geolocation points of the road segment correctly.

While this too could not decode all the polylines(two of 153), we decided to leave the problematic values from the analysis.

**Traffic Speed dataset:** Because of its huge size(approx 3 GB) we decided not to combine monthly traffic speed data into one file, and instead created yearly files for five years from 2015-19 and used them as and when required for analysis. Also we assigned appropriate data types for fields to avoid running into memory issues.

As speed was not part of the crashes dataset, to see the crashes and traffic speed relation, we decided we can combine crashes and traffic speed data. Through earlier processing mentioned for crashes dataset, we could find at which road segment(linkId) crash occurred and now if we combine the two datasets(using timestamp), we can then see the average traffic speed at at crashing time, this can give us insight into crashes vs traffic speed.

To combine these two dataset, we used the pandas - merge\_asof() method, using linkId as column to merge on before checking for combining on timestamp. Also to get traffic speed very close to crash time, we used the tolerance of 5 minutes.

This way finally we were able to combine the crashes and traffic speed dataset, which initially didn't have any common field or info to connect them.

**Traffic Volume Dataset:** While this dataset provides traffic volume at different times of the day, unfortunately it does not provide it for all the days in any year. So we decided to find average traffic volume per day for each road by using pandas group by and computed mean over grouped feature to see any correlation between traffic accidents and different roadways. Also as traffic volume is high compared to crashes data, we computed the log of the volume to better visualize the relation. Further as volume would differ for different road types like expressways/highways, we created a field to differentiate road type - expressway/highway etc

## Analysis and Visualizations:

In this project we wanted to answer several questions related to traffic crashes/fatalities. We expand upon all of these in our notebook but for the purpose of size limitations of the report, we focus on the ones highlighted in bold -

- 1. Most common causes of traffic crashes in New York?**
- 2. How does driving skill (proxied by license cost) affect crashes and their injuries / deaths?**

- 3. Traffic count / time relation with accidents for each zip code**
- 4. Traffic speed / time relation with accidents for each zip code**
- 5. Could lives be saved by quicker medical attention, i.e. hospitals closer to accident sites? Are there any locations, which do not have hospitals nearby enough?**
6. Are there any locations which are more prone to traffic accidents than others?
7. Are there any sections of society which are more prone to accidents, causal/victims like pedestrians, cyclists or bike users in specific locations?
8. Attributes of vehicles with most accidents

## Most Common Causes of Traffic Crashes:

To be able to prevent crashes it is important to first understand the most common causes for crashes.

To visualize this we used Altair library, which provides declarative visualizations in python. It is built on top of the Vega-Lite high-level grammar for interactive graphics.

From this we can clearly see that out of the top 10 causes, the top nine are driver's mistakes. This gave us the insight that perhaps driving skill is the most important aspect that needs to be improved to prevent high number of crashes in NYC (Figure 1).

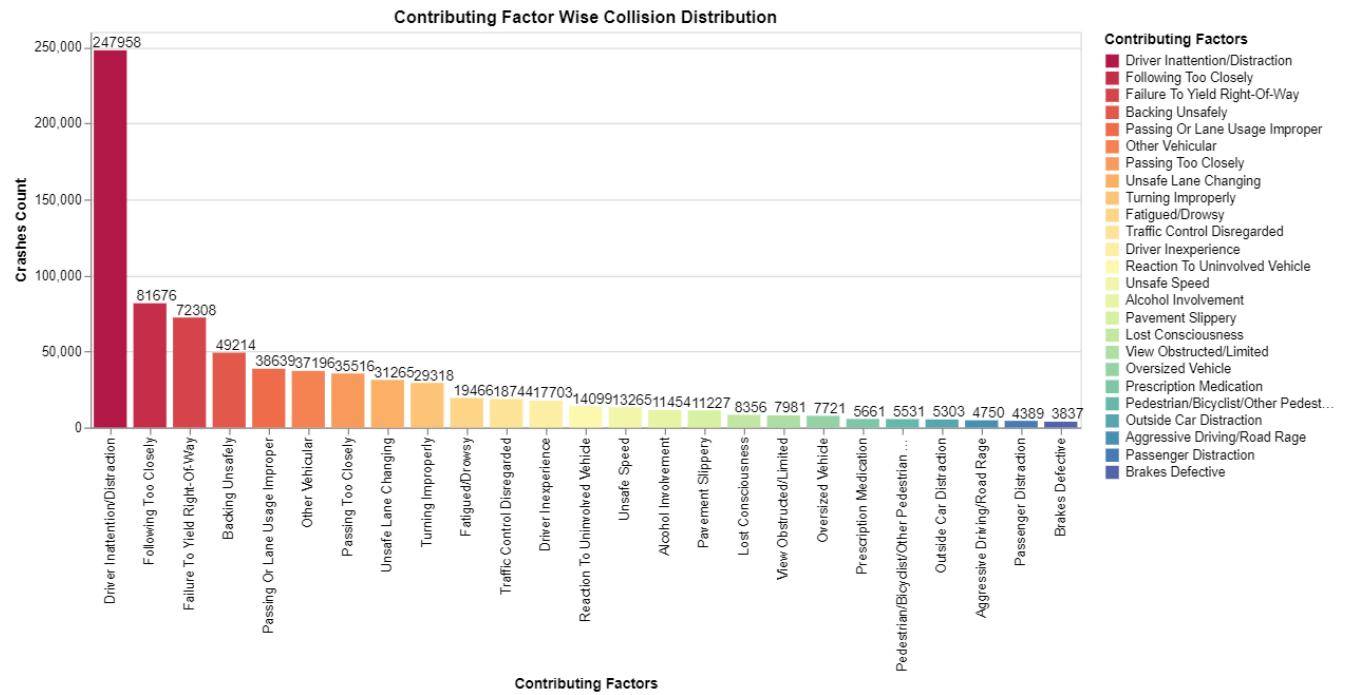


Figure 1

**Recommendation:** Apart from the #1 factor driver inattention, rest of the factors can be improved by defensive driving education in the public and are directly related with techniques that this driving education entails.

## How does driving skill (proxied by license cost) affect crashes:

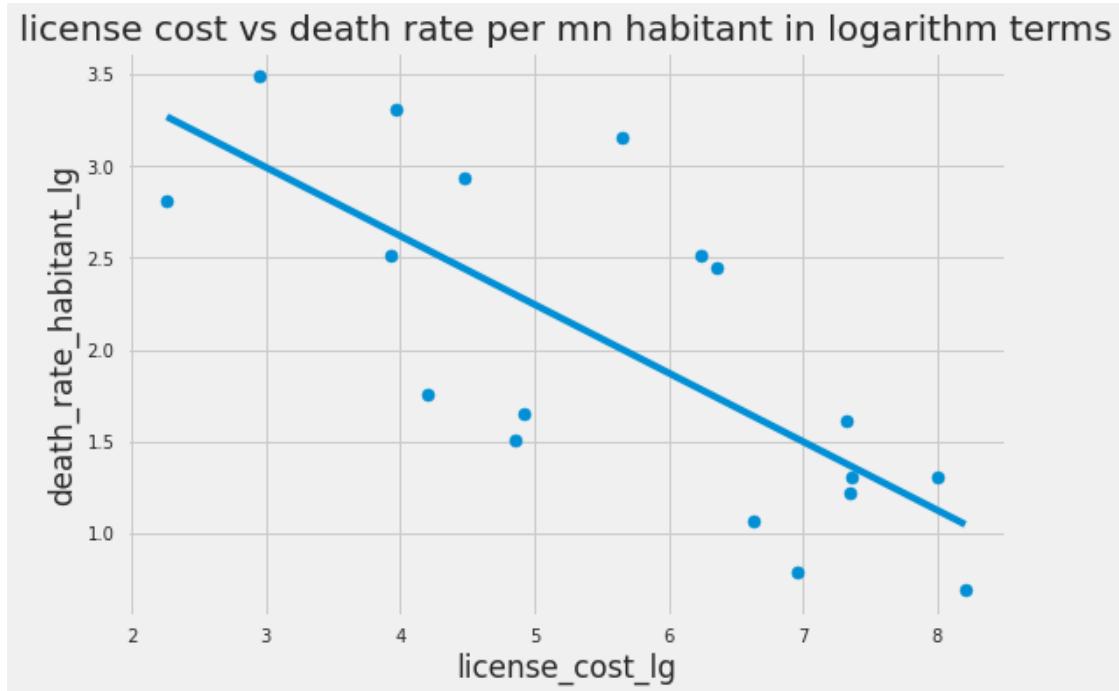
We tried to find data for training skill but could not find it and so decided to use cost for getting a driving license as a proxy for training skill.

Data for license cost and death rate were not on the same scale and were heavily skewed. Hence we transformed both the license\_cost and death rate to logarithm scale and then plotted the scatter plot to see the correlation. We can clearly see the correlation here.

We computed the Pearson correlation(-0.54) and found a negative moderate correlation between license cost and death rate.

We also tested a NULL hypothesis by setting alpha at 0.05 - that is "Death rate is same across different license cost" and got a p-value of 0.0024. That is, the NULL hypothesis is rejected and the death rate is not the same across different license costs. **Now this may be misleading.**

While there is inverse correlation between average license cost and death rate, this does not mean there is causal relationship. But what we can say, is that more training and tests require additional cost and hence cause an increase in license cost which further causes reduction in traffic crashes because of more skilled drivers on the road(Figure 2) .



**Figure 2**

**Recommendation:** We would recommend New York Department of Transportation to look at driving training/test methodology at top five countries with lowest traffic related fatalities(Norway, Sweden, Ireland, Denmark,Germany) and adopt the best practices to achieve the adopted Vision ZERO by NYC.

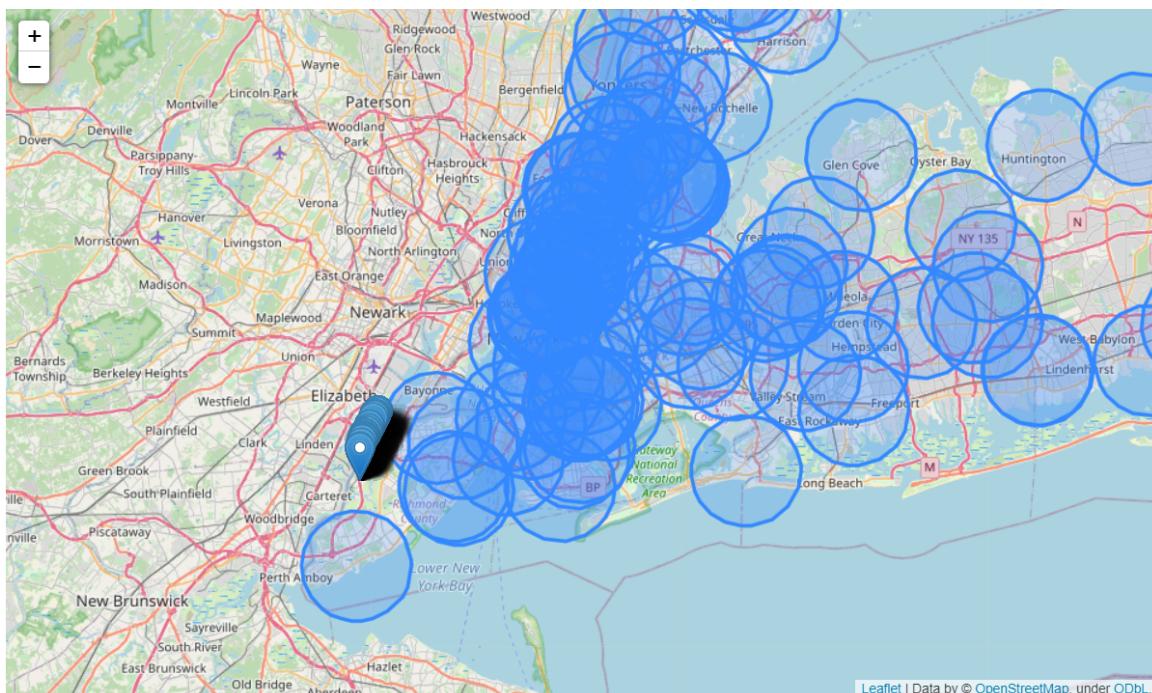
## Are Hospitals Close Enough to Crash Sites:

As per WHO report, much of traffic related fatalities can be prevented by timely treatment of accident victims. In that endeavour, we analyzed the hospital coverage of accidents and checked if new hospitals are required and at what locations, and if this can prevent some casualties.

For this we used Folium library and instead of computing the minimum distance to each hospital, which would be time consuming, we utilized Folium buffer() method, to create a buffer of one mile around the hospital and then combined the connected buffers to create a unary buffer. Post that we checked if a crash location was outside the buffer range. While we appreciate that New York has good coverage of Hospitals, and that's the reason, even on choosing a small buffer of one mile we don't see many crashes outside the buffer range.

**Recommendation:** There still is scope of improvement and 245 crashes, we see happened outside the hospital buffer range in 2015-2019.

This area of *West Shore Expressway* can be further equipped with hospitals/primary care centers to prevent casualties(Figure 3) .



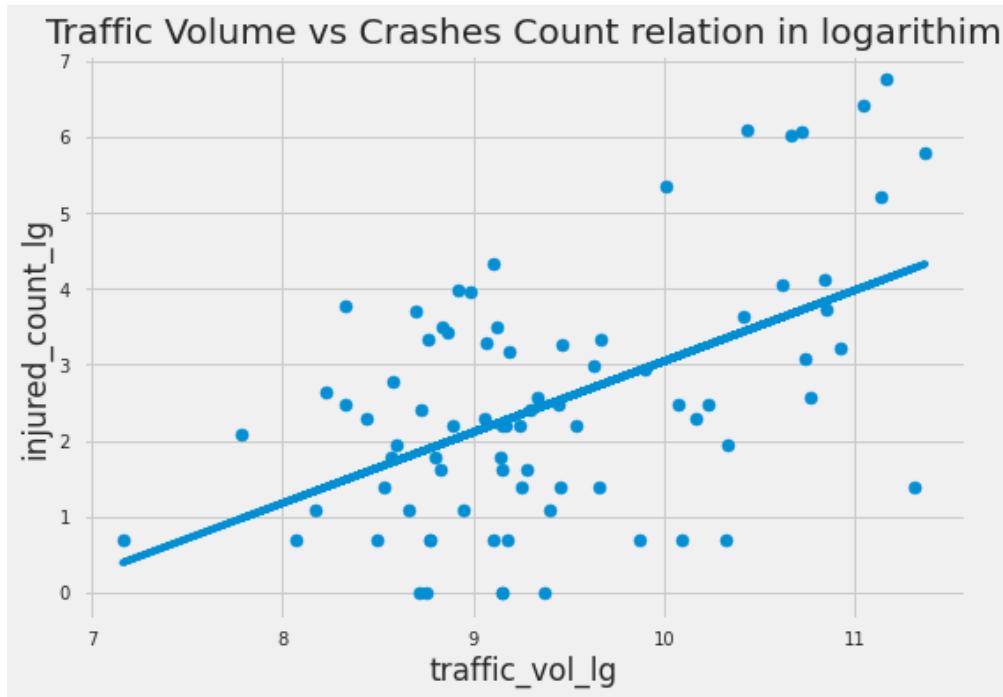
**Figure 3**

## Traffic Volume Relation with Crashes:

We faced multiple hurdles in answering this question. And biggest of this was that traffic volume data provided by [NYC Open Data](#) was not complete. While it provides the data at different times of the day, it doesn't provide for all the days/months of the year. Leaving the data incomplete.

We still tried to get around the problem by computing the average daily traffic volume. Further we saw traffic volume and injured counts are highly skewed and so we transformed both the fields by taking their logarithms. With log of traffic volume and injured count plotted we can clearly see moderate positive correlation of traffic volume with crashes count. We computed the Pearson correlation coefficient and got moderate positive correlation of 0.59(Figure 4) .

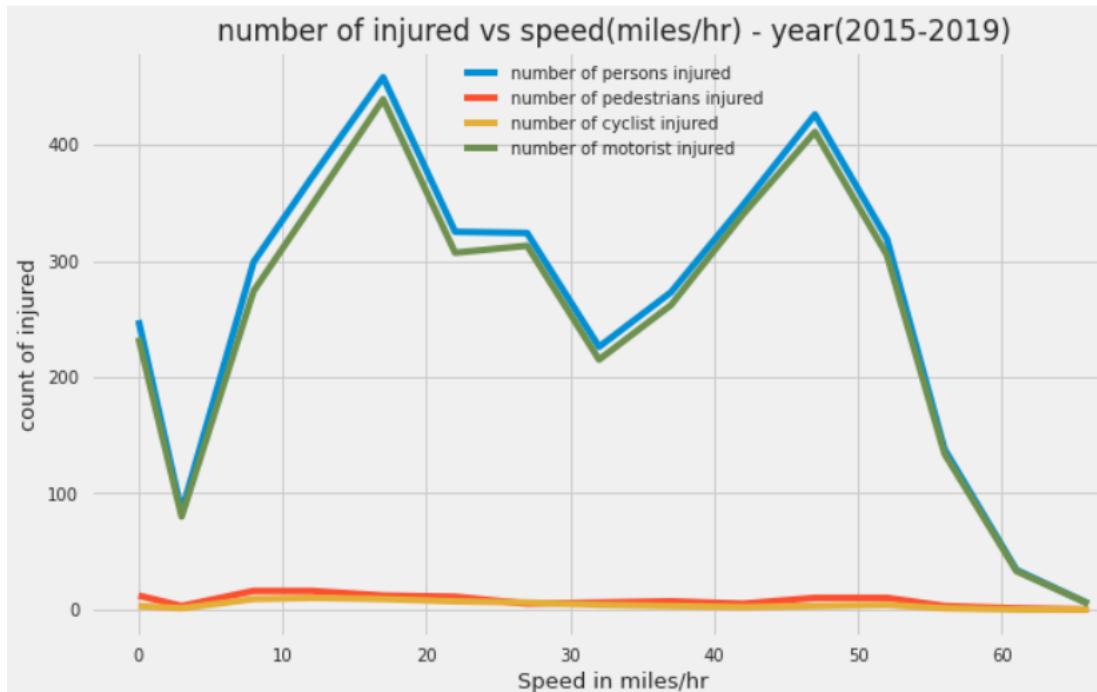
**Recommendation:** We would recommend decongesting the roads with most traffic volume to prevent the crashes.



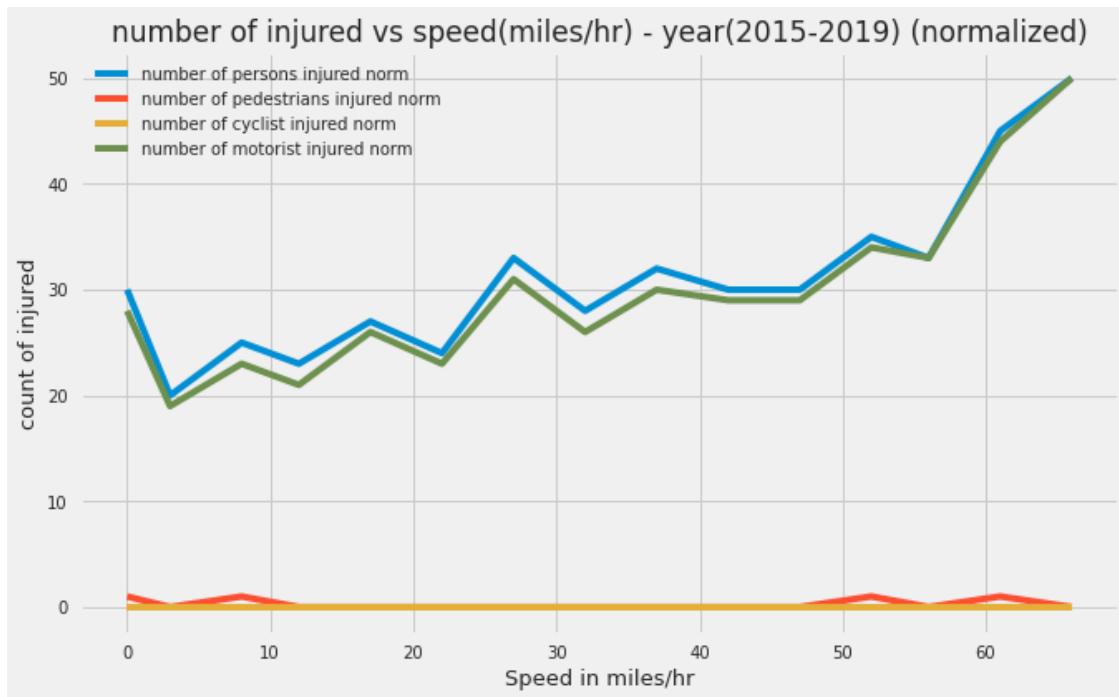
**Figure 4**

**Traffic Speed Relation with Crashes:** Traffic speed dataset was the most challenging to handle in manipulation, we mentioned earlier the methods we used to combine it with crashes dataset.

It was also tricky to interpret the traffic speed and crashes relation too, which we could see on plotting. When crashes/speed plotted directly, it seemed to suggest that the number of injuries tend to decrease after the speed of approx 45 miles/hr, which is contrary to common belief. Obvious point of suspicion was on data volume available for crashes at different speeds(which we had got after merging crashes/traffic speed) (Figure 5).

**Figure 5**

So to nullify this impact, we normalized the crash injuries data by dividing by number of crashes. This seems to solve the problem (Figure 6).

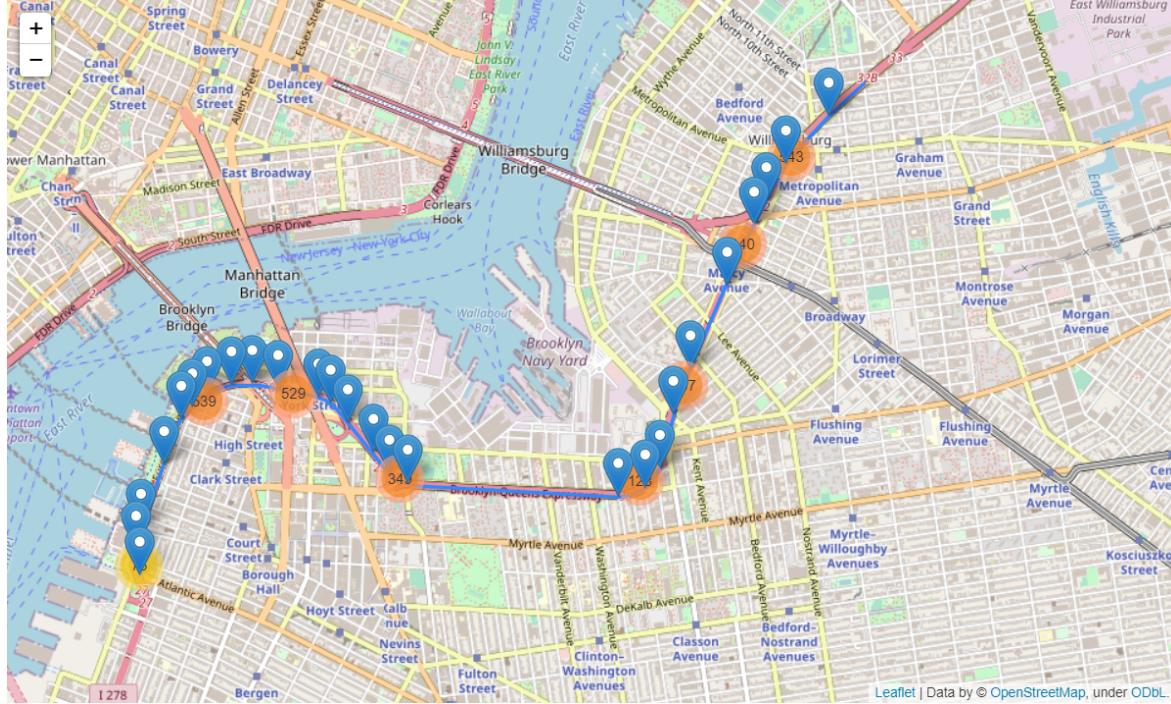
**Figure 6**

One of the most important points we can see here is how injury counts are higher at the speed of ZERO and decrease at further speed initially.

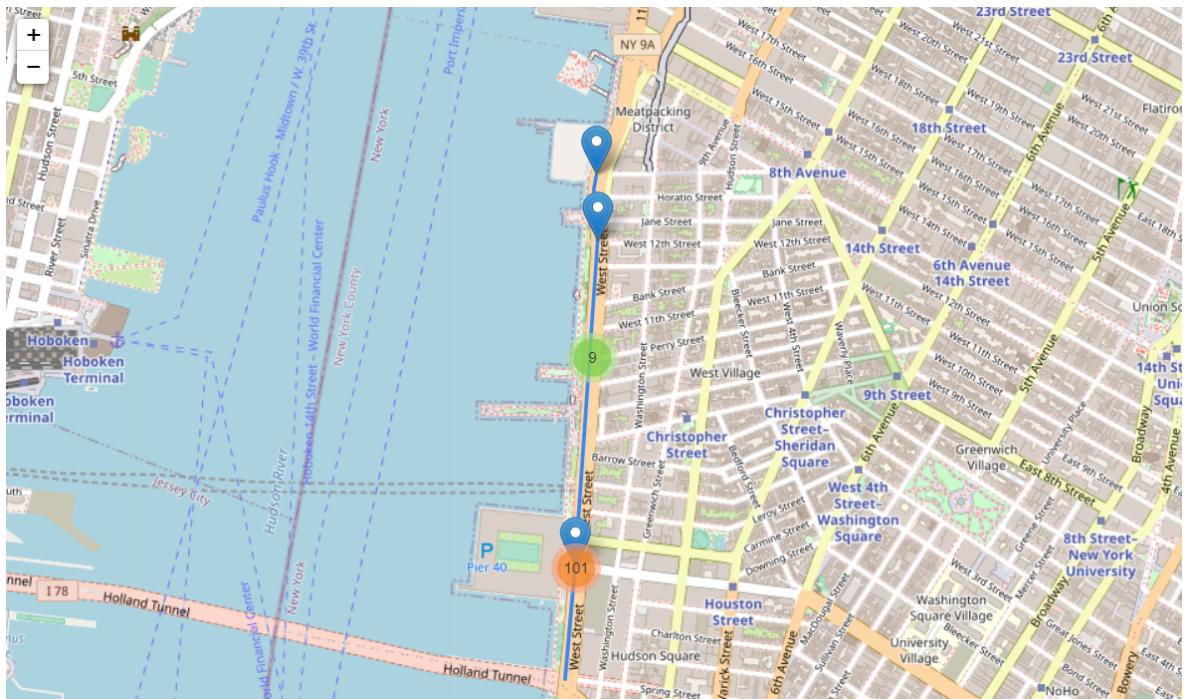
This suggests crashes often happening with static vehicles too, as pointed in the notebook.

Further, we can clearly see how the number of injuries increases steeply after the speed of 55 miles per hour.

**Figure 7**



Leaflet | Data by © OpenStreetMap, under ODbL.



**Figure 8**

Further, we plotted the crash points and corresponding poline road segments, and we could confirm that our crash point to polyline mapping is correct.

Further, we plot this for a few top crash polylines like - Brooklyn Queens Expressway, West Street and Fdr Drive in terms of crashes count.

From the Brooklyn Queens Expressway, we can see that there are more crashes at say bend points/turns and fewer at straight stretches (Figure 7).

Further on plotting the crash/polyline for West Street, we can see that most of the crashes are at intersections (Figure 8).

Similarly, we could see from plotting Fdr Drive, crash/polyline too that most of the crashes are at bend points/intersections.

#### **Recommendation:**

1. Enforce more safety features to prevent crashes at static mode.
2. We further need to look into how to prevent crashes at intersection/bend points.

**Further Work:** While fatalities at these points can be prevented by systematically reducing the speed using speed bumps. Some road design may also be required to see if it can provide more visibility. However this needs further research to confirm what best ways can be to prevent crashes at bend points/intersection.

#### **Statement of Work:**

Amit	Sahil
<ul style="list-style-type: none"> <li>• Knowledge and uses of Folium in spatial maps</li> <li>• Data Source Research and Selection</li> <li>• Data Cleaning</li> <li>• EDA and selection of final visuals</li> <li>• Report writing</li> </ul>	<ul style="list-style-type: none"> <li>• Extensive use and knowledge of leaflets</li> <li>• Data Source Research and Selection</li> <li>• Data Cleaning</li> <li>• EDA and selection of final visuals</li> <li>• Report writing</li> </ul>

#### **References:**

- [1] World Health Organization - global-status-report-road-safety-time-action  
<https://www.afro.who.int/publications/global-status-report-road-safety-time-action>
- [2] World Health Organization - Road Traffic Injuries Jun-21  
<https://www.who.int/news-room/fact-sheets/detail/road-traffic-injuries>
- [3] Ministry of Road Transport and Highways - Road Accidents in India - 2019  
[https://morth.nic.in/sites/default/files/RA\\_Uploading.pdf](https://morth.nic.in/sites/default/files/RA_Uploading.pdf)
- [4] International Cooperative and Mutual Insurance Federation  
<https://www.icmif.org/wp-content/uploads/2020/10/lb-app-road-traffic-accident.jpg>
- [5] The Parrish Law Firm  
<https://www.theparrishlawfirm.com/wp-content/uploads/2020/01/faqs-what-are-the-most-common-causes-of-car-accidents-in-virginia.jpg>
- [6] Proximity Analysis by Alexis Cook <https://www.kaggle.com/alexisbcook/proximity-analysis>