

01

DS-GA 1007 | FALL 2022

NYC Traffic Danger

PLUCKY PYTHONS

Questions

Question 1: What were the most dangerous times of year/week/day?

Question 2: What were the most dangerous areas in NYC in 2021?

Question 3: Which areas had the most bike accidents in NYC in 2021?

Question 4: Are e-bikes more dangerous than regular bikes?

Data

NYC OpenData Motor Vehicle Collisions

Variables of Interest: Crash Date, Crash Time, Zip Code, Latitude, Longitude, Vehicle Type(5), Contributing Factor(5) Persons Injured, Persons Killed, Cyclists Injured, Cyclists Killed

Cleaning:

Initial Data: 1,946,561 rows

Convert date to pandas datetime object

Filtered dataframe to 2021

Dropped rows with missing location values

Dropped unused columns

All Data: 101,757 rows

Bike Data: 7,331 rows

Filled in all missing zip code and boroughs using geopy API with longitude/latitude data

Methodology

Workflow

Question 1: Converted date and time columns into Pandas Datetime index to group number of traffic accidents in specified periods of time and generated line plots for comparison. Compared traffic accidents for various holidays using bar charts.

Question 2: Mapped latitude and longitude on folium maps, grouped by zipcode and graphed with geojson, and clustered with DBSCAN to find most dangerous intersections.

Question 3: Divided data by borough and zip code. Included data on bike commuters in each borough to estimate percentages of injuries and fatalities. Used folium to visualize hot spots of collisions with heat map.

Question 4: Extracted data on injuries and deaths for bikes and e-bikes only. Separated the groups and compared them using bar charts and resampling methods.

05

Original Dataset

NYC OpenData

Motor Vehicle Collisions - Crashes | NYC Open Data

cityofnewyork.us

Contact Us

Blog



Sign In

BETA

Introducing our new data shaping and exploration experience: Filter, group, aggregate, and more!

[Try it now](#)

[Learn more](#)



Motor Vehicle Collisions - Crashes

The Motor Vehicle Collisions crash table contains details on the crash event. Each row represents a crash event.



Find in this Dataset

More Views

Filter

Visualize

Export

Discuss

36

Embed

About

CRASH DATE	CRASH TIME	BOROUGH	ZIP CODE	LATITUDE	LONGITUDE	LOCATION	ON STREET NAME
09/11/2021	2:39						WHITESTONE EXPRESSWAY
03/26/2022	11:45						QUEENSBORO BRIDGE UPPER
06/29/2022	6:55						THROGS NECK BRIDGE
09/11/2021	9:35	BROOKLYN	11208	40.667202	-73.8665	(40.667202°, -73.866...	
12/14/2021	8:13	BROOKLYN	11233	40.683304	-73.917274	(40.683304°, -73.917...	SARATOGA AVENUE
04/14/2021	12:47						MAJOR DEEGAN EXPRESSWAY RAMP
12/14/2021	17:05			40.709183	-73.956825	(40.709183°, -73.956...	BROOKLYN QUEENS EXPRESSWAY
12/14/2021	8:17	BRONX	10475	40.86816	-73.83148	(40.86816°, -73.8314...	
12/14/2021	21:10	BROOKLYN	11207	40.67172	-73.9071	(40.67172°, -73.9071...	

< Previous

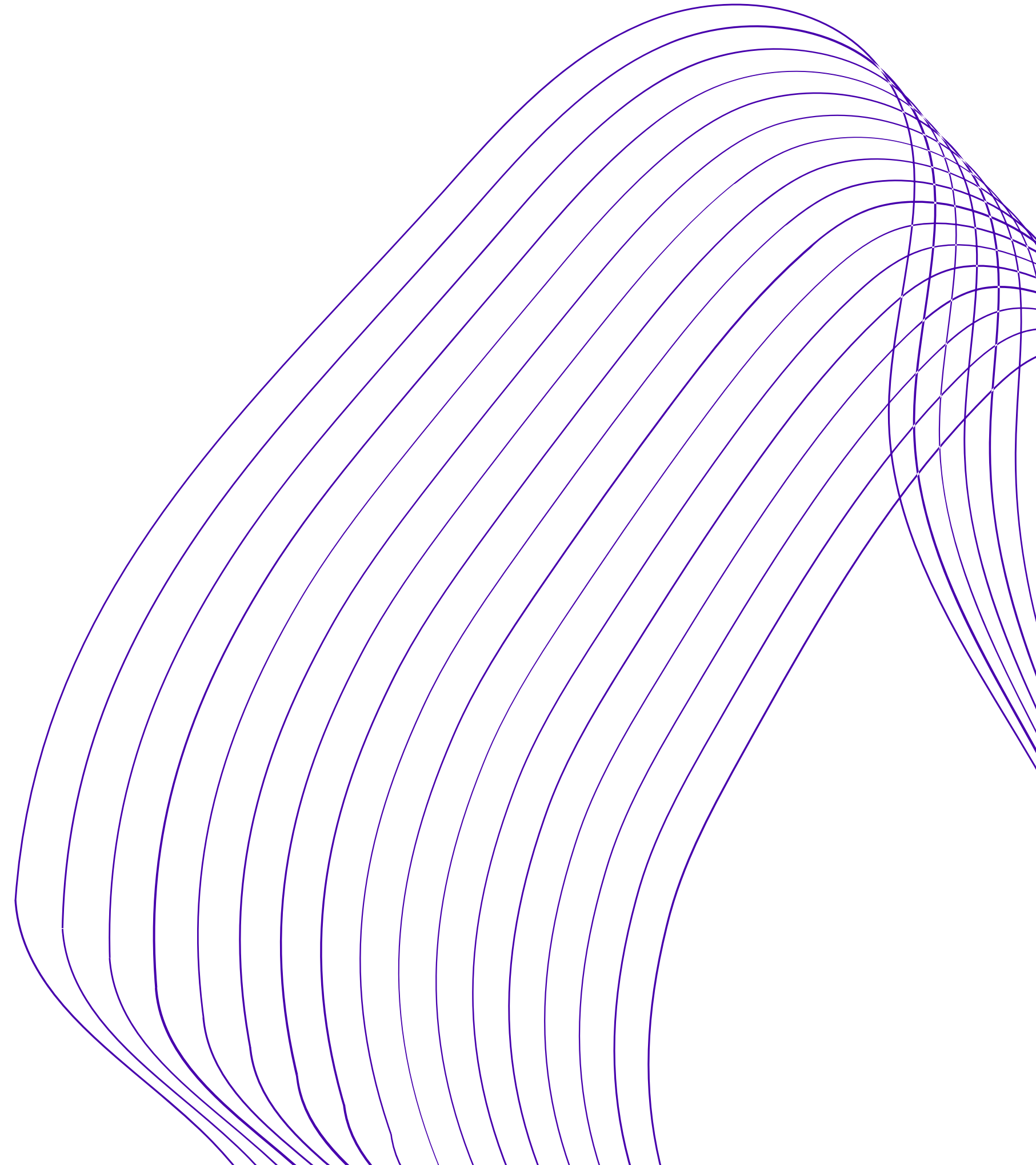
Next >

Showing Motor Vehicle Collisions 1 to 100 out of 1,946,561

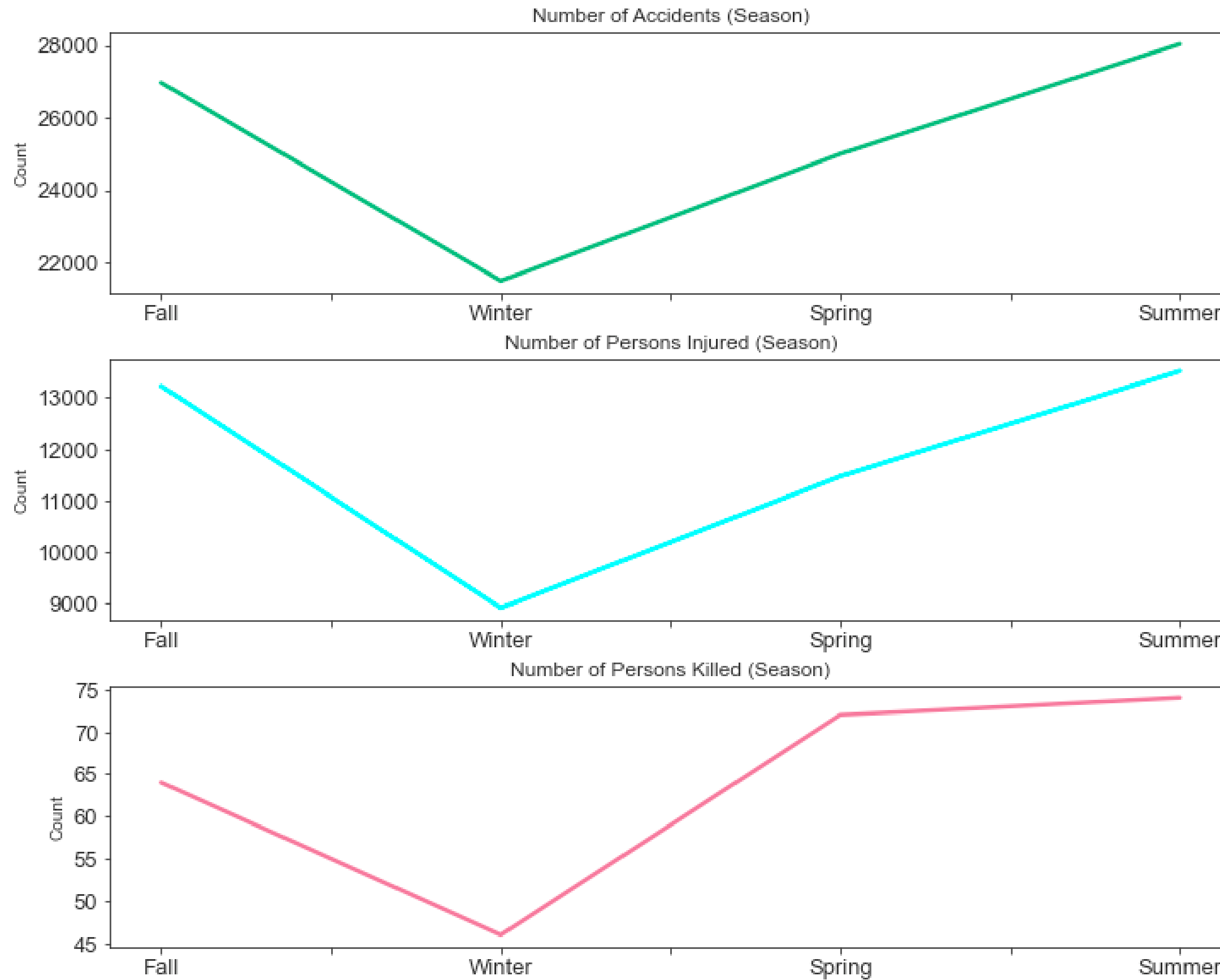
LINK: [HTTPS://DATA.CITYOFNEWYORK.US/PUBLIC-SAFETY/MOTOR-VEHICLE-COLLISIONS-CRASHES/H9GI-NX95](https://data.cityofnewyork.us/public-safety/motor-vehicle-collisions-crashes/h9gi-nx95)

What were the most dangerous times of year/week/day?

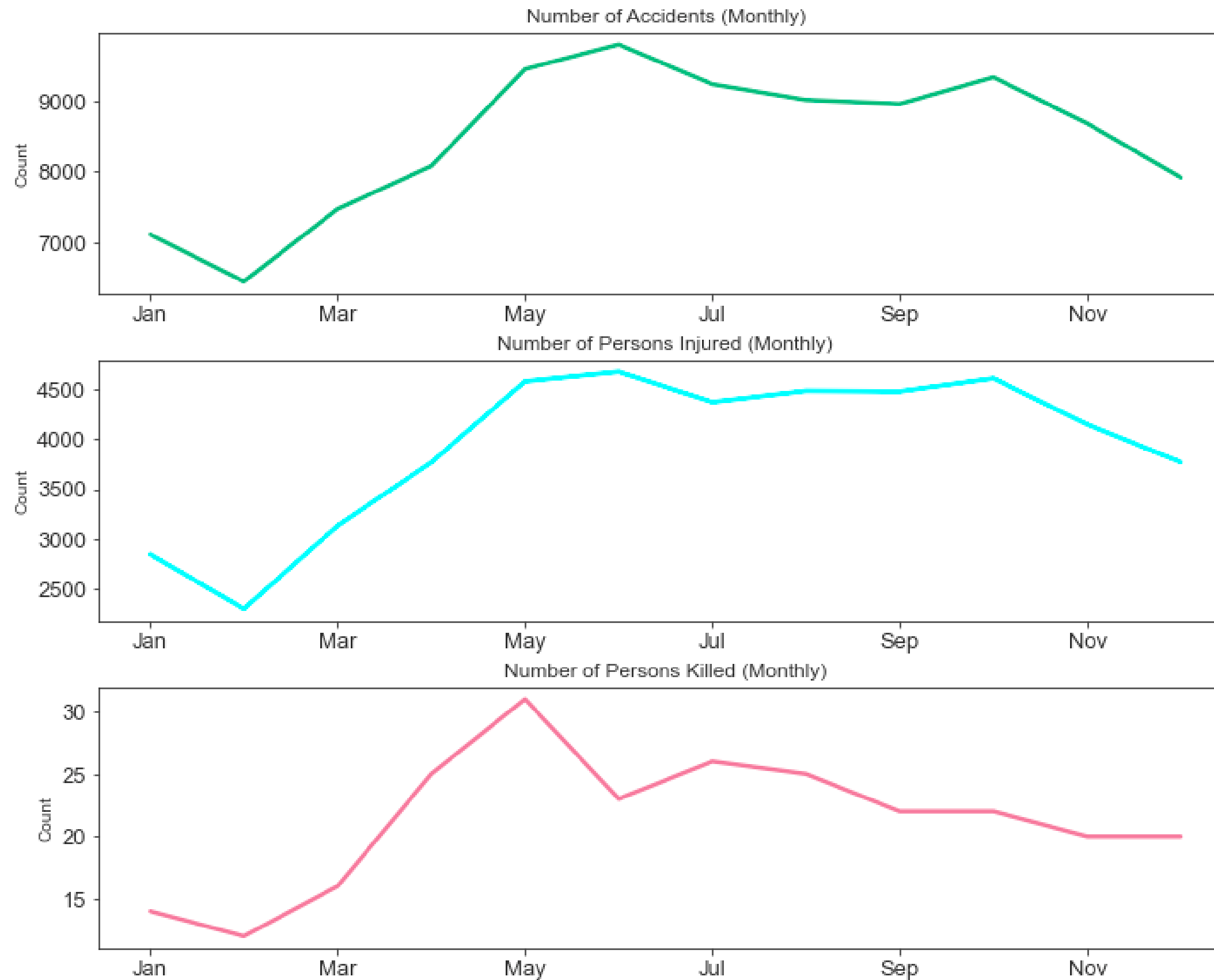
Method: Converted date and time columns into Pandas Datetime index to group number of traffic accidents in specified periods of time and generated line plots for comparison. Compared traffic accidents for various holidays using bar charts.



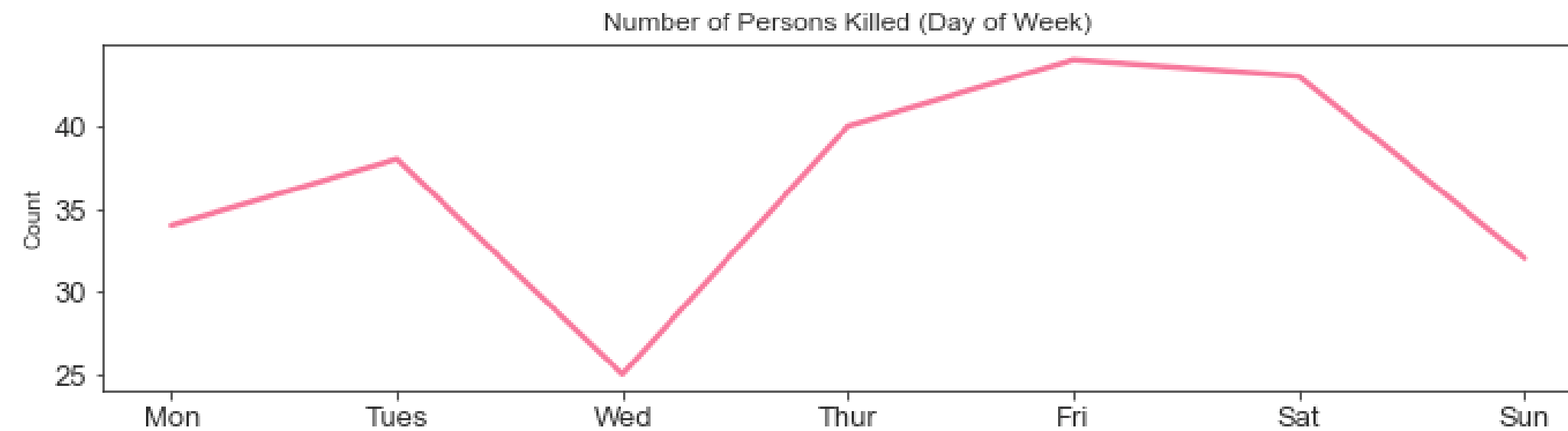
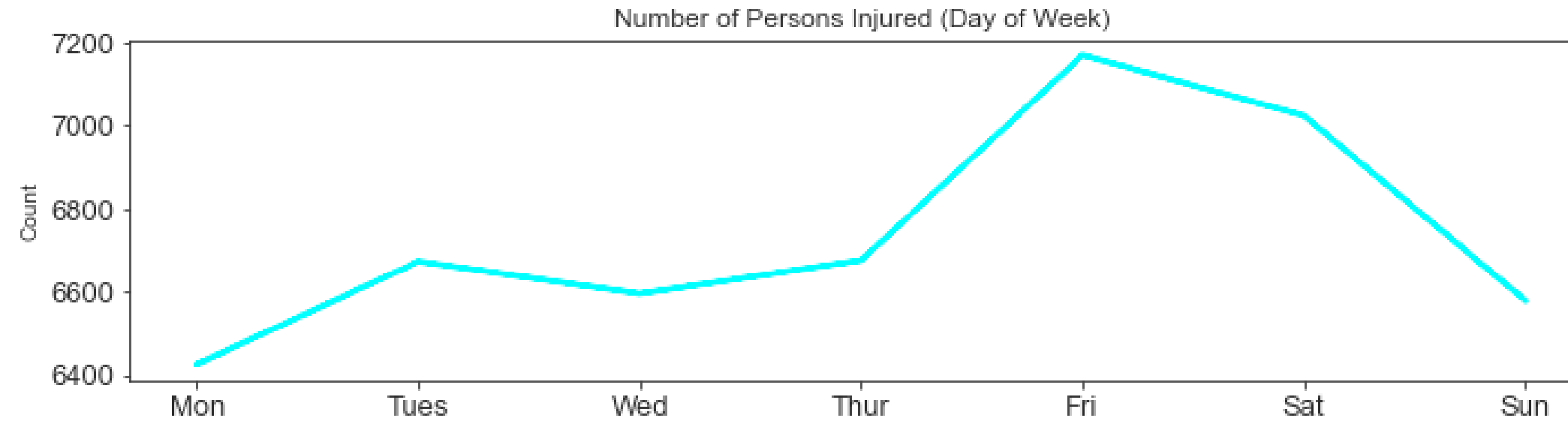
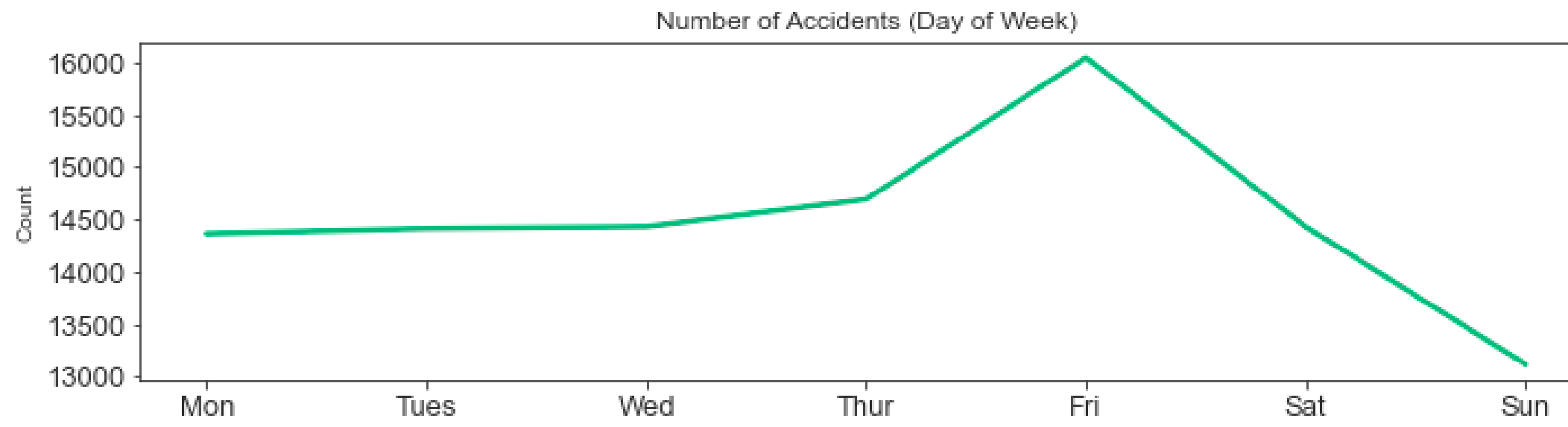
Seasons



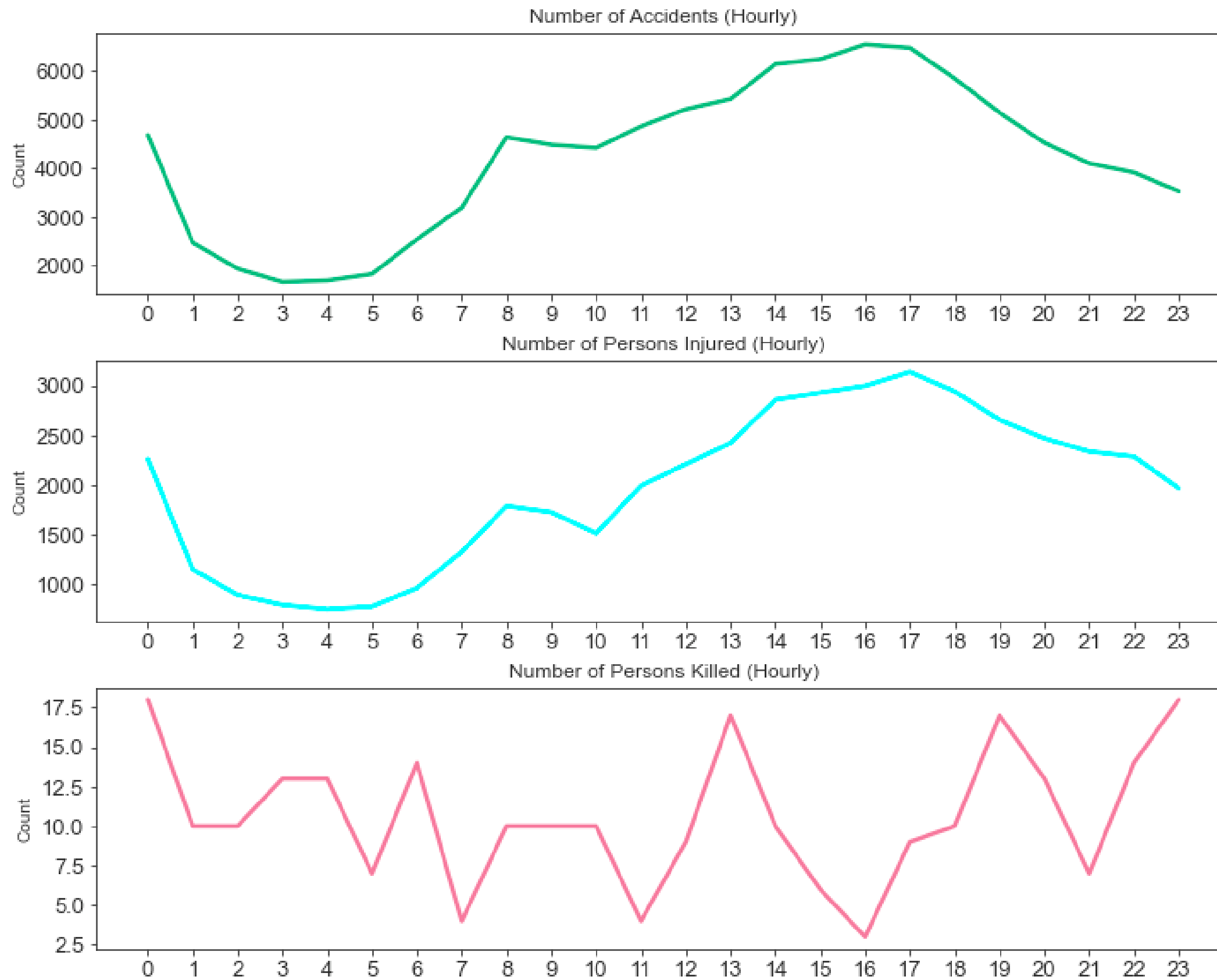
Months



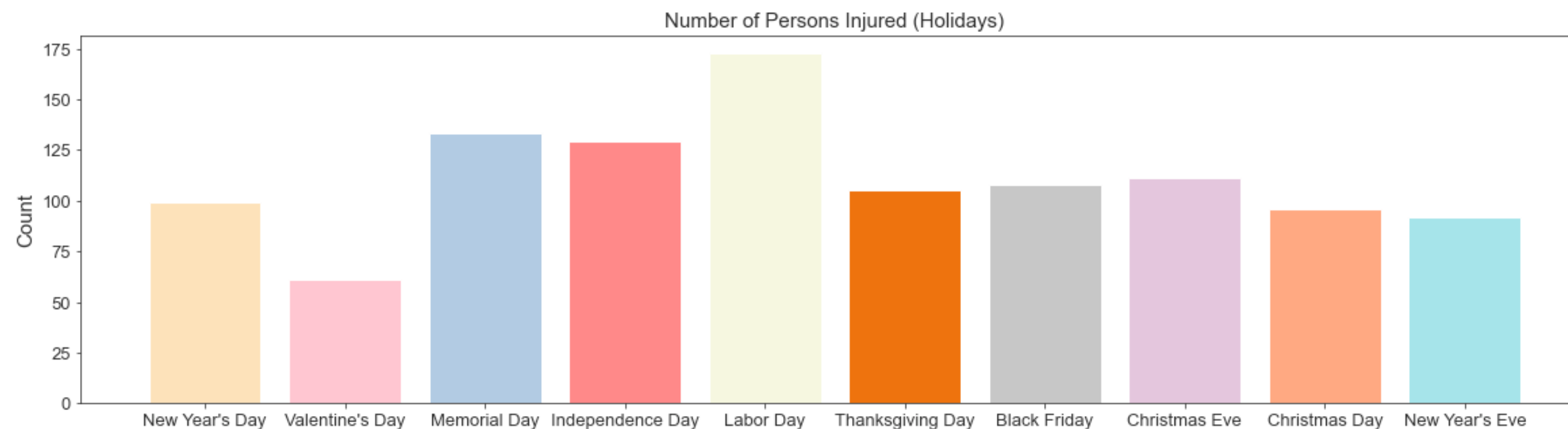
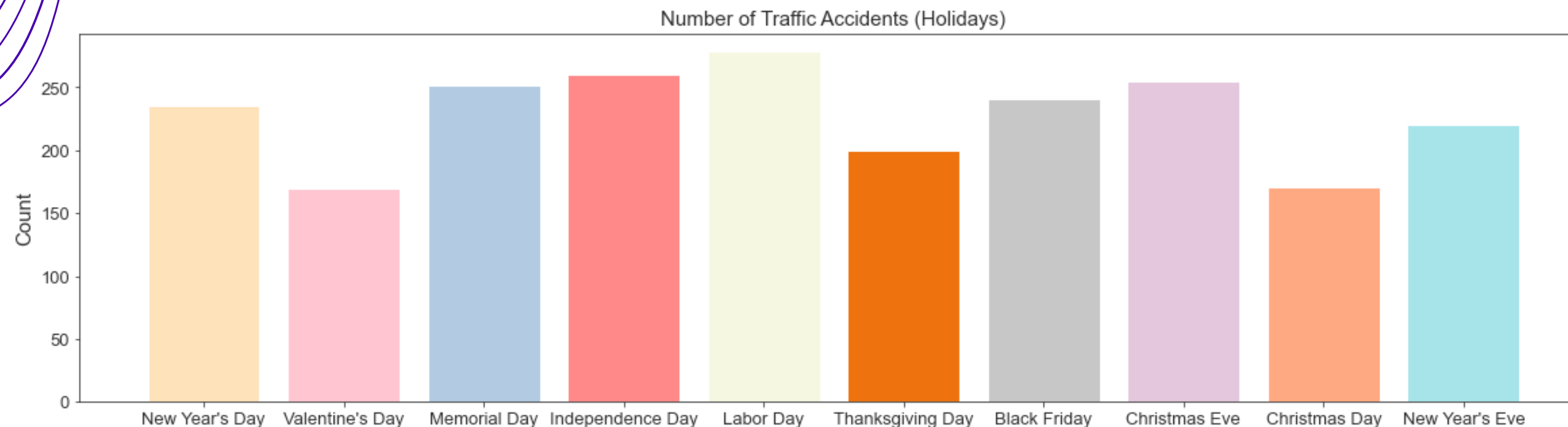
Weekdays



Time of Day

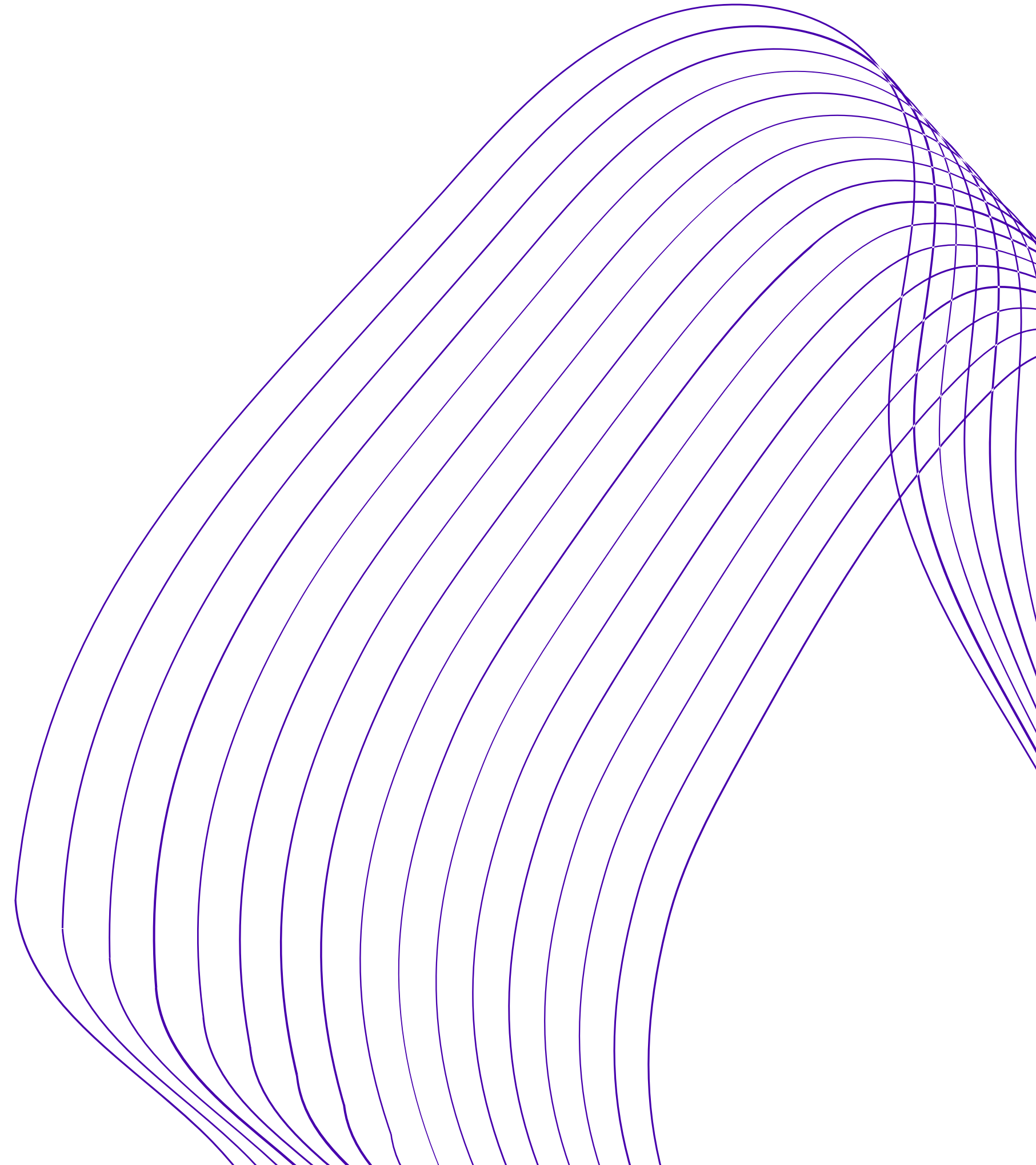


Holidays



What were the most dangerous areas in NYC in 2021?

Methods: Mapped latitude and longitude on folium maps, grouped by zipcode and graphed with geojson, and clustered with DBSCAN to find most dangerous intersections.



All Deaths Folium Map

```
if val == 'all persons':
    val = 9
elif val == 'pedestrians':
    val = 11
elif val == 'cyclists':
    val = 13
elif val == 'motorists':
    val = 15
else:
    print("unacceptable entry, please try again")

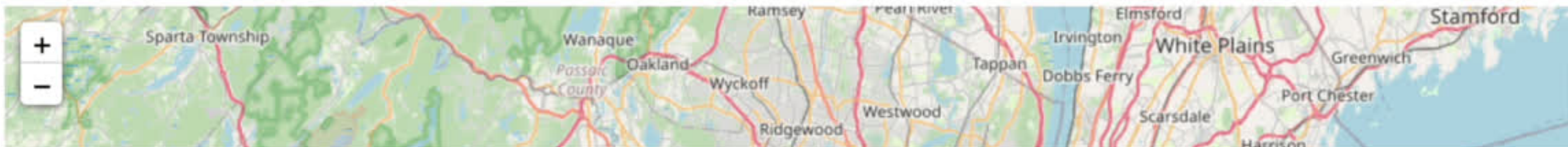
#create death database based on the column number above
death = all_data[all_data.iloc[:,val] > 0]
death

#iteratively create markers in folium map, with lat/long and contributing factor vehicle as popup
for i in range(len(death['ZIP CODE'])):
    folium.Marker(
        location=[death['LATITUDE'].iloc[i],death['LONGITUDE'].iloc[i]],
        popup=str(death['CONTRIBUTING FACTOR VEHICLE 1'].iloc[i]),
        icon=folium.Icon(color="red", icon="info-sign"),
    ).add_to(m)
m
```

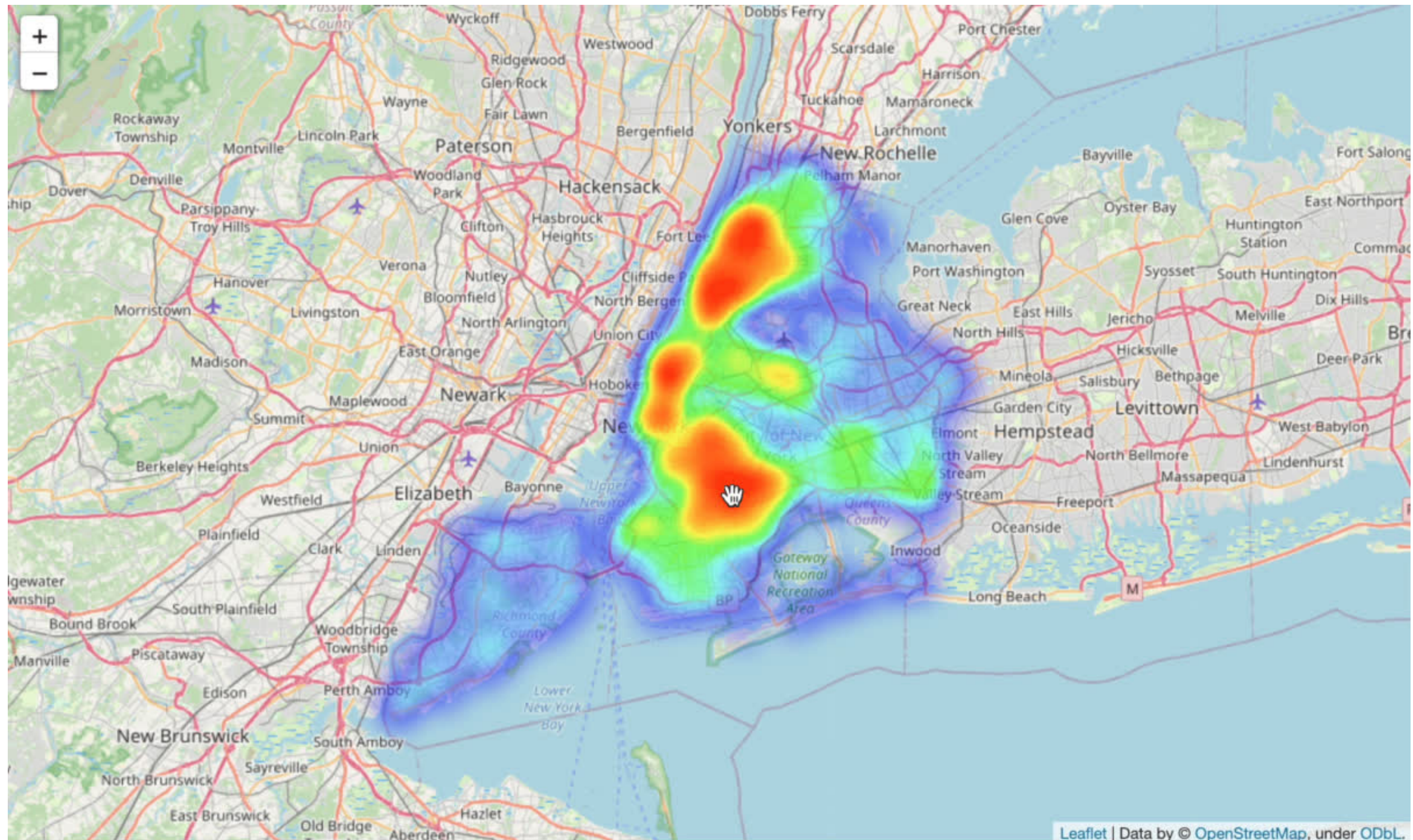
Enter one of the following to see a map of the death:

all persons
pedestrians
cyclists
motorists
motorists

Out[6]:

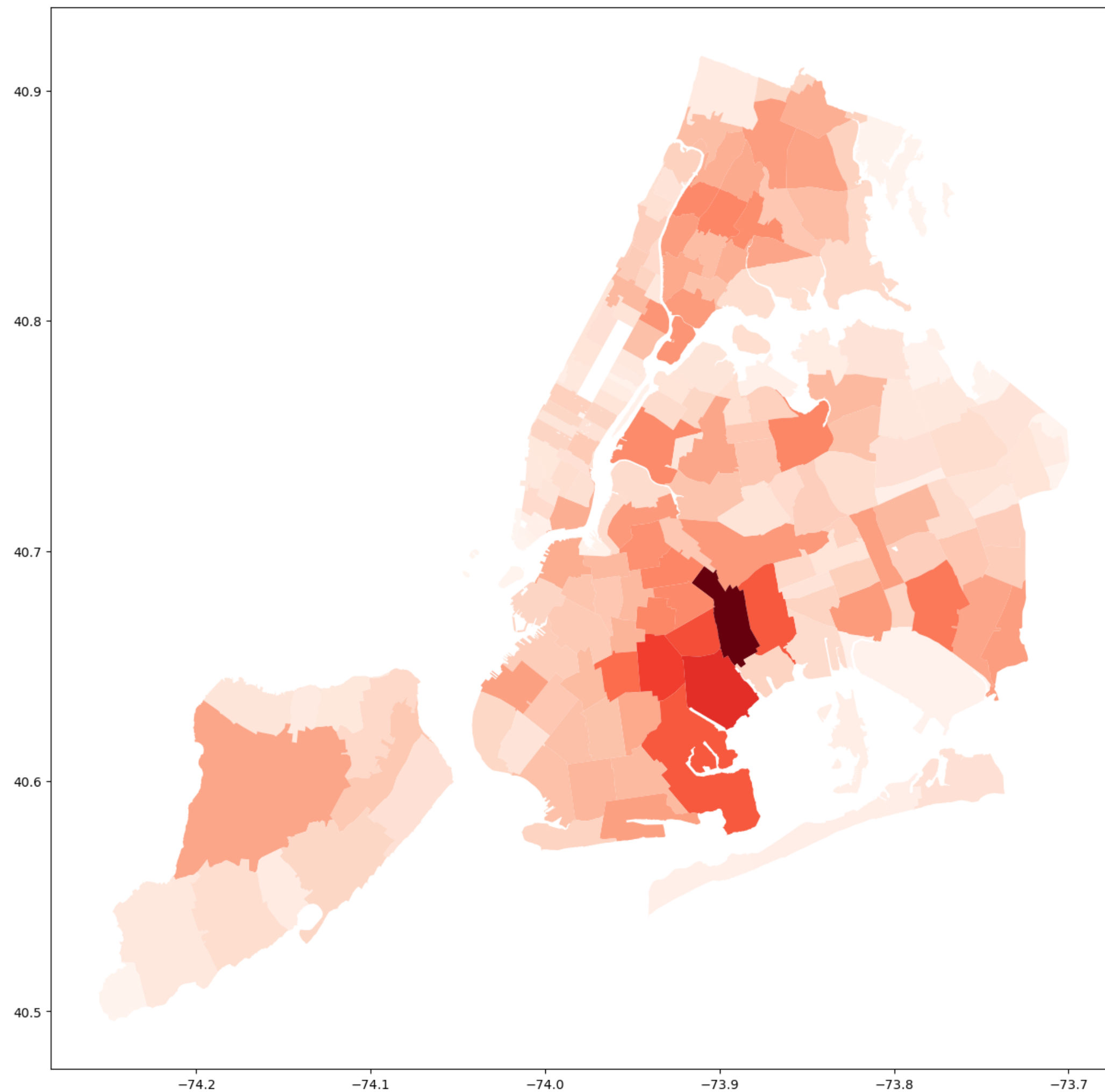


All Injuries Heat Map

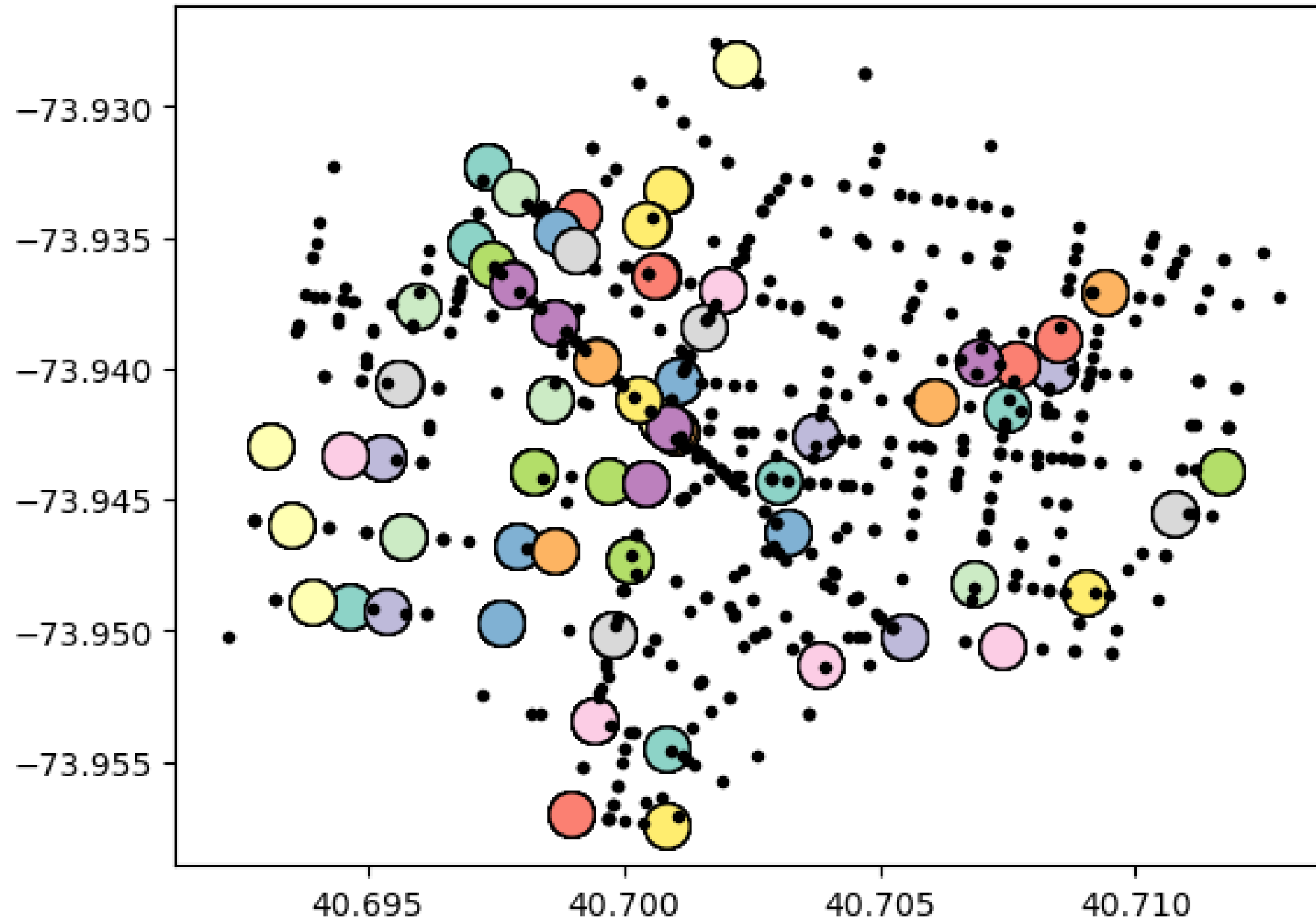


Map of Zipcodes Colorized by Injuries

15



Estimated number of clusters: 61



Plot of
Clusters
Found Using
DBSCAN
Zipcode:11206

Cluster Markers Map for 11206

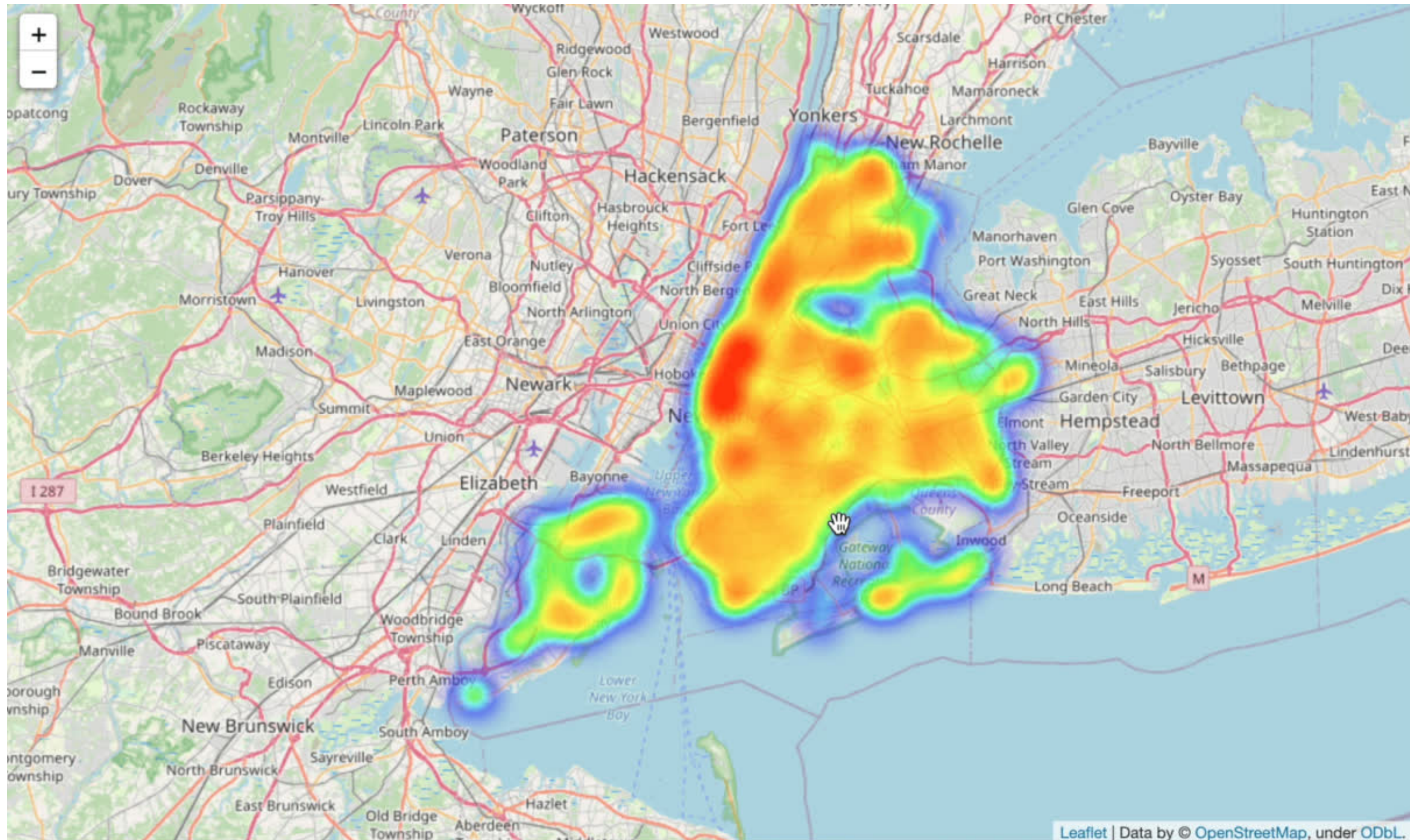


Which areas had the most bike accidents in NYC in 2021?

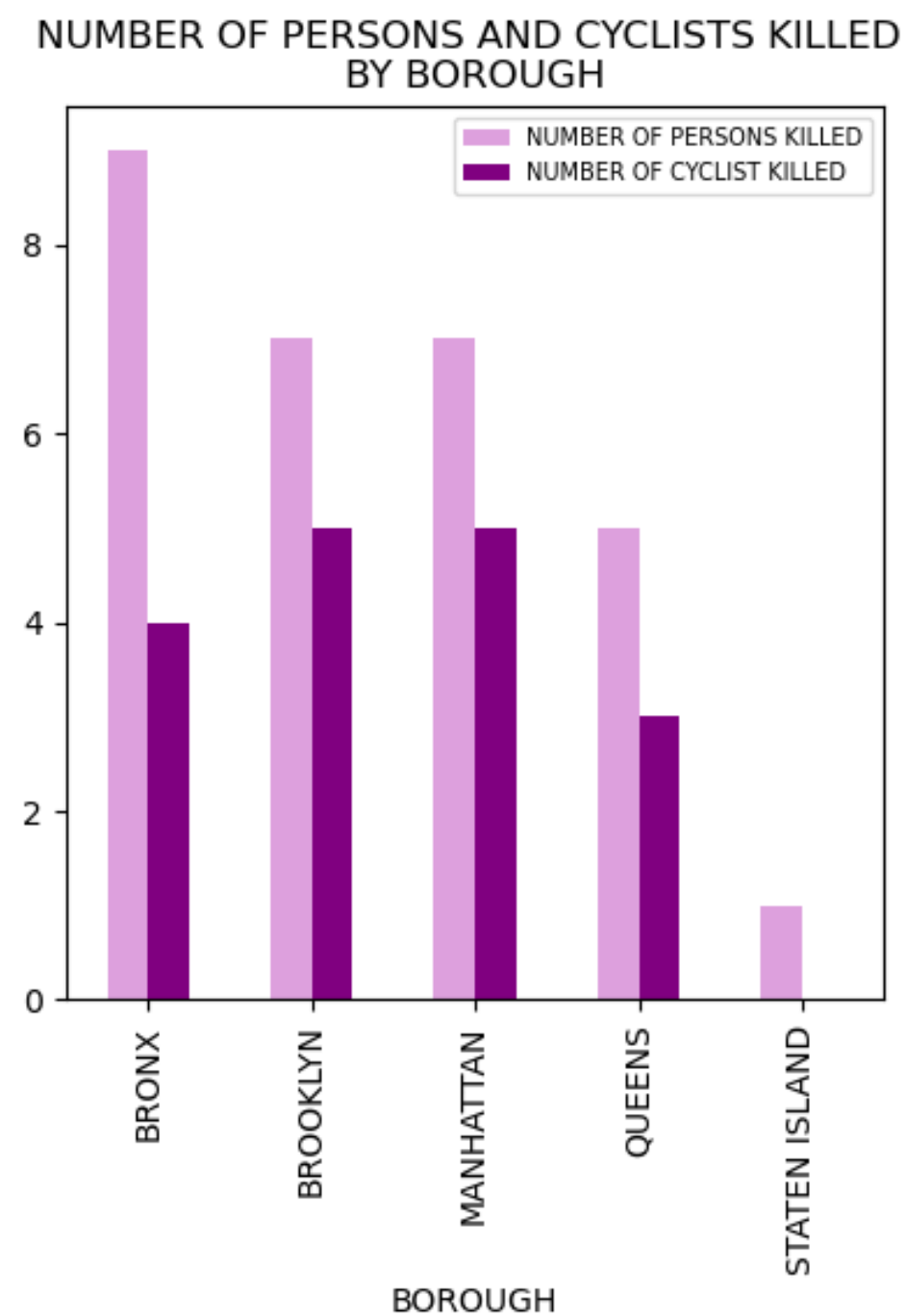
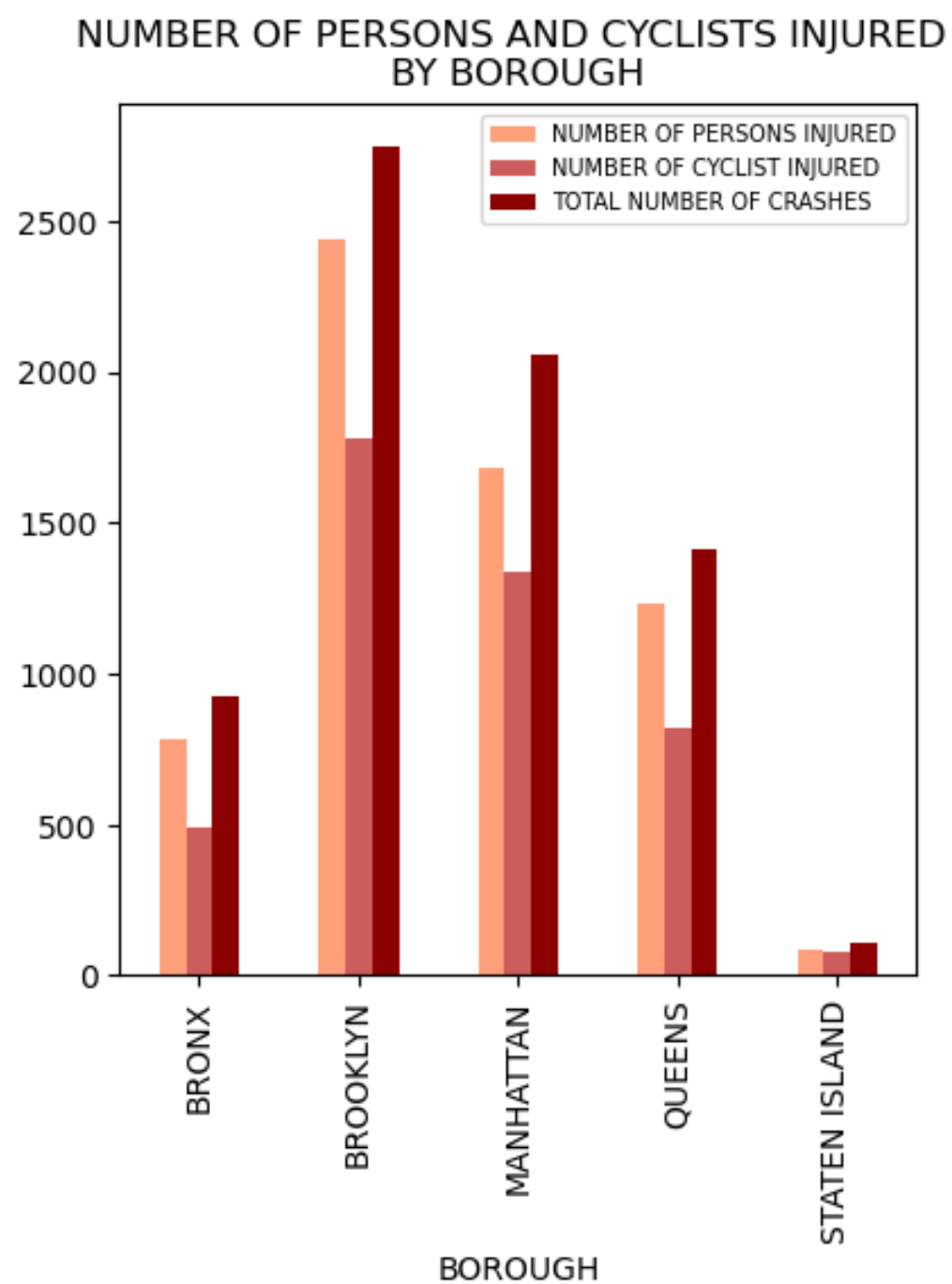
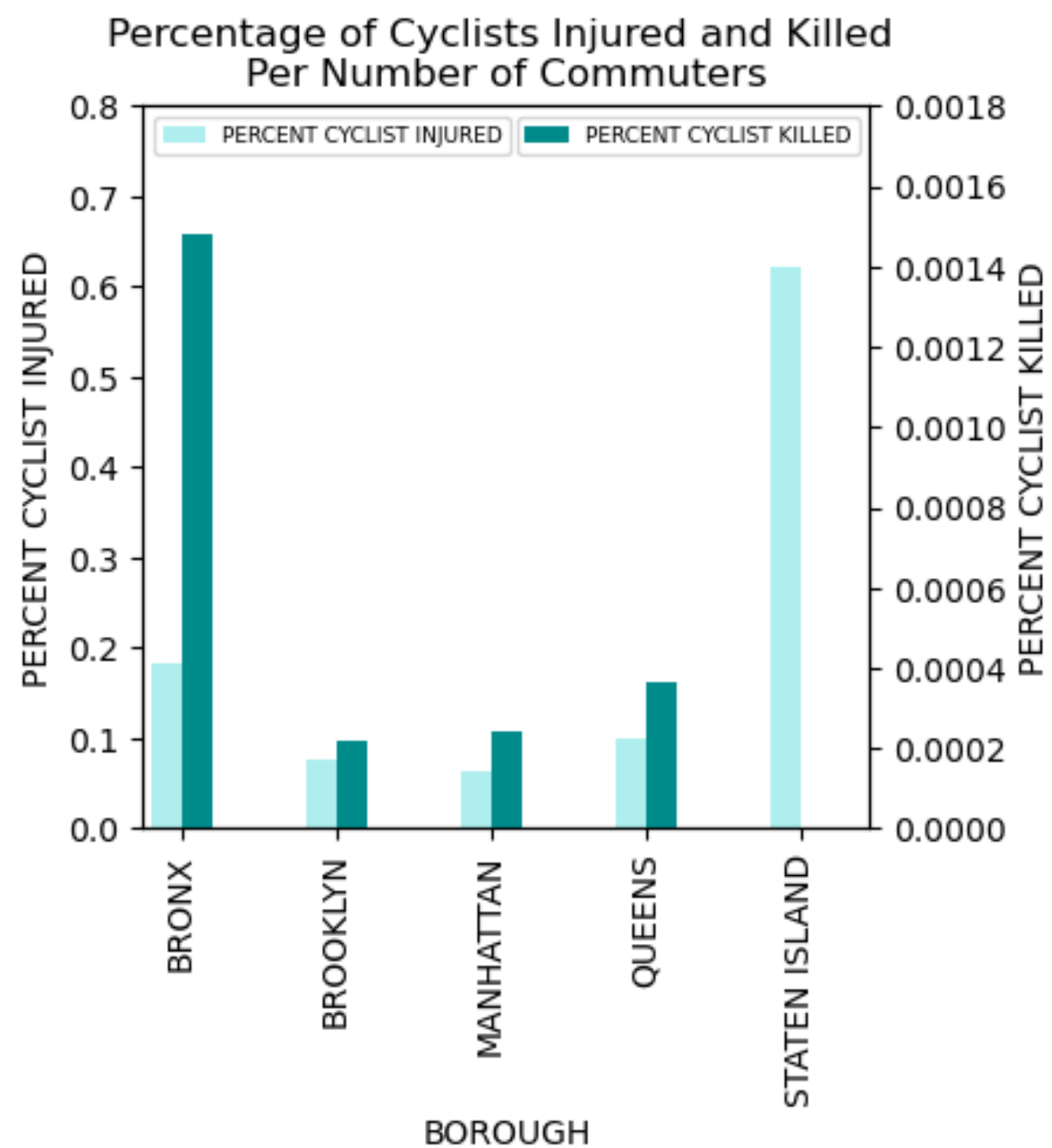
Methods: Divided data by borough and zip code. Included data on bike commuters in each borough to estimate percentages of injuries and fatalities. Used folium to visualize hot spots of collisions with heat map.



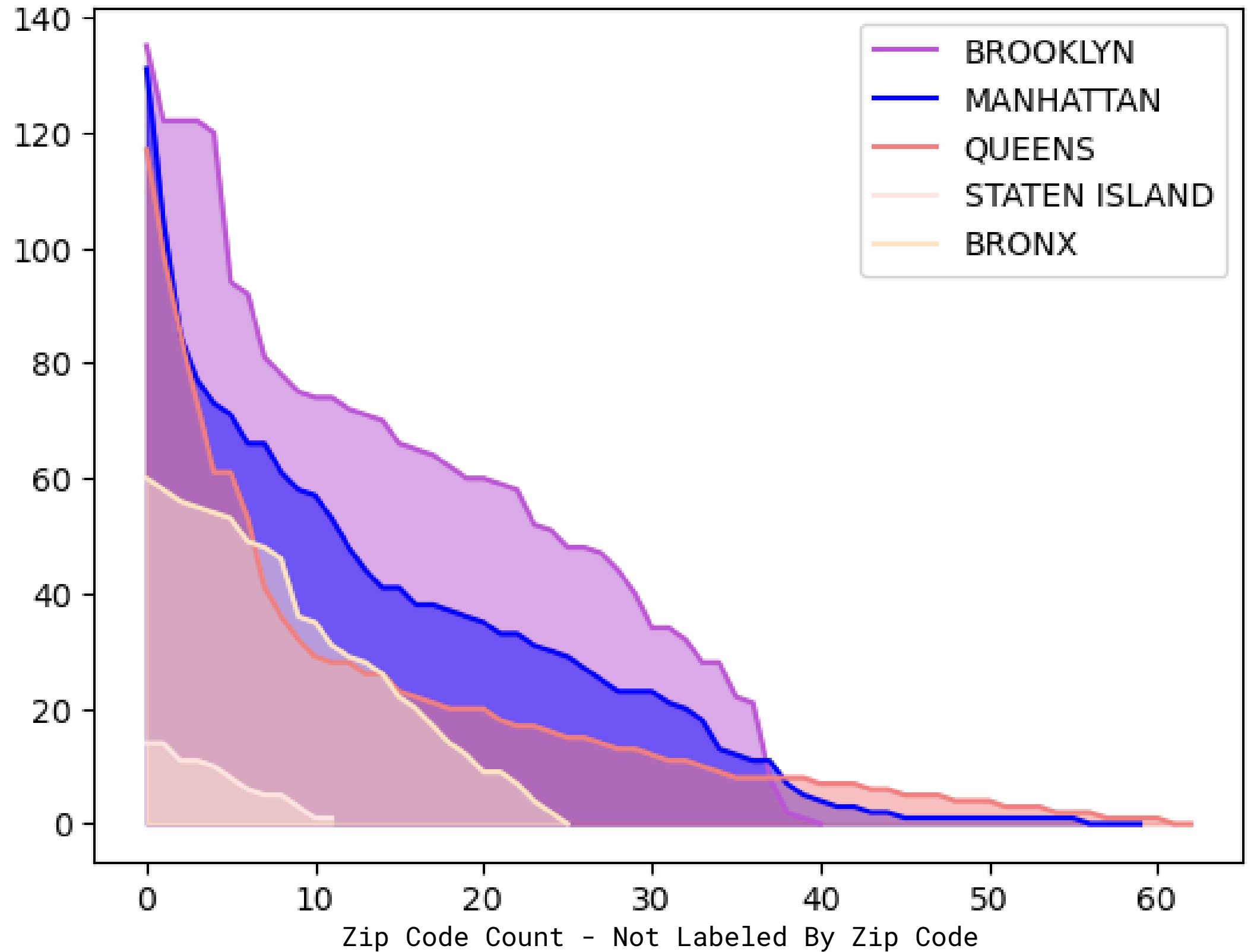
All Crashes Involving Bikes Heat Map



Injuries and Deaths by Borough for Bike Collisions

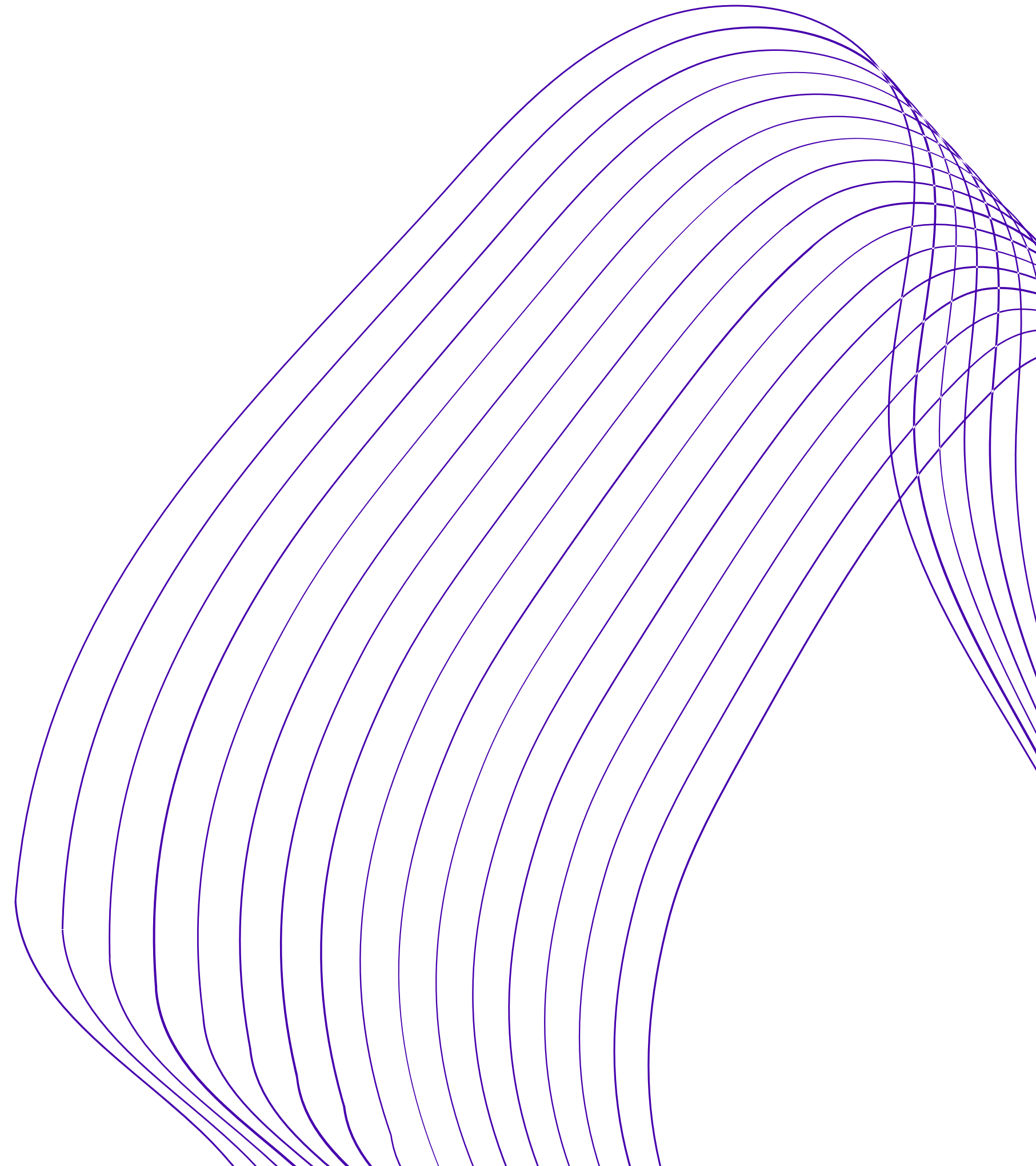


Persons Injured by Zipcode for Bike Collisions



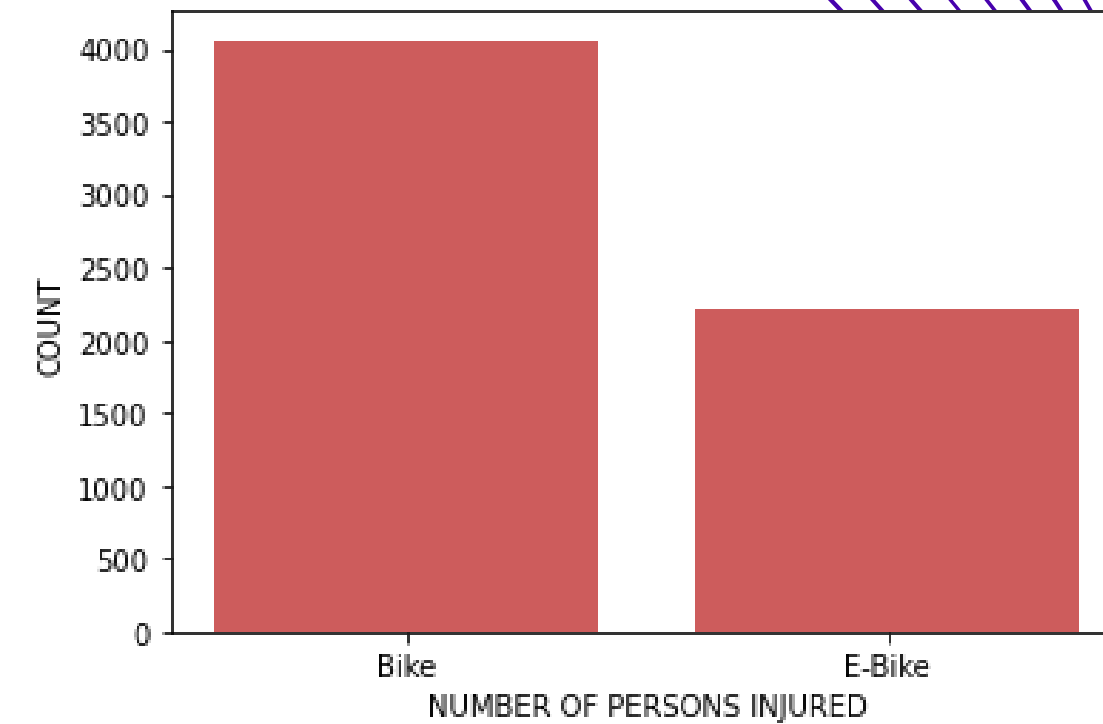
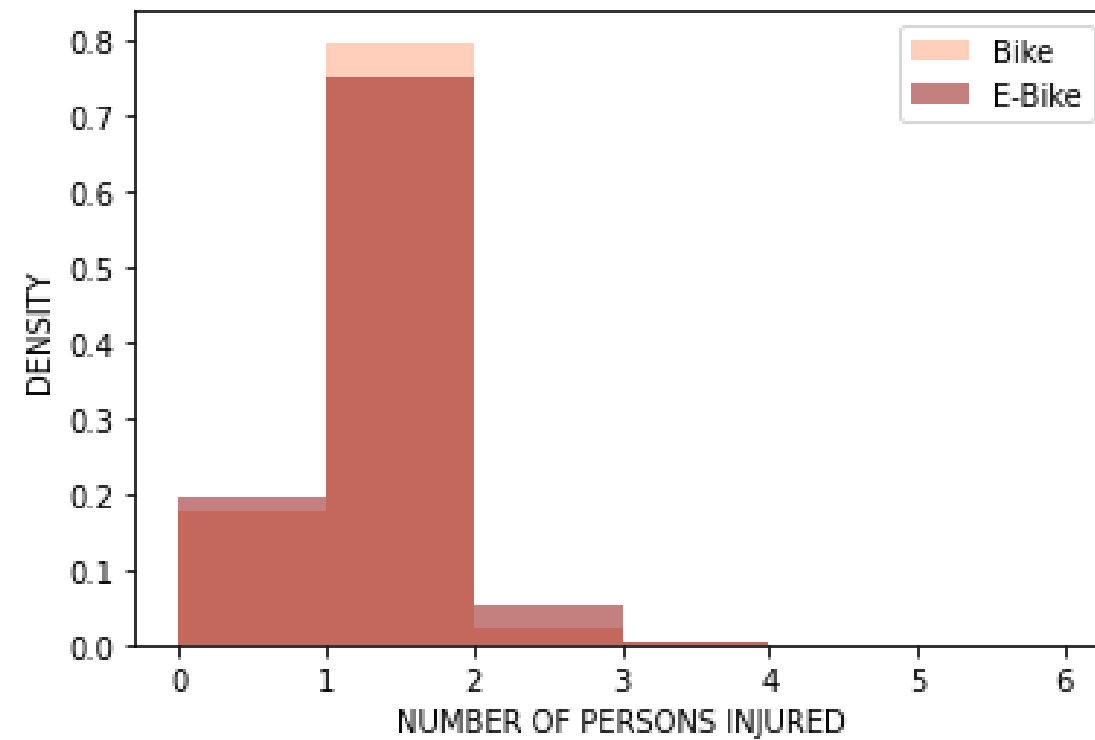
Are e-bikes more dangerous than regular bikes?

Method: Extracted data on injuries and deaths for bikes and e-bikes only. Separated the groups and compared them using bar charts and resampling methods.

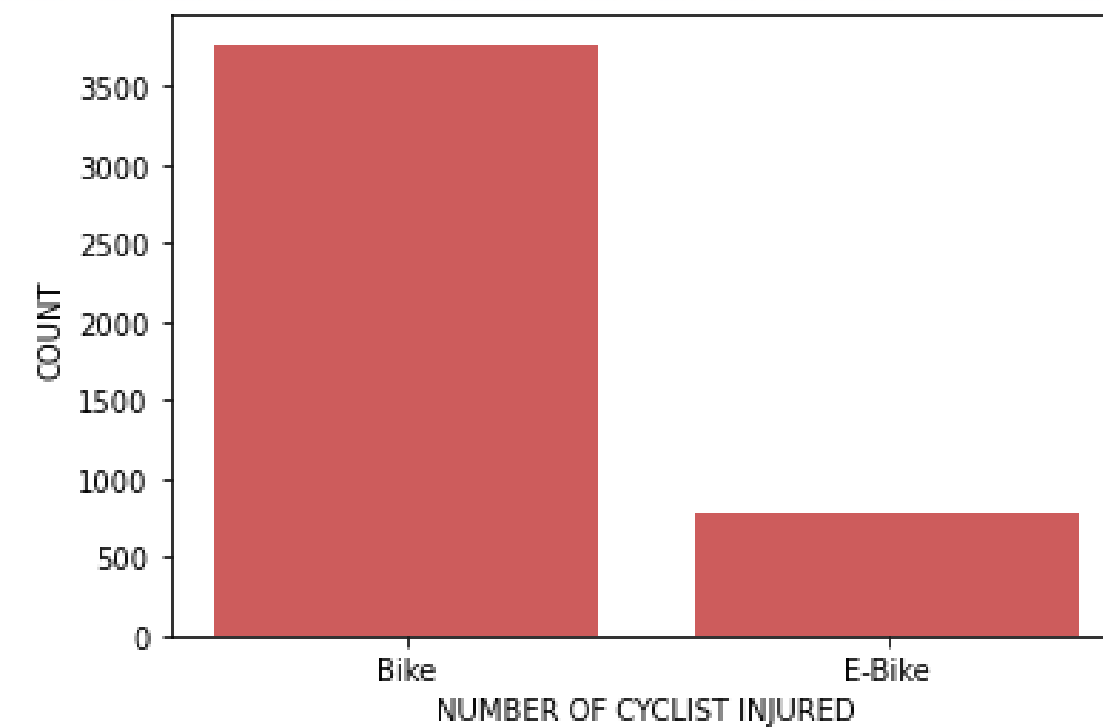
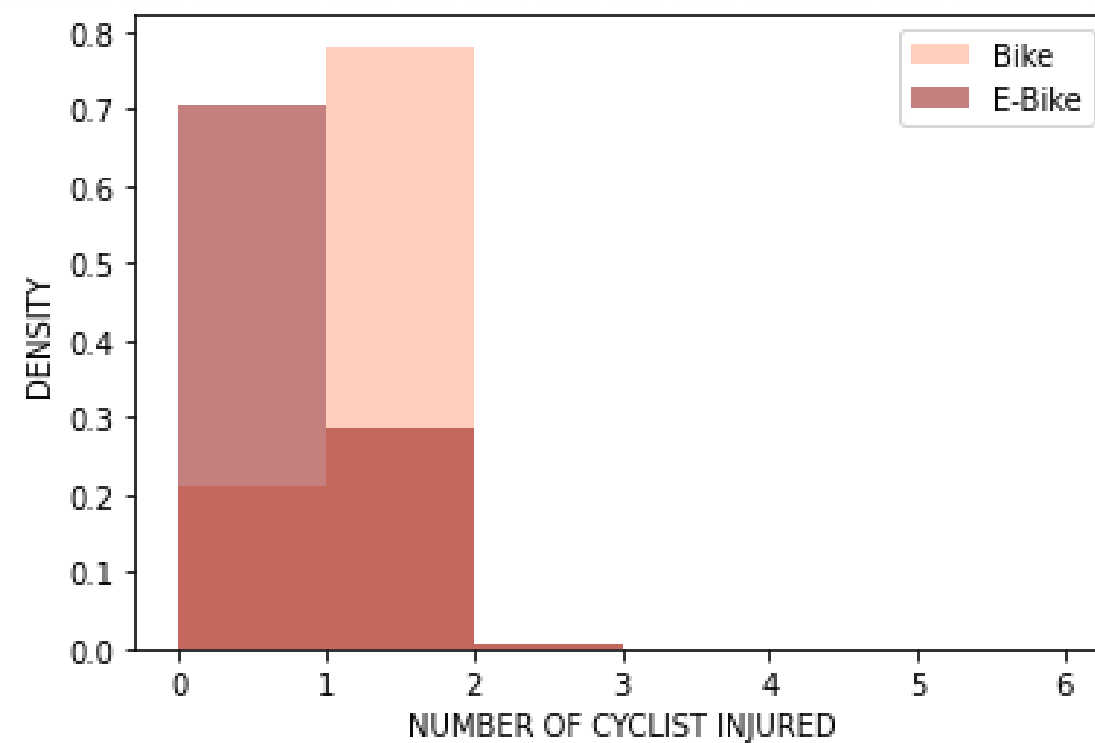


Number Injured

24



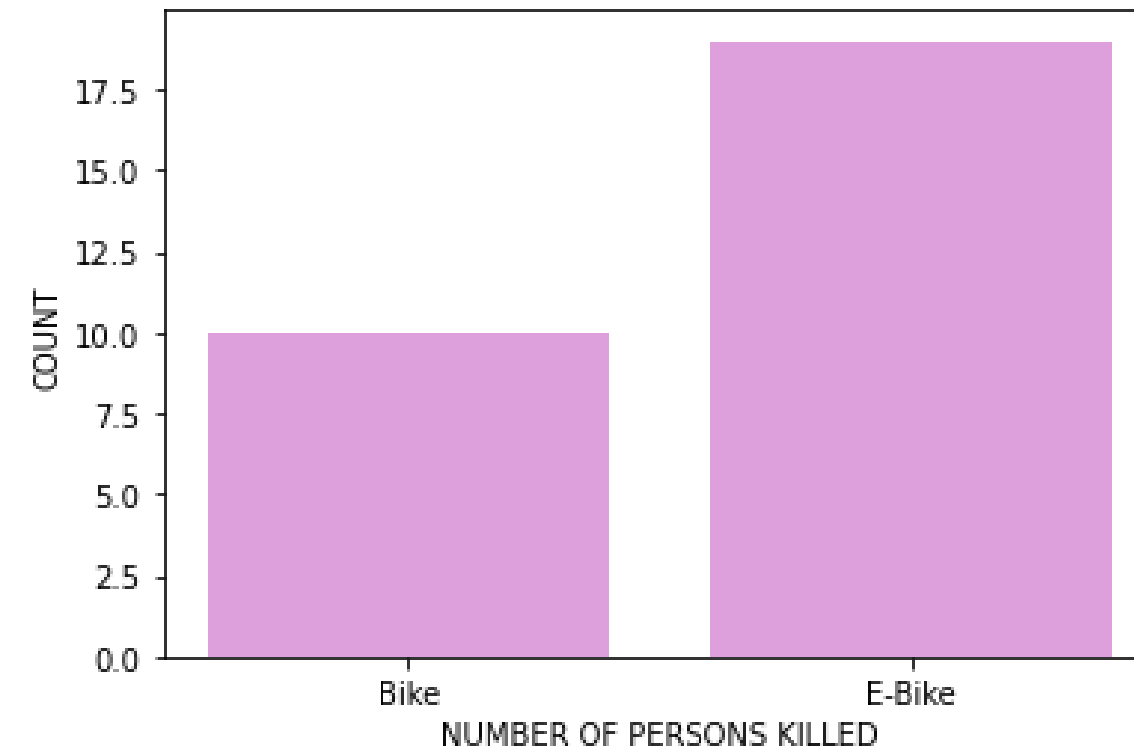
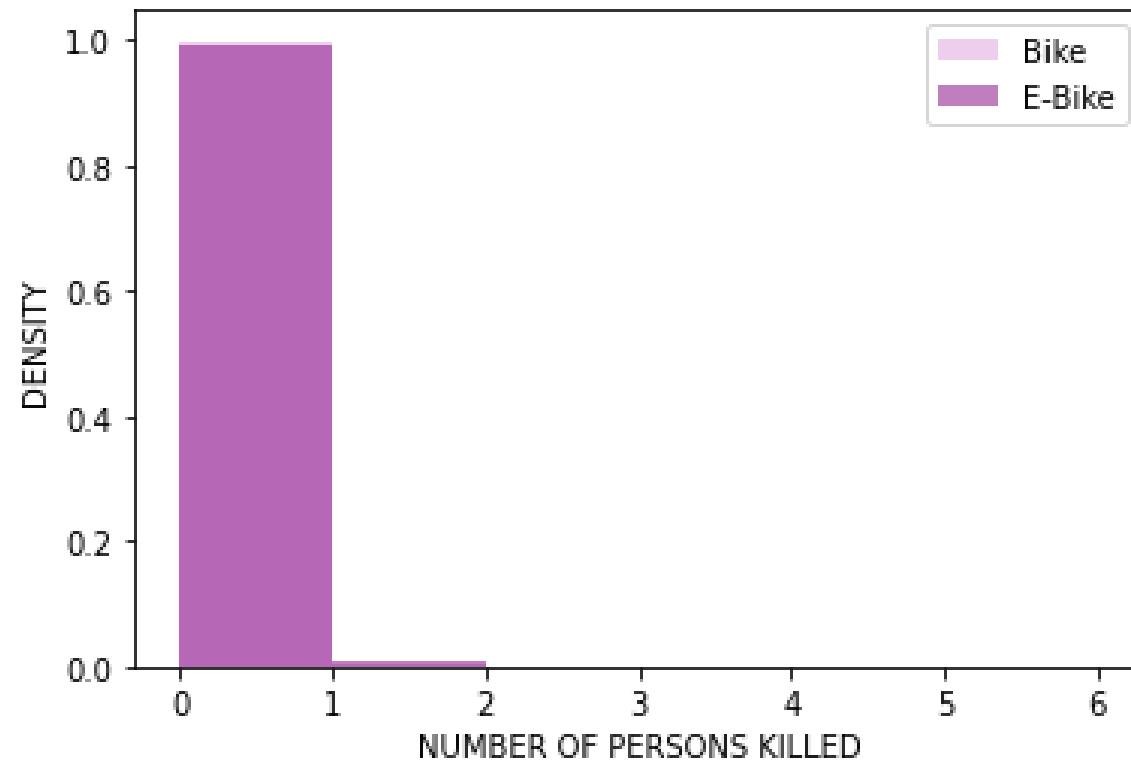
NUMBER OF PERSONS INJURED Welch's t-test p-value: 0.48282999332773846



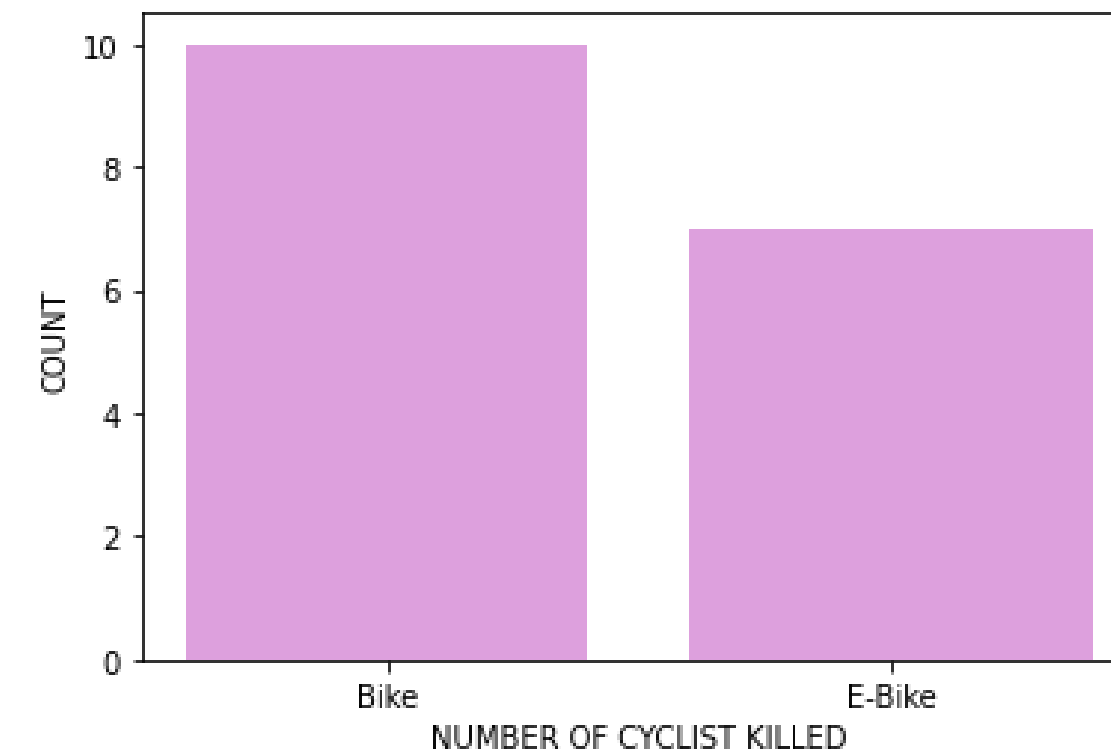
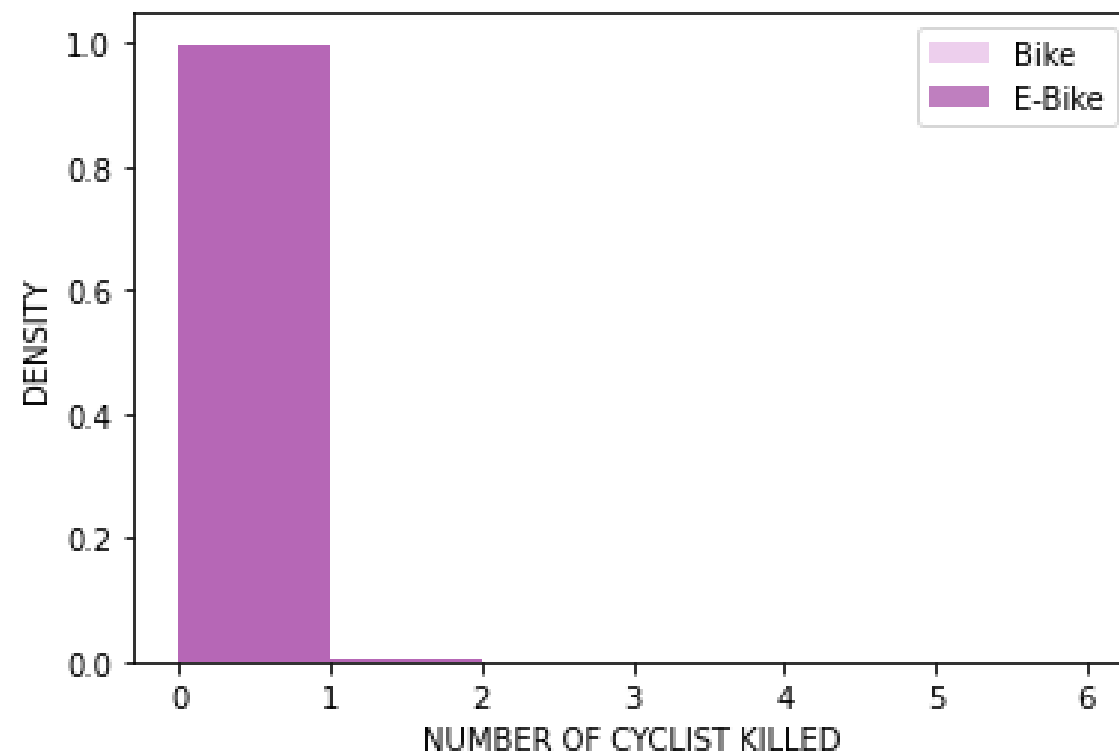
NUMBER OF CYCLIST INJURED Welch's t-test p-value: 0.0

Number Killed

25



NUMBER OF PERSONS KILLED Welch's t-test p-value: 0.005328408716739479



NUMBER OF CYCLIST KILLED Welch's t-test p-value: 0.6116547798822987

Under and Over Sampling Analysis

NUMBER OF PERSONS INJURED:

Undersample count: 3

Undersample significance ratio: 0.003

Oversample count: 16

Oversample significance ratio: 0.016

NUMBER OF PERSONS KILLED:

Undersample count: 981

Undersample significance ratio: 0.981

Oversample count: 1000

Oversample significance ratio: 1.0

NUMBER OF CYCLIST INJURED:

Undersample count: 1000

Undersample significance ratio: 1.0

Oversample count: 1000

Oversample significance ratio: 1.0

NUMBER OF CYCLIST KILLED:

Undersample count: 7

Undersample significance ratio: 0.007

Oversample count: 2

Oversample significance ratio: 0.002

100% of the time there is a significant difference between cyclist injuries for bikes and e-bikes involved in motor vehicle collisions.

Bikes tend to be significantly more 'dangerous' - characterized as having more cyclist injuries per crash - than e-bikes.

Most of the time there is a significant difference between total people killed for bikes versus e-bikes.

E-Bikes tend to be more 'dangerous' - characterized as more total deaths occurring when one is involved in a crash - than bikes.

Thank you!

Questions?