

# The Battle of Neighbourhoods

## Introduction

In the last 5 weeks we have learnt how we can use the location data of a particular neighbourhood and draw interesting conclusions from that which can actually solve critical business problems. We have learnt how to use the Foursquare API to get the locational data of a particular location, the exploratory data analysis using python and pandas has helped us to restructure our data and perform useful operations like data cleaning and handling missing values. The folium library has been used to draw the maps of different geolocations as retrieved from the JSON file created due to the data extraction from Foursquare generated URL.

In week 4 and week 5 of our project we need to assign tasks to ourselves ie defining business problems and from the knowledge we achieved from our previous weeks, we need to put this into action to solve the business queries.

## The Task

This project, and associated files, were produced as to meet the final objectives of Coursera's IBM Applied Data Science Capstone Certificate program. The following sections are the objectives in meeting this criterion.

Section 1: Clearly define a problem or an idea of your choice, where you would need to leverage the Foursquare location data to solve or execute. Remember that data science problems always target an audience and are meant to help a group of stakeholders solve a problem, so make sure that you explicitly describe your audience and why they would care about your problem. This submission will eventually become your Introduction/Business Problem section in your final report. So I recommend that you push the report (having your Introduction/Business Problem section only for now) to your Github repository and submit a link to it.

Section 2: Describe the data that you will be using to solve the problem or execute your idea. Remember that you will need to use the Foursquare location data to solve the problem or execute your idea. You can absolutely use other datasets in combination with the Foursquare location data. So make sure that you provide adequate explanation and discussion, with examples, of the data that you will be using, even if it is only Foursquare location data. This submission will eventually become your Data section in your final report. So I recommend that you push the report (having your Data section) to your Github repository and submit a link to it.

# 1.Problem Description and Background discussion

*Determining suitable locations to set up quarantine centres in New York city due to massive outbreak of COVID-19*

## 1.1.Background:

Due to the outbreak of the highly contagious and deadly corona(COVID-19) virus the world is witnessing deaths ,panics and health hazard.The virus is showing no sign to stop and by the time the vaccine comes out the only way to protect ourselves is to make social distancing among ourselves.The government has set up quarantine centres in New York city to keep the virus affected patients and those having travel history in recent past for 14 days in those locations.

However,the disease being highly contagious, it is very important to set up the quarantine stations in those places where the population is low and and the health facilities like hospitals are also near to those stations.

The neighbourhood data with the geo locations will give us data about the locations where there is a possibility of mass gathering.So when the government decides to start the unlock phase it will be easier for them to decide what places to open and what not to avoid the mass gathering.

## 1.2 Business Problem:

So the main idea behind this project is two-fold. First, it's to analyse the population data of New York city to check what places of the city is less populous so that government can set quarantine centres in those locations and second is to identify the places where there could be possibility of mass gathering so that those places could be taken under surveillance.

## 2. Data Sources

### 1.Population Data of US

To analyse the population data of the New York city we need to analyse the a data set containing the post codes of the city and the population. From this dataset available in Kaggle we get the population data of the year 2010 for the whole country.

LINK: <http://zipatlas.com/us/ny/zip-code-comparison/population-density.htm>

## **2. Post codes of New York**

Now from the population data of the whole United States we need to determine the population of the Zip codes of various locations of New York. So we merge the previous data frame with this data frame on the postcodes.

Link: <https://www.zipcodestogo.com/New%20York/>

## **3.Postcode Geolocation**

Now to access the neighbourhood we need to get the geolocations of the different postcodes of New York city. The data set that gives the geolocations and the postcodes is given below.

Link : [https://cocl.us/new\\_york\\_dataset](https://cocl.us/new_york_dataset)

## **4.Foursquare API data**

Now to access the possible locations where mass gathering can take place, we use the Foursquare API data with our own client-ID. We take a given radius and set the latitude and longitude of New York to get the locations of various points of gathering. The data can be derived from the Foursquare's website.

Link: <https://foursquare.com/>

# Approach

## 1.Importing Important Libraries

```
1 # Importing Important Libraries
2
3 import pandas as pd
4 import numpy as np
5 import requests
6 import json
7 from bs4 import BeautifulSoup
8 import folium
9 from geopy.geocoders import Nominatim
10 import matplotlib.pyplot as plt
11 import matplotlib.cm as cm
12 import matplotlib.colors as colors
13 %matplotlib inline
14
```

## 2.Web Scrapping and Population Dataframe creation

```
1 # Web Scrapping using BeautifulSoup
2
3 with open('NY_Population.html') as html_file:
4     soup=BeautifulSoup(html_file,'lxml')
5
6 table = soup.find('table')
7 table
8 population_df = pd.read_html(str(table))[0]
9 population_df.drop(population_df.index[0:15],inplace=True)
10 population_df=population_df.reset_index(drop=True)
11 population_df.drop(columns=[0,2,3,5,6],inplace=True)
12 population_df=population_df.rename(columns={1:'zipcode',4:'population'})
13 population_df.head()
14 #population_df.shape
15
```

|   | zipcode | population |
|---|---------|------------|
| 0 | 10162   | 1726       |
| 1 | 10028   | 44987      |
| 2 | 10128   | 59856      |
| 3 | 10031   | 60221      |
| 4 | 10009   | 58595      |

### 3.Web Scrapping and Geotag Dataframe creation

```
In [311]: 1 # Extracting Geo location of the neighbourhood datas and Borough and making dataframes
2
3 with open('nyu_2451_34572-geojson.json') as json_data:
4     newyork_data = json.load(json_data)
5     neighborhoods_data = newyork_data['features']
6     neighborhoods_data
7     neighborhoods_data[0]
8     # define the dataframe columns
9     column_names = ['Borough', 'Neighborhood', 'Latitude', 'Longitude']
10
11 # instantiate the dataframe
12 neighborhoods = pd.DataFrame(columns=column_names)
13 neighborhoods
14
15 for data in neighborhoods_data:
16     borough = neighborhood_name = data['properties']['borough']
17     neighborhood_name = data['properties']['name']
18
19     neighborhood_latlon = data['geometry']['coordinates']
20     neighborhood_lat = neighborhood_latlon[1]
21     neighborhood_lon = neighborhood_latlon[0]
22
23     neighborhoods = neighborhoods.append({'Borough': borough,
24                                           'Neighborhood': neighborhood_name,
25                                           'Latitude': neighborhood_lat,
26                                           'Longitude': neighborhood_lon}, ignore_index=True)
27 neighborhoods.head(3)
28 #neighborhoods.shape
```

```
Out[311]:
```

|   | Borough | Neighborhood | Latitude  | Longitude  |
|---|---------|--------------|-----------|------------|
| 0 | Bronx   | Wakefield    | 40.894705 | -73.847201 |
| 1 | Bronx   | Co-op City   | 40.874294 | -73.829939 |
| 2 | Bronx   | Eastchester  | 40.887556 | -73.827806 |

### 4.Merging 2 Dataframes and Missing value handling

```
In [304]: 1 df = pd.merge(df1,neighborhoods,on="Borough", how='left')           #Merging the 2 dataframes from left
2 df.head()
```

```
Out[304]:
```

|   | zipcode | Borough  | population | Neighborhood | Latitude | Longitude |
|---|---------|----------|------------|--------------|----------|-----------|
| 0 | 10001   | New York | 17310      | NaN          | NaN      | NaN       |
| 1 | 10002   | New York | 84870      | NaN          | NaN      | NaN       |
| 2 | 10003   | New York | 53673      | NaN          | NaN      | NaN       |
| 3 | 10005   | New York | 884        | NaN          | NaN      | NaN       |
| 4 | 10006   | New York | 1447       | NaN          | NaN      | NaN       |

```
In [305]: 1 # Handling Missin values
2 # Dropping those rows having Missing Values
3
4 df.isnull().sum()
5 df.dropna( axis=0, how='any',inplace=True)
6 df = df.reset_index(drop=True)           #Resetting Index
```

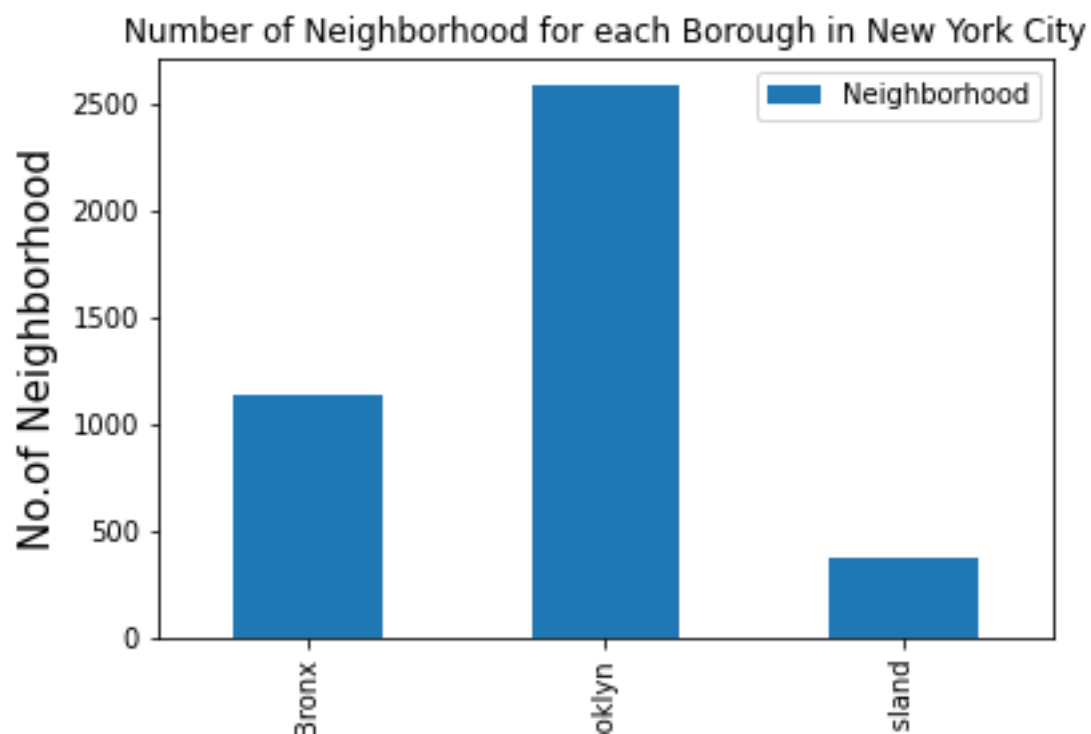
```
In [306]: 1 df.head()
```

```
Out[306]:
```

|   | zipcode | Borough       | population | Neighborhood  | Latitude  | Longitude  |
|---|---------|---------------|------------|---------------|-----------|------------|
| 0 | 10301   | Staten Island | 38805      | St. George    | 40.644982 | -74.079353 |
| 1 | 10301   | Staten Island | 38805      | New Brighton  | 40.640615 | -74.087017 |
| 2 | 10301   | Staten Island | 38805      | Stapleton     | 40.626928 | -74.077902 |
| 3 | 10301   | Staten Island | 38805      | Rosebank      | 40.615305 | -74.069805 |
| 4 | 10301   | Staten Island | 38805      | West Brighton | 40.631879 | -74.107182 |

## 5. Plotting Number of Neighbourhoods of different Borough

```
In [310]: 1 # Plotting which Borough in NY has the most number of Neighbourhoods
2
3 # title
4 plt.title('Number of Neighborhood for each Borough in New York City')
5 #On x-axis
6 plt.xlabel('Borough', fontsize = 15)
7 #On y-axis
8 plt.ylabel('No. of Neighborhood', fontsize=15)
9 #giving a bar plot
10 df.groupby('Borough')['Neighborhood'].count().plot(kind='bar')
11 #Legend
12 plt.legend()
13 #displays the plot
14 plt.savefig('BoroughBarplot.png')
15 print("The Plot Shows Brooklyn has most Number of Neighbourhoods so the population of Brooklyn is likely to be more. So Spe
16 plt.tight_layout()
17
```



The Plot Shows Brooklyn has most Number of Neighbourhoods so the population of Brooklyn is likely to be more. So Special care needs to be taken to control the corona virus spread.

## 6. Analysing Brooklyn Data

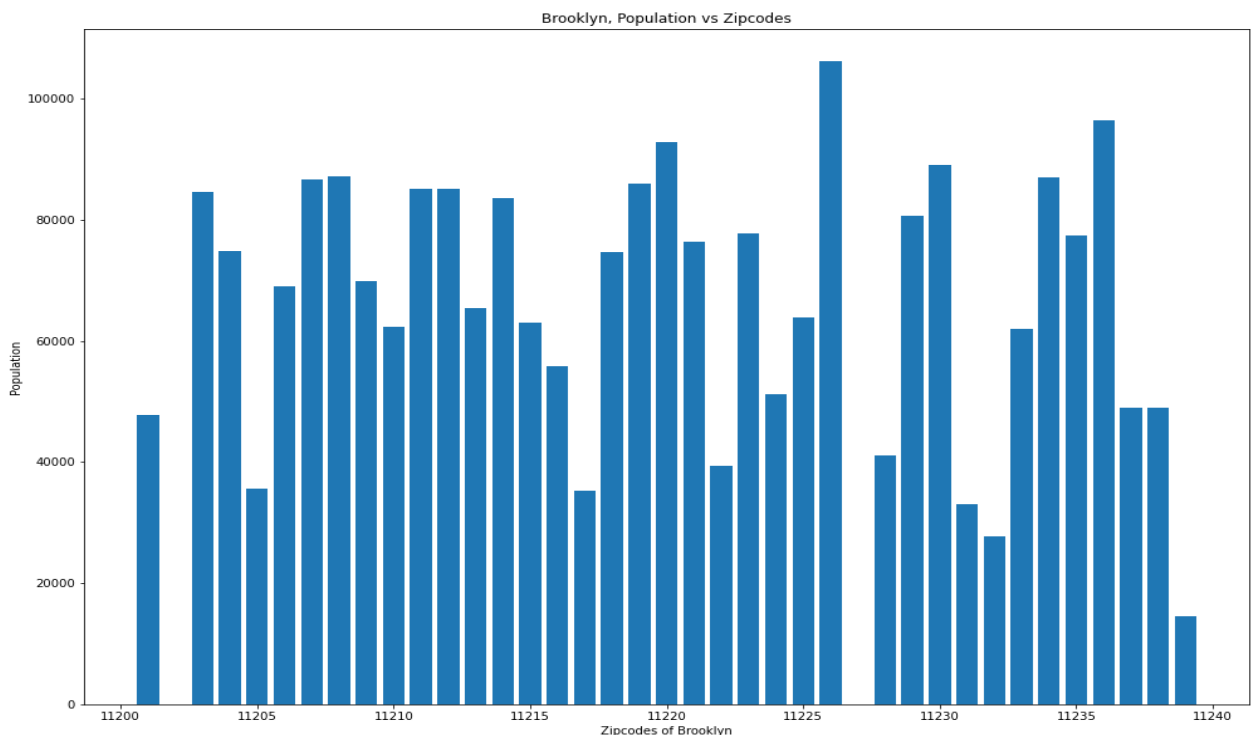
As Brooklyn has most number of Neighborhoods so lets examine the demographic of Brooklyn

```
In [291]: 1 df.population=df.population.astype('int')           # Making population datatype as integer
          2 df.zipcode=df.zipcode.astype('int')           # Making zipcode datatype as integer
          3 filter=df['Borough']=="Brooklyn"               # filter only for Brooklyn Borough
          4 df=df.loc[filter]
          5 df=df.sort_values(by='population',ascending=True) # Sorting population from Low to High
          6 df = df.reset_index(drop=True)
          7 df.head(10)
```

Out[291]:

|   | zipcode | Borough  | population | Neighborhood       | Latitude  | Longitude  |
|---|---------|----------|------------|--------------------|-----------|------------|
| 0 | 11239   | Brooklyn | 14620      | Erasmus            | 40.646926 | -73.948177 |
| 1 | 11239   | Brooklyn | 14620      | Brooklyn Heights   | 40.695864 | -73.993782 |
| 2 | 11239   | Brooklyn | 14620      | Cobble Hill        | 40.687920 | -73.998561 |
| 3 | 11239   | Brooklyn | 14620      | Carroll Gardens    | 40.680540 | -73.994654 |
| 4 | 11239   | Brooklyn | 14620      | Red Hook           | 40.676253 | -74.012759 |
| 5 | 11239   | Brooklyn | 14620      | Gowanus            | 40.673931 | -73.994441 |
| 6 | 11239   | Brooklyn | 14620      | Fort Greene        | 40.688527 | -73.972906 |
| 7 | 11239   | Brooklyn | 14620      | Bedford Stuyvesant | 40.687232 | -73.941785 |
| 8 | 11239   | Brooklyn | 14620      | Park Slope         | 40.672321 | -73.977050 |
| 9 | 11239   | Brooklyn | 14620      | East New York      | 40.669926 | -73.880699 |

```
In [292]: 1 # Plotting Zipcode Vs Population to get an idea about the Population distribution in different places
          2
          3 x=df.zipcode
          4 y=df.population
          5 plt.figure(figsize=(14,10))
          6 plt.xlabel('Zipcodes of Brooklyn')
          7 plt.ylabel('Population')
          8 plt.title("Brooklyn, Population vs Zipcodes")
          9 plt.bar(x,y)
         10 plt.tight_layout()
         11 plt.savefig('PopulationVsZipBarplot.png') #Saving The Image
```



## 7. Brooklyn Neighbourhood In map

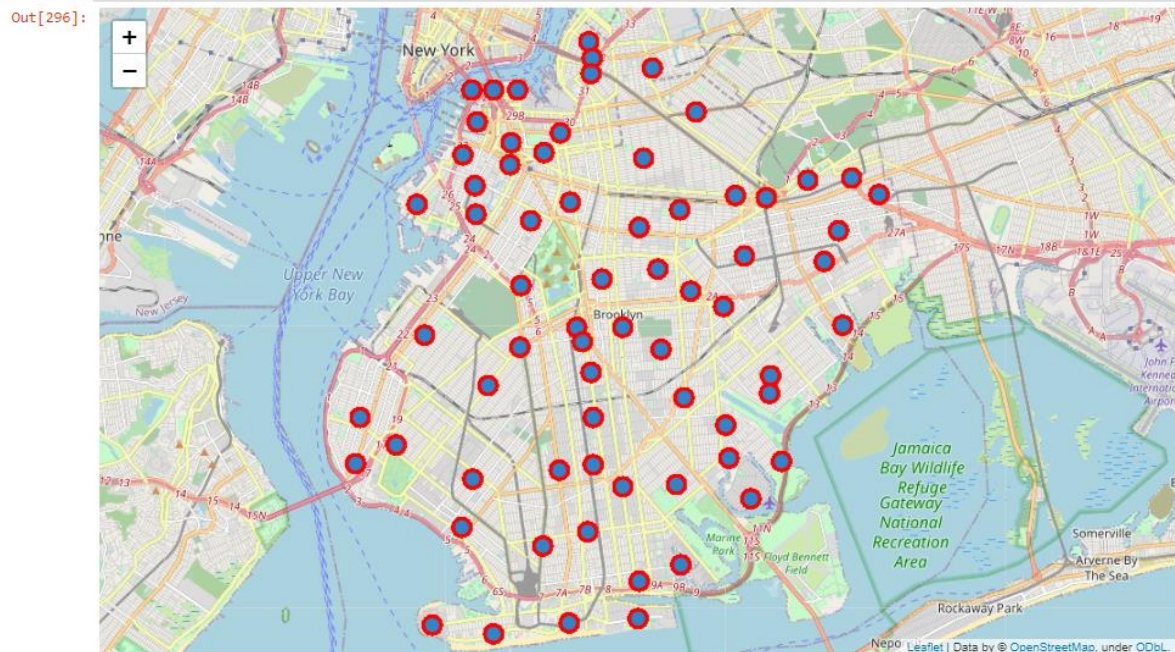
```
In [295]: 1 address = 'Brooklyn, New York, USA' # Address of the Location which is to be examined
2
3 geolocator = Nominatim(user_agent="NY_Traveller")
4 location = geolocator.geocode(address) # To access the Location of the given adress
5 latitude = location.latitude # To access the Latitude of Brooklyn, NY, USA
6 longitude = location.longitude # To access the Longitude of Brooklyn, NY, USA
7
8 print('The geograpical coordinate of NY is : {}, {}'.format(latitude, longitude))
```

The geograpical coordinate of NY is : 40.6501038, -73.9495823.

Observing the Map of Brooklyn

```
In [296]: 1 # create map of Toronto using Latitude and Longitude values
2 map_NY = folium.Map(location=[latitude, longitude], zoom_start=12)
3
4 # add markers to map
5 for lat, lng, borough, neighborhood, population in zip(df['Latitude'], df['Longitude'],
6                                                       df['Borough'], df['Neighborhood'], df['population']):
7     label = '{} , {}'.format(neighborhood, population)
8     label = folium.Popup(label, parse_html=True)
9     folium.CircleMarker(
10         [lat, lng],
11         radius=8,
12         popup=label,
13         color='red',
14         fill=True,
15         fill_color='#3186cc',
16         fill_opacity=0.6,
17         parse_html=False).add_to(map_NY)
18
19 map_NY
```

## Brooklyn Neighbourhood Map





## 8.Foursquare API and URL making

```
In [265]: 1 radius=100000
2 LIMIT=10000
3 neighborhood_latitude=Df.loc[0,"Latitude"]
4 neighborhood_longitude=Df.loc[0,"Longitude"]
5 neighbourhood_name = Df.loc[0, 'Neighborhood']
6
7
8 url = 'https://api.foursquare.com/v2/venues/explore?&client_id={}&client_secret={}&v={}&ll={},{}&radius={}&limit={}'.format(
9     CLIENT_ID,
10    CLIENT_SECRET,
11    VERSION,
12    neighborhood_latitude,
13    neighborhood_longitude,
14    radius,
15    LIMIT)
16
17 print("The Foursquare URL to extract the Erasmus Neighbourhood data is :",url)
```

The Foursquare URL to extract the Erasmus Neighbourhood data is : [https://api.foursquare.com/v2/venues/explore?&client\\_id=LPYCJIALA3DLB8T31E4ES13IPEWMC0YK5U3WBGXWBNQW43AF&client\\_secret=2TRPJHGRCOFPWS1QVERZB1JXZ5DTBYOH0HLDKN2U55XXVPLQ&v=20180604&ll=40.70317632822692,-73.9887528074504&radius=100000&limit=10000](https://api.foursquare.com/v2/venues/explore?&client_id=LPYCJIALA3DLB8T31E4ES13IPEWMC0YK5U3WBGXWBNQW43AF&client_secret=2TRPJHGRCOFPWS1QVERZB1JXZ5DTBYOH0HLDKN2U55XXVPLQ&v=20180604&ll=40.70317632822692,-73.9887528074504&radius=100000&limit=10000)

```
In [266]: 1 #Making Json file from the url made
2 results = requests.get(url).json()
3 #results
4 def get_category_type(row):
5     try:
6         categories_list = row['categories']
7     except:
8         categories_list = row['venue.categories']
9
10    if len(categories_list) == 0:
11        return None
12    else:
13        return categories_list[0]['name']
```

```
In [267]: 1 venues = results['response']['groups'][0]['items']
2 nearby_venues = pd.json_normalize(venues) # flatten JSON
3
4 # filter columns
5 filtered_columns = ['venue.name', 'venue.categories', 'venue.location.lat', 'venue.location.lng']
6 nearby_venues = nearby_venues.loc[:, filtered_columns]
7
8 # filter the category for each row
9 nearby_venues['venue.categories'] = nearby_venues.apply(get_category_type, axis=1)
10
11 # clean columns
12 nearby_venues.columns = [col.split(".")[0] for col in nearby_venues.columns]
13 nearby_venues["Neighbourhood"] = df.Neighborhood
14 nearby_venues["Neighbourhood_lat"] = df.Latitude
15 nearby_venues["Neighbourhood_lng"] = df.Longitude
16 nearby_venues.rename(columns={"lat": "Venue_lat", "lng": "Venue_lng", "name": "Venue"}, inplace=True)
17 nearby_venues.head()
```

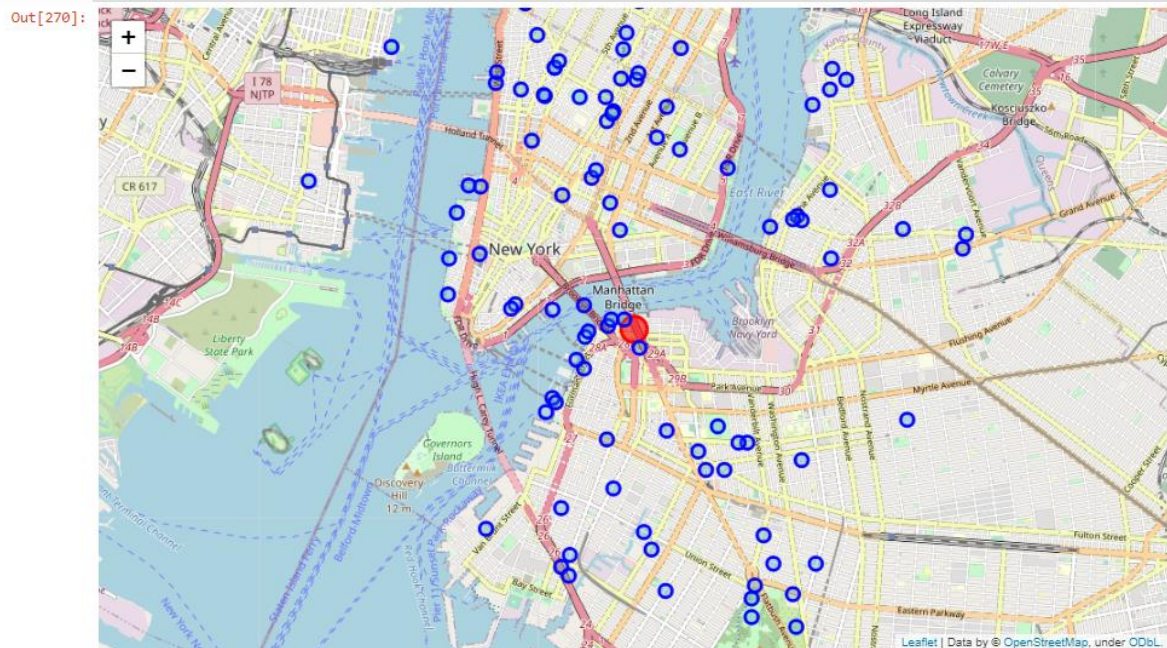
```
Out[267]:
```

|   | Venue                         | categories     | Venue_lat | Venue_lng  | Neighbourhood    | Neighbourhood_lat | Neighbourhood_lng |
|---|-------------------------------|----------------|-----------|------------|------------------|-------------------|-------------------|
| 0 | Pebble Beach                  | Beach          | 40.704329 | -73.990265 | Erasmus          | 40.646926         | -73.948177        |
| 1 | Brooklyn Bridge Park          | Park           | 40.702282 | -73.996456 | Brooklyn Heights | 40.695864         | -73.993782        |
| 2 | Brooklyn Heights Promenade    | Scenic Lookout | 40.698462 | -73.996707 | Cobble Hill      | 40.687920         | -73.998561        |
| 3 | Brooklyn Bridge               | Bridge         | 40.705967 | -73.996707 | Carroll Gardens  | 40.680540         | -73.994654        |
| 4 | Brooklyn Bridge Park - Pier 1 | Park           | 40.702900 | -73.995987 | Red Hook         | 40.676253         | -74.012759        |

This dataframe gives us the data of various venue location so that we can analyse the venues of the Brooklyn datas.

## 9.Mapping Brooklyn Venues

```
In [270]: 1 # create map of Brooklyn using the selected neighbourhood Latitude and Longitude values
2 map_tohood = folium.Map(location=[neighborhood_latitude, neighborhood_longitude], zoom_start=12.5)
3
4 #add a red circle marker to represent the selected neighborhood
5 folium.CircleMarker(
6     [neighborhood_latitude, neighborhood_longitude],
7     radius=12,
8     color='red',
9     popup=Df["Neighborhood"],
10    fill = True,
11    fill_color = 'red',
12    fill_opacity = 0.6
13).add_to(map_tohood)
14
15
16
17 # add markers to map
18 for lat, lng, name, categories in zip(nearby_venues['Venue_lat'], nearby_venues['Venue_lng'], nearby_venues['Venue'], nearby
19 label = '{}', {}'.format(name, categories)
20 label = folium.Popup(label, parse_html=True)
21 folium.CircleMarker(
22     [lat, lng],
23     radius=6,
24     popup=label,
25     color='blue',
26     fill=True,
27     fill_color='#3199cc',
28     fill_opacity=0.3).add_to(map_tohood)
29
30 map_tohood
```



## Observations and conclusions

1. Among the NY Boroughs Brooklyn is Most Populous
2. Brooklyn has 30 Unique zipcode Locations and population is segregated according to those locations
3. The Plot we created shows the Populations densities in those postcode areas
4. Among the 100 limited observations from Foursquare we see what locations of the city Dumbo has the probability of having public gathering
5. We get an idea of how the venues of distributed in the Unlock period which locations can be opened to avoid mass gathering

## Limitations

- 1.The population data we have used is of 2010 census.So the current population will vary from the actual one
- 2.Foursquare API doesnot give us the locations of hospital and healthcare centres