

## ▼ Capstone Project - The Battle of the Neighborhoods

COVID Risk Categories Applied For Reopening Restaraunts

Assignment by Rafael Scarpe Simão

### Table of contents

- [Introduction: Business Problem](#)
- [Data Sources](#)
- [Dataset Structuring](#)
- [Visualization](#)
- [Analysis](#)
- [Results and Discussion](#)
- [Conclusion](#)

## ▼ Introduction: Business Problem

Among the numerous and diverse challenges humanity has faced through modern history, the recent Corona Virus pandemic showed itself as one of the most urgent and formidable one. Not only it brought unparalleled distress to every nation's healthcare system, but inflicted heavy losses in its economies, specially regarding the survival of small businesses - mainly composed of family restaurants and stores. Still, despite minimal improvements in the infection curve scenarios, many countries are starting to reopen their economies due to the aggravating economical factors.

Invariably, the small food business are soon to receive quarantine eager clients, when the pandemic, despite still existing, has began to be treated as a "new normal" situation. In this scenario, this project aims to provide both businesses and clients a dynamic tool for evaluating the risk of attending/providing services at each neighbourhood, based on infection rate, populational density and quantity of stores available (as wait lines and maximum capacity are a crucial factor imposed by health authorities and overall safety practices).

The initial goal of the project is to evaluate risks specifically for the **food service industry**. This model must not be used to analyse the risk of individual restaurants, but the risk associated with its operation in its neighbourhood. It's also expected to establish a relationship between clusters of venues and infection rate over neighborhoods.

This is an easily replicable experiment for every city in the world, but, due to ease of access of data

## ▼ Data Sources

The data needed for the project is listed below (besides the restaurants location data obtained from Foursquare API):

1. Geographical data of buroughs and neighbourhoods of NYC and their latitude/longitude.  
Available at: [https://geo.nyu.edu/catalog/nyu\\_2451\\_34572](https://geo.nyu.edu/catalog/nyu_2451_34572). The data has been conveniently uploaded at the server as: [https://cocl.us/new\\_york\\_dataset](https://cocl.us/new_york_dataset)

Borough Neighborhood Latitude Longitude 35 Manhattan Turtle Bay 40.752042 -73.967708 36  
Manhattan Tudor City 40.746917 -73.971219 37 Manhattan Stuyvesant Town 40.731000  
-73.974052 38 Manhattan Flatiron 40.739673 -73.990947 39 Manhattan Hudson Yards 40.756658  
-74.000111

2. COVID-19 infection data, kill rate and quantity of registered cases. Available at:  
<https://www1.nyc.gov/site/doh/covid/covid-19-data.page> (downloaded on 07/21/2020 as a csv for this project)

ZIP Neighborhood Borough Cases Cases per 100,000 Deaths per 100,000 Percent positive (of people tested)  
10001 Chelsea [NoMad/West](#) Chelsea Manhattan 415 1761,23 101,85 7,78 10002  
Chinatown/Lower East Side Manhattan 1214 1581,65 208,45 11 10003 East Village [Gramercy/Greenwich](#) Village Manhattan 504 936,77 63,2 5,92 10004 Financial District Manhattan 36 986,14 27,39 6,36 10005 Financial District Manhattan 76 905,18 23,82 5,71

3. NYC population data per burough/neighbourhood. Available at  
<https://data.cityofnewyork.us/City-Government/New-York-City-Population-By-Neighborhood-Tabulatio/swpk-hqdp>. Data extracted via SODA API JSON:  
<https://data.cityofnewyork.us/resource/swpk-hqdp.json>

## ▼ Dataset Structuring

## ▼ Libraries

```

1 import numpy as np #Vectorized data
2
3 import pandas as pd #Dataframes
4 pd.set_option('display.max_columns', None)
5 pd.set_option('display.max_rows', None)

```

```

6
7 import json #JSON Files
8 from pandas.io.json import json_normalize # transform JSON file into a pandas dataframe
9
10 import requests #Requests
11
12 #!conda install -c conda-forge geopy --yes
13 from geopy.geocoders import Nominatim #Address into lat/long
14
15 import matplotlib.cm as cm #Plotting and plotting modules
16 import matplotlib.colors as colors
17
18 import matplotlib.pyplot as plt #Plotting additional graphs
19 from matplotlib.ticker import PercentFormatter
20
21 from sklearn.cluster import KMeans #K-Means and Clustering
22
23 #!conda install -c conda-forge folium=0.5.0 --yes
24 import folium #Map Rendering (Folium)
25
26 print('Success!')

```

Success!

## ▼ Data Import

```

1 # NYC Neighborhood Names
2
3 !wget -q -O 'newyork_data.json' https://cocl.us/new_york_dataset #Get JSON
4 with open('newyork_data.json') as json_data:
5     newyork_data = json.load(json_data) #Load data from JSON
6 neighborhoods_data = newyork_data['features']
7 column_names = ['Borough', 'Neighborhood', 'Latitude', 'Longitude'] #Specifying columns for
8 neighborhoods = pd.DataFrame(columns=column_names) #Create DF header
9
10 for data in neighborhoods_data: #Fill DF rows
11     borough = neighborhood_name = data['properties']['borough']
12     neighborhood_name = data['properties']['name']
13
14     neighborhood_latlon = data['geometry']['coordinates']
15     neighborhood_lat = neighborhood_latlon[1]
16     neighborhood_lon = neighborhood_latlon[0]
17
18     neighborhoods = neighborhoods.append({'Borough': borough,
19                                         'Neighborhood': neighborhood_name,
20                                         'Latitude': neighborhood_lat,
21                                         'Longitude': neighborhood_lon}, ignore_index=True)
22
23
24 neighborhoods.drop_duplicates(subset = "Neighborhood", keep = False, inplace = True) #Drop
https://colab.research.google.com/drive/1NKI2o1P-QhrctPVOR9ZBo_u-fnGr7gef#scrollTo=G0YKBUKa7eQU&uniquifier=1&printMode=true

```

```

25 print("Rows:",len(neighborhoods['Neighborhood'].unique())) #Checking size with data source
26
27 neighborhoods.sort_values("Neighborhood",inplace=True) #Sort alphabetically
28 neighborhoods.head()

```

Rows: 298

	Borough	Neighborhood	Latitude	Longitude
298	Bronx	Allerton	40.865788	-73.859319
215	Staten Island	Annadale	40.538114	-74.178549
241	Staten Island	Arden Heights	40.549286	-74.185887
227	Staten Island	Arlington	40.635325	-74.165104
228	Staten Island	Arrochar	40.596313	-74.067124

```

1 #COVID Infection Data
2
3 covid_df=pd.read_csv('https://raw.githubusercontent.com/rafzss/Coursera_Capstone/master/Da
4 covid_df=covid_df.rename(columns={"Cases per 100,000":"Cases per 100k"})
5 covid_df=covid_df.rename(columns={"Deaths per 100,000":"Deaths per 100k"})
6 covid_df=covid_df.rename(columns={"Percent positive<br>(of people tested)":"Percent Positi
7 covid_df.drop(["ZIP"],axis=1,inplace=True)
8 print("Rows before splitting:",len(covid_df['Neighborhood'].unique())) #Checking size with
9
10 #Splitting rows of grouped neighborhoods
11
12 def splitDataFrameList(df,target_column,separator): # Based on: https://gist.github.com/jl
13
14     # df = dataframe to split,
15     # target_column = the column containing the values to split
16     # separator = the symbol used to perform the split
17     # returns: a dataframe with each entry for the target column separated, with each elem
18     #The values in the other columns are duplicated across the newly divided rows.
19
20     row_accumulator = []
21     def splitListToRows(row, separator):
22         split_row = row[target_column].split(separator)
23         for s in split_row:
24             new_row = row.to_dict()
25             new_row[target_column] = s
26             row_accumulator.append(new_row)
27
28     df.apply(splitListToRows, axis=1, args = (separator, ))
29     new_df = pd.DataFrame(row_accumulator)
30     return new_df
31
32 covid_df=splitDataFrameList(covid_df,"Neighborhood","/")
33
34 covid_df=covid_df.sort_values('Cases per 100k').drop_duplicates('Neighborhood',keep='last'

```

```

35 covid_df.sort_values("Neighborhood", inplace=True) #Sort alphabetically
36
37 print("Rows after splitting:", len(covid_df['Neighborhood'].unique())) #Checking size with
38 covid_df.head()

```

Rows before splitting: 162

Rows after splitting: 225

	Neighborhood	Borough	Cases	Cases per 100k	Deaths per 100k	Percent Positive
285	Airport	Queens	1690	4634.93	490.92	23.11
156	Allerton	Bronx	3164	4424.66	496.45	20.93
12	Alphabet City	Manhattan	774	1317.49	117.45	7.94
104	Annadale	Staten Island	1597	2648.23	97.84	16.97

1 # NYC Population

2

```

3 nycpop=pd.read_csv('https://raw.githubusercontent.com/rafzss/Coursera_Capstone/master/Data'
4 print("Latest Year:",nycpop["Year"].max()) #Get latest year of poll

```

```

5 nycpop=nycpop[nycpop['Year']==nycpop["Year"].max()] #Filter dataset by latest year
6

```

7 # This step approximates the population for neighborhoods that were grouped from the loaded

8 # The total population is simply divided by the number of occurrences of neighborhoods in

9

```

10 nycpop=splitDataFrameList(nycpop,"NTA Name","-") #Split rows with more than one neighborho
11 dupsdf=nycpop.pivot_table(index=['NTA Code'],aggfunc='size') #Count duplicates of each nei
12 nycpop=nycpop.merge(dupsdf.rename('Count').to_frame(),on=['NTA Code'],how='left') #Add pre
13 nycpop['Population']=round(nycpop['Population']/nycpop['Count']) #Approximate population o
14

```

```

15 nycpop.drop_duplicates(subset ="NTA Name", keep = False, inplace = True) #Drop duplicates
16 nycpop=nycpop.drop(['FIPS County Code','NTA Code','Count'],axis=1) #Drop unnecessary colum
17 print("Rows:",len(nycpop['NTA Name'].unique())) #Checking size with data source
18

```

```

19 nycpop.head()

```

Latest Year: 2010

Rows: 298

	Borough	Year	NTA Name	Population
0	Bronx	2010	Claremont	15539.0
1	Bronx	2010	Bathgate	15539.0
2	Bronx	2010	Eastchester	11506.0
3	Bronx	2010	Edenwald	11506.0
4	Bronx	2010	Baychester	11506.0

## ▼ Data Preparation

```

1 neighborhoods.set_index('Neighborhood') #Set index as neighborhood
2 nycpop.set_index('NTA Name') #Set index as neighborhood
3 covid_df.set_index('Neighborhood') #Set index as neighborhood
4
5 df1 = pd.merge(nycpop, covid_df, how='inner', left_on='NTA Name', right_on='Neighborhood')
6 df = pd.merge(neighborhoods,df1,how='inner',left_on='Neighborhood',right_on='NTA Name') #M
7 df.drop_duplicates(subset ="Neighborhood_x", keep = False, inplace = True) #Drop duplicate
8 df=df.drop(['Borough_x','Year','NTA Name','Neighborhood_y','Borough_y'],axis=1) #Drop redu
9 df=df.rename(columns={"Neighborhood_x":"Neighborhood"}) #Rename 'neighborhood'
10 df.set_index('Neighborhood') #Set index as neighborhood
11 print(df.shape)
12 df.head()

```

→ (150, 9)

	Borough	Neighborhood	Latitude	Longitude	Population	Cases	Cases per 100k	Deaths per 100k	Per Posi
0	Bronx	Allerton	40.865788	-73.859319	14452.0	3164	4424.66	496.45	2
1	Staten Island	Annadale	40.538114	-74.178549	6942.0	1597	2648.23	97.84	1
2	Staten Island	Arrochar	40.596313	-74.067124	5360.0	1315	3089.68	187.97	1

We see that the dataset went from the original 298 rows of the neighborhood dataset to 150, due to removal of duplicates and wrong matchings of the neighborhood strings of each dataset. Python allows for ease to use matching algoritm through pre-built libraries, such as Levenshtein's algoritm in the 'fuzzywuzzy' library. This will not be applied in this project so as to not overextend it.

## ▼ Data Visualization

```

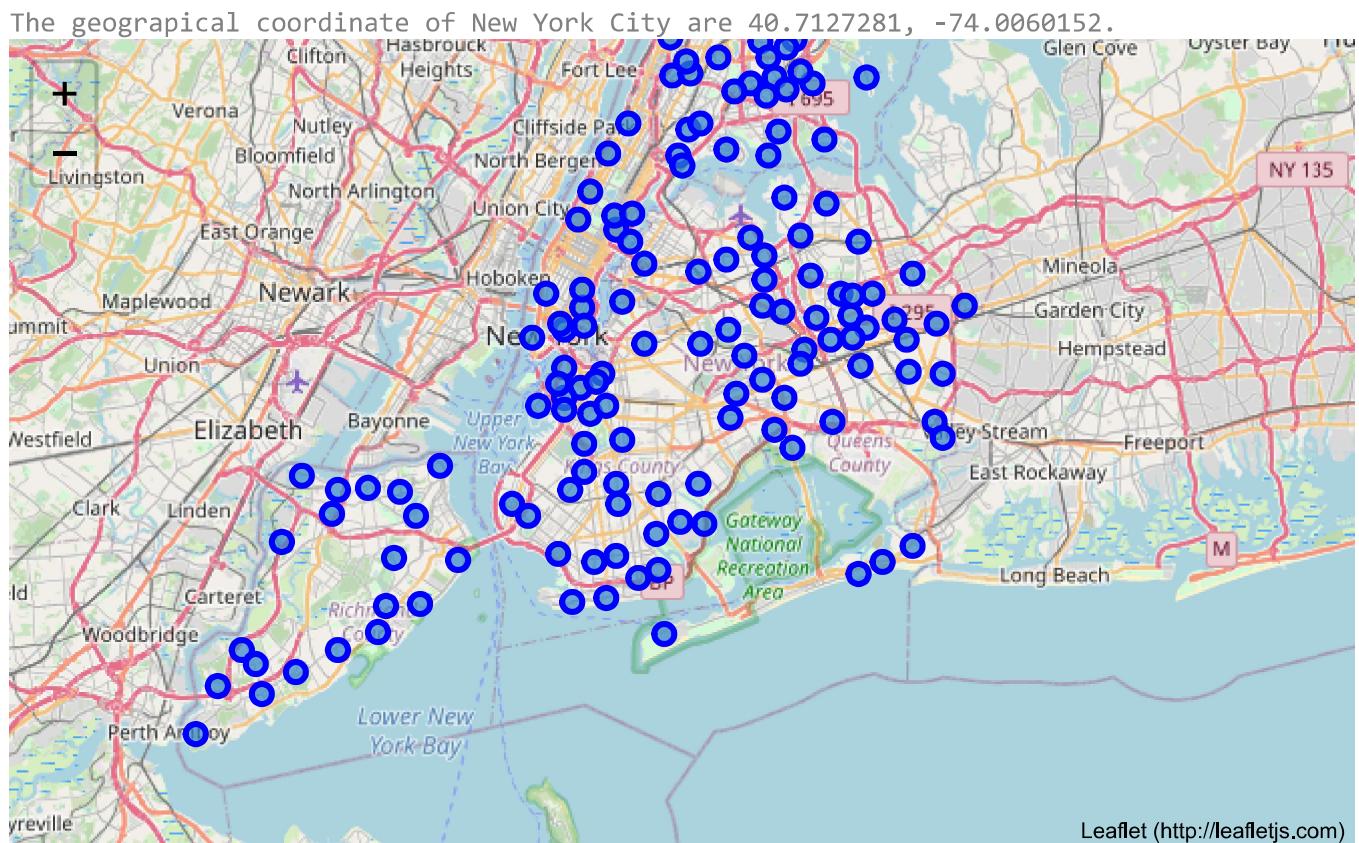
1 #Locate NYC Coordinates
2
3 address = 'New York City, NY'
4 geolocator = Nominatim(user_agent="ny_explorer")
5 location = geolocator.geocode(address)
6 latitude = location.latitude
7 longitude = location.longitude
8 print('The geographical coordinate of New York City are {}, {}'.format(latitude, longitude)
9
10 # Create map of New York using latitude and longitude values
https://colab.research.google.com/drive/1NKI2o1P-QhrctPVOR9ZBo_u-fnGr7gef#scrollTo=G0YKBUKa7eQU&uniquifier=1&printMode=true

```

```
to # Create map of New York using latitude and longitude values
11 map_newyork = folium.Map(location=[latitude, longitude], zoom_start=10)
12
13 # add markers to map
14 for lat, lng, borough, neighborhood in zip(df['Latitude'], df['Longitude'], df['Borough'],
15     label = '{}, {}'.format(neighborhood, borough)
16     label = folium.Popup(label, parse_html=True)
17     folium.CircleMarker(
18         [lat, lng],
19         radius=5,
20         popup=label,
21         color='blue',
22         fill=True,
23         fill_color='#3186cc',
24         fill_opacity=0.7,
25         parse_html=False).add_to(map_newyork)
26
27 map_newyork
```



```
3 def getNearbyVenues(names, latitudes, longitudes, radius=500):
4
5     venues_list=[]
6     for name, lat, lng in zip(names, latitudes, longitudes):
7         print(name)
8
9         # create the API request URL
10        url = 'https://api.foursquare.com/v2/venues/explore?&client_id={}&client_secret={}
11            CLIENT_ID,
12            CLIENT_SECRET,
13            VERSION,
14            lat,
15            lng,
16            radius,
17            LIMIT)
18
19        # make the GET request
20        results = requests.get(url).json()["response"]['groups'][0]['items']
21
22        # return only relevant information for each nearby venue
23        venues_list.append([
24            name,
25            lat,
26            lng,
27            v['venue']['name'],
28            v['venue']['location']['lat'],
29            v['venue']['location']['lng'],
30            v['venue']['categories'][0]['name']) for v in results])
31
32    nearby_venues = pd.DataFrame([item for venue_list in venues_list for item in venue_lis
33    nearby_venues.columns = ['Neighborhood',
34                            'Neighborhood Latitude',
35                            'Neighborhood Longitude',
36                            'Venue',
37                            'Venue Latitude',
38                            'Venue Longitude',
39                            'Venue Category']
40
41    return(nearby_venues)
42
43 # Function that extracts the category of the venue
44
45 def get_category_type(row):
46     try:
47         categories_list = row['categories']
48     except:
49         categories_list = row['venue.categories']
50
51     if len(categories_list) == 0:
52         return None
53     else:
54         return categories_list[0]['name']
```



## ▼ Analysis

### ▼ Foursquare Connection

```

1 CLIENT_ID = 'VLXFKV40KQLKJP3NM3U2VTOWIZYHKJXM3CXWH3SUVILRPA5D' # your Foursquare ID
2 CLIENT_SECRET = 'QJVKBVYQLA52BKUEV3AQGQPUBB0UFP00RLPQTMYXBAPGZMZD' # your Foursquare Secret
3 VERSION = '20180605' # Foursquare API version
4
5 print('Your credentails:')
6 print('CLIENT_ID: ' + CLIENT_ID)
7 print('CLIENT_SECRET:' + CLIENT_SECRET)

```

⇨ Your credentails:  
 CLIENT\_ID: VLXFKV40KQLKJP3NM3U2VTOWIZYHKJXM3CXWH3SUVILRPA5D  
 CLIENT\_SECRET: QJVKBVYQLA52BKUEV3AQGQPUBB0UFP00RLPQTMYXBAPGZMZD

### ▼ Foursquare Venues and Categories

```

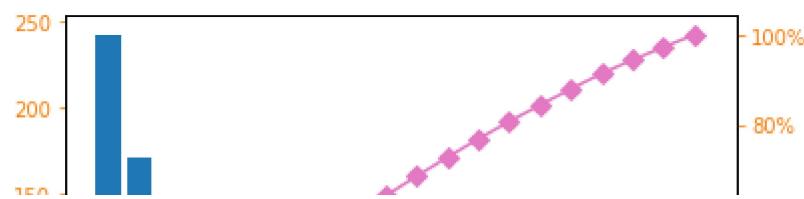
1 #Standard Function for fetching Foursquare data
2

```

12.11.2011 17:31:25 cat\_imagine

```
1 # Fetching Foursquare Data:  
2  
3 LIMIT=100 #Limit of venues returned by Foursquare API (for each neighborhood)  
4 radius=500 #Limit of radius for getting venues  
5  
6 nyc_venues = getNearbyVenues(names=df['Neighborhood'],  
7                                latitudes=df['Latitude'],  
8                                longitudes=df['Longitude'])  
9  
10  
11 print('There are {} uniques categories.'.format(len(nyc_venues['Venue Category'].unique())))  
  
1 # Pareto Chart  
2  
3 df_pareto = nyc_venues.groupby(['Venue Category'])['Venue Category'].count().to_frame(name='Count')  
4 df_pareto.index=df_pareto['Venue Category']  
5 df_pareto = df_pareto.sort_values(by='Count',ascending=False)  
6 df_pareto = df_pareto.nlargest(20,['Count'])  
7 df_pareto["cumulative_percent"] = df_pareto["Count"].cumsum()/df_pareto["Count"].sum()*100  
8 df_pareto  
9  
10 fig, ax = plt.subplots()  
11 ax.bar(df_pareto.index, df_pareto["Count"], color="C0")  
12 ax2 = ax.twinx()  
13 ax2.plot(df_pareto.index, df_pareto["cumulative_percent"], color="C6", marker="D", ms=7)  
14 ax2.yaxis.set_major_formatter(PercentFormatter())  
15  
16 ax.tick_params(axis="y", colors="C1")  
17 ax2.tick_params(axis="y", colors="C1")  
18  
19 for tick in ax.get_xticklabels():  
20     tick.set_rotation(90)  
21     tick.set_color("C0")  
22 plt.show()
```





As imagined, food services places correspond to the vast majority of the venues analysed (at least from the top 20 venues).

50 |  

```
1 # Grouping venues into neighborhoods
2
3 nyc_grouped=nyc_venues.groupby(['Neighborhood'])['Venue Category'].count().to_frame(name='Count')
4 nyc_grouped.sort_values
5 print(nyc_grouped.shape)
6
7
8 # Merging venue data with dataset:
9
10 df2 = pd.merge(df,nyc_grouped,how='left',left_on='Neighborhood',right_on='Neighborhood')
11 print(df2.shape)
12
13 df2.head() #This is the final dataset for analysis
14 # Some columns will not be used, but they provide details for future analysis and comparison
```

	Borough	Neighborhood	Latitude	Longitude	Population	Cases	Cases per 100k	Deaths per 100k	Per Posit.
0	Bronx	Allerton	40.865788	-73.859319	14452.0	3164	4424.66	496.45	2
1	Staten Island	Annadale	40.538114	-74.178549	6942.0	1597	2648.23	97.84	1
2	Staten Island	Arrochar	40.596313	-74.067124	5360.0	1315	3089.68	187.97	1

## ▼ Cluster Neighborhoods

```
1 # K-Means for clustering neighborhoods into 4 risk rates (Low, Average, High and Critical)
2
3 df_kcluster = df2.filter(['Cases per 100k','Venues Qty']) # Dataframe with X
4
5 k=4
6 kmeans = KMeans(n_clusters=k, random_state=0).fit(df_kcluster) # Run k-means clustering
7 kmeans.labels_[0:10] # Check cluster labels generated for each row in the dataframe
8
```

```

9
10 # Create new dataframe with clusters
11
12 # add clustering labels
13 df_kcluster.insert(0, 'Cluster Labels', kmeans.labels_)
14 df_kcluster=df_kcluster.drop(['Cases per 100k','Venues Qty'],axis=1) # Drop redundant colu
15 df_final = df2.merge(df_kcluster, left_index=True,right_index=True) # Merge Cluster Labels w
16
17 print(df_final.shape)
18 df_final.head()

```

↳ (150, 11)

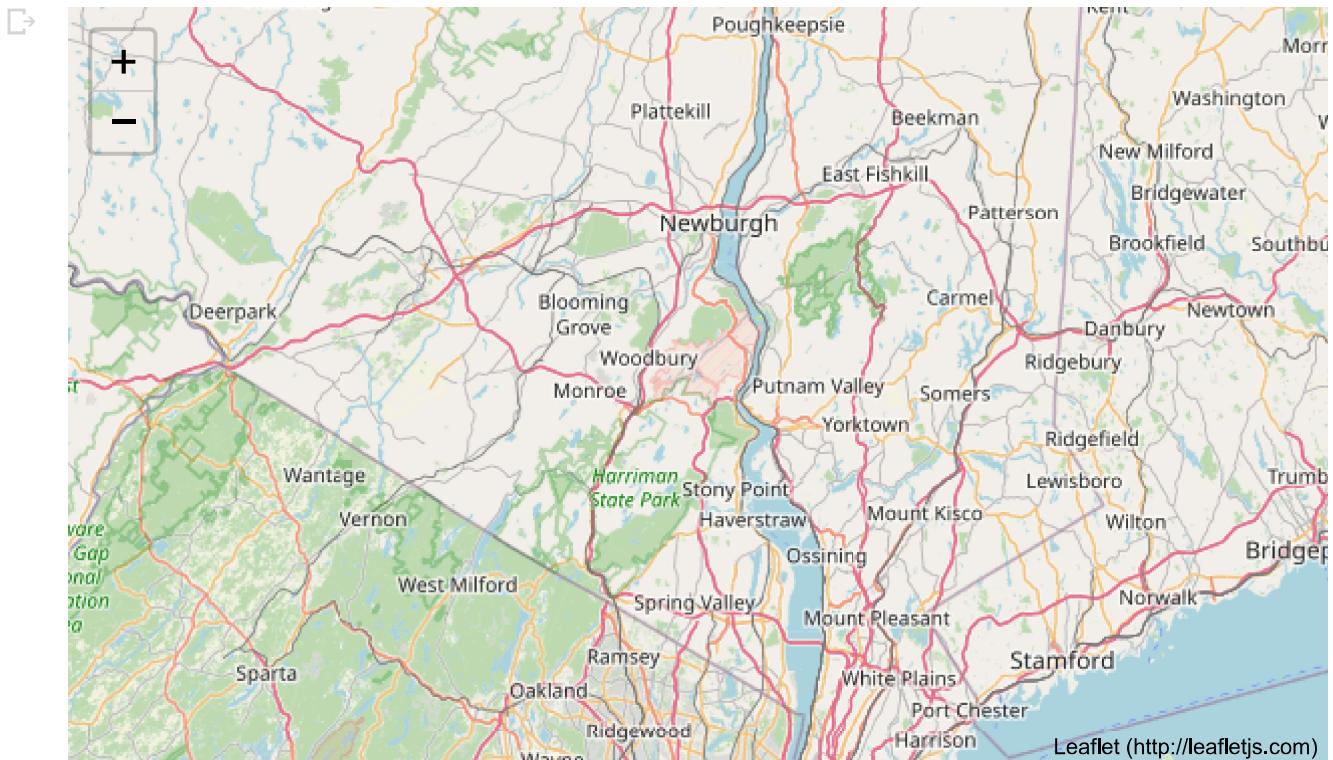
	Borough	Neighborhood	Latitude	Longitude	Population	Cases	Cases per 100k	Deaths per 100k	Per Posi
0	Bronx	Allerton	40.865788	-73.859319	14452.0	3164	4424.66	496.45	2
1	Staten Island	Annadale	40.538114	-74.178549	6942.0	1597	2648.23	97.84	1
2	Staten Island	Arrochar	40.596313	-74.067124	5360.0	1315	3089.68	187.97	1

## ▼ Visualize Clusters

```

1 # create map
2 map_clusters = folium.Map(location=[latitude, longitude], zoom_start=11)
3
4 # set color scheme for the clusters
5 x = np.arange(k)
6 ys = [i + x + (i*x)**2 for i in range(k)]
7 colors_array = cm.rainbow(np.linspace(0, 1, len(ys)))
8 rainbow = [colors.rgb2hex(i) for i in colors_array]
9
10 # add markers to the map
11 markers_colors = []
12 for lat, lon, poi, cluster in zip(df_final['Latitude'], df_final['Longitude'], df_final['N
13     label = folium.Popup(str(poi) + ' Cluster ' + str(cluster), parse_html=True)
14     folium.CircleMarker(
15         [lat, lon],
16         radius=5,
17         popup=label,
18         color=rainbow[cluster-1],
19         fill=True,
20         fill_color=rainbow[cluster-1],
21         fill_opacity=0.7).add_to(map_clusters)
22
23 map_clusters

```



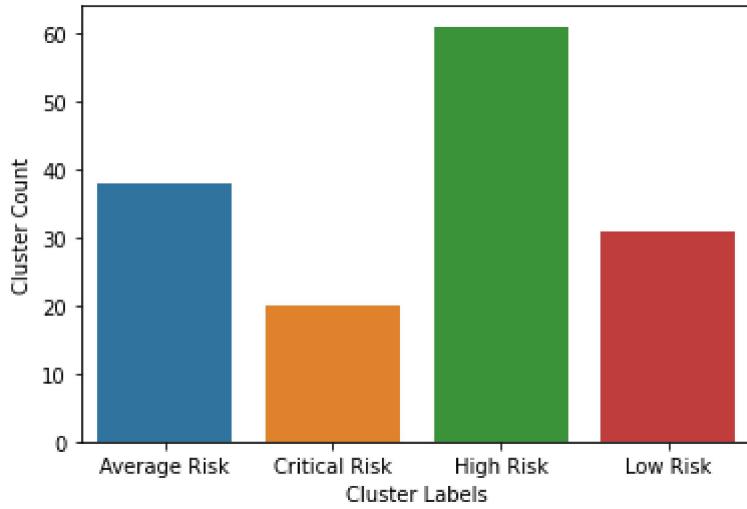
## ▼ Results and Discussion

```

1 import seaborn as sns # Plotting library
2
3 # Bar Plot
4
5 df_plot = df_final.filter(['Cases per 100k','Venues Qty','Cluster Labels'],axis=1)
6 df_plot['Cluster Labels']=df_plot['Cluster Labels'].replace(0,'High Risk').replace(1,'Low')
7 df_count = df_plot.groupby(['Cluster Labels'])['Cluster Labels'].count().to_frame(name='Cl')
8 df_count.head()
9
10 sns.barplot(x='Cluster Labels',y='Cluster Count',data=df_count) # Bar Plot

```

↳ <matplotlib.axes.\_subplots.AxesSubplot at 0x7f55b996b160>



In the actual scenario (07-23-2020), New York City shows a higher amount of neighborhoods with 'High Risk' cases.

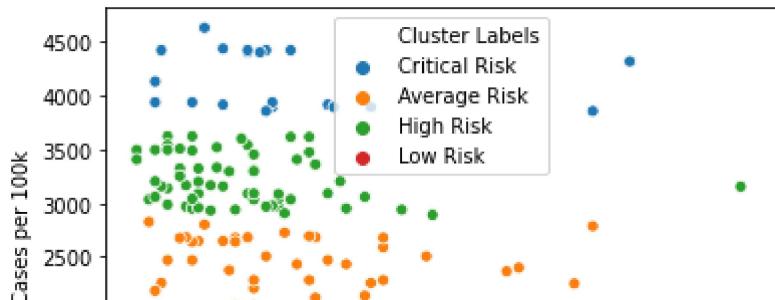
```

1 # Scatterplot of Cluster Labels
2
3 sns.scatterplot(x='Venues Qty',y='Cases per 100k',data=df_plot,hue='Cluster Labels')

```

↳

<matplotlib.axes.\_subplots.AxesSubplot at 0x7f55ba8150b8>



We can see then the classification by clusters that represent the healthcare precautions strictness for each risk group (cluster labels). It's intuitive to think that neighborhoods with a high availability of stores and low infection rate may present lower risks than the opposite scenario. It is noted as well that the clusters have their frontiers in a clear geospatial manner, but this probably represents the correlation between amount of venues, people transit and socioeconomic bias of each region, meaning that the relationship of geography is not of causation.

- Group 0: High Risk
- Group 1: Low Risk
- Group 2: Critical Risk
- Group 3: Average Risk

In summary, the results did not display abnormalities or an excessive amount of outliers. The data showed a rather clear clustering, indexed by the independent variables of the study. It's important to remember that K-means is an unsupervised algorithm and should not be artificially labeled. The previous results show a interpretation of the randomic data generated in this session.

## ▼ Conclusion

Despite discarding an elevated amount of data due to string mismatching of geospatial data definitions and numeric data of pandemic, the model showed itself to clearly define well structured cluster labels, without evident outliers.

The ease of use of Foursquare API and Python's libraries allowed for an agile and assertive approach to such urgent and delicate problem. The complexity of the pandemic scenario which we live on is not denied, but the model could - simply - establish action response levels for the variety of healthcare solutions according to the social context of each neighborhood, where it could serve as a response guide, particularly to small towns that do not rely on a structured core of medical scientists.

Future refinement of this project involve better string comparison tools for merging different sourced datasets and inclusion of other variables to the predicted model.

