

# Coffee Shops Analysis in Manhattan; Segmenting and Clustering Neighborhoods

Capstone Project - Battle of the Neighborhoods: Applied Data Science Capstone by IBM(Coursera)  
Author: Lopamudra Nayak



## Table of Contents

- [1. Introduction](#)
- [2. Download and Explore Dataset](#)
- [3. Explore Neighborhoods in New York City](#)
- [4. Methodology](#)
- [5. Analyze Each Neighborhood](#)
- [6. Cluster Neighborhoods](#)
- [7. Examine Clusters](#)
- [8. Results](#)
- [9. Conclusion](#)

## 1. Introduction

New York is one of the most famous cities in the world also referred as the "Big Apple", this vibrant city is known for its exclusive shops, flashy Broadway performances, and high-flying business tycoons, and it's a city that has long captivated people from all over the world. But, aside from the flashing lights, it's also home to the beloved Statue of Liberty, a symbol of freedom and hope, and Central Park, one of the tidiest and best-kept parks in the world. The proverbial apple of New York State's eye, Manhattan is the most popular tourist destination and one of the best places to visit in New York State. It's also one of the best cities in New York and, arguably, the world. One can find literally everything your heart desires starting from the rarest antique cufflinks to the most delicious entree to grace your lips. An activity that is touristy but also for its residents is drinking coffee. The activity can be done alone, accompanied by friends, and anytime in this city.

Americans drink about 146 billion cups of coffee per year. In 2020, Americans drank an average of 1.87 cups of coffee per day. Between 2011 and 2020, US. coffee consumption per capita per day remained relatively close to two cups, with the exception of 2016, when US survey respondents drank roughly 1.64 cups of coffee per day. In US, people aged seventy and older drank the most coffee per day in 2020. This particular age group drank an average of almost 2.2 cups per capita per day that year. Americans aged eighteen and nineteen drank the least per capita, namely about 0.8 cups a day. In terms of gender, men drank more than women per capita in the United States in the same year. In 2020, Americans drank about 1.3 cups of coffee during breakfast each day, making it the most common time of day for people to drink coffee in the US that year. The least common time for Americans to drink coffee was during other meals. ([total-us-coffee-per-capita-consumption](https://www.statista.com/statistics/456360/us-total-us-coffee-per-capita-consumption/))

In this project, we will try to find an optimal location for a coffee shop. This report will be targeted at stakeholders who want to start from scratch, buy an existing business, or anybody interested in a good coffee in Manhattan, New York City.

We will try to detect locations that are not already crowded with coffee shops. We are also particularly interested in areas with no coffee shops in the vicinity. We would also prefer locations as close to the city center as possible, assuming that the first two conditions are met.

### Import Libraries:

Before we get the data and start exploring it, let's download all the dependencies that we will need.

```
In [1]: import numpy as np # Library to handle data in a vectorized manner
import pandas as pd # Library for data analysis
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)

import json # Library to handle JSON files
#conda install -c conda-forge geopy --yes
from geopy.geocoders import Nominatim # convert an address into latitude and longitude values
import requests # library to handle requests
from pandas.io.json import json_normalize # transform JSON file into a pandas dataframe
# Matplotlib and associated plotting modules
import matplotlib.cm as cm
import matplotlib.colors as colors

# import k-means from clustering stage
from sklearn.cluster import KMeans
#conda install -c conda-forge folium=0.5.0 --yes
import folium # map rendering library
print('Libraries imported.')
```

Libraries imported.

## 2. Download and Explore Dataset

In order to segment the neighborhoods and explore them, we will essentially need a dataset that contains the 5 boroughs and the neighborhoods that exist in each borough as well as the the latitude and longitude coordinates of each neighborhood. I have used the prepared files from IBM cloud server and simply run a `wget` command to access the data.

```
In [2]: !wget -q -O 'newyork_data.json' https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMDeveloperSkillsNetwork-D50701EN-SkillsNetwork/labs/newyork_data.json
print('Data downloaded!')
```

Data downloaded!

### Load and explore the data

Next, let's load the data.

```
In [3]: with open('newyork_data.json') as json_data:
    newyork_data = json.load(json_data)
```

Let's take a quick look at the data.

```
In [4]: newyork_data
{'id': 'nyu_2451_34572_3',
 'geometry': {'type': 'Point',
              'coordinates': [-73.82789644716412, 40.887555677350775]},
 'geometry_name': 'geom',
 'properties': {'name': 'Eastchester',
                'stacked': 1,
                'annoline1': 'Eastchester',
                'annoline2': None,
                'annoline3': None,
                'annoangle': 0.0},
 'borough': 'Bronx',
 'bbox': [-73.82789644716412, 40.887555677350775,
          -73.82789644716412, 40.887555677350775],
 'type': 'Feature",
 'id': 'nyu_2451_34572_4',
 'geometry': {'type': 'Point',
              'coordinates': [-73.9056425951682, 40.89543742690383]},
```

Notice how all the relevant data is in the `features` key, which is basically a list of the neighborhoods. So, let's define a new variable that includes this data.

```
In [5]: neighborhoods_data = newyork_data['features']
```

Let's take a look at the first item in this list.

```
In [6]: neighborhoods_data[0]
Out[6]: {"type": "Feature",
  "id": "nyu_2451_24572_1",
  "geometry": {"type": "Point",
    "coordinates": [-73.84720052054902, 40.89470517661]},
  "properties": {"name": "Wakefield",
    "label": "1",
    "annoline1": "Wakefield",
    "annoline2": None,
    "annoline3": None,
    "borough": "Bronx",
    "bbox": [-73.84720052054902,
     40.89470517661,
     -73.84720052054902,
     40.89470517661]}}
```

Transform the data into a `pandas` dataframe

Now we need to transform this data of nested Python dictionaries into a `pandas` dataframe. So, let's first create an empty dataframe.

```
In [7]: # define the dataframe columns
column_names = ['Borough', 'Neighborhood', 'Latitude', 'Longitude']

# instantiate the dataframe
neighborhoods = pd.DataFrame(columns=column_names)
```

Let's look at the empty dataframe to confirm that the columns are as intended.

```
In [8]: neighborhoods
```

```
Out[8]: Borough Neighborhood Latitude Longitude
```

Then let's loop through the data and fill the dataframe one row at a time.

```
In [9]: for data in neighborhoods_data:
  borough = data['properties']['borough']
  neighborhood_name = data['properties']['name']

  neighborhood_latitude = data['geometry']['coordinates'][0][1]
  neighborhood_lat = neighborhood_latitude[1]
  neighborhood_lon = neighborhood_latitude[0]

  neighborhoods = neighborhoods.append({'Borough': borough,
    'Neighborhood': neighborhood_name,
    'Latitude': neighborhood_lat,
    'Longitude': neighborhood_lon}, ignore_index=True)
```

Now just quickly examine the resulting dataframe.

```
In [10]: neighborhoods.tail()
```

```
Out[10]:
```

Borough	Neighborhood	Latitude	Longitude
301	Manhattan	40.75658	-74.00011
302	Queens	40.58738	-73.805530
303	Queens	40.61122	-73.765968
304	Queens	40.756091	-73.946631
305	Staten Island	40.61731	-74.081740

```
In [11]: neighborhoods.shape
```

```
Out[11]: (306, 4)
```

And make sure that the dataset has all 5 boroughs and 306 neighborhoods.

```
In [12]: print('The dataframe has {} boroughs and {} neighborhoods.'.format(
  len(neighborhoods['Borough'].unique()),
  neighborhoods.shape[0]
))
```

The dataframes has 5 boroughs and 306 neighborhoods.

Use geopy library to get the latitude and longitude values of New York City.

In order to define an instance of the geocoder, we need to define a user\_agent. We will name our agent `ny_explorer`, as shown below.

```
In [13]: address = 'New York City, NY'

geolocator = Nominatim(user_agent="ny_explorer")
location = geolocator.geocode(address)
latitude = location.latitude
longitude = location.longitude
print("The geographical coordinate of New York City are {}, {}".format(latitude, longitude))

The geographical coordinate of New York City are 40.712781, -74.0066152.
```

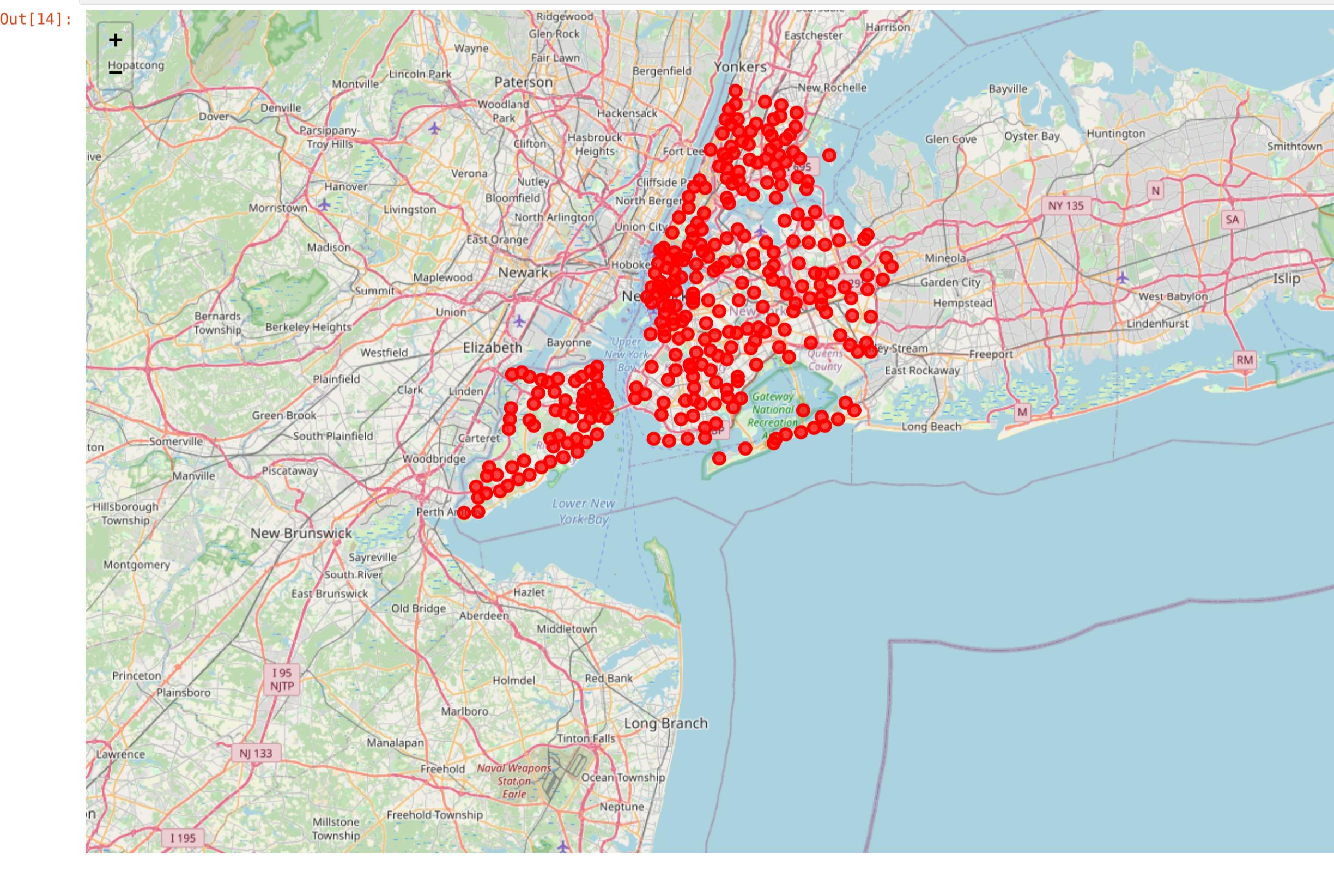
Create a map of New York with neighborhoods superimposed on top.

```
In [14]: # create map of New York using latitude and longitude values

map_newyork = folium.Map(location=[latitude, longitude], zoom_start=15)

# add markers to map
for lat, lng, borough, neighborhood in zip(neighborhoods['Latitude'], neighborhoods['Longitude'], neighborhoods['Borough'], neighborhoods['Neighborhood']):
  label = '{}, {}'.format(neighborhood, borough)
  label = folium.Popup(label, parse_html=True)
  folium.CircleMarker(
    [lat, lng],
    radius=5,
    radius_giant=10,
    color='red',
    fill=True,
    fill_color="#2196cc",
    fill_opacity=0.7,
    parse_html=False).add_to(map_newyork)

map_newyork
```



Leaflet (<http://leafletjs.com>)

Folium is a great visualization library. Feel free to zoom into the above map, and click on each circle mark to reveal the name of the neighborhood and its respective borough.

However, for illustration purposes, let's simplify the above map and segment and cluster only the neighborhoods in Manhattan. So let's slice the original dataframe and create a new dataframe of the Manhattan data.

```
In [15]: manhattan_data = neighborhoods[neighborhoods['Borough'] == 'Manhattan'].reset_index(drop=True)
```

```
Out[15]:
```

Borough	Neighborhood	Latitude	Longitude
0	Marble Hill	40.876551	-73.910660
1	Chinatown	40.715618	-73.994279
2	Washington Heights	40.851903	-73.936900
3	Inwood	40.867884	-73.921210
4	Hamilton Heights	40.823604	-73.949688

```
In [16]: manhattan_data.shape
```

```
Out[16]: (40, 4)
```

Let's get the geographical coordinates of Manhattan.

```
In [17]: address = 'Manhattan'

geocoder = Nominatim(user_agent="ny_explorer")
location = geocoder.geocode(address)
latitude = location.latitude
longitude = location.longitude
print('The geographical coordinate of Manhattan are {}, {}'.format(latitude, longitude))

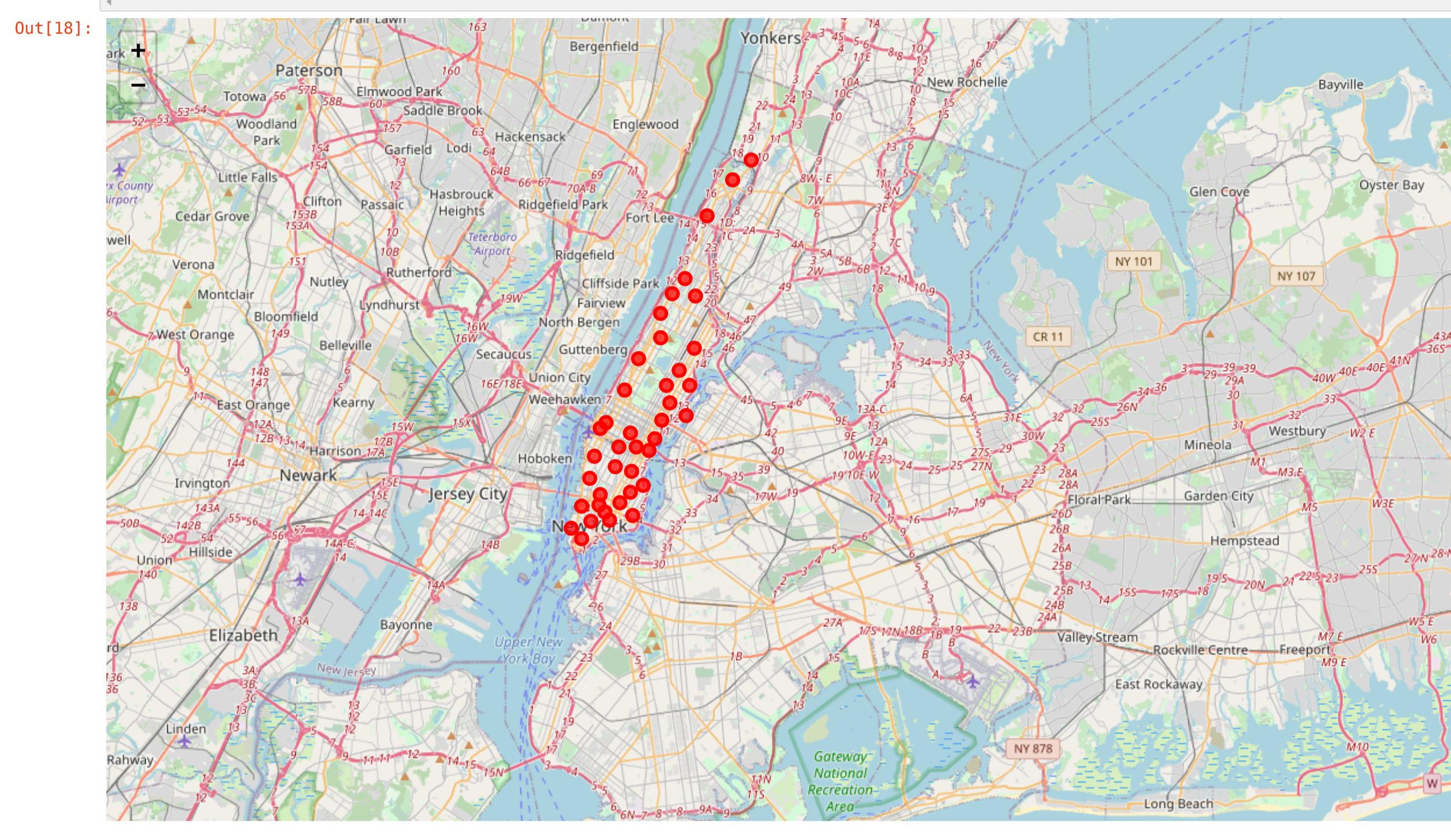
The geographical coordinate of Manhattan are 40.7896239, -73.9598939.
```

As we did with all of New York City, let's visualize Manhattan the neighborhoods in it.

```
In [18]: # create map of Manhattan using latitude and longitude values
map_manhattan = folium.Map(location=[latitude, longitude], zoom_start=11)

# add markers to map
for lat, lng, label in zip(manhattan_data['Latitude'], manhattan_data['Longitude'], manhattan_data['Neighborhood']):
    folium.CircleMarker(
        [lat, lng],
        radius=5,
        popup=label,
        color='red',
        fill=True,
        fill_color="#3186cc",
        fill_opacity=0.7,
        parse_html=False).add_to(map_manhattan)

map_manhattan
```



Leaflet (<http://leafletjs.com>)

Next, we are going to start utilizing the Foursquare API to explore the neighborhoods and segment them.

### 3. Methodology

The objective of this project is to find areas of Manhattan with a low density of coffee shops.

- First, we will check the candidate neighborhoods. It will be created the latitude & longitude coordinates for centroids of our candidate neighborhoods. We will create a grid of cells covering our area of interest which is centered around the Neighborhood Upper West Side.
- Secondly, we will get the top 100 venues that are in Central Park within a radius of 3.5 kilometers and explore their neighborhoods.

```
In [19]: manhattan_data.loc[manhattan_data['Neighborhood']=='Upper West Side']
```

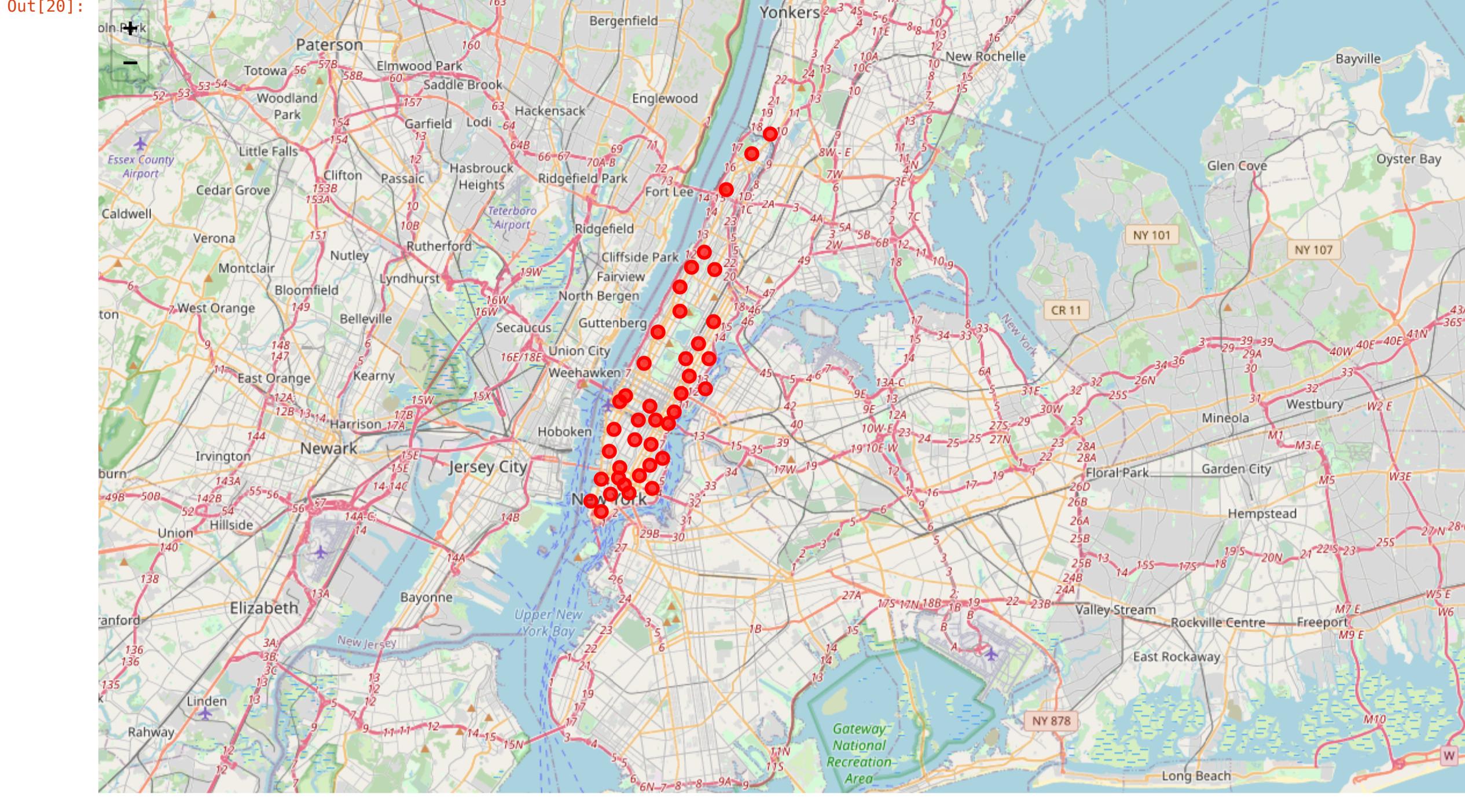
```
Out[19]:
Borough Neighborhood Latitude Longitude
12 Manhattan Upper West Side 40.7870568 -73.977099
```

```
In [20]: # create map of Montreal using latitude and longitude values
latitude = 46.787658
longitude = -73.977059

map_manhattan = folium.Map(location=[latitude, longitude], zoom_start=11)

# add markers to map
for lat, lng, label in zip(manhattan_data['Latitude'], manhattan_data['Longitude'], manhattan_data['Neighborhood']):
    folium.CircleMarker(
        [lat, lng],
        radius=5,
        popup=label,
        color='red',
        fill=True,
        fill_color="#3186cc",
        fill_opacity=0.7,
        parse_html=False).add_to(map_manhattan)

map_manhattan
```



Leaflet (<http://leafletjs.com>)

Now, we are going to start utilizing the Foursquare API to explore the neighborhoods and segment them.

Define Foursquare Credentials and Version

```
In [21]: CLIENT_ID = 'WSXE230WGSJAUOKVYCT4LBPMQ2G4Y1BGKSJUD0GHNPZP' # your Foursquare ID
CLIENT_SECRET = 'CVC44HNF25CREGZB0JBX130453PFYD1U4PG4TIA3ZV0Z' # your Foursquare Secret
VERSION = '20180605' # Foursquare API version
LIMIT = 100 # A default Foursquare API limit value

print('Your credentials:')
print('CLIENT_ID: ' + CLIENT_ID)
print('CLIENT_SECRET: ' + CLIENT_SECRET)

Your credentials:
CLIENT_ID: WSXE230WGSJAUOKVYCT4LBPMQ2G4Y1BGKSJUD0GHNPZP
CLIENT_SECRET:CVC44HNF25CREGZB0JBX130453PFYD1U4PG4TIA3ZV0Z
```

Let's explore the first neighborhood in our dataframe.

Get the neighborhood's name.

```
In [22]: manhattan_data.loc[0, 'Neighborhood']
```

```
Out[22]: 'Marble Hill'
```

Get the neighborhood's latitude and longitude values.

```
In [23]: neighborhood_latitude = manhattan_data.loc[0, 'Latitude'] # neighborhood latitude value
neighborhood_longitude = manhattan_data.loc[0, 'Longitude'] # neighborhood longitude value

neighborhood_name = manhattan_data.loc[0, 'Neighborhood'] # neighborhood name

print('Latitude and longitude values of {} are {}, {}'.format(neighborhood_name,
                                                             neighborhood_latitude,
                                                             neighborhood_longitude))
```

Latitude and longitude values of Marble Hill are 40.87655077879964, -73.91065965862981.

#### 4. Explore Neighborhoods in Manhattan

Let's create a function to get 100 venues of the neighborhoods within a radius of 3.5kms

```
In [24]: def getNearbyVenues(names, latitudes, longitudes, radius=3500):
    venues_list=[]
    for name, lat, lng in zip(names, latitudes, longitudes):
        print(name)

        # create the API request URL
        url = 'https://api.foursquare.com/v2/venues/explore?&client_id={}&client_secret={}&v={}&ll={},{}&radius={}&limit={}'.format(
            CLIENT_ID,
            CLIENT_SECRET,
            VERSION,
            lat,
            lng,
            radius,
            LIMIT)

        # make the GET request
        results = requests.get(url).json()['response']['groups'][0]['items']

        # return only relevant information for each nearby venue
        venues_list.append([{
            'name': name,
            'lat': lat,
            'lng': lng,
            'venue': {
                'name': v['venue']['name'],
                'location': {
                    'lat': v['venue']['location']['lat'],
                    'lng': v['venue']['location']['lng']
                },
                'categories': [v['venue']['categories'][0]['name']] for v in results
            }
        }])

    nearby_venues = pd.DataFrame([item for venue_list in venues_list for item in venue_list])
    nearby_venues.columns = ['Neighborhood',
                            'Neighborhood Latitude',
                            'Neighborhood Longitude',
                            'Venue',
                            'Venue Latitude',
                            'Venue Longitude',
                            'Venue Category']

    return(nearby_venues)
```

In [ ]:

Now write the code to run the above function on each neighborhood and create a new dataframe called `manhattan_venues`.

```
In [25]: # type your answer here
manhattan_venues = getNearbyVenues(names=manhattan_data['Neighborhood'],
                                    latitudes=manhattan_data['Latitude'],
                                    longitudes=manhattan_data['Longitude'])
```

```
Marble Hill
Chinatown
Washington Heights
Hamilton Heights
Manhattanville
Central Harlem
East Harlem
East Side
Yorkville
Lenox Hill
Roosevelt Island
Upper West Side
Columbus Circle
Clinton
Midtown
Murray Hill
Chelsea
Greenwich Village
...+ 44 more
```

Let's check the size of the resulting dataframe



In [34]: manhattan\_part = manhattan\_grouped[['Neighborhood', 'Coffee Shop']]

Out[34]:

Neighborhood	Coffee Shop
0 Battery Park City	0.03
1 Carnegie Hill	0.02
2 Central Harlem	0.07
3 Chelsea	0.03
4 Chinatown	0.02
5 Civic Center	0.04
6 Clinton	0.02
7 East Harlem	0.03
8 East Village	0.01
9 Financial District	0.03
10 Flatiron	0.02
11 Gramercy	0.01
12 Greenwich Village	0.01
13 Hamilton Heights	0.07
14 Hudson Yards	0.02
15 Inwood	0.01
16 Lenox Hill	0.01
17 Lincoln Square	0.02
18 Little Italy	0.02
19 Lower East Side	0.01
20 Manhattan Valley	0.04
21 Manhattanville	0.09
22 Marble Hill	0.02
23 Midtown	0.00
24 Midtown South	0.02
25 Morningside Heights	0.08
26 Murray Hill	0.02
27 Noho	0.01
28 Roosevelt Island	0.01
29 Soho	0.03
30 Stuyvesant Town	0.02
31 Sutton Place	0.01
32 Tribeca	0.03
33 Tudor City	0.00
34 Turtle Bay	0.00
35 Upper East Side	0.03
36 Upper West Side	0.01
37 Washington Heights	0.02
38 West Village	0.04
39 Yorkville	0.02

In [35]: manhattan\_merged = pd.merge(manhattan\_data, manhattan\_part, on='Neighborhood')

Out[35]:

Borough	Neighborhood	Latitude	Longitude	Coffee Shop
0 Manhattan	Marble Hill	40.876551	-73.910660	0.02
1 Manhattan	Chinatown	40.715618	-73.994279	0.02
2 Manhattan	Washington Heights	40.851903	-73.936900	0.02
3 Manhattan	Inwood	40.867684	-73.921210	0.01
4 Manhattan	Hamilton Heights	40.823604	-73.949698	0.07

Bar chart of coffee shops

In [36]:

```
graph = pd.DataFrame(manhattan_onehot.groupby('Neighborhood')['Coffee Shop'].sum())
graph = graph.sort_values(by = 'Coffee Shop', ascending=False)
graph.iloc[:30].plot(kind='bar', figsize=(10,6), color= 'red')
plt.xlabel("Neighborhoods")
plt.ylabel("No of Coffee Shop")
plt.title("Neighborhoods vs No of Coffee Shop")
plt.show()
```

In [37]: graph

Out[37]:

Neighborhood	Coffee Shop
Manhattanville	9
Morningside Heights	8
Hamilton Heights	7
Central Harlem	7
Manhattan Valley	4
West Village	4
Civic Center	4
Upper East Side	3
Tribeca	3
Soho	3
Battery Park City	3
East Harlem	3
Chelsea	3
Financial District	3
Yorkville	2
Murray Hill	2
Clinton	2
Midtown South	2
Chinatown	2
Marble Hill	2
Carnegie Hill	2
Little Italy	2
Lincoln Square	2
Washington Heights	2
Hudson Yards	2
Flatiron	2
Stuyvesant Town	2
Upper West Side	1
Sutton Place	1
Gramercy	1
Roosevelt Island	1
Noho	1
Greenwich Village	1
Lower East Side	1
Lenox Hill	1
Inwood	1
East Village	1
Tudor City	0
Turtle Bay	0
Midtown	0

In [38]: coffee\_df = manhattan\_grouped[['Neighborhood', 'Coffee Shop']]

Out[38]:

Neighborhood	Coffee Shop
0 Battery Park City	0.03
1 Carnegie Hill	0.02
2 Central Harlem	0.07
3 Chelsea	0.03
4 Chinatown	0.02

## 5. Cluster Neighborhoods

In [39]: df\_clustering = coffee\_df.drop('Neighborhood', 1)

In [40]:

```
from sklearn.cluster import KMeans
from yellowbrick.cluster.elbow import kelbow_visualizer
from yellowbrick.datasets.loaders import load_nfl
df_clustering, y = load_nfl()
# Use the quick method and immediately show the figure
kelbow_visualizer(KMeans(random_state=4), df_clustering, k=(2,10))
```

/home/lopa/.local/lib/python3.6/site-packages/sklearn/base.py:213: FutureWarning: From version 0.24, get\_params will raise an AttributeError if a parameter cannot be retrieved as an instance attribute. Previously it would return None.  
FutureWarning)

Out[40]:

```
KElbowVisualizer(ax=matplotlib.axes._subplots.AxesSubplot object at 0x7f0929e2cf60,
estimator=KMeans(n_clusters=9, random_state=4), k=None)
```

Run k-means to cluster the neighborhood into 4 clusters.

In [41]:

```
# set number of clusters
kClusters = 4

coffee_grouped_clustering = coffee_df.drop('Neighborhood', 1)

# run k-means clustering
kmeans = KMeans(n_clusters=kClusters, random_state=0).fit(coffee_grouped_clustering)

# check cluster labels generated for each row in the dataframe
kmeans.labels_[0:10]
```

Out[41]: array([0, 3, 1, 0, 3, 0, 2, 0, 1, dtype=int32)

```
In [42]: # make a copy of manhattan_data
coffee_merged = coffee_df
# add clustering labels
coffee_merged['Cluster Labels'] = kmeans.labels_
# add latitude/longitude for each neighborhood
coffee_merged = coffee_merged.join(manhattan_data.set_index('Neighborhood'), on='Neighborhood')
coffee_merged.head() # check the last columns
```

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
import sys
Out[42]:
Neighborhood Coffee Shop Cluster Labels Borough Latitude Longitude
0 Battery Park City 0.03 0 Manhattan 40.711932 -74.016869
1 Carnegie Hill 0.02 3 Manhattan 40.782683 -73.953256
2 Central Harlem 0.07 1 Manhattan 40.815976 -73.943211
3 Chelsea 0.03 0 Manhattan 40.744035 -74.003116
4 Chinatown 0.02 3 Manhattan 40.715618 -73.994279
```

```
In [43]: # Sort by cluster
coffee_merged.sort_values(['Cluster Labels'], inplace=True)
coffee_merged.head()
```

```
Out[43]:
Neighborhood Coffee Shop Cluster Labels Borough Latitude Longitude
0 Battery Park City 0.03 0 Manhattan 40.711932 -74.016869
29 Soho 0.03 0 Manhattan 40.722184 -74.000657
38 West Village 0.04 0 Manhattan 40.734434 -74.006180
9 Financial District 0.03 0 Manhattan 40.707107 -74.010665
7 East Harlem 0.03 0 Manhattan 40.792249 -73.944182
```

```
In [44]: coffee_merged = coffee_merged.reset_index(drop=True)
```

```
coffee_merged.head()
Out[44]:
Neighborhood Coffee Shop Cluster Labels Borough Latitude Longitude
0 Battery Park City 0.03 0 Manhattan 40.711932 -74.016869
1 Soho 0.03 0 Manhattan 40.722184 -74.000657
2 West Village 0.04 0 Manhattan 40.734434 -74.006180
3 Financial District 0.03 0 Manhattan 40.707107 -74.010665
4 East Harlem 0.03 0 Manhattan 40.792249 -73.944182
```

```
In [45]: coffee_merged.shape
```

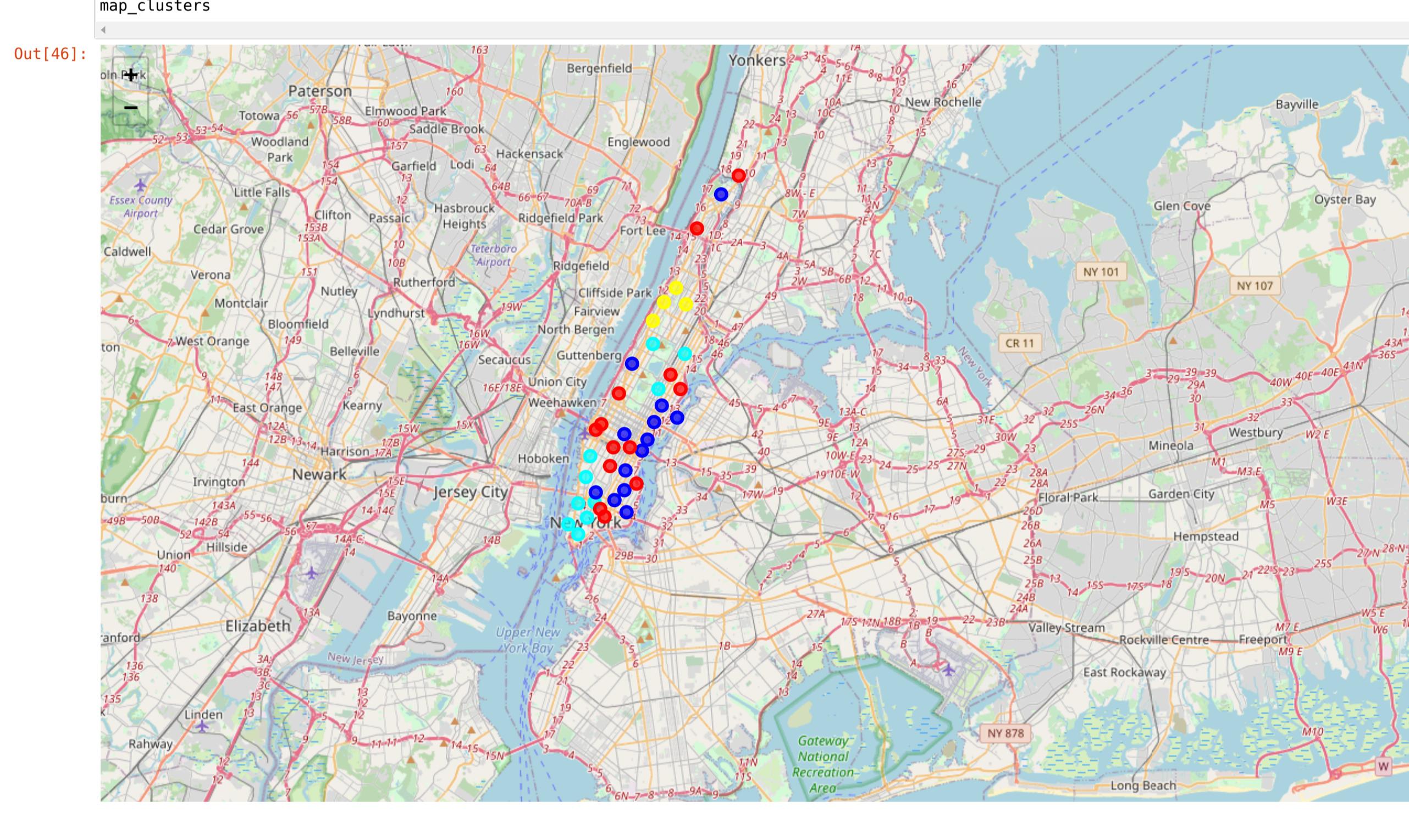
```
Out[45]: (40, 6)
```

Finally, let's visualize the resulting clusters

```
In [46]: # create map
latitude = 40.787658
longitude = -73.977059
map_clusters = folium.Map(location=[latitude, longitude], zoom_start=11)

# set color scheme for the clusters
x = np.arange(kclusters)
colors_array = cm.rainbow(x)
#colors_array = cm.rainbow(np.linspace(0, 1, len(ys)))
#rainbow = (colors.rgb2hex(i) for i in colors_array)
rainbow = ['yellow', 'blue', 'red', 'cyan']

# add markers to the map
markers_colors = []
for lat, lon, poi, cluster in zip(coffee_merged['Latitude'], coffee_merged['Longitude'], coffee_merged['Neighborhood'], coffee_merged['Cluster Labels']):
    label = folium.Popup(str(poi) + ' Cluster ' + str(cluster), parse_html=True)
    folium.CircleMarker(
        [lat, lon],
        radius=10,
        popup=label,
        color=rainbow[cluster-1],
        fill=True,
        fill_color=rainbow[cluster-1],
        fill_opacity=0.7).add_to(map_clusters)
map_clusters
```



## 5. Examine Clusters

```
In [47]: # minimum frequency of coffee shop
coffee_merged['Coffee Shop'].min()
```

```
Out[47]: 0.0
```

```
In [48]: # maximum frequency of coffee shop
coffee_merged['Coffee Shop'].max()
```

```
Out[48]: 0.09
```

Cluster 1

```
In [49]: coffee_merged.loc[coffee_merged['Cluster Labels'] == 0]
```

```
Out[49]:
Neighborhood Coffee Shop Cluster Labels Borough Latitude Longitude
0 Battery Park City 0.03 0 Manhattan 40.711932 -74.016869
1 Soho 0.03 0 Manhattan 40.722184 -74.000657
2 West Village 0.04 0 Manhattan 40.734434 -74.006180
3 Financial District 0.03 0 Manhattan 40.707107 -74.010665
4 East Harlem 0.03 0 Manhattan 40.792249 -73.944182
5 Tribeca 0.03 0 Manhattan 40.771522 -74.010683
6 Upper East Side 0.03 0 Manhattan 40.775608 -73.960508
7 Chelsea 0.03 0 Manhattan 40.744035 -74.003116
8 Civic Center 0.04 0 Manhattan 40.715229 -74.005415
9 Manhattan Valley 0.04 0 Manhattan 40.797307 -73.964286
```

Cluster 2

```
In [50]: coffee_merged.loc[coffee_merged['Cluster Labels'] == 1]
```

```
Out[50]:
Neighborhood Coffee Shop Cluster Labels Borough Latitude Longitude
10 Central Harlem 0.07 1 Manhattan 40.815976 -73.943211
11 Hamilton Heights 0.07 1 Manhattan 40.828304 -73.949688
12 Morningside Heights 0.08 1 Manhattan 40.808000 -73.963896
13 Manhattanville 0.09 1 Manhattan 40.816934 -73.957385
```

Cluster 3

In [51]: coffee\_merged.loc[coffee\_merged['Cluster\_Labels'] == 2]

Out[51]:

Cluster 4

In [52]: coffee\_merged.loc[coffee\_merged['Cluster\_Labels'] == 3]

Out[52]:

In [ ]:

In [53]: # locations that are not already crowded with coffee shops

low\_coffee = (coffee\_merged.loc[coffee\_merged['Cluster\_Labels'] == 2]).reset\_index(drop=True)

low\_coffee

Out[53]:

In [54]: # locations that are crowded with coffee shops

high\_coffee = (coffee\_merged.loc[coffee\_merged['Cluster\_Labels'] == 1]).reset\_index(drop=True)

high\_coffee

Out[54]:

In [ ]:

In [55]: from numpy import cos, sin, arcsin, sqrt

from math import radians

```
def haversine(row):
    lat1 = 40.787658
    lon1 = -73.977059
    lat2 = row['Latitude']
    lon2 = row['Longitude']
    lat1, lon1, lat2, lon2 = map(radians, [lat1, lon1, lat2, lon2])
    dlat = lat2 - lat1
    dlon = lon2 - lon1
    a = sin(dlon/2)**2 + cos(lat1) * cos(lat2) * sin(dlat/2)**2
    c = 2 * arcsin(sqrt(a))
    km = 6367 * c
    return km
```

low\_coffee['Upper West Side (km)'] = low\_coffee.apply(lambda row: haversine(row), axis=1)

low\_coffee

Out[55]:

In [56]: low\_coffee.sort\_values(['Upper West Side (km)'])[:5]

Out[56]:

In [57]: #minimum distance

distance\_min = low\_coffee.loc[low\_coffee['Coffee Shop'] == 0.00].reset\_index(drop=True)

distance\_min['Upper West Side (km)'].min()

Out[57]: 1.1334214594689875

In [58]: distance\_min['Upper West Side (km)'].idxmin()

Out[58]: 2

In [59]: high\_coffee['Upper West Side (km)'] = high\_coffee.apply(lambda row: haversine(row), axis=1)

high\_coffee

Out[59]:

In [60]: high\_coffee.loc[high\_coffee['Coffee Shop'] == 0.09].reset\_index(drop=True)

Out[60]:

## 6. Results

This work analyzed 40 neighborhoods in Manhattan with the objective of detecting places that are no longer full of coffee shops. The area with the lowest density of cafeterias in Montreal is concentrated in the peripheral region of the analyzed area. This region is related to cluster 3 (red dots on the map). In this cluster, 13 neighborhoods have a low or no cafeteria density. The neighborhoods where there is a high incidence of coffee shops are located in cluster 2, and they are appx 2-4 Km away to our central coordinate.

**Midtown** is the closest neighborhood to Upper West Side and has no coffee shop. It is 1.13 km away from the central coordinate. Tudor City and Turtle Bay are also the nearby places closer appx 1.4-1.5 Km with no coffee shops.

## 7. Conclusion

The neighborhoods where there is a high incidence of coffee shops are located in cluster 2, many of them are appx 2-4 Km away to the central coordinate. The neighborhoods with the highest coffee frequencies are Central Harlem, Hamilton Heights, Morningside Heights, Manhattanville. Project requirements include detecting places that are not already crowded with cafes, or areas with no cafes nearby, and as close as possible to the city center. The following three options are suggested that meet the project requirements : **Midtown**, the closest neighborhood and has no coffee shop with 1.13 km away from the central coordinate, Tudor City and Turtle Bay (1.4-1.5 Km).

In [ ]: