

Capstone Project : Battle of the Neighborhoods: Location - Pharmacy

1. Data Acquisition

Import required python libraries

```
In [1]: 1 import numpy as np # library to handle data in a vectorized manner
2
3 import pandas as pd # library for data analysis
4 pd.set_option('display.max_columns', None)
5 pd.set_option('display.max_rows', None)
6
7 import json # library to handle JSON files
8
9 #!conda install -c conda-forge geopy --yes # uncomment this line if you
10 from geopy.geocoders import Nominatim # convert an address into latitude
11
12 import requests # library to handle requests
13 from pandas.io.json import json_normalize # tranform JSON file into a p
14
15 # Matplotlib and associated plotting modules
16 import matplotlib.cm as cm
17 import matplotlib.colors as colors
18
19 # import k-means from clustering stage
20 from sklearn.cluster import KMeans
21
22 #!conda install -c conda-forge folium=0.5.0 --yes # uncomment this line
23 import folium # map rendering library
24
25 #!conda install -c conda-forge beautifulsoup4 --yes
26 import bs4
27 from bs4 import BeautifulSoup
28
29 print('Libraries imported.')
```

Libraries imported.

Clean data

Map and segment and cluster the neighborhoods in one borough. Create a new dataframe of the boroughs data.

```

In [2]: 1 request_data = requests.get('https://en.wikipedia.org/wiki/List_of_post
2 soup = BeautifulSoup(request_data, "html.parser")
3 table=soup.find('table')
4 #dataframe will consist of three columns: PostalCode, Borough, and Neig
5 df = pd.DataFrame(columns = ['PostalCode', 'Borough', 'Neighborhood'])
6 df.shape
7
8 # Search all the postcode, borough, neighborhood
9 for trcelldata in table.find_all('tr'):
10     trdata=[]
11     for tdcelldata in trcelldata.find_all('td'):
12         trdata.append(tdcelldata.text.strip())
13     if len(trdata)==3:
14         df.loc[len(df)] = trdata
15 df.dropna()
16 # drop Borough not assigned
17 df.drop(df[ df['Borough'] == 'Not assigned' ].index, inplace = True)
18
19 # assign Borough to Neighborhood if latter not assigned
20 df.Neighborhood[df.Neighborhood == 'Not assigned'] = df.Borough
21
22 # Create a second frame grouping by postal code with neighborhoods tran
23 # In the dataset this is already done so no need to do this, but doing
24 df1=df.groupby('PostalCode')['Neighborhood'].apply(lambda x: "%s" % ',
25 # Merge , drop the extra column not comma separated and rename to corre
26 df2 = pd.merge(df, df1, on='PostalCode').drop_duplicates().drop(['Neigh
27
28 df2.shape
29 # Rename column to match that of df
30 geospatial_df = pd.read_csv('http://cocl.us/Geospatial_data')
31 geospatial_df.rename(columns={'Postal Code' : 'PostalCode'}, inplace=True)
32 geospatial_df.head()
33 #Merge to get the required dataframe
34 df_neighborhood_geo = pd.merge(df2, geospatial_df, on='PostalCode')
35 df_neighborhood_geo.head()
36

```

Out[2]:

	PostalCode	Borough	Neighborhood	Latitude	Longitude
0	M3A	North York	Parkwoods	43.753259	-79.329656
1	M4A	North York	Victoria Village	43.725882	-79.315572
2	M5A	Downtown Toronto	Regent Park, Harbourfront	43.654260	-79.360636
3	M6A	North York	Lawrence Manor, Lawrence Heights	43.718518	-79.464763
4	M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government	43.662301	-79.389494

Visualize Functions for boroughs

```
In [3]: 1 def get_lat_long_borough_data(borough_name, borough_address):
2         bor_data = neighborhoods[neighborhoods['Borough'].str.contains(bor_
3
4         geolocator = Nominatim(user_agent="t_explorer")
5         location = geolocator.geocode(borough_address)
6         latitude = location.latitude
7         longitude = location.longitude
8         print('The geograpical coordinate of Borough: ' + borough_address +
9         return(bor_data, latitude, longitude)
10
11 # create map of borough using latitude and longitude values
12 def draw_folium_map(borr_data, latitude, longitude):
13     map_bor = folium.Map(location=[latitude, longitude], zoom_start=11)
14
15     # add markers to map
16     for lat, lng, label in zip(borr_data['Latitude'], borr_data['Longitud
17         label = folium.Popup(label, parse_html=True)
18         folium.CircleMarker(
19             [lat, lng],
20             radius=5,
21             popup=label,
22             color='blue',
23             fill=True,
24             fill_color='#3186cc',
25             fill_opacity=0.7,
26             parse_html=False).add_to(map_bor)
27     return map_bor
```

Visualize the borough the neighborhoods in it.

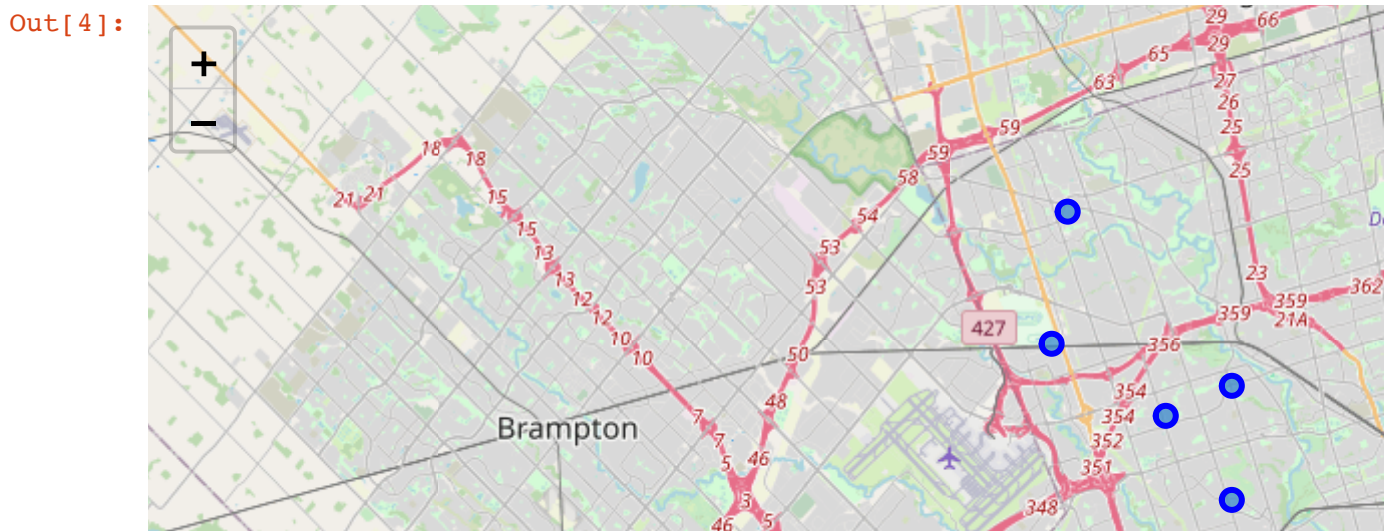
```

In [4]: 1 toronto_data = df_neighborhood_geo[df_neighborhood_geo.Borough.str.contains('Etobicoke')]
2 toronto_data.reset_index(drop=True, inplace=True)
3 toronto_data.head()
4 neighborhoods = toronto_data
5 neighborhoods['Neighborhood']
6 # hoods = neighborhoods['Neighborhood'].unique()
7 # n = 0
8 # for hood in hoods:
9 #     nr = hood.split(",")
10 #     n = n + (len(nr))
11 #     print (n)
12
13
14 hoods = neighborhoods['Borough'].unique()
15 for hood in hoods:
16     print (' ' + hood)
17
18 (bor_data, latitude, longitude ) = get_lat_long_borough_data('Etobicoke')
19
20 bor_data.head()
21 draw_folium_map(bor_data,latitude,longitude)

```

Etobicoke

The geographical coordinate of Borough: Etobicoke, Toronto are 43.6435559, -79.5656326.



Foursquare API Tool: To explore the neighborhoods and segment them.

Define Foursquare Credentials and Version

```

In [5]: 1 CLIENT_ID = 'xxxx' # your Foursquare ID
2 CLIENT_SECRET = 'xxxxx' # your Foursquare Secret for privacy xxxx
3 VERSION = '20180605' # Foursquare API version

```

Functions to help with processing venues

Functions to **get_category_type** and **get_med_nearby** to get medical nearby venues from the Foursquare lab.

```
In [6]: 1 LIMIT = 100 # limit of number of venues returned by Foursquare API
2 radius = 1000 # define radius
3
4 # function that extracts the category of the venue
5 def get_category_type(row):
6     try:
7         categories_list = row['categories']
8     except:
9         categories_list = row['venue.categories']
10
11     if len(categories_list) == 0:
12         return None
13     else:
14         return categories_list[0]['name']
15
16 def get_med_nearby(lat, lon, category, client_id, client_secret, version):
17     categories = category if isinstance(category, str) else ','.join(categories)
18     url = 'https://api.foursquare.com/v2/venues/search?client_id={}&client_secret={}&lat={}&lon={}&radius={}&version={}'.format(
19         client_id, client_secret, lat, lon, radius, version)
20     try:
21         results = requests.get(url).json()['response']['venues']
22         dataframe = json_normalize(results)
23         dataframe.head()
24         # keep only columns that include venue name, and anything that
25         filtered_columns = ['name', 'categories', 'location.lat', 'location.lng']
26         dataframe_filtered = dataframe.loc[:, filtered_columns]
27         # filter the category for each row
28         dataframe_filtered['categories'] = dataframe_filtered.apply(get_category_type, axis=1)
29     except:
30         dataframe_filtered = pd.DataFrame()
31         print(f'\nError on {url}')
32     return dataframe_filtered
33
```

2. Data Processing or Manipulation

Function to repeat the same process to all the neighborhoods in the borough for general and medical venues

```

In [7]: 1 def getNearbyVenues(names, latitudes, longitudes, radius=500):
2
3     venues_list=[]
4     for name, lat, lng in zip(names, latitudes, longitudes):
5         print(name)
6
7         # create the API request URL
8         url = 'https://api.foursquare.com/v2/venues/explore?&client_id=
9             CLIENT_ID,
10            CLIENT_SECRET,
11            VERSION,
12            lat,
13            lng,
14            radius,
15            LIMIT)
16
17         # make the GET request
18         results = requests.get(url).json()["response"]["groups"][0]["it
19
20         # return only relevant information for each nearby venue
21         venues_list.append([ (
22             name,
23             lat,
24             lng,
25             v['venue']['name'],
26             v['venue']['location']['lat'],
27             v['venue']['location']['lng'],
28             v['venue']['categories'][0]['name']) for v in results])
29
30     nearby_venues = pd.DataFrame([item for venue_list in venues_list fo
31     nearby_venues.columns = ['Neighborhood',
32                             'Neighborhood Latitude',
33                             'Neighborhood Longitude',
34                             'Venue',
35                             'Venue Latitude',
36                             'Venue Longitude',
37                             'Venue Category']
38
39     return(nearby_venues)
40
41
42 def getNearbyMedVenues(names, latitudes, longitudes, radius=500):
43     client_id = CLIENT_ID
44     client_secret = CLIENT_SECRET
45     version = VERSION
46     limit = LIMIT
47     categories = '4bf58dd8d48988d104941735,4bf58dd8d48988d10f951735'
48     venues_list=[]
49     for name, lat, lng in zip(names, latitudes, longitudes):
50         url = 'https://api.foursquare.com/v2/venues/search?client_id={}
51         results = requests.get(url).json()['response']['venues']
52         venues_list.append([ (
53             name,
54             lat,
55             lng,
56             v['name'],

```

```

57         v['location']['lat'],
58         v['location']['lng'],
59         v['categories'][0]['name']) for v in results])
60
61
62     nearby_med_venues = pd.DataFrame([item for venue_list in venues_list
63     nearby_med_venues.columns = ['Neighborhood',
64         'Neighborhood Latitude',
65         'Neighborhood Longitude',
66         'Venue',
67         'Venue Latitude',
68         'Venue Longitude',
69         'Venue Category']
70
71     return(nearby_med_venues)

```

Create a new dataframe called **bor_venues**.

```

In [8]: 1 # type your answer here
2
3 bor_venues = getNearbyMedVenues(names=bor_data['Neighborhood'],
4                                 latitudes=bor_data['Latitude'],
5                                 longitudes=bor_data['Longitude']
6                                 )
7 bor_venues

```

Out[8]:

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	Islington Avenue, Humber Valley Village	43.667856	-79.532242	Shoppers Drug Mart	43.663067	-79.531753	Pharmacy
1	Eringate, Bloordale Gardens, Old Burnhamthorpe...	43.643515	-79.577201	Shoppers Drug Mart	43.641312	-79.576924	Pharmacy
2	Eringate, Bloordale Gardens, Old Burnhamthorpe...	43.643515	-79.577201	Burnhamthorpe Health Centre	43.642328	-79.576959	Medical Centre
3	Eringate, Bloordale Gardens, Old Burnhamthorpe...	43.643515	-79.577201	Dr Henry Nirenberg Dental Office	43.641895	-79.578301	Dentist Office
	Eringate,						

Number of unique categories from all the returned venues

```
In [9]: 1 print('There are {} uniques categories.'.format(bor_venues['Venue Category']
2
```

There are ['Pharmacy' 'Medical Center' "Dentist's Office" "Doctor's Office"
 'Veterinarian' 'Chiropractor' 'Medical Lab' 'Supplement Shop' 'Hospital'] uniques categories.

```
In [10]: 1 print('There are {} uniques categories.'.format(len(bor_venues['Venue Category'])
2
```

There are 9 uniques categories.

3. Analyze Each Neighborhood

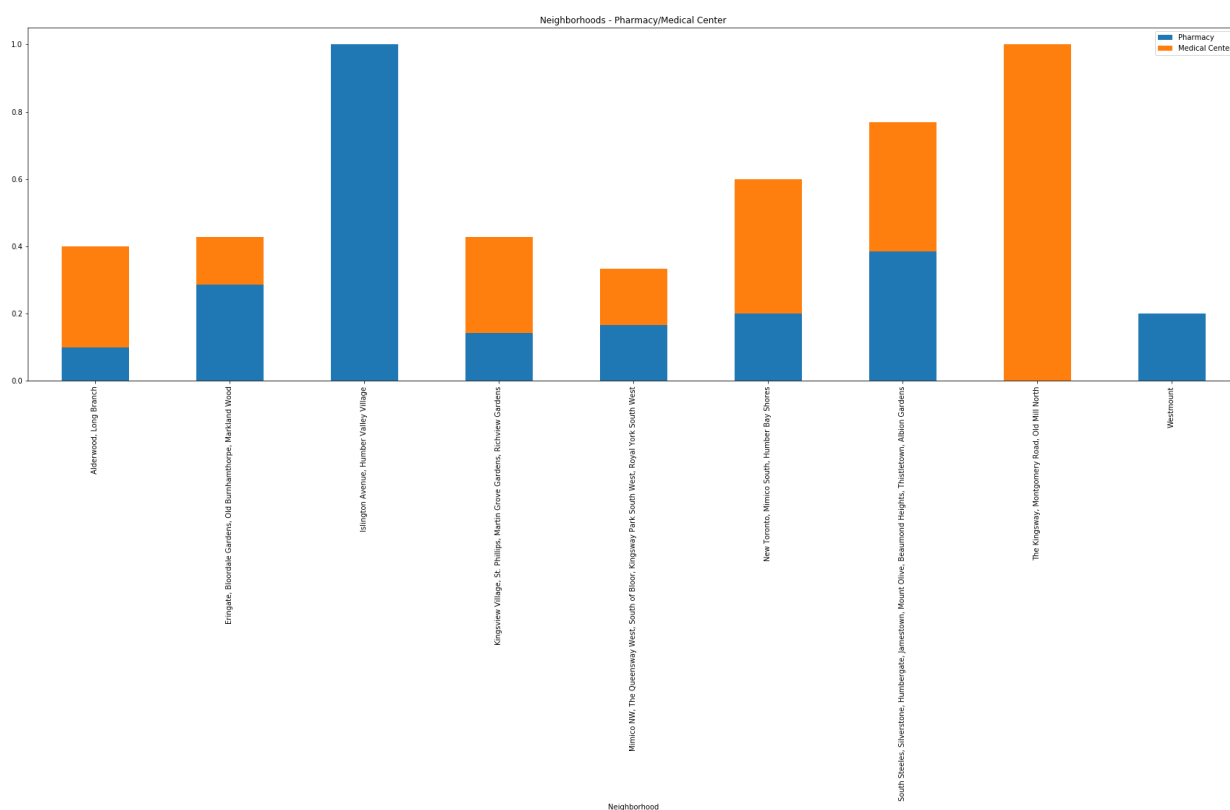
```
In [11]: 1 # one hot encoding
2 bor_allonehot = pd.get_dummies(bor_venues[['Venue Category']], prefix="
3
4 # add neighborhood column back to dataframe
5 bor_allonehot['Neighborhood'] = bor_venues['Neighborhood']
6 bor_allonehot.head()
```

Out[11]:

	Chiropractor	Dentist's Office	Doctor's Office	Hospital	Medical Center	Medical Lab	Pharmacy	Supplement Shop	Veterinarian
0	0	0	0	0	0	0	1	0	
1	0	0	0	0	0	0	1	0	
2	0	0	0	0	1	0	0	0	
3	0	1	0	0	0	0	0	0	
4	0	0	1	0	0	0	0	0	

Group rows by neighborhood and by taking the mean of the frequency of occurrence of each category


```
In [13]: 1 bor_onehot = bor_allonehot[['Neighborhood', 'Pharmacy']]
2 bor_grouped = bor_onehot.groupby('Neighborhood').mean().reset_index()
3 bor_grouped_count = bor_onehot.groupby('Neighborhood').count()
4
5 bor_dpeonehot = bor_allonehot[['Neighborhood', 'Pharmacy', 'Medical Center']]
6 bor_dpegroupped = bor_dpeonehot.groupby('Neighborhood').mean().reset_index()
7 bor_dpegroupped_count = bor_dpeonehot.groupby('Neighborhood').count()
8
9 import matplotlib.pyplot as plt
10 plt.rcParams["figure.figsize"] = [30,9]
11 plt.title('Neighborhoods - Pharmacy/Medical Center')
12 # # gca stands for 'get current axis'
13 ax = plt.gca()
14 bor_dpegroupped.plot(kind='bar',x='Neighborhood', stacked=True, ax=ax)
15 plt.show()
16
```



4. Cluster Neighborhoods

Run *k*-means to cluster the neighborhood into 3 clusters.

```
In [14]: 1 # neighborhoods_venues_sorted.drop('Cluster Labels', inplace=True)
2 neighborhoods_venues_sorted = bor_grouped
3 # neighborhoods_venues_sorted.head()
4 bor_grouped_clustering = bor_grouped.drop('Neighborhood', 1)
5
6 # set number of clusters
7 kclusters = 3
8 # run k-means clustering
9 kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(bor_grouped_c
10
11 # check cluster labels generated for each row in the dataframe
12 kmeans.labels_[0:10]
13 print(kmeans.labels_)

[0 2 1 0 0 0 2 0 0]
```

Create a new dataframe with cluster

```
In [15]: 1 # add clustering labels
2 neighborhoods_venues_sorted.insert(0, 'Cluster Labels', kmeans.labels_)
3
4 bor_merged = bor_data
5 # merge toronto_grouped with toronto_data to add latitude/longitude for
6 bor_merged = bor_merged.join(neighborhoods_venues_sorted.set_index('Nei
7 bor_merged.dropna(subset = ["Cluster Labels"], inplace=True)
```

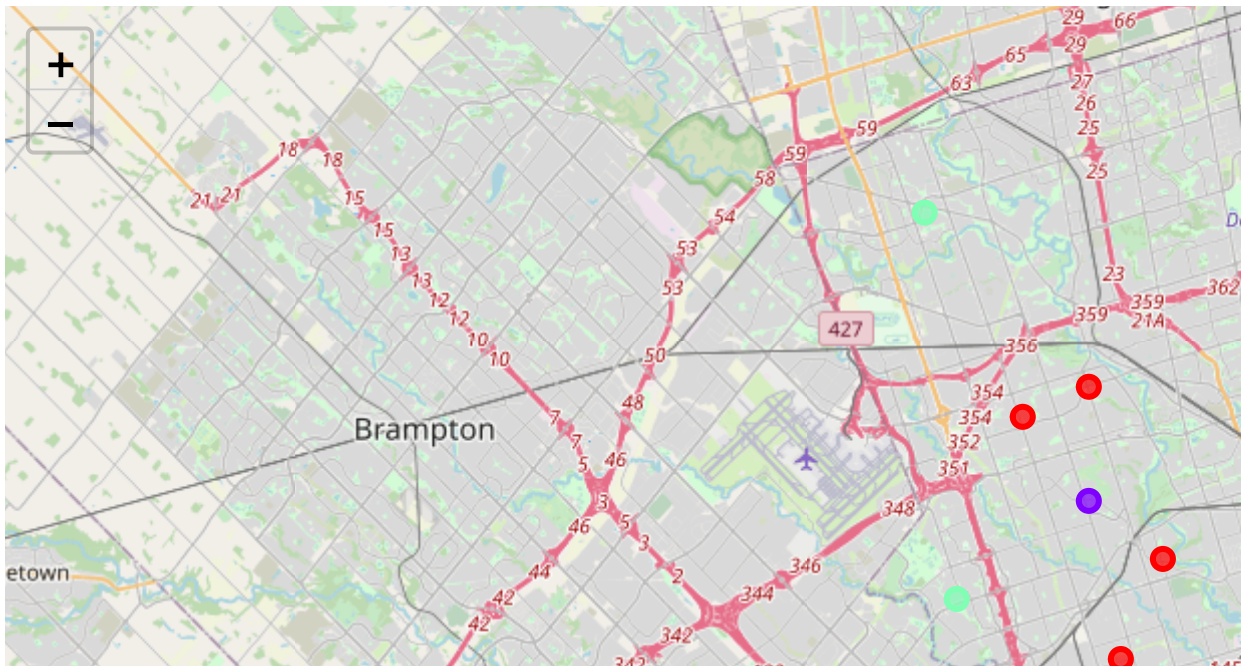
Visualize the resulting clusters

```

In [16]: 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(kclusters)
6 ys = [i + x + (i*x)**2 for i in range(kclusters)]
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, clusterf in zip(bor_merged['Latitude'], bor_merged['
13     cluster = int(clusterf)
14     label = folium.Popup(str(poi) + ' Cluster ' + str(cluster), parse_h
15     folium.CircleMarker(
16         [lat, lon],
17         radius=5,
18         popup=label,
19         color=rainbow[cluster-1],
20         fill=True,
21         fill_color=rainbow[cluster-1],
22         fill_opacity=0.7).add_to(map_clusters)
23
24 map_clusters

```

Out[16]:



5. Examine Clusters

Now, you can examine each cluster and determine the discriminating venue categories that distinguish each cluster. Based on the defining categories, you can then assign a name to each cluster. I will leave this exercise to you.

Cluster 0

```
In [17]: 1 cluster0 = bor_merged.loc[bor_merged['Cluster Labels'] == 0.0, bor_merged]
         2
```

Cluster 1

```
In [18]: 1 cluster1 = bor_merged.loc[bor_merged['Cluster Labels'] == 1.0, bor_merged]
         2 cluster1
```

Out[18]:

	Borough	Cluster Labels	Pharmacy
0	Etobicoke	1.0	1.0

Cluster 2

```
In [19]: 1 cluster2 = bor_merged.loc[bor_merged['Cluster Labels'] == 2.0, bor_merged]
         2
```

```
In [20]: 1 cluster0 = bor_merged.loc[bor_merged['Cluster Labels'] == 0.0, bor_merged]
         2
         3 cluster1 = bor_merged.loc[bor_merged['Cluster Labels'] == 1.0, bor_merged]
         4
         5 cluster2 = bor_merged.loc[bor_merged['Cluster Labels'] == 2.0, bor_merged]
         6
         7 def get_cluster_hoods(bor_merged, clusternum):
         8     return bor_merged.loc[bor_merged['Cluster Labels'] == clusternum]
         9
        10 cluster_stats1 = []
        11 for x in range(kclusters):
        12     dfcluster = get_cluster_hoods(bor_merged, x)
        13     print('Cluster:', x, 'Avg # Pharmacies:', dfcluster['Pharmacy'].mean())
        14     cluster_stats1.append([(
        15         x,
        16         dfcluster['Pharmacy'].mean(),
        17         dfcluster.shape[0])])
        18 cluster_stats1
```

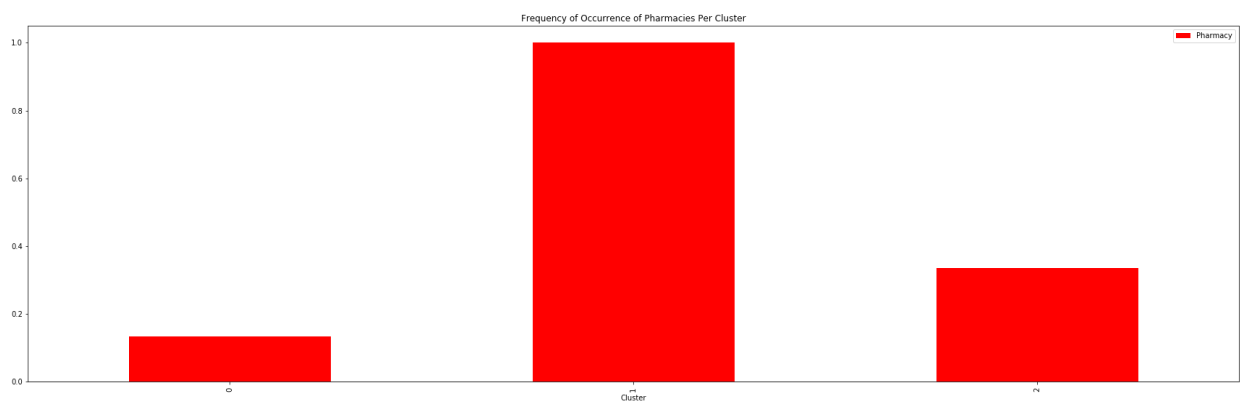
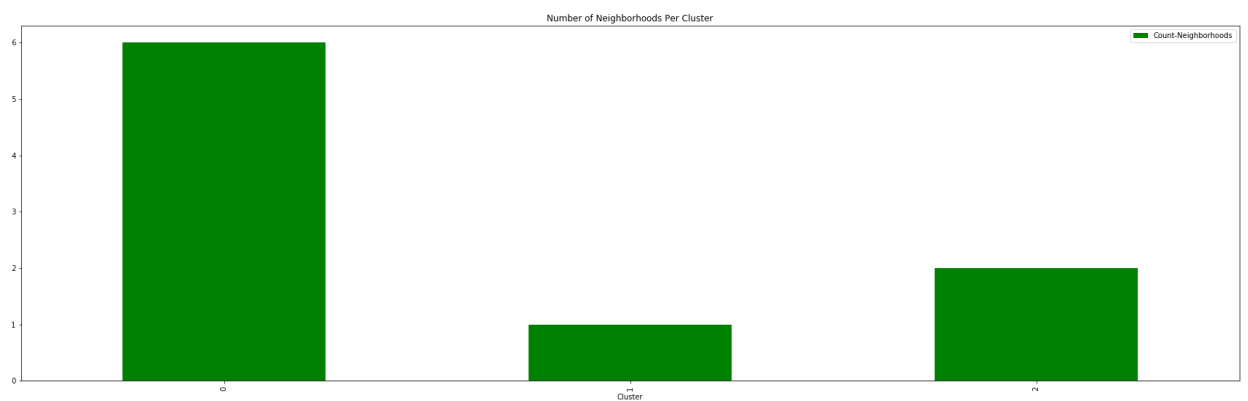
Cluster: 0 Avg # Pharmacies: 0.1349206349206349 #Neighborhoods: 6

Cluster: 1 Avg # Pharmacies: 1.0 #Neighborhoods: 1

Cluster: 2 Avg # Pharmacies: 0.3351648351648352 #Neighborhoods: 2

Out[20]: [[(0, 0.1349206349206349, 6)], [(1, 1.0, 1)], [(2, 0.3351648351648352, 2)]]

```
In [21]: 1 df = pd.DataFrame([item for x in cluster_stats1 for item in x], columns=  
2 df  
3 # gca stands for 'get current axis'  
4 ax = plt.gca()  
5 df.plot(kind='bar', x='Cluster', y='Count-Neighborhoods', stacked=True, color='green',  
6 plt.title('Number of Neighborhoods Per Cluster')  
7  
8 plt.show()  
9  
10 ax = plt.gca()  
11  
12 df.plot(kind='bar', x='Cluster', y='Pharmacy', stacked=True, color='red',  
13 plt.title('Frequency of Occurrence of Pharmacies Per Cluster')  
14  
15 plt.show()  
16
```



```
In [22]: 1 cluster0 = bor_merged.loc[bor_merged['Cluster Labels'] == 0.0]
2 cluster1 = bor_merged.loc[bor_merged['Cluster Labels'] == 1.0]
3 cluster2 = bor_merged.loc[bor_merged['Cluster Labels'] == 2.0]
4 cluster0
```

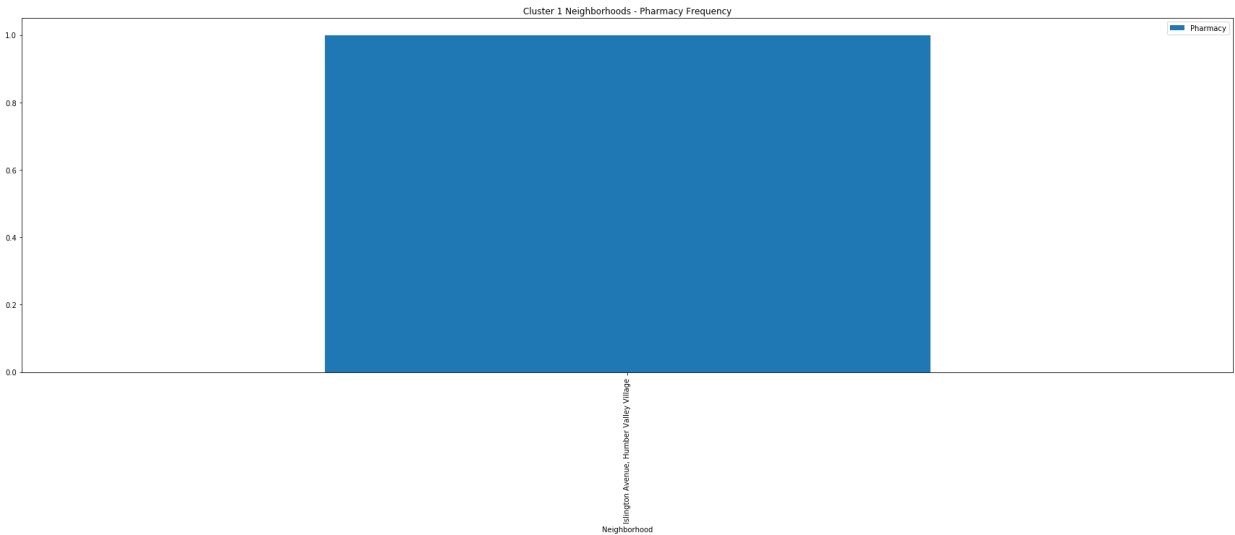
Out[22]:

	PostalCode	Borough	Neighborhood	Latitude	Longitude	Cluster Labels	Pharmacy
3	M9P	Etobicoke	Westmount	43.696319	-79.532242	0.0	0.200000
4	M9R	Etobicoke	Kingsview Village, St. Phillips, Martin Grove ...	43.688905	-79.554724	0.0	0.142857
5	M8V	Etobicoke	New Toronto, Mimico South, Humber Bay Shores	43.605647	-79.501321	0.0	0.200000
7	M8W	Etobicoke	Alderwood, Long Branch	43.602414	-79.543484	0.0	0.100000
9	M8X	Etobicoke	The Kingsway, Montgomery Road, Old Mill North	43.653654	-79.506944	0.0	0.000000
11	M8Z	Etobicoke	Mimico NW, The Queensway West, South of Bloor,...	43.628841	-79.520999	0.0	0.166667

Cluster 0 Neighborhoods (Low in Pharmacy from Modeling)

```
In [23]: 1 cluster1.plot(x ='Neighborhood', y='Pharmacy', kind = 'bar', title='Clu
2
```

Out[23]: <matplotlib.axes._subplots.AxesSubplot at 0x1a20f4a290>



Create data for neighborhoods with pharmacy and medical center in cluster 0

In [24]:

1

2 bor_dpeonehot_merge = pd.merge(bor_dpegrouped, cluster0, on='Neighborhood')

3 bor_dpeonehot_merge

4

Out[24]:

	Neighborhood	Pharmacy_x	Medical Center	PostalCode	Borough	Latitude	Longitude	Cluster Labels	Phi
0	Alderwood, Long Branch	0.100000	0.300000	M8W	Etobicoke	43.602414	-79.543484	0.0	
1	Kingsview Village, St. Phillips, Martin Grove ...	0.142857	0.285714	M9R	Etobicoke	43.688905	-79.554724	0.0	
2	Mimico NW, The Queensway West, South of Bloor,...	0.166667	0.166667	M8Z	Etobicoke	43.628841	-79.520999	0.0	
3	New Toronto, Mimico South, Humber Bay Shores	0.200000	0.400000	M8V	Etobicoke	43.605647	-79.501321	0.0	
4	The Kingsway, Montgomery Road, Old Mill North	0.000000	1.000000	M8X	Etobicoke	43.653654	-79.506944	0.0	
5	Westmount	0.200000	0.000000	M9P	Etobicoke	43.696319	-79.532242	0.0	

In []:

1