

Import libraries

```
In [1]: 1 import pandas as pd
2 import numpy as np
3 import requests
4 import warnings
5 warnings.filterwarnings('ignore')
6 import geopy
7 from geopy.geocoders import Nominatim
8 import geocoder
9 import folium
10 from sklearn.cluster import KMeans
11 import matplotlib.cm as cm
12 import matplotlib.colors as colors
13 from IPython.display import HTML, display
```

Set the url to table location and get the content of the page in variable

```
In [2]: 1 url = 'https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M'
2 page = requests.get(url).content
```

Use the pandas read_html function to read the table. In this case the first table on the page ([0])

```
In [3]: 1 df_raw = pd.read_html(page,header=0)[0]
```

```
In [4]: 1 df_raw.head()
```

Out[4]:

	Postcode	Borough	Neighbourhood
0	M1A	Not assigned	Not assigned
1	M2A	Not assigned	Not assigned
2	M3A	North York	Parkwoods
3	M4A	North York	Victoria Village
4	M5A	Downtown Toronto	Harbourfront

Let's clean up the dataset according to specifications, filtering and replacing values

```
In [5]: 1 df = df_raw[df_raw['Borough'] != 'Not assigned'] #filter
2 df.columns = ['PostalCode', 'Borough', 'Neighborhood'] #set columns
3 df.loc[df['Neighborhood'] == 'Not assigned', ['Neighborhood']] = df['Borough'] #replace Neighborhood with Borough if N
```

Group the dataframe and apply the string concatenation

```
In [6]: 1 df = df.groupby(by=['PostalCode', 'Borough']).agg(lambda x: ', '.join(set(x))).reset_index()
```

Finally let's show the final frame (103 records with 3 columns)

```
In [7]: 1 df.shape
```

Out[7]: (103, 3)

We have the dataframe setup, now let's find the longitude and latitude. I'm using the provided csv as the geocoder is malfunctioning

```
In [8]: 1 file = 'Geospatial_Coordinates.csv'
2 df_geo = pd.read_csv(file)
3 df_geo.rename(columns = {'Postal Code':'PostalCode'}, inplace = True)
```

We now have the location data per postalcode so now it will be joined together with the Toronto dataframe

```
In [9]: 1 df_toronto = pd.merge(df, df_geo, on='PostalCode', how='left')
```

This leaves with clean appended dataframe incl lat and long



```
In [10]: 1 df_toronto.head()
```

Out[10]:

	PostalCode	Borough	Neighborhood	Latitude	Longitude
0	M1B	Scarborough	Rouge, Malvern	43.806686	-79.194353
1	M1C	Scarborough	Rouge Hill, Highland Creek, Port Union	43.784535	-79.160497
2	M1E	Scarborough	Morningside, Guildwood, West Hill	43.763573	-79.188711
3	M1G	Scarborough	Woburn	43.770992	-79.216917
4	M1H	Scarborough	Cedarbrae	43.773136	-79.239476

Let's start with our first map and plot the locations on the Toronto map

We'll define the start zoom location of the folium map with the address and geolocator to get the latitude and longitude

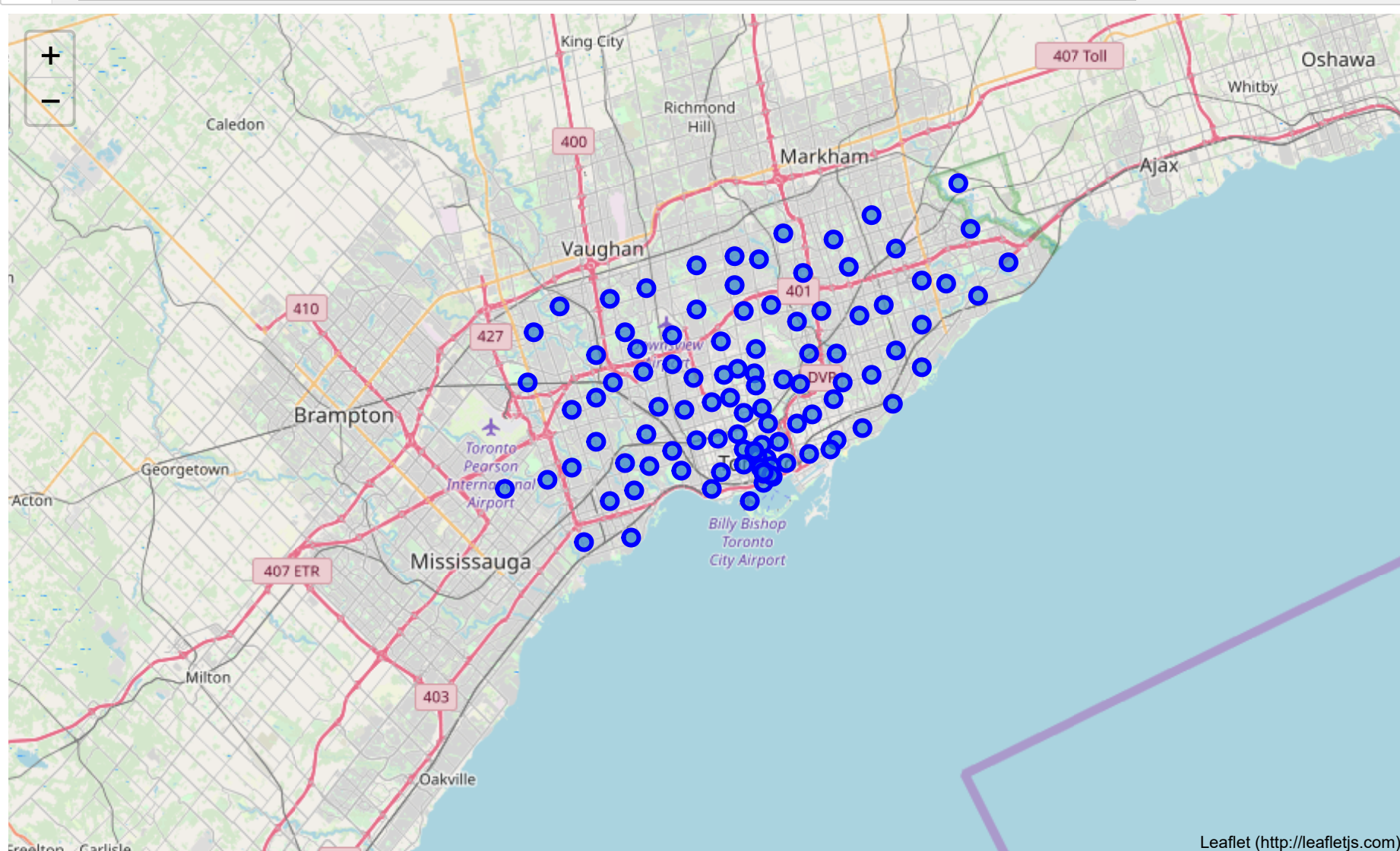
```
In [11]: 1 address = 'Toronto, Ontario'
2
3 geolocator = Nominatim(user_agent="ny_explorer")
4 location = geolocator.geocode(address)
5 latitude = location.latitude
6 longitude = location.longitude
7 print('The geographical coordinates of Toronto are {}, {}'.format(latitude, longitude))
```

The geographical coordinates of Toronto are 43.653963, -79.387207.

Now we plot the areas onto the Toronto map

```
In [12]: 1 # Function required to show folium maps inline
```

```
In [22]: 1 # create map of New York using Latitude and Longitude values
2 map_toronto = folium.Map(location=[latitude, longitude], zoom_start=10)
3
4 # add markers to map
5 for lat, lng, borough, neighborhood in zip(df_toronto['Latitude'], df_toronto['Longitude'], df_toronto['Borough'], df_toronto['Neighborhood']):
6     label = '{} {}'.format(neighborhood, borough)
7     popup = folium.Popup(label, parse_html=True)
8     marker = folium.CircleMarker(
9         [lat, lng],
10        radius=5,
11        popup=popup,
12        color='blue',
13        fill=True,
14        fill_color='blue',
15        fill_opacity=0.7).add_to(map_toronto)
16 map_toronto.save('toronto_map1.html')
17 display(map_toronto)
```



Now we have the initial map setup we continue with clustering. We will start with limiting the dataset as now we have 103 points on the map. We'll limit the set to only show locations for Borough's with the name Toronto in it

```
In [ ]: 1 toronto_data = df_toronto[df_toronto.Borough.str.contains('Toronto')].reset_index(drop=True)
        2 toronto_data.shape
```

This leaves us with 38 Boroughs we wil continue to work with. Let's map them again

```
In [ ]: 1 # create map of New York using Latitude and Longitude values
        2 map_toronto_filtered = folium.Map(location=[latitude, longitude], zoom_start=11)
        3
        4 # add markers to map
        5 for lat, lng, borough, neighborhood in zip(toronto_data['Latitude'], toronto_data['Longitude'], toronto_data['Borough'], toronto_data['Neighborhood']):
        6     label = '{} , {}'.format(neighborhood, borough)
        7     label = folium.Popup(label, parse_html=True)
        8     folium.CircleMarker(
        9         [lat, lng],
        10        radius=5,
        11        popup=label,
        12        color='blue',
        13        fill=True,
        14        fill_color='#3186cc',
        15        fill_opacity=0.7,
        16        parse_html=False).add_to(map_toronto_filtered)
        17
        18 map_toronto_filtered
```

Let's look at venues with foursquare to the locations we have

```
In [ ]: 1 # @hidden cell
        2 CLIENT_ID = '5LAMV3DNDHBER3VMUROMJNYRCJ5S35VSIB5BTJKJW2KHVG55' # your Foursquare ID
        3 CLIENT_SECRET = '3D5FLKF00A41T5XPDDIZTLLWTPIJWAMVQIU3AS5POKUEV1BW' # your Foursquare Secret
        4 VERSION = '20180605' # Foursquare API version
```

```
In [ ]: 1 # some variables we need for the below functions
        2
        3 LIMIT = 100 # limit of number of venues returned by Foursquare API
        4 radius = 500 # define radius
```

Function to get venues to process all the neighborhoods in Toronto

```
In [ ]: 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={}&client_secret={}&v={}&ll={},{}&radius={}&limit={}&nearby={}'.format(
        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]["items"]
        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 for item in venue_list])
        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)
```

```
In [ ]: 1 #Apply the function
2 toronto_venues = getNearbyVenues(names=toronto_data ['Neighborhood'],
3                                   latitudes=toronto_data ['Latitude'],
4                                   longitudes=toronto_data['Longitude'])
```

So we got 1700 records returned with 7 columns

```
In [ ]: 1 print(toronto_venues.shape)
2 toronto_venues.head()
```

We continue to analyse the neighborhoods

```
In [ ]: 1 # one hot encoding
2 toronto_onehot = pd.get_dummies(toronto_venues[['Venue Category']], prefix="", prefix_sep="")
3
4 # add neighborhood column back to dataframe
5 toronto_onehot['Neighborhood'] = toronto_venues['Neighborhood']
6
7 # move neighborhood column to the first column
8 fixed_columns = [toronto_onehot.columns[-1]] + list(toronto_onehot.columns[:-1])
9 toronto_onehot = toronto_onehot[fixed_columns]
10
11 toronto_onehot.head()
```

```
In [ ]: 1 toronto_onehot.shape
```

So got the categories converted to numerical values and transposed them into columns. We have 1700 records with 234 Venues categories

We group them by Neighborhood and that will leave us with 38 neighborhoods with 234 venue categories

```
In [ ]: 1 toronto_grouped = toronto_onehot.groupby('Neighborhood').mean().reset_index()
2 toronto_grouped.shape
```

That is a lot to process so we will get the top 10 venues for each neighborhood

```
In [ ]: 1 def return_most_common_venues(row, num_top_venues):
2     row_categories = row.iloc[1:]
3     row_categories_sorted = row_categories.sort_values(ascending=False)
4     return row_categories_sorted.index.values[0:num_top_venues]
```

```
In [ ]: 1 num_top_venues = 10
2
3 indicators = ['st', 'nd', 'rd']
4
5 # create columns according to number of top venues
6 columns = ['Neighborhood']
7 for ind in np.arange(num_top_venues):
8     try:
9         columns.append('{}{} Most Common Venue'.format(ind+1, indicators[ind]))
10    except:
11        columns.append('{}th Most Common Venue'.format(ind+1))
12
13 # create a new dataframe
14 neighborhoods_venues_sorted = pd.DataFrame(columns=columns)
15 neighborhoods_venues_sorted['Neighborhood'] = toronto_grouped['Neighborhood']
16
17 for ind in np.arange(toronto_grouped.shape[0]):
18     neighborhoods_venues_sorted.iloc[ind, 1:] = return_most_common_venues(toronto_grouped.iloc[ind, :], num_top_venues)
19
20 neighborhoods_venues_sorted.head(3)
```

We will now cluster the neighborhood with k-means into 5 clusters

```
In [ ]: 1 # set number of clusters
2 kclusters = 5
3
4 toronto_grouped_clustering = toronto_grouped.drop('Neighborhood', 1)
5
6 # run k-means clustering
7 kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(toronto_grouped_clustering)
8
9 # check cluster labels generated for each row in the dataframe
10 kmeans.labels_[0:10]
```

```
In [ ]: 1 kmeans_labels = kmeans.labels_
```


Let's add the clusters back to the neighborhoods and venues

```
In [ ]: 1 # add clustering labels
2 neighborhoods_venues_sorted.insert(0, 'Cluster Labels', kmeans_labels)
3
4 toronto_merged = toronto_data
5
6 # merge toronto_grouped with toronto_data to add Latitude/Longitude for each neighborhood
7 toronto_merged = toronto_merged.join(neighborhoods_venues_sorted.set_index('Neighborhood'), on='Neighborhood')
8
9 toronto_merged.head()
```

Finally, let's visualize the resulting clusters

```
In [ ]: 1 # create map
2 map_clusters = folium.Map(location=[latitude, longitude], zoom_start=12)
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, cluster in zip(toronto_merged['Latitude'], toronto_merged['Longitude'], toronto_merged['Neighborhood'], toronto_merged['Cluster Labels']):
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 map_clusters
```

The question then is what do the clusters represent. What is in those various clusters so we can name them better than Cluster 0-4

We will filter the toronto_merged frame into their respective variable so we analyse further

```
In [ ]: 1 label_0 = toronto_merged.loc[toronto_merged['Cluster Labels'] == 0, toronto_merged.columns[[1] + list(range(5, toronto_merged.columns.size-1))]
2 label_1 = toronto_merged.loc[toronto_merged['Cluster Labels'] == 1, toronto_merged.columns[[1] + list(range(5, toronto_merged.columns.size-1))]
3 label_2 = toronto_merged.loc[toronto_merged['Cluster Labels'] == 2, toronto_merged.columns[[1] + list(range(5, toronto_merged.columns.size-1))]
4 label_3 = toronto_merged.loc[toronto_merged['Cluster Labels'] == 3, toronto_merged.columns[[1] + list(range(5, toronto_merged.columns.size-1))]
5 label_4 = toronto_merged.loc[toronto_merged['Cluster Labels'] == 4, toronto_merged.columns[[1] + list(range(5, toronto_merged.columns.size-1))]
```

I'm not exactly sure how to find the common ground in the various clusters

```
In [ ]: 1 venues_columns = neighborhoods_venues_sorted.columns
2 venues_columns = venues_columns.drop(['Cluster Labels', 'Neighborhood'])
```

```
In [ ]: 1 label_1.head()
```

```
In [ ]: 1 label_4 # park
```

```
In [ ]: 1
```