

Capstone Project

This notebook will be used for the capstone project.

PART I. SCRAPPING

```
import pandas as pd
import numpy as np
```

```
"Hello Capstone Project Course!"
Hello Capstone Project Course!
```

```
# Initial dataframe with the postal codes
url = 'https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:M'
pc = pd.read_html(url)

pc_df = pc[0]

pc_df.head()
```

	Postal Code	Borough	Neighbourhood
0	M1A	Not assigned	Not assigned
1	M2A	Not assigned	Not assigned
2	M3A	Not assigned	Parkwoods
3	M4A	North York	Victoria Village
4	M5A	Downtown Toronto	Regent Park, Harbourfront

```
# Drop rows with Borough not assigned
pc_df.drop(pc_df[pc_df['Borough'] == 'Not assigned'].index, inplace=True)
pc_df.reset_index()
pc_df.head()
```

	Postal Code	Borough	Neighbourhood
2	M3A	North York	Parkwoods
3	M4A	North York	Victoria Village
4	M5A	Downtown Toronto	Regent Park, Harbourfront
5	M6A	North York	Lawrence Manor, Lawrence Heights
6	M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government

```
# If the neighbourhood is not assigned, we set its value equal to the district
for index, row in pc_df.iterrows():
    if row['Neighbourhood'] == 'Not assigned':
        pc_df.at[index, 'Neighbourhood'] = row['Borough']
pc_df.head()
```

	Postal Code	Borough	Neighbourhood
2	M3A	North York	Parkwoods
3	M4A	North York	Victoria Village
4	M5A	Downtown Toronto	Regent Park, Harbourfront
5	M6A	North York	Lawrence Manor, Lawrence Heights
6	M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government

```
# Rename the last column:
pc_df.rename(columns={'Neighbourhood': 'Neighborhood'}, inplace=True)
pc_df.head()
```

	Postal Code	Borough	Neighborhood
2	M3A	North York	Parkwoods
3	M4A	North York	Victoria Village
4	M5A	Downtown Toronto	Regent Park, Harbourfront
5	M6A	North York	Lawrence Manor, Lawrence Heights
6	M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government

```
pc_df.shape
```

```
(103, 3)
```

PART II. GEOSPACIAL DATA

```
import urllib.request
url = 'https://codel.us/geospatial_data'
filename = 'Geospatial_data.csv'
urllib.request.urlretrieve(url, filename)

# Geospatial data dataframe
gs_df = pd.read_csv(filename)
gs_df.head()
```

	Postal Code	Latitude	Longitude
0	MTB	43.806686	-79.194353
1	MTC	43.784535	-79.160497
2	MIE	43.763573	-79.188711
3	MIG	43.770992	-79.216917
4	MIH	43.773136	-79.239476

```
# merge the dataframes by Postal Code (inner join)
toronto_data = pd.merge(left=pc_df, right=gs_df, left_on='Postal Code', right_on='Postal Code', how='inner')
toronto_data.head()
```

	Postal Code	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

```
final_df.shape
```

```
(103, 5)
```

PART III. ANALYSIS

1. Explore Dataset

```
# remove districts that do not contain the word Toronto
toronto_data = final_df[final_df['Borough'].str.contains('Toronto', regex=False)].reset_index(drop=True)
toronto_data.head()
```

	Postal Code	Borough	Neighborhood	Latitude	Longitude
0	M5A	Downtown Toronto	Regent Park, Harbourfront	43.654260	-79.360636
1	M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government	43.662301	-79.389494
2	M5B	Downtown Toronto	Garden District, Ryerson	43.657162	-79.378937
3	M5C	Downtown Toronto	St James Town	43.651494	-79.375418
4	M4E	East Toronto	The Beaches	43.676357	-79.293031

```
toronto_data.shape
```

```
(39, 5)
```

Libraries needed

```
import json # library to handle JSON files

#conda install -c conda-forge geoppy --yes # uncomment this line if you haven't compiled
from geoppy.geocoders import Geocoder # convert an address into latitude and longitude

import requests # library to handle requests
from pandas.io.json import 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 # uncomment this line if you haven't
import folium # map rendering library
```

geographic coordinates of Toronto

```
address = 'Toronto, Canada'

geolocator = Nominatim(user_agent='toronto_explorer')
location = geolocator.geocode(address)
latitude = location.latitude
longitude = location.longitude
print('The geograpical coordinate of Toronto are {}, {}'.format(latitude, longitude))
```

The geographical coordinate of Toronto are 43.6534817, -79.3893947.

Neighbourhoods in Toronto

```
# create map of Toronto using latitude and longitude values
map_toronto = folium.Map(location=[latitude, longitude], zoom_start=11)

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

map_toronto
```

Make this Notebook Trusted to load map: File -> Trust Notebook

Credentials and version for the Foursquare API.

```
CLIENT_ID = '' # your Foursquare ID
CLIENT_SECRET = '' # your Foursquare Secret
VERSION = '20180605' # Foursquare API version
LIMIT = 100 # default Foursquare API limit value
```

2. Explore Neighborhoods in Toronto

```
def getNearbyVenues(names, latitudes, longitudes, radius=500):

    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={}&version={}&limit={}&nearby={}&radius={}&types={}&query={}'.format(
            CLIENT_ID,
            CLIENT_SECRET,
            VERSION,
            lat,
            lng,
            radius,
            LIMIT)

        # make the GET request
        results = requests.get(url).json()

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

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

    return(nearby_venues)
```

Create a new dataframe with the venue data

```
toronto_venues = getNearbyVenues(names=toronto_data['Neighborhood'],
                                 latitudes=toronto_data['Latitude'],
                                 longitudes=toronto_data['Longitude'])
```

```
toronto_venues.head()
```

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	Regent Park, Harbourfront	43.65426	-79.360636	Roselle Desserts	43.653447	-79.362017	Bakery
1	Regent Park, Harbourfront	43.65426	-79.360636	Tandem Coffee	43.653559	-79.361809	Coffee Shop
2	Regent Park, Harbourfront	43.65426	-79.360636	Cooper Koo Family YMCA	43.653249	-79.358008	Distribution Center
3	Regent Park, Harbourfront	43.65426	-79.360636	Body Blitz Spa East	43.654735	-79.359874	Spa
4	Regent Park, Harbourfront	43.65426	-79.360636	Impact Kitchen	43.656369	-79.356980	Restaurant

Size of the new dataframe

```
toronto_venues.shape
```

```
(1608, 7)
```

How many venues were returned for each neighborhood:

```
toronto_venues.groupby('Neighborhood').count()
```

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
	Neighborhood						
	Brocton, Parkdale Village, Exhibition Place	56	56	56	56	56	56
	Business reply mail Processing Centre, South Central Letter Processing Plant	22	22	22	22	22	22
	CN Tower, King and Spadina, Railway Lands, Harbourfront West, Bathurst Quay, South Niagara, Island airport	16	16	16	16	16	16
	Central Bay Street	62	62	62	62	62	62
	Christie	16	16	16	16	16	16
	Church and Wellesley	79	79	79	79	79	79
	Commerce Court, Victoria Hotel	100	100	100	100	100	100
	Davisville	36	36	36	36	36	36
	Davisville North	9	9	9	9	9	9
	Dufferin, Dovercourt Village	15	15	15	15	15	15
	First Canadian Place, Underground city	100	100	100	100	100	100
	Forest Hill North & West, Forest Hill Road Park	4	4	4	4	4	4
	Garden District, Ryerson	100	100	100	100	100	100
	Harbourfront East, Union Station, Toronto Islands	100	100	100	100	100	100
	High Park, The Junction South	26	26	26	26	26	26
	India Bazaar, The Beaches West	20	20	20	20	20	20
	Kensington Market, Chinatown, Grange Park	62	62	62	62	62	62
	Lawrence Park	4	4	4	4	4	4
	Little Portugal, Trinity	43	43	43	43	43	43
	Moore Park, Summerhill East	3	3	3	3	3	3
	North Toronto West, Lawrence Park	17	17	17	17	17	17
	Parkdale, Roncesvalles	14	14	14	14	14	14
	Queen's Park, Ontario Provincial Government	36	36	36	36	36	36
	Regent Park, Harbourfront	46	46	46	46	46	46
	Richmond, Adelaide, King	96	96	96	96	96	96
	Rosedale	4	4	4	4	4	4
	Roselawn	3	3	3	3	3	3
	Runnymede, Swansea	33	33	33	33	33	33
	St. James Town	79	79	79	79	79	79
	St. James Town, Cabbagetown	46	46	46	46	46	46
	Stn A PO Boxes	96	96	96	96	96	96
	Studio District	37	37	37	37	37	37
	Summerhill West, Rathnelly, South Hill, Forest Hill SE, Deer Park	14	14	14	14	14	14
	The Annex, North Midtown, Yorkville	19	19	19	19	19	19
	The Beaches	5	5	5	5	5	5
	The Danforth West, Riverdale	42	42	42	42	42	42
	Toronto Dominion Centre, Design Exchange	100	100	100	100	100	100
	University of Toronto, Harbord	32	32	32	32	32	32

How many unique categories can be curated from all the returned venues

```
print('There are {} unique categories.'.format(len(toronto_venues['Venue Category'].unique()))
There are 235 unique categories.
```

3. Analyze Each Neighborhood

```
# one hot encoding
toronto_onehot = pd.get_dummies(toronto_venues[['Venue Category']], prefix="", prefix_sep="",
                                columns=toronto_venues['Venue Category'].unique(),
                                dtype=int, drop_first=False)

# add neighborhood column back to dataframe
toronto_onehot['Neighborhood'] = toronto_venues['Neighborhood']

# move neighborhood column to the first column
fixed_columns = toronto_onehot.columns[1:] + list(toronto_onehot.columns[0])
toronto_onehot = toronto_onehot[fixed_columns]

toronto_onehot.head()
```

	Yoga Studio	Adult Boutique	Airport	Airport Food Court	Airport Gate	Airport Lounge	Airport Service	Airport Terminal	American Restaurant	Antique Shop	Theater
0	0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0	0

5 rows x 235 columns

Size of the dataframe

```
toronto_onehot.shape
```

```
(1608, 235)
```

Next, let's group rows by neighborhood and by taking the mean of the frequency of occurrence of each category

```
toronto_grouped = toronto_onehot.groupby('Neighborhood').mean().reset_index()
toronto_grouped
```

37	Dominion Centre, Exchange	0.000000	0.000000	0.0000	0.0000	0.0000	0.000	0.000	0.000	0.000	0.000	0.000
38	University of Toronto, Harbord	0.031250	0.000000	0.0000	0.0000	0.0000	0.0000	0.000	0.000	0.000	0.000	0.000

39 rows × 235 columns

```
toronto_grouped.shape
```

(39, 235)

New dataframe selecting only the relevant characteristics

Visualization

Number of elements per cluster and average of gyms per cluster

```
In [82]: %matplotlib inline
import matplotlib.pyplot as plt

listaMedias = []
listaCantidades = []
for n in range(0,6):
    listaMedias.append(final_gym_df.loc[final_gym_df['Cluster Labels'] == n]['Oym'].mean())
    listaCantidades.append(len(final_gym_df.loc[final_gym_df['Cluster Labels'] == n]))

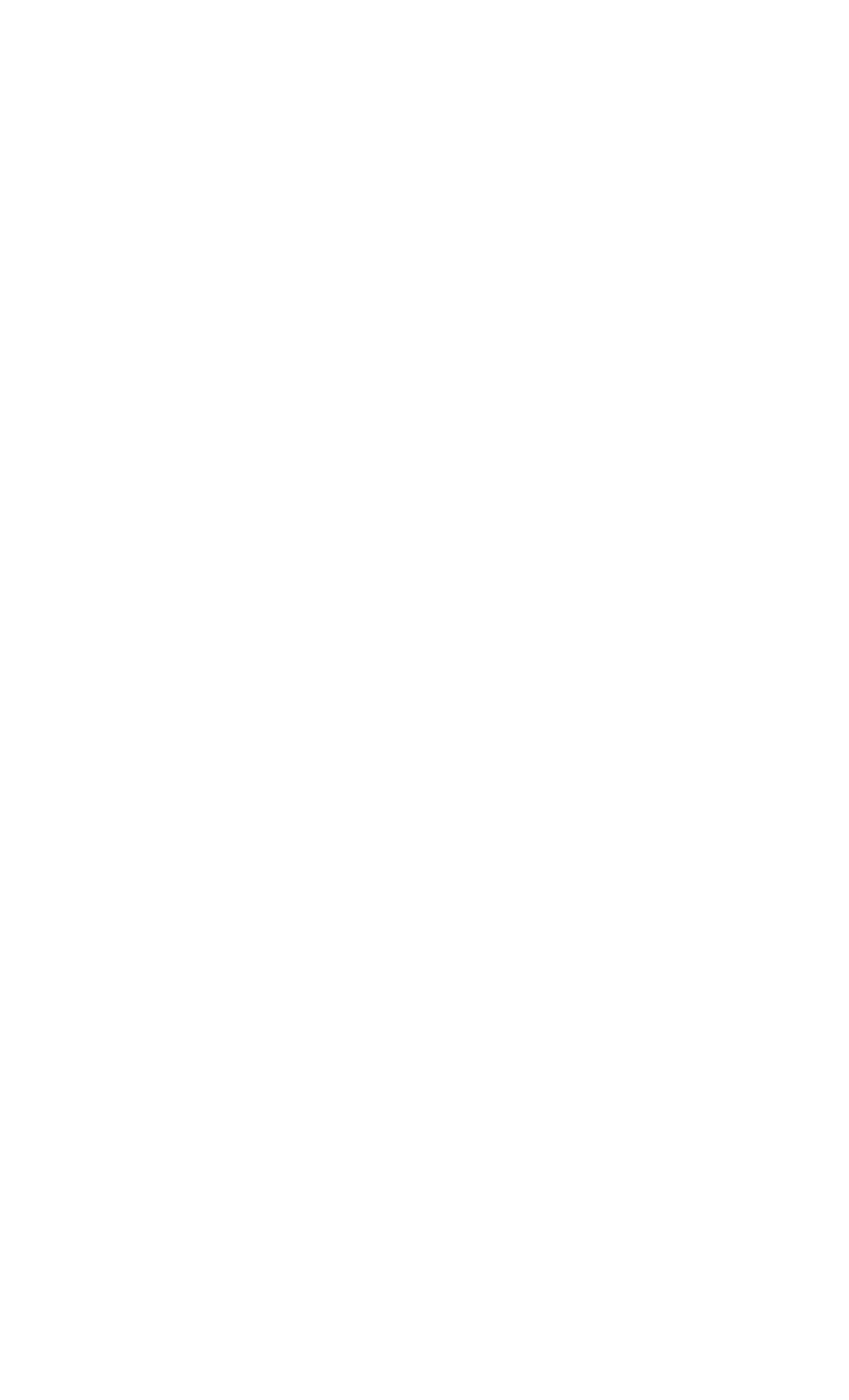
print(listaMedias)
print(listaCantidades)

[0.0, 0.029776936026936027, 0.04181818181818182, 0.01, 0.02307489451476793, 0.05277777777777778]
[27, 3, 3, 2, 2, 2]
```

Plot the average number of gyms per cluster



Plot the number of elements per cluster



In [] :