



Capstone Project - The Battle of Neighbourhoods (Week 2)

By Phiraphon Phankoson

Introduction

The travel agency would like to open a travel package to bring tourist into Toronto City so the travel agency have to understand in this city for creating travel package and finding area where able to be response to tourist by that area should have restaurants, cafe, museum and shopping mall.

For traveler that want to visit, travel or person who have to be in Toronto.

Data

To consider the objective, I listed the below data sources used for the analysis.

a) Toronto Neighborhood Data: The following Wikipedia page was scraped to pull out the necessary information: https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M

The information obtained i.e. the table of postal codes was transformed into a pandas data frame for further analysis.

b) Coordinate data for each Neighborhood in Toronto: The following csv file gave us the geographical coordinates of each postal code: http://cocl.us/Geospatial_data

Methodology

1. Load Toronto Neighborhood Data and transformed into a pandas data frame for further analysis.

```
1 import pandas as pd
2 import numpy as np
3 import requests
4 from bs4 import BeautifulSoup
```

Load Toronto Neighborhood Data

```
1 url = "https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M"
2 source = requests.get(url).text
3 Canada_data = BeautifulSoup(source, 'lxml')
```

Transform to pandas data frame

```
1 column_names = ['Postalcode', 'Borough', 'Neighborhood']
2 toronto = pd.DataFrame(columns = column_names)
```

```
1 content = Canada_data.find('div', class_='mw-parser-output')
2 table = content.table.tbody
3 postcode = 0
4 borough = 0
5 neighborhood = 0
6
7 for tr in table.find_all('tr'):
8     i = 0
9     for td in tr.find_all('td'):
10         if i == 0:
11             postcode = td.text
12             i = i + 1
13         elif i == 1:
14             borough = td.text
15             i = i + 1
16         elif i == 2:
17             neighborhood = td.text.strip('\n').replace(']', '')
18             toronto = toronto.append({'Postalcode': postcode, 'Borough': borough, 'Neighborhood': neighborhood})
19
20 toronto = toronto[toronto.Borough != 'Not assigned']
21 toronto = toronto[toronto.Borough != 0]
22 toronto.reset_index(drop = True, inplace = True)
23 i = 0
24 for i in range(0,toronto.shape[0]):
25     if toronto.iloc[i][2] == 'Not assigned':
26         toronto.iloc[i][2] = toronto.iloc[i][1]
27         i = i+1
28
29 df = toronto.groupby(['Postalcode', 'Borough'])['Neighborhood'].apply(', '.join).reset_index()
30 df.head()
```

	Postalcode	Borough	Neighborhood
0	M1B	Scarborough	Rouge, Malvern
1	M1C	Scarborough	Highland Creek, Rouge Hill, Port Union
2	M1E	Scarborough	Guildwood, Morningside, West Hill
3	M1G	Scarborough	Woburn
4	M1H	Scarborough	Cedarbrae

Methodology

2. Load Coordinate data for each Neighborhood in Toronto and put them to data frame.

```
1 file = 'http://cocl.us/Geospatial_data'
2 df_lo = pd.read_csv(file)
3 df_lo.head()
```

	Postal Code	Latitude	Longitude
0	M1B	43.806686	-79.194353
1	M1C	43.784535	-79.160497
2	M1E	43.763573	-79.188711
3	M1G	43.770992	-79.216917
4	M1H	43.773136	-79.239476

3. Merge 2 data frame together

```
1 df.rename(columns={'Postalcode':'Postal Code'}, inplace=True)
2 df_lo.set_index("Postal Code")
3 df.set_index("Postal Code")
4 toronto_data=pd.merge(df,df_lo)
5 toronto_data.head()
```

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

Methodology

4. Explore and cluster the neighborhoods in Toronto

```
1 from geopy.geocoders import Nominatim
2 import matplotlib.cm as cm
3 import matplotlib.colors as colors
4 from sklearn.cluster import KMeans
5
6 !conda install -c conda-forge folium=0.5.0 --yes
7 import folium

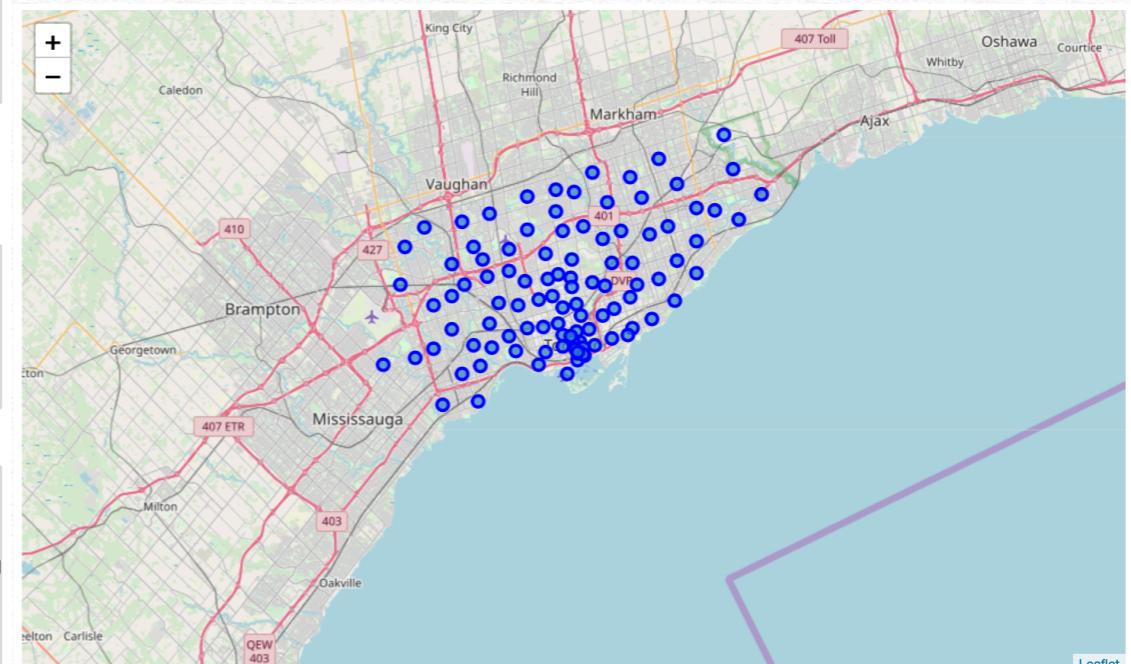
Collecting package metadata (current_repotdata.json): done
Solving environment: done

# All requested packages already installed.
```

```
1 address = 'Toronto, ON, Canada'
2
3 geolocator = Nominatim(user_agent="to_explorer")
4 location = geolocator.geocode(address)
5 latitude = location.latitude
6 longitude = location.longitude
7 print('The geographical coordinate of Toronto, ON, Canada are {}, {}'.format(latitude, longitude))

The geographical coordinate of Toronto, ON, Canada are 43.653963, -79.387207.
```

```
1 # create map of Toronto 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 in zip(toonto_data['Latitude'], toonto_data['Longitude'], toonto_data['Borough']):
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)
17
18 map_toronto
```



Methodology

5. Define Foursquare.

```
1 # defining radius and limit of venues to get
2 radius=500
3 LIMIT=100

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_se
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 ve
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

1 toronto_venues = getNearbyVenues(names=toronto_data['Neighborhood'],
2                                 latitudes=toronto_data['Latitude'],
3                                 longitudes=toronto_data['Longitude']
4                                 )
```

Let's check the size of the resulting dataframe

```
1 print(toronto_venues.shape)
2 toronto_venues.head()
```

(2220, 7)

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	Rouge, Malvern	43.806686	-79.194353	Wendy's	43.807448	-79.199056	Fast Food Restaurant
1	Highland Creek, Rouge Hill, Port Union	43.784535	-79.160497	Royal Canadian Legion	43.782533	-79.163085	Bar
2	Highland Creek, Rouge Hill, Port Union	43.784535	-79.160497	Scarborough Historical Society	43.788755	-79.162438	History Museum
3	Guildwood, Morningside, West Hill	43.763573	-79.188711	Swiss Chalet Rotisserie & Grill	43.767697	-79.189914	Pizza Place
4	Guildwood, Morningside, West Hill	43.763573	-79.188711	G & G Electronics	43.765309	-79.191537	Electronics Store

Methodology

6. Analysing Each Neighbourhood

```

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

```

	Accessories Store	Afghan Restaurant	Airport	Airport Food Court	Airport Lounge	Airport Service	Airport Terminal	American Restaurant	Antique Shop	Aquarium
0	0	0	0	0	0	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0
2	0	0	0	0	0	0	0	0	0	0
3	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0

And let's examine the new dataframe size.

```

1 toronto_onehot.shape
(2220, 268)

```

```

1 toronto_grouped = toronto_onehot.groupby('Neighborhood').mean().reset_index()
2 toronto_grouped.head()

```

Neighborhood	Accessories Store	Afghan Restaurant	Airport	Airport Food Court	Airport Lounge	Airport Service	Airport Terminal	American Restaurant
0 Adelaide, King, Richmond	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.02
1 Agincourt	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00
2 Agincourt North, L'Amoreaux East, Milliken, St...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00
3 Albion Gardens, Beaumont Heights, Humbergate, ...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00
4 Alderwood, Long Branch	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.00

Let's print each neighborhood along with the top 5 most common venues

```

1 num_top_venues = 5
2
3 for hood in toronto_grouped['Neighborhood']:
4     print("----"+hood+"----")
5     temp = toronto_grouped[toronto_grouped['Neighborhood'] == hood].T.reset_index()
6     temp.columns = ['venue','freq']
7     temp = temp.iloc[1:]
8     temp['freq'] = temp['freq'].astype(float)
9     temp = temp.round({'freq': 2})
10    print(temp.sort_values('freq', ascending=False).reset_index(drop=True).head(num_top_venues))
11    print('\n')

```

```

----Adelaide, King, Richmond----
      venue freq
0   Coffee Shop 0.06
1       Bar 0.04
2      Café 0.04
3  Steakhouse 0.04
4  Thai Restaurant 0.04

```

Methodology

7. Take result to data frame.

```
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
5     return row_categories_sorted.index.values[0:num_top_venues]
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
```

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

	Neighborhood	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue	6th Most Common Venue	7th Most Common Venue	8th Most Common Venue
0	Adelaide, King, Richmond	Coffee Shop	Bar	Café	Thai Restaurant	Steakhouse	Restaurant	Bakery	Re
1	Agincourt	Clothing Store	Latin American Restaurant	Skating Rink	Lounge	Breakfast Spot	Falafel Restaurant	Event Space	
2	Agincourt North, L'Amoreaux East, Milliken, St...	Playground	Park	Yoga Studio	Dumpling Restaurant	Discount Store	Dog Run	Doner Restaurant	
3	Albion Gardens, Beaumont Heights, Humbergate, ...	Grocery Store	Pizza Place	Fried Chicken Joint	Coffee Shop	Beer Store	Pharmacy	Fast Food Restaurant	S
4	Alderwood, Long Branch	Pizza Place	Pub	Pharmacy	Gym	Sandwich Place	Coffee Shop	Skating Rink	

Methodology

8. Answer the question of this project.

Result

The travel agency can create the travel package from catogorie of 1st Most Common Venue

```
1 toronto_merged['1st Most Common Venue'].unique()

array(['Fast Food Restaurant', 'History Museum', 'Spa', 'Coffee Shop',
       'Caribbean Restaurant', 'Playground', 'Discount Store', 'Bakery',
       'American Restaurant', 'College Stadium', 'Indian Restaurant',
       'Clothing Store', 'Pizza Place', 'Dog Run', 'Caf  ', 'Cafeteria',
       'Piano Bar', 'Ramen Restaurant', 'Park', 'Gym',
       'Miscellaneous Shop', 'Airport', 'Grocery Store', 'Baseball Field',
       'Liquor Store', 'Pub', 'Sporting Goods Shop', 'Greek Restaurant',
       'Pet Store', 'Dessert Shop', 'Summer Camp', 'Sandwich Place',
       'Garden', 'Mexican Restaurant', 'Bar', 'Airport Lounge',
       'Metro Station', 'Tennis Court', 'Gift Shop', 'Light Rail Station',
       'Furniture / Home Store', 'Rental Car Location'], dtype=object)
```

I assume that the travel package is travel Museum

```
1 toronto_merged.loc[toronto_merged['1st Most Common Venue'] == 'History Museum']
```

Postal Code	Borough	Neighborhood	Latitude	Longitude	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue
1 M1C	Scarborough	Highland Creek, Rouge Hill, Port Union	43.784535	-79.160497	0	History Museum	Bar	Yoga Studio

Result

If travel package is travel Museum, the travel agency bring tourist to Scarborough and walk around the Highland Creek, Rouge Hill and Port Union.

This project show about if someone would like to use data to creating the travel package, you can do that!

Conclusion

Have to set travel package first and then will able to plan path for tourist

The accuracy of data depends purely depends on the data provided by FourSquare

Thank you.