

Capstone Project

The Battle of Neighborhoods Opening
Chinese Restaurant in Toronto

Carol Wang

Business Understanding/Problem Description

- **Objective:** In this project, we will find the best neighborhood in Toronto for opening a new Chinese restaurant.
- There are some questions we need to address:
 - People live in which neighborhood likely eat in this kind of restaurant?
 - How many "similar" restaurants are available nearby?
 - Do the "similar" restaurants cost more? If so, what specialty do that have?

Data Source

- List of postal code for communities in Toronto from Wikipedia:
- https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M
- Geospatial Coordinates CSV file for Toronto postal codes from http://cocl.us/Geospatial_data
- Foursquare API venues explore methods to get venues of given neighborhoods
- Foursquare API venues methods to get ranks and likes of restaurants by given venue id

Methodology

- I will use of demographic data by neighborhoods, which allows me to gauge people's taste in each neighborhood. For example, I assume areas with a majority of Chinese people would be prefer to go to Chinese restaurant.
- I will use the Foursquare API to explore neighborhoods in Toronto and apply the explore function to get the most common venue categories in each of the neighborhoods.
- I will then use this feature to group the neighborhoods into clusters and use the k-means clustering algorithm to refine the groupings.
- Finally, I will use the Folium library to visualize the neighborhoods in Toronto and their emerging clusters as well as push the right neighborhood for setting up the restaurant business.

Data Cleaning:

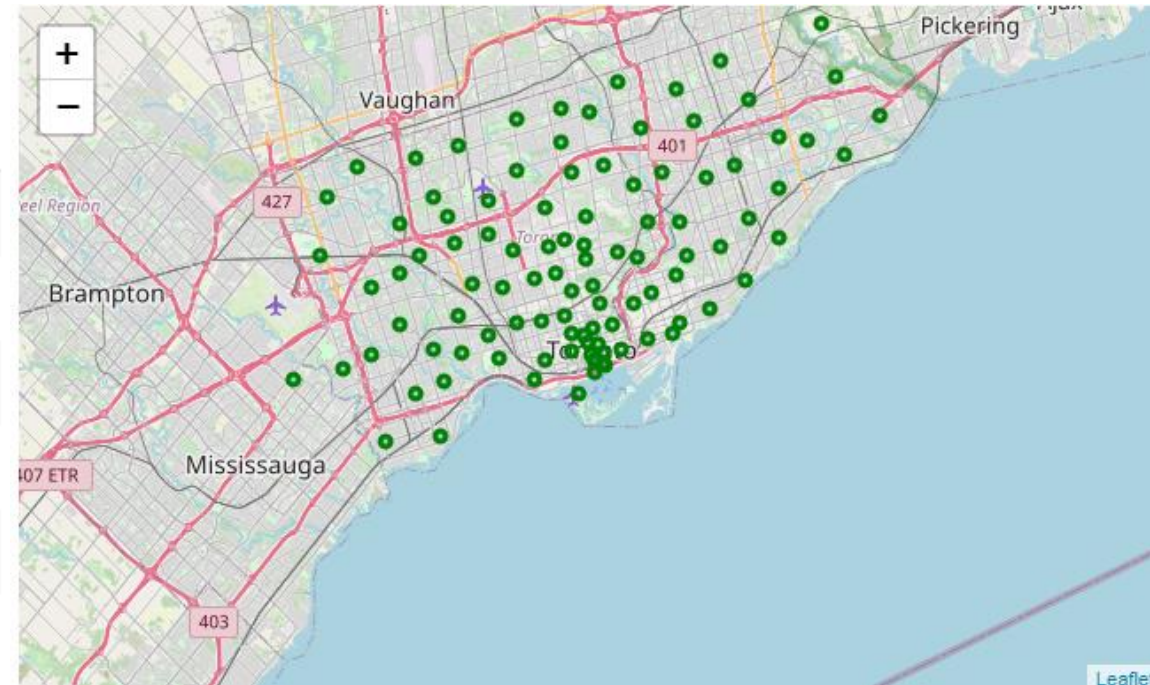
- I will use BeautifulSoup library to scrape Wikipedia to list of postal codes for each community. Remove the empty data in each row and group each zip code with all neighborhoods:

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

Data Cleaning:

- Then geospatial data (coordinates) is read into pandas dataframe. Then joined with the previous dataframe to get a new dataframe with zip code, neighborhood, latitude, and longitude coordinates.

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



Foursquare API

- I will use Foursquare API to explore neighborhoods and merge it to dataframe including neighborhoods and top-rated venues.

	neighborhood	neighborhood latitude	neighborhood longitude	venue	venue latitude	venue longitude	venue category
0	Malvern, Rouge	43.806686	-79.194353	Wendy's	43.807448	-79.199056	Fast Food Restaurant
1	Rouge Hill, Port Union, Highland Creek	43.784535	-79.160497	Royal Canadian Legion	43.782533	-79.163085	Bar
2	Rouge Hill, Port Union, Highland Creek	43.784535	-79.160497	SEBS Engineering Inc. (Sustainable Energy and ...	43.782371	-79.156820	Construction & Landscaping
3	Guildwood, Morningside, West Hill	43.763573	-79.188711	RBC Royal Bank	43.766790	-79.191151	Bank
4	Guildwood, Morningside, West Hill	43.763573	-79.188711	G & G Electronics	43.765309	-79.191537	Electronics Store

Features:

- When extracting features from a dataset, it is often useful to transform categorical features into vectors. Then group them by neighborhood by taking the mean of the frequency of occurrence.

	neighborhood	Accessories Store	Afghan Restaurant	Airport	Airport Food Court	Airport Gate	Airport Lounge	Airport Service	Airport Terminal	American Restaurant	Antique Shop	Ac
0	Agincourt	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	
1	Alderwood, Long Branch	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	
2	Bathurst Manor, Wilson Heights, Downsview North	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	
3	Bayview Village	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.000000	0.0	
4	Bedford Park, Lawrence Manor East	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.038462	0.0	

	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	9th Most Common Venue	10th Most Common Venue
0	Agincourt	Lounge	Latin American Restaurant	Skating Rink	Clothing Store	Breakfast Spot	Doner Restaurant	Dim Sum Restaurant	Diner	Discount Store	Distribut Cen
1	Alderwood, Long Branch	Pizza Place	Coffee Shop	Gym	Skating Rink	Pharmacy	Pub	Sandwich Place	Pool	Yoga Studio	Dess Sh
2	Bathurst Manor, Wilson Heights, Downsview North	Bank	Coffee Shop	Pharmacy	Pizza Place	Restaurant	Deli / Bodega	Ice Cream Shop	Fried Chicken Joint	Frozen Yogurt Shop	Mid East Restaurant
3	Bayview Village	Café	Bank	Japanese Restaurant	Chinese Restaurant	Deli / Bodega	Dessert Shop	Dim Sum Restaurant	Diner	Discount Store	Distribut Cen
4	Bedford Park, Lawrence Manor East	Sandwich Place	Coffee Shop	Italian Restaurant	Restaurant	Thai Restaurant	Pharmacy	Pizza Place	Indian Restaurant	Pub	C

Clustering

- In clustering, I try to find homogeneous subgroup within our data such that data points in each cluster are as similar as possible according to a similarity measure such as Euclidean-based distance or correlation-based distance.
- Clustering is considered an unsupervised learning method since we don't have the ground truth to compare the output of the clustering algorithm to the true labels to evaluate its performance. In this project, I use K-means clustering, which is a type of unsupervised learning, which is used when have unlabeled data. The K-means algorithm is an iterative algorithm that tries to partition the dataset into K pre-defined distinct non-overlapping subgroups where each data point belongs to only one group. In this case, the clustering is based on types of nearby venues.

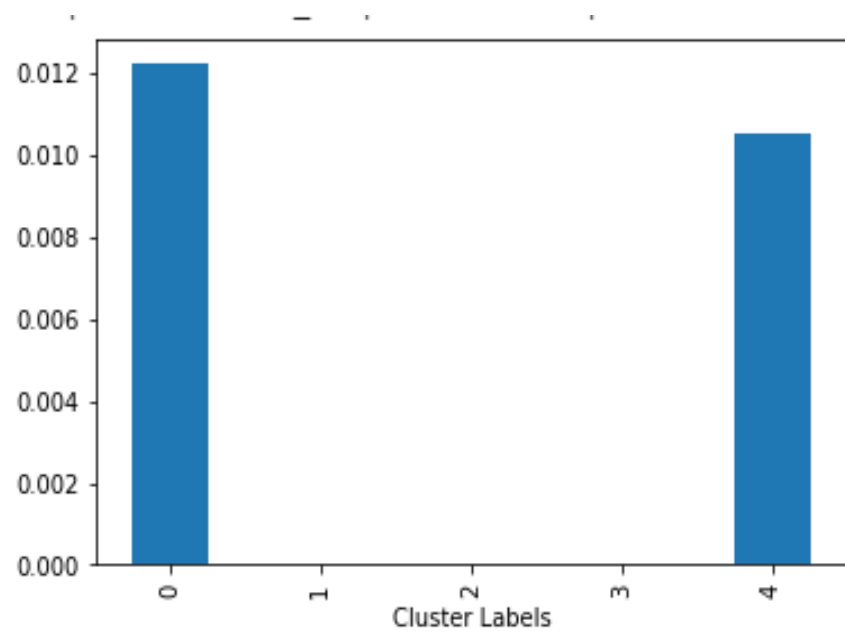
Clustering

- Then data is grouped into 5 clusters. The clusters and the top 10 venues for each neighborhood are shown below:

	Borough	Neighbourhood	Postcode	Latitude	Longitude	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
0	Scarborough	Malvern, Rouge	M1B	43.806686	-79.194353	4.0	Fast Food Restaurant	Yoga Studio	Dog Run	Department Store	
1	Scarborough	Rouge Hill, Port Union, Highland Creek	M1C	43.784535	-79.160497	4.0	Bar	Construction & Landscaping	Yoga Studio	Doner Restaurant	Re
2	Scarborough	Guildwood, Morningside, West Hill	M1E	43.763573	-79.188711	4.0	Electronics Store	Mexican Restaurant	Restaurant	Rental Car Location	E
3	Scarborough	Woburn	M1G	43.770992	-79.216917	1.0	Coffee Shop	Korean Restaurant	Yoga Studio	Dog Run	
4	Scarborough	Cedarbrae	M1H	43.773136	-79.239476	4.0	Fried Chicken Joint	Gas Station	Hakka Restaurant	Bakery	



- Then dataset is grouped by clusters to get the means of the frequency occurrence of each cluster, and examined by a heatmap.



Conclusion

- Conclusions can be taken from the previous maps, tables, and exploratory data analysis. For example, cluster 1 and 4 has high frequency of Chinese restaurants, so try not to open a new Chinese restaurant in cluster 1 or 4 neighborhoods.

