

Capstone Project - The Battle of Neighborhoods

Danilo Evangelista (<https://github.com/danilodn/>)

1. Notes

Capstone project report for the “Applied Data Science Capstone” course from the “IBM Data Science Professional Certificate” on Coursera (<https://www.coursera.org/professional-certificates/ibm-data-science>).

2. Introduction/Business Problem

Stakeholders want to open a new coffee shop branch either on Vancouver or Toronto and need the subsidiary information on the city/competitors to decide about the expanding strategy.

Their intention is to open one Café shop right away and elaborate a long-term strategy to allocate several stores over the next years in these two cities.

The plan should consider the competitors’ location and neighborhood characteristics.

3. Dataset

The first set of data used was the list of boroughs and neighborhoods of Toronto and Vancouver, extracted from Wikipedia (https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M and https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_V), as seen in the following images:



Postal Code	Borough	Neighborhood
M1A	Not assigned	Not assigned
M2A	Not assigned	Not assigned
M3A	North York	Parkwoods
M4A	North York	Victoria Village
M5A	Downtown Toronto	Regent Park, Harbourfront
M6A	North York	Lawrence Manor, Lawrence Heights
M7A	Downtown Toronto	Queen's Park, Ontario Provincial Government
M8A	Not assigned	Not assigned
M9A	Etobicoke	Islington Avenue, Humber Valley Village
M1B	Scarborough	Malvern, Rouge
M2B	Not assigned	Not assigned

Toronto's list of Postal Cods, Boroughs and Neighborhoods

Vancouver's list of Postal Cods, Boroughs and Neighborhoods

We also retrieved the list of venues for both cities from Foursquare API, each limited to 100 occurrences from 1000 meters of each neighborhood centroid, which resulted in 1772 different category venues for Vancouver and 4951 for Toronto.

From Wikipedia we managed to get data through Pandas method “pd.read_html”.

```
In [237]: link_V = "https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_V"
df_V = pd.read_html(link_V)[0]
df_V.head()
```

2

After that, we used Pgeocode to get the coordinates for each postal code/neighborhood, resulting in the following DataFrame for each city:

```
[201]:
```

	PostalCode	Neighborhood	latitude	longitude
0	M3A	North York (York Heights / Victoria Village / ...	43.7545	-79.3300
1	M4A	North York (Sweeney Park / Wigmore Park)	43.7276	-79.3148
2	M5A	Downtown Toronto (Regent Park / Port of Toronto)	43.6555	-79.3626
3	M6A	North York (Lawrence Manor / Lawrence Heights)	43.7223	-79.4504
4	M7A	Queen's Park Ontario Provincial Government	43.6641	-79.3889

Then we run Foursquare API, getting the top 100 venues from 1000 meters from each postal code/neighborhood centroid coordinate, along with its name, category, and coordinates, resulting in another DataFrame for each city:

```
[178]:
```

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	North York (York Heights / Victoria Village / ...	43.7545	-79.33	Allwyn's Bakery	43.759840	-79.324719	Caribbean Restaurant
1	North York (York Heights / Victoria Village / ...	43.7545	-79.33	Brookbanks Park	43.751976	-79.332140	Park
2	North York (York Heights / Victoria Village / ...	43.7545	-79.33	Tim Hortons	43.760668	-79.326368	Café
3	North York (York Heights / Victoria Village / ...	43.7545	-79.33	A&W	43.760643	-79.326865	Fast Food Restaurant
4	North York (York Heights / Victoria Village / ...	43.7545	-79.33	Shoppers Drug Mart	43.760857	-79.324961	Pharmacy

As our focus was Coffee Shops and Cafés, we created one more DataFrame with only these venue types.

The Toronto DataFrames resulted in 4951 rows (all venues) and 601 rows (coffee shops and cafés) and the Vancouver resulted in 1779 rows (all venues) and 162 rows (coffee shops and cafés).

```
[202]: print(toronto_venues.shape)
print(toronto_coffee.shape)

(4951, 7)
(601, 7)
```

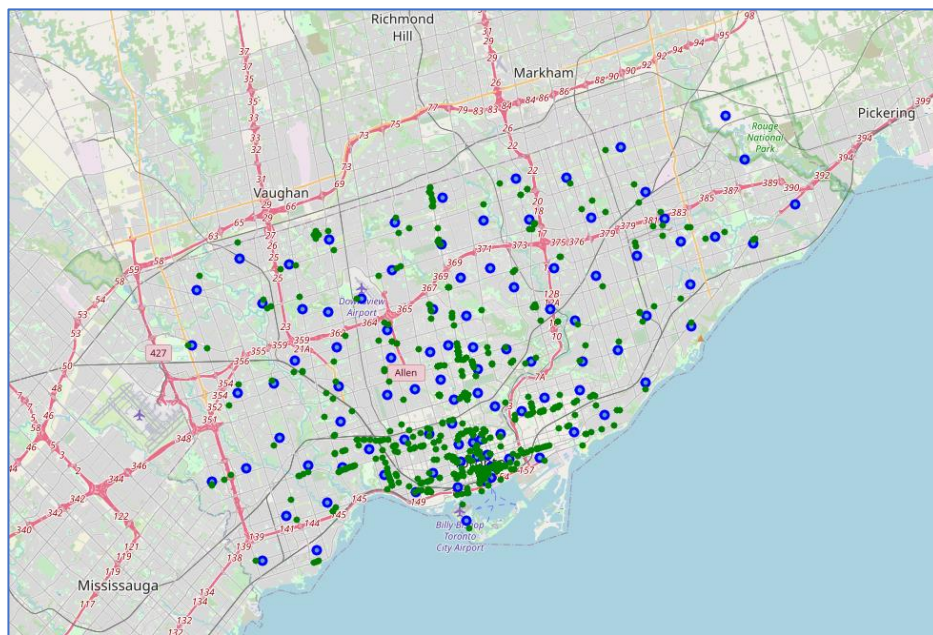
```
[203]: print(vancouver_venues.shape)
print(vancouver_coffee.shape)

(1779, 7)
(162, 7)
```

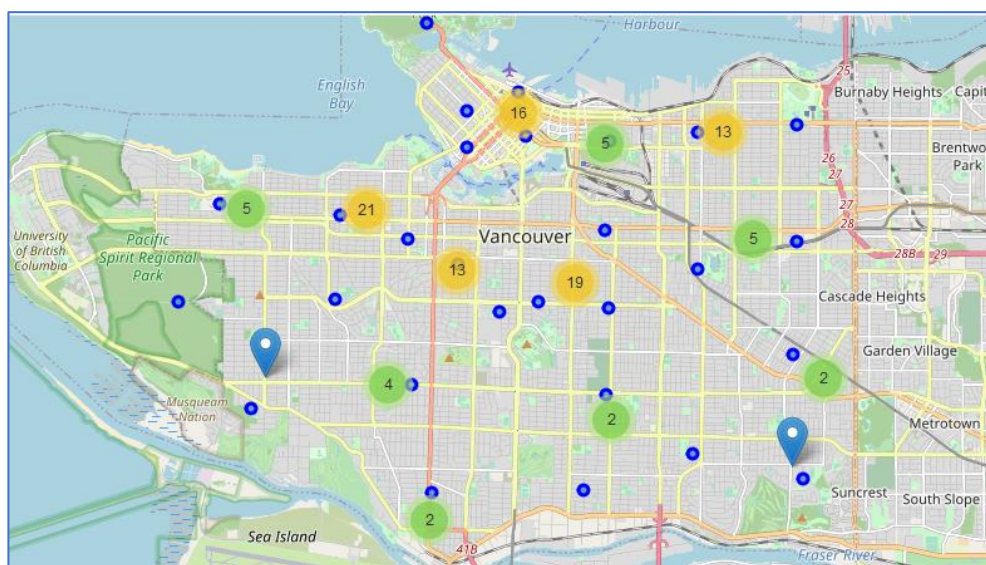
4. Methodology

With the DataFrames ready, first we used Python Folium to plot the neighborhoods in a map, to get a perspective of the cities and data.

Then we managed to superimpose the coffee shops and cafés in the same maps, to get a grip of the distribution of these venues over the cities, as we can see:



Also, we created a dynamic map using MarkerCluster, from Folium, which expands the venues when we zoom in the map and concentrate it when we zoom out, showing the number of venues per region. For example:



4.1. Clustering

We used the K-means clustering Classification algorithm, done with Python scikit-learn, to classify the neighborhoods in clusters of similarity or dissimilarity, to understand their characteristics, based on the venues populated there.

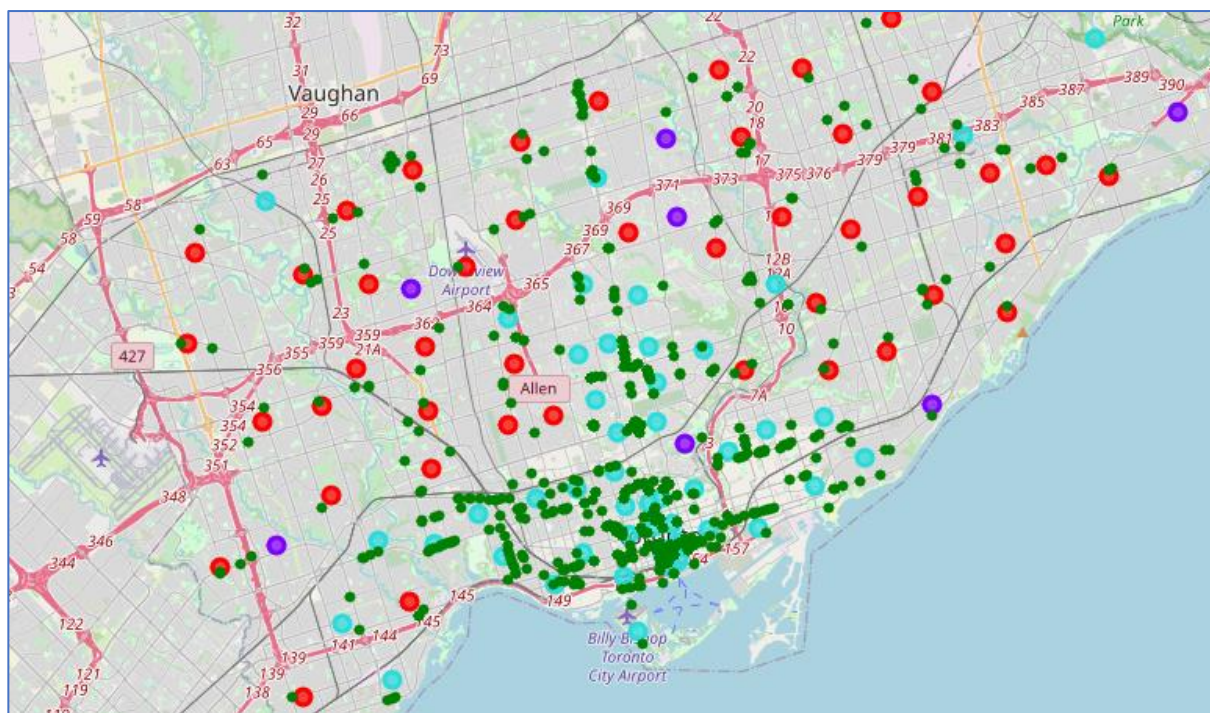
To do that, we used one-hot encoding to get dummy data and grouped neighborhoods, getting the top 10 venues for each neighborhood, for each city.

[124]:

	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	Central Toronto (Davisville North)	Coffee Shop	Italian Restaurant	Pizza Place	Café	Park	Sushi Restaurant	Restaurant	Yoga Studio	Diner	Pub
1	Central Toronto (Davisville)	Sushi Restaurant	Italian Restaurant	Pizza Place	Indian Restaurant	Café	Coffee Shop	Bakery	Sandwich Place	Bank	Bar
2	Central Toronto (Forest Hill North & West)	Park	Coffee Shop	Café	Bank	Sushi Restaurant	Italian Restaurant	Pharmacy	Gym / Fitness Center	Trail	Bagel Shop
3	Central Toronto (Lawrence Park East)	Sushi Restaurant	Italian Restaurant	Coffee Shop	Bakery	Bus Line	Café	Fast Food Restaurant	Pub	Bank	Asian Restaurant
4	Central Toronto (Moore Park / Summerhill East)	Coffee Shop	Sushi Restaurant	Italian Restaurant	Grocery Store	Thai Restaurant	Gym	Park	Gastropub	Spa	Bank

Trying some values for k, the best values were 3 for Toronto and 4 for Vancouver, grouping this way the neighborhoods according to its venue's categories.

Then, we created maps with the clustered neighborhoods and finally maps with the clusters and coffee/café's venues superimposed on it (green dots), which gave a better sense of all data together.

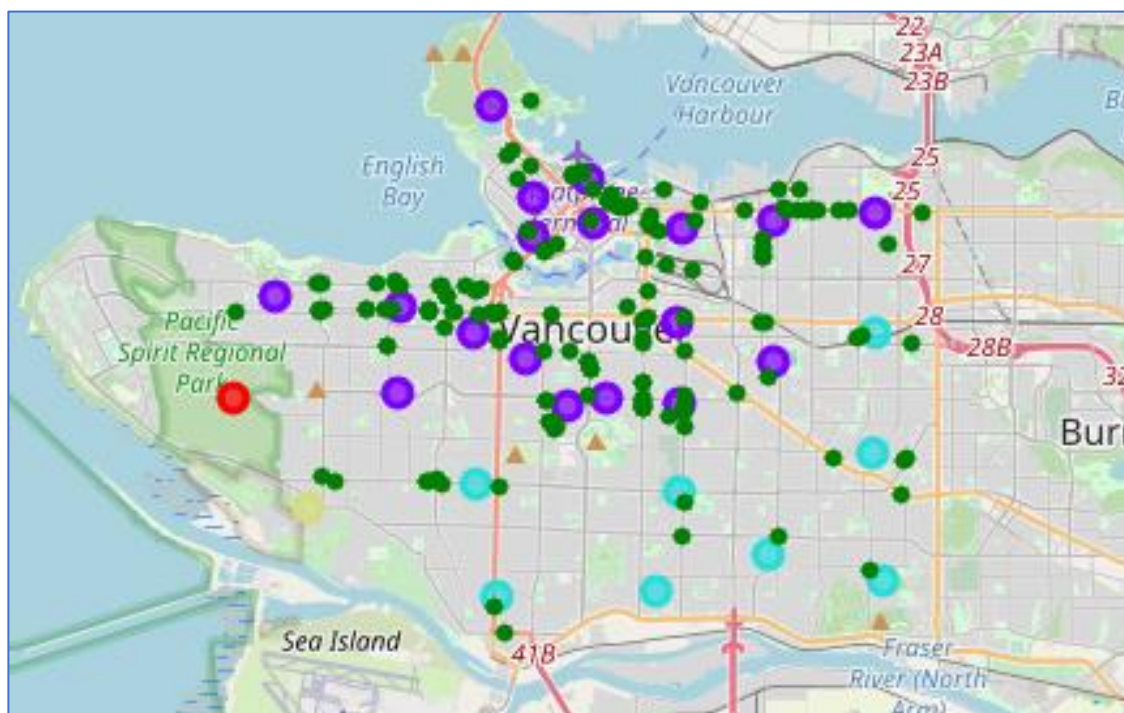


Resulting clusters for Toronto are:

- 0 (red): residence neighborhoods, predominantly with bus stops, parks, grocery stores, pharmacies, and banks, along a few restaurants, malls, and cafés.
- 1 (purple): a few sparse neighborhoods with distinct venues, as parks, zoo, soccer field, baseball field, golf course, theater, etc.
- 2 (cyan): neighbors with downtown characteristics with a lot of coffee shops/cafés, restaurants, and hotels.

Vancouver:

- 0 (red): one neighborhood situating the Pacific Spirit Regional Park.
- 1 (purple): neighbors with downtown characteristics with a lot of coffee shops/cafés, restaurants, and a few hotels (less than Toronto).
- 2 (cyan): residence neighborhoods, predominantly with bus stops, parks, restaurants, and other sparse venues.
- 3 (yellow): one neighborhood situating a golf course.



5. Results

From this Project we managed to understand Toronto and Vancouver neighborhood's characteristics, according to their venue vocation, and the distribution of coffee shops/cafés in the cities.

We classified satisfactorily the neighborhoods in clusters of similarity based on its venues and venue categories.

Also provided dynamic maps to the stakeholders, making the results visually appealing and easily understandable, making it very appropriate for the business decision making.

Must we say that this project can be used in the future too, updating the status of the neighborhoods and its venues at any time needed.

6. Discussion

The scope of this report could be extended by researching population of neighborhoods to be used together with the number of venues for each one.

It could have considered individual income from each city too, but it escapes the predetermined objective.

Also, we could have done some histograms, bar charts, area plots, etc. to help in the visualization of the information.

7. Conclusion

It is an open field and Foursquare API can be used for an infinite number of applications, since government policies to businesses decision to individual interests.

The bigger challenges are to get enough good data to use, and have the time to explore everything, but we must have in mind the objective, without escaping focus of the project.