

IBM – Coursera
Data Science Capstone Course

Capstone Project – Final Report
Neighborhood clustering in Toronto, Manhattan and Paris

Mauricio Stefanel - 2019

I. Introduction and Business Problem

This report is made for the final course of the applied Data Science Specialization, this has 4 courses created by IBM and Coursera. In this project the learner has the right to decide the topic where would be needed to leverage the Foursquare location data to solve or execute.

The main goal will be exploring the neighborhoods of New York City, Toronto city and Paris in order to define with neighborhoods are related bearing in mind the most common venues close to each neighborhood.

The idea comes from a family whose want to move having similar venue close to their home. This will work for a family that wants to move from New York, Toronto or Paris, to any of the neighborhoods of these cities.

So, can we define which neighborhoods are more alike to move on without much changes in the closest venues? If so, which are the neighborhoods more likely to have the same venues

The target audience for this report are:

- Potential families in New York, Toronto, or Paris that want to move to a similar neighborhood in those cities.
- Entrepreneurs who wants to open a new venue having in mind the lack or excess of similar venues in each neighborhood.
- Learners who would be interested in Clustering and a location foursquare application

II. Data description

To consider the problem we can list the data as below:

- The Toronto Neighborhood Coordinates are obtained from the List of Postal Codes of Canada (https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M) and the coordinates of each postal code (https://cocl.us/Geospatial_data), clean the data that don't have Borough or Neighborhood, group by Boroughs and Postal code. Then, the Postal Code is matched between the first data and the postal code coordinates, to finally get the coordinates of each neighborhood. And then finally takes only the Center of Toronto

```
In [68]: toronto_data.head()
```

Out[68]:

	PostalCode	Borough	Neighborhood	Latitude	Longitude
0	M4E	East Toronto	The Beaches	43.676357	-79.293031
1	M4K	East Toronto	The Danforth West,Riverdale	43.679557	-79.352188
2	M4L	East Toronto	The Beaches West,India Bazaar	43.668999	-79.315572
3	M4M	East Toronto	Studio District	43.659526	-79.340923
4	M4N	Central Toronto	Lawrence Park	43.728020	-79.388790

- The New York Neighborhood Coordinates was downloaded from the Json file (https://cocl.us/new_york_dataset), then capturing the Borough, Neighborhoods and the location (latitude and longitude), and finally takes the Manhattan's information. And then finally takes only the Center of Manhattan.

```
In [71]: manhattan_data.head()
```

Out[71]:

	Borough	Neighborhood	Latitude	Longitude
0	Manhattan	Marble Hill	40.876551	-73.910660
1	Manhattan	Chinatown	40.715618	-73.994279
2	Manhattan	Washington Heights	40.851903	-73.936900
3	Manhattan	Inwood	40.867684	-73.921210
4	Manhattan	Hamilton Heights	40.823604	-73.949688

- The Paris Neighborhood Coordinates is loaded from the CSV file (<https://opendata.paris.fr/explore/dataset/arrondissements/download>) after the information is cleaned and structured in a data frame. And then finally takes only the Center of Paris.

```
In [69]: paris_data.head()
```

Out[69]:

	Neighborhood	Latitude	Longitude
0	Bourse	48.8682792225	2.34280254689
1	Temple	48.86287238	2.3600009859
2	Reuilly	48.8349743815	2.42132490078
3	Louvre	48.8625627018	2.33644336205
4	Hôtel-de-Ville	48.8543414263	2.35762962032

- The venues of each neighborhood are obtained by Foursquare with a limit of 100 venues and a radius of 500 of the center of each neighborhood. For those venues is saved the latitude, longitude and venue category

```
In [72]: venues_merged.head()
```

Out[72]:

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	The Beaches	43.6764	-79.293	The Big Carrot Natural Food Market	43.678879	-79.297734	Health Food Store
1	The Beaches	43.6764	-79.293	Grover Pub and Grub	43.679181	-79.297215	Pub
2	The Beaches	43.6764	-79.293	Guru Raghavendra Ji	43.680187	-79.292337	Astrologer
3	The Beaches	43.6764	-79.293	Upper Beaches	43.680563	-79.292869	Neighborhood
4	The Danforth West/Riverdale	43.6796	-79.3522	Pantheon	43.677621	-79.351434	Greek Restaurant

III. Methodology

Firstly, we need to get the list of the location of each neighborhood within Toronto, Manhattan and Paris, the coordinates of the center of each city and a list of venues close to each neighborhood with the location and the Venue Category. In the previous section was explained the process of how this data was obtained.

After getting ready the data, it is important to understand the data and check for any inconsistency or improve opportunities. In this case, at evaluate the venues category, we figured it out that Coffee Shops can be obtained as “Coffee Shops” and “Café”, and Gym can be obtained also as “Gym / Fitness Center”.

Once the data is ready to model, we first create a DataFrame in which defines the recurrence of each venue category on each neighborhood with help of `get_dummies` for pandas libraries, and then calculate per neighborhood how many venues there are.

Then is used the Kmeans function to define which are the main segments of neighborhood due to the recurrence of venues. With trial and error, I defined the number of clusters of 6. After this, the segments are mapped for each city to take a view of how it is distributed and then summarized the information of each cluster for the 5 venues category with the most recurrence.

IV. Results

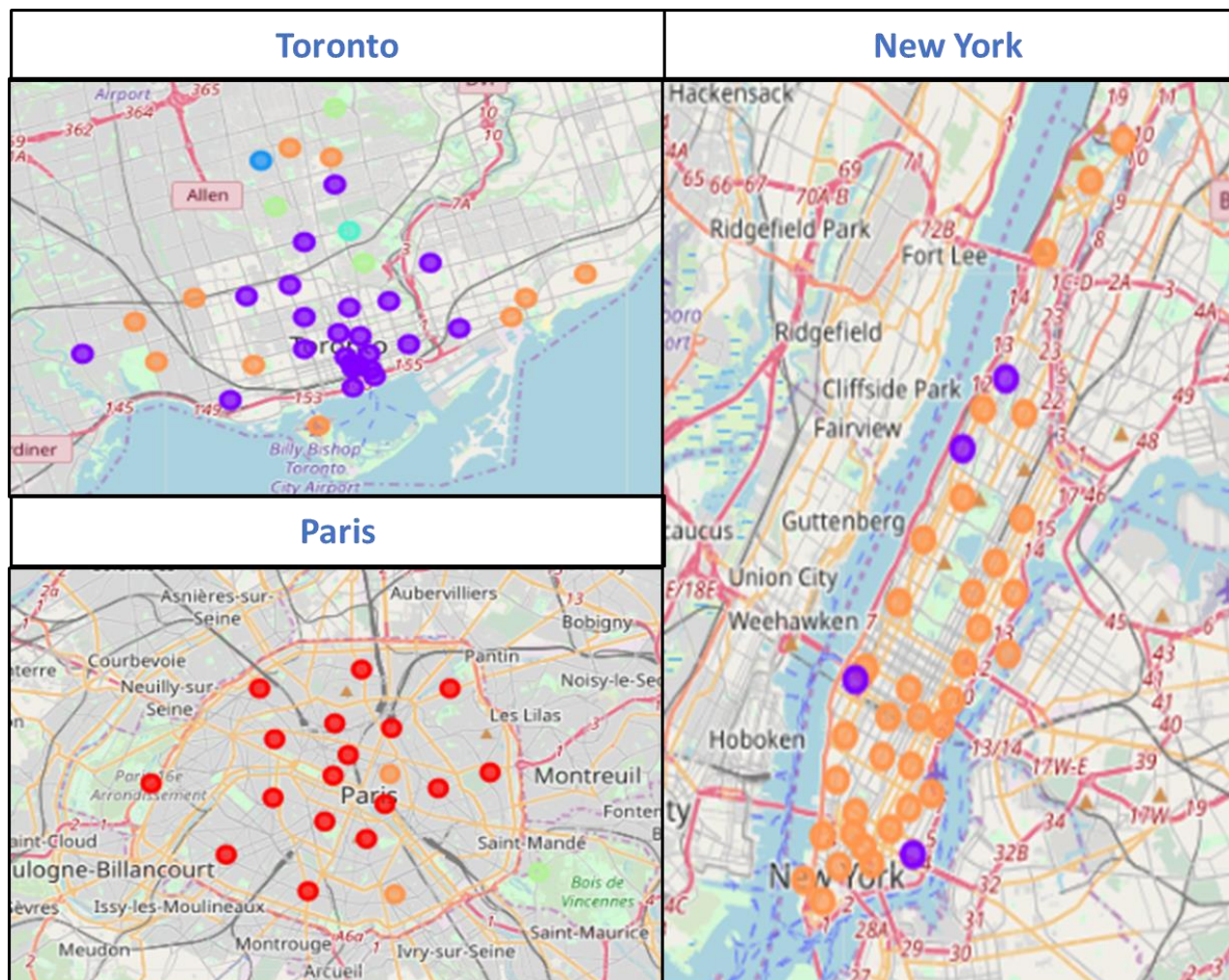
The results from the k-means clustering shows that we can categorized neighborhoods into 6 cluster based on the frequency of occurrence for each venue category:

	N. of neighborhoods	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
Cluster						
0	17	French Restaurant	Hotel	Coffee Shop	Italian Restaurant	Plaza
1	31	Coffee Shop	Italian Restaurant	Park	Pizza Place	Restaurant
2	1	Garden	Home Service	Zoo	Farmers Market	Duty-free Shop
3	4	Park	Trail	Bus Line	Swim School	Jewelry Store
4	44	Coffee Shop	Gym	Italian Restaurant	Park	Bakery
5	1	Summer Camp	Playground	Zoo	Farmers Market	Duty-free Shop

As coffee shop and Italian Restaurant are the most common place the conclusion would not consider much of these categories. Also, we can see that the cluster 0, 1 and 4 are the main cluster, the 2, 3 and 5 are cluster with neighborhoods that are almost unique.

The first Cluster are neighborhoods, above Coffee Shops and Italian Restaurant, with French Restaurants and Hotels, the second have Parks, Pizza places, and the fourth with Gyms, parks and Bakeries.

The following maps shows the clustering for each city. The red points are from cluster 0, purple for cluster 1 orange for cluster 4



As graph shows, there is a strong relationship between the city and the cluster, Toronto for cluster 1, Manhattan for cluster 4 and Paris for cluster 0. Even thou, there are some neighborhoods that can be compared.

V. Discussion

As observations shows from the maps in the Result section, there are a strong pattern for each country, which drive the clustering algorithm. Also, the frequency of Coffee Shops, Italian Restaurant and French Restaurant are the most commons venues which is also a driven for the clustering.

There are three cluster with a low number of neighborhoods within. This can be explained because in those places none of the most frequent category are this neighborhood.

VI. Conclusion

If anyone from any country would like to move for another neighborhood with a similitud in the most common places I would recommend to look forward in the same country. However, in some cases there are similar neighborhoods in another city to move.

