# COURSERA CAPSTONE FINAL REPORT

April 4, 2020

Prepared by:

Mohammed Saidul Islam

# Contents

# 1 BUSINESS PROBLEM DESCRIPTION

The business problem that is addressed in this notebook is that, **if a person wants to open a new coffee shop** in a city in Canada, then what are the things that he/she has to look into before opening the shop. Here, by analyzing and exploring all of the Neighborhoods in the **Boroughs(North York, East York and York)** in the city **Vaughan**, he can get useful insights about the **venues** present in the neighborhoods. If he/she can find a neighborhood where no coffee shop is present currently he/she could try to establish one in that neighborhood. Also, he/she has to explore the neighboring neighborhoods to get more better insights for his/her business.

In this case, the stakeholders are himself/herself and the people in the neighborhoods. As he/she will be the **owner** of the coffee shop, and he/she wants to make profit off of it, he/she needs to analyze all the neighborhoods near the city. So, he/she will be the **internal stakeholder**.

And the customer will be the consumers. The popularity and prosperity of his/her business will very much depend of the customers' mood, whether they like the coffee shop or not, whether they like the services given by the employees or not. So, the **customers** will be the **external stakeholder** of the business.

# 2 OVERVIEW OF THE DATASET

The **dataset** that I am working on is the **Neighborhood data of Canada** according to their **postal codes**. It has been downloaded from the wikipedia page: Canada Postal codes. To scrape the webpage, I have used the **"beautifulsoup4"** library. The dataset consists of **three columns**, namely, **PostalCode** ==>refers to the postal code of each of the Neighborhood, **Borough** ==>the Borough in which the Neighborhood is situated, and **Neighborhood** ==>the name of the Neighborhood. To explore each of the Neighborhoods, where all of the **coffee shops, parks, restaurants** and **other venues**, the **Foursquare API** has been used. To use the Foursquare API I needed the **latitude** and the **longitude** values of each of the Neighborhoods. The latitude and the longitude values are collected from this website.

# 3  METHODOLOGY

As the business problem revolves around opening a coffee shop in a neighborhood in city of Vaughan in Canada, at first step the relevant **boroughs** are selected. The boroughs are: **North York, East York and York**.

In the second step, **all the neighborhoods** that resides in the boroughs selected have been figured out. After that, using the **foursquare API**, the **venues** that are residing in those neighborhoods are found out.

In the next step, **filtering** of the neighborhoods have been done based on the criteria on the absence of coffee shops. This results in the neighborhoods in those boroughs that does not have any coffee shops in them.

Finally, a **clustering technique (k-means clustering)** was used to find the clusters of similar neighborhoods. The clustering gives the necessary insight that is needed to find a place where if the coffee shop is established would result in **higher profit and customer satisfaction** for the owner.

# 4  ANALYSIS ON THE DATA

At first the selected neighborhoods are one hot encoded based on the data collected from the foursquare API.
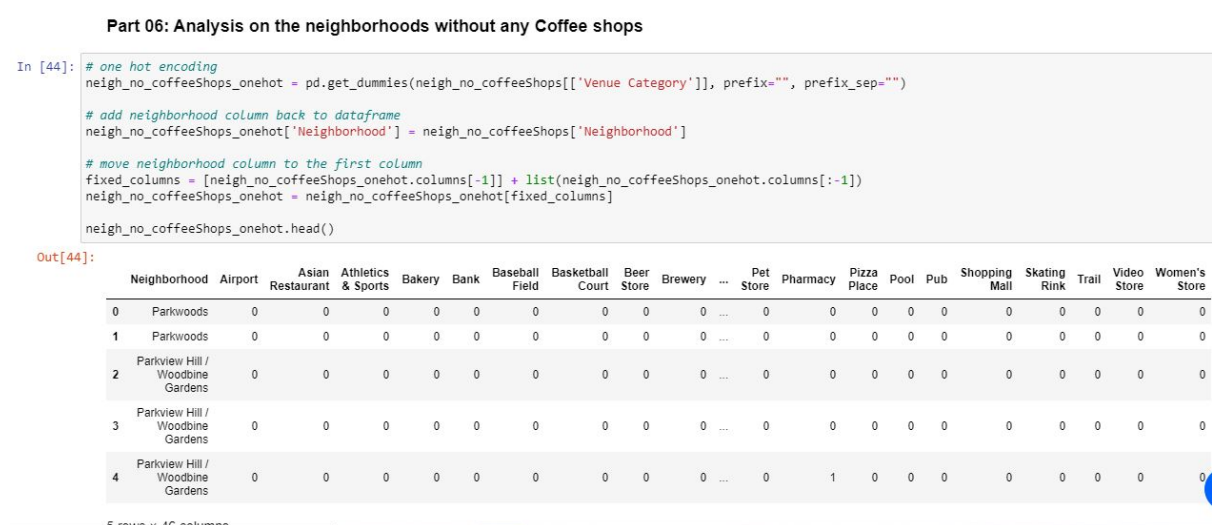


**Figure 1:** One hot encoded data

Then the neighborhoods are grouped by the mean of the one-hot values.



**Figure 2:** Group by data

In the next step, top five venues of each neighborhoods are generated.



**Figure 3:** Top 5 venues in neighborhoods

In the following step, the neighborhoods are merged and the venues are sorted in descending order based on their frequency values.
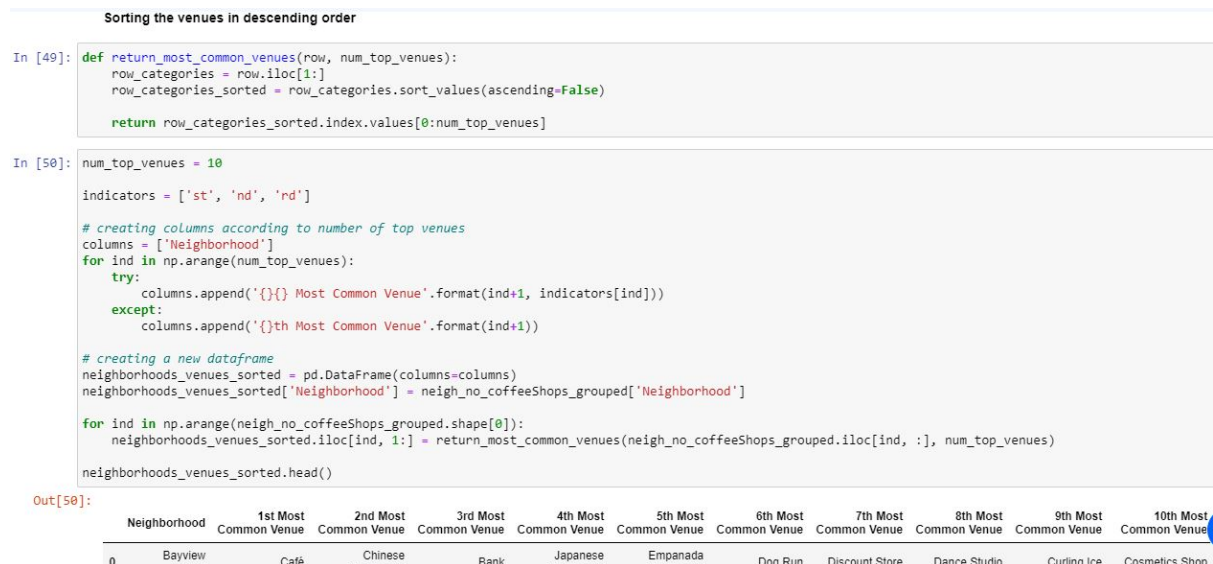


**Figure 4:** Sorted venues

Then, a clustering technique is used to cluster the neighborhoods and a map is generated based on the clusters.



**Figure 5:** Clustering the neighborhoods

**Figure 6:** Result of the clustering



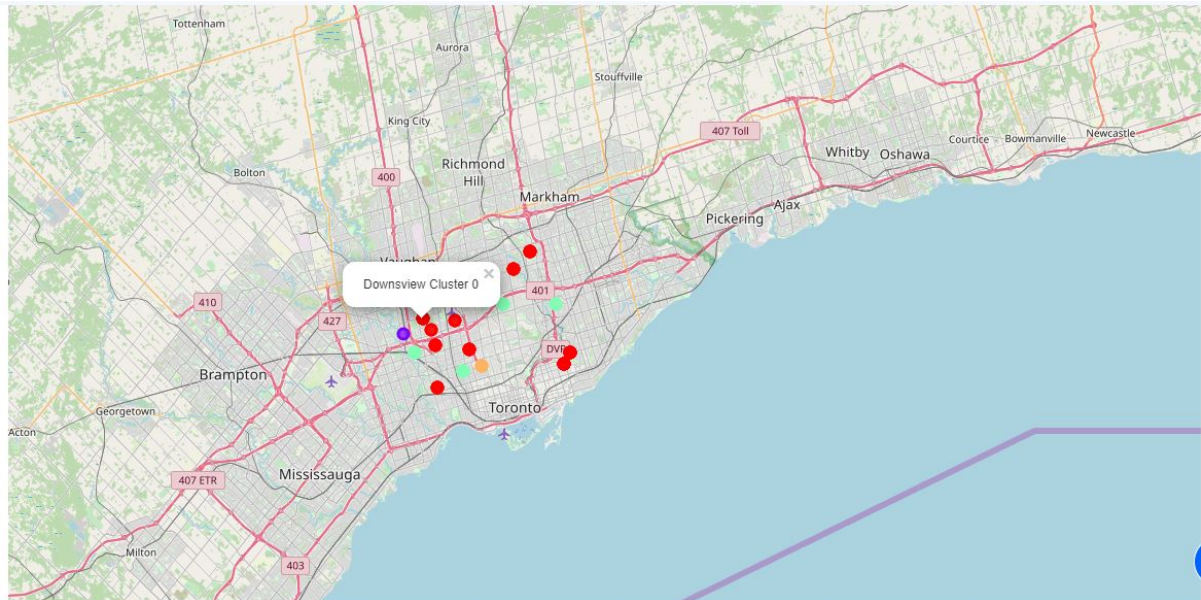**Figure 7:** Map showing different neighborhood clusters

Finally, the clusters are examined and based on the result, clusters of neighborhoods are chosen based on their business potential.

**Part 08: Examining the clusters**

*Cluster 1*

In [54]: `neigh_no_coffeeShops_merged.loc[neigh_no_coffeeShops_merged['Cluster Labels'] == 0, neigh_no_coffeeShops_merged.columns[[1] + list(range(5, neigh_no_cof`

Out[54]:

| | Neighborhood Latitude | Venue Longitude | Venue Category | Cluster Labels | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 43.706397 | -79.312913 | Gastropub | 0 | Pizza Place | Bus Line | Fast Food Restaurant | Athletics & Sports | Bank | Pharmacy | Pet Store | Gastropub | Gym / Fitness Center | Intersection |
| 3 | 43.706397 | -79.309279 | Gym / Fitness Center | 0 | Pizza Place | Bus Line | Fast Food Restaurant | Athletics & Sports | Bank | Pharmacy | Pet Store | Gastropub | Gym / Fitness Center | Intersection |
| 4 | 43.706397 | -79.312825 | Pharmacy | 0 | Pizza Place | Bus Line | Fast Food Restaurant | Athletics & Sports | Bank | Pharmacy | Pet Store | Gastropub | Gym / Fitness Center | Intersection |
| 5 | 43.706397 | -79.312270 | Bank | 0 | Pizza Place | Bus Line | Fast Food Restaurant | Athletics & Sports | Bank | Pharmacy | Pet Store | Gastropub | Gym / Fitness Center | Intersection |
| 6 | 43.706397 | -79.313130 | Pizza Place | 0 | Pizza Place | Bus Line | Fast Food Restaurant | Athletics & Sports | Bank | Pharmacy | Pet Store | Gastropub | Gym / Fitness Center | Intersection |

**Figure 8:** Cluster 1 data

*Cluster 2*

In [55]: `neigh_no_coffeeShops_merged.loc[neigh_no_coffeeShops_merged['Cluster Labels'] == 1, neigh_no_coffeeShops_merged.columns[[1] + list(range(5, neigh_no_coffe`

Out[55]:

| | Neighborhood Latitude | Venue Longitude | Venue Category | Cluster Labels | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 63 | 43.724766 | -79.532854 | Baseball Field | 1 | Baseball Field | Women's Store | Chinese Restaurant | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store |

**Figure 9:** Cluster 2 data

*Cluster 3*

In [56]: `neigh_no_coffeeShops_merged.loc[neigh_no_coffeeShops_merged['Cluster Labels'] == 2, neigh_no_coffeeShops_merged.columns[[1] + list(range(5, neigh_no_coffe`

Out[56]:

| | Neighborhood Latitude | Venue Longitude | Venue Category | Cluster Labels | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 62 | 43.756303 | -79.570637 | Empanada Restaurant | 2 | Empanada Restaurant | Women's Store | Field | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store | Construction & Landscaping |

**Figure 10:** Cluster 3 data

**Cluster 4**

```
In [57]: d.loc[neigh_no_coffeeShops_merged['Cluster Labels'] == 3, neigh_no_coffeeShops_merged.columns[[1] + list(range(5, neigh_no_coffeeShops_merged.shape[1]))]]
```

Out[57]:

| | Neighborhood Latitude | Venue Longitude | Venue Category | Cluster Labels | 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 | 43.753259 | -79.332140 | Park | 3 | Food & Drink Shop | Park | Café | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store |
| 1 | 43.753259 | -79.333114 | Food & Drink Shop | 3 | Food & Drink Shop | Park | Café | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store |
| 31 | 43.689026 | -79.456300 | Park | 3 | Park | Women's Store | Market | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store | Construction & Landscaping |
| 32 | 43.689026 | -79.456333 | Women's Store | 3 | Park | Women's Store | Market | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store | Construction & Landscaping |
| 33 | 43.689026 | -79.456317 | Market | 3 | Park | Women's Store | Market | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop | Convenience Store | Construction & Landscaping |

**Figure 11:** Cluster 4 data

**Cluster 5**

```
In [58]: neigh_no_coffeeShops_merged.loc[neigh_no_coffeeShops_merged['Cluster Labels'] == 4, neigh_no_coffeeShops_merged.columns[[1] + list(range(5, neigh_no_coffe
```

Out[58]:

| | Neighborhood Latitude | Venue Longitude | Venue Category | Cluster Labels | 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 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 28 | 43.693781 | -79.428705 | Field | 4 | Field | Trail | Hockey Arena | Chinese Restaurant | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop |
| 29 | 43.693781 | -79.426106 | Trail | 4 | Field | Trail | Hockey Arena | Chinese Restaurant | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop |
| 30 | 43.693781 | -79.431761 | Hockey Arena | 4 | Field | Trail | Hockey Arena | Chinese Restaurant | Empanada Restaurant | Dog Run | Discount Store | Dance Studio | Curling Ice | Cosmetics Shop |

**Figure 12:** Cluster 5 data

From our cluster analysis, we can see that the neighborhoods that falls in **cluster 0** and **cluster 3** has more venues in them than the other clusters. So, those neighborhoods might have more **potential customers** for any business.

```
In [62]: #Finding out the neighborhoods of interest
         neighborhoods_of_interest_1 = neigh_no_coffeeShops_merged[neigh_no_coffeeShops_merged['Cluster Labels'] == 0].Neighborhood
         neighborhoods_of_interest_2 = neigh_no_coffeeShops_merged[neigh_no_coffeeShops_merged['Cluster Labels'] == 3].Neighborhood
```

```
In [65]: #Neighborhoods of interest: 01
         print(neighborhoods_of_interest_1.unique())

         ['Parkview Hill / Woodbine Gardens' 'Glencairn' 'Woodbine Heights'
          'Hillcrest Village' 'Bayview Village' 'Downsview'
          'North Park / Maple Leaf Park / Upwood Park'
          'Runnymede / The Junction North']
```

```
In [66]: #Neighborhoods of interest: 02
         print(neighborhoods_of_interest_2.unique())

         ['Parkwoods' 'Caledonia-Fairbanks' 'Weston' 'York Mills West']
```

**Figure 13:** Potential clusters of neighborhoods

# 5 RESULTS AND DISCUSSION

So the cluster analysis results in 5 clusters of neighborhoods present in the boroughs of: North York, East York and York. To select the neighborhoods that would be perfect for opening a coffee shop two neighborhoods clusters have been selected, namely **cluster 0** and **cluster 3**.
In cluster 0, the neighborhoods present are: 'Parkview Hill / Woodbine Gardens', 'Glencairn', 'Woodbine Heights', 'Hillcrest Village', 'Bayview Village', 'Downsview', 'North Park / Maple Leaf Park / Upwood Park', 'Runnymede / The Junction North'.
In cluster 3, the neighborhoods present are: 'Parkwoods', 'Caledonia-Fairbanks', 'Weston', 'York Mills West'.
Although they fall in the same cluster, the distance between neighborhoods in cluster 3 is much greater than the neighborhoods in cluster 0.
So neighborhoods in cluster 0 would be a good choice for a potential neighborhood to open a coffee shop based on business perspective. Remember, the data that have been worked on, consists only of the neighborhoods that does not have any coffee shops in them. From the map analysis of the clusters it is found that the **Downsview** neighborhood might be the best choice in cluster 0.

# 6 CONCLUSION

Although the dataset consists of neighborhood data of every city in Canada and the foursquare API has been used to find out all the venues residing in those neighborhoods, but lack of population data, population density data in the neighborhoods certainly limit the capability to get a proper analysis of the business potential of each neighborhood. But, based on the current data, it can be said that, **Downsview** is a good choice to open a coffee shop in the city of Vaughan.