# Applied Data Science Capstone
## THE BATTLE OF NEIGHBORHOODS – REPORT

"Where is the best place to setup a new restaurant in Toronto"
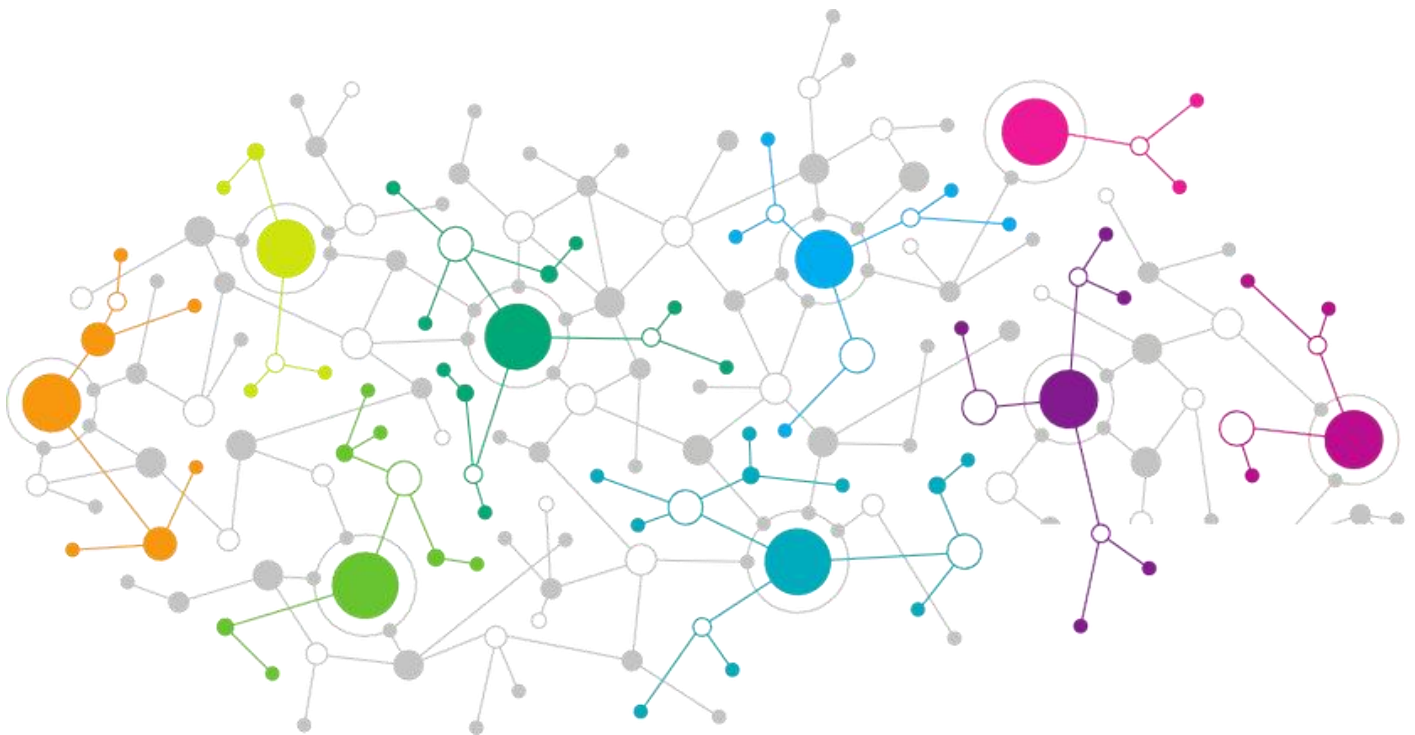
# Table of contents

# Table of Tables

# Table of Figures

# Introduction

This study will try to provide an optimum location for opening a new business in a specific city, based on:

- A neighborhoods property, e.g. the second most common language spoken (after English) in the neighborhood
- The number of competitors in the neighborhood
- The population density in the neighborhood
- The average income of the neighborhood

Let's investigate in the city of Toronto (Canada), and propose the best possible place for opening a new restaurant with ethnic cuisine.

Assume that preferably we would like the new restaurant to be in a neighborhood with a high degree of the same ethnic characteristics, i.e. assume the languages spoken in that neighborhood, so to make advantage of the cultural element of the area to try to increase the number of potential clients. To sustain the new business, there should be a lot of population, the less number of competitors possible. The restaurant should be of middle class and above to generate higher incomes.

The results could be highly usable for people having ethnic cooking skills or restaurant-businessmen, who want to open an ethnic restaurant in a neighborhood having some degree of the same ethnic culture in Toronto.

# Input data

## DATA DESCRIPTION

The datasets used in this study describes Toronto's neighborhoods and their main characteristics like localization (latitude and longitude), population, average income, ethnicity, etc.

The first dataset comes from Wikipedia website and is about Toronto's demographic information:

- https://en.wikipedia.org/wiki/Demographics_of_Toronto_neighbourhoods

The geographic coordinates of the neighborhoods of Toronto are taken from:

- File 'List_neihgborhood_Toronto.csv' which the latitude and longitude per neighborhood in Toronto area (see Appendix 1 for the details)

Combining the above data sets, we get demographic information, focused in Toronto's neighborhoods and the exact coordinates per neighborhood. Then by using the Foursquare API, we can retrieve further information for venues, venue categories and venue coordinates for every area. The Foursquare data set combined with the neighborhood's data set with demographic information will be the main data set that we will be used for the analysis. Visualization of the results via maps and graphs, where possible, will help to explain the data.

Based on the language spoken (second language spoken after 'English'), the neighborhood's population, the level of wealth, and the number of ethnic restaurants (restaurants with ethnicity common with the language spoken) the best possible set of candidate neighborhoods can be retrieved. Then by using k-means algorithm the candidate neighborhoods will be further analyzed. The results, via tables and maps will conclude on finding the best neighborhood to start an ethnic restaurant in an ethnic-cultural neighborhood, show any existing patterns and similarities between ethnic restaurants and ethnic populated neighborhoods in Toronto area.

## LIBRARIES USED

This paragraph gives the list of the libraries used and their version:

```
NUMPY: 1.15.4
PANDAS: 0.23.4
REQUESTS: 2.20.1
JSON: 2.0.9
BEAUTIFULSOUP: 4.6.3
MATPLOTLIB: 2.2.2
SKLEARN: 0.20.1
```

All the calculations have been done using Jupyter and Python 3:

https://labs.cognitiveclass.ai/tools/jupyterlab/

# Methodology

## DATA SCRAPPING FROM WIKIPEDIA

For first step, the information in the Wikipedia link must be transformed in a suitable form that enables further dataframe analysis. The link that provides the demographic data is the following:

[https://en.wikipedia.org/wiki/Demographics_of_Toronto_neighbourhoods](https://en.wikipedia.org/wiki/Demographics_of_Toronto_neighbourhoods)

By using 'BeautifulSoup' we retrieve the json data and fetch the wanted tags. We then clear the data via regular expressions for unwanted characters (e.g. remove '\n', empty spaces, etc.), remove non-meaningful data, rename index and columns and the dataframe with the demographic data is as shown below:

| | Neighborhood | Population | Density | Average income | Second language after English, % | Second language after English, name | Second language population |
|---|---|---|---|---|---|---|---|
| 0 | Agincourt | 44577 | 3580 | 25750 | 19.3 | Cantonese | 8603 |
| 1 | Alderwood | 11656 | 2360 | 35239 | 6.2 | Polish | 722 |
| 2 | Alexandra Park | 4355 | 13609 | 19687 | 17.9 | Cantonese | 779 |
| 3 | Allenby | 2513 | 4333 | 245592 | 1.4 | Russian | 35 |
| 4 | Amesbury | 17318 | 4934 | 27546 | 6.1 | Spanish | 1056 |
| 5 | Armour Heights | 4384 | 1914 | 116651 | 9.4 | Russian | 412 |
| 6 | Banbury | 6641 | 2442 | 92319 | 5.1 | Chinese | 338 |
| 7 | Bathurst Manor | 14945 | 3187 | 34169 | 9.5 | Russian | 1419 |
| 8 | Bay Street Corridor | 4787 | 43518 | 40598 | 9.6 | Mandarin | 459 |
| 9 | Bayview Village | 12280 | 2966 | 46752 | 8.4 | Cantonese | 1031 |

*Table 1: Toronto neighborhood demographic information*

Each neighborhood is depicted via its population, density, average income, the most common language after English spoken in the area (assume it as named as 'language' from now on), the name of the language, the population speaking that language (assume it as 'ethnic population').

## COMPLETE DEMOGRAPHIC DATAFRAME WITH NEIGHBORHOOD COORDINATES

It is possible to add latitude and longitude data, by merging the demographic dataframe and the neighborhoods of Old Toronto only. File 'oldToronto.csv' contains the coordinates of Old Toronto. The updated dataframe is as shown below:

| | Neighborhood | Population | Density | Average income | Percentage | Language | Second language population | Latitude | Longitude |
|---|---|---|---|---|---|---|---|---|---|
| 0 | Alexandra Park | 4355 | 13609 | 19687 | 17.9 | Cantonese | 779 | 43.71627 | -79.40555 |
| 1 | Allenby | 2513 | 4333 | 245592 | 1.4 | Russian | 35 | 43.71275 | -79.54746 |
| 2 | Bay Street Corridor | 4787 | 43518 | 40598 | 9.6 | Mandarin | 459 | 43.65777 | -79.38619 |
| 3 | Bedford Park | 13749 | 6057 | 80827 | 0.7 | Greek | 96 | 43.73138 | -79.42116 |
| 4 | Bloor West Village | 5175 | 6993 | 55578 | 3.6 | Ukrainian | 186 | 43.65936 | -79.48543 |
| 5 | Bracondale Hill | 5343 | 8618 | 41605 | 4.8 | Greek | 256 | 43.67600 | -79.42803 |
| 6 | Brockton | 9039 | 8217 | 27260 | 19.9 | Portuguese | 1798 | 43.66055 | -79.40531 |
| 7 | Cabbagetown | 11120 | 7943 | 50398 | 1.6 | Chinese | 177 | 43.66763 | -79.36606 |
| 8 | Carleton Village | 6544 | 8843 | 23301 | 17.0 | Portuguese | 1112 | 43.67200 | -79.45700 |
| 9 | Casa Loma | 3597 | 5369 | 82203 | 1.8 | Korean | 64 | 43.67000 | -79.41000 |

*Table 2: Head of Toronto Dataframe*

We can visualize the map of Toronto, with the neighborhoods, the language used (after English) – which depicts the ethnic group, the percentage of the ethnic group. Different color is used for each language (e.g. 'pink' is used for 'Portuguese').



*Figure 1: Ethnies represenetd on Toronto map*

## FOURSQUARE API

**Note**: Foursquare API version is: 20180506

Now that we have the demographic information per neighborhood in Toronto, lets collect all venues within 1 km radius from the center of each neighborhood (limit to 100 venues) and store the results in a dataframe. We are only interested in 'Restaurants' so we filter the venue category by this type:

| | Neighborhood | Neighborhood Latitude | Neighborhood Longitude | Venue | Venue Latitude | Venue Longitude | Venue Category |
|---|---|---|---|---|---|---|---|
| 0 | Alexandra Park | 43.71627 | -79.40555 | Cibo Wine Bar | 43.711464 | -79.399570 | Italian Restaurant |
| 1 | Alexandra Park | 43.71627 | -79.40555 | La Vecchia Ristorante | 43.710167 | -79.399086 | Italian Restaurant |
| 2 | Alexandra Park | 43.71627 | -79.40555 | Grazie Ristorante | 43.709329 | -79.398823 | Italian Restaurant |
| 3 | Alexandra Park | 43.71627 | -79.40555 | Tio's Urban Mexican | 43.714630 | -79.400000 | Mexican Restaurant |
| 4 | Alexandra Park | 43.71627 | -79.40555 | Sushi Shop | 43.713609 | -79.399844 | Restaurant |
| 5 | Alexandra Park | 43.71627 | -79.40555 | Banh Mi Boys | 43.709217 | -79.398777 | Fast Food Restaurant |
| 6 | Alexandra Park | 43.71627 | -79.40555 | Sushi Rock Café | 43.709089 | -79.398641 | Sushi Restaurant |
| 7 | Alexandra Park | 43.71627 | -79.40555 | Mai Thai Restaurant | 43.708779 | -79.398720 | Thai Restaurant |
| 8 | Allenby | 43.71275 | -79.54746 | Faley Restaurant | 43.713817 | -79.558676 | Asian Restaurant |
| 9 | Allenby | 43.71275 | -79.54746 | Mcdonald's in Walmart | 43.714250 | -79.553289 | Fast Food Restaurant |

*Table 3: List of restaurants per neighborhood*

Then we calculate the sum of ethnic restaurants per neighborhood (note ethnic assumed the ethnic group speaking the second most common language after English in the area:

| | Neighborhood | Population | Density | Average income | Percentage | Language | Second language population | Latitude | Longitude | Total Restaurants |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Alexandra Park | 4355 | 13609 | 19687 | 17.9 | Cantonese | 779 | 43.71627 | -79.40555 | 8 |
| 1 | Allenby | 2513 | 4333 | 245592 | 1.4 | Russian | 35 | 43.71275 | -79.54746 | 3 |
| 2 | Bay Street Corridor | 4787 | 43518 | 40598 | 9.6 | Mandarin | 459 | 43.65777 | -79.38619 | 24 |
| 3 | Bedford Park | 13749 | 6057 | 80827 | 0.7 | Greek | 96 | 43.73138 | -79.42116 | 16 |
| 4 | Bloor West Village | 5175 | 6993 | 55578 | 3.6 | Ukrainian | 186 | 43.65936 | -79.48543 | 4 |
| 5 | Bracondale Hill | 5343 | 8618 | 41605 | 4.8 | Greek | 256 | 43.67600 | -79.42803 | 30 |
| 6 | Brockton | 9039 | 8217 | 27260 | 19.9 | Portuguese | 1798 | 43.66055 | -79.40531 | 32 |
| 7 | Cabbagetown | 11120 | 7943 | 50398 | 1.6 | Chinese | 177 | 43.66763 | -79.36606 | 13 |
| 8 | Carleton Village | 6544 | 8843 | 23301 | 17.0 | Portuguese | 1112 | 43.67200 | -79.45700 | 19 |
| 9 | Casa Loma | 3597 | 5369 | 82203 | 1.8 | Korean | 64 | 43.67000 | -79.41000 | 38 |

*Table 4: Number of restaurants per neighborhood*

Since a language can be spoken by more than one country (and represent more than one cuisines), the following speaking groups are formed: For Portuguese assume common ethnic group for Brazilian and Portuguese Restaurants. For Japanese assume common ethnic group for Sushi and Japanese Restaurants. For Cantonese assumed common ethnic group for Thai, Taiwanese, Vietnamese, Cantonese, Indonesian Restaurants. For Mandarin as Chinese for Chinese Restaurant. This is important as it enables to differentiate the ethnic restaurants per neighborhood. These restaurants will be the competitors if we want to open a new ethnic restaurant. We can then visualize in a map the number of competitors for each neighborhood:

*Figure 2: Map of languages*

In this map, each color corresponds to a language.

## WEIGHT FACTORS

Based on the initial requirements, to have the less competition possible, it is needed to introduce a new factor to depict the "Density of ethnic restaurants (same as language) out of total restaurants" (less is best). According to this factor the neighborhoods that should be avoided due to high number of competitors (ethnic restaurants same as the language spoken), is as follows:
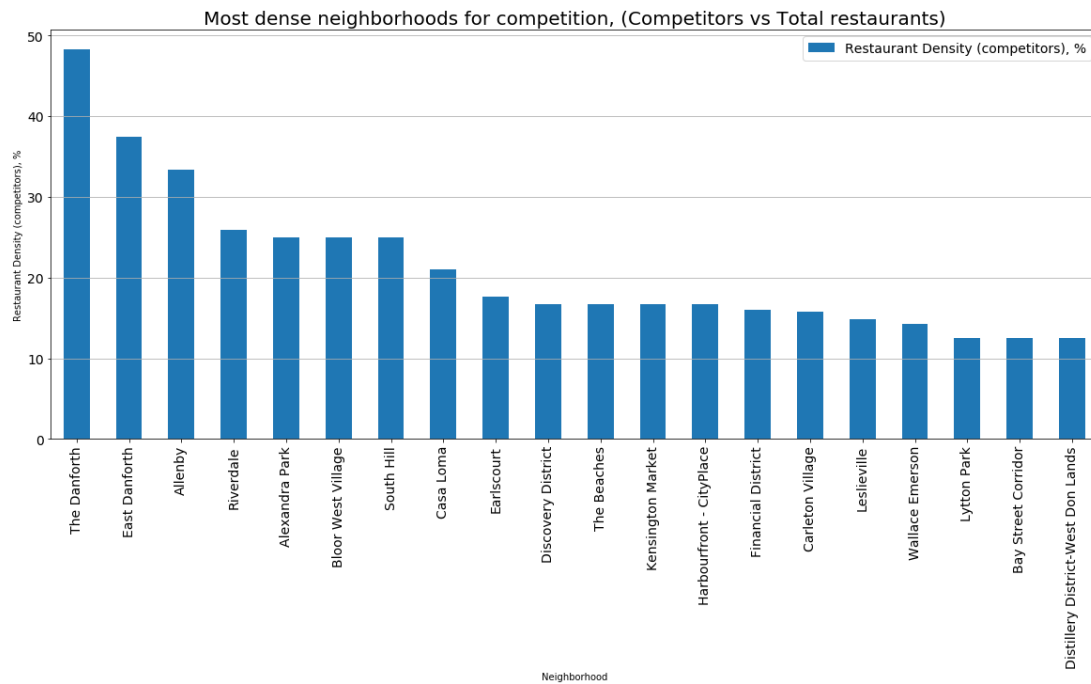


*Figure 3: Most dense neighborhood for competition*

We can see that "The Danforth", "East Danforth" and "Allenby" are highly competitive for this kind of business and better to be avoided.
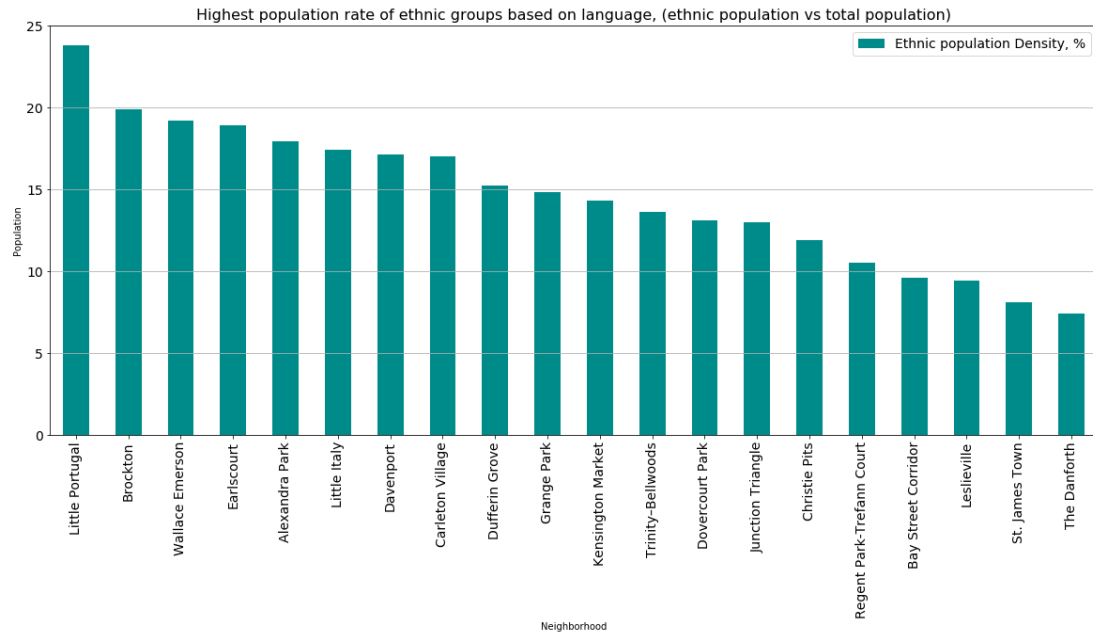
*Figure 4: Highest population rate of ethnic groups based on language*

Similarly, for the requirement to have strong ethnic presence in a neighborhood, the density of ethnic population over the total population is introduced (large is best). Neighborhoods such "Little Portugal", 'Brockton', "Wallace Emerson" have a high degree of the dominant language-ethnic groups. They should be considered in relation to average income and population later if they are good candidates.

We combine the 2 density factors with the main dataframe to a final dataframe, as shown below:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Earlscourt | Portuguese | 17240 | 3258 | 17 | 2 | 43.678000 | -79.449000 | 17.65 | 18.90 |
| 1 | Leslieville | Cantonese | 23567 | 2215 | 27 | 3 | 43.661927 | -79.332039 | 14.81 | 9.40 |
| 2 | Riverdale | Cantonese | 31007 | 2077 | 27 | 6 | 43.667750 | -79.349610 | 25.93 | 6.70 |
| 3 | Wallace Emerson | Portuguese | 10338 | 1984 | 21 | 2 | 43.663000 | -79.441000 | 14.29 | 19.20 |
| 4 | Brockton | Portuguese | 9039 | 1798 | 32 | 0 | 43.660550 | -79.405310 | 3.12 | 19.90 |
| 5 | Davenport | Portuguese | 8781 | 1501 | 19 | 1 | 43.673000 | -79.428000 | 10.53 | 17.11 |
| 6 | Dufferin Grove | Portuguese | 9875 | 1501 | 29 | 1 | 43.657000 | -79.428000 | 6.90 | 15.21 |
| 7 | Little Italy | Portuguese | 7917 | 1377 | 30 | 1 | 43.655000 | -79.413000 | 6.67 | 17.41 |
| 8 | Grange Park | Chinese | 9007 | 1333 | 27 | 2 | 43.653000 | -79.393000 | 11.11 | 14.81 |
| 9 | Little Portugal | Portuguese | 5013 | 1193 | 24 | 1 | 43.650000 | -79.435556 | 8.33 | 23.82 |

*Table 5: Final dataframe*

# Results

## MANUAL SELECTION OF BEST LOCATIONS

We could manually try to search for the optimum location base on the following criteria:

1) Assume middle-class and above neighborhoods only, i.e. merge the dataframe with the 'average-income' available from demographic information. Find the average income and filter neighborhoods above the man value, i.e. middle and above class.

2) Large population, so to attract as many people as possible. Further filter the above dataframe for neighborhoods with population above the mean population value

3) Large ethnic community, so to have significant cultural characteristics. Consider for ethnic group significant high, i.e. above average number of all ethnic groups.

4) Less number of competitors, so to avoid competition as much as possible. Based on restaurant density we keep only the neighborhoods where the competition is below the average number of competitors.

We obtain the following dataframe:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 27 | Davisville | Persian | 23727 | 355 | 35 | 0 | 43.701000 | -79.389000 | 2.86 | 1.5 | 55735 |
| 35 | The Annex | Spanish | 15602 | 202 | 35 | 0 | 43.670000 | -79.404000 | 2.86 | 1.3 | 63636 |
| 38 | Deer Park | Russian | 15165 | 166 | 16 | 0 | 43.688056 | -79.394028 | 6.25 | 1.1 | 80704 |
| 28 | Swansea | Polish | 11133 | 333 | 13 | 0 | 43.643889 | -79.477778 | 7.69 | 3.0 | 58681 |
| 22 | Forest Hill | Russian | 24056 | 577 | 11 | 0 | 43.700000 | -79.416667 | 9.09 | 2.4 | 101631 |

*Table 6: Results for handmade calculations*

"The Annex" and "Davisville" are the first candidates, even if "Davisville" has larger population and similar restaurant density, i.e. "Davisville" is better. "Deer Park" and "Swansea" are respectively in second and third position. "Deer Park" has less ethnic population and a slightly larger restaurant density. **The optimal place seems to be "Davisville" for opening a Persian, medium-upper class restaurant.**

# ANALYSIS OF RESULTS WITH K-MEANS ALGORITHM

Let us add to the main dataframe "Toronto_restaurants_final" the information about the average income, since this is relative with the type of restaurant that will open, i.e. lower, middle, upper, high class and then try to apply the k-means algorithm and see the results. For the clustering algorithm, a cluster of 5 groups (k=5) will be sufficient for the analysis. The following clusters are formed:

## Cluster 0:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors) | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | The Annex | Spanish | 15602 | 202 | 35 | 0 | 43.670000 | -79.404000 | 2.86 | 1.30 | 63636 | 0 |
| 1 | Fashion District | Portuguese | 4642 | 51 | 32 | 0 | 43.645000 | -79.398000 | 3.12 | 1.12 | 63282 | 0 |
| 2 | Summerhill | Chinese | 5100 | 56 | 29 | 0 | 43.683000 | -79.390000 | 3.45 | 1.12 | 88937 | 0 |
| 3 | Chaplin Estates | French | 4906 | 58 | 33 | 1 | 43.700000 | -79.400000 | 6.06 | 1.20 | 81288 | 0 |
| 4 | Deer Park | Russian | 15165 | 166 | 16 | 0 | 43.688056 | -79.394028 | 6.25 | 1.10 | 80704 | 0 |
| 5 | Bedford Park | Greek | 13749 | 96 | 16 | 1 | 43.731380 | -79.421160 | 12.50 | 0.71 | 80827 | 0 |
| 6 | Financial District | Japanese | 548 | 9 | 25 | 3 | 43.647935 | -79.381752 | 16.00 | 1.82 | 63952 | 0 |
| 7 | The Beaches | Cantonese | 20416 | 142 | 12 | 1 | 43.667266 | -79.297128 | 16.67 | 0.70 | 67536 | 0 |
| 8 | Casa Loma | Korean | 3597 | 64 | 38 | 7 | 43.670000 | -79.410000 | 21.05 | 1.81 | 82203 | 0 |

## Cluster 1:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Seaton Village | Portuguese | 5259 | 262 | 35 | 0 | 43.668000 | -79.416000 | 2.86 | 5.00 | 41506 | 1 |
| 1 | Davisville | Persian | 23727 | 355 | 35 | 0 | 43.701000 | -79.389000 | 2.86 | 1.50 | 55735 | 1 |
| 2 | Harbord Village | Portuguese | 5906 | 242 | 33 | 0 | 43.661000 | -79.406000 | 3.03 | 4.11 | 45792 | 1 |
| 3 | Playter Estates | Chinese | 3968 | 71 | 32 | 0 | 43.678056 | -79.355556 | 3.12 | 1.81 | 44557 | 1 |
| 4 | Bracondale Hill | Greek | 5343 | 256 | 30 | 0 | 43.676000 | -79.428030 | 3.33 | 4.81 | 41605 | 1 |
| 5 | Niagara | Portuguese | 6524 | 260 | 26 | 0 | 43.643000 | -79.408000 | 3.85 | 4.00 | 44611 | 1 |
| 6 | Upper Beaches | Cantonese | 19830 | 138 | 25 | 0 | 43.646667 | -79.408333 | 4.00 | 0.70 | 44346 | 1 |
| 7 | Roncesvalles | Polish | 15996 | 703 | 23 | 0 | 43.646231 | -79.449048 | 4.35 | 4.40 | 46820 | 1 |
| 8 | High Park North | Polish | 22746 | 682 | 18 | 0 | 43.656000 | -79.475000 | 5.56 | 3.00 | 46437 | 1 |
| 9 | Wychwood | Portuguese | 4182 | 112 | 31 | 1 | 43.676200 | -79.424400 | 6.45 | 2.70 | 53613 | 1 |
| 10 | Swansea | Polish | 5175 | 333 | 13 | 0 | 43.643889 | -79.477778 | 7.69 | 3.00 | 58681 | 1 |
| 11 | Cabbagetown | Chinese | 11120 | 177 | 13 | 0 | 43.667630 | -79.366060 | 7.69 | 1.60 | 50398 | 1 |
| 12 | Corktown | Spanish | 4484 | 94 | 22 | 1 | 43.655518 | -79.359712 | 9.09 | 2.12 | 54681 | 1 |
| 13 | Bay Street Corridor | Mandarin | 4787 | 459 | 24 | 2 | 43.657770 | -79.386190 | 12.50 | 9.61 | 40598 | 1 |
| 14 | Discovery District | Chinese | 7262 | 472 | 24 | 3 | 43.658000 | -79.388000 | 16.67 | 6.51 | 41998 | 1 |
| 15 | Bloor West Village | Ukrainian | 5175 | 186 | 4 | 0 | 43.659360 | -79.485430 | 25.00 | 3.61 | 55578 | 1 |
| 16 | Riverdale | Cantonese | 31007 | 2077 | 27 | 6 | 43.667750 | -79.349610 | 25.93 | 6.70 | 40139 | 1 |
| 17 | The Danforth | Greek | 7849 | 580 | 29 | 13 | 43.678472 | -79.347222 | 48.28 | 7.40 | 44979 | 1 |

## Cluster 2:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Rosedale | Chinese | 7672 | 76 | 23 | 0 | 43.646231 | -79.449048 | 4.35 | 1.00 | 213941 | 2 |
| 1 | Allenby | Russian | 2513 | 35 | 3 | 0 | 43.712750 | -79.547460 | 33.33 | 1.43 | 245592 | 2 |

# Cluster 3:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Christie Pits | Portuguese | 5124 | 609 | 38 | 0 | 43.664722 | -79.420833 | 2.63 | 11.90 | 30556 | 3 |
| 1 | Brockton | Portuguese | 9039 | 1798 | 32 | 0 | 43.660550 | -79.405310 | 3.12 | 19.90 | 27260 | 3 |
| 2 | Regal Heights | Spanish | 2719 | 149 | 31 | 0 | 43.676200 | -79.424400 | 3.23 | 5.52 | 36652 | 3 |
| 3 | Church and Wellesley | Spanish | 13397 | 241 | 28 | 0 | 43.665694 | -79.380956 | 3.57 | 1.81 | 37653 | 3 |
| 4 | Parkdale | Polish | 28367 | 822 | 26 | 0 | 43.640454 | -79.436731 | 3.85 | 2.90 | 26314 | 3 |
| 5 | Trinity–Bellwoods | Portuguese | 8687 | 1181 | 25 | 0 | 43.646667 | -79.408333 | 4.00 | 13.61 | 31106 | 3 |
| 6 | Dovercourt Park | Portuguese | 8497 | 1113 | 22 | 0 | 43.665000 | -79.432000 | 4.55 | 13.11 | 28311 | 3 |
| 7 | The Junction | Portuguese | 11391 | 467 | 20 | 0 | 43.665556 | -79.464444 | 5.00 | 4.11 | 34906 | 3 |
| 8 | Little Italy | Portuguese | 7917 | 1377 | 30 | 1 | 43.655000 | -79.413000 | 6.67 | 17.41 | 31231 | 3 |
| 9 | Dufferin Grove | Portuguese | 9875 | 1501 | 29 | 1 | 43.657000 | -79.428000 | 6.90 | 15.21 | 27961 | 3 |
| 10 | Garden District | Chinese | 8240 | 247 | 25 | 1 | 43.658500 | -79.375800 | 8.00 | 3.01 | 37614 | 3 |
| 11 | St. James Town | Filipino | 14666 | 1187 | 25 | 1 | 43.669167 | -79.372778 | 8.00 | 8.10 | 22341 | 3 |
| 12 | Little Portugal | Portuguese | 5013 | 1193 | 24 | 1 | 43.650000 | -79.435556 | 8.33 | 23.82 | 29224 | 3 |
| 13 | Davenport | Portuguese | 8781 | 1501 | 19 | 1 | 43.673000 | -79.428000 | 10.53 | 17.11 | 28335 | 3 |
| 14 | Grange Park | Chinese | 9007 | 1333 | 27 | 2 | 43.653000 | -79.393000 | 11.11 | 14.81 | 35277 | 3 |
| 15 | Junction Triangle | Portuguese | 6666 | 866 | 25 | 2 | 43.659000 | -79.446000 | 12.00 | 13.01 | 28067 | 3 |
| 16 | Wallace Emerson | Portuguese | 10338 | 1984 | 21 | 2 | 43.663000 | -79.441000 | 14.29 | 19.20 | 25029 | 3 |
| 17 | Leslieville | Cantonese | 23567 | 2215 | 27 | 3 | 43.661927 | -79.332039 | 14.81 | 9.40 | 30886 | 3 |
| 18 | Carleton Village | Portuguese | 6544 | 1112 | 19 | 2 | 43.672000 | -79.457000 | 15.79 | 17.01 | 23301 | 3 |
| 19 | Kensington Market | Cantonese | 3740 | 534 | 30 | 4 | 43.654772 | -79.400678 | 16.67 | 14.30 | 23335 | 3 |
| 20 | Earlscourt | Portuguese | 17240 | 3258 | 17 | 2 | 43.678000 | -79.449000 | 17.65 | 18.90 | 26672 | 3 |
| 21 | Alexandra Park | Cantonese | 4355 | 779 | 8 | 1 | 43.716270 | -79.405550 | 25.00 | 17.91 | 19687 | 3 |
| 22 | East Danforth | Cantonese | 21440 | 900 | 8 | 2 | 43.688056 | -79.301944 | 37.50 | 4.20 | 33847 | 3 |
| 23 | Port Lands | Mandarin | 571 | 19 | 1 | 0 | 43.648056 | -79.338333 | 100.00 | 3.50 | 36243 | 3 |

# Cluster 4:

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Forest Hill | Russian | 24056 | 577 | 11 | 0 | 43.700000 | -79.416667 | 9.09 | 2.40 | 101631 | 4 |
| 1 | Yorkville | French | 6045 | 114 | 38 | 3 | 43.670278 | -79.391111 | 10.53 | 1.90 | 105239 | 4 |
| 2 | Lytton Park | Serbian | 6494 | 58 | 8 | 0 | 43.716000 | -79.406000 | 12.50 | 0.91 | 127356 | 4 |
| 3 | South Hill | French | 6218 | 62 | 16 | 3 | 43.681000 | -79.404000 | 25.00 | 1.01 | 120453 | 4 |

# Discussion

We can see from the results, that the main factor that the machine learning algorithm has used to divide the neighborhoods is the 'Average income' data. This property proved to be the more decisive from all other properties of the neighborhoods. In more details per cluster, we can see the following:

**Cluster 0**, has the upper-class population (60k – 90k). All neighborhoods either have small ethnic group, or small population relative to neighborhoods of other clusters. "The Annex" (Spanish) seems to be the best option for this group.

**Cluster 1**, has middle class areas (40k – 50k). We could say that being at the average class, both low-level and high-level income citizens can be attracted, i.e. this is the most representative group of neighborhoods. Let us further filter for population more than the average of the cluster (Cluster1_final):

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 6 | Upper Beaches | Cantonese | 19830 | 138 | 25 | 0 | 43.646667 | -79.408333 | 4.00 | 0.7 | 44346 | 1 |
| 1 | Davisville | Persian | 23727 | 355 | 35 | 0 | 43.701000 | -79.389000 | 2.86 | 1.5 | 55735 | 1 |
| 11 | Cabbagetown | Chinese | 11120 | 177 | 13 | 0 | 43.667630 | -79.366060 | 7.69 | 1.6 | 50398 | 1 |
| 8 | High Park North | Polish | 22746 | 682 | 18 | 0 | 43.656000 | -79.475000 | 5.56 | 3.0 | 46437 | 1 |
| 10 | Swansea | Polish | 11133 | 333 | 13 | 0 | 43.643889 | -79.477778 | 7.69 | 3.0 | 58681 | 1 |
| 7 | Roncesvalles | Polish | 15996 | 703 | 23 | 0 | 43.646231 | -79.449048 | 4.35 | 4.4 | 46820 | 1 |
| 16 | Riverdale | Cantonese | 31007 | 2077 | 27 | 6 | 43.667750 | -79.349610 | 25.93 | 6.7 | 40139 | 1 |

*Table 7: Second filtering on Cluster 1 results*

The strongest ethnic groups are at "Davisville", "High Park North", "Roncesvalles" and "Riverdale" (max). From all the above "Davisville" (Persian) has the lowest competition (1.5%) and the second largest population after "Riverdale". So, for this cluster and overall clusters, "Davisville" (Persian) is the best option for opening a new Persian restaurant (middle-class).

**Cluster 2**, has the most expensive areas (>200k), although the population at these areas is small and the ethnic group small, i.e. The areas do not represent a high a cultural neighborhood. Not efficient enough.

**Cluster 3**, has the low-class areas (<35k). At these areas there is very high competition for almost the half of the neighborhoods. Best of all seems to be "Parkdale" (Polish) with very high population, very strong Polish group representative, no competition for other ethnic restaurants and relatively low competition from other types of restaurants. For low-level class "Parkdale" (Polish) is the clear winner.

**Cluster 4**, has high-class areas (>100k). Small ethnic groups relatively to population and not many restaurants in the area. "Forest Hill" (Russian) seems the exception and for this cluster is the best option. For high-class restaurant "Forest Hill" (Russian) is the best option.

# Conclusion

We have analyzed the neighborhoods of Toronto with respect to:

- Population
- Competition
- 14 Ethnic group presences (language oriented)
- Average income

The results from observation are the same as the ones from applying the k-means algorithm. The 'Average income' was the most distinctive property for the neighborhoods, more important than other significant properties such as the population. Below the best candidates, based on "Average income":

| | Neighborhood | Language | Population | Second language population | Total Restaurants | Number of Competitors | Latitude | Longitude | Restaurant Density (competitors), % | Ethnic population Density, % | Average income | Cluster Labels |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Parkdale | Polish | 28367 | 822 | 26 | 0 | 43.640454 | -79.436731 | 3.85 | 2.9 | 26314 | 3 |
| 1 | Davisville | Persian | 23727 | 355 | 35 | 0 | 43.701000 | -79.389000 | 2.86 | 1.5 | 55735 | 1 |
| 2 | The Annex | Spanish | 15602 | 202 | 35 | 0 | 43.670000 | -79.404000 | 2.86 | 1.3 | 63636 | 0 |
| 3 | Forest Hill | Russian | 24056 | 577 | 11 | 0 | 43.700000 | -79.416667 | 9.09 | 2.4 | 101631 | 4 |

*Table 8: Best places*

# Appendix 1

Content of "List_neighborhood_Toronto.csv" file:

| Neighborhood | Latitude | Longitude |
|---|---|---|
| Alexandra Park | 43.71627 | -79.40555 |
| Allenby | 43.71275 | -79.54746 |
| Bay Street Corridor | 43.65777 | -79.38619 |
| Bedford Park | 43.73138 | -79.42116 |
| Bloor West Village | 43.65936 | -79.48543 |
| Bracondale Hill | 43.676 | -79.42803 |
| Brockton | 43.66055 | -79.40531 |
| Cabbagetown | 43.66763 | -79.36606 |
| Carleton Village | 43.672 | -79.457 |
| Casa Loma | 43.67 | -79.41 |
| Chaplin Estates | 43.7 | -79.4 |
| Christie Pits | 43.664722 | -79.420833 |
| Church and Wellesley | 43.665694 | -79.380956 |
| Corktown | 43.655518 | -79.359712 |
| Davenport | 43.673 | -79.428 |
| Davisville | 43.701 | -79.389 |
| Deer Park | 43.688056 | -79.394028 |
| Discovery District | 43.658 | -79.388 |
| Distillery District/West Don Lands | 43.655 | -79.353 |
| Dovercourt Park | 43.665 | -79.432 |
| Dufferin Grove | 43.657 | -79.428 |
| Earlscourt | 43.678 | -79.449 |
| East Danforth | 43.688056 | -79.301944 |
| Fashion District | 43.645 | -79.398 |
| Financial District | 43.647935 | -79.381752 |
| Forest Hill | 43.7 | -79.416667 |
| Fort York/Liberty Village | 43.637 | -79.422 |
| Garden District | 43.6585 | -79.3758 |
| Grange Park | 43.653 | -79.393 |
| Harbord Village | 43.661 | -79.406 |
| Harbourfront / CityPlace | 43.638 | -79.385 |
| High Park North | 43.656 | -79.475 |
| Junction Triangle | 43.659 | -79.446 |
| Kensington Market | 43.654772 | -79.400678 |
| Lawrence Park | 43.722 | -79.388 |
| Leslieville | 43.661927 | -79.332039 |
| Little Italy | 43.655 | -79.413 |
| Little Portugal | 43.65 | -79.435556 |
| Lytton Park | 43.716 | -79.406 |
| Moore Park | 43.691 | -79.377 |
| Niagara | 43.643 | -79.408 |
| Parkdale | 43.640454 | -79.436731 |
| Playter Estates | 43.678056 | -79.355556 |
| Port Lands | 43.648056 | -79.338333 |
| Regal Heights | 43.6762 | -79.4244 |
| Regent Park/Trefann Court | 43.656548 | -79.36201 |
| Riverdale | 43.66775 | -79.34961 |
| Roncesvalles | 43.646231 | -79.449048 |
| Rosedale | 43.646231 | -79.449048 |
| Seaton Village | 43.668 | -79.416 |
| South Hill | 43.681 | -79.404 |
| St. James Town | 43.669167 | -79.372778 |
| Summerhill | 43.683 | -79.39 |
| Swansea | 43.643889 | -79.477778 |
| The Annex | 43.67 | -79.404 |
| The Beaches | 43.667266 | -79.297128 |
| The Danforth | 43.678472 | -79.347222 |
| The Junction | 43.665556 | -79.464444 |
| Toronto Islands | 43.620833 | -79.378611 |
| Trinity–Bellwoods | 43.646667 | -79.408333 |
| Upper Beaches | 43.646667 | -79.408333 |
| Wallace Emerson | 43.663 | -79.441 |
| Wychwood | 43.6762 | -79.4244 |
| Yorkville | 43.670278 | -79.391111 |