# Yahuza in Toronto- Choosing the best neighbourhood

## Introduction

My client, Yahuza Suya, a Nigerian-Style barbeque brand are planning to open a new spot in Toronto, to cater to the growing Nigerian population in Canada and enter new markets. Their product, being mostly chicken and beef barbeque, will appeal to a wider market since chicken and beef is eaten mostly around the world.

## Problem

While they are well funded, Yahuza Suya will like to open their spot where there is a higher probability of success. They want to ensure the neighbourhood will be a safe one, with less crime likely. Another important factor is that the pricing of the product must be within the income range as they do not want to be priced out of the market.

## Stakeholders

Yahuza

Other businesspersons wanting to do business in Toronto

People looking for a home and want to understand the neighbourhoods

Looking for how to fit into Toronto.

## DATA

1. We will need demographics data to get a snapshot of the neighbourhoods and their income ranges. This will help us sort and get the neighbourhoods with the highest incomes and those that follow.
2. We will need the crime data which will enable us sort and filter out the safest neighbourhoods.
3. List of Metro stations to understand the metro station distribution
4. Foursquare data showing us the neighbourhoods and the venues present there. This will help us understand the characteristics of the neighbourhoods and how our spot will fit in there.

### Sources

General Toronto data was obtained from

https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M

The demographic data was obtained from
https://en.wikipedia.org/wiki/Demographics_of_Toronto_neighbourhoods

The crime data was obtained from

http://data.torontopolice.on.ca/datasets?q=crime

The list of Metro stations was obtained from

https://en.wikipedia.org/wiki/List_of_Toronto_subway_stations
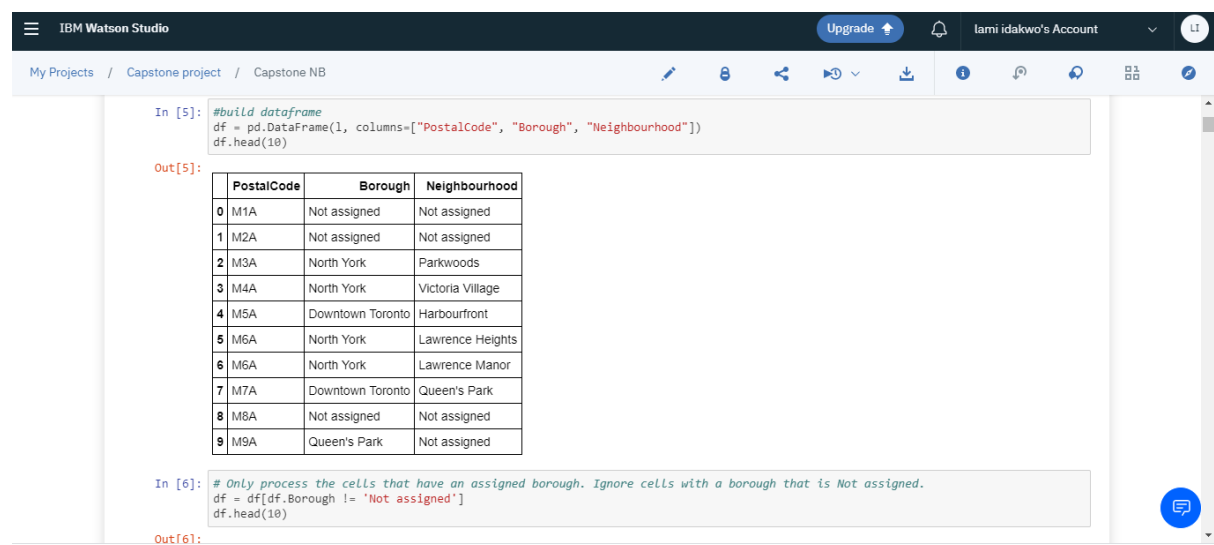
**Methodology**

We first import all necessary libraries. We ignore libraries like folium and sklearn until we need them so that it does not slow down the processing

**Neigbourhood Data**

At the basis of all the data is the list of boroughs and neighbourhoods with their postal codes. I used BeautifulSoup to scrape the data from https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M.

I then used it to build a dataframe which looked like this:



To clean it up and ensure I don't have empty cells, I processed only the cells with assigned boroughs. I used the 'dropna()' function.
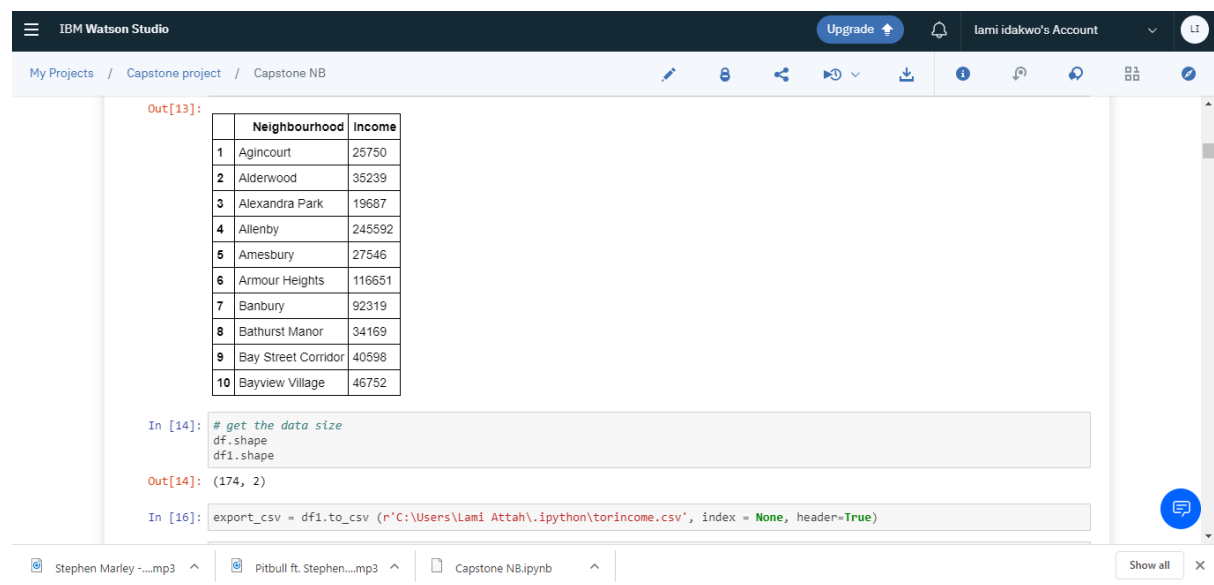
There were some boroughs with no neighbourhood. In such cases, I assigned the borough to the neighbourhoods. For neighbourhoods within the same borough and postal code, we grouped them together.

After doing this, we have a full dataframe with 103 rows and 3 columns.

**Income data**

To get the income data, I used BeautifulSoup again to scrape
https://en.wikipedia.org/wiki/Demographics_of_Toronto_neighbourhoods and obtain the
demographic data.

To use only the income, I dropped all other columns except the income column. This gave a
dataframe that looked like this:



To get the neighbourhoods with the higher incomes, I sorted the dataframe to show the a
descending order of the income. We can see the top ten (10) neighbourhoods with the highest
average incomes. The dataframe looks like this:

**Crime data**

The crime data was not importing from my computer so I decided to clean it up in Excel, then upload it to a website I own, and use read_csv to get the table.

I then sort the dataframe according to the assault average(assault is the most common form of crime in Toronto).

The dataframe looks like this :



We see the top ten safest neighbourhoods.

To complement the borough and neighbourhood dataframe and enrich our data, we obtain the coordinates and merge the two dataframes. We now have a dataframe that has the postal codes, neighbourhood, boroughs, and coordinates (latitude and longitude).

The dataframe looks like this:

Next, we obtained the list of Subway stations to understand their spread across Toronto. We know that the airport is a central point and that most of Toronto is not far from Pearson YYZ.

We clean the data in Excel and add the coordinate data. We have a dataframe that looks like:



Next, we import folium and use it to visualize the Toronto map.

It looks like this:



We also visualize the Subway stations in Toronto to see the spread. We can see that the subway stations are well spread but we take note of a central point where we have most stations. The visualization is below:

**FOURSQUARE API:**

We then use Foursquare API to obtain the venues in Toronto and observe. We decided to look for general venues and not the metro stations because we can see that the subway stations are well-spread. In addition, the metro stations form a part of the foursquare results and it will be important to use the other venues to get a full understanding of the neighbourhoods.

Results:

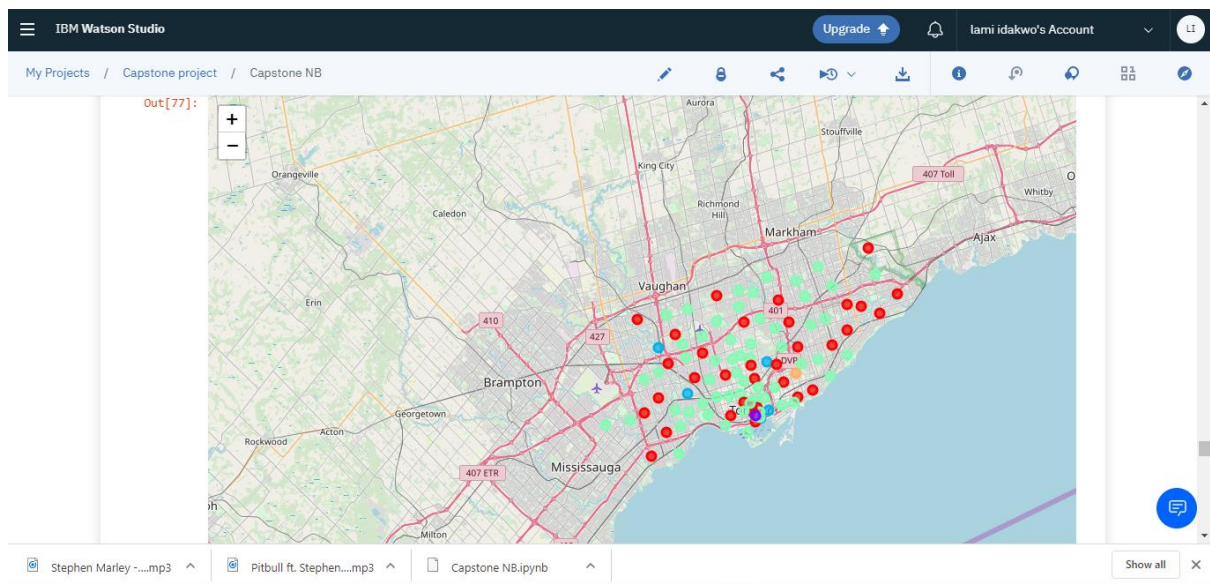In order to perform k-means clustering, we used one-hot encoding.

We then print each neighbourhood according to the most popular venues. We then obtain the 1st ten(10) most common venues for each neighbourhood.

**K-means**

To perform k-means, we take our k=5. We will be observing 5 clusters. We are using k-means so that we can have a statistical visual of the characteristics of the neighbourhood. Such a view will bring certain things to light which we may have omitted in looking at income and crime data.

After perfoming k-means, we visualize the 5 clusters shown below:



We then begin to observe each clusters:

**Cluster 0-** We see that cluster 0 has very many spots. It also has bus stations, trains stations and metro stations among its top ten most common place. It appears to be favourable for our business line. However, it may be saturated and there may be too much competition.

**Cluster 1** – Cluster one is just one neighbourhood which has foreign food restaurants and other **services. However, it does not have transportation services in its top ten common places.**

**Cluster 2-** While cluster two has a few varieties of venues, it seems to be quite a homely cluster, most likely family-oriented. It also does not have transportation services as common venues.

**Cluster 3-** Cluster 3 appears to have many spots, including transportation services, it also appears to cater to a large number of foreign nationals with the amount of foreign food restaurants around there.

**Cluster 4**- Cluster 4 is a bit similar to cluster 1 and appears to be focused on a certain type of market.

**Discussion:**

Based on the results obtained, I closed in on cluster 3, which clusters neighbourhoods with the kind of venues we will expect around our own. It also has transportation services as many of the top ten most common venues which is one of our main factors. From cluster 3, I look out for neighbourhoods which satisfy our factors- (safety and income average).

Cluster 3 is shown below:



On observation, I pick these six(6) neighbourhoods:

**Lawrence Park**

**Rosedale**

**Parkwoods**

**Willowdale West**

**Riverdale West**

**Mimico**

To narrow it down to three and have more specific results, I look closely and see that:

1.**Lawrence Park** satisfies all the three requirements and is close to a college. This will attract young adults who have an active nightlife.

2. **Willowdale West** is also a good idea as it satisfies the requirements.

3. The third option will be **Parkwoods**, which also satisfies the requirement of safety and income. Further search reveals there are many immigrants.

**Conclusion**

Choosing a good neighbourhood usually depends on more than one factor. In most cases, it is usually a healthy combination of important factors.

For this task, while safety and pricing were major factors, it would have been counterproductive to choose based on these two factors alone. The use of foursquare data and the venues provided gave us a more descriptive view of these neighbourhoods and we can make a more robust decision.

Such tasks show the important of Data Science as it enables fact-based decision-making, even if it means coming to the same assumption.