

# **Most Prospective Business in Toronto**

**Capstone Project - The Battle of  
Neighborhoods (Week 2)**

Jeffrey Chai Zhi Yang  
9 February 2021

# Table of Contents

- +Introduction
- +Data
- +Methodology
- +Results & Discussion
- +Conclusion

# Introduction (1/2)

## +Background

As an entrepreneur, John wishes to start his own business in Toronto, but he needs ideas on what is the most popular business to do in the city at the moment. He has been doing his own primary market research around the neighbourhood but has deemed it as a very time inefficient method of doing so.

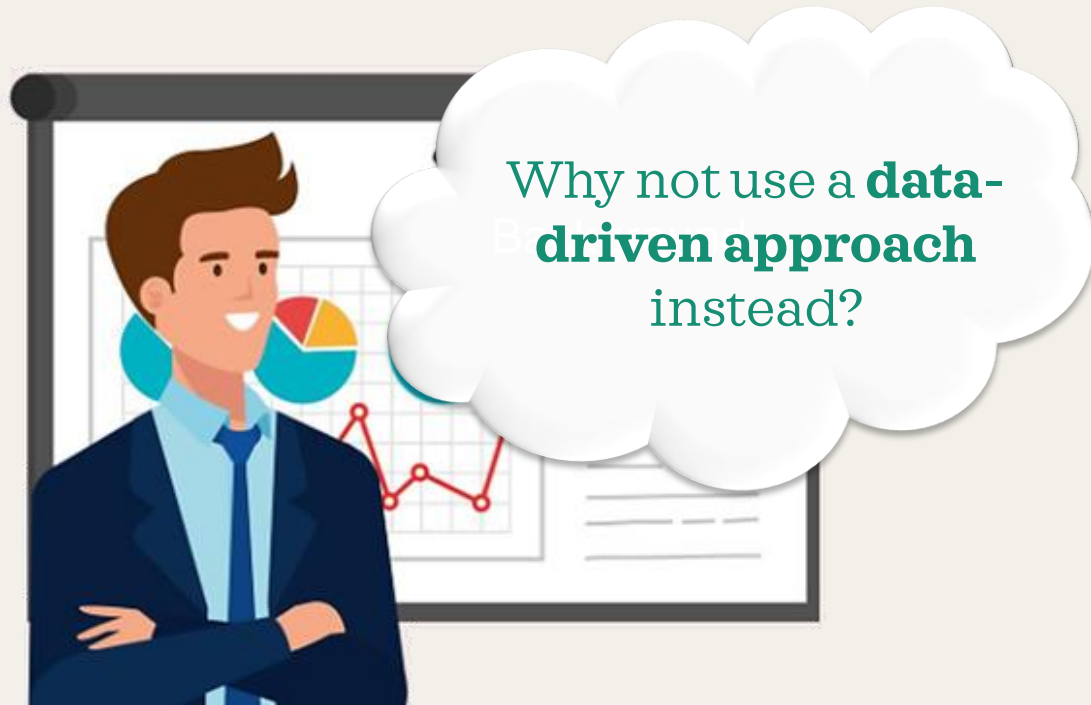
## +Business Problem

As an entrepreneur, John wishes to start his own business in Toronto, but he needs ideas on what is the most popular business to do in the city at the moment. He has been doing his own primary market research around the neighbourhood but has deemed it as a very time inefficient method of doing so.

# Introduction (2/2)

## + Motivation

The project aims to find out what is the most visited venues in the majority of the neighbourhoods (namely postal code M) in Toronto so that this can give John some form of confidence, driven by data analytics, that he could probably set up a similar business given the highest level of demand for that given establishment.



# Data

## + Data Sources & Pre-processing



### Source

# 1

List of Postal Codes (only M) of Canada

Web-Scraping

Convert to dataframe

### Source

# 2

Geographical coordinates of each postal code (only M) in Toronto

Read CSV

Convert to dataframe

### Source

# 3

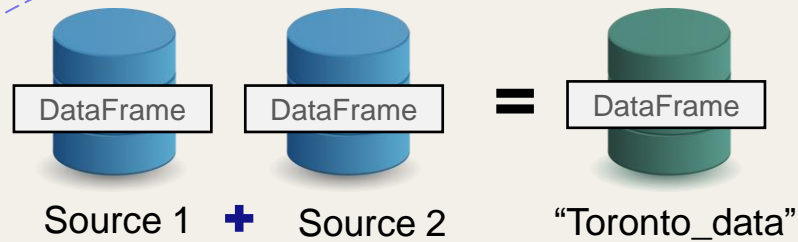
Information on each Neighbourhood in Toronto

API Call

Convert to dataframe

# Methodology

## + Overview



Inputs to function:  
['Neighbourhood'],  
['Latitude'],  
['Longitude'],



API Call to Foursquare

```
# Explore neighborhoods in Toronto
def getNearbyVenues(names, latitude, longitude, radius):
    venue_list = []
    for name, lat, lng in zip(names, latitude, longitude):
        print(name)

    # create the API request URL
    url = "https://api.foursquare.com/v2/venues/explore?llat={}&llong={}&radius={}&client_id={}&client_secret={}"
    CLIENT_ID, CLIENT_SECRET, VENUE_ID, PUBLIC, radius, VENUE_ID

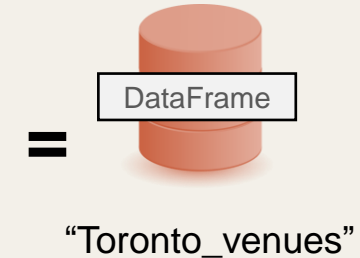
    # make the GET request
    results = requests.get(url).json()[{"response": "venues"}][{"items": []}]

    # return only relevant information for each nearby venue
    venue_list.append([
        name,
        lat,
        lng,
        venue["name"],
        venue["location"]["lat"],
        venue["location"]["lng"],
        venue["categories"][0]["name"]
    ])

    nearby_venues = pd.DataFrame([item for venue_list in venue_list for item in venue_list])
    nearby_venues.columns = ["Neighbourhood",
                             "Neighbourhood Latitude",
                             "Neighbourhood Longitude",
                             "Venue Name",
                             "Venue Latitude",
                             "Venue Longitude",
                             "Venue Category"]

    return(nearby_venues)
```

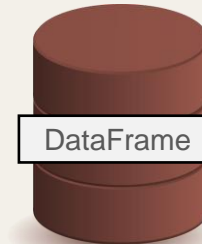
Defined Function for  
API Call:  
"getNearbyVenues"



+

A series of data  
manipulation steps...\*

- On this final aggregated table, perform count function on the '1<sup>st</sup> Most Visited Venue'<sup>#</sup> Column.
- The venue with the most number of counts of venue = **Most Prospective Business!**

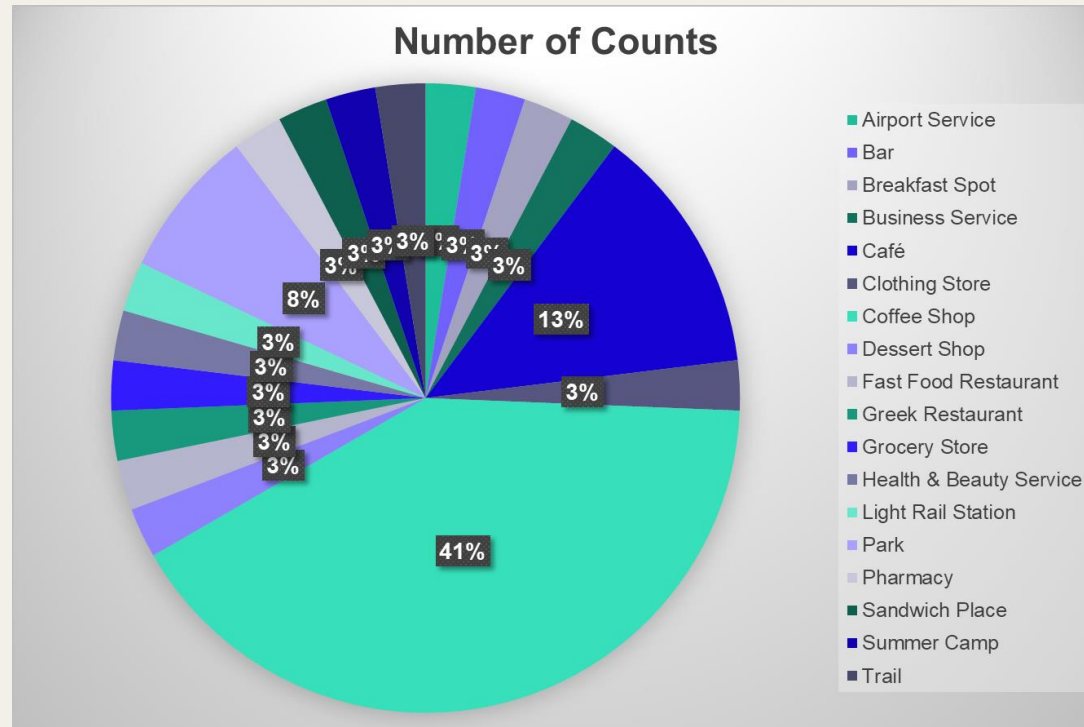


"neighborhoods\_venues\_  
sorted"

\*For more details on the data manipulation steps, please refer to the Jupyter notebook in the repository  
# One of the columns in 'neighborhoods\_venues\_sorted' dataframe.

# Results & Discussion

+ Most Number of Counts in '1<sup>st</sup> Most Visited Venue'



**COFFEE SHOP (41%)**



# Conclusion

+ To recommend John to try 'Coffee Shop' for his first entrepreneur business!

## + Future Work for additional Analysis

- Where should John better allocate his resources such that he can strategically set up his first set of coffee shop outlets?
- How can he discover and understand more about the market trends of individual neighbourhoods?
- Is he able to identify neighborhoods that generally loves coffee shops more? Or loves the similar common venues should John wish to venture into our sort of businesses in the future?



# Future Work (Cont.)

+ Possible Data Science Method to find similar/dissimilar market trends across neighborhoods in Toronto – **K-Means Clustering**<sup>^</sup>

```
In [16]: # clustering neighborhoods
from sklearn.cluster import KMeans
# set number of clusters
N_clusters = 5
toronto_grouped_clustering = toronto_grouped.drop('Neighborhood', 1)
# run k-means clustering
kmeans = KMeans(n_clusters=N_clusters, random_state=0).fit(toronto_grouped_clustering)
# check cluster labels generated for each row in the dataset
kmeans.labels_[0:10]

Out[16]: array([2, 2, 2, 2, 2, 2, 2, 2, 2, 2])

In [17]: # add clustering labels
neighborhood_names_sorted.insert(0, 'Cluster Labels', kmeans.labels_)
toronto_merged = toronto_data
toronto_merged.rename(columns={'Neighborhood': 'Neighborhood'}, inplace=True)
# merge neighborhood_grouped with neighborhood_data to add latitude/longitude for each neighborhood
toronto_merged = toronto_merged.join(neighborhood_names_sorted.set_index('Neighborhood'), on='Neighborhood')
toronto_merged.head() # check the last column!

Out[17]:
```

	Postal Code	Neighborhood	Neighborhood	Latitude	Longitude	Cluster Label	1st Most Common Value	2nd Most Common Value	3rd Most Common Value	4th Most Common Value	5th Most Common Value	6th Most Common Value	7th Most Common Value	8th Most Common Value
0	M5A	Downtown Toronto	Rogers Park, Yeshuofan	43.654000	-79.383000	2	Coffee Shop	Park	Park	Park	Cafe	Theater	Wine Shop	
1	M5A	Downtown Toronto	Queen's Park, Ontario Provincial Government	43.652001	-79.384004	2	Coffee Shop	College	College	College	College	College	College	College
2	M5B	Downtown Toronto	Queen's Park, Ontario Provincial Government	43.651702	-79.379007	2	Coffee Shop	Clothing Store	Cafe	Japanese Restaurant	Italian Restaurant	Chinese Restaurant	Indian Restaurant	Subs/Teli Shop
3	M5C	Downtown Toronto	St. James Town	43.651004	-79.375010	2	Coffee Shop	Cafe	Coffee Shop	Grocery	American Restaurant	Chinese Restaurant	Indian Restaurant	
4	M5C	East Toronto	The Beaches	43.670007	-79.260001	0	Trail	Health Food Store	Pub	Wine Shop	Cuban Restaurant	Diner Restaurant	Day Run	

```
In [18]: # Matplotlib and associated plotting modules
import matplotlib.pyplot as plt
import matplotlib.colors as mcolors

# create map
map_clusters = folium.Map(location=[latitude, longitude], zoom_start=11)

# set color scheme for the clusters
c = mcolors.to_rgba_array(kmeans.labels_)
c = c[:, :3] # (R, G, B) for i in range(kmeans.labels_.size)
cmap = mcolors.LinearSegmentedColormap.from_list('toronto_clusters', c)

# add markers to the map
markers = []
for lat, lon, pos, cluster in zip(toronto_merged['latitude'], toronto_merged['longitude'], toronto_merged['Neighborhood'], toronto_merged['Cluster Labels']):
    pos = folium.Popup(pos)
    folium.CircleMarker(
        [lat, lon],
        radius=6,
        popup=pos,
        color=cmap(cluster),
        fill=True,
        fill_color=cmap(cluster),
        fill_opacity=0.7).add_to(map_clusters)

map_clusters
```



<sup>^</sup>Since this is an additional piece of analysis, I will not be elaborating further in this report but I have included in my repository an example of how I have performed such clustering technique on the dataset, where I used K=5 as an example.