# Coursera Capstone

## IBM Applied Data Science Capstone

Opening a New Indian Restaurant in Mississauga, Canada

By Sameer Pradeep

March 2020

# Introduction

Canada has been one of the top destinations for Indian immigrants for a long time. Many of them have settled there and bought their culture with them including Indian restaurants. So, the aim of this project to suggest a neighborhood in the city of Mississauga in Canada to start a new Indian restaurant without facing much competition from the existing ones.

# Business Problem

The objective of this capstone project is to analyse and select the best locations in the city of Mississauga, Canada to open a new shopping mall. Using data science methodology and machine learning techniques like clustering, this project aims to provide solutions to answer the business question: In the city of Mississauga, Canada, if a property developer is looking to open a new Indian Restaurant, where would you recommend that they open it?

# Target Audience of this project

This project is particularly useful to property developers and investors looking to open or invest in new Indian Restaurant in the city of Mississauga.

To solve the problem, we will need the following data:

• List of neighbourhoods in Mississauga. This defines the scope of this project which is confined to the city of Mississauga in Canada.

• Latitude and longitude coordinates of those neighbourhoods. This is required in order to plot the map and also to get the venue data.

• Venue data, particularly data related to Indian Restaurants. We will use this data to perform clustering on the neighbourhoods.

# Sources of data and methods to extract them

This Wikipedia page (https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Mississauga) contains a list of neighbourhoods in Mississauga, with a total of 17 neighbourhoods. We will use web scraping techniques to extract the data from the Wikipedia page, with the help of Python requests and beautifulsoup packages. Then we will get the geographical coordinates of the neighbourhoods using Python Geocoder package which will give us the latitude and longitude coordinates of the neighbourhoods.

After that, we will use Foursquare API to get the venue data for those neighbourhoods. Foursquare has one of the largest database of 105+ million places and is used by over 125,000 developers.

Foursquare API will provide many categories of the venue data, we are particularly interested in the Shopping Mall category in order to help us to solve the business problem put forward. This is a project that will make use of many data science skills, from web scraping (Wikipedia), working with API (Foursquare), data cleaning, data wrangling, to machine learning (K-means clustering) and map visualization (Folium). In the next section, we will present the Methodology section where we will discuss the steps taken in this project, the data analysis that we did and the machine learning technique that was used.

Firstly, we need to get the list of neighbourhoods in the city of Kuala Lumpur. Fortunately, the list is available in the Wikipedia page (https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Mississauga). We will do web scraping using Python requests and beautifulsoup packages to extract the list of neighbourhoods data. However, this is just a list of names. We need to get the geographical coordinates in the form of latitude and longitude in order to be able to use Foursquare API. To do so, we will use the wonderful Geocoder package that will allow us to convert address into geographical coordinates in the form of latitude and longitude. After gathering the data, we will populate the data into a pandas DataFrame and then visualize the neighbourhoods in a map using Folium package. This allows us to perform a sanity check to make sure that the geographical coordinates data returned by Geocoder are correctly plotted in the city of Mississauga.
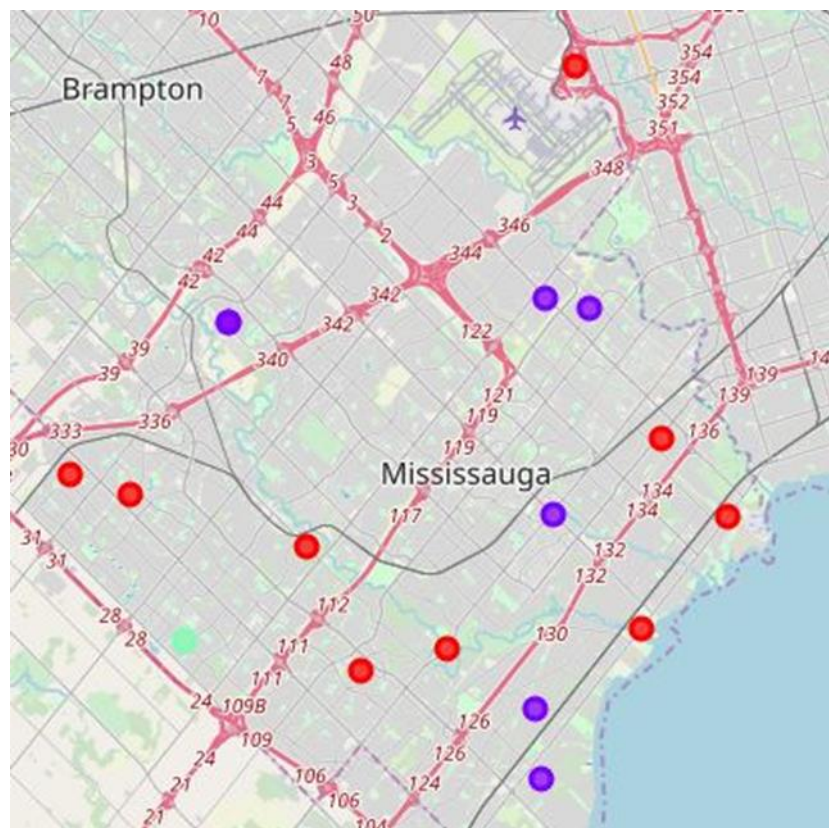
Next, we will use Foursquare API to get the top 100 venues that are within a radius of 2000 meters. We need to register a Foursquare Developer Account in order to obtain the Foursquare ID and Foursquare secret key. We then make API calls to Foursquare passing in the geographical coordinates of the neighbourhoods in a Python loop. Foursquare will return the venue data in JSON format and we will extract the venue name, venue category, venue latitude and longitude. With the data, we can check how many venues were returned for each neighbourhood and examine how many unique categories can be curated from all the returned venues. Then, we will analyse each neighbourhood by grouping the rows by neighbourhood and taking the mean of the frequency of occurrence of each venue category. By doing so, we are also preparing the data for use in clustering. Since we are analysing the "Indian Restaurant" data, we will filter the "Shopping Mall" as venue category for the neighbourhoods.

Lastly, we will perform clustering on the data by using k-means clustering. K-means clustering algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible. It is one of the simplest and popular unsupervised machine learning algorithms and is particularly suited to solve the problem for this project. We will cluster the neighbourhoods into 3 clusters based on their frequency of occurrence for "Indian Restaurant". The results will allow us to identify which neighbourhoods have higher concentration of Indian Restaurants while which neighbourhoods have fewer number of Indian Restaurants. Based on the occurrence of Indian Restaurants in different neighbourhoods, it will help us to answer the question as to which neighbourhoods are most suitable to open new Indian Restaurants.

The results from the k-means clustering show that we can categorize the neighbourhoods into 3 clusters based on the frequency of occurrence for "Indian Restaurants":

•        Cluster 0: Neighbourhoods with high concentration of Indian Restaurants

•        Cluster 1: Neighbourhoods with moderate number of Indian Restaurants

•        Cluster 2: Neighbourhoods with low number to no existence of Indian Restaurants

The results of the clustering are visualized in the map below with cluster 0 in red colour, cluster 1 in purple colour, and cluster 2 in mint green colour.

As observations noted from the map in the Results section, cluster 2 has very low number to no shopping mall in the neighbourhood of Churchill Meadows. This represents a great opportunity and high potential areas to open new Indian Restaurants as there is very little to no competition from the existing ones.

# Limitations and Suggestions for Future Research

In this project, we only consider one factor i.e. frequency of occurrence of Indian Restaurants, there are other factors such as population and income of residents that could influence the location decision of a new shopping mall. However, to the best knowledge of this researcher such data are not available to the neighbourhood level required by this project. Future research could devise a methodology to estimate such data to be used in the clustering algorithm to determine the preferred locations to open a new Indian Restaurant. In addition, this project made use of the free Sandbox Tier Account of Foursquare API that came with limitations as to the number of API calls and results returned. Future research could make use of paid account to bypass these limitations and obtain more results.

In this project, we have gone through the process of identifying the business problem, specifying the data required, extracting and preparing the data, performing machine learning by clustering the data into 3 clusters based on their similarities, and lastly providing recommendations to the relevant stakeholders i.e. property developers and investors regarding the best locations to open a new Indian Restaurant. To answer the business question that was raised in the introduction section, the answer proposed by this project is: The neighbourhood of Churchill Meadows is the most preferred location to open a new Indian restaurant. The findings of this project will help the relevant stakeholders to capitalize on the opportunities on high potential locations while avoiding crowded areas in their decisions to open a new Indian Restaurant.