Coursera Capstone Project IBM Data Science Specialization

Opening an Indian Cuisine Restaurant in New York City, USA

By: Tushar Singhal

August 2020

1. Introduction

- •While opening a restaurant can be a very lucrative business, a lack of demand causes many restaurants to close within the first year of opening.
- •There are many different factors that can account for a restaurant's success such as location, competition and quality of the food. This is an important question that every business owner must face when choosing whether to open a restaurant or not, as well as location of the business.
- •To demonstrate the process of picking a location for a client opening a business, the project will focus on answering were to open the restaurant.

1.2 Business Problem

•The following question: "If the client wanted to open an Indian Restaurant in New York City, what areas are the best options to open the restaurant?" For an Indian Restaurant, the location and competition are both determined by where the restaurant is opened.

1.3 Target audience of this project

•This project is mainly useful to an anyone who wishes to open an Indian Cuisine Restaurant in New York City. The insight from the project will be helpful for determining the best possible location of the restaurant. It will help it understanding whether there is a lot or little competition in a given neighbourhood of the city. Accordingly, the person who is interested in opening it can take the necessary decision with the help of the gathered intel.

2. Data

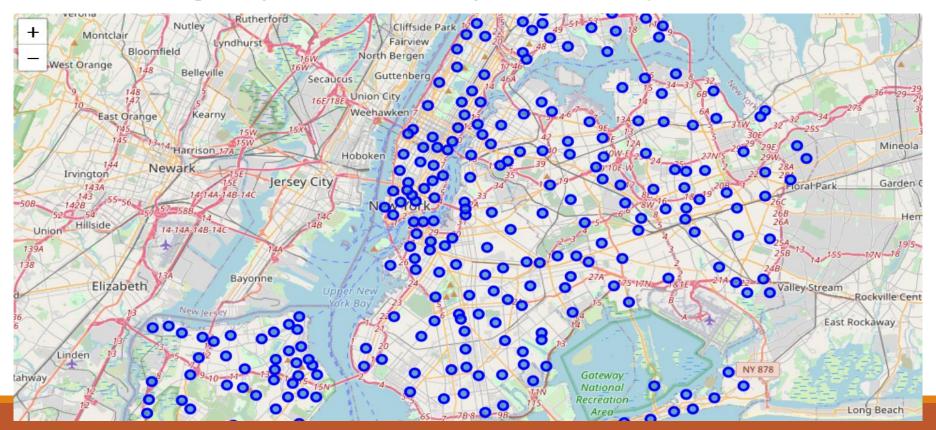
To solve the problem, the following data is needed:

- List of neighbourhoods in New York City. This defines the scope of this project which is confined to the city of New York.
- 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.

The data for the neighbourhoods of New York is obtained from previous week of the course. 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.

3. Methodology

- •First the json file containing the New York neighbourhood data set was loaded.
- •In total there are 306 unique neighbourhoods and using the folium library this is visualized as shown below



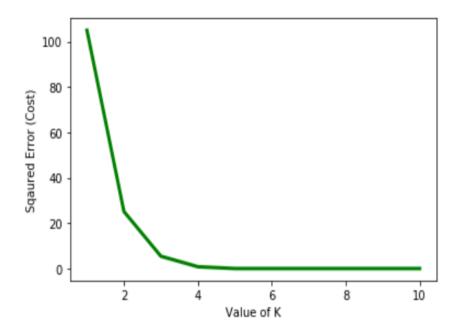
- Using the Foursquare API up to 100 venues in a radius of 500m in each neighbourhood is obtained. This is done with the help of REST API. The data is stored in nyc_venues data frame.
- In total there are 429 unique venue categories. The data frame is then onehot encoded according to venues.

	Neighborhood	Neighborhood Latitude	Neighborhood Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	Wakefield	40.894705	-73.847201	Lollipops Gelato	40.894123	- 73.845892	Dessert Shop
1	Wakefield	40.894705	-73.847201	Carvel Ice Cream	40.890487	-73.848568	Ice Cream Shop
2	Wakefield	40.894705	-73.847201	Walgreens	40.896528	- 73.844700	Pharmacy
3	Wakefield	40.894705	-73.847201	Rite Aid	40.896649	- 73.844846	Pharmacy
4	Wakefield	40.894705	-73.847201	Shell	40.894187	-73.845862	Gas Station

• Then the data frame is grouped by neighbourhood and summed. A new data frame consisting of only the Indian cuisine restaurants is created and the three different kinds of Indian Restaurants are summed to a single total.

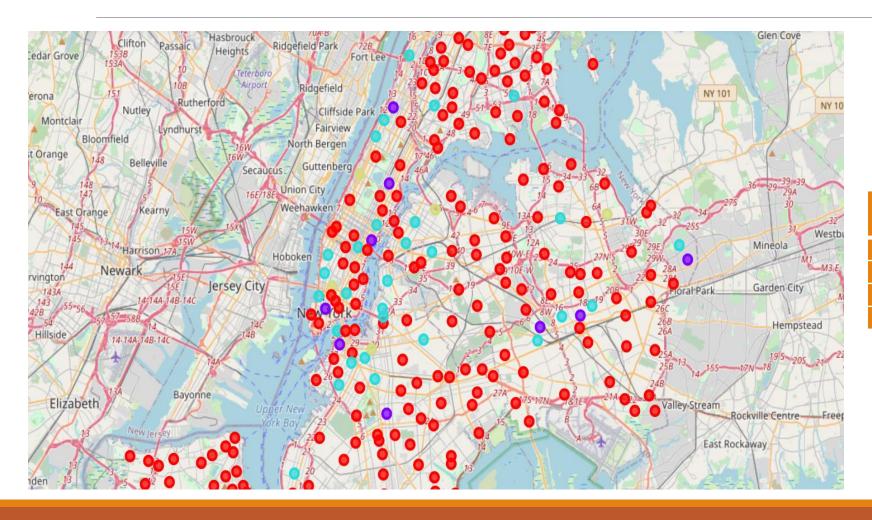
	Neighborhood	Total
0	Allerton	0
1	Annadale	0
2	Arden Heights	0
3	Arlington	0
4	Arrochar	0
5	Arverne	0
6	Astoria	3
7	Astoria Heights	0
8	Auburndale	0
9	Bath Beach	0
10	Battery Park City	0

- Used k-means clustering to identify the areas which have either a smaller number of restaurants or where they are clustered to find the optimal location for a new one.
- To determine the value of k a graph of cost Squared error vs k is plotted using matplot.pyplot library.



The elbow point is at 4 we can choose k as 4. After choosing k clustering is performed on the data to obtain the cluster labels. The obtained labels is merged with neighbourhood name and location data.

4. Result



Cluster Number	Colour	Number of
		Restaurants
0	Red	0
1	Purple	2
2	Teal	1
3	Gold	3 or more

5. Discussions

From the results of the clustering it can be inferred that the best neighbourhoods would be those which have either 0 or 1 Indian restaurant at the most. These neighbourhoods are represented by Red and Teal colour respectively. On the other hand red coloured clusters are predominantly in southern part of New York these neighbourhoods could be suburbs, so these could be avoided. The best areas are either the northern or Southern Manhattan were there are 0 or 1 restaurants. Even Brooklyn or Queens have good places to setup. At any cost Purple and Gold coloured clusters should be avoided as they won't lead to higher profits due to competition already existing there.

6. Conclusion

Opening a restaurant is a complex task that can lead to a large monetary loss if not done properly. Thus, extensive research about the area would greatly increase the likelihood of the restaurant succeeding. From the project above, I demonstrated the workflow necessary for a client to determine what area the restaurant should open. This project can further be improved upon by taking demography and income of the households, this will give an idea of the spending power in that neighbourhood.