# Data Acquisition and Cleaning

## Data Sources

This project works with two datasets. The first one is the data related to New York's Neighborhood and the data is available in json format. And the second dataset is the data related to Toronto's Neighborhood and this data is obtained from Wikipedia and the data is Scrapped using Web Scrapping Tool BeautifulSoup.

## Data Cleaning

The first data source as mentioned above is in json format. It initially consisted of many different classes of data. Upon examining them, the data that we are interested in was found under 'features' category. Therefore, we extracted the data and stored the data in a dataframe using pandas library. The second data source is a Wikipedia page that contains Postal Code of the city of Toronto in a wiki table. To scrape the data from the URL, BeautifulSoup has been used to extract the table data. After going through a few more steps, the dataframe was obtained which consists of: Postal Code, Borough and Neighborhood. But the problem with this dataframe was, it has some values under the column 'Borough' which were not assigned in the first place. So, the rows with no assigned value in the 'Borough' column were dropped.

## Feature Selection

Now that we have obtained the different neighborhoods and their respective geometric coordinates for the city of New York and Toronto, it is time to come up with different venues that the different venues have to offer. We use Foursquare API to do this job. Foursquare API provides with an access to an enormous database consisting of venues from all around the world including rich variety of information such as addresses, tips, photos and comments. Having signed up for a Foursquare developer, using the Client ID and Client Secret, it is possible to make API requests in order in order to retrieve venue information. By feeding a function with Neighborhood name and its geometric coordinates, using Foursquare API different venues (Restaurants, Coffee shops, etc) were extracted. After performing One-HotEncoding and grouping together the rows by neighborhoods, the NY dataset and Toronto dataset. Both the dataframes were combined into a single dataframe in order to perform clustering operation.