

Clustering of Similar Neighborhoods in Different Cities

New York City and Toronto March 2020

Introduction

- Purpose: Explore and investigate the distribution of different venues in New York City and Toronto.
- Analysis to help new visitors to stay in better place.
- New York welcomes 65.2 million Tourists
- Toronto welcomes 27.5 million Tourists
- This will help understand the 2 cities better.



Data Preprocessing



FILES(GEOLOCATION DATA)

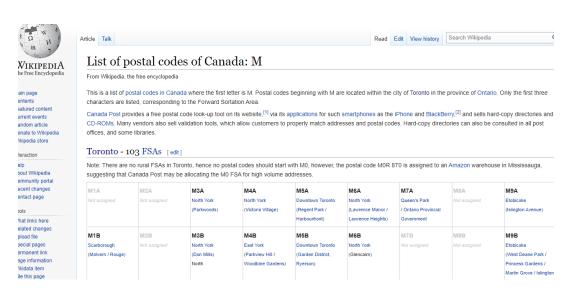


HTTPS://COCL.US/GEOSPATIAL DATA (.CSV) -- USED IN TORONTO.IPYNB



HTTPS://COCL.US/NEW_YORK_DATASET (.JSON) -- USED IN NEWYORKCITY.IPYNB

Data Preprocessing





	Postcode	Borough	Neighborhood
0	МЗА	North York	Parkwoods
1	M4A	North York	Victoria Village
2	M5A	Downtown Toronto	Harbourfront
3	M6A	North York	Lawrence Heights
4	M6A	North York	Lawrence Manor
5	M7A	Downtown Toronto	Queen's Park
6	M9A	Etobicoke	Islington Avenue
7	M1B	Scarborough	Rouge
8	M1B	Scarborough	Malvern
9	МЗВ	North York	Don Mills North
10	M4B	East York	Woodbine Gardens
11	M4B	East York	Parkview Hill
12	M5B	Downtown Toronto	Ryerson
13	M5B	Downtown Toronto	Garden District
14	M6B	North York	Glencairn
15	M9B	Etobicoke	Cloverdale
16	M9B	Etobicoke	Islington
17	M9B	Etobicoke	Martin Grove
18	M9B	Etobicoke	Princess Gardens
19	M9B	Etobicoke	West Deane Park
20	M1C	Scarborough	Highland Creek
21	M1C	Scarborough	Rouge Hill
22	M1C	Scarborough	Port Union

Data sets

Geolocation Data

- NY_Neighborhood_Locations.csv
- Toronto_Neighborhood_Locations.csv

One hot coded with Neighbourhood and nearby venues

- NY_grouped.csv
- Toronto_grouped.csv

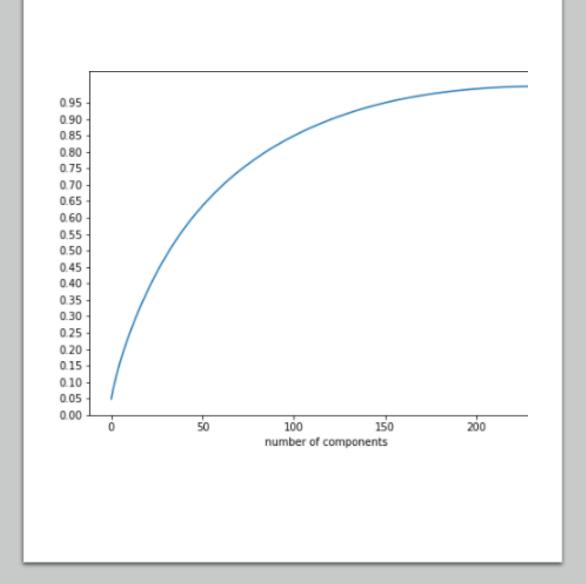
Packages used for Analysis

```
In [1]: import pandas as pd
import numpy as np

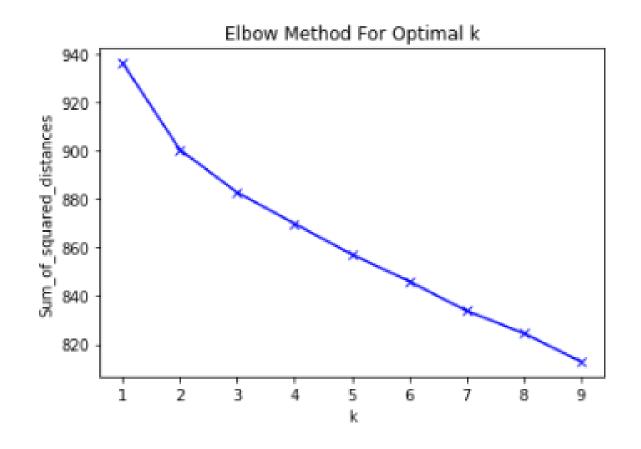
import sklearn
from geopy.geocoders import Nominatim
import folium
import matplotlib.pyplot as plt
from sklearn.cluster import KMeans
from sklearn.decomposition import PCA
from sklearn.preprocessing import MinMaxScaler
from sklearn.metrics import silhouette_samples
from sklearn.metrics import silhouette_score
```

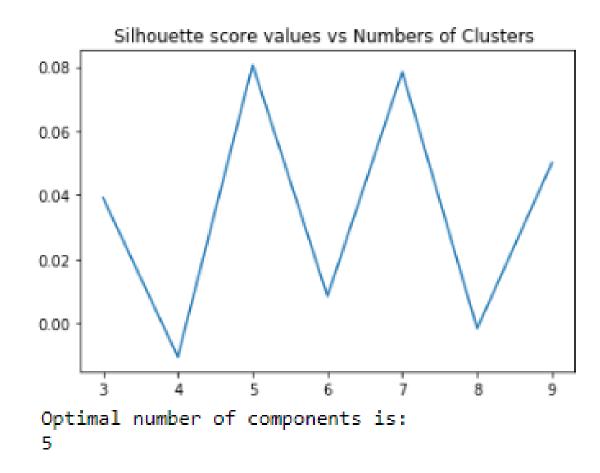
PCA

- Principal Component Analysis is a dimension reduction tool that can be used to reduce a large set of variables to a smaller set by grouping the similar variables that explains the whole dataset.
- With the use of MinMaxScalar functionality of Scikit learn preprocessing, it scales the dataset such as all features values are in the range {0,1}. This retains the variance of the dataset.
- With this it was safe to reduce the number of components to 175 retaining the variance.



K-Means(Optimal K clusters)





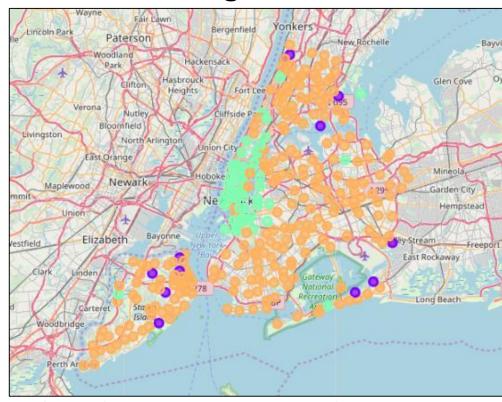
Results

NY before clustering



Neighborhood Locations

NY after clustering



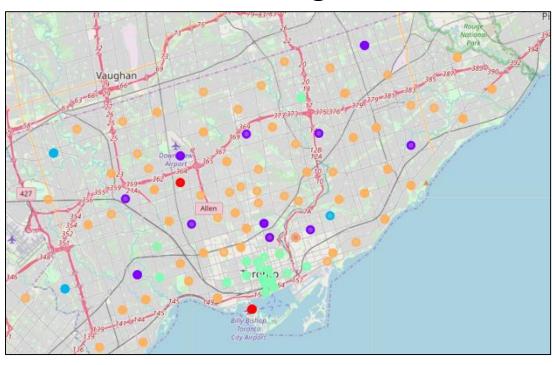
Similar colors represent same clusters

Results

Toronto before clustering



Toronto after clustering



Neighborhood Locations

Similar colors represent same clusters

Results



We can see how the same clusters in Toronto and New York City represent how these places are similar being in different cities



We can also see there are few clusters which are not available in New York city which can be unique to Toronto

- People travelling from one country to another can choose to stay in places that they are habitual
- It would benefit them with convenience and understanding of where to stay next.

Conclusion