

Capstone Project

The Battle of Neighborhoods Reports.

Zaheer Habib, Coursera final project submission 2019

Introduction.

Description of the background of the Project

New York is the popular city of the United States with an estimate of around population of 8,398,748 which is distributed over 302.6 square miles. It is also the most densely populated city which is in the Southern tip and city is the center of NY metropolitan area.

NYC consists of five boroughs – Brooklyn, Queens, Manhattan, The Bronx, and Staten Island. Many districts and landmarks in NYC are famous; including three of the world's ten most visited tourist attractions in 2013, around 62.8 million tourists visited in 2017.

< https://en.wikipedia.org/wiki/New_York_City >

"New York City is multicultural. About 36% of the city's population is foreign-born,^[23] one of the highest among US cities. The eleven nations constituting the largest sources of modern immigration, to New York City are to New York City are the Dominican Republic, China, Jamaica, Guyana, Mexico, Ecuador, Brazil, Haiti, Trinidad and Tobago, Colombia, Russia and El Salvador".^[24]

https://en.wikipedia.org/wiki/Demographics_of_New_York_City

Requirement analysis.

An entrepreneur from China likes to invest in New York City (NYC); he is a very enthusiastic businessman having sound business implementation knowledge. So, he contacts with the local technology company to explore the best area in New York which he will invest so that he gets quick return on investment (ROI). Subsequently, to get maximum ROI he pushes developer to find the area in NYC where he will establish his business but only constraint is, he likes to invest in the business which is traditionally belongs to Chinese culture. Example Chinese restaurant, SPA or retail shop which deal in Chinese herbs. So below are the two requirements which need consideration.

- Investment should need to make in NYC best potential borough and its neighborhood area to get maximum ROI
- Technology advise entrepreneur in related to Chinese culture like Chinese restaurant, SPA or retail herbal shop etc.

Problem

The Data that might contribute to determining the explore information to get best location which will further helps to identify the metrics that best describe what kind of business he supposed to put his capital. This project aims to predict the area and type of business which the entrepreneur should needs to invest to get maximum ROI.

Interest.

This project aims to predict the area and type of business which the entrepreneur should needs to invest to get maximum ROI.

Data acquisition and cleaning.

The download Json file is processed accordingly as require and all outliers, irrelevant data are filter out before use for next stage of the analysis to make process more efficient and error free.

Target variables, lat & long is created as require and further filter out at later stages when this is not further required. Tabular data is created throughout the project for analysis data and better understating of the derived information which helps to further drilled down as required.

To consider the problem we get through the below sites to get data

- Forsquare API to get the most common venues of given Borough and neighborhood.
- Coursera Lab NYC borough json data: https://geo.nyu.edu/catalog/nyu_2451_34572

Methodology & Exploratory data analysis

To come up with the finding we were using different concepts, tools and methods which is mention below.

- GitHub repository is used to deposit files and share our work with community
- Tabular data which has main component consist of Borough,
- Tabular data which has main component consist of Venue Frequency
- Forsquare API to explore the neighborhood of the borough
- **folium** library to visualize geographic details through choropleth plotting.
- Geolocator is used to get geocode(address) latitude and longitude.

In summary section, one of our aim is to visualize the requirement to plot in choropleth style; when reviewing all the basic requirement we were consider all these problems, we created a initial map using choropleth style to plot our existence borough and further NEW choropleth style map is created based on derived data which is done according to clusters which is generated through k-means algorithm where this clustered is highlighted according to the venue density in respective brought. This derived information have the information which is help to analysis and created our final map and different tabular information.

- Borough name
- Cluster name
- Venue
- Frequency
- Lat & Long

Neighborhood has a total of 5 boroughs and 306 neighborhoods. In order to segment the neighborhoods and explore them, we will essentially need a dataset that contains the 5 boroughs and the neighborhoods that exist in each borough as well as the latitude and longitude coordinates of each neighborhood.

Exploring newyork_data json file through features method

```
{'type': 'Feature',
'id': 'nyu_2451_34572.1',
'geometry': {'type': 'Point',
'coordinates': [-73.84720052054902, 40.89470517661]},
'geometry_name': 'geom',
'properties': {'name': 'Wakefield',
'stacked': 1,
'annoline1': 'Wakefield',
'annoline2': None,
'annoline3': None,
'annoangle': 0.0,
'borough': 'Bronx',
'bbox': [-73.84720052054902,
40.89470517661,
-73.84720052054902,
40.89470517661]}}
```

Exploring neighborhoods: total 5 boroughs and 306 neighborhoods hit

[21]:	Borough	Neighborhood	Latitude	Longitude
0	Bronx	Wakefield	40.894705	-73.847201
1	Bronx	Co-op City	40.874294	-73.829939
2	Bronx	Eastchester	40.887556	-73.827806
3	Bronx	Fieldston	40.895437	-73.905643
4	Bronx	Riverdale	40.890834	-73.912585

Get credential from Foursquare Credentials and Version

```
Your credentails:
CLIENT_ID: SGSD5IWJSV43PQ42KAK4BQXDNKQA2R2X2HP3LAV4FGQJFA0A
CLIENT_SECRET: 0XBMMNUB0PGQUMPKNZWXEZXMZ3BVMETQR5FSY4YAQSXJ5SFU
```

Get Latitude and longitude values of Marble Hill are 40.87655077879964, -73.91065965862981 and generate map 100 topmost neighborhood near 500 radius. When send request to foursquare we get 25 venue.

```
]:
```

	name	categories	lat	lng
0	Arturo's	None	40.874412	-73.910271
1	Bikram Yoga	None	40.876844	-73.906204
2	Tibbett Diner	None	40.880404	-73.908937
3	Starbucks	None	40.877531	-73.905582
4	Dunkin'	None	40.877136	-73.906666

And how many venues were returned by Foursquare?

```
]:
```

```
print('{} venues were returned by Foursquare.'.format(nearby_venues.shape[0]))
```

25 venues were returned by Foursquare.

Visualized the cluster through k-means.

Finally, we break our analysis into five clusters and each cluster is color coded for the ease of presentation to understand their neighborhood. We can see that majority of the neighborhood represents different colors for different clusters. These neighborhoods have their own cluster (Blue, Red, Purple and Yellow). This color scheme helps to name our cluster based on the venue and neighborhood.

Among five clusters below are the details.

- Cluster 4 is the largest cluster which has 20 neighborhoods while Cluster 1 is the second largest having 17.
- Cluster 2 and 5 only have 1 neighborhood
- Cluster 3 has 12 neighborhoods.