

Coursera Capstone

Stewart Foodservice Inc. POS in Beirut, Lebanon

Detailed Report

By Ayham Mhd

Contents

| 1. | Introduction |
|-------|---------------------------|
| 1.1. | Stake holder overview |
| ••••• | |
| 1.2. | Problem Definition |
| 1.3. | Initial thinking |
| 2. | Description of data |
| 3. | Methodology |
| 3.1. | Preparing the environment |
| 3.2. | Exploring Beirut City |
| 3.3. | Exploring Neighborhoods |
| 3.4. | Clustering Neighborhoods |
| 4. | The Result8 |

1. Introduction

This report illustrates Coursera Capstone related to IBM Data science course.

1.1. Stake holder overview

Stewart Foodservice is owned by industry veteran Aubrey MacMillan. He started out in the business over 50 years ago at Burgess Wholesale, then an egg-grading station, working his way up through deliveries, sales and management to owner of what has become one of the leading Canadian owned foodservice distributors.



1.2. Problem Definition

Stewart Foodservice Inc, who supplies restaurants and so on, decided to open a store or a point of sale POS in Beirut, Lebanon, and asked me an engineer having some background in data science that collected from IBM data science courses on Coursera, to suggest the place where to install this POS.

1.3. Initial thinking

As the company is the first point in food industry supply chain, and because the distance between the supplier and the client can shorten the time to market (TTM) which is a cost, in addition, the more this distance the more the cost of delivery, so I thought that selecting the right closest place of POS to the targeted market (restaurant, coffee shops and so forth) could be the most important factor.

p.s. TTM is the length of time it takes from a product being conceived until its being available for sale (Wikipedia).

2. Description of data

I've collected data from Forsquare API using my credentials like the following:

```
CLIENT_ID = 'My Foursquare ID'
CLIENT_SECRET = 'Foursquare Secret'
VERSION = '20180604'
LIMIT = 100
RADIUS = 2000

request_parameters = {
    "client_id": CLIENT_ID,
    "client_secret": CLIENT_SECRET,
    "v": '20180605',
    "section": "restaurant",
    "near": "Beirut",
    "radius": RADIUS,
    "limit": LIMIT}

data = requests.get("https://api.foursquare.com/v2/venues/explore", params=request parameters)
```

This gave me a comprehansive data enabling me to extract many meaningful pointers, I extracted the data filtered as I want certain categories of venues, i.e. only related to targeted customers who are restaurants, coffee shops and so forth, and I chose the following useful fields: name of venue, category, longitude and latitude as the following:

| Ing | lat | categories | name |
|-----------|-----------|---------------------------|----------------------|
| 35.478505 | 33.896835 | Bar | Moscow Mule |
| 35.479906 | 33.895523 | Café | Cafe Younes |
| 35.477883 | 33.896285 | Burger Joint | Classic Burger Joint |
| 35.479539 | 33.896052 | Diner | Roadster Diner |
| 35.478478 | 33.897890 | Middle Eastern Restaurant | Socrate |

3. Methodology

It was a descriptive research consists of the following steps after introductory part:

3.1. Preparing the environment

I used *numpy* to handle data in a vectorized manner, and will use *pandas* library for data analysis and json library to handle JSON files, and will use geopy to be able to use *geocoders.Nominatim* to convert an address into latitude and longitude values, and *requests* library to handle requests, and *json normalize* that will transform JSON file into a pandas data frame.

I used *Matplotlib* and associated plotting modules, and *k-means* for clustering stage.

Then I used *folium* map rendering library.

First importing preregisites

```
import numpy as np # library to handle data in a vectorized manner
import pandas as pd # library for data analsysis
pd.set_option('display.max_columns', None)
pd.set_option('display.max_rows', None)
import json # library to handle JSON files
!conda install -c conda-forge geopy --yes
from geopy.geocoders import Nominatim # convert an address into Latitude and Longitude values
import requests # library to handle requests
from pandas.io.json import json normalize # tranform JSON file into a pandas dataframe
# MatplotLib and associated plotting modules
import matplotlib.cm as cm
import matplotlib.colors as colors
#import k-means from clustering stage
from sklearn.cluster import KMeans
!conda install -c conda-forge folium=0.5.0 --yes
import folium # map rendering library
print('Prerequisit Packages were Loaded...')
```

3.2. Exploring Beirut City

I used geopy library to get the latitude and longitude values of Beirut City,

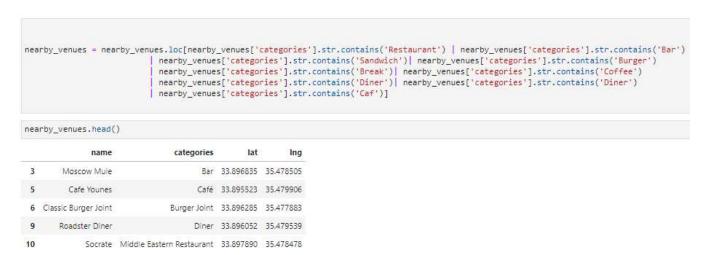
```
address = 'Beirut'
geolocator = Nominatim(user_agent="ny_explorer")
location = geolocator.geocode(address)
latitude = location.latitude
longitude = location.longitude
print('The geograpical coordinate of Beirut City are {}, {}.'.format(latitude, longitude))
The geograpical coordinate of Beirut City are 33.8959203, 35.47843.
```

then using my Forsquare credentials I got the top 100 venues that are in Beirut within a radius of 500 meters.

Then I used a function that extracts the category of the venue then clean the json and structure it into a pandas data frame.

```
#From Foursquare Lab all the information is in the items key. Before we proceed, let's borrow the get_category_type function from the Foursquare Lab.
# function that extracts the category of the venue
def get_category_type(row):
        categories_list = row['categories']
    except:
        categories_list = row['venue.categories']
    if len(categories_list) == 0:
        return None
        return categories_list[0]['name']
#Now we are ready to clean the json and structure it into a pandas dataframe. venues = results['response']['groups'][0]['items']
nearby_venues = pd.json_normalize(venues) # flatten JSON
# filter columns
filtered_columns = ['venue.name', 'venue.categories', 'venue.location.lat', 'venue.location.lng']
nearby_venues =nearby_venues.loc[:, filtered_columns]
# filter the category for each row
nearby_venues['venue.categories'] = nearby_venues.apply(get_category_type, axis=1)
nearby_venues.columns = [col.split(".")[-1] for col in nearby_venues.columns]
nearby_venues.head()
```

Then I filtered the result categories to match our target (Restaurants and etc.)



3.3. Exploring Neighborhoods

I used a function to repeat the same process to all the neighborhoods in Beirut using Foursquare API as well, creating a new data frame which also must be filtered.

| | Neighborhood | Neighborhood Latitude | Neighborhood Longitude | Venue | Venue Latitude | Venue Longitude | Venue Category |
|---|--------------|-----------------------|------------------------|----------------------|----------------|-----------------|---------------------------|
| 3 | Moscow Mule | 33.896835 | 35.478505 | Moscow Mule | 33.896835 | 35,478505 | Bar |
| 5 | Moscow Mule | 33.896835 | 35.478505 | Cafe Younes | 33.895523 | 35.479906 | Café |
| 6 | Moscow Mule | 33.896835 | 35.478505 | Fiber | 33.898009 | 35.479450 | Restaurant |
| 7 | Moscow Mule | 33.896835 | 35,478505 | Classic Burger Joint | 33.896285 | 35,477883 | Burger Joint |
| 8 | Moscow Mule | 33.896835 | 35,478505 | Socrate | 33.897890 | 35,478478 | Middle Eastern Restaurant |

Then I analyzed each neighborhood by one hot encoding.

Analyze Each Neighborhood



Then grouping rows by neighborhood by taking the mean of the frequency of occurrence of each category, then putting that into a pandas data frame.



3.4. Clustering Neighborhoods

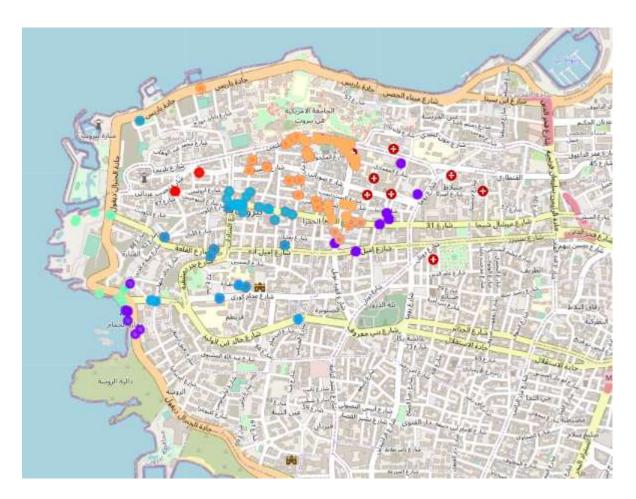
I clustered neighborhoods using Kmean as following:

Cluster Neighborhoods

```
t_grouped_clustering = t_grouped.drop('Neighborhood', 1)
# run k-means clustering
kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(t_grouped_clustering)
# check cluster labels generated for each row in the dataframe kmeans.labels_[8:10]
array([0, 3, 1, 4, 2, 0, 0, 4, 1, 0], dtype=int32)
#Let's create a new dataframe that includes the cluster as well as the top 10 venues for each neighborhood. 
# add clustering labels 
neighborhoods_venues_sorted.insert(0, 'Cluster Labels', kmeans.labels_)
t merged = t venues
# merge toronto_grouped with toronto_data to add Latitude/Longitude for each neighborhood
t_merged = t_merged.join(neighborhoods_venues_sorted.set_index('Neighborhood'), on='Neighborhood')
t_merged.head() # check the Last columns!
                                                                                                                                                                                                                                                                                       10th Most
Common
Venue
                                                                                                                                                                3rd Most
Common
Venue
   Neighborhood Neighborhood Latitude Longitude
                                                           Venue Venue Venue Latitude Longitude
                                                          Moscow
Mule 33.896835 35.478505
                          33.896835
                                            35.478505
                                                           Cafe 33.895523 35.479906
                          33.896835
                                            35.478505
                                            35.478505
                                                                                                                                                                                                                                                                             Bar Breakfast Spot
                                                                                                                                                                                                                  Fast Food
7 Moscow Mule 33.896835 35.478505 Burger 33.896285 35.477883 Burger Joint 0 Restaurant
```

and visualize the resulting clusters using Folium library:

```
map_clusters = folium.Map(location=[latitude, longitude], zoom_start=15)
# set color scheme for the clusters
x = np.arange(kclusters)
ys = [i + x + (i*x)**2 for i in range(kclusters)]
colors_array = cm.rainbow(np.linspace(0, 1, len(ys)))
rainbow = [colors.rgb2hex(i) for i in colors_array]
# add markers to the map
markers_colors = []
for lat, lon, poi, cluster in zip(t_merged['Venue Latitude'], t_merged['Venue Longitude'], t_merged['Neighborhood'], t_merged['Cluster Labels']):
label = folium.Popup(str(poi) + 'Cluster' + str(cluster), parse_html=True)
     folium.CircleMarker(
         [lat, lon],
         radius=5.
         popup=label,
         color=rainbow[cluster-1],
         fill=True,
         fill_color=rainbow[cluster-1],
         fill_opacity=0.7).add_to(map_clusters)
map_clusters
```



4. The Result

I chose the cluster that have the largest number of neighborhoods, get the average coordinates to plot the suggested POS location, this way it will be close to maximum number of the customers, and by this way I can consider that I had chosen the optimal place for the POS because of reducing TTM and delivery and shipping fees.

| _merged.groupby(t_merged['Cluster Labels']).count() | | | | | | | | | | | | | | | | | |
|---|--------------|--------------------------|---------------------------|-------|-------------------|--------------------|-------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|-----------------------------|---------------------------|
| | Neighborhood | Neighborhood Latitude | Neighborhood Longitude | Venue | Venue Latitude | Venue Longitude | Venue Category | 1st Most Common Venue | 2nd Most Common Venue | 3rd Most Common Venue | 4th Most Common Venue | 5th Most Common Venue | 6th Most Common Venue | 7th Most Common Venue | 8th Most Common Venue | 9th Most Common Venue | 10th Mos Commo Venu |
| Cluster Labels | | | | | | | | | | | | | | | | | |
| 0 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 835 | 83 |
| 1 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 | 534 |
| 2 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 133 | 13 |
| 3 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 29 | 2 |
| 4 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 | 654 |

```
Lat = cluster0['Venue Latitude'].mean()
Lng = cluster0['Venue Longitude'].mean()
print('Latitude is {}, and Longitude is {}'.format(Lat, Lng))
Latitude is 33.89664370405219, and Longitude is 35.479637994214215
```

Let's add a circle marker where suggested POS place

```
folium.CircleMarker(
    [Lat, Lng],
    radius=15,
    popup="Suggested POS Place!",
    color='green',
    fill=True,
    fill_color='#3186cc',
    fill_opacity=0.7,
    parse_html=False).add_to(map_clusters)
map_clusters
```

