

IBM Data Science

Analyzing Neighborhoods of Manchester For Starting A New Restaurant.

Capstone Project

Karan Bhatia

27/02/2022

Contents

Contents.....	1
Introduction	2
Business Problem.....	2
Data	2
Neighbourhoods Data	2
Geographical Coordinates.....	3
Venue Data from FourSquare	3
Methodology.....	3
Feature Extraction.....	3
Unsupervised Learning	4
Plotting.....	4
Results.....	5
Discussion.....	6
Conclusion.....	6

Introduction

Manchester is one of the biggest cities in the UK and also one of the most populous cities of United Kingdom. It lies within the United Kingdom's second-most populous urban area, with a population of 2.9 million and third-most populous metropolitan area, with a population of 3.3 million.

The population of Manchester comprises of people of various ethnicities from all over the world. The city is full of restaurants serving thousands of hungry customers, every day. The diversity in the population of the city has brought in a vast diversity in food habits of people. In this project we will study the neighbourhoods and make recommendations accordingly.

Business Problem

Our client is an investor who is interested in investing in a restaurant in Manchester. They have approached us to study the market and suggest them a location in one of the neighbourhoods which would be in best interest of the business. Our main objectives of this project would be to extract and analyse right data about various neighbourhoods of Manchester using various data science techniques and suggest our client a fitting location for their restaurant.

Data

In order to achieve our final goal we will need the following data:

- Neighbourhoods of Manchester.
- Geographical coordinates of the neighbourhoods.
- Venue data from FourSquare.

Neighbourhoods Data

This data was extracted from Areas of Manchester Wikipedia page (https://en.wikipedia.org/wiki/Category:Areas_of_Manchester) using web scraping with BeautifulSoup library in Python. This will give us a detailed list of neighbourhoods present in Manchester.

Geographical Coordinates

Later, the geographical coordinates of various neighbourhoods were extracted using GeoPy library in Python. Geographical coordinates are necessary for plotting maps during the project for visualizing our data. After using GeoPy we added two columns to our dataframe with latitude and longitude information of each neighbourhood as shown below:

	Neighbourhood	Latitude	Longitude
0	Baguley	53.399090	-2.285610
1	Barlow Moor	53.422164	-2.245970
2	Belle Vue, Manchester	42.955859	-71.459019
3	Benchill	53.381730	-2.261250
4	Beswick, Manchester	53.478390	-2.200320

Venue Data from FourSquare

Later we extracted venue data using FourSquare API. This venue data was used to study the venues in various neighbourhoods in Manchester. This data provided important details of various restaurants in the area and helped us understand the competition. This data was very important because it helped us draw the main conclusion of the project.

Methodology

Feature Extraction

Feature extraction was carried out through One Hot Encoding. In this method, each feature is a category that belongs to a venue which is then converted into binary, this means that 1 means this category is found in the venue and 0 means the opposite. Then, all the venues are grouped by the neighbourhoods, computing at the same time the mean. This will give us a venue for each row and each column will contain the frequency of occurrence of that particular category.

```
man_1hot = pd.get_dummies(explore_man[['Venue Category']], prefix="", prefix_sep="")

# Add neighbourhood column back to dataframe
man_1hot['Neighbourhood'] = explore_man['Neighbourhood']

# Move neighbourhood column to the first column
fixed_columns = [man_1hot.columns[-1]] + man_1hot.columns[:-1].values.tolist()
man_1hot = man_1hot[fixed_columns]

man_1hot.head()
```

Unsupervised Learning

Unsupervised learning was carried out in order to find out the similarities between found similarities between neighbourhoods. K-Means, a clustering algorithm, was implemented. In this case K-Means is used due to its simplicity and its similarity approach to find patterns.

- **K-Means:** K-Means is a clustering algorithm. This algorithm search clusters within the data and the main objective function is to minimize the data dispersion for each cluster. Thus, each group found represents a set of data with a pattern inside the multi-dimensional features. It is necessary for this algorithm to have a prior idea about the number of clusters since it is considered an input of this algorithm. For this reason, the elbow method is implemented. A chart that compares error vs number of cluster is done and the elbow is selected. Then, further analysis of each cluster is done.

```
max_range = 15 #Max range 15 (number of clusters)

from sklearn.metrics import silhouette_samples, silhouette_score

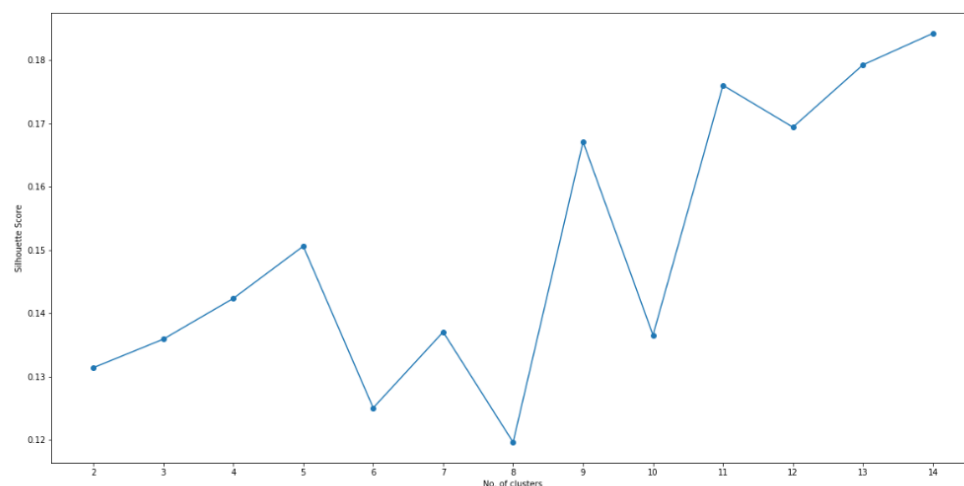
indices = []
scores = []

for man_clusters in range(2, max_range) :

    # Run k-means clustering
    man_gc = man_grouped_clustering
    kmeans = KMeans(n_clusters = man_clusters, init = 'k-means++', random_state = 0).fit_predict(man_gc)

    # Gets the score for the clustering operation performed
    score = silhouette_score(man_gc, kmeans)

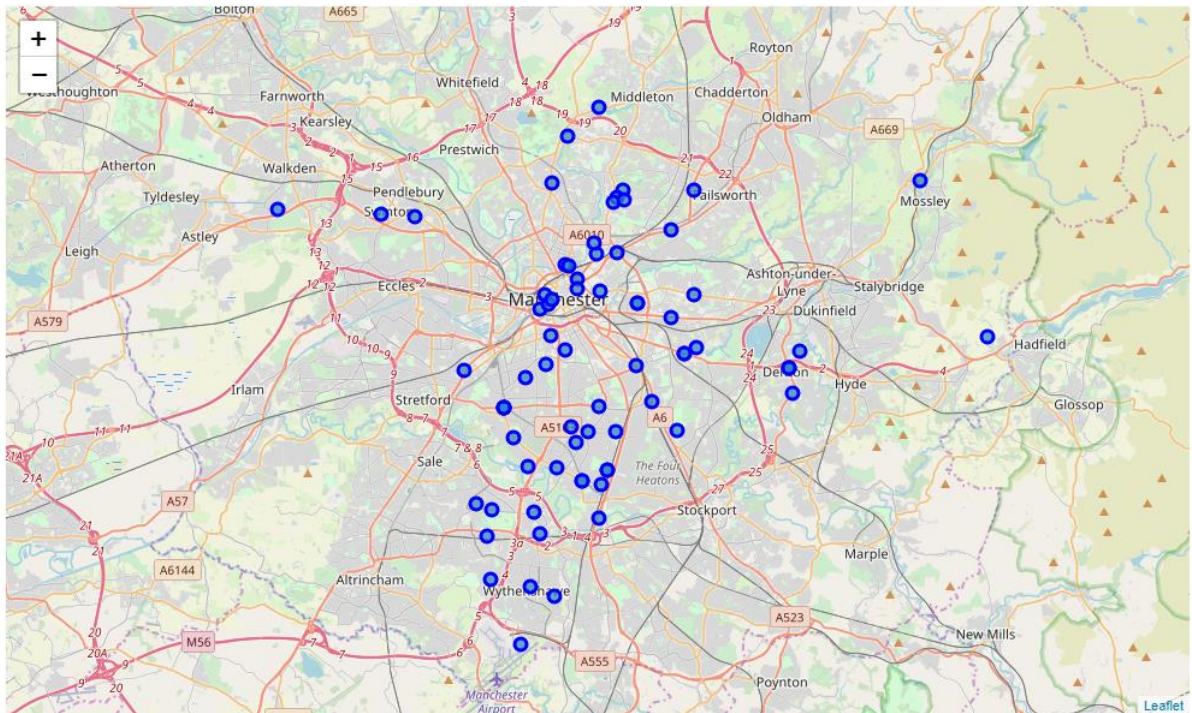
    # Appending the index and score to the respective Lists
    indices.append(man_clusters)
    scores.append(score)
```



Plotting

Various plotting techniques we used as well in order to visualize the data. Visualizing data often gives a clear understanding of the data as it is easier to spot patterns in a visualized data as compares to quantitative data.

- Folium: Folium library was used to plot maps of Manchester city as well as neighbourhoods. Folium was also used to visualize the cluster data.



Results

The above mentioned, K-Means clustering method was applied to the dataframe of neighbourhoods of Manchester city. As mentioned earlier the number of clusters that was derived from elbow method was 8. The code as well as plotting of clusters can be seen below:

```
map_clusters = folium.Map(location=[latitude, longitude], zoom_start=11)

# Setup color scheme for different clusters
x = np.arange(man_clusters)
ys = [i + x + (i*x)**2 for i in range(man_clusters)]
colors_array = cm.rainbow(np.linspace(0, 1, len(ys)))
rainbow = [colors.rgb2hex(i) for i in colors_array]

markers_colors = []
for lat, lon, poi, cluster in zip(man_final['Latitude'], man_final['Longitude'], man_final['Neighbourhood'],
                                  man_final['Cluster Labels']):
    label = folium.Popup(str(poi) + ' (Cluster ' + str(cluster + 1) + ')', parse_html=True)
    map_clusters.add_child(
        folium.features.CircleMarker(
            [lat, lon],
            radius=5,
            popup=label,
            color=rainbow[cluster-1],
            fill=True,
            fill_color=rainbow[cluster-1],
            fill_opacity=0.7))

map_clusters
```