

A thick black L-shaped frame is positioned around the text. It starts at the top left, goes right, then down, then right again, forming a partial rectangular border around the central text.

HOW TO SUPPORT DIFFERENT
VISITORS TO LIST ISTANBUL
DISTRICTS THAT FIT THEIR NEEDS
IN TERM OF CULINARY/ FOOD
VENUES

Discussion of the background

- Istanbul has a lot of culinary types that offer a variety of domestic food and also a choice of food from around the world and places to eat in 39 different districts consisting of 25 districts sit on the European side, and 14 rest on the Asian side. Because of the enormous size of these districts, to get to know them all would take a lifetime so it is very difficult for newcomers especially those who want to go to restaurants to choose the type of restaurant that suits their tastes from the wide variety of information in the media.
- In this Capstone Project, the problem solving steps will be carried out on how to utilize Foursquare location data and use one of the machine learning techniques, namely clustering to make an analysis and decision to determine which neighborhoods are in accordance with the tastes of customers in the city of Istanbul when the big event was held.

Data

In this Capstone Project we need the data needed to be processed and analyzed, including

- **List of districts of Istanbul**

Source : https://en.wikipedia.org/wiki/List_of_districts_of_Istanbul

The list of districts of Istanbul data in the form of this table will be taken from Wikipedia through the scrapping method then that table will be cleaned, explored, and processed then to add coordinates in each district in Istanbul can automatically use the geocoder class of Geopy client.

- **Food Venue/Restaurants in each neighborhoods of Istanbul City**

Source : Foursquare APIs

By using this Foursquare APIs after logging in using registered credentials, we will get all the closest places from all neighborhoods in Istanbul City so that we can get restaurants by filtering the data of the closest places from all the neighborhoods.

Data Preparation

1. Scraping List of Districts of Istanbul Table from Wikipedia

3.1 Use pandas to transform the data in the table on the Wikipedia page into a dataframe

```
[2]: df = pd.read_html('https://en.wikipedia.org/wiki/List_of_districts_of_Istanbul')[0] #[0] means the first table on the website  
df
```

Data Preparation

2. Performed Cleaning and Manipulation of the Data

3.2 Cleaning and Manipulation the Data

```
[3]: # rename column Population(2019) with Population
df.rename(columns={"Population (2019)": "Population"}, inplace=True)

# delete the rows with label 39,40,41,42
df.drop([39, 40, 41, 42], axis=0, inplace=True)

df
```

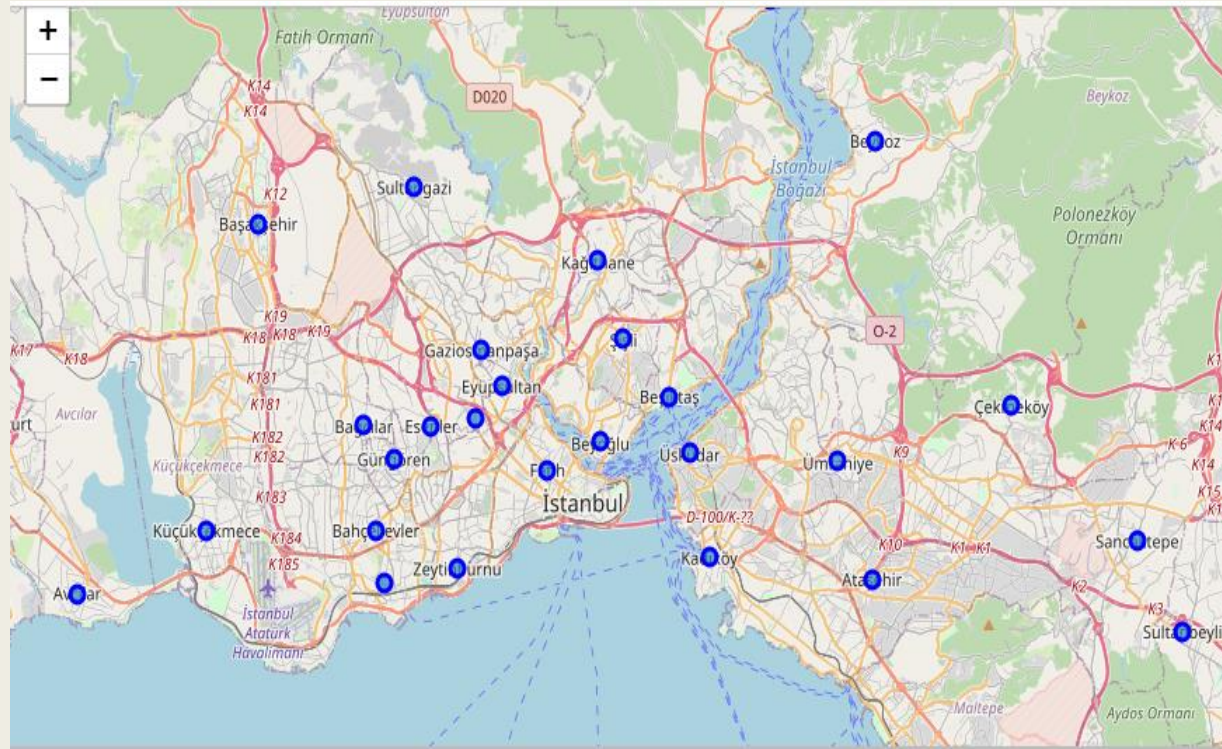
	District	Population	Area (km ²)	Density (per km ²)
0	Adalar	15238	11.05	1379
1	Arnavutköy	282488	450.35	627
2	Ataşehir	425094	25.23	16849
3	Avcılar	448882	42.01	10685
4	Bağcılar	745125	22.36	33324
5	Bahçelievler	611059	16.62	36766
6	Bakırköy	229239	29.64	7734
7	Başakşehir	460259	104.30	4413
8	Bayrampaşa	274735	9.61	28588
9	Beşiktaş	182649	18.01	10142
10	Beykoz	248260	310.36	800
11	Beylikdüzü	352412	37.78	9328
12	Beyoğlu	233323	8.91	26187
13	Büyükdere	254103	139.17	1826
14	Çatalca	73718	1115.13	66
15	Çekmeköy	264508	148.09	1786

**AFTER LITTLE
MANIPULATION, THE DATA-
FRAME IS OBTAINED**

Data Preparation

3. Getting Coordinates of Districts using Geopy Client

```
# module geocoder class from geopy client to convert an address each neighborhood into latitude and longitude values  
from geopy.geocoders import Nominatim  
  
# In order to define an instance of the geocoder, we need to define a user_agent.  
geolocator = Nominatim(user_agent="Istanbul_explorer")  
  
df['Latitude'] = df['District'].apply(geolocator.geocode).apply(lambda x: (x.latitude))  
df['Longitude'] = df['District'].apply(geolocator.geocode).apply(lambda x: (x.longitude))
```



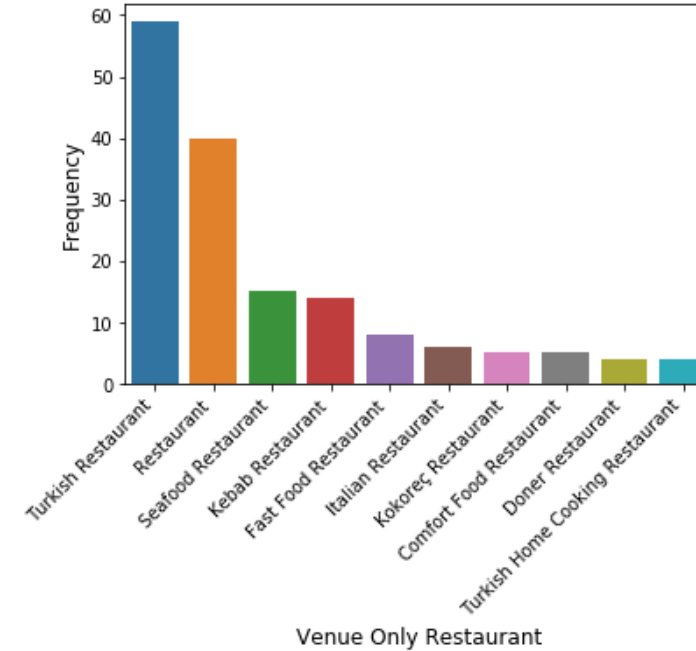
**USED
PYTHON FOLIUM
LIBRARY TO
VISUALIZE
GEOGRAPHIC
DETAILS OF
ISTANBUL AND ITS
39 DISTRICTS**

Exploratory Data Analysis

1. Using Foursquare Location Data

- will concentrate in Restaurant Category only and explore all the 39 districts in Istanbul.
- We find out 18 unique venue categories that only restaurant in Istanbul and Turkish Restaurants top the charts as we can see in the plot below

10 Most Frequently Occuring Venues Only Restaurant in 39 Districts of Istanbul



2. Analyze each neighborhood to know about the top 5 venues of each one

◀ ▶

[illegible]

Create a data-frame with pandas one hot encoding for the venue categories that only contain restaurant

```
# dataframe istanbul_venues_only_restaurant contains column Neighborhood, Neighborhood Latitude, Neighborhood Longitude, Venue Category, and Venue Name
# one hot encoding
istanbul_onehot = pd.get_dummies(istanbul_venues_only_restaurant[['Venue Category']], prefix="", prefix_sep="")

# add neighborhood column back to dataframe
istanbul_onehot['Neighborhood'] = istanbul_venues_only_restaurant['Neighborhood']

# move neighborhood column to the first column
fixed_columns = [istanbul_onehot.columns[-1]] + list(istanbul_onehot.columns[:-1])
istanbul_onehot = istanbul_onehot[fixed_columns]

istanbul_onehot.head()
```

[illegible]

Use pandas groupby on neighborhood column and calculate the mean of the frequency of occurrence of each venue category

```
# reset_index() use for reset index each row after grouping
istanbul_grouped = istanbul_onehot.groupby('Neighborhood').mean().reset_index()
istanbul_grouped
```

	Neighborhood	Chinese Restaurant	Comfort Food Restaurant	Doner Restaurant	Eastern European Restaurant	Falafel Restaurant	Fast Food Restaurant	Italian Restaurant	Kebab Restaurant	Kokoreç Restaurant	Kumpir Restaurant	Mediterranean Restaurant
0	Arnavutköy	0.000	0.00	0.00	0.0	0.000000	0.166667	0.000000	0.000000	0.000000	0.000000	
1	Ataşehir	0.000	0.00	0.25	0.0	0.000000	0.000000	0.250000	0.125000	0.000000	0.000000	
2	Avcılar	0.000	0.00	0.00	0.0	0.000000	0.000000	0.000000	0.000000	1.000000	0.000000	
3	Bahçelievler	0.000	0.00	0.00	0.0	0.000000	0.000000	0.000000	0.000000	0.250000	0.000000	
4	Bakırköy	0.000	0.00	0.00	0.0	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	
5	Bayrampaşa	0.000	0.20	0.00	0.0	0.000000	0.000000	0.000000	0.000000	0.200000	0.000000	

Output each neighborhood along with the top 5 most common venues categories that only restaurant in Istanbul

```
num_top_venues = 5

for hood in istanbul_grouped['Neighborhood']:
    print("----"+hood+"----")
    temp = istanbul_grouped[istanbul_grouped['Neighborhood'] == hood].T.reset_index()
    temp.columns = ['venue', 'freq of occurrence']
    temp = temp.iloc[1:]
    temp['freq of occurrence'] = temp['freq of occurrence'].astype(float)
    temp = temp.round({'freq': 3})
    print(temp.sort_values('freq of occurrence', ascending=False).reset_index(drop=True).head(num_top_venues))
    print('\n')
```

----Arnavutköy----

	venue	freq of occurrence
0	Restaurant	0.500000
1	Turkish Restaurant	0.333333
2	Fast Food Restaurant	0.166667
3	Mediterranean Restaurant	0.000000
4	Turkish Home Cooking Restaurant	0.000000

----Ataşehir----

	venue	freq of occurrence
0	Restaurant	0.375
1	Doner Restaurant	0.250
2	Italian Restaurant	0.250
3	Kebab Restaurant	0.125
4	Chinese Restaurant	0.000

Use clustering method to clustering these 34 districts
based on the venue categories that only contain
restaurant

Cluster Neighborhoods

```
# import k-means from clustering stage
from sklearn.cluster import KMeans

# Run k-means to cluster the neighborhood into 5 clusters
# set number of clusters
kclusters = 5

istanbul_grouped_clustering = istanbul_grouped.drop('Neighborhood', 1)

# run k-means clustering
kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(istanbul_grouped_clustering)

# check cluster labels generated for each row in the dataframe
kmeans.labels_[0:10]

array([1, 3, 3, 3, 2, 2, 3, 1, 0, 2], dtype=int32)
```

Let's create a new dataframe that includes the cluster label as well as the top 10 venues for each neighborhood.

```
# add clustering labels to first column of dataframe neighborhoods_venues_sorted
neighborhoods_venues_sorted.insert(0, 'Cluster Labels', kmeans.labels_)

istanbul_merged = df

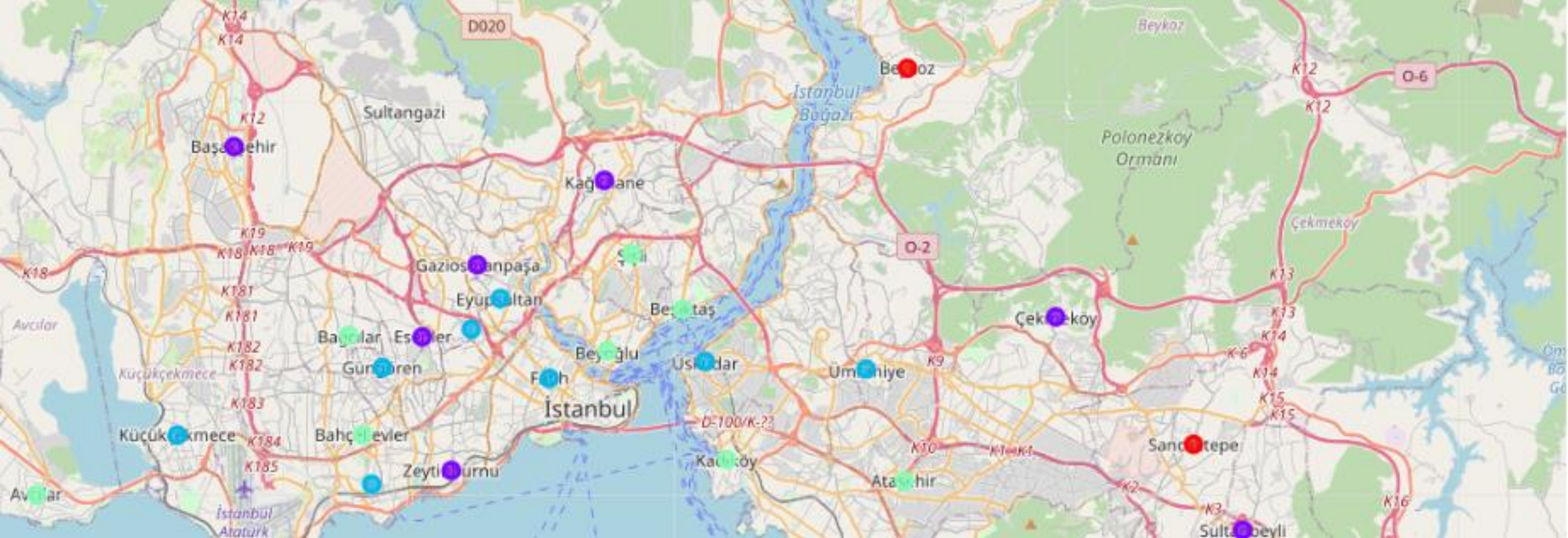
istanbul_merged.rename(columns={'District': 'Neighborhood'}, inplace=True)

# merge neighborhoods_venues_sorted with df to add latitude/longitude for each neighborhood
# istanbul_merged = istanbul_merged.join(neighborhoods_venues_sorted.set_index('Neighborhood'), on='Neighborhood')

# with merge function, istanbul_merged and neighborhoods_venues_sorted in this case `Neighborhood` is the only column name in both dataframes
merged_inner = pd.merge(left=istanbul_merged, right=neighborhoods_venues_sorted, left_on='Neighborhood', right_on='Neighborhood')

merged_inner.head() # check the last columns!
```

	Neighborhood	Population	Area (km ²)	Density (per km ²)	Latitude	Longitude	Cluster Labels	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
0	Arnavutköy	282488	450.35	627	41.184182	28.740729	1	Restaurant	Turkish Restaurant	Fast Food Restaurant	Kebab Restaurant	Comfort Food Restaurant	Doner Restaurant	Eastern European Restaurant	Falafe Restaurant
1	Ataşehir	425094	25.23	16849	40.984749	29.106720	3	Restaurant	Doner Restaurant	Italian Restaurant	Kebab Restaurant	Turkish Restaurant	Comfort Food Restaurant	Eastern European Restaurant	Falafe Restaurant



WE CAN VISUALIZE THE RESULTING CLUSTERS BY REPRESENT THESE 5 CLUSTERS USING FOLIUM LIBRARY AS BELOW

Results

- Turkish Restaurant is the most frequently occurring venues that only restaurant in the 39 Districts of Istanbul.
- The types of restaurants outside Turkey such as Thai Restaurant and Chinese Restaurant are the least frequently occurring venues that only restaurant in the 39 Districts of Istanbul.
- Neighborhood Güngören and Silivri have the highest number of restaurants/food venues in Istanbul.
- Neighborhood Avcılar, Bakırköy, Bağcılar, Beykoz, Kartal, and Tuzla has the least number of restaurants in Istanbul.
- Cluster 1, cluster 2, and cluster 3 which are marked with a circle maker in red, purple and blue in folium map are in the downtown area of Istanbul so they can choose from many types of restaurants, especially Turkish restaurants.
- Istanbul's strategic position has a variety of land lines and many rail lines that connect between neighborhoods making it easy for every tourist to travel, especially to find a place to eat.
- If tourists want to find a place to eat with a taste of Turkey outside such as European restaurants, they can travel in Clusters 4 and Clusters 5 a bit far from downtown Istanbul.

Discussion

The clustering is completely based on the most common venues obtained from Foursquare data, especially in this clustering is done on the most common venues data that contains the word restaurant in each neighborhood in Istanbul. In this clustering analysis, we make a number of assumptions, including food venue/restaurant data that appears in each neighborhood based on the closest distance from the center of the neighborhood not from the Atatürk Olympic Stadium, ignoring the price range of each restaurant, the cleanliness and hygiene of food from each restaurant, restaurant services and etc. Since we don't have such data and it would be difficult to farm it for a small exploratory study like ours. Hence, our analysis only helps tourists to get an overview of food venue/restaurants distribution by categories in the 39 districts of Istanbul.

Conclusion

- Many problems in real life where the data associated with these problems can be used to find solutions. For example the problem above we want to help every visitor from various countries to list and visualize Istanbul districts that fit their needs in terms of culinary/food venues. Existing data is used for segmenting and clustering each neighborhood in Istanbul based on the most common venue that contains a word restaurant. From this data will help each tourist determine which neighborhood has a restaurant that suits their interests.
- Data manipulation starts from scrapping tables from wikipedia, cleaning and manipulation of the data, getting coordinates of districts using geopy client, using foursquare location data to explore the neighborhoods to get venues that only restaurants and finally use clustering methods to clustering these 34 districts based on the venue categories that only contain restaurant and we can visualize the results using the folium leaflet map.
- From the results of clustering 34 neighborhoods in Istanbul we grouped them into five clusters. The first, second, and third clusters located in the center of Istanbul have many types of restaurants, especially the highest number of Turkish restaurants and also foreign restaurants such as Fast Food Restaurant and Italian Restaurant. These three clusters are also closer to the Atatürk Olympic Stadium. The fourth and fifth clusters are located a bit far from downtown Istanbul if tourists want to find a place to eat with a taste of Turkey outside such as European restaurants.