

The Battle of the Neighborhoods in São Paulo City, Brazil.

Final project of the Courseira IBM Data Science Course



1. Introduction: Business Problem

In this project we will try to find another ways to make a new marketing campaign by a segmentation method. This report will be targeted to people interested in suggest marketing contents, products advertising and etc. to the people that lives in the Districts of São Paulo's Capital, Brazil.

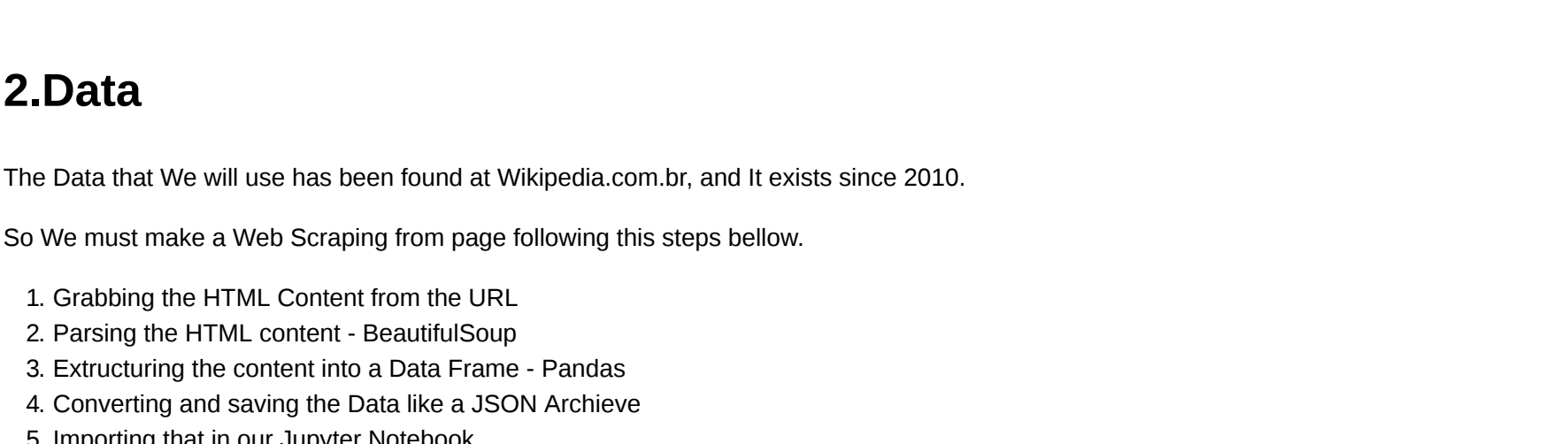
All entrepreneurs looks a way for increase the selling. Therefore, different marketing strategies are made of. We are talking about keeping a good presence in the social medias, to optimize the mechanisms of search, investing in customer service. All of those are so good actions and also bring great results, but We are able to make more.

We know that are lots of venues in the São Paulo City, therefore We will try figure out the business profile of each borough. Assuming that São Paulo is an alpha global city, We must to know what are the subject most reached by our costumers.

And finally, help them to make better buying decisions. We will use all power of data science to get a success we need.

Districts of São Paulo City

We are talking about 96 Districts



2.Data

The Data that We will use has been found at Wikipédia.com.br, and it exists since 2010.

So we must make a Web Scraping from page following this steps below.

- Grabbing the HTML Content from the URL.
- Parsing the HTML content - BeautifulSoup
- Extracting the content into a Data Frame - Pandas
- Converting and saving the Data like a JSON Archive
- Importing that in our Jupyter Notebook

How?

For that We will use some libs in Python language.

Like that:

```
import time
import requests
import pandas as pd
from bs4 import BeautifulSoup
from selenium import webdriver
from selenium.webdriver.firefox.options import Options
from selenium.webdriver.common.by import By
import json
```

URL: https://pt.wikipedia.org/wiki/Lista_de_distritos_de_S%C3%A3o_Paulo_por_nome

	Districts	Population
0	Grajaú	360.787
1	Jardim Ângela	295.434
2	Sapopemba	284.524
3	Capão Redondo	268.729
4	Jardim São Luís	267.871
...
92	Sã	23.651
93	Parí	17.299
94	Barra Funda	14.383
95	Marsilac	8.258
96	None	NaN

97 rows x 2 columns

3. Methodology

This survey is a Clustering problem. In this case We will use a machine learning method.

K-Means machine learning method will be used for acquiring information concerning the venues located close to each area the data platform Foursquare was used. In this project our purpose is to create an appropriate number of clusters based on the Geolocalization and the quantity and the kind of venues located in each District in São Paulo City.

In the end of this project we will be able to see the different types of public that we may consider important according to the kind of business We have. Also, We will figure out what kind of destinations each venue has an interest. In general we acquire a clean perception of the most visitable venues in each District. From that We will see what kind of costumers we may offer like recommendations to each client.

4. Preparing Data

4.1 Manipulation of the dataframe to discover the coordinates of the districts. For that we must make a tratament of our Dataset.

The first step is to eliminate the NaN data's. The second one consists on create a new column with the Address. (This is so important to get the coordinates of each Distric of São Paulo)

	Districts	Population	Address	Latitude	Longitude
0	Grajaú	360.787	Grajaú São Paulo Brasil	-23.785907	-46.869197
1	Jardim Ângela	295.434	Jardim Ângela São Paulo Brasil	-23.712528	-46.768720
2	Sapopemba	284.524	Sapopemba São Paulo Brasil	-23.694326	-46.590885
3	Capão Redondo	268.729	Capão Redondo São Paulo Brasil	-23.671903	-46.779435
4	Jardim São Luís	267.871	Jardim São Luís São Paulo Brasil	-23.683573	-46.737752
...
91	Jaguara	24.895	Jaguara São Paulo Brasil	-23.507448	-46.755315
92	Sã	23.651	Sã São Paulo Brasil	-23.550051	-46.633362
93	Parí	17.299	Parí São Paulo Brasil	-23.522979	-46.616486
94	Barra Funda	14.383	Barra Funda São Paulo Brasil	-23.595462	-46.695913
95	Marsilac	8.258	Marsilac São Paulo Brasil	-23.637142	-46.710236

96 rows x 5 columns

The Coordinates had been gotten by The Library GeoPy, that fetchs the Latitude and Longitude from the Address, how We said lastly.

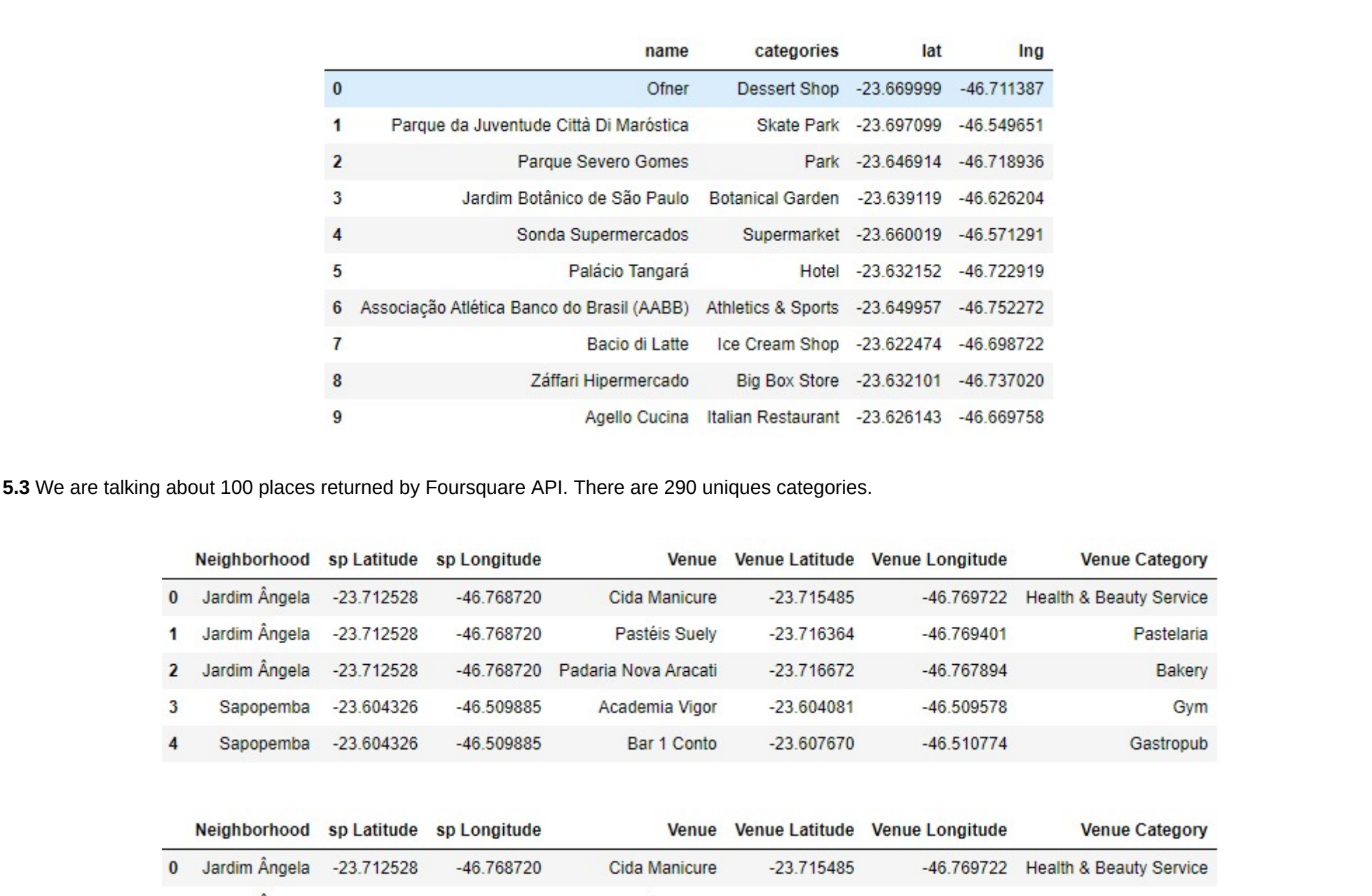
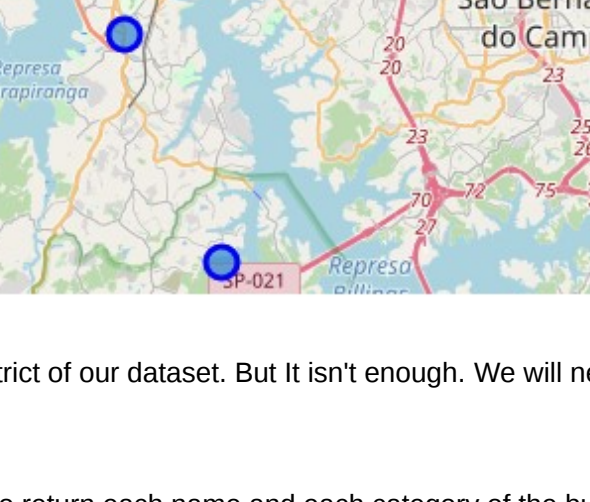
from geopy.geocoders import Nominatim

Districts	Population	Address	Latitude	Longitude
0	Grajaú	Grajaú São Paulo Brasil	-23.785907	-46.869197
1	Jardim Ângela	Jardim Ângela São Paulo Brasil	-23.712528	-46.768720
2	Sapopemba	Sapopemba São Paulo Brasil	-23.694326	-46.590885
3	Capão Redondo	Capão Redondo São Paulo Brasil	-23.671903	-46.779435
4	Jardim São Luís	Jardim São Luís São Paulo Brasil	-23.683573	-46.737752
...
91	Jaguara	Jaguara São Paulo Brasil	-23.507448	-46.755315
92	Sã	Sã São Paulo Brasil	-23.550051	-46.633362
93	Parí	Parí São Paulo Brasil	-23.522979	-46.616486
94	Barra Funda	Barra Funda São Paulo Brasil	-23.595462	-46.695913
95	Marsilac	Marsilac São Paulo Brasil	-23.637142	-46.710236

96 rows x 5 columns

5. Visualizing Data

5.1 Now We will utilize the Foursquare API to figure out the Latitude and Longitude. The next step is to plot the Maps through the Folium Library



How We can see, each blue point represents each District of our dataset. But it isn't enough. We will need to figure out all kind of Venues we can find from each reference points.

5.2 It is awesome thinking the Foursquare API is able to return each name and each category of the business places.

	name	categories	lat	lng
0	Parque da Juventude CIBD Di Maristota	Dessert Shop	-23.697099	-46.711387
1	Parque Severo Gomes	Park	-23.648914	-46.718936
2	Jardim Botânico de São Paulo	Botanical Garden	-23.638119	-46.626204
3	Sonda Supermercados	Supermarket	-23.660019	-46.571291
4	Palácio Tangará	Hotel	-23.632152	-46.722919
5	Associação Atlético Banco do Brasil (AABB)	Athletics & Sports	-23.649697	-46.752272
6	Baca & Leite	Ice Cream Shop	-23.652474	-46.688722
7	Zaffari HomeMERCADO	Big Box Store	-23.632101	-46.737020
8	Agello Cucina	Italian Restaurant	-23.626143	-46.669758

5.3 We are talking about 100 places returned by Foursquare API. There are 290 unqiues categories.

Neighborhood	sp Latitude	sp Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category	
0	Jardim Ângela	-23.712528	-46.768720	Cida Mancure	-23.715485	-46.769722	Health & Beauty Service
1	Jardim Ângela	-23.712528	-46.768720	Pastelaria Suely	-23.710364	-46.768401	Pasteleria
2	Jardim Ângela	-23.712528	-46.768720	Pastaria Nova Aracati	-23.716672	-46.767694	Bakery
3	Sapopemba	-23.694326	-46.590885	Academia Vitor	-23.694081	-46.590578	Gym
4	Sapopemba	-23.694326	-46.590885	Bar 1 Conto	-23.697570	-46.510774	Gastropub

We figure out the Best K number for a Clustering Algorithm. In this case, We are using K-Means algorithm and the number of clusters may use k for.

We used the Elbow Method. This Approach is very common used with the K-means Clustering. So we have got that K is equal to 5.

Number of Clusters	Distortion Score (Approximate)
1	0.550
2	0.535
3	0.525
4	0.515
5	0.505
6	0.500
7	0.495
8	0.490
9	0.485
10	0.480

5.4 How we know, to work with Machine Learning, it is necessary to convert the categorical variables into numerical variables.

For that we have used The One Hot Encoding method.

Neighborhood	Acad House	Accessories Store	African Restaurant	Antique Shop	Arcade	Argentinian Restaurant	Art Gallery	Art Museum	Art Store	Vietnamese Restaurant	Volleyball Court	Warehouse Store
0	Jardim Ângela	0	0	0	0	0	0	0	0	0	0	0
1	Jardim Ângela	0	0	0	0	0	0	0	0	0	0	0
2	Jardim Ângela	0	0	0	0	0	0	0	0	0	0	0
3	Sapopemba	0	0	0	0	0	0	0	0	0	0	0
4	Sapopemba	0	0	0	0	0	0	0	0	0	0	0
...
95	Itaim Paulista	0	0	0	0	0	0	0	0	0	0	0
96	Itaim Paulista	0	0	0	0	0	0	0	0	0	0	0
97	Jabaquara	0	0	0	0	0	0	0	0	0	0	0
98	Jabaquara	0	0	0	0	0	0	0	0	0	0	0
99	Jabaquara	0	0	0	0	0	0	0	0	0	0	0

100 rows x 291 columns

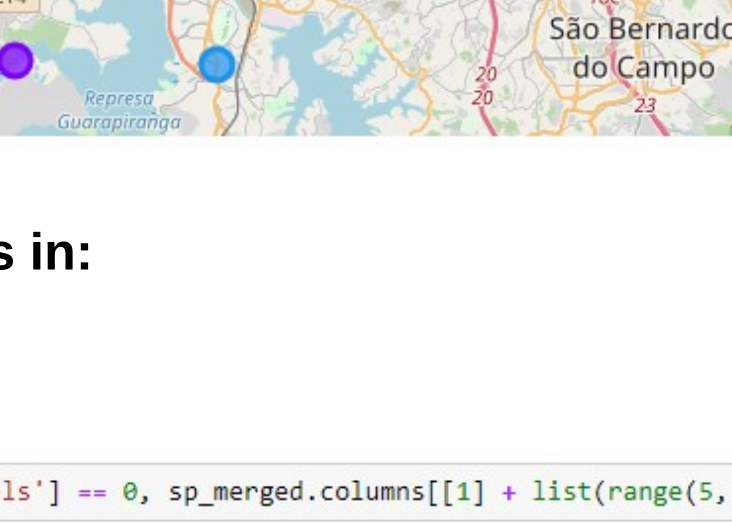
5.5 In case of this business problem We want to discover the most common Venue of São Paulo's District. All of this driven by location.

Neighborhood	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 Most Common Venue
0	Alto de Pinheiros	Plaza	Dog Run	Bike Rental / Bike Share	Cafe	Gym / Fitness Center	Trail	Spa	Department Store	Fish & Chips Restaurant
1	Arhangara	Gym / Fitness Center	Pet Store	Ice Cream Shop	Pizza Place	Convenience Store	Grocery Store	Pet Store	Event Space	Escape Room
2	Aracandara	Bakery	Clothing Store	Gym / Fitness Center	Candy Store	Grocery Store	Yoga Studio	Flea Market	Farmers Market	Fish & Chips Restaurant
3	Altus Azim	Pizza Place	Department Store	Nearstand	Grocery Store	Beer Garden	Bar	Bakery	Gymnastics Gym	Pasteleria
4	Barra Funda	Music Venue	Restaurant	Standart Restaurant	Cafe	Sandwich Place	Chocolate Shop	Courty	Dance Club	Nightclub

6. Modeling

6.1 It's interesting We figure out the Best K number for a Clustering Algorithm. In this case, We are using K-Means algorithm and the number of K represents the best quantities of Clusters may use for.

So, We might have used the Elbow Method. This Approach is very common used with the K-means Clustering. So We have got that K is equal 6.



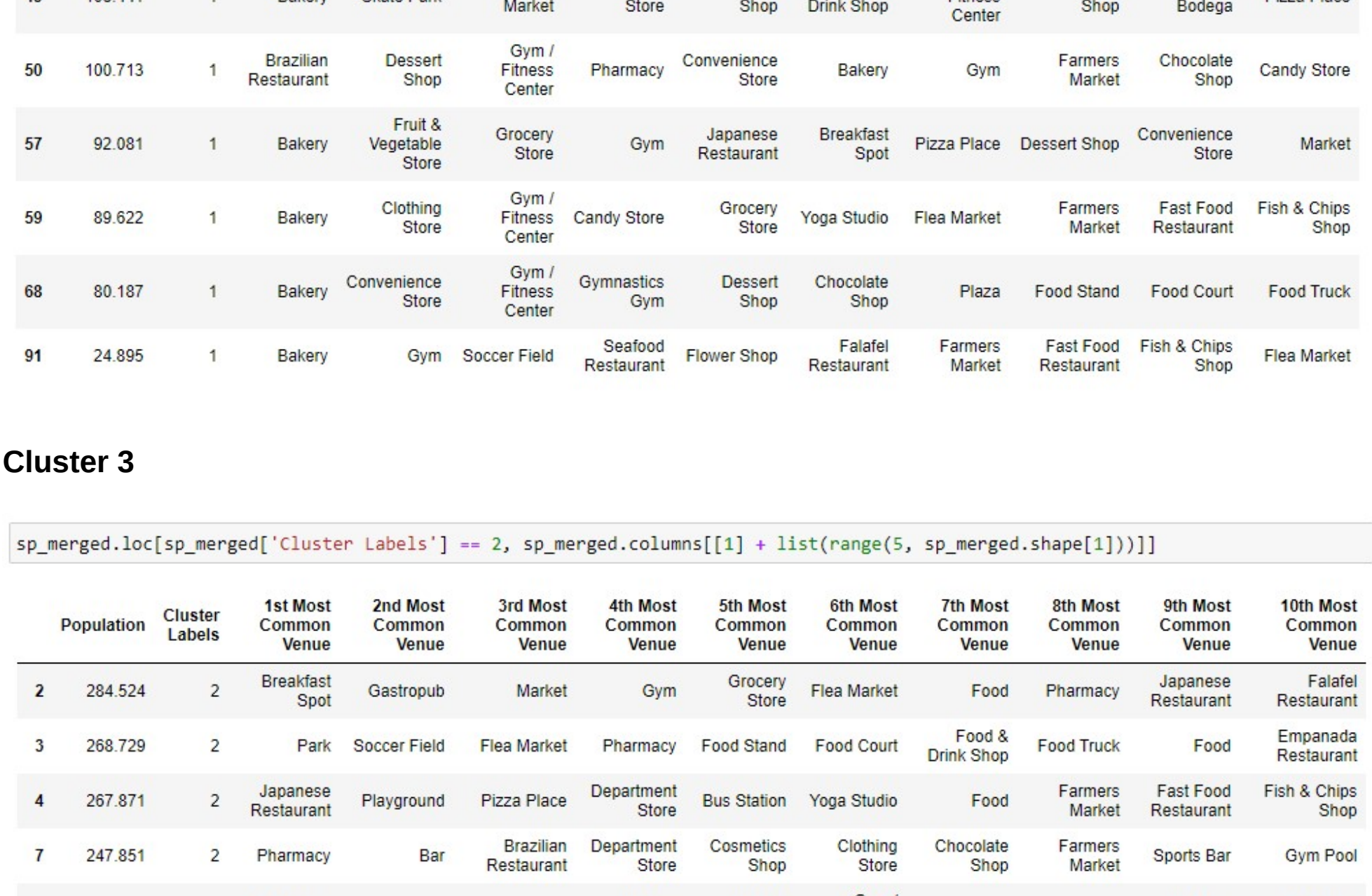
```
In [ ]: # set number of clusters
kclusters = 6

# run k-means clustering
kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(sp_grouped_clustering)
# check Cluster Labels generated for each row in the dataframe
kmeans.labels_[0:10]
```

Districts	Population	Address	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	9th Most Common Venue	10th Most Common Venue
0	Grajaú	Grajaú São Paulo Brasil	-23.785907	-46.869197	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
1	Jardim Ângela	Jardim Ângela São Paulo Brasil	-23.712528	-46.768720	1.0	Bakery	Health & Beauty Service	Pasteleria	Flea Market	Escape Room	Event Space	Faifal Restaurant	Farmers Market		
2	Sapopemba	Sapopemba São Paulo Brasil	-23.694326	-46.590885	2.0	Breakfast Spot	Gastropub	Market	Gym	Grocery Store	Flea Market	Food	Pharmacy		
3	Capão Redondo	Capão Redondo São Paulo Brasil	-23.671903	-46.779435	2.0	Park	Soccer Field	Flea Market	Pharmacy	Food Court	Food Court	Food & Drink Shop	Food & Drink Shop		
4	Jardim São Luís	Jardim São Luís São Paulo Brasil	-23.683573	-46.737752	2.0	Japanese Restaurant	Playground	Pizza Place	Department Store	Bus Station	Yoga Studio	Food	Farmers Market		

6.2 After merging the data, We had needed to make a new treatment of the dataset.

6.3 Now We have a new map with the clusters separated by colors. It's simply wonderful! The 6 clusterings separated by Colors!



The Clustering consists in:

Cluster 1

Population	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	9th Most Common Venue	10th Most Common Venue
34	127.015	0	Gym / Fitness Center	Pet Store	Market	Gym	Residential Building (Apartment / Condo)	Event Space	Faifal Restaurant	Farmers Market	Fast Food Restaurant
73	65.859	0	Gym / Fitness Center	Plaza	Ice Cream Shop	Bakery	Pizza Place	Convenience Store	Grocery Store	Pet Store	Fish & Chips Shop

Cluster 2

Population	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	9th Most Common Venue	10th Most Common Venue
1	295.434	1	Bakery	Health & Beauty Service	Pasteleria	Flea Market	Escape Room	Event Space	Faifal Restaurant	Farmers Market	Fast Food Restaurant
5	265.681	1	Music Venue	Fried Chicken Joint	Bakery	Gymnastics Gym	Flower Shop	Faifal Restaurant	Farmers Market	Fast Food Restaurant	Fish & Chips Shop
15	184.616	1	Bakery	Convenience Store	Grocery Store	Gym / Fitness Center	Farmers Market	Arts & Crafts Store	Brazilian Restaurant	Italian Restaurant	Diner
26	136.623	1	Bakery	Chocolate Shop	Gym / Fitness Center	Cafe	Pharmacy	Food Stand	Food Court	Food & Drink Shop	Food
21	131.823	1	Brazilian Restaurant	Nature Preserve	Snack Place	Cafe	French Restaurant	Food Stand	Fried Chicken Joint	Food Stand	Food Truck
42	109.088	1	Brazilian Restaurant	Bakery	Market	Food & Drink Shop	Farmers Market	Gym	Gym / Fitness Center	Furniture / Home Store	Pizza Place
43	108.441	1	Bakery	Shate Park	Farmers Market	Grocery Store	Gourmet Shop	Food & Drink Shop	Gym / Fitness Center	Motorcycle Shop	Pizza Place
50	100.713	1	Brazilian Restaurant	Dessert Shop	Gym / Fitness Center	Pharmacy	Convenience Store	Bakery	Gym	Farmers Market	Chocolate Shop
57	92.081	1	Bakery	Fruit & Vegetable Store	Grocery Store	Gym	Japanese Restaurant	Breakfast Spot	Pizza Place	Dessert Shop	Convenience Store
59	89.622	1	Bakery	Clothing Store	Gym / Fitness Center	Candy Store	Grocery Store	Yoga Studio	Flea Market	Farmers Market	Fast Food Restaurant
61	86.187	1	Bakery	Convenience Store	Gym / Fitness Center	Gymnastics Gym	Dessert Shop	Chocolate Shop	Farmers Market	Fast Food Restaurant	Food Court
98	24.895	1	Bakery	Gym	Soccer Field	Seafood Restaurant	Flower Shop	Faifal Restaurant	Pharmacy	Fast Food Restaurant	Fish & Chips Shop

Cluster 3

Population	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	9th Most Common Venue	10th Most Common Venue
2	284.524	2	Breakfast Spot	Gastropub	Market	Gym	Convenience Store	Food Court	Food & Drink Shop	Food	Food Truck
3	268.729	2	Park	Soccer Field	Flea Market	Pharmacy	Food Stand	Food Court	Food & Drink Shop	Food	Food Truck
4	267.871	2	Japanese Restaurant	Playground	Pizza Place	Department Store	Bus Station	Yoga Studio	Food	Farmers Market	Fast Food Restaurant
7	247.851	2	Pharmacy	Bar	Brazilian Restaurant	Department Store	Cosmetics Shop	Clothing Store	Chocolate Shop	Farmers Market	Sports Bar
8	224.074	2	Japanese Restaurant	Bakery	Food Truck	Boutique Alley	Dessert Shop	Gym / Fitness Center	Pizza Place	Brewery	Brazilian Restaurant
...
89	33.882	2	Korean Restaurant	Brazilian Restaurant	Women's Store	Dessert Shop	Cafe	Bar	Coffee Shop	Grocery Store	Bakery
90	29.265	2	Brazilian Restaurant	Clothing Store	Pizza Place	Park	Burger Joint	Furniture / Home Store	Gaming Cafe	Cafe	Shoe Store
92	23.651	2	Brazilian Restaurant	Bookstore	Miscellaneous Shop	Cosmetics Shop	Bakery	Music Venue	Chocolate Shop	Japanese Restaurant	Arts & Crafts Store
93	17.299	2	Clothing Store	Brazilian Restaurant	Middle Eastern Restaurant	Restaurant	Shopping Mall	Cafe	Faifal Restaurant	Winehouse	Bar
94	14.383	2	Music Venue	Restaurant	Brazilian Restaurant	Cafe	Sandwich Place	Chocolate Shop	Country Dance Club	Nightclub	Office

Cluster 4

Population	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	9th Most Common Venue	10th Most Common Venue
21	143.932	3	Bar	Park	Brewery	Food	Farmers Market	Fast Food Restaurant	Fish & Chips Shop	Flea Market	Flower Shop
63	84.843	3	Bar	Chinese Restaurant	Gym / Fitness Center	Burger Joint	Diner	Food & Drink Shop	Flea Market	Fast Food Restaurant	Fish & Chips Shop
72	68.258	3	Brazilian Restaurant	Planetarium	IT Services	Farmers Market					