

IBM Data Science Professional Certificate

IBM Capstone Project: Applied Data Science Capstone

Final report

1. Introduction/Business Problem section

This project aims to utilize data science concepts and machine learning tools to solve a popular problem for restaurant owners in Madrid, Spain: where is the best neighborhood to open a tapas restaurant? In this project, I will go through the processes of problem definition, data preparation, and use machine learning to improve business decisions.

Madrid is the capital and most-populous city of Spain. The Madrid urban agglomeration has the third-largest GDP in the European Union and its influence in politics, education, entertainment, environment, media, fashion, science, culture, and the arts all contribute to its status as one of the world's major global cities.

Madrid is reputed to have a "vibrant nightlife". It is one of the city's main attractions with tapas bars, cocktail bars, clubs, jazz lounges, live music venues and flamenco theatres. That makes Madrid a great place for entrepreneurs to start and grow their business.

1.2. Target Audience:

Customer is a company, that want to start a tapas restaurant. We should help to make the best choice, to recommend the correct borough and explain why this place is better than others.

2. Description of the Data:

For solving the problem, we need different type of data. Which are:

- List of neighborhoods of Madrid
- Car Parking list and their geolocations of Madrid.

- Light Rail/Metro/Train/Tram/Bus Station list and their geolocations of Madrid

For the first point we use Google Geoencoding API to find the approximate coordinates of the neighborhoods. And for the second and third datasets we get the datasets from the Foursquare API.

The important disadvantage of API is the response limit of queries. In every query, the response of venues list limited as 100 result. So, we tried to send query with sub-categoryIDs of venues' (which is exist in Foursquare web site) and update the dataframe with "for loop" easily. we explained this detail later also.

3. Methodology section

3.1. Business Understanding:

Our main goal is to find optimum borough for a new restaurant.

3.2. Analytic approach:

As we know, choosing a neighborhood and location for a restaurant is an important and difficult step for the entrepreneurs. Most of time it can be as crucial factor as menu quality and service for the restaurant's success. The neighborhood is not only an important factor for customers but also it affects the restaurant quality, including the menu, interior design, style of furniture.

Then, to choose the location for the restaurant, we divide the project into two parts:

1. we take into account the suggestions of the article 4 Important Factors When Choosing a Location to Open a Restaurant which are
 - Parking. A restaurant should have its own parking lot. If that isn't an option (for example, in a major cities), consider partnering with a hotel in the area that has its own parking options. Also, it has an acceptable distance to public parking lots.
 - Accessibility. There's a reason that major restaurant chains are often located near highway exits: It makes them accessible for customers. Also, it is a good option to close distance to public transportations. Light Rail Station, Metro Station, Train/Tram Station, Bus Stop around the restaurant is provide convenience for customers.
 - Visibility. This goes along with accessibility and is very important for new restaurant locations. People have to know the restaurant is there, either in person or on their mobile devices.
 - Population Base. There need to be enough people who live in or pass through the area regularly to keep the restaurant busy. So "the particular area's population base" is another important factor about suitable place for a new restaurant. Due to Madrid has a very crowded city with huge number of tourists,

the Population Base is negligible. The Visibility is out of the scope of this project. It is hard to keep in view, because it needs different type of data like usage density of every streets, frequently preferred routes and the preferred reason, etc. May be a field research is also needed. Anyway, because of the difficulties, we focused on Parking and Accessibility and collected data about them.

2. From the first part we take the best thirty neighborhoods with better parking and accessibility factor. Then, we will use Foursquare API to find the top 100 venues within a radius of 1000 meters of the center of that neighborhoods. Finally, we will run k-means clustering on the data to cluster neighborhoods.

3.3. Exploratory Data Analysis:

1. Download csv-file from <https://datos.madrid.es/egob/catalogo/200078-1-distritos-barrios.csv> and Convert csv-file to DataFrame by pandas.
2. Add the latitude and longitude coordinates from google map API

	Neighborhood	Lat	Lon
0	PALACIO , Madrid, Spain, Madrid, S...	40.4163	-3.71432
1	IMPERIAL , Madrid, Spain, Madrid, S...	40.4068	-3.72009
2	PACIFICO , Madrid, Spain, Madrid, S...	40.4041	-3.67794
3	RECOLETOS , Madrid, Spain, Madrid, S...	40.4214	-3.68956
4	EL VISO , Madrid, Spain, Madrid, S...	40.449	-3.68668
...
126	CANILLEJAS , Madrid, Spain, Madrid, S...	40.4438	-3.61466
127	EL GOLOSO , Madrid, Spain, Madrid, S...	40.5562	-3.70402
128	ATALAYA , Madrid, Spain, Madrid, S...	40.4659	-3.66658
129	EL SALVADOR , Madrid, Spain, Madrid, S...	40.4444	-3.63065
130	COSTILLARES , Madrid, Spain, Madrid, S...	40.4737	-3.67183

131 rows × 3 columns

Fig. 1. A part of clear borough list from The London Datastore.

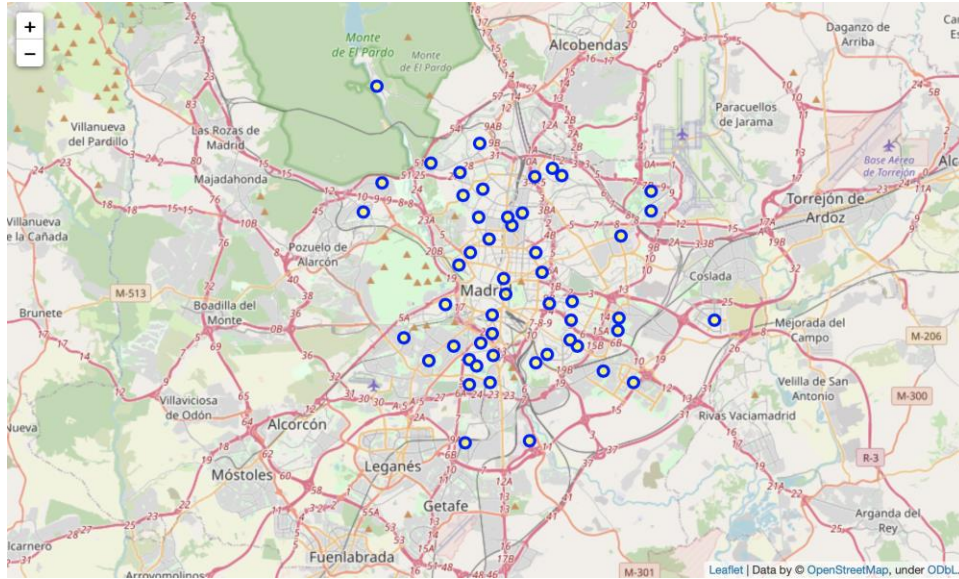


Fig. 2. Madrid Neighborhoods.

3.4 Create DataSets For Comparison

We need 2 datasets for comparison between the neighborhoods

- Car Parking list and their geolocations of Madrid
- Light Rail/Metro/Train/Tram/Bus Station list and their geolocations of Madrid

We can get these data via Foursquare Api. Firstly, we establish the Foursquare API.

	Neighborhood	Lat	Lon	CarPark_Dist	Close_Bus/Train
0	PALACIO , Madrid, Spain, Madrid, S...	40.416339	-3.714318	1008.48	13
1	IMPERIAL , Madrid, Spain, Madrid, S...	40.406845	-3.720088	1483.29	0
2	RECOLETOS , Madrid, Spain, Madrid, S...	40.421446	-3.689557	1752.70	10
3	PACIFICO , Madrid, Spain, Madrid, S...	40.404145	-3.677941	1902.74	0
4	EL VISO , Madrid, Spain, Madrid, S...	40.449021	-3.686681	2141.15	0

Fig. 3. Data set with the parking lot distance and with the number of public transport stops.

3.5 Comparison of Neighborhoods

In this step we can start to comparison of neighborhoods to determine the best suitable location for a new restaurant in Madrid. We sort the neighborhoods from the average distance of parking lots and create a new dataframe as df with first 30 neighborhoods.

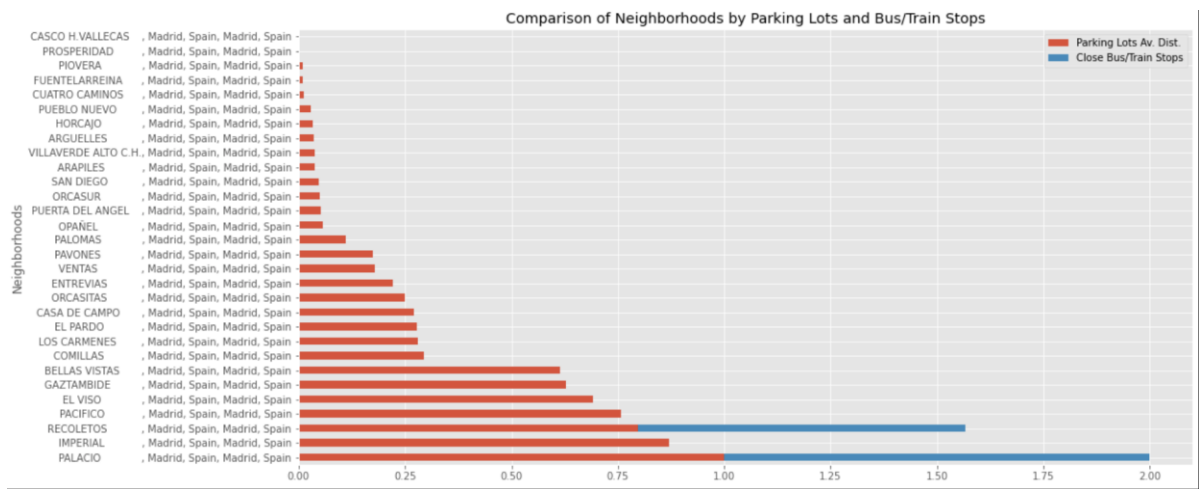


Fig. 4. Comparison among the first 30 neighborhoods.

We find that the top three neighborhoods are: 1 PALACIO, 2 RECOLETOS y 3 IMPERIAL.

We show a map with all the best neighborhood

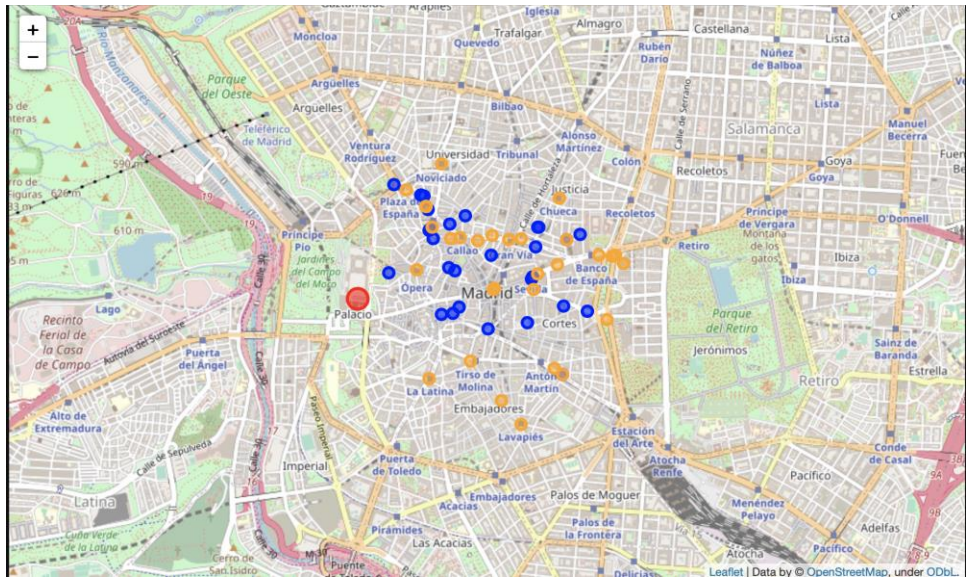


Fig. 5. Palacio, best neighborhood.

The result (best neighborhood) as RED circle marker, Parking Lots as BLUE circle markers and Metro/Train/Tram/Bus Stations as ORANGE circle markers.

According with this preliminary analysis, the best neighborhood to open a restaurant is **Palacio**. Now, we will explore the top 100 venues that are in the neighborhoods within a radius of 1500 meters from the first thirty neighborhood with better parking lot and bus/train stations.

3.6 Explore Restaurants in Madrid

Before we get into clustering, I'm creating a new dataframe with all restaurants data that was returned by Foursquare API. Since some of these venues were double counted, I will drop them in order to make a map of these restaurants.

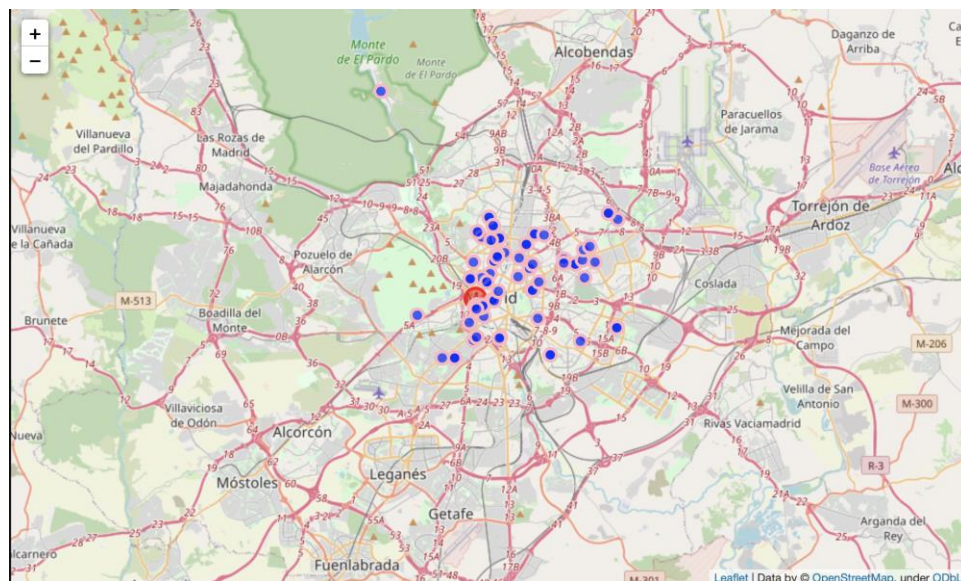


Fig. 6. Tapas restaurant for each borough in blue point and red the Palace neighborhood is shown.

4.Results

4.1 Machine Learning - Cluster Neighborhoods

We are ready to get into machine learning! For this project, we using k-means clustering. To begin with, a cluster is a collection of data points aggregated together based on their similarities. Using machine learning algorithms, we can cluster the neighborhoods based on their similarities with each other. K-means algorithm, in particular, first identifies k number of centroids, and then allocates every data point to the cluster, in a way that the data point is closer to that cluster's centroid than any other centroid. K-means algorithm runs this in a repetitive fashion before the centroids are stabilized and the clusters are formed. I am using this method because it is an

unsupervised learning method meaning that the algorithm will find the similarities between the data points for us given, we don't know them to begin with.

4.1.1 Find best K

One limitation of k-means clustering is that the algorithm does not decide how many clusters to form on its own and we need to find the best K to make clustering more accurate. The Elbow Method is one of the most popular methods to determine this optimal value of k. We iterate the values of k from 1 to 10 and calculate the distortion and inertia values for each value of k in the given range. Distortion is the average of the squared distances from the cluster centers of the respective clusters while inertia is the sum of squared distances of samples to their closest cluster center.

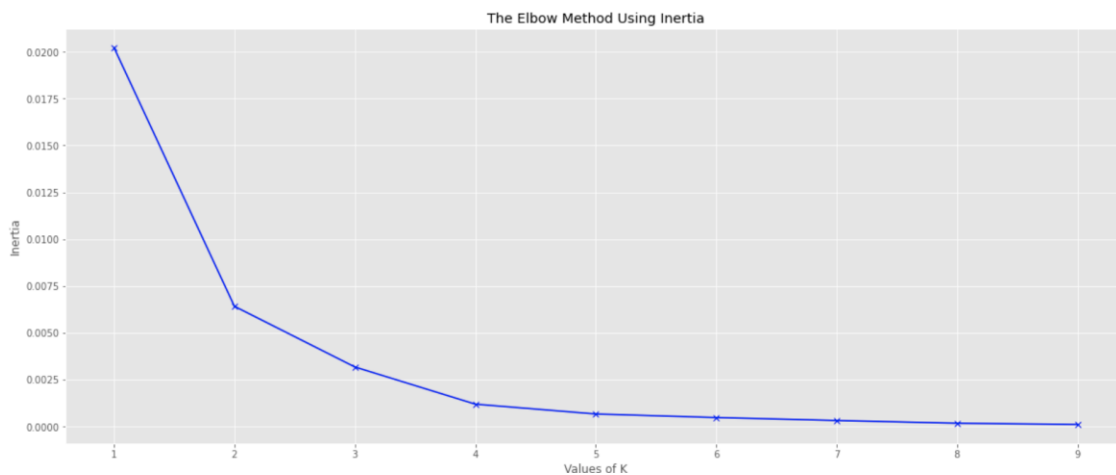


Fig. 7 The elbow method using inertia.

To determine the optimal number of K, we select the value of k at the “elbow” of the plots, the point after which the distortion/inertia starts decreasing in a linear fashion. Given these plots, we conclude that 4-6 clusters would work best for our data. Ultimately, I decided to go with 5.

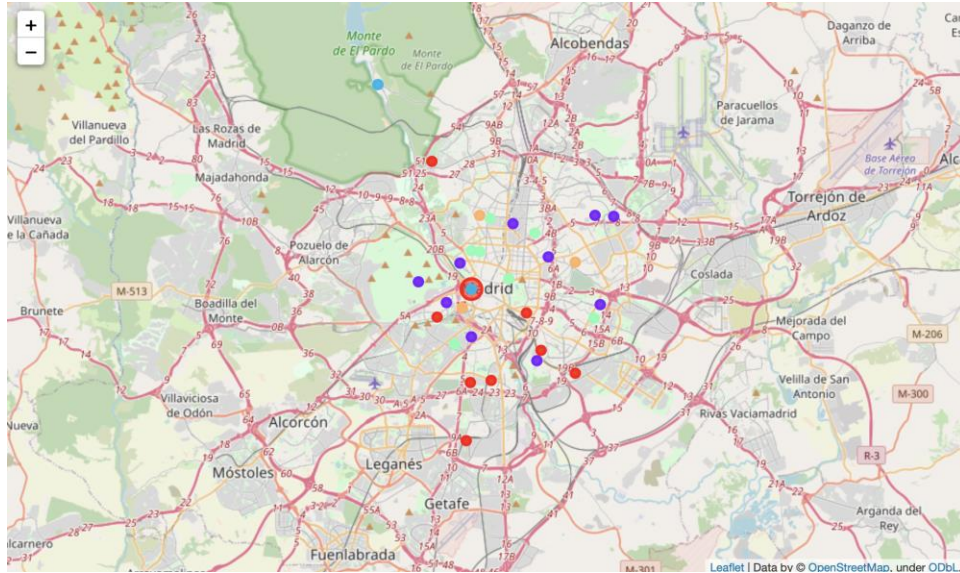


Fig. 8. Five clusters

The cluster are:

- Cluster 0 = Orange. It has an average frequency of Tapas Restaurants of 0.06666666666666667.
- Cluster 1 = Red. It has an average frequency of Tapas Restaurants of 0.0027380952380952383.
- Cluster 2 = Purple. It has average frequency of Tapas Restaurants of 0.027064030267977636.
- Cluster 3 = Blue. It has an average frequency of Tapas Restaurants of 0.09166666666666667.
- Cluster 4 = Green. It has an average frequency of Tapas Restaurants of 0.04817986611090059.
- Palacio Neighborhood = big red point

From Fig. 8 we can see that the Palacio Neighborhood is in the cluster 3 which has the higher average frequency of Tapas Restaurants. These two are good indicator that the best place to open a tapas restaurant in the Palacio's Neighborhood.

5. Discussion and Conclusions

In this project, we want to determine the best neighborhood in Madrid for the entrepreneurs to open a tapas restaurant. For this aim, we get the neighborhoods list of Madrid and then we use the Google Geoencoding API to find the approximate coordinates of the neighborhoods.

For comparison between the neighborhoods, we need "Parking Lots" list and "Metro/Train/Tram/Bus Station list" with their geolocations. We get these data via Foursquare Api. In Foursquare web site, we learned that the categoryId of our venues' (Parking, Metro Station, Tram Station and Bus Station) categoryIds and use these Ids to get and create the datasets. As we know, Foursquare API has a response limit (the number of venues returned by Foursquare API has a maximum limit is 100) of queries. Because of there are more than 100 Bus/Train Stops in Madrid, we tried to send query with sub-categoryIDs of venues and update the dataframe with "for loop" with 1000 meters radius easily.

Then, we calculate the average distance of all "parking lot" for each neighborhood and calculate the total "Metro/Train/Tram/Bus Station number" below than 1000 meters for each neighborhood.

After compare the neighborhoods with these two factor, we reach that The best neighborhood for a new restaurant in Madrid is: PALACIO.

Palacio has:

- 13 public transport (Metro/Train/Tram/Bus) Stations in 1 km of distance.
- The average distance for 17 Parking Lots is around 1 km.

During our cluster analysis, we found that Cluster 1 has the lowest average frequency of tapas restaurants while Cluster 3 has the highest. However, Cluster 4 has the highest number of tapas restaurants, with a lower average frequency, which might be because Cluster 4 has a high number of neighborhoods compared to other clusters, and that there are other common venues in the neighborhood which makes the frequency of tapas restaurants lower. On the other hand, cluster 3 has only two neighborhoods, these are: Palacio and El Pardo.

Based on the previous result, we can conclude that the better place to open a new tapas restaurant is in the Palacio's neighborhood. That is because it has the closest public transport station as well as the closest parking lots.

6. Final comments

It is only an exploratory project and it is not claimed that the result of the project is 100% proper for trade and investing industry. The project framework needs more thinking and planning with systematically. And also, the project needs more datasets, field and market research.

we should also note some limitations to this analysis. To start with, the Foursquare API limit to return only the top 100 venues within the radius of 1500 meters. But the neighborhoods are very different in shapes and sizes. Some neighborhoods are much larger but less populated while others are more densely populated with a smaller area. Thus, the Foursquare API might not have been able to capture all the tapas restaurants in each neighborhood. However, we calculated the frequency of these restaurants within the 1500-meter radius, which could still reflect the average frequency of tapas restaurants within that neighborhood. Thus, before opening a restaurant, it might be better to do some research on that specific neighborhood, for example, on its commercial pricing, consumers, competitors, and take other factors into consideration. Here a list of some upgrades to do:

3. We need different type of data like usage density of every streets, frequently preferred routes and the preferred reason.
4. Field research is also important. For instance, a restaurant has 10 parking lots in too close distance but they haven't got enough capacity so most of time it has no free parking lot. So here 10 close parking lots couldn't add positive effect to make a decision.
5. Taxes are another an important factor in determining the best location. So, we need extra datasets.
6. More data about the neighborhood. For instances, bad smell in the area, construction in the street, closely a loud venue.

So, the project has insufficiencies.

