

# The Battle of Neighborhoods Gym in Madrid

**T**he target of this project is to assess an investor to decide where would be the best location to place a Gym in Madrid. Using Foursquare data we want to cluster the neighborhoods based on their similarity in on the “Outdoors and Recreation” category so we can understand which are the sport offers each cluster offers to its population.



*New gym in Madrid is about to open*

Based on this premise, we will be able to find the neighborhoods where a Gym / Fitness Center has the better chance to get more potential clients.

# Data

The information he information needed about districts, neighborhoods and population can be found in the website “Portal de datos abiertos del Ayuntamiento de Madrid”: <https://datos.madrid.es/portal/site/egob>

From this website we can download an excel file with several indicators about (“panel\_indicadores\_distritos\_barrios\_2019.xls”) the population classified by neighborhood from which the necessary data will be extracted, such as district and neighborhood codes and their names that will be used for it.

We will additionally use two .CSV files downloaded from the same webpage containing information about which neighborhood belongs to each borough both of them you can find them in the project repository.

The information about districts and neighborhoods will be used to obtain geographical coordinates (latitude and longitude) where these are located by using the geocoder library.

The Foursquare API will be used to collect information about the venues and possible competitors in the neighborhoods of Madrid.

With the venues obtained from Foursquare it will be possible to classify them by category and finally establish a clustering for neighborhood by means of KMeans algorithm.

# Methodology

These are the sequential steps necessary to identify the top ten common venues for each neighborhood in Madrid classified by category and located in clusters in a map.

- The first step is to obtain a dataframe with codes and descriptions about districts and neighborhoods of Madrid city. I used two CSV files described in the “Data” section:
- Next, it is necessary to clean the dataframe, define columns, drop duplicates and unnecessary information like “district surface”
- Then, call argcis from geocoder library to obtain the latitude and longitude coordinates for each neighborhood in the dataframe. This will be necessary to find the Madrid venues by means of the Foursquare API.

borough code	borough	neighborhood
0	1	CENTRO
1	1	EMBAJADORES
2	1	CORTES
3	1	JUSTICIA
4	1	UNIVERSIDAD
...	...	...
126	21	BARAJAS ALAMEDA DE OSUNA
127	21	BARAJAS AEROPUERTO
128	21	BARAJAS CASCO H.BARAJAS
129	21	BARAJAS TIMON
130	21	BARAJAS CORRALEJOS

borough code	borough	neighborhood	Latitude	Longitude
0	1	CENTRO	40.40958	-3.88009
1	1	EMBAJADORES	40.39152	-3.69289
2	1	CORTES	40.41607	-3.69893
3	1	JUSTICIA	40.42479	-3.69308
4	1	UNIVERSIDAD	40.42565	-3.70726
...	...	...	...	...

- Once we have identified all neighborhoods and its location coordinates, we can use the Foursquare API to find the venues existing in each neighborhood. We will search for venues within 1km from the neighborhood center. As we are looking for Gyms, we will filter the venues to obtain only those ones related to the “Outdoors and Recreation” category.

neighborhood	neighborhood	neighborhood	neighborhood	neighborhood	venue	venue	venue	venue	venue
neighborhood	neighborhood	neighborhood	neighborhood	neighborhood	venue	venue	venue	venue	venue
0	PALACIO	40.40958	-3.88009		BeOne	40.404435	-3.885560	Gym / Fitness Center	
1	PALACIO	40.40958	-3.88009	Skatepark de Boadilla	40.410489	-3.888021		Skate Park	
2	PALACIO	40.40958	-3.88009		Fitness19	40.409288	-3.880009	Gym / Fitness Center	
3	PALACIO	40.40958	-3.88009	Parque Sofía de Grecia II	40.409460	-3.882442		Plaza	
4	PALACIO	40.40958	-3.88009	Parque Santillana del Mar	40.407508	-3.882004		Park	
...	...	...	...	...	...	...	...	...	...
2930	CORRALEJOS	40.46540	-3.61164	Recinto Para Perros Parque JCI	40.462782	-3.608200		Dog Run	
2931	CORRALEJOS	40.46540	-3.61164	Piscina de Novotel	40.462584	-3.614951		Pool	
2932	CORRALEJOS	40.46540	-3.61164	Fitness Pullman Hotel	40.462311	-3.614936		Gym / Fitness Center	
2933	CORRALEJOS	40.46540	-3.61164	Lago del Parque Juan Carlos I	40.460311	-3.609649		Lake	
2934	CORRALEJOS	40.46540	-3.61164	Donut Juan Carlos I	40.459320	-3.606842		Sculpture Garden	

## Coursera Capstone Project

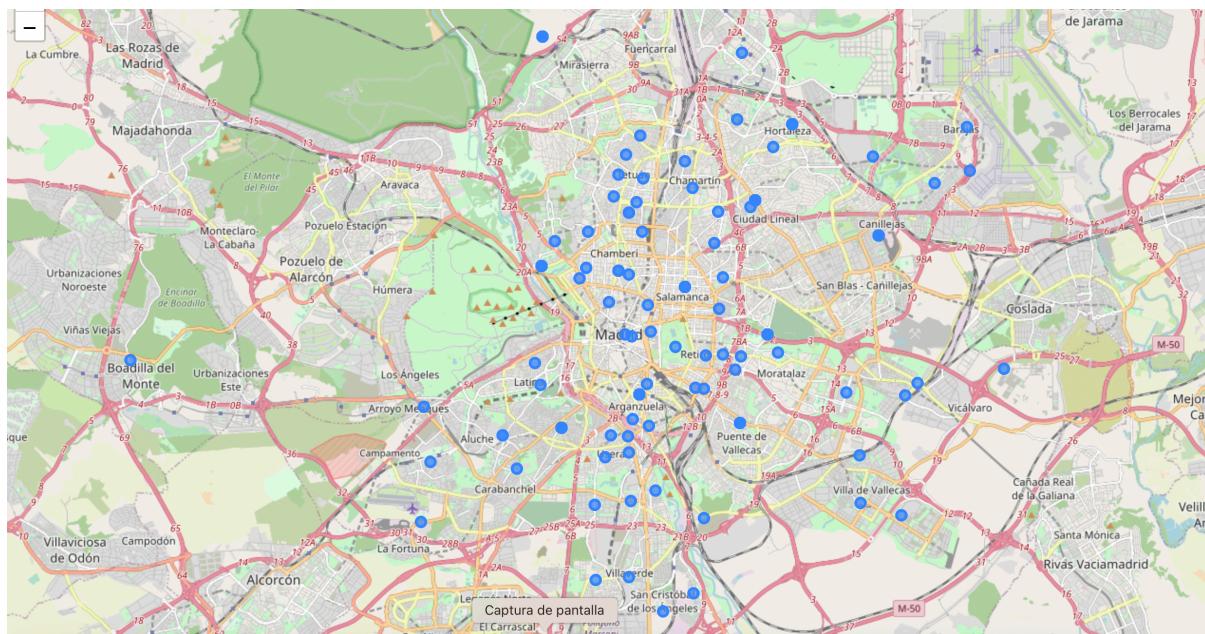
- With this new data frame we will group the venues to get the top 14 most common venues in each neighborhood based in the number of repetitions

	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue
0	ABRANTES	Plaza	Gym / Fitness Center	Athletics & Sports	Park	Playground	Roof Deck	Mountain	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym	Gun Range
1	ACACIAS	Plaza	Park	Gym	Gym / Fitness Center	Yoga Studio	Garden	Playground	Pool	Gym Pool	Gymnastics Gym	Lake	Outdoors & Recreation	Pedestrian Plaza	Athletics & Sports
2	ADELFA'S	Gym	Plaza	Gym / Fitness Center	Athletics & Sports	Park	Harbor / Marina	Martial Arts School	Skate Park	Soccer Field	Gymnastics Gym	Gym Pool	Lake	Basketball Court	Indoor Play Area
3	AEROPUERTO	Sculpture Garden	Playground	Plaza	Soccer Field	Garden	Gym	Yoga Studio	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym / Fitness Center	Gun Range
4	ALAMEDA DE OSUNA	Park	Plaza	Gym	Gym / Fitness Center	Tennis Court	Garden	Scenic Lookout	Yoga Studio	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gun Range

- Now, through K-Means Clustering unsupervised algorithm, it divides the data into K non-overlapping clusters grouping similar venues. It will be used for K=5. We will then merge Madrid data containing neighborhoods and coordinates, with the neighborhoods venues clustered.

borough code	borough	neighborhood	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	11th Most Common Venue	12th Most Common Venue
0	1 CENTRO	PALACIO	40.40958	-3.88009	0.0	Gym / Fitness Center	Plaza	Athletics & Sports	Park	Pool	Skate Park	Cycle Studio	Dog Run	Lake	Basketball Court	Indoor Play Area	Harbor / Marina
1	1 CENTRO	EMBAJADORES	40.39152	-3.69289	3.0	Park	Plaza	Gym	Playground	Gym / Fitness Center	Soccer Field	Yoga Studio	Garden	Pool	Skating Rink	Skate Park	Fountain
2	1 CENTRO	CORTES	40.41607	-3.69893	0.0	Plaza	Gym	Gym / Fitness Center	Athletics & Sports	Park	Garden	Roof Deck	Fountain	Yoga Studio	Sculpture Garden	Pedestrian Plaza	Lake
3	1 CENTRO	JUSTICIA	40.42479	-3.69308	0.0	Plaza	Park	Yoga Studio	Gym	Gym / Fitness Center	Athletics & Sports	Gymnastics Gym	Fountain	Playground	Garden	Gym Pool	Outdoors & Recreation
4	1 CENTRO	UNIVERSIDAD	40.42565	-3.70726	0.0	Plaza	Gym	Gym / Fitness Center	Park	Athletics & Sports	Garden	Gym Pool	Gymnastics Gym	Outdoors & Recreation	Roof Deck	Pedestrian Plaza	Playground

- Through the folium library show a map with the clusters.



# Coursera Capstone Project

- Now is time to analyze the different clusters to identify the best cluster to invest in a gym based on how frequently you can find them

## Cluster 1

madrid_merged.loc[madrid_merged['Cluster Labels'] == 0, madrid_merged.columns[[1] + list(range(5, madrid_merged.shape[1]))]]																
borough	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue	
0	CENTRO	0.0	Gym / Fitness Center	Plaza	Athletics & Sports	Park	Pool	Skate Park	Cycle Studio	Dog Run	Lake	Basketball Court	Indoor Play Area	Harbor / Marina	Botanical Garden	Boxing Gym
2	CENTRO	0.0	Plaza	Gym	Gym / Fitness Center	Park	Athletics & Sports	Garden	Roof Deck	Fountain	Yoga Studio	Sculpture Garden	Pedestrian Plaza	Lake	Playground	Tree

## Cluster 2

madrid_merged.loc[madrid_merged['Cluster Labels'] == 1, madrid_merged.columns[[1] + list(range(5, madrid_merged.shape[1]))]]																
borough	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue	
16	RETIRO	1.0	Park	Plaza	Garden	Fountain	Gym	Gym / Fitness Center	Lake	Yoga Studio	Harbor / Marina	Athletics & Sports	Playground	Sports Club	Scenic Lookout	Outdoors & Recreation
60	LATINA	1.0	Park	Athletics & Sports	Plaza	Gym	Playground	Soccer Field	Gym / Fitness Center	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gun Range	Mountain

## Cluster 3

madrid_merged.loc[madrid_merged['Cluster Labels'] == 2, madrid_merged.columns[[1] + list(range(5, madrid_merged.shape[1]))]]																
borough	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue	
43	FUENCARRAL-EL PARDO	2.0	Scenic Lookout	Yoga Studio	Outdoors & Recreation	Martial Arts School	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym / Fitness Center	Gym	Gun Range	Golf Course	Garden
45	FUENCARRAL-EL PARDO	2.0	Scenic Lookout	Yoga Studio	Outdoors & Recreation	Martial Arts School	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym / Fitness Center	Gym	Gun Range	Golf Course	Garden

## Cluster 4

madrid_merged.loc[madrid_merged['Cluster Labels'] == 3, madrid_merged.columns[[1] + list(range(5, madrid_merged.shape[1]))]]																
borough	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue	
1	CENTRO	3.0	Park	Plaza	Gym	Playground	Gym / Fitness Center	Soccer Field	Yoga Studio	Garden	Pool	Skating Rink	Skate Park	Fountain	Lake	Golf Course
13	RETIRO	3.0	Gym	Plaza	Park	Athletics & Sports	Martial Arts School	Harbor / Marina	Gym / Fitness Center	Playground	Gymnastics Gym	Skate Park	Soccer Field	Dog Run	Garden	Golf Course

## Cluster 5

madrid_merged.loc[madrid_merged['Cluster Labels'] == 4, madrid_merged.columns[[1] + list(range(5, madrid_merged.shape[1]))]]																
borough	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	11th Most Common Venue	12th Most Common Venue	13th Most Common Venue	14th Most Common Venue	
44	FUENCARRAL-EL PARDO	4.0	Park	Mountain	Dog Run	Gym	Martial Arts School	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym / Fitness Center	Yoga Studio	Outdoors & Recreation	Golf Course
46	FUENCARRAL-EL PARDO	4.0	Park	Mountain	Dog Run	Gym	Martial Arts School	Lake	Indoor Play Area	Harbor / Marina	Gymnastics Gym	Gym Pool	Gym / Fitness Center	Yoga Studio	Outdoors & Recreation	Golf Course

## Results

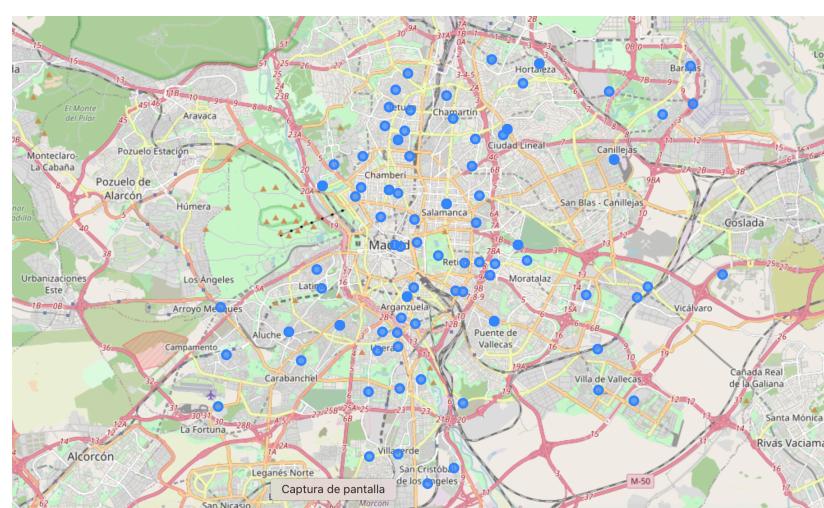
With the data now ready, we run k-means to cluster the neighborhoods into five clusters. The cluster number was established after multiple samplings and iterations. With our clusters established, this dataframe is merged with the total scores data to provide us with our final pieces of criteria in selecting the appropriate neighborhood(s).

Based on the analysis, the cluster where gym is lower in the “top 14 most common venues” list is the 2nd cluster, where we can observe the average position is about 10th position.

From the cluster 2 neighborhoods, the selection must be done based on the population and the lowest possible average age to maximize the number of potential clients

<b>borough</b>	<b>Population</b>	<b>Average age</b>
16 RETIRO	119379	47,13
60 LATINA	238154	46,55
74 USERA	139501	42,34
79 PUENTE DE VALLECAS	234770	43,3
85 MORATALAZ	94609	47,61
91 CIUDAD LINEAL	216270	45,87
107 VILLAVERDE	148883	42,15
113 VILLA DE VALLECAS	110436	38,9
114 VICALVARO	72126	40,83

## On the map:



# Discussion

From the results discovered and presented, the following observations and recommendations can be made:

- Based on the criteria given by the client and the cluster data, the main recommendation for a new gym would be Villa de Vallecas, followed by Puente de Vallecas (as both has a good potential and they are nearby).
- Villaverde and Usera are also good locations based on the population/age ratio
- Latina, even if it is the biggest in population, the average age of its population is quite high, reducing the amount of potential clients as the younger ones are more interested in high intensity activities, so my recommendation would be not to choose it at first instance.

# Conclusion

In conclusion, the scope of this of the analysis is somewhat limited. The venues to do business is ever changing, and the information available may be outdated as it relays on user information via Foursquare. Overall though, the model created can easily be replicated again and again with monitored data via the Foursquare API and the data from the forthcoming census in 2021.

With the data analyzed and scoring system established by the investor group, we stand by the recommendations made.