# Report: The Battle of Neighborhoods
## Coursera Capstone

## Kenyi Josué Ramírez

Data Sceince Course: Applied Data Science Capstone

# Overview

# The Problem

I currently live in Guadadalaja, Jalisco, México. I'm looking for a job as Data Scientist but my chances of finding a job increase if a look in a big City. Besides, I would like to find a similar city like Guadalajara.

# The Problem

The main characteristics of this question are:

- Average dairy income for each city in México.
- Population size, because this determine a better life quality.
- Venues around 5 km to the downtown.

# Proposal of solution

the detailed problem could be seen as a clustering problem in the first approach, we try to identify similar cities to Guadalajara. The derivation will be similar to the third assignment of this course with a lite variation in the data. In this case, we shall aggregate information about population size and average daily income for each city.

# Proposal of solution

The main tools for this Data Analysis will be:

### Tools

1. Pandas.
2. KMeans from Sklearn.
3. Numpy.

# Sources

**Sources**

1. Mexican Government Page
2. Wikipedia
3. https://foursquare.com/

The information to solve this problem is available in many sources. We will take the info from the 2020 Census made by INEGI (The National Institute of Statistics and Geography), this data is placed in Wikipedia and the official page of the government. The geospatial data and venue will be taken from Foursquare.

# Sources



Figure: HTML Table from Wikipedia.

# Sources



Figure: CVS sample (Is saved in GitHub as CSV file).

# Sources

The preprocessing of this data was made using pandas to create a unique table for clustering process.

# Clustering Model

We use K-Mean algorithm for clustering. In our data each variable represent a dimension in the space and we cosider each sample of data as a part of a data cloud (Cluster). This represent similarities between samples.

# Clustering Model

| | City | Nominal daily income | 2020 Census | Zoo Exhibit | Accessories Store | Airport | Airport Lounge | Airport Terminal | American Restaurant | Amp |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | Acapulco | 0.134864 | 0.056454 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.01 | 0.0 |
| 1 | Aguascalientes | 0.392071 | 0.079105 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 |
| 2 | Buenavista | 0.436511 | 0.007703 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 |
| 3 | Cabo San Lucas | 0.301220 | 0.006149 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 |
| 4 | Campeche | 0.735280 | 0.011327 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.00 | 0.0 |

Figure: Data table before the clustering.

# Clustering Model

We use the elbow method to choose the best k-value for the model. In this case we take k=5.



The Elbow Method using Distortion

# Clustering Model

We proceed to fit the model and we get this segmentation:

# Clustering Model

We see that Guadalajara lies in the zero cluster, that is the set of similar cities to Guadalajara.

```
1 clustered_data.loc[clustered_data['Cluster Labels'] == 0, clustered_data.columns[0]]
```

| | |
|---|---|
| 1 | Aguascalientes |
| 2 | Buenavista |
| 4 | Campeche |
| 7 | Chalco |
| 9 | Chicoloapan |
| 10 | Chihuahua |
| 12 | Chimalhuacán |
| 13 | Ciudad Acuña |
| 14 | Ciudad Apodaca |
| 15 | Ciudad Benito Juárez |
| 16 | Ciudad Juárez |
| 17 | Ciudad López Mateos |
| 18 | Ciudad Madero |
| 19 | Ciudad Nezahualcóyotl |
| 20 | Ciudad Nicolás Romero |
| 22 | Ciudad Victoria |
| 24 | Ciudad del Carmen |
| 25 | Coatzacoalcos |
| 27 | Cuautitlán Izcalli |
| 28 | Cuautla |
| 29 | Cuernavaca |
| 32 | Ecatepec |
| 33 | Ensenada |
| 34 | García |
| 35 | General Escobedo |
| 36 | Guadalajara |
| 37 | Guadalupe |
| 42 | Ixtapaluca |
| 43 | Jiutepec |
| 48 | Matamoros |

# Clustering Model

Finally we can choose the city in the first cluster with the highest Nominal Dairy income, in this case that corresponds to Querétaro City.

```
[34]  1 clustered_data[clustered_data['Cluster Labels'] == 0].sort_values("Nominal daily income", as
```
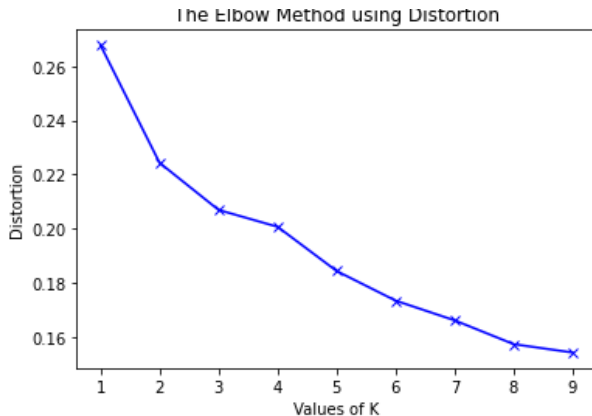
| | City | Nominal daily income | 2020 Census | Zoo Exhibit | Accessories Store | Airport | Airport Lounge | Airport Terminal | American Restaurant | A |
|---|---|---|---|---|---|---|---|---|---|---|
| 67 | Querétaro | 0.745298 | 0.071480 | 0.0 | 0.01 | 0.0 | 0.0 | 0.0 | 0.0 | |
| 73 | San Juan del Río | 0.745298 | 0.003393 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | |
| 24 | Ciudad del Carmen | 0.735280 | 0.004885 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | |
| 4 | Campeche | 0.735280 | 0.011327 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | |
| 15 | Ciudad Benito Juárez | 0.725316 | 0.017800 | 0.0 | 0.00 | 0.0 | 0.0 | 0.0 | 0.0 | |

5 rows × 345 columns

# Clustering Model

This conclusion makes sense because it is a growing city with a lot of new jobs and it is very similar to Guadalajara in the kind of venues[1].

---

[1] https://www.diariodequeretaro.com.mx/local/imparable-mancha-urbana-crece-2.9-5127464.html

# References

📄 Elbow Method for optimal value of k in KMeans
https://www.geeksforgeeks.org/elbow-method-for-optimal-value-of-k-in-kmeans/

📄 Panda Documentation
https://pandas.pydata.org/docs/

📄 SKlearn KMeans documentation
https://scikit-learn.org/stable/modules/generated/sklearn.cluster.KMeans.html

📄 INEGI
https://www.inegi.org.mx/

# The End