

# Cousera Capstone Project: Applied Data Science

Ayodeji Henry Ajisegiri  
Lagos, Nigeria

henryajisegiri@gmail.com

# Table of Content

<b>1.0</b>	<b>Introduction</b>
<b>2.0</b>	<b>Business Problem</b>
<b>3.0</b>	<b>Data</b>
<b>4.0</b>	<b>Methodology</b>
<b>5.0</b>	<b>Analysis</b>
<b>6.0</b>	<b>Results &amp; Discussion</b>
<b>7.0</b>	<b>Conclusion</b>

# Introduction

In 2019, the global tourism industry was worth over \$5 Trillion dollars according to the United Nations World Tourism Organization (UNTWO), with the United States contributing a whopping \$580 Billion (>11%).

UNTWO defines a tourist as someone who travels at least 80km from his or her home for at least 24 hours for business, leisure and/or other reasons. From the regal streets of Grand Bazaar in Istanbul to Nakamise street of Sensoji Temple in Tokyo, over 1.5 billion people thronged to different travel destinations in 2019.

# Introduction

Tourism is a great economic contributor and its impact can be felt across the following industries:

- Accommodation
- Food and Beverage Services
- Recreation and Entertainment
- Transportation
- Travel Services
- Retail Trade (souvenirs and the like)

For the aforementioned reasons, developing countries are looking to standardize their current tourism sites in order to attract international, continental and local tourists.

# Business Problem

My client is a West African country with breathtaking rainforests and a vast variety of wildlife in its savannah alongside other historical sites. The objective of this project is to investigate the ancillary infrastructure surrounding the best tourist sites in the world with the view of strategically replicating such infrastructure to ensure the best experience for potential tourists in order to maximize the impact on the local economy.

# Data

The data for this project was retrieved and process through multiple sources, however, the core data required are in two segments;

Longitude and Latitude of the 10 tourist locations

Location exploration using the above data points

The longitude and latitude were obtained manually because of the randomness of the locations while the location exploration will be obtained using Foursquare API.

# Data

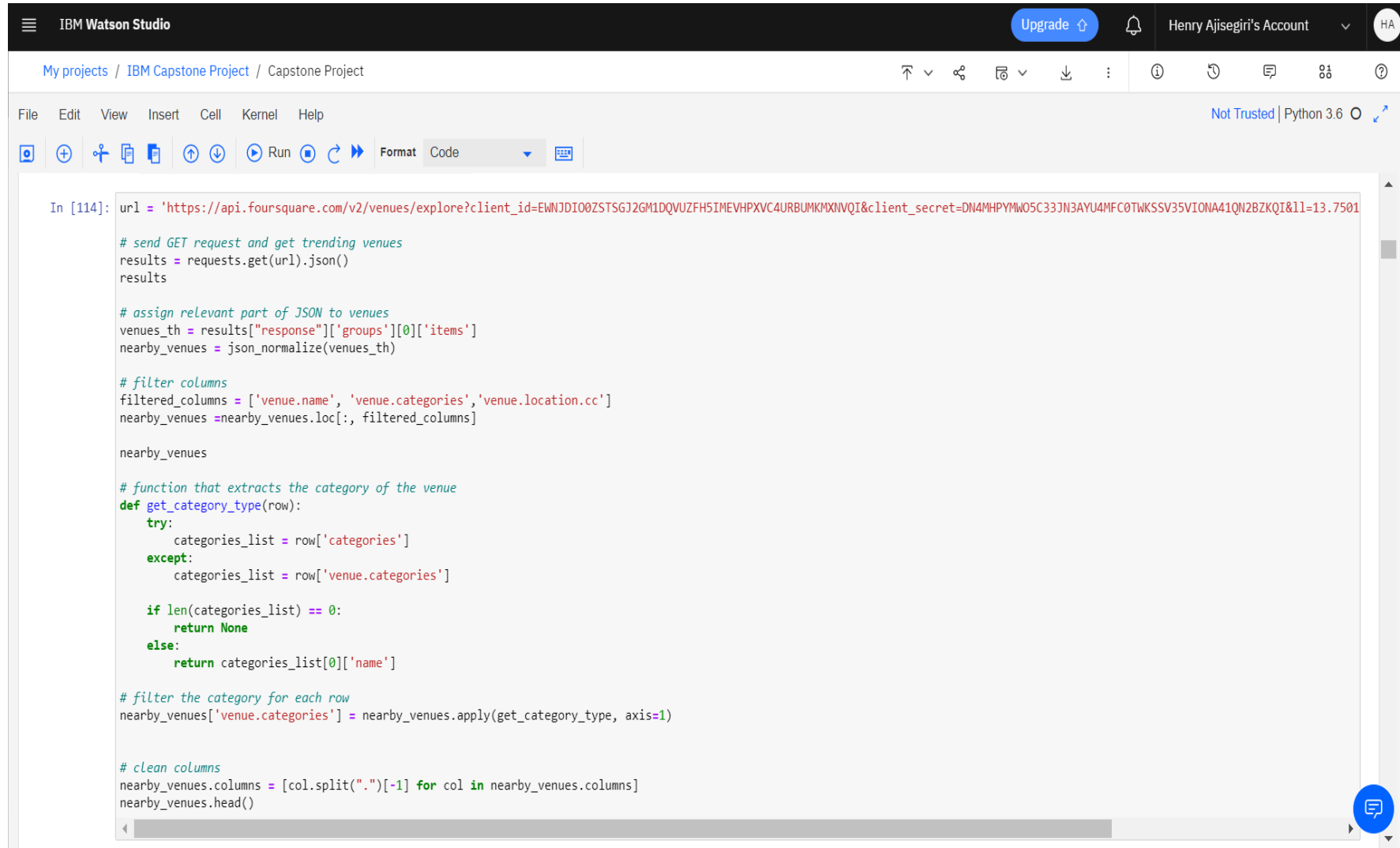
The major criteria for selecting the tourist sites are as follows:

- The site must receive atleast 1m visitors yearly
- The availability of the site's data on foursquare's database

The selected tourist sites are in the table below;

Tourist Site	Country	Country Code
The Forbidden City	China	CN
The Grand Palace	Thailand	TH
The Grand Bazaar	Turkey	TR
Sacre-Coeur Basilica	France	FR
St. Peter's Basilica	Vatican/Italy	VT/IT
Taj Mahal	India	IN
The Acropolis	Greece	GR
Eiffel Tower	France	FR
Sensoji Temple	Japan	JP
Burj Khalifa	UAE	AE

# Methodology



The screenshot displays the IBM Watson Studio web interface. At the top, there's a navigation bar with the IBM Watson Studio logo, an 'Upgrade' button, a notification bell, and a user profile for 'Henry Ajisegiri's Account'. Below this is a breadcrumb trail: 'My projects / IBM Capstone Project / Capstone Project'. A toolbar contains various icons for file operations, running, and formatting. The main area is a Jupyter Notebook with a Python 3.6 kernel. The code in the notebook performs the following steps:

- 1. Defines a URL to the Foursquare API for exploring venues.
- 2. Sends a GET request and parses the JSON response.
- 3. Extracts the relevant 'items' from the JSON response.
- 4. Normalizes the JSON data.
- 5. Filters the columns to include venue name, categories, and location.
- 6. Defines a function to extract the category type from the venue data.
- 7. Applies the function to filter the category for each row.
- 8. Cleans the columns by splitting them at the first period.

```
In [114]: url = 'https://api.foursquare.com/v2/venues/explore?client_id=EWNJDIO0ZSTSGJ2GM1DQVUFH5IMEVHPXVC4URBUMKMXNVQI&client_secret=DN4MHPYMW05C33JN3AYU4MFC0TWKSSV35VIONA41QN2BZKQI&l1=13.7501'

# send GET request and get trending venues
results = requests.get(url).json()
results

# assign relevant part of JSON to venues
venues_th = results["response"]["groups"][0]['items']
nearby_venues = json_normalize(venues_th)

# filter columns
filtered_columns = ['venue.name', 'venue.categories', 'venue.location.cc']
nearby_venues = nearby_venues.loc[:, filtered_columns]

nearby_venues

# function that extracts the category of the venue
def get_category_type(row):
    try:
        categories_list = row['categories']
    except:
        categories_list = row['venue.categories']

    if len(categories_list) == 0:
        return None
    else:
        return categories_list[0]['name']

# filter the category for each row
nearby_venues['venue.categories'] = nearby_venues.apply(get_category_type, axis=1)

# clean columns
nearby_venues.columns = [col.split(".")[1] for col in nearby_venues.columns]
nearby_venues.head()
```



# Methodology

## **One Hot Encoding**

One hot encoding is the process through which categorical variables are converted into a form (usually numerical) that can be provided to Machine Learning algorithms in order to perform efficiently during fitting and prediction. As stated earlier, K Means Clustering algorithm was used for this project and all unique items in the 'categories' column was one hot encoded.

## **Ten (10) Most Common Venues**

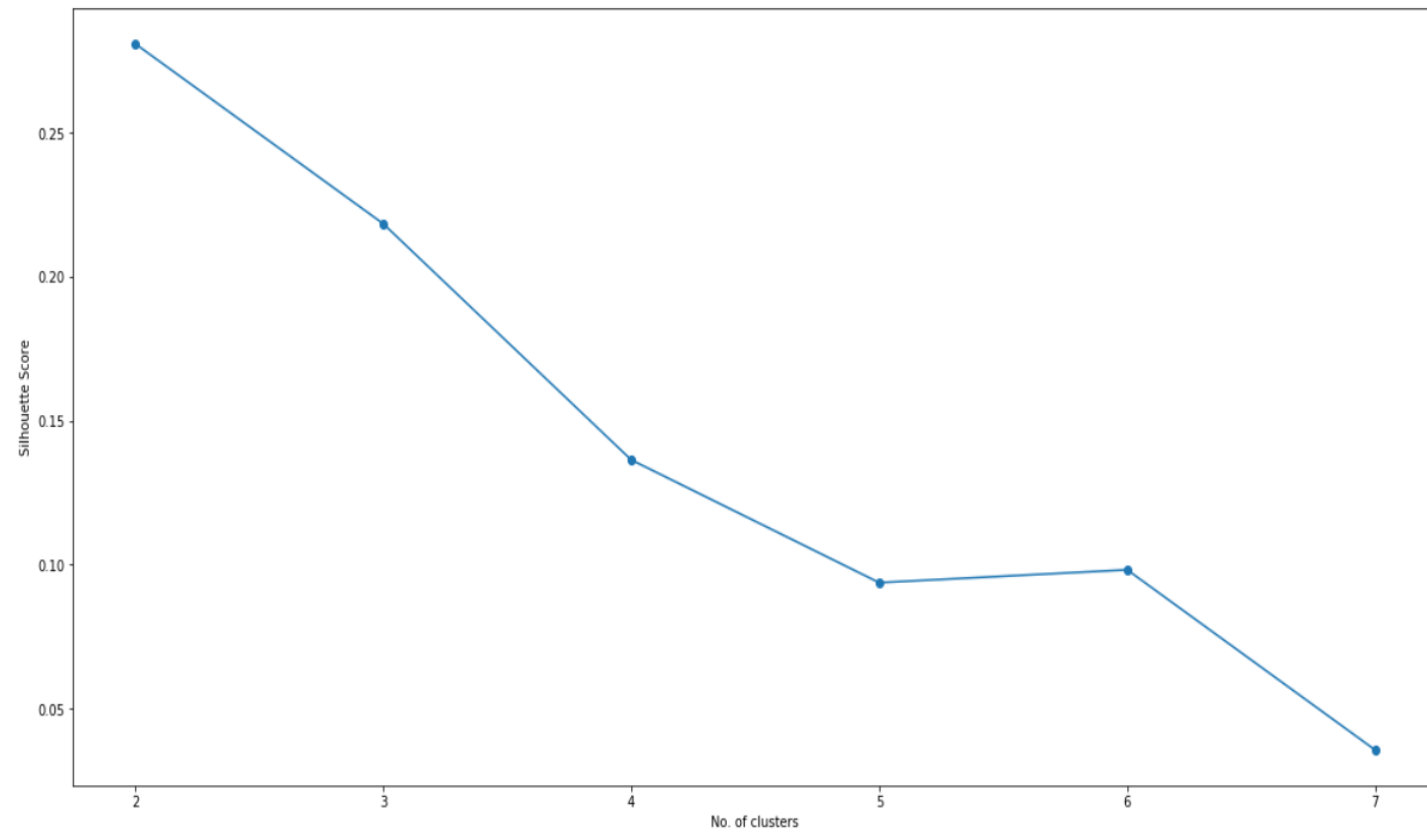
From over 900 locations and 191 unique categories, the 10 most popular venues were selected into a new dataframe for the purpose of training the K Means Clustering Algorithm.

# Methodology

## **Silhouette Score**

This is a measure of how similar an object is to its own cluster (cohesion) compared to other clusters (separation). The silhouette ranges from -1 to +1, where a high value indicates that the object is well matched to its own cluster and poorly matched to neighbouring clusters. Based on the Silhouette Score of various clusters below 20, the optimal number of clusters was determined.

# Methodology



*Silhouette score vs no. of clusters*

# Methodology

## **K Means Clustering**

The data obtained from the tourist sites' exploration was trained using the K Means Clustering algorithm in order to group the variation of ancillary tourist attractions into clusters which will serve as insight(s) to help the client take the best approach in creating an enabling environment for tourism to thrive in the country.

K Means was selected because of the size of the variables as K Means is computationally faster than other clustering algorithms.

## Results

My projects / IBM Capstone Project / Capstone Project

File Edit View Insert Cell Kernel Help
Not Trusted | Python 3.6

	Cluster Labels	country	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	
Out[164]:	0	1	AE	Boutique	Hotel	Café	Coffee Shop	Restaurant	Bakery	Clothing Store	Ice Cream Shop	Jewelry Store	Cosmetics Shop
	1	0	CN	Historic Site	Hotel	Café	Chinese Restaurant	Park	Hostel	Scenic Lookout	French Restaurant	Dive Bar	Sculpture Garden
	2	1	FR	French Restaurant	Italian Restaurant	Hotel	Bakery	Plaza	Bistro	Art Museum	Bar	Pizza Place	Restaurant
	3	1	GR	Historic Site	Bar	Greek Restaurant	Ice Cream Shop	Hotel	Coffee Shop	Café	Plaza	Dessert Shop	Wine Bar
	4	3	IN	Hotel	Indian Restaurant	Multicuisine Indian Restaurant	Resort	Fast Food Restaurant	Historic Site	Café	Market	Pizza Place	Coffee Shop
	5	1	IT	Italian Restaurant	Ice Cream Shop	Pizza Place	Hotel	Sandwich Place	Restaurant	Gourmet Shop	Art Gallery	Seafood Restaurant	Dessert Shop
	6	1	JP	Japanese Restaurant	Café	Wagashi Place	Unagi Restaurant	Dessert Shop	BBQ Joint	Soba Restaurant	Hostel	Sukiyaki Restaurant	Coffee Shop
	7	1	TH	Café	Thai Restaurant	Noodle House	Bar	Dessert Shop	Historic Site	Palace	History Museum	Hotel	Asian Restaurant
	8	1	TR	Hotel	Turkish Restaurant	Jewelry Store	Historic Site	Restaurant	Coffee Shop	Dessert Shop	Café	Art Gallery	Plaza
	9	2	VA	Church	Fountain	Museum	Plaza	Art Gallery	Art Museum	Monument / Landmark	French Restaurant	Grilled Meat Restaurant	Greek Restaurant

# Discussions

From the results, most (UAE, France, Greece, Italy, Japan, Thailand and Turkey) of the countries have similar infrastructure supporting their tourism industry. China, India and the Vatican are all in different clusters.

It is clear that the client should study the tourist sites in cluster 1 as they are among the top visited tourist attractions in the world and the size of the cluster also shows that a lot of countries have shaped their infrastructure in similar fashion. It is also important to study the uniqueness of China, India and the Vatican.

From the data exploration phase, it was observed that Asian restaurants are well distributed in all tourist attractions with Italian restaurant coming a distant second. This may be a response to the huge number of Asians visiting several sites around the world or the fact that the world loves Asian food, I choose the former.

# Conclusion

Tourism in many developing and least developed countries is the most viable and sustainable economic development option, and in some countries, the main source of foreign exchange earnings. For the client, tourism is a great option for diversification and enabling the support structure similar to cluster 1 countries will lead to more jobs via small to medium businesses, improved GDP and increased foreign exchange earnings.