

## 1.0 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 atleast 80km from his or her home for atleast 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.

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

1. Accommodation
2. Food and Beverage Services
3. Recreation and Entertainment
4. Transportation
5. Travel Services
6. 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.

## 2.0 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.

This objective will be achieved by randomly selecting 10 popular tourist attractions across 5 continents, exploring existing outlets within 600 metres radius of the tourist site by using Foursquare location API and clustering similar outlets using K-Means Clustering algorithm.

The table below contains the biggest tourist spenders of 2018;

Country	Amount Spent (\$)
China	277 billion
United States	144 billion
Germany	94 billion
United Kingdom	76 billion
France	48 billion
Australia	37 billion

Russia	35 billion
Canada	33 billion
South Korea	32 billion
Italy	30 billion

How does this West African country build its tourism infrastructure to attract these billions of dollars?

### 3.0 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.

### 3.1 Tourist Site Selection

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

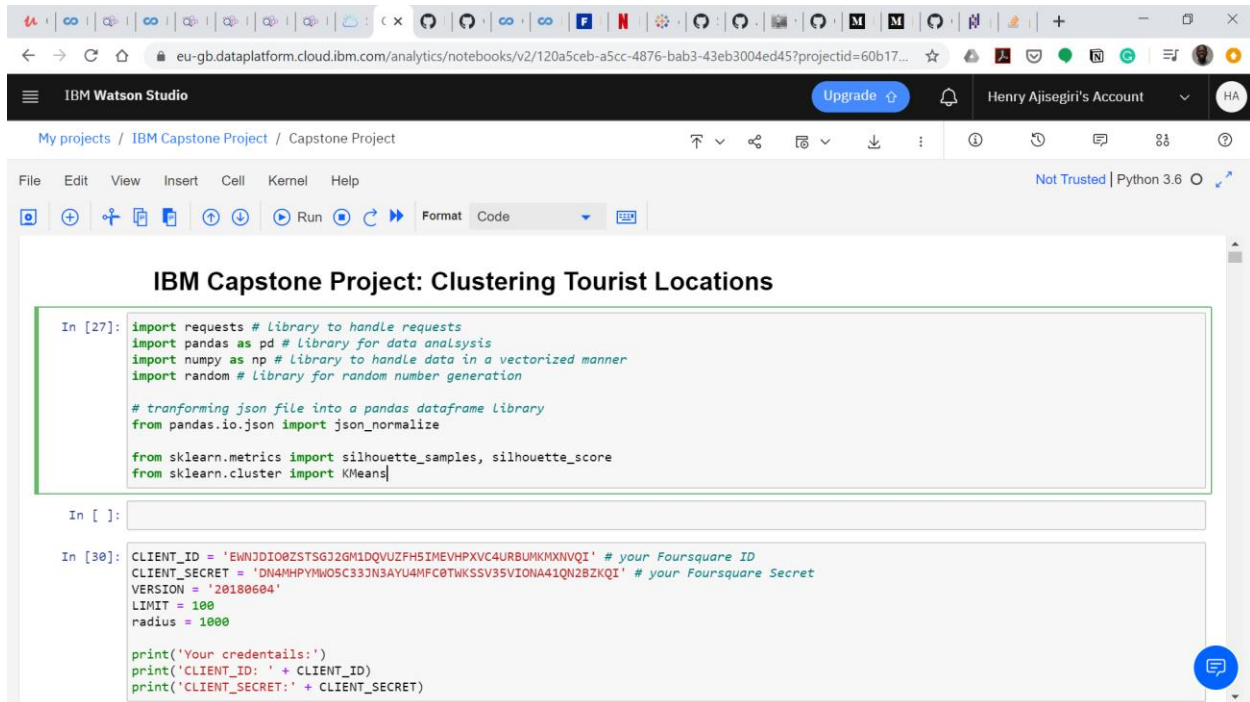
- The site must receive atleast 1m visitors yearly
- The availability of the site's data on foursquare's database

The selected tourist sites are in the table below;

<b>Tourist Site</b>	<b>Country</b>	<b>Country Code</b>
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

## 4.0 METHODOLOGY

The entire analysis was performed using Python Language and Machine Learning algorithm (K Means) in a Jupyter IDE hosted on IBM Watson Studio. The libraries used are shown in the image below;



```
In [27]: import requests # Library to handle requests
import pandas as pd # Library for data analysis
import numpy as np # Library to handle data in a vectorized manner
import random # Library for random number generation

# transforming json file into a pandas dataframe library
from pandas.io.json import json_normalize

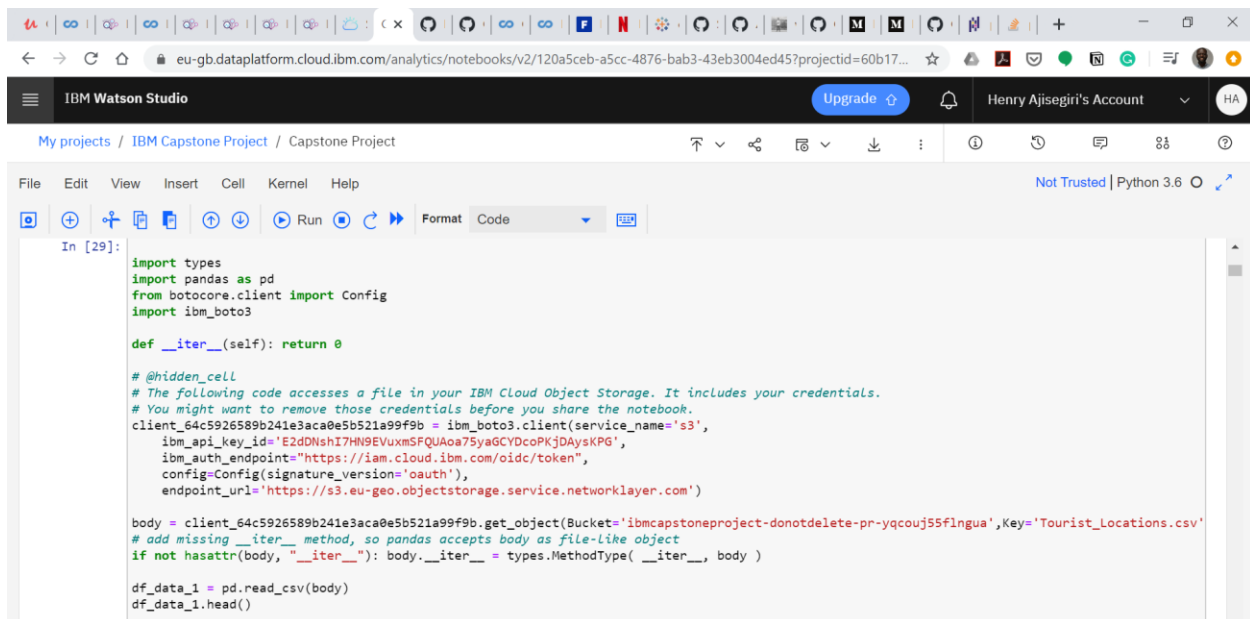
from sklearn.metrics import silhouette_samples, silhouette_score
from sklearn.cluster import KMeans

In [ ]:

In [30]: CLIENT_ID = 'EWNJDI0BZSTSGJ2GM1DQVUFH5IMEVHPXVC4URBWMKMXNVQI' # your Foursquare ID
CLIENT_SECRET = 'DN4MHPYMW05C33JN3AYU4MFC0TNKSSV35VIONA44QN2BZKQI' # your Foursquare Secret
VERSION = '20180604'
LIMIT = 100
radius = 1000

print('Your credentials:')
print('CLIENT_ID: ' + CLIENT_ID)
print('CLIENT_SECRET: ' + CLIENT_SECRET)
```

The manually created csv file was loaded directly into Watson studio as shown below;



```
In [29]: import types
import pandas as pd
from botocore.client import Config
import ibm_boto3

def __iter__(self): return 0

# @hidden_cell
# The following code accesses a file in your IBM Cloud Object Storage. It includes your credentials.
# You might want to remove those credentials before you share the notebook.
client_64c5926589b241e3aca0e5b521a99f9b = ibm_boto3.client(service_name='s3',
    ibm_api_key_id='E2dDNshI7HN9EVuxmSFQUAoa75yaGCDcoPKjDAYSKPG',
    ibm_auth_endpoint="https://iam.cloud.ibm.com/oidc/token",
    config=Config(signature_version='oauth'),
    endpoint_url='https://s3.eu-geo.objectstorage.service.networklayer.com')

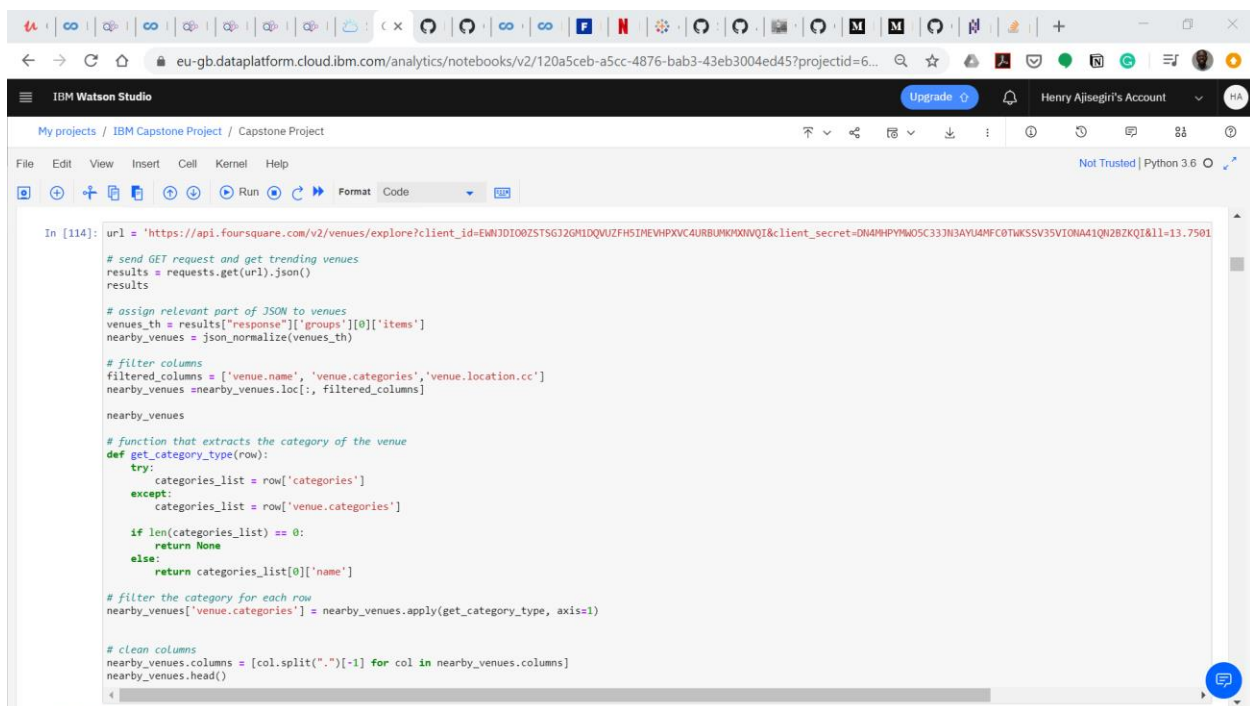
body = client_64c5926589b241e3aca0e5b521a99f9b.get_object(Bucket='ibmcapstoneproject-donotdelete-pr-yqcouj55flngua',Key='Tourist_Locations.csv')
# add missing __iter__ method, so pandas accepts body as file-like object
if not hasattr(body, "__iter__"): body.__iter__ = types.MethodType(__iter__, body)

df_data_1 = pd.read_csv(body)
df_data_1.head()
```

## 5.0 ANALYSIS

### 5.1 Foursquare

Foursquare's API was used to extract the details of the most popular spots around the tourist sites earlier mentioned. The specified radius coverage is 1000m and the number of spots extracted was limited to 100 per site. The image below shows the Foursquare API call used to extract the required information and the data wrangling process to make the data presentable; this code was run 10 times for the 10 different locations selected.



```
In [114]: url = 'https://api.foursquare.com/v2/venues/explore?client_id=EMW3DIO8Z5TSG3ZGMDQVZFHS1MEVHPXVC4URBUMQ00VQI&client_secret=DN4MPYMW05C33JN3AYU4MFC0TWKSSV3SVIONA41QN2BZKQI&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)

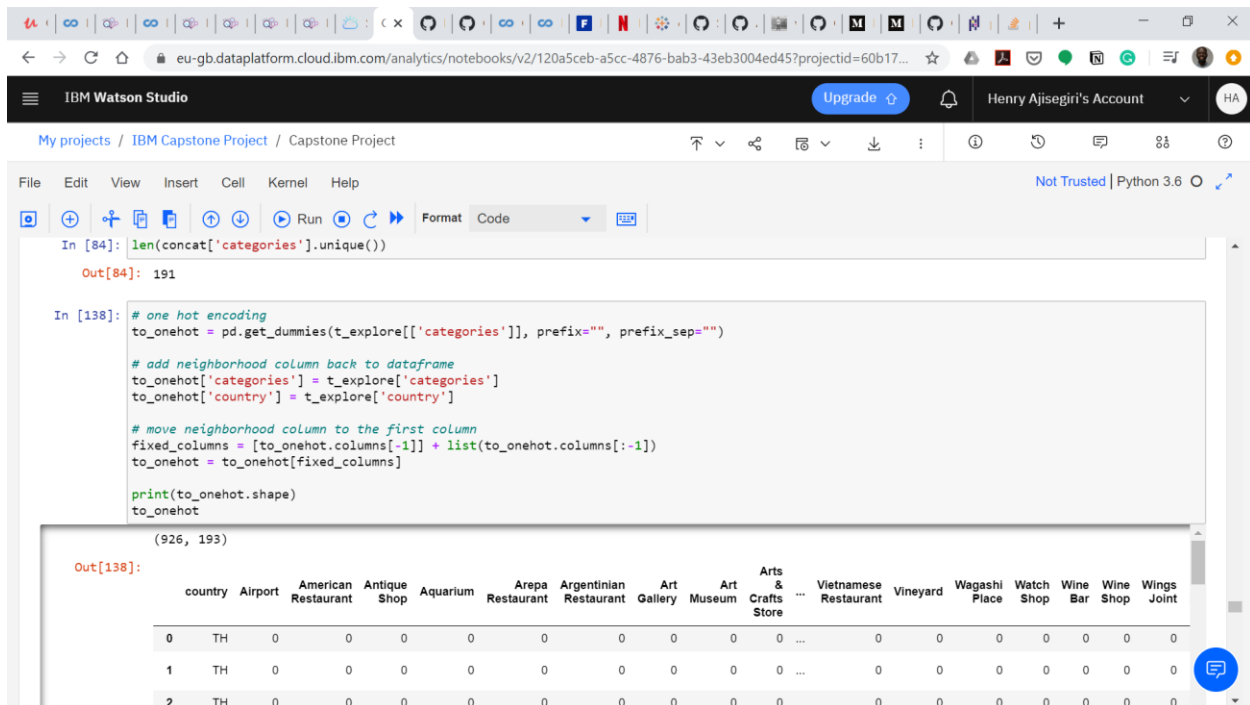
# clean columns
nearby_venues.columns = [col.split(".")[1] for col in nearby_venues.columns]
nearby_venues.head()
```

#### Key steps

- The required authentication details were filled into the API url
- Necessary information was read from the json file into a dataframe
- The dataframe was cleaned up for better presentation and easier operations down the line

### 5.2 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.



The screenshot shows an IBM Watson Studio notebook interface. The top bar includes the IBM Watson Studio logo, an 'Upgrade' button, and the user's account 'Henry Ajisegiri's Account'. The notebook is titled 'My projects / IBM Capstone Project / Capstone Project'. The code editor shows the following Python code:

```
In [84]: len(concat['categories'].unique())
Out[84]: 191

In [138]: # one hot encoding
to_onehot = pd.get_dummies(t_explore[['categories']], prefix="", prefix_sep="")

# add neighborhood column back to dataframe
to_onehot['categories'] = t_explore['categories']
to_onehot['country'] = t_explore['country']

# move neighborhood column to the first column
fixed_columns = [to_onehot.columns[-1]] + list(to_onehot.columns[:-1])
to_onehot = to_onehot[fixed_columns]

print(to_onehot.shape)
to_onehot
```

The output of the code is a DataFrame with 926 rows and 193 columns. The first three rows are shown:

	country	Airport	American Restaurant	Antique Shop	Aquarium	Arepa Restaurant	Argentinian Restaurant	Art Gallery	Art Museum	Arts & Crafts Store	...	Vietnamese Restaurant	Vineyard	Wagashi Place	Watch Shop	Wine Bar	Wine Shop	Wings Joint
0	TH	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0
1	TH	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0
2	TH	0	0	0	0	0	0	0	0	0	...	0	0	0	0	0	0	0

### 5.3 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.

```
In [147]: def return_most_common_venues(row, num_top_venues):
row_categories = row.iloc[1:]
row_categories_sorted = row_categories.sort_values(ascending=False)
return row_categories_sorted.index.values[0:num_top_venues]

In [149]: num_top_venues = 10

indicators = ['st', 'nd', 'rd']

# create columns according to number of top venues
columns = ['country']
for ind in np.arange(num_top_venues):
    try:
        columns.append('{} {} Most Common Venue'.format(ind+1, indicators[ind]))
    except:
        columns.append('{}th Most Common Venue'.format(ind+1))

# create a new dataframe
t_explore_sorted = pd.DataFrame(columns=columns)
t_explore_sorted['country'] = to_grouped['country']

for ind in np.arange(to_grouped.shape[0]):
    t_explore_sorted.iloc[ind, 1:] = return_most_common_venues(to_grouped.iloc[ind, :], num_top_venues)
```

## 5.4 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.

```
In [143]: max_range = 8

tourism_grouped_clustering = to_grouped.drop('country', 1)

from sklearn.metrics import silhouette_samples, silhouette_score

indices = []
scores = []

for kclusters in range(2, max_range):

    # Run k-means clustering
    kmc = tourism_grouped_clustering
    kmeans = KMeans(n_clusters = kclusters, init = 'k-means++', random_state = 0).fit_predict(kmc)

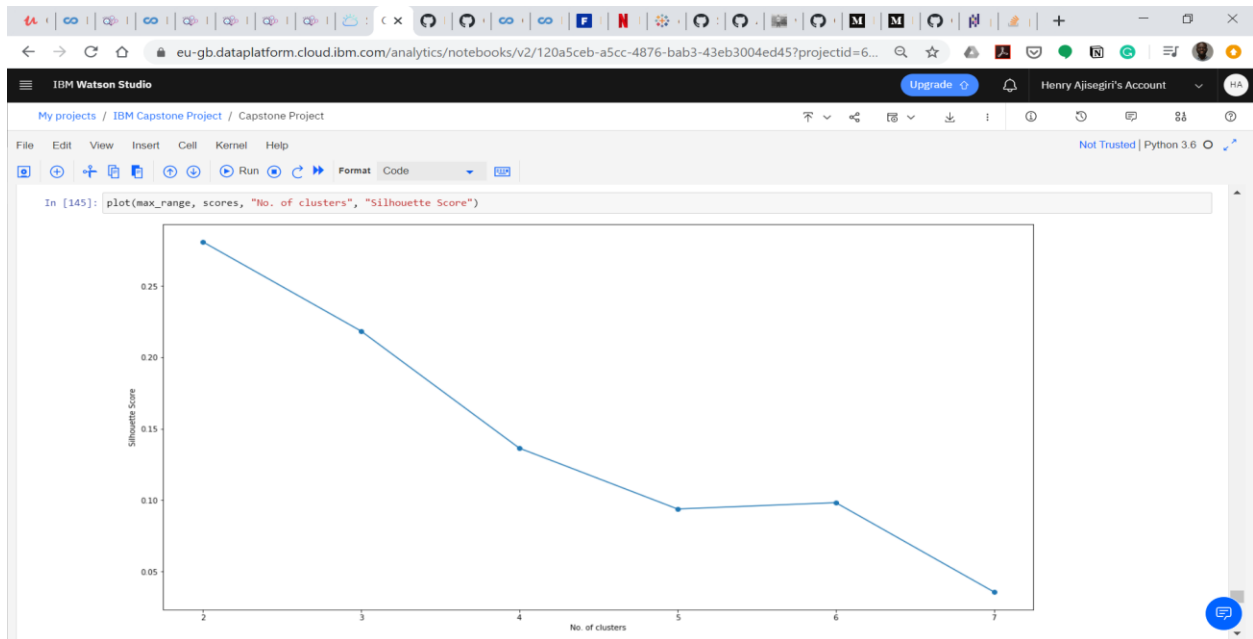
    # Gets the score for the clustering operation performed
    score = silhouette_score(kmc, kmeans)

    # Appending the index and score to the respective lists
    indices.append(kclusters)
    scores.append(score)

In [144]: import matplotlib.pyplot as plt
%matplotlib inline

def plot(x, y, xlabel, ylabel):
    plt.figure(figsize=(20,10))
    plt.plot(np.arange(2, x), y, 'o-')
    plt.xlabel(xlabel)
    plt.ylabel(ylabel)
    plt.xticks(np.arange(2, x))
    plt.show()
```

With the codes above, the graph below was plotted;

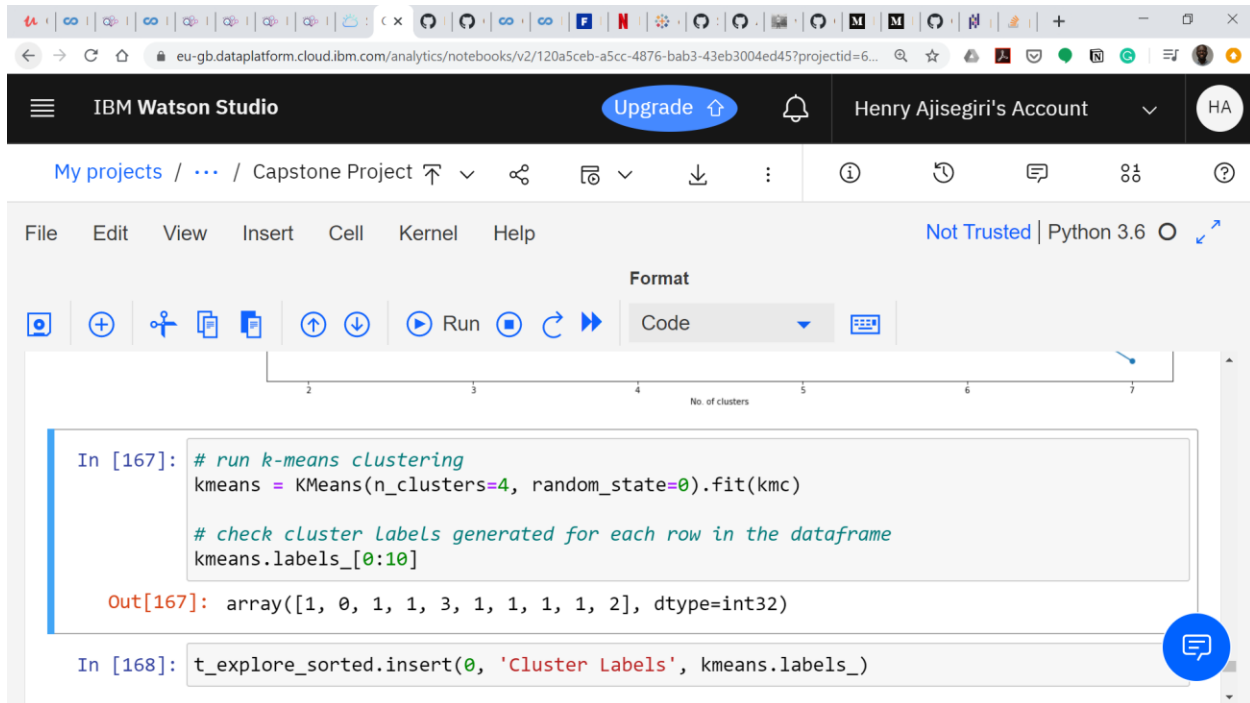


From the Silhouette Score, the optimal number of clusters is two (2). However, since the business problem is not strongly based on high precision and accuracy, the K Mean Clustering algorithm was performed with four (4) clusters.

## 5.5 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.



```
In [167]: # run k-means clustering
kmeans = KMeans(n_clusters=4, random_state=0).fit(kmc)

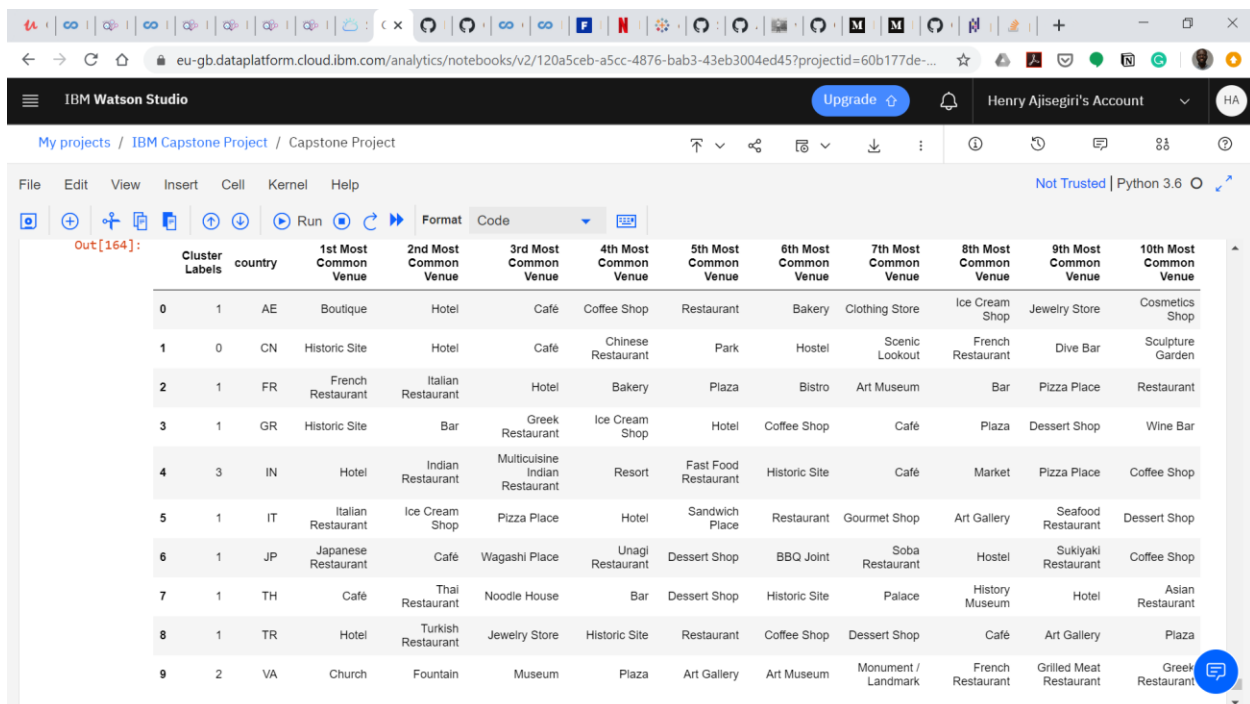
# check cluster labels generated for each row in the dataframe
kmeans.labels_[0:10]

Out[167]: array([1, 0, 1, 1, 3, 1, 1, 1, 2], dtype=int32)

In [168]: t_explore_sorted.insert(0, 'Cluster Labels', kmeans.labels_)
```

## 6.0 RESULTS & DISCUSSION

The countries were divided into 4 clusters but the clusters cannot be visualized with a map because the locations selected are vastly dispersed across the world. However, the result can be viewed on a table.



	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
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



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.

## **7.0 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.