# Coursera Capstone

## IBM Applied Data Science Capstone

## *Opening a New SPA in Mumbai*

## Introduction:

Identify the location for setting up new SPA business in the Mumbai city.

## Business Problem:

The objective of this capstone project is to analyze and select the best locations in the city of Mumbai to open a new SPA. Using data science methodology and machine learning techniques like clustering, this project aims to provide solutions to answer the business question: In the city of Mumbai if an investor or individual is looking to open a new SPA, where would you recommend that they open it?

# Data:

**To solve the problem, we will need the following data:**

• List of neighborhoods in Mumbai. This defines the scope of this project which is confined to the Mumbai, the financial capital city of the country of India.

• Latitude and longitude coordinates of those neighborhoods. This is required in order to plot the map and also to get the venue data.

• Venue data, particularly data related to SPA. We will use this data to perform clustering on the neighborhoods.

**Sources of data and methods to extract them**

This Wikipedia page (https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Mumbai) contains a list of neighborhoods in Mumbai, with a total of 136 neighborhoods. We will use web scraping techniques to extract the data from the Wikipedia page, with the help of Python requests and beautifulsoup4 packages. Then we will get the geographical coordinates of the neighborhoods using Python Geocoder package which will give us the latitude and longitude coordinates of the neighborhoods.
After that, we will use Foursquare API to get the venue data for those neighbourhoods. Foursquare API will provide many categories of the venue data, we are particularly interested in the SPA category in order to help us to solve the business problem put forward. This is a project that will make use of many data science skills, from web scraping (Wikipedia), working with API (Foursquare), data cleaning, data wrangling, to machine learning (K-means clustering) and map visualization (Folium). In the next section, we will present the Methodology section where we will discuss the steps taken in this project, the data analysis that we did and the machine learning technique that was used.
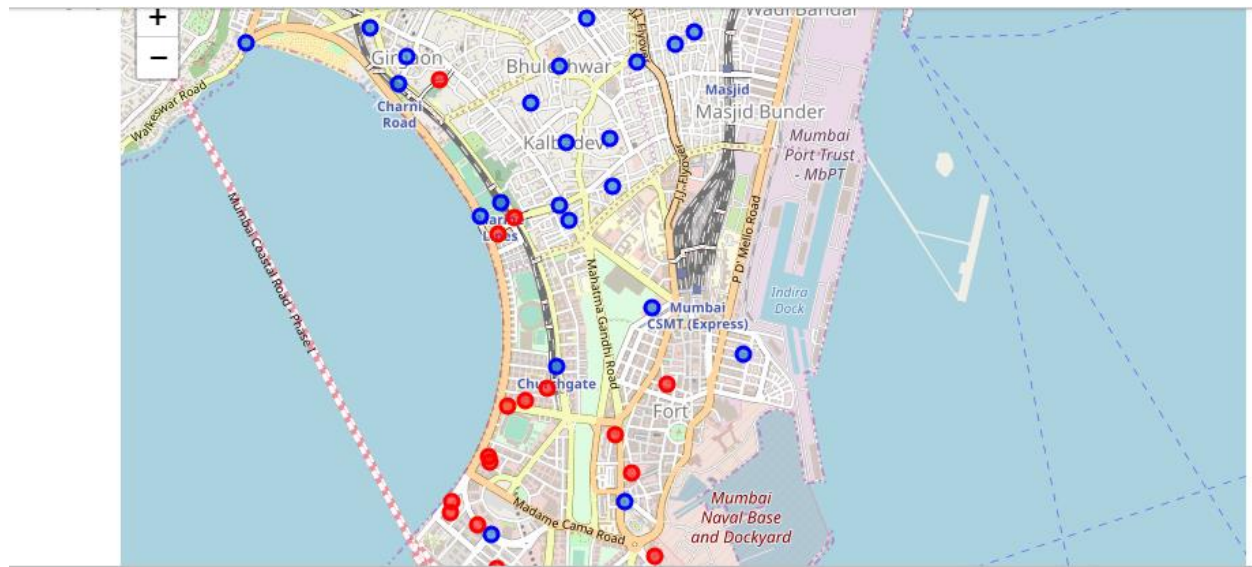
# Methodology

Firstly, we need to get the list of neighbourhoods in the city of Mumbai. Fortunately, the list is available in the Wikipedia page (https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Mumbai). We will do web scraping using Python requests and beautifulsoup4 packages to extract the list of neighbourhoods data. However, this is just a list of names. We need to get the geographical coordinates in the form of latitude and longitude in order to be able to use Foursquare API. To do so, we will use the wonderful Geocoder package that will allow us to convert address into geographical coordinates in the form of latitude and longitude. After gathering the data, we will populate the data into a pandas DataFrame and then visualize the neighbourhoods in a map using Folium package. This allows us to perform a sanity check to make sure that the geographical coordinates data returned by Geocoder are correctly plotted in the city of Mumbai.
Next, we will use Foursquare API to get the top 50 venues that are within a radius of 2000 meters. We need to register a Foursquare Developer Account in order to obtain the Foursquare ID and Foursquare secret key. We then make API calls to Foursquare passing in the geographical coordinates of the

neighbourhoods in a Python loop. Foursquare will return the venue data in JSON format and we will extract the venue name, venue category, venue latitude and longitude. With the data, we can check how many venues were returned for each neighbourhood and examine how many unique categories can be curated from all the returned venues. Since we are analysing the "SPA" data, we will filter the "SPA" as venue category for the neighbourhoods.

Lastly, we will perform clustering on the data by using k-means clustering. K-means clustering algorithm identifies k number of centroids, and then allocates every data point to the nearest cluster, while keeping the centroids as small as possible. It is one of the simplest and popular unsupervised machine learning algorithms and is particularly suited to solve the problem for this project. We will cluster the neighbourhoods into 3 clusters based on their frequency of occurrence for "SPA". The results will allow us to identify which neighbourhoods have higher concentration of SPA while which neighbourhoods have fewer number of SPA. Based on the occurrence of SPA in different neighbourhoods, it will help us to answer the question as to which neighbourhoods are most suitable to open new SPA.

# Results

The results from the k-means clustering show that we can categorize the neighbourhoods.
The results of the clustering are visualized in the map below with RED most of the SPA are already present in the South Mumbai the central and north Mumbai section is completely open for the business.

# Limitations and Suggestions for Future Research

In this project, we only consider one factor i.e. frequency of occurrence of SPA, there are other factors such as population and income of residents that could influence the location decision of a new SPA. However, to the best knowledge of this researcher such data are not available to the neighbourhood level required by this project. Future research could devise a methodology to estimate such data to be used in the clustering algorithm to determine the preferred locations to open a new SPA. In addition, this project made use of the free Sandbox Tier Account of Foursquare API that came with limitations as to the number of API calls and results returned. Future research could make use of paid account to bypass these limitations and obtain more results.

# Conclusion

In this project, we have gone through the process of identifying the business problem, specifying the data required, extracting and preparing the data, performing machine learning by clustering the data based on their similarities, and lastly providing recommendations to the relevant stakeholders i.e. investors, Individual business owners regarding the best locations to open a new SPA center. The findings of this project will help the relevant stakeholders to capitalize on the opportunities on high potential locations while avoiding overcrowded areas in their decisions to open a new SPA center.