

# Coursera Capstone

## IBM Applied Data Science Capstone

Opening a new Indian Restaurant in Hyderabad.

By Rajkumar

June 2020



### Introduction:

For many foodies, visiting restaurants is a great way to relax and enjoy themselves during weekends and holidays. A **restaurant** (sometimes known as a **diner**) is a place where cooked food is sold to the public, and where people sit down to eat it. It is also a place where people go to enjoy the time and to eat a meal.

**Indian cuisine** consists of a variety of regional and traditional cuisines native to the Indian Continent. Given the diversity in soil, climate, culture, ethnic groups, and occupations, these cuisines vary substantially and use locally available spices, herbs, vegetables, and fruits. Indian food is also heavily influenced by religion, in particular Hinduism, cultural choices and traditions. The cuisine is also influenced by centuries of Islamic rule, particularly the Mughal rule. Samosas and pilaffs are examples

Of course, as with any business decision, opening a new Restaurant requires serious consideration and is a lot more complicated than it seems. Particularly, the location of the restaurant is one of the most important decisions that will determine whether the restaurant will be a success or a failure.

## **Business Problem:**

The objective of this capstone project is to analyse and select the best locations in the city of Hyderabad, India to open an Indian Restaurant. 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 Hyderabad, India, if a property developer is looking to open a new Indian Restaurant, where would you recommend that they open it?

## **Target Audience of this project:**

This project is particularly useful to property developers and investors looking to open or invest in new Indian Restaurants in the capital city of India i.e. Hyderabad. This project is timely as the city is currently suffering from oversupply of Indian restaurants. Data from the National Property Information Centre (NAPIC) released last year showed that an additional 15 per cent will be added to existing restaurant space, and the agency predicted that total occupancy may dip below 86 per cent.

## **Data:**

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

- List of neighbourhoods in Hyderabad. This defines the scope of this project which is confined to the city of Hyderabad, the capital of Telangana.
- Latitude and longitude coordinates of these neighbourhoods. This is required in order to plot the map and also to get the venue data.
- Venue data, particularly data related to restaurants. We will use this data to perform clustering on the neighbourhoods.

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

This Wikipedia page

([https://en.wikipedia.org/wiki/Category:Neighbourhoods\\_in\\_Hyderabad, Ind](https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Hyderabad,_Ind)

[ia](#)) contains a list of neighbourhoods in Hyderabad, with a total of 200 neighbourhoods. We will use web scraping techniques to extract the data from the Wikipedia page, with the help of Python requests and beautiful soup packages. Then we will get the geographical coordinates of the neighbourhoods using Python Geocoder package which will give us the latitude and longitude coordinates of the neighbourhoods.

After that, we will use Foursquare API to get the venue data for those neighbourhoods. Foursquare has one of the largest databases of 105+ million places and is used by over 125,000 developers.

Foursquare API will provide many categories of the venue data, we are particularly interested in the Restaurant 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 Hyderabad. Fortunately, the list is in Wikipedia page ([https://en.wikipedia.org/wiki/Category:Neighbourhoods\\_in\\_Hyderabad,\\_In](https://en.wikipedia.org/wiki/Category:Neighbourhoods_in_Hyderabad,_In)). We will do web scraping using Python requests and beautiful soup packages to extract the list of neighbourhood's 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 Data Frame 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 Hyderabad. Next, we will use Foursquare API to get the top 100 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

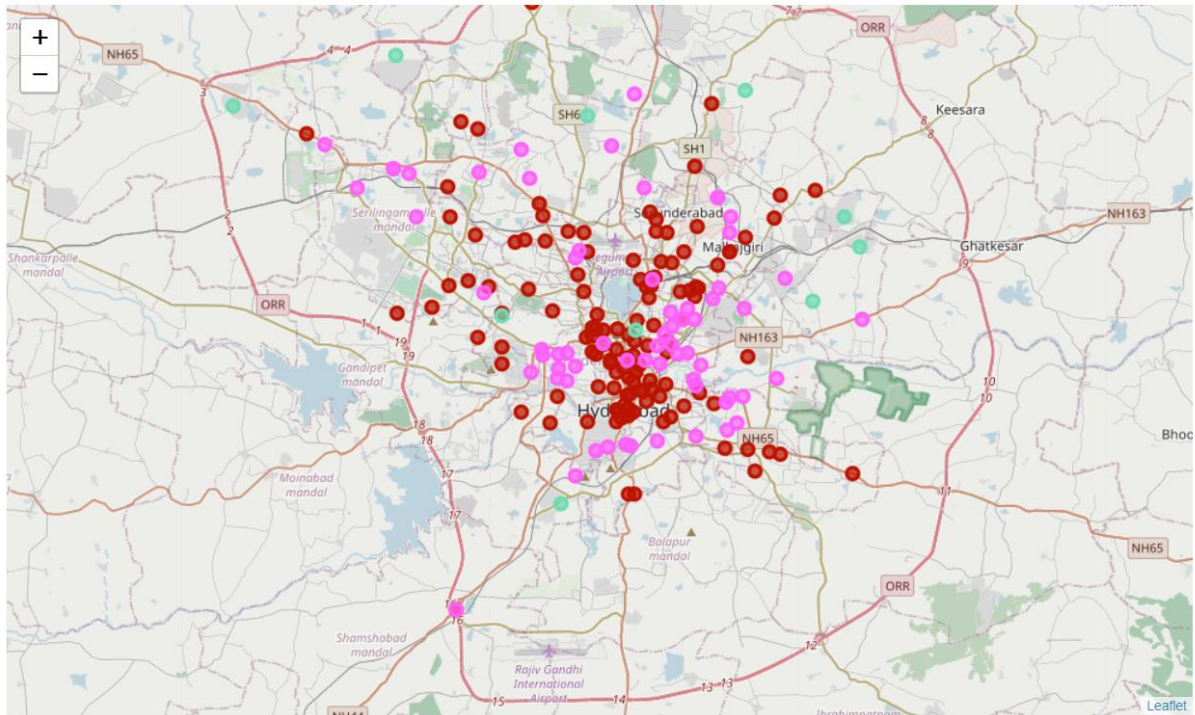
curated from all the returned venues. Then, we will analyse each neighbourhood by grouping the rows by neighbourhood and taking the mean of the frequency of occurrence of each venue category. By doing so, we are also preparing the data for use in clustering. Since we are analysing the “Indian Restaurants” data, we will filter the “Indian Restaurants” 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 “Indian Restaurants”. The results will allow us to identify which neighbourhoods have higher concentration of Indian Restaurants while which neighbourhoods have fewer number of Indian Restaurants. Based on the occurrence of Indian Restaurants in different neighbourhoods, it will help us to answer the question as to which neighbourhoods are most suitable to open new Indian Restaurants.

## **Results**

The results from the k-means clustering show that we can categorize the neighbourhoods into 3 clusters based on the frequency of occurrence for “Indian Restaurants”

- Cluster 1: Neighbourhoods with high concentration of Indian Restaurants
- Cluster 2: Neighbourhoods with moderate number of Indian Restaurants
- Cluster 3: Neighbourhoods with low number to no existence of Indian Restaurants.

The results of the clustering are visualized in the map below with cluster 1 in red colour, cluster 2 in light pink colour, and cluster 3 in mint green colour.



## Discussion

As observations noted from the map in the Results section, most of the Indian Restaurants are concentrated in the central area of Hyderabad city, with the highest number in cluster 1 and moderate number in cluster 2. On the other hand, cluster 3 has very low number to no Indian Restaurants in the neighbourhoods. This represents a great opportunity and high potential areas to open Indian Restaurants as there is very little to no competition from existing Restaurants. Meanwhile, Indian Restaurants in cluster 1 are likely suffering from intense competition due to oversupply and high concentration of Indian Restaurants. From another perspective, the results also show that the oversupply of Indian Restaurants mostly happened in the central area of the city, with the suburb area still have very few Indian Restaurants. Therefore, this project recommends property developers to capitalize on these findings to open new Indian Restaurants in neighbourhoods in cluster 3 with little to no competition. Property developers with unique selling propositions to stand out from the competition can also open new Indian Restaurants in neighbourhoods in cluster 2 with moderate competition. Lastly, property developers are advised to avoid neighbourhoods in cluster 1 which already have high concentration of shopping malls and suffering from intense competition.

## Limitations and Suggestions for Future Research

In this project, we only consider one factor i.e. frequency of occurrence of Indian Restaurants, there are other factors such as population and

income of residents that could influence the location decision of a new Indian Restaurants. 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 Indian Restaurants. 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 into 3 clusters based on their similarities, and lastly providing recommendations to the relevant stakeholders i.e. property developers and investors regarding the best locations to open a new Indian Restaurant.. To answer the business question that was raised in the introduction section, the answer proposed by this project is: The neighbourhoods in cluster 3 are the most preferred locations to open a new Indian Restaurant. 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 Indian Restaurant.