

Capstone Project - The Battle of  
Neighborhoods  
Geolocation analysis for opening shopping  
mall in Bangkok  
Phiphat Prapapanpong  
January 12, 2020



# Business Problem

---

Bangkok is the capital and most populous city of Thailand. Over 8 million people lived within Bangkok and this city is the most visited city in the world with approximately 23 million tourists per year. The shopping malls are one of the places that Thai people and tourists usually visit and spend a lot of money. Therefore, if we are able to help decision making on finding the best location to build a shopping mall. It can reduce the risk for investors and also able to attract more tourists to visit and spend. Businesses can make more informed decisions that can improve both efficiency and effectiveness

In this project, I will try to understand the preferences of each district by leveraging venue data from Foursquare's 'Places API' and 'k-means clustering' machine learning algorithm.

# Data Acquisition and Cleansing

- List of districts in Bangkok. This defines the scope of this project which is confined to the city of Bangkok from [https://en.wikipedia.org/wiki/List\\_of\\_districts\\_of\\_Bangkok](https://en.wikipedia.org/wiki/List_of_districts_of_Bangkok), In this data, we can see district code, district name, population, latitude and longitude
- The list of Latitude and longitude coordinates of these districts is required in order to plot the map and also to get the nearby venue data from Foursquare
- Venue data from Foursquare API, particularly data related to shopping malls. I will use this data to perform clustering on the districts.



# Data extraction and data cleansing 1

- Web scraping technique “pandas.io.html.read\_html” was used to extract the data from the Wikipedia page “[https://en.wikipedia.org/wiki/List\\_of\\_districts\\_of\\_Bangkok](https://en.wikipedia.org/wiki/List_of_districts_of_Bangkok)”,

```
1 from pandas.io.html import read_html
2 url='https://en.wikipedia.org/wiki/List_of_districts_of_Bangkok'
3 wikipable = read_html(url, attrs={"class":"wikitable"},header =0)
4 wiki_df = pd.DataFrame(wikipable[0])
5 wiki_df= wiki_df.rename(columns={'District(Khet)': 'District'})
6 wiki_df= wiki_df.rename(columns={'No. ofSubdistricts(Khwaeng)': 'No_of_Subdistricts'})
7 wiki_df= wiki_df.rename(columns={'Thai': 'DistrictThai'})
8 wiki_df['Latitude'].fillna(0, inplace=True)
9 wiki_df['Longitude'].fillna(0, inplace=True)
10 wiki_df.head(11)
```

# Data extraction and data cleansing 2

- This data provides information of Bangkok subdivided into 50 districts including district name, district code, latitude, longitude and also provide population of each districts.
- However, some records do not provide data as show in the figure below;

	District	Code	DistrictThai	Population	No_of_Subdistricts	Latitude	Longitude
0	Bang Bon	50	บางบอน	105161	4	0.0	0.0
19	Khan Na Yao	43	คันนายาว	88678	2	0.0	0.0
44	Thawi Watthana	48	ทวีวัฒนา	76351	2	0.0	0.0
46	Thung Khru	49	ทุ่งครุ	116473	2	0.0	0.0
47	Wang Thonglang	45	วังทองหลาง	114768	4	0.0	0.0



# Data extraction and data cleansing 3

- Python Geocoder package used to extract the latitude and longitude coordinates of these records  
After that, we will use Foursquare API to get the venue data for those districts

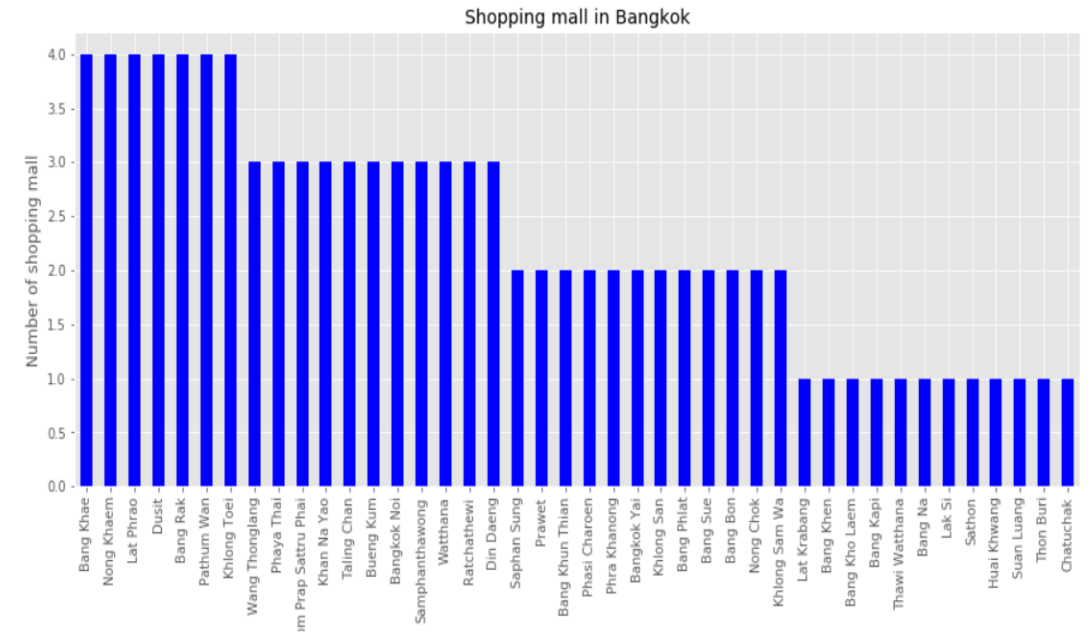
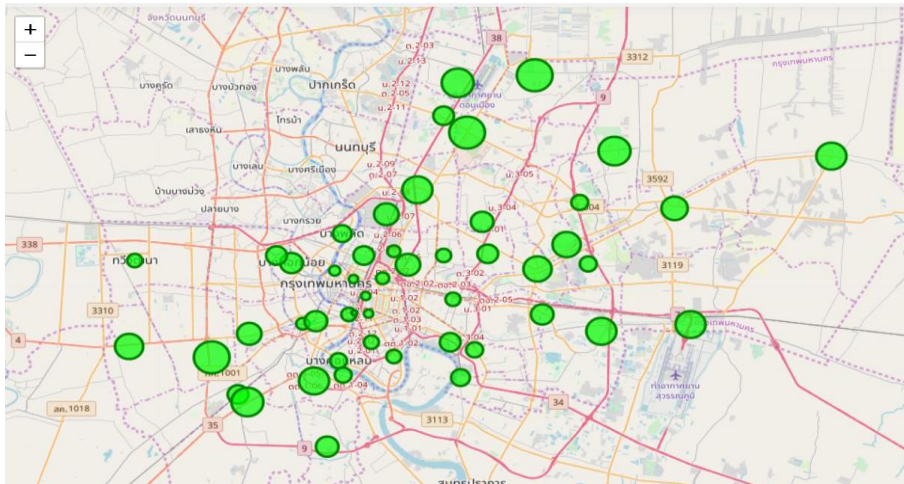


# FOURSQUARE

VenueName	VenueId	VenueLatitude	VenueLongitude	VenueCategory
ขนมจีนเทวดา บิ๊บบ้านสดๆ	559b4bfc498eb645065670b6	13.659428	100.433692	Noodle House
ไฟทองโภชนา	4ca8499a44a8224b303f1640	13.662101	100.435264	Thai Restaurant
UNIQLO (ยูนิโคล์)	528ec195498ec899efc22903	13.663285	100.439450	Clothing Store
Starbucks Reserve (สตาร์บัคส์ รีเสิร์ฟ)	4bc0040d461576b003b07932	13.663825	100.437668	Coffee Shop
Tops Market (ท็อปส์ มาร์เก็ต)	4bdeb42c0ee3a593b28631b0	13.662781	100.437410	Supermarket
อ๋องสเด็ก & เย็นตาโฟ	4bf382bc706e20a16256a898	13.670467	100.419624	Thai Restaurant

# Exploratory Data Analysis

- Map visualization, we found that people in Bangkok lived in suburb area as show in the figure below. Size of circle represent population in each district.

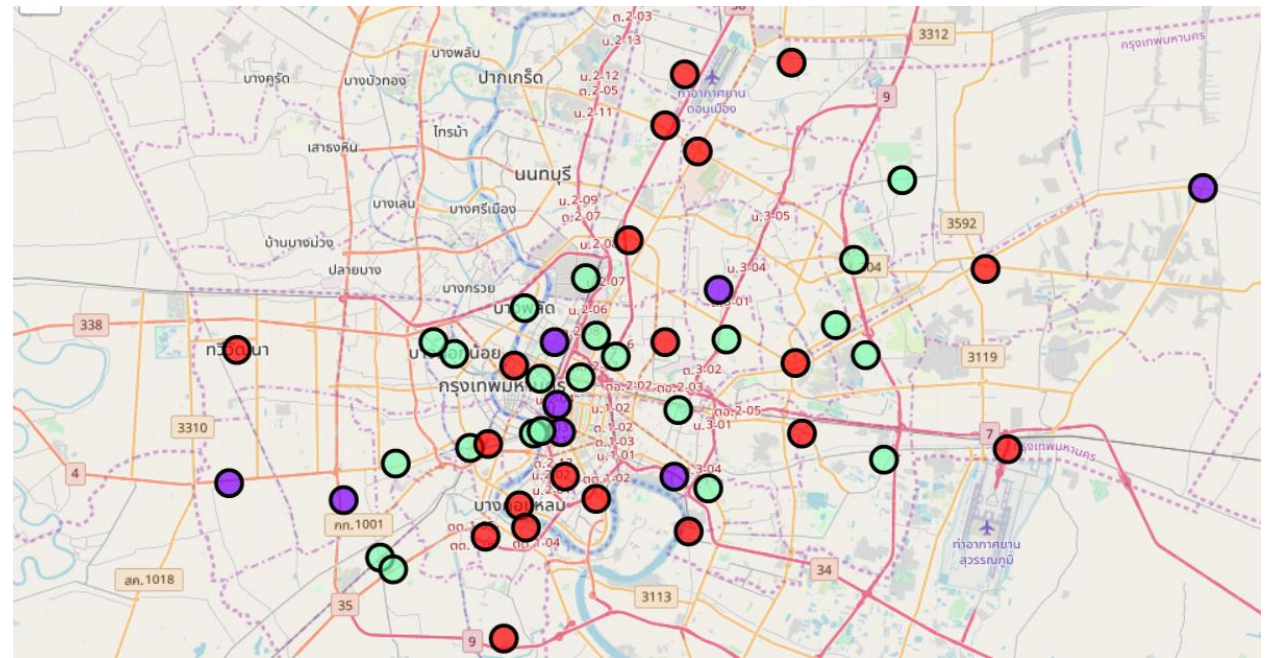


- The total number of shopping mall in each district shows in the bar chart

# Machine Learning (K-Mean)

- We will cluster the district into 3 clusters based on their frequency of occurrence for “Shopping Mall”. Before using K-Mean, we have to transform data using one-hot encoder. The results will allow us to identify which districts have a higher concentration of shopping malls while which districts have a fewer number of shopping malls by using k-means clustering. Based on the occurrence of shopping malls in a different district, it will help us to plan which district we should acquire land for developing new shopping mall.

```
1 # set number of clusters
2 kclusters = 3
3
4 df_clustering = df_mall.drop(["District"], 1)
5
6 # run k-means clustering
7 kmeans = KMeans(n_clusters=kclusters, random_state=0).fit(df_clustering)
8
```





# Results

- From the results, we found that the shopping malls were clustered using number of shopping malls in each district, so it can be classified to high, medium and low number fo shopping mall in each area. However, there are other factors which can be used to cluster such as number of tourists, income of people around that area, number of people traveled around these areas. These information can make clustering more accurate and will generate more information for decision making.

	District	Shopping Mall	Cluster	Code	DistrictThai	Population	No_of_Subdistricts	Latitude	Longitude
1	Bang Kapi	0.01	0	6	บางกะปิ	148465	2	13.765833	100.647778
3	Bang Khen	0.01	0	5	บางเขน	189539	2	13.873889	100.596389
4	Bang Kho Laem	0.01	0	31	บางคอแหลม	94956	3	13.693333	100.502500
6	Bang Na	0.01	0	47	บางนา	95912	2	13.680081	100.591800
13	Chatuchak	0.01	0	30	จตุจักร	160906	5	13.828611	100.559722

	District	Shopping Mall	Cluster	Code	DistrictThai	Population	No_of_Subdistricts	Latitude	Longitude
0	Bang Bon	0.02	2	50	บางบอน	105161	4	13.666503	100.428859
5	Bang Khun Thian	0.02	2	21	บางขุนเทียน	165491	2	13.660833	100.435833
7	Bang Phlat	0.02	2	25	บางพลัด	99273	4	13.793889	100.505000
9	Bang Sue	0.02	2	29	บางซื่อ	132234	2	13.809722	100.537222
10	Bangkok Noi	0.03	2	20	บางกอกน้อย	117793	5	13.770867	100.467933

	District	Shopping Mall	Cluster	Code	DistrictThai	Population	No_of_Subdistricts	Latitude	Longitude
2	Bang Khae	0.04	1	40	บางแค	191781	4	13.696111	100.409444
8	Bang Rak	0.04	1	4	บางรัก	45875	5	13.730833	100.524167
17	Dusit	0.04	1	2	ดุสิต	107655	5	13.776944	100.520556
22	Khlong Toei	0.04	1	33	คลองเตย	109041	3	13.708056	100.583889
25	Lat Phrao	0.04	1	38	ลาดพร้าว	122182	2	13.803611	100.607500

# Conclusion

- According to objective of this project we would like to support decision making on planning to develop new shopping mall in Bangkok, 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 shopping mall.



# Reference

- List of districts in Bangkok

[https://en.wikipedia.org/wiki/List\\_of\\_districts\\_of\\_Bangkok](https://en.wikipedia.org/wiki/List_of_districts_of_Bangkok)

- Foursquare Developers Documentation. Foursquare. Retrieved from

<https://developer.foursquare.com/docs>