

# Battle of the Neighborhoods

Coursera IBM Capstone Project

## 1. Introduction

Moscow is the economic centre with the high density of population and non-declining number of incoming labour from overall Russia as well as other countries appeals new businesses in both traditional and innovative sectors.

We majorly address problem to the local small or middle-size businesses that are concentrated on 'venue'-type businesses and where questions as following can arise:

1. Based on average income of the area, which city do we choose?
2. What is the situation in borrows and whom can we target?
3. What venues types are in the area and which ones are most popular?

... and so many more. Therefore, we to help answer questions that certainly arise when there exists a consideration of establishing a new venue to enrich strategical aims based on data.

What assumptions we should make? Probably the most convenient ones would be:

1. Feasible and availability pricing for commercial renting
2. Not oversaturated area with this type of business
3. Not necessary but we can look at the areas with higher house renting prices and assume that income there can be measured through house renting prices are representative
4. We still try to remain in a populated area
5. Closer to the metro is a plus

## 2. Data used

1. Forsquare API to get venues and details of a chosen city
2. As up-to-date granular data for income distribution in Moscow was not available, therefore we use average house renting prices per m<sup>2</sup> from Domofond.ru within districts
3. We use Wikipedia page on Moscow population data by Districts (wikipedia.com)
4. We use GIS-Lab website to get Moscow geodata (<https://gis-lab.info/qa/moscow-atd.html>)
5. When the venue type is chosen, we analyse the use sorted CIAN data on available commercial areas to rent (cian.ru)

## 3. Data cleansing

If we generalize, following operations were performed:

1. Data was scraped from multiple sources and cleaned for each source sequentially and further was combined into one table.
2. First and foremost, we used Domofond data on renting houses, cleaned and sorted it to the useful values and then mapped with the geospatial data we had for Moscow

3. We removed areas with n/a values and cleaned values with extra spacing and unnecessary symbols
4. We also selected value in a certain range of values for the pharmacy case and cleaned all unused data
5. We removed in some cells information that contained descriptions and maintained only values needed in those cells
6. We transferred some data values from strings to floats and vice versa to allow data processing and analysis

	Borough	District Name
0	Западный	Филёвский Парк
1	Зеленоградский	Матушкино
2	Западный	Внуково
3	Зеленоградский	Савёлки
4	Зеленоградский	Силино
...	...	...

	Borough	District Name	Rent	Latitude	Longitude	Area	Population	Population density
108	Южный	Зябликово	789	55.621900	37.742900	438	133278.0	30428.77
120	Восточный	Новокосино	765	55.740167	37.861725	360	107907.0	29974.17
11	Юго-Западный	Ломоносовский	976	55.678778	37.533239	334	88320.0	26443.11
52	Северный	Восточное Дегунино	780	55.880100	37.557600	377	98923.0	26239.52
69	Северо-Восточный	Бибирево	789	55.893800	37.611100	645	160163.0	24831.47
46	Северный	Бескудниковский	826	55.865900	37.553300	330	79603.0	24122.12
55	Юго-Западный	Зюзино	889	55.653100	37.598400	545	126815.0	23268.81
78	Северо-Восточный	Северное Медведково	833	55.888000	37.645500	566	127819.0	22582.86
111	Восточный	Новогиреево	799	55.748154	37.804108	445	98415.0	22115.73
76	Северо-Восточный	Южное Медведково	778	55.871000	37.638300	388	85698.0	22087.11
53	Северный	Савёловский	988	55.801600	37.564700	270	59287.0	21958.15
39	Юго-Западный	Коньково	921	55.643414	37.530588	718	156389.0	21781.20
101	Юго-Восточный	Марьино	782	55.652664	37.744774	1191	253943.0	21321.83
100	Южный	Орехово-Борисово Южное	913	55.604264	37.733132	694	147789.0	21295.24
64	Южный	Чертаново Северное	857	55.634300	37.603500	540	114548.0	21212.59

(2483, 7)							
	District Name	Latitude	Longitude	Venue	Venue Latitude	Venue Longitude	Venue Category
0	Филёвский Парк	55.75120	37.508200	RollerShop (Роллермар)	55.748076	37.513905	Board Shop
1	Филёвский Парк	55.75120	37.508200	HL   Holy Land	55.749521	37.502167	Cosmetics Shop
2	Филёвский Парк	55.75120	37.508200	Ели сацебели	55.755320	37.507718	Café
3	Филёвский Парк	55.75120	37.508200	Листок	55.753128	37.502470	Café
4	Филёвский Парк	55.75120	37.508200	НикОль	55.749222	37.501394	Cosmetics Shop
...	...	...	...	...	...	...	...
2478	Кунцево	55.74261	37.398482	Остановка «ЦКБ»	55.745189	37.395215	Bus Stop
2479	Кунцево	55.74261	37.398482	Парк	55.746161	37.399371	Park
2480	Кунцево	55.74261	37.398482	Синегория	55.740114	37.403173	Dance Studio
2481	Кунцево	55.74261	37.398482	Остановка «Стадион»	55.738309	37.397978	Bus Stop
2482	Кунцево	55.74261	37.398482	шиномонтаж "на колесах"	55.746240	37.394133	Auto Workshop

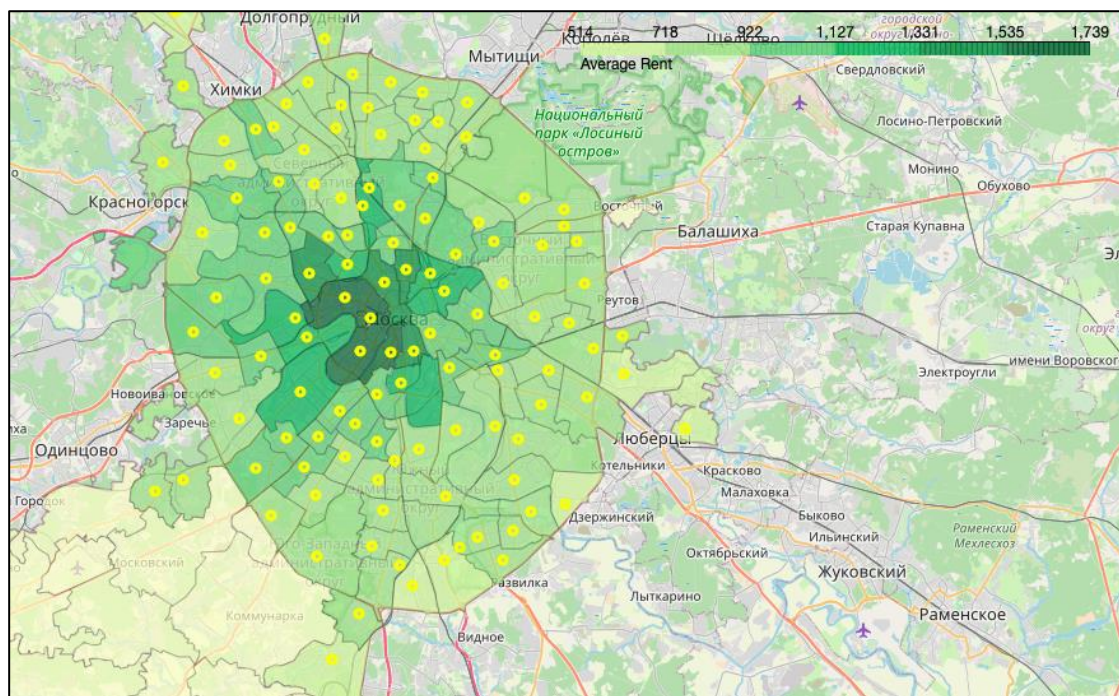
## 4. Methodology

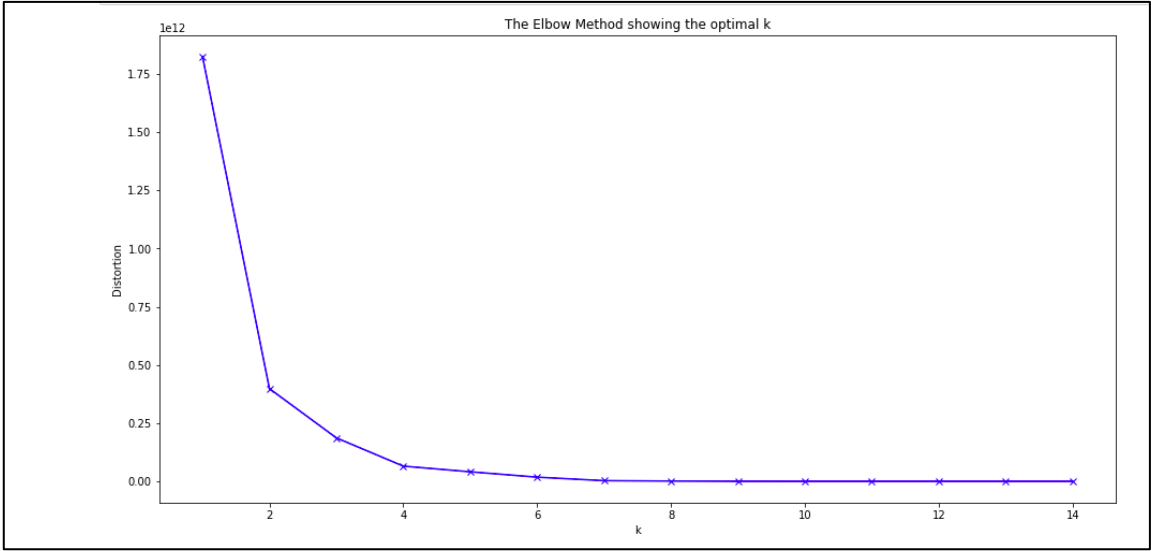
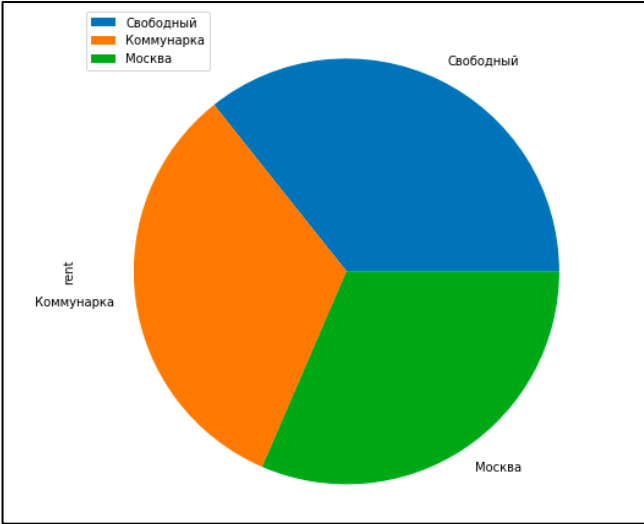
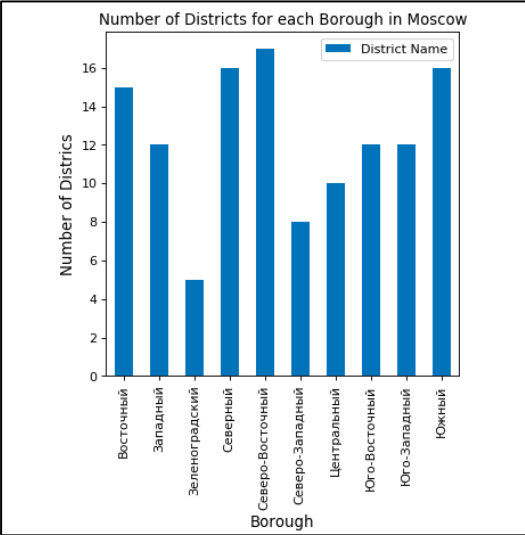
Idea is to get overall data on renting for living and analyse it as an alternative to the income data. Then we want to pick the most lucrative city for analysis. We concentrate on Moscow (3<sup>rd</sup> place) as the two highest seem to be ranked so due to some values in the average that upscaled the overall pricing. For the sake of better understanding we want check how many districts are there in each borough. We further proceed with the highest average renting price per borough.

Then we change the granularity to the level of districts to be more precise in our analysis and draw a heatmap with average renting prices per district. Then we explore population density to get a better understanding of how many customers we can serve

We get venue data and first look at the overall picture, e.g. what is present more frequently among. We proceed with clustering of these big groups.

At the end, we choose one particular venue based on analysis and own interest and add pricing of the commercial renting for the business building and analyse which location would be the optimal for opening a pharmacy based on given data and conditions. We cluster and build population density heatmap with these clusters. Optimal k for clustering is obtained via elbow-method.





	District Name	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
0	Академический	Coffee Shop	Pharmacy	Health Food Store	Sporting Goods Shop	Beer Store	Dance Studio	Park	Optical Shop
1	Алексеевский	Recording Studio	Supermarket	Pizza Place	Auto Workshop	Historic Site	Gym / Fitness Center	Mobile Phone Shop	Food & Drink Shop
2	Алтуфьевский	Supermarket	Lake	Gym / Fitness Center	Hotel	Pub	Dry Cleaner	Café	Grocery Store
3	Арбат	Coffee Shop	Hotel	Yoga Studio	Hostel	Bakery	Museum	Caucasian Restaurant	Gym / Fitness Center
4	Аэропорт	Coffee Shop	Salon / Barbershop	Martial Arts Dojo	Smoke Shop	Gourmet Shop	Clothing Store	Pet Store	Pharmacy
...	...	...	...	...	...	...	...	...	...
88	Южное Тушино	Pizza Place	Flower Shop	Sushi Restaurant	Convenience Store	Park	Pub	Coffee Shop	Food & Drink Shop
89	Южнопортовый	Film Studio	Electronics Store	Dance Studio	Boutique	Bakery	Baby Store	Fried Chicken Joint	Café
90	Якиманка	Coffee Shop	Bakery	Gym / Fitness Center	Shoe Store	Bridal Shop	Donut Shop	Pub	Mobile Phone Shop
91	Ярославский	Park	Fountain	Pizza Place	Shopping Mall	Café	Auto Workshop	Coffee Shop	Photography Studio
92	Ясенево	Fast Food Restaurant	Blini House	Pharmacy	Supermarket	Bookstore	Sporting Goods Shop	Burger Joint	Shopping Mall

93 rows x 11 columns

	District Name	Area	Population	Population density	Cluster Labels	1st Most Common Venue	2nd Most Common Venue	3rd Most Common Venue	4th Most Common Venue	5th Most Common Venue
0	Филевский Парк	962	94323.0	9804.89	1.0	Cosmetics Shop	Café	Insurance Office	Bus Stop	Asian Restaurant
4	Крюково	1049	100098.0	9542.23	1.0	Paintball Field	Scenic Lookout	Hockey Arena	Eastern European Restaurant	Exhibit
5	Щукино	769	111203.0	14460.73	1.0	Bus Stop	Gym / Fitness Center	Pizza Place	Stadium	Fountain
6	Нагатинский Затон	980	120954.0	12342.24	1.0	Gym / Fitness Center	Candy Store	Wine Shop	Food	Burger Joint
9	Тёплый Стан	750	134321.0	17909.47	1.0	Supermarket	Basketball Court	Health Food Store	Shopping Mall	Sushi Restaurant
...	...	...	...	...	...	...	...	...	...	...
116	Выхино-Жулебино	1497	225154.0	15040.35	1.0	Restaurant	Convenience Store	Grocery Store	Bar	Nightclub
117	Восточное Измайлово	385	78154.0	20299.74	1.0	Shawarma Place	Exhibit	Furniture / Home Store	Supermarket	Garden
119	Косино-Ухтомский	1505	82267.0	5466.25	1.0	Soccer Field	Harbor / Marina	Print Shop	Fish Market	Exhibit
120	Новокосино	360	107907.0	29974.17	1.0	Park	Health Food Store	Bookstore	Supermarket	Shopping Mall
122	Кунцево	1656	152364.0	9200.72	1.0	Spa	Bus Stop	Park	Dance Studio	Convenience Store

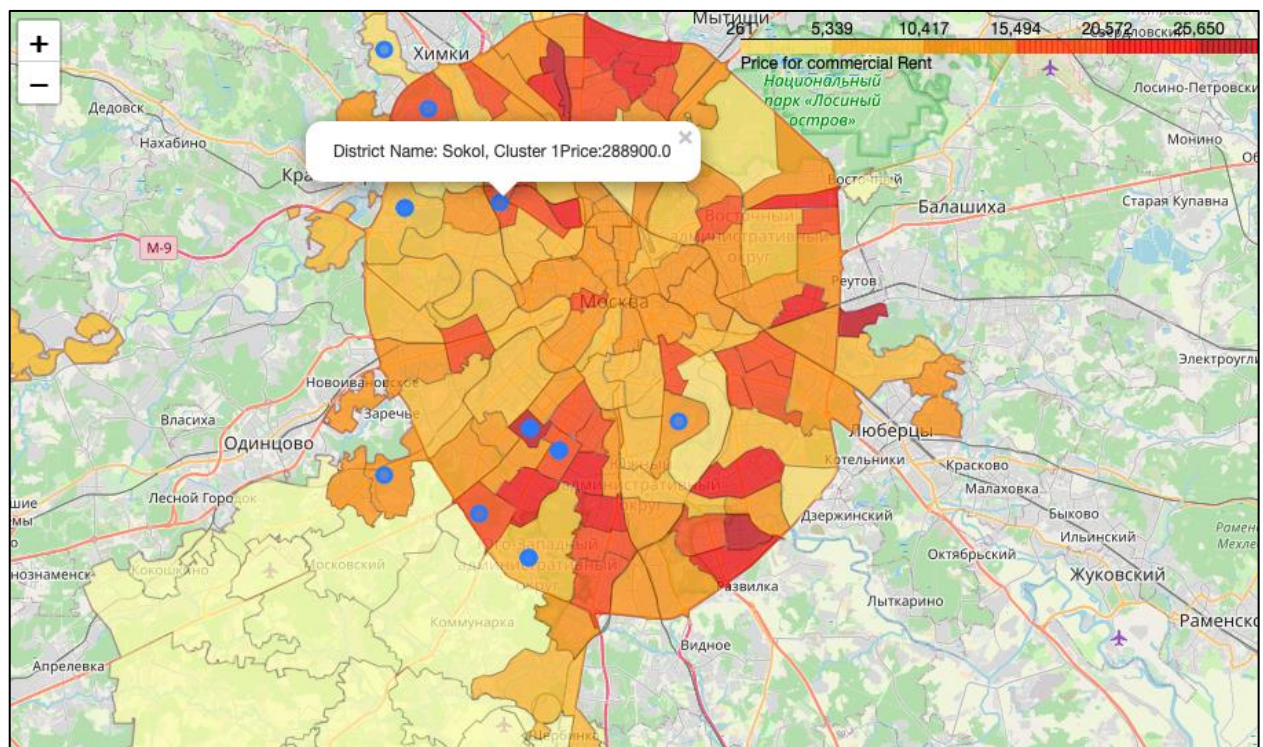
72 rows x 15 columns



## 4. Results

By analyzing data obtained we can see that the pharmacy is present in both clusters (**Cluster 0**: Shops&Malls&All for living, **Cluster 1**: Being active and socializing). However, based on given data market is not oversaturated with Pharmacies, therefore we target it further as the business to open.

Then, we obtain two clusters for pharmacies, **first - with expensive prices and second - with average prices on commercial renting**. We plot a map and look at overall Moscow data we have obtained. The best suit seems to be the district Sokol as it has **high population density, pricing is reasonable**, and **area of the commercial part is spacious**. Sokol is also the district where the **mean rent for living is slightly above average**, therefore we assume that ability to pay of customers is also there. District is **not oversaturated** with pharmacies and seems to be the perfect match for the new spot and the distance to the metro is only **3 min walk**.



## Cluster 1:

	District Name	Venue Latitude	Venue Longitude	Area of Com. Rent	Metro	Price	Borough	Rent	Area	Population	Population density
5	Тверской	55.773321	37.598090	106.0	м. Лубянка (4 мин пешком)	1100001.0	Центральный	1500	727	77864.0	10710.32
7	Тверской	55.773321	37.598090	114.0	м. Китай-город (7 мин пешком)	848999.0	Центральный	1500	727	77864.0	10710.32
10	Тверской	55.776630	37.607171	106.0	м. Лубянка (4 мин пешком)	1100001.0	Центральный	1500	727	77864.0	10710.32
12	Тверской	55.776630	37.607171	114.0	м. Китай-город (7 мин пешком)	848999.0	Центральный	1500	727	77864.0	10710.32
14	Якиманка	55.732885	37.613876	78.0	м. Новокузнецкая (2 мин пешком)	1250001.0	Центральный	1562	480	27672.0	5765.00
15	Якиманка	55.732885	37.613876	145.0	м. Третьяковская (7 мин пешком)	899000.0	Центральный	1562	480	27672.0	5765.00

## Cluster 2

	District Name	Venue Latitude	Venue Longitude	Area of Com. Rent	Metro	Price	Borough	Rent	Area	Population	Population density
0	Нагатинский Затон	55.683017	37.681513	85.0	м. Коломенская (3 мин пешком)	450005.0	Южный	918	980	120954.0	12342.24
1	Северное Тушино	55.859719	37.432788	65.0	м. Первомайская (14 мин пешком)	124995.0	Северо-Западный	827	940	165762.0	17634.26
2	Сокол	55.803670	37.509224	81.0	м. Войковская (3 мин на машине)	149000.0	Северный	952	372	59507.0	15996.51
3	Дорогомилово	55.737393	37.524963	100.0	м. Кутузовская (5 мин на машине)	449000.0	Западный	1172	795	76093.0	9571.45
4	Тверской	55.773321	37.598090	72.0	м. Белорусская (2 мин пешком)	480000.0	Центральный	1500	727	77864.0	10710.32
6	Тверской	55.773321	37.598090	86.0	м. Белорусская (2 мин пешком)	599899.0	Центральный	1500	727	77864.0	10710.32
8	Тверской	55.773321	37.598090	133.0	м. Менделеевская (2 мин пешком)	288900.0	Центральный	1500	727	77864.0	10710.32
9	Тверской	55.776630	37.607171	72.0	м. Белорусская (2 мин пешком)	480000.0	Центральный	1500	727	77864.0	10710.32
11	Тверской	55.776630	37.607171	86.0	м. Белорусская (2 мин пешком)	599899.0	Центральный	1500	727	77864.0	10710.32
13	Тверской	55.776630	37.607171	133.0	м. Менделеевская (2 мин пешком)	288900.0	Центральный	1500	727	77864.0	10710.32

## **5. Discussion**

First and foremost obtained results are highly dependent of the quality and completeness of data which in some cases is under a big question mark. Although we obtained a neat and nice result, it could have been better whether we used dynamic data to show the trends in, for instance, changes in rents and therefore, could have made a better prediction. Moreover, whether we had access to fiscal data of the pharmacies we could have done a more rigorous competitor analysis and may have developed some set of better actions.

The granularity of districts although is not at an extremely high level but is still not enough to make a precise analysis. One could have used rankings for the pharmacies to obtain understanding of performances nearby

K-means is although a good clustering algorithm but there exist others that can help to obtain results. For example, Mean shift algorithm that finds and adapts centroids based on the density of examples in the feature space.

## **6. Conclusion**

Overall this project overlooked at the initial steps to be undertaken for conservative local businesses to halt usage of purely the 'guts' feeling but also to utilize the technologies to enforce a better decision making based on data which is sooner or later will become inevitable part of each business