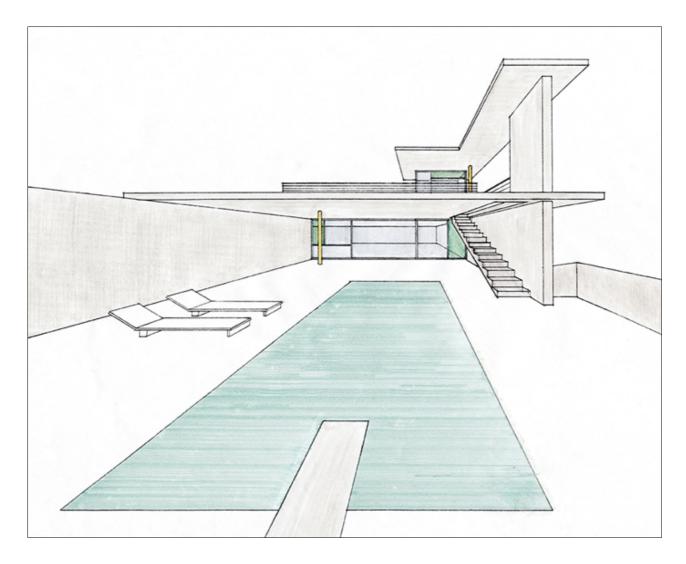# Intelligent new home finder

Prepared for: coursa IBM data science professional certificate Capstone project

Prepared by: Caifang Cai

8 May 2020

# INTRODUCTION

### Objective

The main objective of this project is to provided an intelligent new home location finder for families relocated for a new job. This is based on the fact that the client family has already have a reference home and looking for a similar place to install the family. The criteria taken into account including:

1. Transport (bus, subway, local trains);
2. Closeness to commerce (bakery, supermarkets, etc. );
3. Nearby culture environment. Such as museums, galleries etc.
4. Public parks and sport centres.

The maximum distance should be limited in a radius of 12km to the new working place.

This major goal is to provide an intelligent  tool to take into account of the above 4 major criteria to find a new home location similar to a reference home address. I give the top three recommendations locations for a new home location.

### Stat of the art

Relocating home to unfamiliar region is indeed a painful  task. It takes too much humain energy.

So far, when we need to relocate the home following a new job offre.  The first step we need to do is to check on internets for days and days or to ask friends for all kinds of information. The second step is to process all these informations and make a 'not bad' choice barely based on humain juge. Both steps are not reliable which often leads to a 'not good' choice.

### Interests

This project interests all people who need to relocate a family or to find a new home location in an unfamiliar region. The only condition to use is to have a reference home address. This reference home address can be your current home address that you are satisfied  with or any address that you know well to be a good location for your home.

# DATA ACQUISITION AND CLEANING

Table 1: office and reference home addresses used in homefinder

| | |
|---|---|
| **Office address** | 100 Avenue de Paris, 91344 Massy, France |
| **Reference home address** | 102 Boulevard Richard Wallace, 92800 Puteaux |

## Dara source and demo example

We will use Foursquare data as our source data, an example in Paris region as our final result for demonstration. The reference home address and the new office address around which we are looking to install a family are given as follows in Table 1.

More specifically, we will use Foursquare explore end-point to explore the areas of concern.

## Data cleaning

The raw data obtained from Foursquare is a json file which basically as following:

```
{'meta': {'code': 200, 'requestId': '5ef6780030567d545e80ec19'},
 'response': {'suggestedFilters': {'header': 'Tap to show:',
   'filters': [{'name': 'Open now', 'key': 'openNow'}]},
  'headerLocation': 'Issy-les-Moulineaux',
  'headerFullLocation': 'Issy-les-Moulineaux, Paris',
  'headerLocationGranularity': 'neighborhood',
  'totalResults': 73,
  'suggestedBounds': {'ne': {'lat': 48.83141760900001,
    'lng': 2.271401296161224},
   'sw': {'lat': 48.81341759099999, 'lng': 2.244113103838776}},
  'groups': [{'type': 'Recommended Places',
    'name': 'recommended',
    'items': [{'reasons': {'count': 0,
      'items': [{'summary': 'This spot is popular',
        'type': 'general',
        'reasonName': 'globalInteractionReason'}]},
     'venue': {'id': '541485c1498efe97bd3efa43',
      'name': 'La Passerelle',
      'location': {'address': '172 quai de Stalingrad',
       'lat': 48.82526662415111,
       'lng': 2.2575543895488503,
       'labeledLatLngs': [{'label': 'display',
         'lat': 48.82526662415111,
         'lng': 2.2575543895488503}],
       'distance': 317,
       'postalCode': '92130',
       'cc': 'FR',
       'city': 'Issy-les-Moulineaux',
       'state': 'Île-de-France',
          'country': 'France',}
```

The first we need to do is the extract the useful information related to our home
relocating later.

Here is a list of extracted values used later for recommendations of new home location.

['venue.categories','venue.location.distance','venue.location.formattedAddress','venue.
location.lat','venue.location.lng','venue.name']

After extraction, the data is not real for analyse yet. Here is a short example

| | venue.cate gories | venue.location.dist ance | venue.location.format tedAddress | venue.locat ion.lat | venue.locat ion.lng | venue.n ame |
|---|---|---|---|---|---|---|
| 0 | [{'id': '52e81612bcbc57f1 066b79f9', 'name': 'M... | 317 | [172 quai de Stalingrad, 92130 Issy-les-Moulin... | 48.825267 | 2.25755 4 | La Passer elle |
| 1 | [{'id': '4bf58dd8d48988d1 63941735', 'name': 'P... | 454 | [170 quai de Stalingrad, 92130 Issy-les-Moulin... | 48.826249 | 2.25563 0 | Île Saint-Germa in |
| 2 | [{'id': '4bf58dd8d48988d1 0c941735', 'name': 'F... | 434 | [113 bis avenue de Verdun, 92130 Issy-les-Moul... | 48.818860 | 2.25530 4 | Issy Guing uette |
| 3 | [{'id': '4bf58dd8d48988d1 75941735', 'name': 'G... | 374 | [1-6 boulevard Garibaldi, 92130 Issy-les-Mouli... | 48.823489 | 2.26260 4 | MurMu r |

The contents are too long and can not be interpreted directly by humain. I then filtered the contents to get only
the quantitive numbers and short, readable informations. Here gives an example after data cleaning.

| | category | distance | address | latitude | longitude | name | |
|---|---|---|---|---|---|---|---|
| 0 | Modern European Restaurant | 317 | 172 quai de Stalingrad, 92130 Issy-les-Mouline... | 48.825267 | 2.257554 | La Passerelle | |
| 1 | Park | 454 | 170 quai de Stalingrad, 92130 Issy-les-Mouline... | 48.826249 | 2.255630 | Île Saint-Germain | |
| 2 | French Restaurant | 434 | 113 bis avenue de Verdun, 92130 Issy-les-Mouli... | 48.818860 | 2.255304 | Issy Guinguette | |
| 3 | Gym / Fitness Center | 374 | 1-6 boulevard Garibaldi, 92130 Issy-les-Moulin... | 48.823489 | 2.262604 | MurMur | |

## Feature preparation

I need to convert all these data into features that can be used for quantitative analysis. One important thing is to keep these features reflecting the four criteria that we discussed in the introduction. These four criteria are the rules for us to carry out further recommendations. I classify all the items into four categories commerce, transport, culture and environment (park and sport centres ). I then count the numbers items and distances to office in each category. In each category, I use two values to describe it. One is the mean distance to home. The other is the density of numbers. In order to avoid influences of very small and very large numbers. I use a category value to describe the density of number. For example, for transport, the *num_transport* will have values of None, Few and Many. At the end, a concerned address will has the following shape of the feature vector.

| dis_commerce | num_commerce | dis_culture | num_culture | dis_transport | num_transport | dis_environement | num_environement | |
|---|---|---|---|---|---|---|---|---|
| 0.67481 | high | 0.659 | median | 0.7702 | many | 0.5485 | low | |

# NEW HOME RECOMMENDATION METHOD