## Task:

to explore something interesting in neighborhood description.

## Idea:

Given an example below, I first explore what information tends to be included in the neighborhood description.

"Very safe and quiet residential neighborhood, you can walk around anytime day or night. Great access to **Harvard Square** and all it has to offer - **Gourmet restaurants**, **bars & pubs**, **cultural events**, **historical sites**, shopping, entertainment, etc. Easy access to **MIT**, **downtown Boston**, **area hospitals**, all in about 12 - 15 minutes with Uber or Lyft "

It can be seen that several locations (marked bold) are referred in the sentences. The locations might attract readers most because those words are what they care about. We should include information with reference to these locations if another description is auto generated by our system which describes the nearby property in this neighborhood.

After we extract the locations, we also interest in the descriptive words for this location. For example,

"Very safe and quiet residential neighborhood, you can walk around anytime day or night.

Great access to Harvard Square and all it has to offer - Gourmet restaurants, bars & pubs, cultural events, historical sites, shopping, entertainment, etc. Easy access to MIT, downtown Boston, area hospitals, all in about 12 - 15 minutes with Uber or Lyft"

Those words marked wave underline are to describe the locations. We should extract them too. Combined locations with descriptive words, the basic element of a sentence can be generated.

Take Boston for example. Boston has several neighborhoods, e.g. Allston, Brighton, Chinatown. I am going to extract popular locations for every neighborhood in Boston and compare the frequency of word occurrence. I wrote four python source files.

The python library I used is Spacy. (<a href="https://spacy.io/docs/">https://spacy.io/docs/</a>) Some features of this library might help my work:

- Tokenization
- pos tagging
- sentence segmentation
- dependency parsing
- entity recognition
- word vector (Glove)
- sentiment analysis

The function of 4 source files:

generateData.py

paras: [city]

generate a dict (python structure) which contains lists of neighborhood-related descriptions stored in neighborhood keys. Write data of one neighborhood into one file.

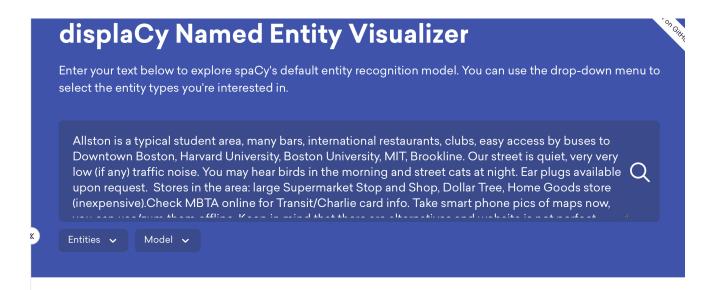
findNameEntites.py

paras: [city] [neighborhood]

generate keywords of one neighborhood of one city. Choose NameEntity as location-related

features.

A good visualizer: https://demos.explosion.ai/displacy-ent/



Allston ORG is a typical student area, many bars, international restaurants, clubs, easy access by buses to Downtown Boston ORG, Harvard University ORG, Boston University ORG, MIT ORG, Brookline OPE. Our street is quiet, very very low (if any) traffic noise. You may hear birds in the morning and street cats at night. Ear plugs available upon request. Stores in the area: large Supermarket Stop and Shop, Dollar Tree, Home Goods ORG store (inexpensive). Check MBTA online for Transit/Charlie card info. Take smart phone pics of maps now, you can use/zum them offline. Keep in mind that there are alternatives and website is not perfect.

findAdj.py

paras: [city] [neighbourhood] [keyPlace]

find descriptive words for a specific key place in one neighborhood of one city. I expect to get adj. words, e.g. far, near, but it does not perform as I expected. Still under construction.

findKeywordsInCity.py

paras: [city]

find all the location-related keywords in one city. Draw a heat map for visualization which shows how frequent a location appears in the description of one neighborhood in one city.

## What I can get from the data visualization

- The heat map is sparse. only several locations are frequently appear in all neighborhoods.
- Some locations only appear in one neighborhood. It means that that location must be described in this neighborhood.
- If one location appears in several neighborhoods, it is suggested that those neighborhoods might be closed to each other.
- City name is always the most frequent word.
- It seldom exists many locations which all appear very frequently in one neighborhood. Seems good because we don't have to struggle in selection.

|        | \ <i>'</i> ' | 1.5      |
|--------|--------------|----------|
| I Nata | Merial       | lization |
| Dala   | visuai       | nzanon   |

see below

