# Applied Data Science Capstone – The Battle of Neighborhoods I

## The housing price at Great Toronto Area

By Daniel Chen

➢ **Motivation**

With the growth of population at Canada, demands of houses has been growing. By doing the research of housing price at great Toronto area will help people who want to invest on estate or buy for living. For example, if there is an investor wants to do the investment for buying a house at Toronto; it is crucial to having the information about what kinds of venue around it. In advance research, the investor can do the prediction of pricing on the specific area. The other advantage benefits to people like me; plan to buy a house that contents personal living requirements. By doing the research on neighborhood at Toronto is helping me to decide where will be the best fit for me to live. In addition, including the budget can assist me to choose the house is affordable and the area is suitable to me.

➢ **Data**

Housing prices – *Zoocasa*

- Scraped by Beautifulsoup

Location – *Nominatim of GeoPy*

Venues – *Foursquare*

- ➤ **Code** – web scraping housing prices from zoocasa
  - ▪ Import libraries

```python
1  import os
2  from bs4 import BeautifulSoup
3  from selenium import webdriver
4  from selenium.webdriver.common.keys import Keys
5  from selenium.webdriver.common.by import By
6
7  from geopy.geocoders import Nominatim
8
9  import time
10 import sys
11 import numpy as np
12 import pandas as pd
13 import regex as re
```

  - ▪ Applying Chrome drive

```python
1  chromedriver = "" # path to the chromedriver executable
2  chromedriver = os.path.expanduser(chromedriver)
3  print('chromedriver path: {}'.format(chromedriver))
4  sys.path.append(chromedriver)
5  driver = webdriver.Chrome(chromedriver)
```

  - ▪ Scraping

```python
1  house_infos = []
2  zoocasa_url = "https://www.zoocasa.com/search?bedrooms=3&latitude=43.72400463181717&longitude=-79.39055077202798&zoom=12"
3  pages = 22
4  for i in range(pages):
5      print("It is scriping at " + i+1 + " pages...")
6      driver.get(zoocasa_url)
7      listing_houses = soup.find_all("div", class_="card-wrapper")
8
9      for house in listing_houses:
10
11         house_address = house.find_all("div", {"class": "street"})
12         house_details = house.find_all("div", {"class": "details"})
13         house_price = house.find_all("div", {"class": "price"})
14
15         house_info = [house_address, house_details, house_price]
16         house_infos.append(list(map(lambda x: [info.getText() for info in x][0].strip("\n"),house_info)))
17      time.sleep(5)
18      if i < pages-1:
19          zoocasa_url = "https://www.zoocasa.com" + soup.find_all("a", class_="icon-arrow-right-open active")[0]['href']
```

- ➢ **Making Dataframe** – *Including Address, Home Type, Price, Latitude, and Longitude*
  - ▪ Dataframe from Zoocasa – House Postal code, Address, type, and price

```
1  df = pd.DataFrame(data = house_infos)
2  df.columns = ['Address', 'Home Type', 'Price']
3
4  #Make price value to be integer
5  df["Price"] = [str(x[2].replace("$","").replace("\n","").replace(",", "").replace(" ","")) for x in house_infos]
6
7  #Make adress to be available searching location on Nominatim
8  addr_series = df["Address"].copy()
9  for idx, address in enumerate(address_series_latlon):
10     if "-" in address:
11         minusidx = address.index('-')
12         addr_series.set_value(idx, addr_series[idx][minusidx+1::])
13
14 #Getting postal code from geocoder.canadapost
15 postalcodes = []
16 for address in addr_series:
17     g = geocoder.canadapost(address + ", Toronto")
18     postalcodes.append(g.postal[0:3])
19 df.insert(0, "Postal Code", postalcodes)
```

  - ▪ Determine the postal code belongs to Toronto

```
1  !wget -q -O 'Geospatial_Coordinates.csv' http://cocl.us/Geospatial_data
2  unique_postal = pd.DataFrame({'Postal Code':df["Postal Code"].unique()})
3
4  lat_lng_coords = pd.read_csv('Geospatial_Coordinates.csv')
5  #Delete the postal code not belong to Toronto
6  postal_notToronto = []
7  for postalcode in unique_postal['Postal Code']:
8      if postalcode not in lat_lng_coords['Postal Code'].tolist():
9          unique_postal = unique_postal[unique_postal['Postal Code'] != postalcode].reset_index()
10         postal_notToronto.append(postalcode)
11 unique_postal
```

| | index | Postal Code |
|---|---|---|
| 0 | 0 | M4J |
| 1 | 1 | M1K |
| 2 | 2 | M1P |
| 3 | 3 | M5P |
| 4 | 4 | M4L |
| 5 | 5 | M5J |
| 6 | 6 | M5B |
| 7 | 7 | M1J |
| 8 | 8 | M5T |
| 9 | 9 | M6P |
| 10 | 10 | M1L |
| 11 | 11 | M5S |
| 12 | 12 | M4R |
| 13 | 13 | M6H |
| 14 | 14 | M5V |
| 15 | 16 | M4A |

  - ▪ Clean the data of house not in Toronto

```
1  for postalcode in postal_notToronto:
2      df = df[df['Postal Code'] != postalcode].reset_index()
```

- Dataframe from Nominatim – House Coordinate (latitude and longitude)

```
1  geolocator = Nominatim()
2  latlon = []
3  for address in addr_series:
4      location = geolocator.geocode(address + ", Toronto")
5      latlon.append([location.latitude,location.longitude])
6  lat_lon_df = pd.DataFrame(data = latlon)
7  lat_lon_df.columns = ["Latitude", "Longitude"]
```

➢ **Foursquare** – Nearby Venues Cataloged by Postal Codes
- Postal Code Coordinates

```
1  lag_lng = pd.Series(map(lambda x: lat_lng_coords[['Latitude', 'Longitude']][lat_lng_coords['Postal Code'] == x].values.tolis
2  lag_lng_df = pd.DataFrame(columns=['Latitude', 'Longitude'], data = list(lag_lng))
3  lag_lng_df = pd.concat([unique_postal, lag_lng_df], axis=1, sort=False)[['Postal Code', 'Latitude', 'Longitude']]
4  lag_lng_df.head()
```

|   | Postal Code | Latitude | Longitude |
|---|---|---|---|
| 0 | M4J | 43.685347 | -79.338106 |
| 1 | M1K | 43.727929 | -79.262029 |
| 2 | M1P | 43.757410 | -79.273304 |
| 3 | M5P | 43.696948 | -79.411307 |
| 4 | M4L | 43.668999 | -79.315572 |

- Venues Dataframe

```
1  #Foursquare ID and Scecretc
2  CLIENT_ID = '5JKMDJDBIRGHDRTH0AX4XEEMAAMKNROAHHEUBA014MYU2C0J' # your Foursquare ID
3  CLIENT_SECRET = 'E45TSRJVHBBEQENESBV2A0K1PPEOZZSAE3KJX0C4XTJ03VAT' # your Foursquare Secret
4  VERSION = '20180605' # Foursquare API version
5
6  LIMIT = 150
7
8  Toronto_venues = getNearbyVenues(names=unique_postal['Postal Code'],
9                                   latitudes=lat_lng_df['Latitude'],
10                                  longitudes=lat_lng_df['Longitude'],
11                                  CLIENT_ID=CLIENT_ID,
12                                  CLIENT_SECRET=CLIENT_SECRET,
13                                  VERSION=VERSION,
14                                  LIMIT=LIMIT
15                                  )
16 Toronto_venues.to_csv("Toronto_Venues")
```

|   | Neighborhood | Neighborhood Latitude | Neighborhood Longitude | Venue | Venue Latitude | Venue Longitude | Venue Category |
|---|---|---|---|---|---|---|---|
| 0 | M4J | 43.685347 | -79.338106 | The Only Cafe | 43.680409 | -79.337898 | Beer Bar |
| 1 | M4J | 43.685347 | -79.338106 | efes bar&grill | 43.680242 | -79.338346 | Gastropub |
| 2 | M4J | 43.685347 | -79.338106 | Serano Bakery | 43.683139 | -79.346531 | Bakery |
| 3 | M4J | 43.685347 | -79.338106 | Mr. Pide | 43.679635 | -79.341530 | Turkish Restaurant |
| 4 | M4J | 43.685347 | -79.338106 | Sakawa Coffee | 43.679906 | -79.339807 | Café |