## *Toronto Housing Market*

Housing market in Toronto is hot. During this pandemic, I have heard demand for housing has been increased and prices has gone up, so it's not a bad idea to check this. I was not able to obtain recent data and only data I was able to get was from Kaggle which is provided to public by Toronto Real Estate Board. I created this analysis so once TREBB release new dataset I have my code ready.

```
1
```

## Set up - Imports and Reading

```
1   # Importing Libraries
2
3   import pandas as pd
4   import numpy as np
5
6   import matplotlib.pyplot as plt
7   %matplotlib inline
8   import seaborn as sns
9   sns.set(color_codes=True)
10  import itertools
11  import folium
12  from folium import Map
```

```
1   # Reading database
```

```
2  database=pd.read_csv('properties.csv')
3  database.head()
4
```

|   | Unnamed: 0 | Address | AreaName | Price ($) | lat | lng |
|---|---|---|---|---|---|---|
| **0** | 0 | 86 Waterford Dr Toronto, ON | Richview | 999888 | 43.679882 | -79.544266 |
| **1** | 1 | #80 - 100 BEDDOE DR Hamilton, ON | Chedoke Park B | 399900 | 43.250000 | -79.904396 |
| **2** | 2 | 213 Bowman Street Hamilton, ON | Ainslie Wood East | 479000 | 43.251690 | -79.919357 |
| **3** | 3 | 102 NEIL Avenue Hamilton, ON | Greenford | 285900 | 43.227161 | -79.767403 |
| **4** | 6 | #1409 - 230 King St Toronto, ON | Downtown | 362000 | 43.651478 | -79.368118 |

## Clean up

```
1  # Cleaning data with unacceptable latitude and longitude values, and reasonable lower price
2  database=database.drop(database[abs(database.lat>90)].index)
3  database=database.drop(database[abs(database.lng>180)].index)
4  database=database.dropna()
5  database=database.drop(database[database['Price ($)']<100000].index)
```

```
1  print(database.size)
2  print(database.shape)
3  database['Price ($)'].max()
```
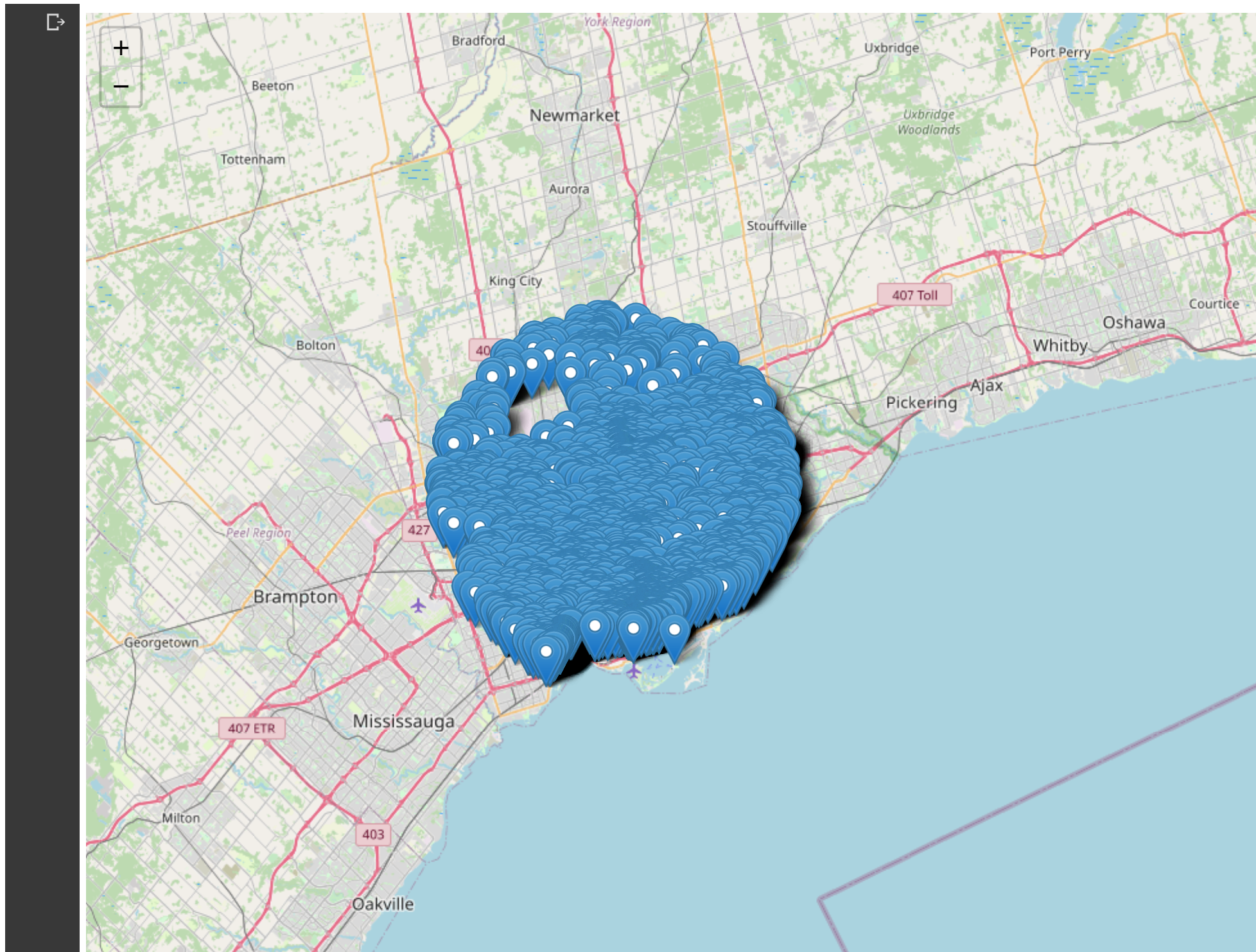
```
133716
(22286, 6)
32500000
```

The dataset seems to have areaname, but I don't know the all the areas in Greater Toronto Area, so I decided to clean it up with lat and lng information

```
1    # Finding distance from the centre of Toronto. I pick Allen road and eEnglinton as the centre
2    # and I am going to find 40km range away from this point. This point coordinate
3    # are : lat (43.73976) and long (-79.42126)
4
5
6    # here a distance finder function has been created:
7
8    from geopy.distance import geodesic
9    def distance_finder (lat, lng):
10     centre = (43.73976, -79.42126)
11     point= (lat,lng)
12     distance= geodesic(centre, point).km
13     return distance
14
15   # lambda operation on Database
16   database['Distance']=database.apply(lambda x: distance_finder(x['lat'],x['lng']), axis=1)
17   # df is the area of my interest
18   df=database[database['Distance']<15]
19
20
```

```
1    # since the data set includes all the transaction in Ontario, I would like to check there is no record outside of the
2    # region of my interest
3
4    locations = df[['lat', 'lng']]
5    locationlist = locations.values.tolist()
6    len(locationlist)
7
8
9    map = folium.Map(location=[43.73976, -79.42126], zoom_start=10)
10   for i in range(0, len(locationlist)):
11       folium.Marker(locationlist[i]).add_to(map)
12   map
13
14
15
```
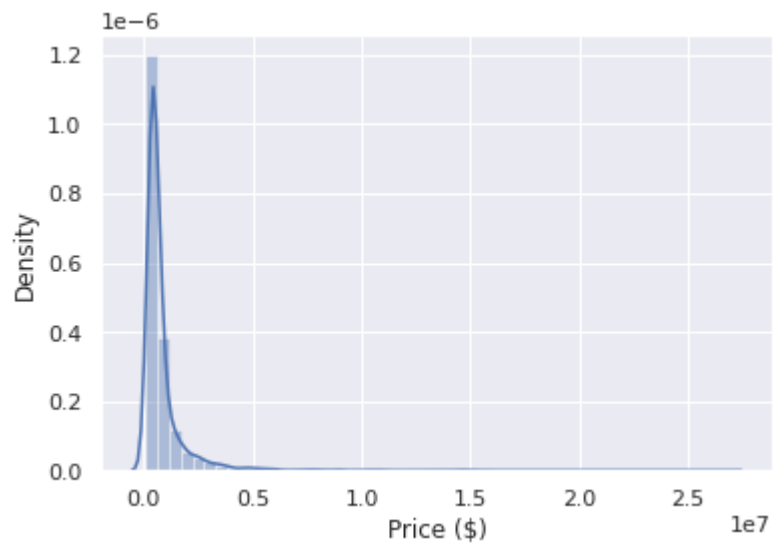
## ▾ Analysis

```
1  #quick sns distribution graph
2  ax = sns.distplot(df['Price ($)'])
3  # Show data is skewed
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated functio
  warnings.warn(msg, FutureWarning)
```

```
1   #checking dataset
2   df.describe()
```

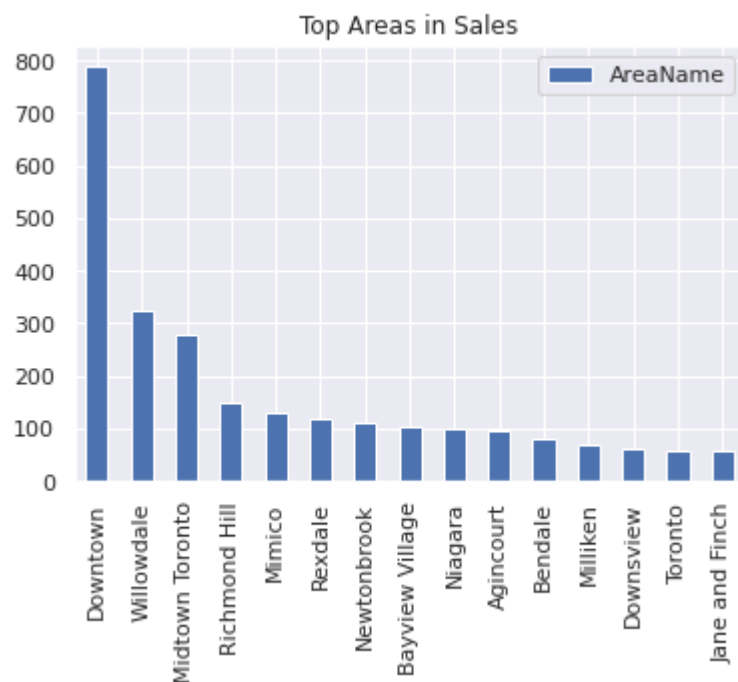|  | Unnamed: 0 | Price ($) | lat | lng | Distance |
|---|---|---|---|---|---|
| count | 4590.000000 | 4.590000e+03 | 4590.000000 | 4590.000000 | 4590.000000 |
| mean | 39078.920915 | 8.258964e+05 | 43.715426 | -79.405040 | 9.271488 |
| std | 29697.845972 | 1.179302e+06 | 0.063987 | 0.077623 | 3.561755 |
| min | 0.000000 | 1.000000e+05 | 43.615078 | -79.605835 | 0.229304 |
| 25% | 6175.250000 | 3.499000e+05 | 43.656661 | -79.447968 | 6.692606 |
| 50% | 36149.500000 | 4.990000e+05 | 43.708008 | -79.399605 | 10.194123 |
| 75% | 74069.750000 | 7.990000e+05 | 43.767937 | -79.369514 | 11.833238 |
| max | 120600.000000 | 2.680000e+07 | 43.874107 | -79.235939 | 14.988489 |

```
1   # I am focusing near price mean region.
2
3
4   df_1=df[((df['Price ($)']<5.299000e+05) & (df['Price ($)']>2.5e5))]
5   ax = sns.distplot(df_1['Price ($)'])
```

```
/usr/local/lib/python3.7/dist-packages/seaborn/distributions.py:2557: FutureWarning: `distplot` is a deprecated functi
  warnings.warn(msg, FutureWarning)
```



```
1    # I would like to now what are the hottest neighbourhoods in sale.
2
3    Area= pd.DataFrame(df['AreaName'].value_counts())
4    Area.sort_values(by='AreaName')
5    Area[0:15].plot(kind='bar', title='Top Areas in Sales')
```

<matplotlib.axes._subplots.AxesSubplot at 0x7f158a658a50>
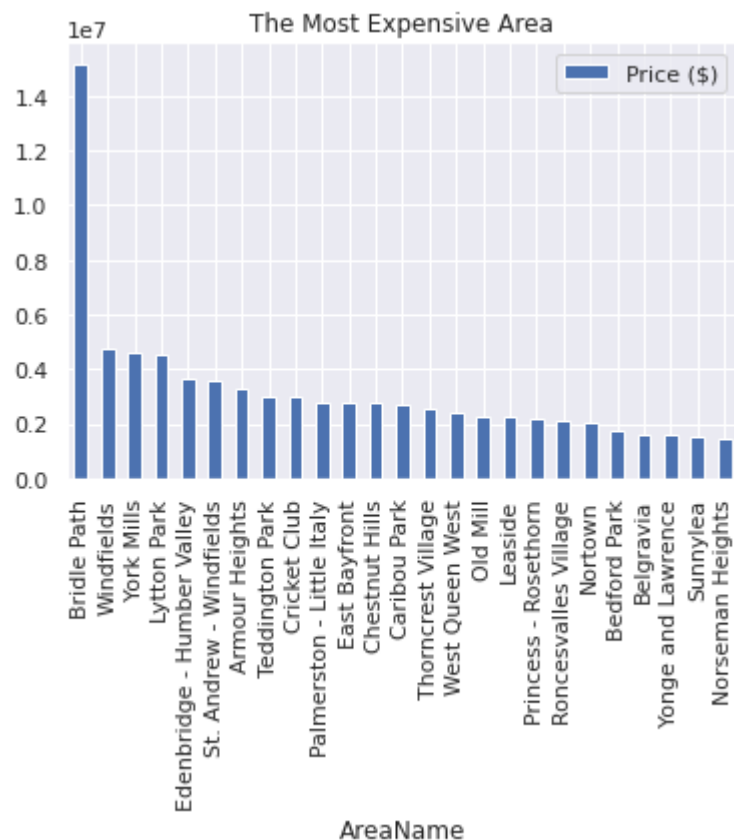


```
1    # I would like to know what are the most expensive neighbourhoods.
2
3    AreaPrice=pd.DataFrame(df.groupby('AreaName')['Price ($)'].mean())
4    AreaPrice=AreaPrice.sort_values(by='Price ($)', ascending=False)
5    AreaPrice[0:25].plot(kind='bar', title='The Most Expensive Area')
6
```

7



```
<matplotlib.axes._subplots.AxesSubplot at 0x7f158a521390>
```

```
1    # Is it possible to buy a cheap house in expensive area?
2
3    # Here I am creating df1 as data frame for 5 to 15 most expensive areas (Top 5 are too expensive).
4    a=AreaPrice[5:15].index
5    df1=pd.DataFrame()
6    df1 = df[df.AreaName.isin(a)]
```

```
1    # From following box graph, I can see Leaside and Nortown are expensive neighbourhoods with few affordable houses.
2
3
4    box = sns.boxplot(x='AreaName', y='Price ($)', data=df1)
```
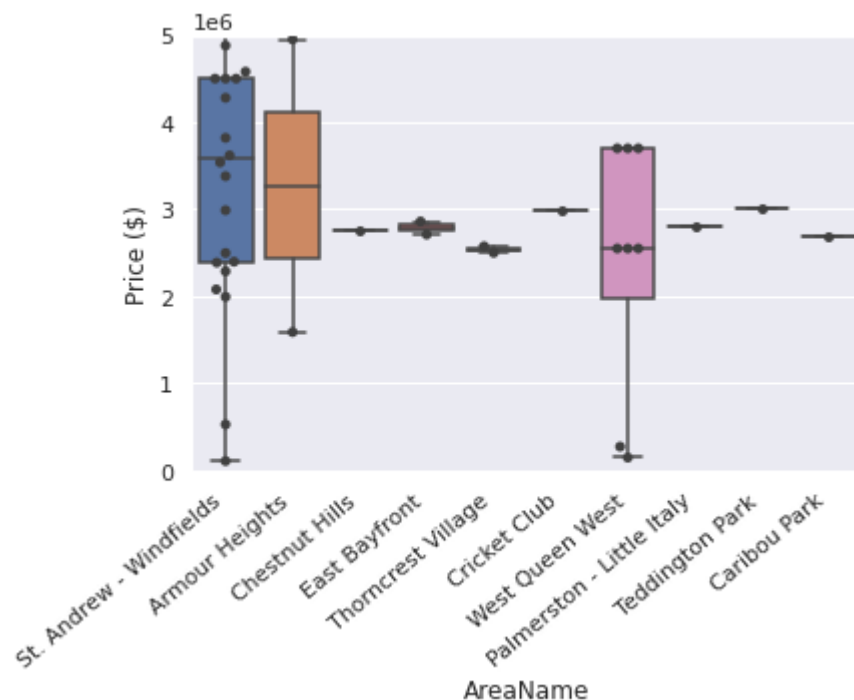
```
5   box = sns.swarmplot(x='AreaName', y='Price ($)', data=df1, color=".25")
6   box.set_ylim([0, 0.5e7])
7   box.set_xticklabels(box.get_xticklabels(), rotation=40, ha="right")
```

```
[Text(0, 0, 'St. Andrew - Windfields'),
 Text(0, 0, 'Armour Heights'),
 Text(0, 0, 'Chestnut Hills'),
 Text(0, 0, 'East Bayfront'),
 Text(0, 0, 'Thorncrest Village'),
 Text(0, 0, 'Cricket Club'),
 Text(0, 0, 'West Queen West'),
 Text(0, 0, 'Palmerston - Little Italy'),
 Text(0, 0, 'Teddington Park'),
 Text(0, 0, 'Caribou Park')]
```



```
1   # I created a table with average price and number of sales for each area, called result1
2   Area['Area']=Area.index
3   Area=Area.rename({'AreaName':'#Sale'}, axis=1)
4   AreaPrice['AreaName']=Area.index
5   result = pd.merge(Area, AreaPrice, on=Area['Area'])
6   result1=result[['#Sale', 'Price ($)']]
```

```
 6   result1=result[[ "#Sale", Price ($)" ]]
 7
 8
 9
10
```

```
 1   # this graph shows number of sales vs price
 2   cc = sns.cubehelix_palette(rot=-.2, as_cmap=True)
 3   cv = sns.relplot(
 4       data=result1,
 5       x='#Sale', y='Price ($)',
 6       palette=cc, sizes=(5, 100),
 7   )
 8
 9   cv.ax.xaxis.grid(True, "minor", linewidth=.05)
10   cv.ax.yaxis.grid(True, "minor", linewidth=.05)
11   cv.set(xscale='log', yscale='log')
```

<seaborn.axisgrid.FacetGrid at 0x7f158ef1ed50>