

## Segmenting and Clustering Toronto Neighbourhoods

### Objective: Segment and cluster the Toronto neighbourhoods based on post codes

Download the table of post codes for neighbourhoods in Toronto from Wikipedia. [https://en.wikipedia.org/wiki/List\\_of\\_postal\\_codes\\_of\\_Canada:\\_M](https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M) ([https://en.wikipedia.org/wiki/List\\_of\\_postal\\_codes\\_of\\_Canada:\\_M](https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M)).

I have Excel 2016, which can hold 1,048,576 rows. The post code table has only 287 rows. I tried using BeautifulSoup first. But when I was half-way through and dealing with all those issues of cleaning the data, I realized that for such a small table, directly downloading it into Excel, then uploading it into Jupyter Notebook will be much easier and faster as well. So, I scraped all my BeautifulSoup codes.

```
In [1]: import pandas as pd
        postcode_df=pd.read_excel('Toronto Post Codes.xlsx')
        postcode_df.head()
```

```
Out[1]:
```

|   | Postcode | Borough          | Neighbourhood    |
|---|----------|------------------|------------------|
| 0 | M1A      | Not assigned     | Not assigned     |
| 1 | M2A      | Not assigned     | Not assigned     |
| 2 | M3A      | North York       | Parkwoods        |
| 3 | M4A      | North York       | Victoria Village |
| 4 | M5A      | Downtown Toronto | Harbourfront     |

```
In [2]: # size of the table
        postcode_df.shape
```

```
Out[2]: (287, 3)
```

```
In [3]: # number of unique Boroughs
        import numpy as np
        postcode_df["Borough"].value_counts()
```

```
Out[3]:
```

|                  |    |
|------------------|----|
| Not assigned     | 77 |
| Etobicoke        | 45 |
| North York       | 38 |
| Downtown Toronto | 37 |
| Scarborough      | 37 |
| Central Toronto  | 17 |
| West Toronto     | 13 |
| York             | 9  |
| East Toronto     | 7  |
| East York        | 6  |
| Mississauga      | 1  |

Name: Borough, dtype: int64

Okay. So there are 77 "Not assigned" in the Borough column. I could use `drop.duplicates()` to drop those "Not assigned". However, that will drop the other boroughs as well, since there are boroughs that have more than one neighbourhoods assigned to it. So, drop by "Borough" won't work.

Let's check the Neighbourhood column. The neighbourhoods should be unique, except for the "Not assigned".

```
In [4]: postcode_df["Neighbourhood"].value_counts()
```

```
Out[4]: Not assigned          77
St. James Town           2
Runnymede                 2
Yorkville                 1
Glencairn                 1
..
Exhibition Place         1
Morningside              1
Parkdale                 1
Bathurst Manor           1
Scarborough Village West 1
Name: Neighbourhood, Length: 209, dtype: int64
```

All the neighbourhoods are unique, except for 1) the "Not assigned" and 2) "Runnymede" and "St. James Town" where there are two copies.

```
In [5]: #Remove the "Not assigned" from the Borough column.
postcode_df.drop_duplicates(subset="Neighbourhood",keep=False, inplace=True)
postcode_df
```

```
Out[5]:
```

|     | Postcode | Borough          | Neighbourhood            |
|-----|----------|------------------|--------------------------|
| 2   | M3A      | North York       | Parkwoods                |
| 3   | M4A      | North York       | Victoria Village         |
| 4   | M5A      | Downtown Toronto | Harbourfront             |
| 5   | M6A      | North York       | Lawrence Heights         |
| 6   | M6A      | North York       | Lawrence Manor           |
| ... | ...      | ...              | ...                      |
| 281 | M8Z      | Etobicoke        | Kingsway Park South West |
| 282 | M8Z      | Etobicoke        | Mimico NW                |
| 283 | M8Z      | Etobicoke        | The Queensway West       |
| 284 | M8Z      | Etobicoke        | Royal York South West    |
| 285 | M8Z      | Etobicoke        | South of Bloor           |

206 rows × 3 columns

All the "Not assigned" in the Neighbourhood column has been dropped. However, since the parameter "keep" was set to False. The process above dropped the two neighbourhoods "Runnymede" and "St. James Town" as well. But I would like to keep those two neighbourhoods, since each copy belong to two different boroughs.

Put the two neighbourhoods back.

```
In [6]: append1=pd.DataFrame({"Postcode":["M5C", "M6N", "M6S", "M4X"], "Borough":["Downtown Toronto", "York", "West Toronto", "Downtown Tor
postcode_df.append(append1,ignore_index=False)
```

```
Out[6]:
```

|     | Postcode | Borough          | Neighbourhood    |
|-----|----------|------------------|------------------|
| 2   | M3A      | North York       | Parkwoods        |
| 3   | M4A      | North York       | Victoria Village |
| 4   | M5A      | Downtown Toronto | Harbourfront     |
| 5   | M6A      | North York       | Lawrence Heights |
| 6   | M6A      | North York       | Lawrence Manor   |
| ... | ...      | ...              | ...              |
| 285 | M8Z      | Etobicoke        | South of Bloor   |
| 0   | M5C      | Downtown Toronto | St. James Town   |
| 1   | M6N      | York             | Runnymede        |
| 2   | M6S      | West Toronto     | Runnymede        |
| 3   | M4X      | Downtown Toronto | St. James Town   |

210 rows × 3 columns

```
In [12]: postcode_df=postcode_df.append(append1,ignore_index=False)
```

```
In [13]: print("The cleaned table has", len(postcode_df["Postcode"]), "row")
```

The cleaned table has 210 row

As shown above,the original table has 287 rows in total and 77 "Not assigned" in the Borough column. After the 77 "Not assigned" rows are dropped, the cleaned table should have 287-77=210 rows. The result of the above code block shows 210 rows.

Let's confirm that there is no more "Not assigned" in Borough column.

```
In [14]: postcode_df["Borough"].value_counts()
```

```
Out[14]: Etobicoke      45
North York    38
Downtown Toronto  37
Scarborough   37
Central Toronto  17
West Toronto   13
York           9
East Toronto   7
East York      6
Mississauga     1
Name: Borough, dtype: int64
```

Cool. There is no more "Not assigned" in the Borough column.

The table is cleaned now. I can start working with it. The next thing I am going to do is to put the neighbourhoods that have the same postcode in the same row. Let's see how many unique postcodes are in the cleaned table.

```
In [16]: postcode_df["Postcode"].value_counts()
```

```
Out[16]: M8Y      8
M9V      8
M5V      7
M8Z      5
M9B      5
..
M4A      1
M4N      1
M5G      1
M5N      1
M6G      1
Name: Postcode, Length: 103, dtype: int64
```

There are 103 unique postcodes. Most of the postcodes has only 1 neighbourhoods associated with it. However, postcodeS "M8Y" and "M9V" have 8 neighbourhoods associated with it, "M5V" has 7, "M4V" and "M8Z" have 5. I need to concatenate all those neighbourhoods under the postcode they are associated with, in one row.

```
In [17]: postcode_df=pd.DataFrame(postcode_df.groupby(["Postcode","Borough"])["Neighbourhood"].apply(lambda x: ','.join(x)))
postcode_df
```

```
Out[17]:
```

|          |             |  | Neighbourhood                                     |
|----------|-------------|--|---|
| Postcode | Borough     |  |   |
| M1B      | Scarborough |  | Rouge,Malvern                                     |
| M1C      | Scarborough |  | Highland Creek,Rouge Hill,Port Union              |
| M1E      | Scarborough |  | Guildwood,Morningside,West Hill                   |
| M1G      | Scarborough |  | Woburn  |
| M1H      | Scarborough |  | Cedarbrae   |
| ...      | ...         |  | ...   |
| M9N      | York        |  | Weston  |
| M9P      | Etobicoke   |  | Westmount   |
| M9R      | Etobicoke   |  | Kingsview Village,Martin Grove Gardens,Richvie... |
| M9V      | Etobicoke   |  | Albion Gardens,Beaumont Heights,Humbergate,Jam... |
| M9W      | Etobicoke   |  | Northwest   |

103 rows × 4 columns

All the neighbourhoods are now concatenated under the postcode they are associated with. However, instead of regular index, the table has "Postcode" and "Borough" as its multi-column index. I am going to fix it below.

In [18]:

postcode\_df.reset\_index(inplace=True)  
postcode\_df

Out[18]:

|     | Postcode | Borough     | Neighbourhood                                     |
|-----|----------|-------------|---|
| 0   | M1B      | Scarborough | Rouge,Malvern                                     |
| 1   | M1C      | Scarborough | Highland Creek,Rouge Hill,Port Union              |
| 2   | M1E      | Scarborough | Guildwood,Morningside,West Hill                   |
| 3   | M1G      | Scarborough | Woburn  |
| 4   | M1H      | Scarborough | Cedarbrae   |
| ... | ...      | ...         | ...   |
| 98  | M9N      | York        | Weston  |
| 99  | M9P      | Etobicoke   | Westmount   |
| 100 | M9R      | Etobicoke   | Kingsview Village,Martin Grove Gardens,Richvie... |
| 101 | M9V      | Etobicoke   | Albion Gardens,Beaumont Heights,Humbergate,Jam... |
| 102 | M9W      | Etobicoke   | Northwest   |

103 rows × 3 columns

Get the size of the cleaned table.

In [19]:

print(postcode\_df.shape)  
  
(103, 3)

Next, find the latitude and longitude coordinates for each postcode.

In [20]:

# Load the "Geospatial\_Coordinates.csv" dataset  
Geo\_df=pd.read\_csv("Geospatial\_Coordinates.csv")  
Geo\_df.head()

Out[20]:

|   | Postcode | Latitude  | Longitude  |
|---|----------|-----------|------------|
| 0 | M1B      | 43.806686 | -79.194353 |
| 1 | M1C      | 43.784535 | -79.160497 |
| 2 | M1E      | 43.763573 | -79.188711 |
| 3 | M1G      | 43.770992 | -79.216917 |
| 4 | M1H      | 43.773136 | -79.239476 |

In [21]:

# Merge "postcode\_df" and "Geo\_df".  
# Use "left" join so that all records from "postcode\_df" and matched records from "Geo\_df" are returned.  
Geomerge\_df=pd.merge(postcode\_df,Geo\_df, how="left")  
Geomerge\_df

Out[21]:

|     | Postcode | Borough     | Neighbourhood                                     | Latitude  | Longitude  |
|-----|----------|-------------|---|-----------|------------|
| 0   | M1B      | Scarborough | Rouge,Malvern                                     | 43.806686 | -79.194353 |
| 1   | M1C      | Scarborough | Highland Creek,Rouge Hill,Port Union              | 43.784535 | -79.160497 |
| 2   | M1E      | Scarborough | Guildwood,Morningside,West Hill                   | 43.763573 | -79.188711 |
| 3   | M1G      | Scarborough | Woburn  | 43.770992 | -79.216917 |
| 4   | M1H      | Scarborough | Cedarbrae   | 43.773136 | -79.239476 |
| ... | ...      | ...         | ...   | ...       | ...        |
| 98  | M9N      | York        | Weston  | 43.706876 | -79.518188 |
| 99  | M9P      | Etobicoke   | Westmount   | 43.696319 | -79.532242 |
| 100 | M9R      | Etobicoke   | Kingsview Village,Martin Grove Gardens,Richvie... | 43.688905 | -79.554724 |
| 101 | M9V      | Etobicoke   | Albion Gardens,Beaumont Heights,Humbergate,Jam... | 43.739416 | -79.588437 |
| 102 | M9W      | Etobicoke   | Northwest   | 43.706748 | -79.594054 |

103 rows × 5 columns

Cluster the neighbourhoods and create visualizations.