# Untitled

January 10, 2019

```python
In [1]: import pandas as pd
        import numpy as np
        import requests

        from bs4 import BeautifulSoup


        source = requests.get('https://en.wikipedia.org/wiki/List_of_postal_codes_of_Canada:_M')

        soup = BeautifulSoup(source, 'html5lib')

        postal_codes_dict = {} # initialize an empty dictionary to save the data in
        for table_cell in soup.find_all('td'):
            try:
                postal_code = table_cell.p.b.text # get the postal code
                postal_code_investigate = table_cell.span.text
                neighborhoods_data = table_cell.span.text # get the rest of the data in the cell
                borough = neighborhoods_data.split('(')[0] # get the borough in the cell

                # if the cell is not assigned then ignore it
                if neighborhoods_data == 'Not assigned':
                    neighborhoods = []
                # else process the data and add it to the dictionary
                else:
                    postal_codes_dict[postal_code] = {}

                    try:
                        neighborhoods = neighborhoods_data.split('(')[1]

                        # remove parantheses from neighborhoods string
                        neighborhoods = neighborhoods.replace('(', ' ')
                        neighborhoods = neighborhoods.replace(')', ' ')

                        neighborhoods_names = neighborhoods.split('/')
                        neighborhoods_clean = ', '.join([name.strip() for name in neighborhoods_
                    except:
                        borough = borough.strip('\n')
                        neighborhoods_clean = borough
```

```python
            # add borough and neighborhood to dictionary
            postal_codes_dict[postal_code]['borough'] = borough
            postal_codes_dict[postal_code]['neighborhoods'] = neighborhoods_clean
        except:
            pass

    # create an empty dataframe
    columns = ['PostalCode', 'Borough', 'Neighborhood']
    toronto_data = pd.DataFrame(columns=columns)
    toronto_data

    # populate dataframe with data from dictionary
    for ind, postal_code in enumerate(postal_codes_dict):
        borough = postal_codes_dict[postal_code]['borough']
        neighborhood = postal_codes_dict[postal_code]['neighborhoods']
        toronto_data = toronto_data.append({"PostalCode": postal_code,
                                            "Borough": borough,
                                            "Neighborhood": neighborhood},
                                            ignore_index=True)

    # print number of rows of dataframe
    toronto_data.shape[0]
```

Out[1]: 0

In [ ]: