NEIGHBORING TORONTO CITY EVERYTHING IS IN ONE NOTEBOOK ITSELF START OF WEEK 3 ASSIGNMENT **IMPORT LIBRARIES** pd.set option('display.max columns', None) pd.set_option('display.max rows', None) from geopy.geocoders import Nominatim # convert an address into latitude and longitude value from pandas.io.json import json normalize # tranform JSON file into a pandas dataframe from sklearn.cluster import DBSCAN from sklearn.preprocessing import StandardScaler GET THE DATA OF TORONTO CITY INTO A DATAFRAME url = 'https://en.wikipedia.org/wiki/List of postal codes of Canada: M' data = pd.read html(url)[0] data head() **Postal Code** 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 Regent Park, Harbourfront **CLEANING UP THE DATA** Toronto data = data.where(data['Borough'] != 'Not assigned').dropna() Toronto data head() **Postal Code Borough** Neighbourhood **2** M3A North York Parkwoods Victoria Village **3** M4A North York 4 M5A Downtown Toronto Regent Park, Harbourfront **5** M6A North York Lawrence Manor, Lawrence Heights **6** M7A Downtown Toronto Queen's Park, Ontario Provincial Government Toronto data.sort values('Postal Code', ascending = True, inplace = True) Toronto data = Toronto data reset index().drop('index', axis = 1) Toronto data head() **Postal Code** Neighbourhood **Borough 0** M1B Scarborough Malvern, Rouge Scarborough Rouge Hill, Port Union, Highland Creek 1 M1C **2** M1E Scarborough Guildwood, Morningside, West Hill **3** M1G Scarborough Woburn **4** M1H Scarborough Cedarbrae for index, row in Toronto data iterrows(): row['Neighbourhood'] == row['Borough'] Toronto data head() **Postal Code** Neighbourhood **Borough 0** M1B Scarborough Malvern, Rouge 1 M1C Scarborough Rouge Hill, Port Union, Highland Creek **2** M1E Scarborough Guildwood, Morningside, West Hill **3** M1G Scarborough Woburn 4 M1H Scarborough Cedarbrae Toronto data shape TRIED FINDING THE LATITUDE AND LONGITUDE OF THE POSTAL CODE (USING GEOPY) Toronto data['Latitude'] = None Toronto data['Longitude'] = None geolocator = Nominatim(user agent="to explorer") for index, row in Toronto data iterrows(): Neighbourhood = row['Neighbourhood'].split(',')[0] address = '{}, Toronto, Ontario'.format(Neighbourhood) location = geolocator.geocode(address) if location is None: latitude = location.latitude longitude = location longitude row['Latitude'] = latitude row['Longitude'] = longitude TRIED FINDING THE LATITUDE AND LONGITUDE OF THE POSTAL CODE (USING GEOCODER) Toronto data['Latitude'] = None Toronto data['Longitude'] = None coordinates = None for index, row in Toronto_data iterrows(): Postal code = row['Postal code'] while (coordinates == None): g = geocoder.google('{}, Toronto, Ontario'.format(Postal code)) coordinates = g.latlng latitude = coordinates[0] longitude = coordinates[1] row['Latitude'] = latitude row['Longitude'] = longitude ~\anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_value(self, series, key) except IndexError: pandas_libs\index.pyx in pandas._libs.index.get_value_at() pandas_libs\util.pxd in pandas._libs.util.get_value_at() TypeError: 'str' object cannot be interpreted as an integer <ipython-input-8-0170bf99ale9> in <module> 12 for index, row in Toronto_data.iterrows(): Postal_code = row['Postal code'] ~\anaconda3\lib\site-packages\pandas\core\series.py in __getitem__(self, key) result = self.index.get_value(self, key) ~\anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_value(self, series, key) raise InvalidIndexError(key) else: except Exception: ~\anaconda3\lib\site-packages\pandas\core\indexes\base.py in get_value(self, series, key) k = self._convert_scalar_indexer(k, kind="getitem") return self._engine.get_value(s, k, tz=getattr(series.dtype, "tz", None)) except KeyError as e1: pandas_libs\index.pyx in pandas._libs.index.IndexEngine.get_value() pandas_libs\index.pyx in pandas._libs.index.IndexEngine.get_value() pandas_libs\index.pyx in pandas. libs.index.IndexEngine.get loc() pandas_libs\hashtable_class_helper.pxi in pandas._libs.hashtable.PyObjectHashTable.get_item() pandas\ libs\hashtable_class helper.pxi in pandas. libs.hashtable.PyObjectHashTable.get item() USING THE GIVEN COORDINATES FILE coordinates df = pd read csv('Geospatial Coordinates.csv') coordinates df head() Postal Code 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 coordinates df shape Toronto data['Latitude'] = coordinates df['Latitude'] e'] = coordinates df[': Toronto_data head() **Postal Code** Borough Neighbourhood Latitude Longitude **0** M1B Scarborough Malvern, Rouge 43.806686 -79.194353 **1** M1C Scarborough Rouge Hill, Port Union, Highland Creek 43.784535 -79.160497 Scarborough Guildwood, Morningside, West Hill **2** M1E 43.763573 -79.188711 **3** M1G 43.770992 -79.216917 Scarborough Woburn **4** M1H Scarborough Cedarbrae 43.773136 -79.239476 **END OF WEEK 3 ASSIGNMENT** START OF WEEK 4 ASSIGNMENT **Problem Statement** You want to open a new coffee shop in Toronto, CA. You want to find the optimal venue to set up your cafe to maximize the profit. 1) We select the top 10 neighborhood having maximum number of offices, colleges and schools. 2) Make a union list of all these neighborhood. 3) For all these neighborhood, we find which one of these has the minimum number of coffee shops. 4) Once we find the neighborhood where we want to set up our shop, we search for trending venues in the neighborhood and set up a shop there. 5) We then combine all the offices, colleges and schools data points to cluster them into groups to find outliers 6) Visualize the colleges and schools in the area and select the area where the data points are densed. 7) If not happy with the result, select the second neighborhood with minimum number of coffee shops. We will be using the Toranto_data that we created earlier to get the neighbouhood data and the Foursquare API for location data **END OF WEEK 4 ASSIGNMENT** START OF WEEK 5 ASSIGNMENT FIND MAXIMUM NUMBER OF OFFICES USING FOURSQUARE API Toronto data['offices'] = None Toronto data['colleges'] = None Toronto data['schools'] = Non RADIUS = 500LIMIT = 100QUERY = 'office' CLIENT ID = 'HNA CLIENT SECRET ACCESS TOKEN = 'CE VERSION = '20180604' office list = [] college list = [] school list = [] for index, row in Toronto data iterrows(): LATITUDE = row['Latitude'] LONGITUDE = row['Longitude'] es/search?client id={}&client secret={}&v={}&l url = 'https://api.foursquare.com/v2/ve l={},{}&radius={}&oauth token={}&query={}&limit={}' format(CLIENT ID,CLIENT SECRET,VERSION, LATITUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests.get(url).json() offices = results['response']['venues'] temp dataframe = json normalize(offices) Number of offices = temp dataframe.shape[0] office list append (Number of offices) C:\Users\Tejas\anaconda3\lib\site-packages\ipykernel_launcher.py:23: FutureWarning: pandas.io.json_json_normalize is deprecated, Toronto data['offices'] = office list Toronto_data head() Postal Code Borough Neighbourhood Latitude Longitude offices colleges schools **0** M1B 43.806686 -79.194353 1 Scarborough Malvern, Rouge None None **1** M1C Scarborough Rouge Hill, Port Union, Highland Creek 43.784535 -79.160497 1 None None **2** M1E Scarborough Guildwood, Morningside, West Hill 43.763573 -79.188711 0 None None **3** M1G Scarborough Woburn 43.770992 -79.216917 0 None None **4** M1H Scarborough Cedarbrae 43.773136 -79.239476 2 None None FIND MAXIMUM NUMBER OF COLLEGS USING FOURSQUARE API QUERY = 'college' for index, row in Toronto data iterrows(): LATITUDE = row['Latitude'] LONGITUDE = row['Longitude'] url = 'https://api.foursquare.com/v2/venues/search?client id={}&client secret={}&v={}&l l={},{}&radius={}&oauth token={}&query={}&limit={}' format(CLIENT ID,CLIENT SECRET,VERSION, LATITUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() colleges = results['response']['venues'] temp dataframe = json normalize(colleges) Number of colleges = temp dataframe.shape[0] college list append(Number of colleges) # This is added back by InteractiveShellApp.init_path() Toronto data['colleges'] = college list Toronto data head() Neighbourhood **Postal Code** Borough Latitude Longitude offices colleges schools **0** M1B Scarborough Malvern, Rouge 43.806686 -79.194353 0 None 1 M1C Scarborough Rouge Hill, Port Union, Highland Creek 43.784535 -79.160497 0 None Scarborough Guildwood, Morningside, West Hill 0 0 **2** M1E 43.763573 -79.188711 None 43.770992 -79.216917 0 0 **3** M1G Scarborough Woburn None 4 M1H 43.773136 -79.239476 0 Scarborough Cedarbrae None FIND MAXIMUM NUMBER OF SCHOOLS USING FOURSQUARE API QUERY = 'school' for index, row in Toronto data iterrows(): LATITUDE = row[' LONGITUDE = row['Longitude'] url = 'https://api.four .={},{}&radius={}&o token={}&query={}&limit={}' format(CLIENT ID,CLIENT SECRET,VERSION, LATITUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() schools = results['response']['venues'] temp dataframe = json normalize(schools) Number of schools = temp dataframe.shape[0] school list append (Number of schools) # This is added back by InteractiveShellApp.init_path() Toronto data['schools'] = school list Toronto data head() **Postal Code** Neighbourhood Latitude Longitude offices colleges schools **Borough 0** M1B Scarborough Malvern, Rouge 43.806686 -79.194353 1 M1C Scarborough Rouge Hill, Port Union, Highland Creek 43.784535 -79.160497 0 **2** M1E Scarborough Guildwood, Morningside, West Hill 43.763573 -79.188711 2 **3** M1G 43.770992 -79.216917 Scarborough Woburn 43.773136 -79.239476 4 M1H Scarborough Cedarbrae ARRINGING THE DATASET WITH RESPECT TO MAXIMUM NUMBER OF OFFICES, COLLEGES AND SCHOOLS Toronto data['sum'] = Toronto data['offices'] + Toronto data['colleges'] + Toronto data[' Toronto_data head() **Postal Code Borough** Neighbourhood Latitude Longitude offices colleges schools Scarborough Malvern, Rouge **0** M1B 43.806686 -79.194353 0 2 **1** M1C Scarborough Rouge Hill, Port Union, Highland Creek 43.784535 -79.160497 0 0 **2** M1E Scarborough Guildwood, Morningside, West Hill 43.763573 -79.188711 0 0 2 2 **3** M1G Scarborough Woburn 43.770992 -79.216917 0 Scarborough Cedarbrae 43.773136 -79.239476 2 **4** M1H arranged Toronto = Toronto data.sort values('sum', ascending = False) arranged Toronto head() **Postal Code Borough** Neighbourhood Latitude Longitude offices colleges schools sum 43.657162 -79.378937 **54** M5B 27 Downtown Toronto Garden District, Ryerson 50 127 **57** M5G 43.657952 -79.387383 23 Downtown Toronto Central Bay Street 50 50 123 **52** M4Y Downtown Toronto Church and Wellesley 43.665860 -79.383160 47 46 10 103 **85** M7A Downtown Toronto Queen's Park, Ontario Provincial Government 43.662301 -79.389494 28 50 10 88 43.662696 -79.400049 **66** M5S Downtown Toronto University of Toronto, Harbord 13 23 50 86 WE SELECT THE TOP 3 RESULTS FOR FURTHER ANALYSIS trimmed Toronto = arranged Toronto.head(3) trimmed Toronto head() **Postal Code** Neighbourhood Latitude Longitude offices colleges schools sum **54** M5B Downtown Toronto Garden District, Ryerson 43.657162 -79.378937 27 127 50 Downtown Toronto Central Bay Street **57** M5G 43.657952 -79.387383 50 50 23 123 **52** M4Y Downtown Toronto Church and Wellesley 43.665860 -79.383160 103 46 10 FOR THESE LOCATIONS WE WILL CHECK WHICH LOCATION HAS THE LEAST NUMBER OF **COFFEE SHOPS** QUERY = 'coffee' coffee list = [] for index, row in trimmed Toronto iterrows(): LATITUDE = row['Latitude'] LONGITUDE = row['Longitude'] search?client id={}&client secret={}&v={}&l url = 'https://api.foursquare l={},{}&radius={}&oauth token={}&guery={}&limit={}' format(CLIENT ID,CLIENT SECRET,VERSION, LATITUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() coffee = results['response']['venues'] temp dataframe = json normalize(coffee) Number of coffee = temp dataframe.shape[0] print(temp dataframe.shape[0]) coffee list append (Number of coffee) C:\Users\Tejas\anaconda3\lib\site-packages\ipykernel_launcher.py:12: FutureWarning: pandas.io.json_json_normalize is deprecated, AS WE CAN SEE, THE THIRD LOCATION HAS COMPARATIVELY LESS COFFEE SHOPS AND TOTAL 103 OFFICEC, COLLEGES AND SCHOOLS. LOOK FOR TRENDING PLACES IN THE NEIGHBORHOOD url = 'https://api.foursquare.com/v2/venues/trending?client id={}&client secret={}&ll={},{} &v={}' format(CLIENT ID, CLIENT SECRET, LATITUDE, LONGITUDE, VERSION) results = requests get(url) json() results if len(results['response']['venues']) == 0: trending venues df = None else: trending venues = results['response']['venues'] trending venues df = json normalize(trending venues) columns filtered = ['name', 'categories'] + ['location.distance', 'location.city', 'loc trending venues df = trending venues df.loc[:, columns filtered] trending venues df GET COORDINATES OF OFFICES IN THE NEIGHBORHOOD LATITUDE = trimmed Toronto.at[52, 'Latitude'] LONGITUDE = trimmed Toronto.at[52, 'Longitude'] OUERY = 'office' target office list = [] target college list = [] target school list = [] url = 'https://api.foursquare.com/v2/venues/search?client id={}&client secret={}&v={}&ll={} , {}&radius={}&oauth token={}&query={}&limit={}'.format(CLIENT ID, CLIENT SECRET, VERSION, LATI TUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() offices = results['response']['venues'] temp dataframe = json normalize(offices) temp dataframe = temp dataframe[['name', 'location.lat', 'location.lng']] office data = temp dataframe print(office data.shape[0]) office data head() name location.lat location.lng **0** Angela's Office, 11th floor OMP 43.670198 -79.379034 **1** Martina's Office - BBDO 43.670275 -79.387297 2 CUPE 2191 Office 43.662360 -79.381737 3 Le Office 43.669748 -79.380067 4 Ken Chan Campaign Office 43.663729 -79.383380 for index, row in office data iterrows(): lat = row['location.lat'] lng = row['location.lng'] target office list append([lat, lng]) en(target_office_list) GET COORDINATES OF COLLEGES IN THE NEIGHBORHOOD QUERY = 'college' ient id={}&client secret={}&v={}&ll={} , {}&radius={}&oauth_token={}&query={}&limit={}'.format(CLIENT_ID,CLIENT_SECRET,VERSION,LATI TUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() colleges = results['response']['venues'] temp dataframe = json normalize(colleges) temp dataframe = temp dataframe[['name', 'location.lat', 'location.lng']] college data = temp dataframe print(college data.shape[0]) college data head() C:\Users\Tejas\anaconda3\lib\site-packages\ipykernel_launcher.py:8: FutureWarning: pandas.io.json_json_normalize is deprecated, name location.lat location.lng 0 College Park 43.661237 -79.383603 1 College Subway Station 43.661311 -79.382819 2 Yonge & College 43.661371 -79.383184 **3** TTC Streetcar #506 College 43.664858 -79.377675 4 College Cleaners 43.662731 -79.381139 for index, row in college data iterrows(): lat = row['location.lat'] lng = row['location.lng'] target college list append([lat, lng]) en(target college list) GET COORDINATES OF SCHOOLS IN THE NEIGHBORHOOD QUERY = 'school' rch?client id={}&client secret={}&v={}&ll={} , {}&radius={}&oauth token={}&query={}&limit={}'.format(CLIENT ID, CLIENT SECRET, VERSION, LATI TUDE, LONGITUDE, RADIUS, ACCESS TOKEN, QUERY, LIMIT) results = requests get(url) json() schools = results['response']['venues'] temp dataframe = json normalize(schools) temp_dataframe = temp_dataframe[['name', 'location.lat', 'location.lng']] school data = temp dataframe print(school data shape[0]) school data head() name location.lat location.lng 0 Church St junior Public School 43.663915 -79.379760 1 St Josephs College Secondary School 43.664259 -79.388681 2 Church Street Junior Public School 43.663483 -79.379778 3 Sol School 43.664743 -79.377530 4 Church Street Public School 43.663703 -79.380324 for index, row in school data.iterrows(): lat = row[lng = row['location.lng'] target school list append([lat, lng]) en(target school list) COMBINING THE OFFICE, COLLEGE AND SCHOOL DATAFRAME temp1 = office data append(college data) combine data = temp1 append(school data) combine data['cluster'] = None print(combine data shape[0]) combine data head() name location.lat location.lng cluster **0** Angela's Office, 11th floor OMP 43.670198 -79.379034 None 1 Martina's Office - BBDO 43.670275 -79.387297 None 2 CUPE 2191 Office 43.662360 -79.381737 None 3 Le Office 43.669748 -79.380067 None 4 Ken Chan Campaign Office 43.663729 -79.383380 None COMBINE THE LIST OF COORDINATES TO PREPARE FOR CLUSTERING for value in target college list: target office list append(value) for value in target school list: target office list append(value) print(len(target office list)) print(combine data.shape[0]) target list = target office list **CLUSTERING USING DBSCAN**

