

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

Emigma Housing prize Dataset

```
In [2]: 1 #Emigma housing prize dataset
        2 df=pd.read_csv("19b7b17c-e655-4013-b4d1-7daa9727225f_NewYorkCityPrope
```

```
In [3]: 1 df.head()
```

Out[3]:

	borough_code	borough_code_definition	neighborhood	building_class_category_code	building_
0	1	Manhattan	ALPHABET CITY	07	RE
1	1	Manhattan	ALPHABET CITY	07	RE
2	1	Manhattan	ALPHABET CITY	07	RE
3	1	Manhattan	ALPHABET CITY	07	RE
4	1	Manhattan	ALPHABET CITY	07	RE

5 rows x 62 columns

NYC Demographic dataset

```
In [58]: 1 #Demographics dataset
        2 df_socio_economic=pd.read_csv("Demographic_Statistics_By_Zip_Code.csv")
        3 df_socio_economic.head()
```

Out[58]:

	JURISDICTION NAME	COUNT PARTICIPANTS	COUNT FEMALE	PERCENT FEMALE	COUNT MALE	PERCENT MALE	COUNT GENDER UNKNOWN	PERCE GEND UNKNOI
0	10001	44	22	0.50	22	0.50	0	
1	10002	35	19	0.54	16	0.46	0	
2	10003	1	1	1.00	0	0.00	0	
3	10004	0	0	0.00	0	0.00	0	
4	10005	2	2	1.00	0	0.00	0	

5 rows x 46 columns

```
In [6]: 1 df_socio_economic.columns
```

```
Out[6]: Index(['JURISDICTION NAME', 'COUNT PARTICIPANTS', 'COUNT FEMALE',
              'PERCENT FEMALE', 'COUNT MALE', 'PERCENT MALE', 'COUNT GENDER U
              NKNOWN',
              'PERCENT GENDER UNKNOWN', 'COUNT GENDER TOTAL', 'PERCENT GENDER
              TOTAL',
              'COUNT PACIFIC ISLANDER', 'PERCENT PACIFIC ISLANDER',
              'COUNT HISPANIC LATINO', 'PERCENT HISPANIC LATINO',
              'COUNT AMERICAN INDIAN', 'PERCENT AMERICAN INDIAN',
              'COUNT ASIAN NON HISPANIC', 'PERCENT ASIAN NON HISPANIC',
              'COUNT WHITE NON HISPANIC', 'PERCENT WHITE NON HISPANIC',
              'COUNT BLACK NON HISPANIC', 'PERCENT BLACK NON HISPANIC',
              'COUNT OTHER ETHNICITY', 'PERCENT OTHER ETHNICITY',
              'COUNT ETHNICITY UNKNOWN', 'PERCENT ETHNICITY UNKNOWN',
              'COUNT ETHNICITY TOTAL', 'PERCENT ETHNICITY TOTAL',
              'COUNT PERMANENT RESIDENT ALIEN', 'PERCENT PERMANENT RESIDENT A
              LIEN',
              'COUNT US CITIZEN', 'PERCENT US CITIZEN', 'COUNT OTHER CITIZEN
              STATUS',
              'PERCENT OTHER CITIZEN STATUS', 'COUNT CITIZEN STATUS UNKNOWN',
              'PERCENT CITIZEN STATUS UNKNOWN', 'COUNT CITIZEN STATUS TOTAL',
              'PERCENT CITIZEN STATUS TOTAL', 'COUNT RECEIVES PUBLIC ASSISTAN
              CE',
              'PERCENT RECEIVES PUBLIC ASSISTANCE',
              'COUNT NRECEIVES PUBLIC ASSISTANCE',
              'PERCENT NRECEIVES PUBLIC ASSISTANCE',
              'COUNT PUBLIC ASSISTANCE UNKNOWN', 'PERCENT PUBLIC ASSISTANCE U
              NKNOWN',
              'COUNT PUBLIC ASSISTANCE TOTAL', 'PERCENT PUBLIC ASSISTANCE TOT
              AL'],
              dtype='object')
```

```
In [60]: 1 processed_socio_economic=df_socio_economic.drop(['PERCENT FEMALE','PE
```

Foursquare Dataset

```
In [7]: 1 df_four=pd.read_table("new york_anon_locationData_newcrawl.txt")
```

In [8]:

```

1 #Foursquare dataset
2 import re
3 handle=open('new_york_anon_locationData_newcrawl.txt')
4 foursquaredf=pd.DataFrame()
5 idlist=[]
6 restlist=[]
7 translator = str.maketrans('', '', '*(())')
8 for line in handle:
9     line=line.strip()
10    words=line.split(';')
11    #print(words[0])
12    #print(words[1])
13    words[0]=int(words[0].translate(translator))
14    words[1]=words[1].translate(translator)
15    words[1]=words[1].split(',')
16    #print(words[0])
17    idlist.append(words[0])
18    restlist.append(words[1])
19    #print(words[1].split(', '))
20 foursquaredf=pd.DataFrame(restlist)
21 foursquaredf=foursquaredf.iloc[:, :6]
22 foursquaredf['id']=idlist

```

In [9]:

```

1 #Work on Foursquare dataset to bring into the required format
2 foursquaredf_final=pd.DataFrame()
3 foursquaredf_final['latitude']=foursquaredf[0].astype(float)
4 foursquaredf_final['longitude']=foursquaredf[1].astype(float)
5 foursquaredf_final['type']=foursquaredf[2].astype(str)
6 foursquaredf_final['num1']=foursquaredf[3]
7 foursquaredf_final['num2']=foursquaredf[5]
8 foursquaredf_final['id']=foursquaredf['id']

```

In [10]:

```
1 foursquaredf_final.head()
```

Out[10]:

	latitude	longitude	type	num1	num2	id
0	40.760265	-73.989105	'Italian'	'217'	'Ristorante Da Rosina'	42889
1	40.780704	-73.954704	'Bakery'	'1291'	'Le Pain Quotidien'	57489
2	40.663925	-74.118230	'Video Store'	'60'	'Blockbuster'	42890
3	40.739840	-73.991770	"Corporate ' Office"	'6'	'Day & Night Office'	74771
4	40.786987	-73.807530	'Other - Food'	'29'	"Erin's Isle"	42891

```
In [11]: 1 #map each of FOursquare venues to one of 7 categories
2 from pydoc import deque
3 import json
4
5 class _Node(object):
6
7     def __init__(self, parent, data):
8         self.parent = parent
9         self.data = data
10        self.children = []
```

```
In [12]: 1 #map each of FOursquare venues to one of 7 categories
2 class Categories(object):
3
4     def __init__(self):
5         self._categories = None
6         # self._resources = Resources()
7
8         self._root = None
9         self._name_category_map = {}
10        self._short_name_category_map = {}
11
12        self._load()
13        self._load_colors()
14
15    def _load(self):
16
17        json_data=open('categories.json')
18        data = json.load(json_data)
19        json_data.close()
20        self._categories = data
21
22        self._root = _Node(None, None)
23        nodeQueue = deque([])
24        nodeQueue.append(self._root)
25
26        objectQueue = deque([])
27        objectQueue.append(self._categories)
28
29        while(len(nodeQueue) > 0):
30            node = nodeQueue.pop()
31            category = objectQueue.pop()
32            if not (node.data is None):
33                self._name_category_map[node.data['name']] = node
34                self._short_name_category_map[node.data['shortName']]
35
36            categories = category.get('categories')
37            if not (categories is None):
38                for child in categories:
```

```
39         childNode = _Node(node, { 'name': child['pluralName'],
40                                     'shortName': child['shortName'] })
41         node.children.append(childNode)
42
43         nodeQueue.append(childNode)
44         objectQueue.append(child)
45
46     def get_parent(self, category):
47         node = self._name_category_map.get(category)
48         if node is None:
49             node = self._short_name_category_map.get(category)
50             if node is None:
51                 return None
52
53         while(node.parent != None and node.parent.data != None):
54             node = node.parent
55         break
56
57         return node.data['name']
58
59
60
61     def get_top_parent(self, category):
62         node = self._name_category_map.get(category)
63         if node is None:
64             node = self._short_name_category_map.get(category)
65             if node is None:
66                 return None
67
68         while(node.parent != None and node.parent.data != None):
69             node = node.parent
70
71         return node.data['name']
72     def _load_colors(self):
73         self._colors = { 'Arts & Entertainment' : 'r', 'College & University' : 'g',
74                         'Food' : 'b', 'Nightlife Spot' : '#DDFF00', 'Outdoor Recreation' : 'm',
75                         'Professional & Other Places' : 'c', 'Residence' : 'p',
76                         'Shop & Service' : '#800000', 'Travel & Transport' : 'k' }
77
78
79
80
81
```

```
In [13]: 1 #map each of FOursquare venues to one of 7 categories
2 import re
3 dog=Categories()
4 cate=[]
5 for i,row in foursquaredf_final.iterrows():
6     tr=re.sub('[\\\'"]', '', row['type'])
7     ssss=re.sub('^[ ]', '',tr)
8     ssss=re.sub(' ',' ',ssss)
9     cate.append(dog.get_top_parent(ssss))
```

```
In [14]: 1 #map each of FOursquare venues to one of 7 categories
2 four_reg=foursquaredf_final
3 four_reg['category']=cate
4
5 four_reg.category = pd.Categorical(four_reg.category)
6 four_reg['category_code'] = four_reg.category.cat.codes
```

```
In [15]: 1 four_reg.head()
```

Out[15]:

	latitude	longitude	type	num1	num2	id	category	category_code
0	40.760265	-73.989105	'Italian'	'217'	'Ristorante Da Rosina'	42889	Food	2
1	40.780704	-73.954704	'Bakery'	'1291'	'Le Pain Quotidien'	57489	Food	2
2	40.663925	-74.118230	'Video Store'	'60'	'Blockbuster'	42890	Shops & Services	7
3	40.739840	-73.991770	"Corporate ' Office"	'6'	'Day & Night Office'	74771	Professional & Other Places	5
4	40.786987	-73.807530	'Other - Food'	'29'	"Erin's Isle"	42891	NaN	-1

```
In [16]: 1 #import the USZipcode api to map longitude and latitude to zip code
2 from uszipcode import ZipcodeSearchEngine
3 search = ZipcodeSearchEngine()
4 res = search.by_coordinate(40.760265, -73.989105, radius=1, returns=1)
```

```
In [17]: 1 #See what does the API return
        2 res
```

```
Out[17]: [{"City": "New York", "Density": 56161.36363636363, "HouseOfUnits": 17
          958, "LandArea": 0.44, "Latitude": 40.7602619, "Longitude": -73.993287
          2, "NEBoundLatitude": 40.768738, "NEBoundLongitude": -73.9781161, "Pop
          ulation": 24711, "SWBoundLatitude": 40.723624900000004, "SWBoundLongit
          ude": -74.004786, "State": "NY", "TotalWages": 1686575064.0, "WaterAre
          a": 0.0, "Wealthy": 68251.99562947675, "Zipcode": "10036", "ZipcodeTyp
          e": "Standard"}]
```

```
In [18]: 1 #The API returns all information about zip code ranging from socio ec
        2 #data to demographic data which could also be used for analysis of hc
        3 zipcodes=[]
        4 density=[]
        5 houseofunit=[]
        6 landarea=[]
        7 pop=[]
        8 totwage=[]
        9 water_area=[]
       10 wealthy=[]
       11 zipped_df=pd.DataFrame()
       12 for i,row in four_reg.iterrows():
       13     lat=float(row['latitude'])
       14     lon=float(row['longitude'])
       15     #print(lat,lon)
       16     try:
       17         res=search.by_coordinate(lat, lon, radius=1, returns=1)[0]
       18         zipp=res['Zipcode']
       19         densityl=res['Density']
       20         houseofunitl=res['HouseOfUnits']
       21         landareal=res['LandArea']
       22         popl=res['Population']
       23         totwagel=res['TotalWages']
       24         water_areal=res['WaterArea']
       25         wealthyl=res['Wealthy']
       26     except:
       27         zipp=-1
       28         densityl=-1
       29         houseofunitl=-1
       30         landareal=-1
       31         popl=-1
       32         totwagel=-1
       33         water_areal=-1
       34         wealthyl=-1
       35     #print(res)
       36     zipcodes.append(zipp)
       37     density.append(densityl)
       38     houseofunit.append(houseofunitl)
       39     landarea.append(landareal)
```



```

39     landarea.append(landareal)
40     pop.append(popl)
41     totwage.append(totwagel)
42     water_area.append(water_areal)
43     wealthy.append(wealthyl)
44     four_reg['zipcode']=zipcodes
45     zipped_df['zipcode']=zipcodes
46     zipped_df['density']=density
47     zipped_df['houseofunits']=houseofunit
48     zipped_df['landarea']=landarea
49     zipped_df['pop']=pop
50     zipped_df['totwage']=totwage
51     zipped_df['water_area']=water_area
52     zipped_df['wealthy']=wealthy

```

In [19]: 1 four_reg.head()

Out[19]:

	latitude	longitude	type	num1	num2	id	category	category_code	zip
0	40.760265	-73.989105	'Italian'	'217'	'Ristorante Da Rosina'	42889	Food	2	-
1	40.780704	-73.954704	'Bakery'	'1291'	'Le Pain Quotidien'	57489	Food	2	-
2	40.663925	-74.118230	'Video Store'	'60'	'Blockbuster'	42890	Shops & Services	7	(
3	40.739840	-73.991770	"Corporate 'Office'"	'6'	'Day & Night Office'	74771	Professional & Other Places	5	-
4	40.786987	-73.807530	'Other - Food'	'29'	"Erin's Isle"	42891	NaN	-1	-

In [76]: 1 zipped_df.head()
2 zipped_df.to_csv("zip_api_data.csv")

```

In [55]: 1 #Count the venues in each of the zip codes
2 zip_codes=four_reg['zipcode'].unique()
3 list_dicts=[]
4 for j,z in enumerate(zip_codes):
5     dict_count=dict.fromkeys(list(four_reg['category'].unique()),0)
6     df_zip_part=four_reg[four_reg['zipcode']==z]
7     for i,row in df_zip_part.iterrows():
8         dict_count[row['category']] +=1
9     list_dicts.append(dict_count)

```

In [56]: 1 zipped_foursquare=pd.DataFrame.from_dict(list_dicts)
2 zipped_foursquare['zipcode']=zip_codes

```
In [75]: 1 #save the resulting datasets to csv file
          2 zipped_foursquare.head()
          3 zipped_foursquare.to_csv("foursquare.csv")
```

```
In [74]: 1 processed_socio_economic.head()
          2 processed_socio_economic.to_csv("Socio_economic.csv")
```

NYC Taxi Data

```
In [68]: 1 df_taxi_jan=pd.read_csv('yellow_tripdata_2016-01.csv')
```

```
In [69]: 1 df_taxi_jan.head()
```

Out[69]:

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	pickup
0	2	2016-01-01 00:00:00	2016-01-01 00:00:00	2	1.10	
1	2	2016-01-01 00:00:00	2016-01-01 00:00:00	5	4.90	
2	2	2016-01-01 00:00:00	2016-01-01 00:00:00	1	10.54	
3	2	2016-01-01 00:00:00	2016-01-01 00:00:00	1	4.75	
4	2	2016-01-01 00:00:00	2016-01-01 00:00:00	3	1.76	

```
In [71]: 1 df_taxi_jan=df_taxi_jan[['pickup_longitude','pickup_latitude','dropoff_longitude','dropoff_latitude']]
```

```
In [72]: 1 df_taxi_jan.head()
```

Out[72]:

	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	-73.990372	40.734695	-73.981842	40.732407
1	-73.980782	40.729912	-73.944473	40.716679
2	-73.984550	40.679565	-73.950272	40.788925
3	-73.993469	40.718990	-73.962242	40.657333
4	-73.960625	40.781330	-73.977264	40.758514

```
In [ ]: 1 df_taxi_jan.head()
```

```
In [78]: 1 sample_taxi=df_taxi_jan.sample(100000)
```

In [85]: 1 sample_taxi.head()

Out[85]:

	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
4231934	-73.963196	40.793873	-73.979362	40.776642
4279558	-73.991035	40.751293	-73.979820	40.746521
6514632	-73.984375	40.722385	-73.981468	40.716774
4841223	-73.982559	40.739552	-73.967133	40.757561
10219474	-73.971558	40.758381	-73.972870	40.795879

```

In [ ]: 1 #Map latitude and longitude to ZIP code
2 dict_pick_count=dict()
3 dict_drop_count=dict()
4 for i,row in sample_taxi.iterrows():
5     pick_lon=float(row['pickup_longitude'])
6     pick_lat=float(row['pickup_latitude'])
7
8     drop_lon=float(row['dropoff_longitude'])
9     drop_lat=float(row['dropoff_latitude'])
10
11     print("hi")
12     try:
13         res_pick=search.by_coordinate(pick_lat, pick_lon, radius=1, r
14         print(res_pick)
15         res_pick_zip=res_pick['Zipcode']
16         res_drop=search.by_coordinate(drop_lat,drop_lon,radius=1,retu
17     except:
18         pass
19     print(res_pick)
20
21
22     if res_pick in dict_pick_count.keys():
23         dict_pick_count[res_pick]+=1
24     else:
25         dict_pick_count[res_pick]=0
26     if res_drop in dict_drop_count.keys():
27         dict_drop_count[res_drop]+=1
28     else:
29         dict_drop_count[res_drop]=0
30

```

```

hi
{
    "City": "New York",
    "Density": 126133.33333333331,
    "HouseOfUnits": 47617,
    "Longitude": -73.963196

```

```
LandArea": 0.75,  
"Latitude": 40.7999209,  
"Longitude": -73.96831019999998,  
"NEBoundLatitude": 40.8108,  
"NEBoundLongitude": -73.95699309999998,  
"Population": 94600,  
"SWBoundLatitude": 40.7878149,  
"SWBoundLongitude": -73.98440699999998,  
"State": "NY",  
"TotalWages": 1675737238.0,  
"WaterArea": 0.0,  
"Wealthy": 17713.924291754756,  
"Zipcode": "10025",  
"ZipcodeType": "Standard"  
,
```

In []:

1