

Public Recycling Bins in the NYC Boroughs

Background:

Recycling helps protect the environment, ensures environmental sustainability and helps to tackle climate change. The New York Department of Sanitation (DSNY) is responsible for collecting trash, cleaning the city of New York and more relevantly, they're responsible for public recycling bins.

Throughout this project, I'll explore 3 datasets to find the ratio of people per bin in each New York borough. By doing so, I'll be able to provide suggestions to further improve the recycling opportunities in New York.

The 3 datasets are:

- MTA Turnstiles:** This dataset will be used to find the daily traffic in each station which will help get a gauge of how many people pass by each day.
- MTA Stations:** This dataset will be used to link the stations to the New York City boroughs.
- Public Recycling Bins:** This dataset will be used to get the number of recycling bins that are in each New York Borough.

[Dataset #1: MTA Turnstiles] (<http://web.mta.info/developers/turnstile.html>)

Importing the MTA Turnstiles Dataset

```
In [1]: import pandas as pd
import astext
import requests
import matplotlib.pyplot as plt
```

```
In [2]: def get_data(weeks):
    url = "http://web.mta.info/developers/data/nyct/turnstile/turnstile_{}.txt"
    dfs = []
    for week in weeks:
        file_url = url.format(week)
        dfs.append(pd.read_csv(file_url))
    return pd.concat(dfs)
```

```
weeks = [200201, 200208, 200215, 200222,
         200307, 200314, 200321,
         200328, 200404, 200411, 200418]
```

```
# February 1st to April 25th
turnstiles_df = get_data(weeks)
```

In [3]: turnstiles_df

C/A	UNIT	SCP	STATION	LINENAME	DIVISION	DATE	TIME	DESC	ENTRIES	EXITS
0	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	03:00:00	REGULAR	7356095 2493703
1	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	07:00:00	REGULAR	7356105 2493714
2	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	11:00:00	REGULAR	7356170 2493761
3	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	15:00:00	REGULAR	7356333 2493812
4	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	19:00:00	REGULAR	7356581 2493862
...
206171	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	05:00:00	REGULAR	5554 514
206172	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	09:00:00	REGULAR	5554 514
206173	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	13:00:00	REGULAR	5554 514
206174	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	17:00:00	REGULAR	5554 514
206175	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	21:00:00	REGULAR	5554 514

2472168 rows x 11 columns

Cleaning the Turnstiles Data and Extracting Features

Removing the white space from column names:

```
In [4]: turnstiles_df.columns
```

```
Out[4]: Index(['C/A', 'UNIT', 'SCP', 'STATION', 'LINENAME', 'DIVISION', 'DATE', 'TIME',
       'DESC', 'ENTRIES', 'EXITS'],
       dtype='object')
```

```
In [5]: turnstiles_df.columns = turnstiles_df.columns.str.strip()
```

```
turnstiles_df.columns
```

```
Out[5]: Index(['C/A', 'UNIT', 'SCP', 'STATION', 'LINENAME', 'DIVISION', 'DATE', 'TIME',
       'DESC', 'ENTRIES', 'EXITS'],
       dtype='object')
```

Converting date and time columns to a single datetime column:

```
In [6]: turnstiles_df["DATE_TIME"] = pd.to_datetime(turnstiles_df["DATE"] + " " + turnstiles_df["TIME"],
                                                format="%m/%d/%Y %H:%M:%S")
```

```
turnstiles_df
```

C/A	UNIT	SCP	STATION	LINENAME	DIVISION	DATE	TIME	DESC	ENTRIES	EXITS	DATE TIME
0	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	03:00:00	REGULAR	7356095 2493703	2020-01-25 03:00:00
1	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	07:00:00	REGULAR	7356105 2493714	2020-01-25 07:00:00
2	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	11:00:00	REGULAR	7356170 2493761	2020-01-25 11:00:00
3	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	15:00:00	REGULAR	7356333 2493812	2020-01-25 15:00:00
4	A002	R051	02-00-00	59 ST	NQR456W	BMT	01/25/2020	19:00:00	REGULAR	7356581 2493862	2020-01-25 19:00:00
...
206171	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	05:00:00	REGULAR	5554 514	2020-04-17 05:00:00
206172	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	09:00:00	REGULAR	5554 514	2020-04-17 09:00:00
206173	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	13:00:00	REGULAR	5554 514	2020-04-17 13:00:00
206174	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	17:00:00	REGULAR	5554 514	2020-04-17 17:00:00
206175	TRAM2	R469	00-05-01	RIT-ROOSEVELT		R	04/17/2020	21:00:00	REGULAR	5554 514	2020-04-17 21:00:00

2472168 rows x 12 columns

Checking for duplicates:

```
In [7]: temp = turnstiles_df
temp['ENTRIES'] = temp['ENTRIES'].str.replace(' ', '')
temp['EXITS'] = temp['EXITS'].str.replace(' ', '')
temp['DATE'] = temp['DATE'].str.replace(' ', '')
```

```
temp['ENTRIES'].unique()
```

```
Out[7]: array([2, 1], dtype=int64)
```

The result above indicates that some values are duplicated.

```
In [8]: turnstiles_df.sort_values(["C/A", "UNIT", "SCP", "STATION", "DATE_TIME"], inplace=True)
# dropping duplicated rows:
turnstiles_df.drop_duplicates(subset=["C/A", "UNIT", "SCP", "STATION", "DATE_TIME"], inplace=True)
```

```
In [9]: temp = turnstiles_df
temp['ENTRIES'] = temp['ENTRIES'].str.replace(' ', '')
temp['EXITS'] = temp['EXITS'].str.replace(' ', '')
temp['DATE'] = temp['DATE'].str.replace(' ', '')
```

```
temp['ENTRIES'].unique()
```

```
Out[9]: array([1], dtype=int64)
```

Now, there are no duplicated values.

```
In [10]: temp = turnstiles_df
temp['C/A'] = temp['C/A'].str.replace(' ', '')
temp['UNIT'] = temp['UNIT'].str.replace(' ', '')
temp['SCP'] = temp['SCP'].str.replace(' ', '')
temp['STATION'] = temp['STATION'].str.replace(' ', '')
temp['DATE'] = temp['DATE'].str.replace(' ', '')
temp['TIME'] = temp['TIME'].str.replace(' ', '')
```

```
temp['ENTRIES'] = temp['ENTRIES'].str.replace(' ', '')
```

```
Out[10]: array([1], dtype=int64)
```

Sanity check with the EXITS column, there are no duplicates.

```
In [11]: temp = turnstiles_df
temp['C/A'] = temp['C/A'].str.replace(' ', '')
temp['UNIT'] = temp['UNIT'].str.replace(' ', '')
temp['SCP'] = temp['SCP'].str.replace(' ', '')
temp['STATION'] = temp['STATION'].str.replace(' ', '')
temp['DATE'] = temp['DATE'].str.replace(' ', '')
temp['TIME'] = temp['TIME'].str.replace(' ', '')
```

```
temp['EXITS'] = temp['EXITS'].str.replace(' ', '')
```

```
Out[11]: array([1], dtype=int64)
```

Now, we have a dataframe that contains each station, where is it located (borough) and its average daily traffic.

[Dataset #2: MTA Stations] (<http://web.mta.info/developers/data/nyct/subway/Stations.csv>)

Importing the MTA Stations dataset:

```
In [12]: stations = pd.read_csv("http://web.mta.info/developers/data/nyct/subway/Stations.csv")
```

```
stations
```

```
Out[12]: StationID      ComplexID      GTFS_StopID      Division      Line      StopName      Latitude      Longitude      Structure      GTFS_StopName      GTFS_StopID      NorthLabel      SouthLabel      ADA
0      34 ST-PENN STA          1           R01          BMT      Astoria-Ditmars Blvd      Q      NW      Elevated      40.775036      -73.912034      NaN      Manhattan      0
1      GRD CNTRL-42 ST          2           R03          BMT      Astoria      Q      NW      Elevated      40.770258      -73.917843      Ditmars Blvd      Manhattan      1
2      34 ST-HERALD SQ          3           R04          BMT      Astoria      30 Av      Q      NW      Elevated      40.766779      -73.921477      Astoria-Ditmars Blvd      Manhattan      0
3      34 ST-HERALD SQ          4           R05          BMT      Astoria      Broadway      Q      NW      Elevated      40.761820      -73.925508      Astoria-Ditmars Blvd      Manhattan      0
4      14 ST-UNION SQ          5           R06          BMT      Astoria      36 Av      Q      NW      Elevated      40.756804      -73.929575      Astoria-Ditmars Blvd      Manhattan      0
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

```
...
```

