

```
In [1]: # importing libraries
import pandas as pd
import numpy as np
```

```
In [40]: import seaborn as sns
pd.options.mode.chained_assignment = None
```

```
In [3]: import matplotlib.pyplot as plt
matplotlib inline
```

Import a 311 NYC service request.

```
In [4]: df = pd.read_csv('311_dataset.csv',low_memory=False)
```

```
In [5]: df.head()
```

	Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Incident Address	Bridge Highway Name	Id
0	32310363	12/31/2015 11:59:45 PM	01-01-16 0:55	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	10034.0	71 VERMILY AVE	...	NaN
1	32309934	12/31/2015 11:59:44 PM	01-01-16 1:26	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	11105.0	27-07 23 AVENUE	...	NaN
2	32309159	12/31/2015 11:59:28 PM	01-01-16 4:51	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	10458.0	2897 VALENTI AVENUE	...	NaN
3	32305098	12/31/2015 11:57:46 PM	01-01-16 7:43	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	10461.0	2940 BAISLEY AVENUE	...	NaN
4	32306529	12/31/2015 11:56:58 PM	01-01-16 3:24	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	11373.0	87-14 57 ROAD	...	NaN

5 rows × 53 columns

```
In [6]: # shape of the dataframe
df.shape
```

```
Out[6]: (300698, 53)
```

```
In [7]: # columns in df
df.columns
```

```
Out[7]: Index(['Unique Key', 'Created Date', 'Closed Date', 'Agency', 'Agency Name',
        'Complaint Type', 'Descriptor', 'Location Type', 'Incident Zip',
        'Incident Address', 'Street Name', 'Cross Street 1', 'Cross Street 2',
        'Intersection Street 1', 'Intersection Street 2', 'Address Type',
        'City', 'Landmark', 'Facility Type', 'Status', 'Due Date',
        'Resolution Description', 'Resolution Action Updated Date',
        'Community Board', 'Borough', 'X Coordinate (State Plane)',
        'Y Coordinate (State Plane)', 'Park Facility Name', 'Park Borough',
        'School Name', 'School Number', 'School Region', 'School Code',
        'School Phone Number', 'School Address', 'School City', 'School State',
        'School Zip', 'School Not Found', 'School or Citywide Complaint',
        'Vehicle Type', 'Taxi Company Borough', 'Taxi Pick Up Location',
        'Bridge Highway Name', 'Bridge Highway Direction', 'Road Ramp',
        'Bridge Highway Segment', 'Garage Lot Name', 'Ferry Direction',
        'Ferry Terminal Name', 'Latitude', 'Longitude', 'Location'],
        dtype='object')
```

```
In [8]: # converting the columns into lower case and replacing the space with '_'.
```

```
df.columns = df.columns.str.lower().str.replace(' ', '_')
```

```
In [9]: # checking for NaN or null values in df
df.isna().sum()
```

Out[9]:	unique_key	created_date	closed_date	agency	agency_name	complaint_type	descriptor	location_type	incident_zip	incident_address	bridge_highway_name	id
	0	0	0	0	2164	0	0	0	0	0	0	0
	created_date	0	0	0	0	0	0	0	0	0	0	0
	closed_date	0	0	0	0	0	0	0	0	0	0	0
	agency	0	0	0	0	0	0	0	0	0	0	0
	agency_name	0	0	0	0	0	0	0	0	0	0	0
	complaint_type	0	0	0	0	0	0	0	0	0	0	0
	descriptor	0	0	0	0	0	0	0	0	0	0	0
	location_type	0	0	0	0	0	0	0	0	0	0	0
	incident_zip	0	0	0	0	0	0	0	0	0	0	0
	incident_address	0	0	0	0	0	0	0	0	0	0	0
	street_name	0	0	0	0	0	0	0	0	0	0	0
	cross_street_1	0	0	0	0	0	0	0	0	0	0	0
	cross_street_2	0	0	0	0	0	0	0	0	0	0	0
	intersection_street_1	0	0	0	0	0	0	0	0	0	0	0
	intersection_street_2	0	0	0	0	0	0	0	0	0	0	0
	address_type	0	0	0	0	0	0	0	0	0	0	0
	city	0	0	0	0	0	0	0	0	0	0	0
	landmark	0	0	0	0	0	0	0	0	0	0	0
	facility_type	0	0	0	0	0	0	0	0	0	0	0
	status	0	0	0	0	0	0	0	0	0	0	0
	due_date	0	0	0	0	0	0	0	0	0	0	0
	resolution_description	0	0	0	0	0	0	0	0	0	0	0
	resolution_action_updated_date	0	0	0	0	0	0	0	0	0	0	0
	community_board	0	0	0	0	0	0	0	0	0	0	0
	borough	0	0	0	0	0	0	0	0	0	0	0
	x_coordinate_state_plane	0	0	0	0	0	0	0	0	0	0	0
	y_coordinate_state_plane	0	0	0	0	0	0	0	0	0	0	0
	park_facility_name	0	0	0	0	0	0	0	0	0	0	0
	park_borough	0	0	0	0	0	0	0	0	0	0	0
	school_name	0	0	0	0	0	0	0	0	0	0	0
	school_number	0	0	0	0	0	0	0	0	0	0	0
	school_region	0	0	0	0	0	0	0	0	0	0	0
	school_code	0	0	0	0	0	0	0	0	0	0	0
	school_zip	0	0	0	0	0	0	0	0	0	0	0
	school_not_found	0	0	0	0	0	0	0	0	0	0	0
	school_or_citywide_complaint	0	0	0	0	0	0	0	0	0	0	0
	vehicle_type	0	0	0	0	0	0	0	0	0	0	0
	taxi_company_borough	0	0	0	0	0	0	0	0	0	0	0
	taxi_pick_up_location	0	0	0	0	0	0	0	0	0	0	0
	bridge_highway_name	0	0	0	0	0	0	0	0	0	0	0
	bridge_highway_direction	0	0	0	0	0	0	0	0	0	0	0
	road_ramp	0	0	0	0	0	0	0	0	0	0	0
	bridge_highway_segment	0	0	0	0	0	0	0	0	0	0	0
	garage_lot_name	0	0	0	0	0	0	0	0	0	0	0
	ferry_direction	0	0	0	0	0	0	0	0	0	0	0
	ferry_terminal_name	0	0	0	0	0	0	0	0	0	0	0
	latitude	0	0	0	0	0	0	0	0	0	0	0
	longitude	0	0	0	0	0	0	0	0	0	0	0
	location	0	0	0	0	0	0	0	0	0	0	0
	dtype: int64											

Read or convert the columns 'Created Date' and 'Closed Date' to datetime datatype and create a new column 'Request\_Closing\_Time' as the time elapsed between request creation and request closing. (Hint: Explore the package/module datetime)

```
In [10]: df['closed_date'] = pd.to_datetime(df['closed_date'])
```

```
In [11]: df['created_date'] = pd.to_datetime(df['created_date'])
```

```
In [12]: requested_closing_time = df['closed_date'] - df['created_date']
```

```
In [13]: df['requested_closing_time'] = requested_closing_time
```

```
In [14]: # converting the date difference into minutes
```

```
requested_closing_time_min = df['requested_closing_time']/np.timedelta64(1,'m')
```

```
In [15]: df['requested_closing_time_min'] = requested_closing_time_min
```

```
In [16]: df.head()
```

	unique_key	created_date	closed_date	agency	agency_name	complaint_type	descriptor	location_type	incident_zip	incident_address	bridge_highway_name	id
0	32310363	2015-12-31 23:59:45	2016-01-01 00:55:00	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	10034.0	71 VERMILY AVENUE	...	NaN
1	32309934	2015-12-31 23:59:44	2016-01-01 01:26:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	11105.0	27-07 23 AVENUE	...	NaN
2	32309159	2015-12-31 23:59:29	2016-01-01 04:51:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	10458.0	2897 VALENTI AVENUE	...	NaN
3	32305098	2015-12-31 23:57:46	2016-01-01 07:43:00	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	10461.0	2940 BAISLEY AVENUE	...	NaN
4	32306529	2015-12-31 23:56:58	2016-01-01 03:24:00	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	11373.0	87-14 57 ROAD	...	NaN

5 rows × 55 columns

```
In [17]: df.columns
```

Out[17]:	Index(['unique_key', 'created_date', 'closed_date', 'agency', 'agency_name', 'complaint_type', 'descriptor', 'location_type', 'incident_zip', 'incident_address', 'street_name', 'cross_street_1', 'cross_street_2', 'intersection_street_1', 'intersection_street_2', 'address_type', 'city', 'landmark', 'facility_type', 'status', 'due_date', 'resolution_description', 'resolution_action_updated_date', 'community_board', 'borough', 'x_coordinate_state_plane', 'y_coordinate_state_plane', 'park_facility_name', 'park_borough', 'school_name', 'school_number', 'school_region', 'school_code', 'school_phone_number', 'school_not_found', 'school_or_citywide_complaint', 'vehicle_type', 'taxi_company_borough', 'taxi_pick_up_location', 'bridge_highway_name', 'bridge_highway_direction', 'road_ramp', 'bridge_highway_segment', 'garage_lot_name', 'ferry_direction', 'ferry_terminal_name', 'latitude', 'longitude', 'location', 'requested_closing_time', 'requested_closing_time_min'], dtype='object')
----------	--

```
In [18]: # creating new dataframe
new_df = df[['unique_key','created_date','closed_date','agency', 'agency_name','complaint_type','descriptor','location_type','address_type','city','requested_closing_time','requested_closing_time_min']]
```

```
In [19]: new_df.head()
```

	unique_key	created_date	closed_date	agency	agency_name	complaint_type	descriptor	location_type	address_type	city
0	32310363	2015-12-31 23:59:45	2016-01-01 00:55:00	NYPD	New York City Police Department	Noise - Street/Sidewalk	Loud Music/Party	Street/Sidewalk	ADDRESS	NEW YORK
1	32309934	2015-12-31 23:59:44	2016-01-01 01:26:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	ADDRESS	ASTORIA
2	32309159	2015-12-31 23:59:29	2016-01-01 04:51:00	NYPD	New York City Police Department	Blocked Driveway	No Access	Street/Sidewalk	ADDRESS	BRONX
3	32305098	2015-12-31 23:57:46	2016-01-01 07:43:00	NYPD	New York City Police Department	Illegal Parking	Commercial Overnight Parking	Street/Sidewalk	ADDRESS	BRONX
4	32306529	2015-12-31 23:56:58	2016-01-01 03:24:00	NYPD	New York City Police Department	Illegal Parking	Blocked Sidewalk	Street/Sidewalk	ADDRESS	ELMHURST

Converting the series city and borough into lower case to avoid redundancies

```
In [20]: new_df['city'] = new_df['city'].str.lower()
```

```
pd.options.mode.chained_assignment = None
```

```
In [21]: new_df['borough'] = new_df['borough'].str.lower()
```

```
pd.options.mode.chained_assignment = None
```

```
In [22]: new_df.head()
```

```
new_df.shape
```

```
Out[22]: (300698, 17)
```

```
In [23]: new_df.duplicated().sum()
```

```
Out[23]: 0
```

```
In [24]: new_df.isnull().sum()
```

Out[24]:	unique_key	created_date	closed_date	agency	agency_name	complaint_type	descriptor	location_type	address_type	city
	0	0	0	0	0	0	0	0	0	0
	created_date	0	0	0	0	0	0	0	0	0
	closed_date	0	0	0	0	0	0	0	0	0
	agency	0	0	0	0	0	0	0	0	0
	agency_name	0	0	0	0	0	0	0	0	0
	complaint_type	0	0	0	0	0	0	0	0	0
	descriptor	0	0	0	0	0	0	0	0	0
	location_type	0	0	0	0	0	0	0	0	0
	address_type	0	0	0	0	0	0	0	0	0
	city	0	0	0	0	0	0	0	0	0
	facility_type	0	0	0	0	0	0	0	0	0
	status	0	0	0	0	0	0	0	0	0
	due_date	0	0	0	0	0	0	0	0	0
	borough	0	0	0	0	0	0	0	0	0
	location	0	0	0	0	0	0	0	0	0
	requested_closing_time	0	0	0	0	0	0	0	0	0
	requested_closing_time_min	0	0	0	0	0	0	0	0	0
	dtype: int64									

```
In [25]: pd.set_option('mode.chained_assignment', None)
new_df.dropna(inplace=True)
```

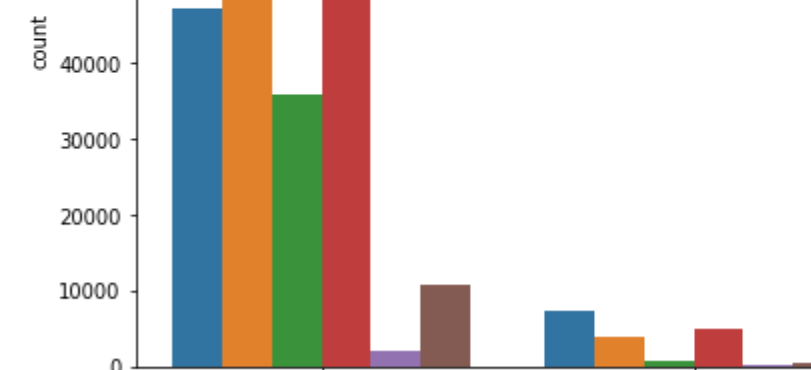
```
In [26]: new_df.isnull().sum()
```

Out[26]:	unique_key	created_date	closed_date	agency	agency_name	complaint_type	descriptor	location_type	address_type	city
	0	0	0	0	0	0	0	0	0	0
	created_date	0	0	0	0	0	0	0	0	0
	closed_date	0	0	0	0	0	0	0	0	0
	agency	0	0	0	0	0	0	0	0	0
	agency_name	0	0	0	0	0	0	0	0	0
	complaint_type	0	0	0	0	0	0	0	0	0
	descriptor	0	0	0	0	0	0	0	0	0
	location_type	0	0	0	0	0	0	0	0	0
	address_type	0	0	0	0	0	0	0	0	0
	city	0	0	0	0	0	0	0	0	0
	facility_type	0	0	0	0	0	0	0	0	0
	status	0	0	0	0	0	0	0	0	0
	due_date	0	0	0	0	0	0	0	0	0
	borough	0	0	0	0	0	0	0	0	0
	location	0	0	0	0	0	0	0	0	0
	requested_closing_time	0	0	0	0	0	0	0	0	0
	requested_closing_time_min	0	0	0	0	0	0	0	0	0
	dtype: object									

Provide major insights/patterns that you can offer in a visual format (graphs or tables); at least 4 major conclusions that you can come up with after generic data mining.

1)

```
In [28]: df['location_type'].value_counts()
sns.countplot(df['location_type'])
plt.set_xticklabels(plot.get_xticklabels(),rotation=90)
plt.show()
```



Insight:

- Street/ SideWalk is the most used Location Type by New York people. From the plot, we can see it out-counts most other types

2)

```
In [29]: plt.figure(figsize=(12,4))
plt2=sns.countplot(df['descriptor'])
plt2.set_xticklabels(plot2.get_xticklabels(),rotation=90)
plt.show()
```