

Problem objective:

Perform a service request data analysis of New York City 311 calls. You will focus on the data wrangling techniques to understand the pattern in the data and also visualize the major complaint types.

In [] : `new importing necessary libraries of python`

In [3] : `import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns`

Task 1: Importing the 311 NYC Service request dataset using pd.read_csv

In [2] : `data = pd.read_csv('311_Service_Requests_from_2010_to_Present.csv')
#User's\shivam\conda3\lib\site-packages\IPythonCore\interactiveshell.py:3146: DtypeWarning: Columns (48,49) have mixed types.Specify dtype option on import or set LowMemory=False
has_raised = await self.run_ast_nodes(code_ast.body, cell_name,`

In [3] : `data.head(10)`

Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Incident Address	Bridge Highway Segment	Bridge Direction	Road Ramp	Bridge Segment	Garage Name	Ferry Direction	Ferry Terminal Name	Latitude	Longitude
0	3233063	12/23/2015 11:59:40 AM	01-01	NYPD	New York City Police Department	Noise - StreetSidewalk	Local	10034.0	VERMILION AVENUE	71	...	NaN	NaN	NaN	NaN	NaN	40.86566	-73.92301
1	3230984	11/29/14 11:59:44 PM	01-01	NYPD	New York City Police Department	Blocked Driveway	No Access	StreetSidewalk	11105.0	2707 23 AVENUE	...	NaN	NaN	NaN	NaN	NaN	40.77946	-73.91504
2	3230919	11/29/15 11:59:25 PM	01-01	NYPD	New York City Police Department	Blocked Driveway	No Access	StreetSidewalk	10458.0	VALENTINE AVENUE	...	NaN	NaN	NaN	NaN	NaN	40.87025	-73.88825
3	3230908	11/29/2015 11:57:40 PM	01-01	NYPD	New York City Police Department	Blocked Driveway	No Access	StreetSidewalk	10461.0	2540 BAMELEY AVENUE	...	NaN	NaN	NaN	NaN	NaN	40.83994	-73.82879
4	3230529	12/1/2015 11:56:33 PM	01-01	NYPD	New York City Police Department	Illegal Parking	Blocked Driveway	StreetSidewalk	11373.0	87-147 ROAD	...	NaN	NaN	NaN	NaN	NaN	40.73300	-73.87410
5	3230558	12/1/2015 11:57:30 PM	01-01	NYPD	New York City Police Department	Illegal Parking	Blocked Driveway	StreetSidewalk	11215.0	200 21 STREET	...	NaN	NaN	NaN	NaN	NaN	40.60823	-73.89258
6	3230559	12/1/2015 11:56:30 PM	01-01	NYPD	New York City Police Department	Illegal Parking	Blocked Driveway	StreetSidewalk	10022.0	524 WEST 169 STREET	...	NaN	NaN	NaN	NaN	NaN	40.84048	-73.82735
7	3230709	11/24/15 11:54:01 AM	01-01	NYPD	New York City Police Department	Blocked Driveway	No Access	StreetSidewalk	10457.0	501 EAST 171 STREET	...	NaN	NaN	NaN	NaN	NaN	40.83763	-73.90205
8	3230551	12/1/2015 11:53:20 AM	01-01	NYPD	New York City Police Department	Illegal Parking	Blocked Driveway	StreetSidewalk	11415.0	83-44 BOULEVARD	...	NaN	NaN	NaN	NaN	NaN	40.70497	-73.82005
9	3230901	12/1/2015 11:53:30 AM	01-01	NYPD	New York City Police Department	Blocked Driveway	No Access	StreetSidewalk	11219.0	1408 66 STREET	...	NaN	NaN	NaN	NaN	NaN	40.82793	-73.99938

10 rows x 19 columns

In [4] : `data.shape`

Out[4] : (398698, 19)

There are 390698 rows and 19 columns

In [9] : `data.describe()`

Out[9] :

Unique Key	Incident Zip	X Coordinate (State Plane)	Y Coordinate (State Plane)	School or Citywide Complaint	Vehicle Type	Taxi Company Borough	Taxi Pick Up Location	Garage Lot Name	Latitude	Longitude
count	3.00000e+05	2.98000e+000	2.97150e+000	297150.000000	0.0	0.0	0.0	0.0	297150.00000	297150.00000
mean	3.13005e+07	1.084388e+05	1.00485e+06	203754.534410	NaN	NaN	NaN	NaN	40.72085	-73.92630
std	3.78867e+07	60.1000e3	2.17330e+04	29800.000000	NaN	NaN	NaN	NaN	0.30022	0.07664
min	3.02704e+07	83.000000	9.13670e+05	121210.000000	NaN	NaN	NaN	NaN	40.49915	-74.25487
25%	3.06033e+07	1.0310.000000	9.91870e+05	183340.000000	NaN	NaN	NaN	NaN	40.66976	-73.87342
50%	3.13005e+07	1.02000.00000	1.00315e+06	201110.000000	NaN	NaN	NaN	NaN	40.71861	-73.93211
75%	3.17846e+07	1.128.000000	1.01837e+06	224125.200000	NaN	NaN	NaN	NaN	40.78140	-73.87605
max	3.23106e+07	1.1697.000000	1.067173e+06	271870.000000	NaN	NaN	NaN	NaN	40.82869	-73.70070

In [24] : `data.info()`

<class 'pandas.core.frame.DataFrame'>
Int64Index: 398697 entries, 0 to 398697
Data columns (total 19 columns):
Column Non-Null Count Dtype
-- -- --
0 unique Key 398697 non-null int64
1 Created Date 398697 non-null object
2 Closed Date 398693 non-null object
3 Agency 398697 non-null object
4 Agency Name 398697 non-null object
5 Complaint Type 398697 non-null object
6 Descriptor 398693 non-null object
7 Location Type 398698 non-null object
8 Incident Zip 398693 non-null float64
9 Incident Address 398698 non-null object
10 Street Name 256288 non-null object
11 Cross Street 1 21413 non-null object
12 Cross Street 2 250919 non-null object
13 Intersection Street 1 43362 non-null object
14 Intersection Street 2 43362 non-null object
15 Address Type 397883 non-null object
16 City 398693 non-null object
17 Landmark 3449 non-null object
18 Facility Type 298827 non-null object
19 Status 398697 non-null object
20 Due Date 398694 non-null object
21 Resolution Description 398697 non-null object
22 Resolution Action Updated Date 298818 non-null object
23 Community Board 398697 non-null object
24 Borough 398697 non-null object
25 Y Coordinate (State Plane) 297150 non-null float64
26 X Coordinate (State Plane) 297150 non-null float64
27 Park Borough 398697 non-null object
28 Park Facility Name 398697 non-null object
29 School Name 398697 non-null object
30 School Number 398697 non-null object
31 School Region 398697 non-null object
32 School Code 398697 non-null object
33 School Phone Number 398697 non-null object
34 School Address 398697 non-null object
35 School City 398697 non-null object
36 School State 398697 non-null object
37 School Not Found 398697 non-null object
38 School or Citywide Complaint 398697 non-null object
39 Vehicle Type 0 non-null float64
40 Taxi Company Borough 0 non-null float64
41 Taxi Pick Up Location 0 non-null float64
42 Bridge Highway Name 243 non-null object
43 Bridge Highway Direction 243 non-null object
44 Bridge Highway Segment 243 non-null object
45 Road Ramp 213 non-null object
46 Bridge Highway Segment 213 non-null object
47 Garage Lot Name 0 non-null float64
48 Ferry Direction 2 non-null object
49 Ferry Terminal Name 398697 non-null object
50 Latitude 297158 non-null float64
51 Longitude 297158 non-null float64
52 Location 297158 non-null object
dtypes: float64(18), int64(1), object(42)
memory usage: 123.4+ MB

In [12] : `data.columns`

Out[12] :

Index(['Unique Key', 'Created Date', 'Closed Date', 'Agency', 'Agency Name', 'Complaint Type', 'Descriptor', 'Location Type', 'Incident Zip', 'Incident Address', 'Bridge Highway Segment', 'Bridge Direction', 'Road Ramp', 'Bridge Segment', 'Garage Name', 'Ferry Direction', 'Ferry Terminal Name', 'Latitude', 'Longitude'],
dtype='object')

In [13] : `data.isnull().sum()`

Out[13] :

Unique Key 0
Created Date 2164
Closed Date 0
Agency 0
Agency Name 0
Complaint Type 0
Descriptor 5914
Location Type 13
Incident Zip 2635
Incident Address 256288
Street Name 44430
Cross Street 1 49779
Cross Street 2 25540
Intersection Street 1 26738
Intersection Street 2 26738
Address Type 2615
City 2614
Landmark 38849
Facility Type 2271
Status 3
Due Date 3
Resolution Description 398697
Resolution Action Updated Date 2387
Community Board 0
Borough 0
Y Coordinate (State Plane) 3540
X Coordinate (State Plane) 3540
Park Borough 0
Park Facility Name 0
School Name 0
School Number 0
School Region 1
School Code 1
School Phone Number 0
School Address 0
School City 0
School State 0
School Zip 0
School Not Found 0
School or Citywide Complaint 398697
Vehicle Type 0
Taxi Company Borough 0
Taxi Pick Up Location 0
Bridge Highway Name 243
Bridge Highway Direction 243
Road Ramp 213
Bridge Highway Segment 213
Garage Lot Name 0
Ferry Direction 2
Ferry Terminal Name 398697
Latitude 3540
Longitude 3540
Location 3540
dtype: int64

In [16] : `data[Complaint Type].unique()`

Out[16] :

array(['Noise - StreetSidewalk', 'Blocked Driveway', 'Illegal Parking', 'Derelict Vehicle', 'Noise - Commercial', 'Noise - House of Worship', 'Noise - Vehicle', 'Noise - Abuse', 'Vandalism', 'Traffic', 'Drinking', 'Bike/Broller/Skate Chronic', 'Panhandling', 'Noise - Park', 'Homeless Encampment', 'Urinating in Public', 'Graffiti', 'Disorderly Youth', 'Illegal Fireworks', 'Ferry Complaints', 'Agency Issues', 'Sequece', 'Animal in a Park'], dtype=object)

In [18] : `data[Complaint Type].value_counts()`

Out[18] :

Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Incident Address	Bridge Highway Segment	Bridge Direction	Road Ramp	Bridge Segment	Garage Name	Ferry Direction	Ferry Terminal Name	Latitude	Longitude	Location
28332	2047220	04/28/2015 10:05:20 AM	05-05	NYPD	New York City Department of Transportation	Animal in a Park	Park	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN

1 rows x 19 columns

In [21] : `data.drop(labels=[28332], axis=0, inplace=True)`

In [22] : `data.shape`

Out[22] : (398697, 19)

Earlier it was 390698 rows and 19 columns and after dropping one row, now it is 390697 rows and 19 columns.

Task 2: 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

In [26] : `data['Created Date'] = pd.to_datetime(data['Created Date'])
data['Closed Date'] = pd.to_datetime(data['Closed Date'])`

In [27] : `data.info()`

<class 'pandas.core.frame.DataFrame'>
Int64Index: 398697 entries, 0 to 398697
Data columns (total 19 columns):
Column Non-Null Count Dtype
-- -- --
0 unique Key 398697 non-null int64
1 Created Date 398697 non-null datetime64[ns]
2 Closed Date 398693 non-null datetime64[ns]
3 Agency 398697 non-null object
4 Agency Name 398697 non-null object
5 Complaint Type 398697 non-null object
6 Descriptor 398693 non-null object
7 Location Type 398698 non-null object
8 Incident Zip 398693 non-null float64
9 Incident Address 398698 non-null object
10 Street Name 256288 non-null object
11 Cross Street 1 21413 non-null object
12 Cross Street 2 250919 non-null object
13 Intersection Street 1 43362 non-null object
14 Intersection Street 2 43362 non-null object
15 Address Type 397883 non-null object
16 City 398693 non-null object
17 Landmark 3449 non-null object
18 Facility Type 298827 non-null object
19 Status 398697 non-null object
20 Due Date 398694 non-null object
21 Resolution Description 398697 non-null object
22 Resolution Action Updated Date 298818 non-null object
23 Community Board 398697 non-null object
24 Borough 398697 non-null object
25 Y Coordinate (State Plane) 297150 non-null float64
26 X Coordinate (State Plane) 297150 non-null float64
27 Park Borough 398697 non-null object
28 Park Facility Name 398697 non-null object
29 School Name 398697 non-null object
30 School Number 398697 non-null object
31 School Region 398697 non-null object
32 School Code 398697 non-null object
33 School Phone Number 398697 non-null object
34 School Address 398697 non-null object
35 School City 398697 non-null object
36 School State 398697 non-null object
37 School Not Found 398697 non-null object
38 School or Citywide Complaint 398697 non-null object
39 Vehicle Type 0 non-null float64
40 Taxi Company Borough 0 non-null float64
41 Taxi Pick Up Location 0 non-null float64
42 Bridge Highway Name 243 non-null object
43 Bridge Highway Direction 243 non-null object
44 Bridge Highway Segment 243 non-null object
45 Road Ramp 213 non-null object
46 Bridge Highway Segment 213 non-null object
47 Garage Lot Name 0 non-null float64
48 Ferry Direction 2 non-null object
49 Ferry Terminal Name 398697 non-null object
50 Latitude 297158 non-null float64
51 Longitude 297158 non-null float64
52 Location 297158 non-null object
dtypes: datetime64[ns](2), float64(18), int64(1), object(49)
memory usage: 123.9+ MB

As you can see, data column 1 and 2 i.e. created date and closed date earlier it was object type and now they are of datetime datatype.

In [30] : `#Creating a new column Request_Closing_Time
data['Request_Closing_Time'] = data['Closed Date'] - data['Created Date']`

In [31] : `data['Request_Closing_Time'].min()`

Out[31] :

data['Request_Closing_Time'].min()

In [32] : `data.head()`

Out[32] :

Unique Key	Created Date	Closed Date	Agency	Agency Name	Complaint Type	Descriptor	Location Type	Incident Zip	Incident Address	Bridge Highway Segment	Bridge Direction	Road Ramp	Bridge Segment	Garage Name	Ferry Direction	Ferry Terminal Name	Latitude	Longitude	Location
0	3233063	2015-12-23 11:59:40	01-01	NYPD	New York City Police Department	Noise - StreetSidewalk	Local	10034.0	VERMILION AVENUE	71	...	NaN	NaN	NaN	NaN	NaN	40.86566	-73.92301	40.86566156303767, -73.92300959371144

1 rows x 19 columns

In [34] : `data.drop(labels=[28332], axis=0, inplace=True)`

In [35] : `data.shape`

Out[35] : (398697, 19)

As you can see here earlier there is 390697 rows and 19 columns and now it is 390697 rows and 19 columns after adding Request_Closing_Time and Request_Closing_Time.

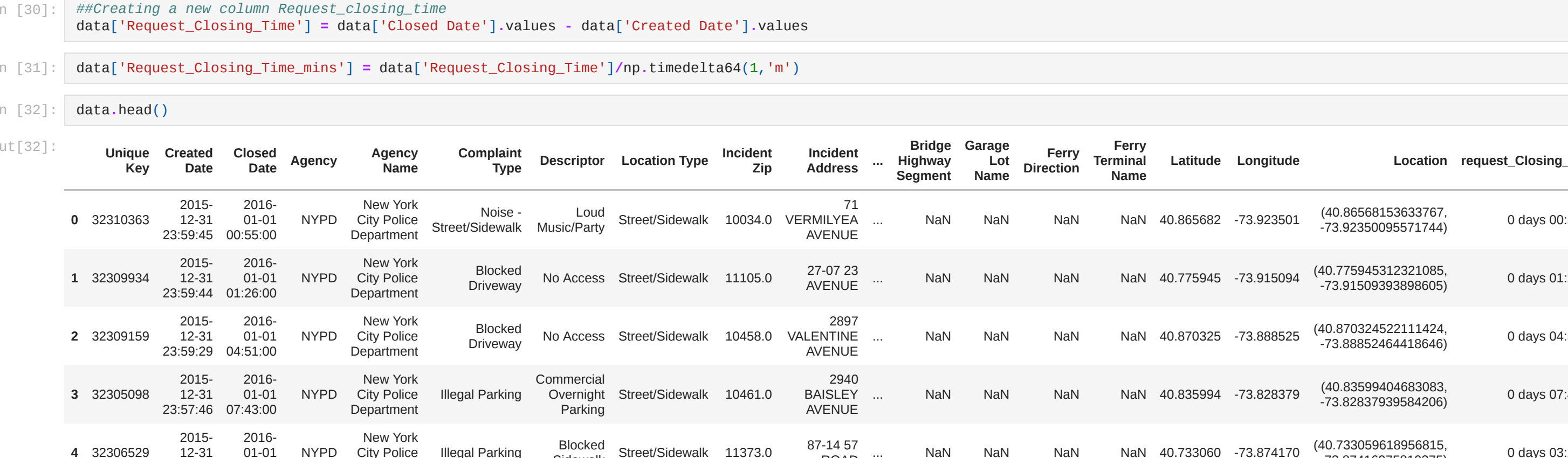
Task 3: Provide major insights/patterns that you can offer in a visual format (graphs or tables) at least 4 major conclusions that you can come up with after getting data mining.

In [40] : `#Visualizing the dataset with suitable graphs.
Visualization 1`

Out[40] :

data.groupby(['City', 'Complaint Type']).size().unstack().fillna(0)

In [39] : `data.plot.bar(figsize=(15,20), stacked=True)
plt.title('Number of complaints in different cities')
plt.xlabel('City')
plt.ylabel('Number of Complaints')
plt.show()`



Conclusion 1: City with maximum number of complaints was found to be Brooklyn having different types of complaints.

In [42] : `#Visualization 2`

Out[42] :

data.groupby(['City', 'Complaint Type']).size().unstack().fillna(0)

In [43] : `data.groupby(['City', 'Complaint Type']).size().unstack().fillna(0)`

Out[43] :

data.groupby(['City', 'Complaint Type']).size().unstack().fillna(0)

Conclusion 2: As we can see here that Blocked Driveway is the most lodged complaints followed by Illegal Parking.

In [61] : `#Checking the status of cases in the top five cities.
data_3 = data.groupby(['City', 'Status']).size().unstack().fillna(0)`

Out[61] :

data_3 = data.groupby(['City', 'Status']).size().unstack().fillna(0)

In [58] : `data_3.sort_values(by='Open', ascending=False).head()`

Out[58] :

Status	Assigned	Closed	Open	Unresolved	percentage
NEW YORK	15.0	0.0	1.0	27.602308	
BROOKLYN	15.0	0.0	1.0	18.461538	
JAMAICA	1.0	0.0	0.0	13.846154	
STATEN ISLAND	2.0	0.0	0.0	9.230769	
BROXN	6.0	0.0	0.0	7.602308	

In [59] : `data_3['Unresolved_percentage'] = data_3['Open'] / data_3['Open'] * 100`

Out[59] :

data_3['Unresolved_percentage'] = data_3['Open'] / data_3['Open'] * 100

In [60] : `data_3.sort_values(by='Unresolved_percentage', ascending=False).head().sum()`

Out[60] :

data_3.sort_values(by='Unresolved_percentage', ascending=False).head().sum()

Conclusion 3: As we can see here that the maximum number of pending cases were found in New York. Maximum Number of complaints has been resolved and the maximum unresolved cases were found in the top 5 cities having maximum number of complaints. Total Percentage of unresolved cases is found to be 79.92.

In [68] : `#Visualization 3
data['Location Type'].fillna(value='StreetSidewalk', inplace=True)`

Out[68] :

data['Location Type'].fillna(value='StreetSidewalk', inplace=True)

In [71] : `plt.figure(figsize=(10,5))
plt.title('Average Request Closing Time For Boroughs')
plt.xlabel('Average Request Closing Time For Boroughs')
plt.ylabel('Request Closing Time (mins)')
plt.show()`

Out[71] :

plt.figure(figsize=(10,5))

Conclusion 4: The maximum cases were located to StreetSidewalk and it can be resolved if the traffic police strictly look into this matter.

In [93] : `#Visualization 4`

Out[93] :

plt.figure(figsize=(8,7))

Conclusion 5: As we observed here that, the maximum average resolving time of complaint for different Boroughs was found in Unspecified location.

In [95] : `data_4 = data.groupby(['Complaint Type', 'Request_Closing_Time_mins']).size().unstack().fillna(0)`

Out[95] :

data_4 = data.groupby(['Complaint Type', 'Request_Closing_Time_mins']).size().unstack().fillna(0)

In [96] : `data_4.head()`

Out[96] :

Complaint Type	Request_Closing_Time_mins	size
Agency Issues	315.034545	
Animal Abuse	112.796009	
Blocked Driveway	225.987778	
Derelict Vehicle	284.454331	

Conclusion 6: As we can see here average request closing time for different complaints.

Task 4: Order the complaint types based on the average 'Request_Closing_Time', grouping them for different locations

In [100] : `data_5 = data.groupby(['City', 'Complaint Type']).size().unstack().fillna(0)`

Out[100] :

Complaint Type	Animal Abuse	Blocked Driveway	Derelict Vehicle	Disorderly Youth	Drinking	Graffiti	Homeless Encampment	Illegal Parking	Commercial	Noise - House of Worship	Noise - StreetSidewalk	Noise - Vehicle	Panhandling	Urinating in Public	Vandalism	Bike/Broller/Skate Chronic
City																
MANHATTAN	176.00210	215.47000	14.216667	92.000000	108.887500	138.303889	137.233333	92.733333

Conclusion 5: As we observed here that, the maximum average resolving time of complaint for different Boroughs was found in Unspecified location.

In [105] : `data_5 = data.groupby(['Complaint Type', 'Request_Closing_Time_mins']).size().unstack().fillna(0)`

Out[105] :

data_5 = data.groupby(['Complaint Type', 'Request_Closing_Time_mins']).size().unstack().fillna(0)

In [106] : `data_5.head()`

Out[106] :

Complaint Type	Request_Closing_Time_mins	size
Agency Issues	315.034545	
Animal Abuse	112.796009	
Blocked Driveway	225.987778	
Derelict Vehicle	284.454331	

Conclusion 6: As we can see here average request closing time for different complaints.

Task 5: Perform a statistical test for the following Please note: For the below statements you need to state the Null and Alternate and then provide a statistical test to accept or reject the Null Hypothesis along with the corresponding p-value.

• Whether the average response time across complaint types is similar or not (overall)
• Are the types of complaint or service requested also location related?

1)ANOVA Analysis (Checking for top 5 complaints)

• Null Hypothesis: The average response time across complaint types is not similar.
• Alternate Hypothesis: The average response time across complaint types is similar.