

INTRODUCTION

```
# Name :- Rohit Yadav
# Py. No. :- 8311655005
# Email :- rohit.yadav.py94@gmail.com
```

Performing a service request data analysis of New York City 311 calls

In [130]# Importing Some Libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
import scipy.stats as stats

In [131]nyc_ser = pd.read_csv('311_Service_Requests_from_2010_to_Present.csv') # Importing the dataset
Data Columns (total 53 columns):
Column Non-Null Count Dtype
0 Unique Key 300698 non-null object
1 Created Date 300698 non-null object
2 Closed Date 300698 non-null object
3 Agency 300698 non-null object
4 Agency Name 300698 non-null object
5 Complaint Type 300698 non-null object
6 Descriptor 300698 non-null object
7 Location Type 300698 non-null object
8 Incident Zip 300698 non-null float64
9 Incident Address 300698 non-null object
10 ... 100698 non-null object
11 Bridge Highway Direction 300698 non-null object
12 Bridge Highway Segment 300698 non-null object
13 Road Ramp 300698 non-null object
14 Bridge Highway Segment 300698 non-null object
15 Garage Lot Name 300698 non-null object
16 Ferry Direction 300698 non-null object
17 Ferry Terminal Name 300698 non-null object
18 Latitude 300698 non-null float64
19 Longitude 300698 non-null float64
20 Location 300698 non-null object
21 Request_Closing_Time 300698 non-null object
dtype: object

In [132]nyc_ser.head() # shows first 5 rows

In [133]nyc_ser.columns # shows all columns

In [134]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', 'Community Board', 'X Coordinate (State Plane)', 'Y Coordinate (State Plane)', 'Park Facility Name', 'Park Borough', 'School Name', 'School Number', 'School Address', 'School City', 'School State', 'School Zip', 'School Not Found', 'School Region', 'School Code', 'Vehicle Type', 'Taxi Company Borough', 'Taxi Pick Up Location', 'Bridge Highway Name', 'Bridge Highway Direction', 'Road Ramp', 'Ferry Direction', 'Ferry Terminal Name', 'Garage Lot Name', 'Ferry Direction', 'Request_Closing_Time'], dtype='object')

In [135]nyc_ser.info() # to show the datatype and non-null values of all columns

In [136]# Using Pandas converting Data type
nyc_ser['Created Date'] = pd.to_datetime(nyc_ser['Created Date'])
nyc_ser['Closed Date'] = pd.to_datetime(nyc_ser['Closed Date'])

In [137]# Creating new column
nyc_ser['Request_Closing_Time'] = (nyc_ser['Closed Date'] - nyc_ser['Created Date']).dt.total_seconds()

In [138]nyc_ser.head()

In [139]nyc_ser.describe(include='all')

In [140]# showing Unique values in Complaint type
nyc_ser['Complaint Type'].unique()

In [141]array(['Noise - Street/Sidewalk', 'Blocked Driveway', 'Illegal Parking', 'Derelict Vehicle', 'Noise - Commercial', 'Noise - House of Worship', 'Posting Advertisement', 'Noise - Vehicle', 'Animal Abuse', 'Vending', 'Traffic', 'Graffiti', 'Bike/Roller/Skate Chronic', 'Panhandling', 'Noise - Park', 'Homeless Encampment', 'Urinating in Public', 'Disorderly Youth', 'Illegal Fireworks', 'Ferry Complaint', 'Ferry Complaint', 'Agency Issues', 'Squeegee', 'Animal in a Park'], dtype=object)

In [142]# Making Bar plot
value_counts = nyc_ser['Complaint Type'].value_counts()
plt.figure(figsize=(16,5))
value_counts.plot(kind='bar')
plt.title('Most common complaints',fontsize=28)

In [143]Text(0.5, 1.8, 'Most common complaints')

In [144]# From above it is clear that Maximum complaints are from Complaint type (Blocked Driveway)

In [145]# Showing Top 10 Complaints
nyc_ser['Complaint Type'].value_counts().nlargest(10)

In [146]Blocked Driveway 77044
Illegal Parking 76981
Noise - Street/Sidewalk 48612
Noise - Commercial 35917
Derelict Vehicle 17718
Noise - Vehicle 17083
Animal Abuse 7778
Traffic 4498
Homeless Encampment 4416
Noise - Park 4042
Name: Complaint Type, dtype: int64

In [147]# showing null values in City
nyc_ser['City'].isnull().sum()

In [148]2614

In [149]nyc_ser['City'].unique()

In [150]array(['NEW YORK', 'ASTORIA', 'BROOK', 'ELMHURST', 'BROOKLYN', 'KING ALBANY', 'JANICA', 'SOUTH RICHMOND HILL', 'NEW RICHMOND', 'HOWARD BEACH', 'FOREST HILLS', 'STATEN ISLAND', 'ZONE PARK', 'RICHMOND HILL', 'WOODHAVEN', 'FLUSHING', 'CORONA', 'QUEENS VILLAGE', 'OAKLAND GARDENS', 'WELLS', 'WASHT', 'EAST ELMHURST', 'SOUTH ZONE PARK', 'WOODSIDE', 'FRESH MEADOWS', 'LONG ISLAND CITY', 'ROCKAWAY PARK', 'SPRINGFIELD GARDENS', 'COLLEGE POINT', 'BAYSIDE', 'GLEN DAKS', 'FAR ROCKAWAY', 'BELLEROSE', 'LITTLE NECK', 'CAMBRIA HEIGHTS', 'ROSELAND', 'SUNNYSIDE', 'WHITESTONE', 'ARVERNE', 'FLORAL PARK', 'NEW YORK PARK', 'CENTRAL PARK', 'BREEZY POINT', 'QUEENS', 'ASTORIA', 'Long Island City', 'Woodside', 'East Elmhurst', 'Howard Beach'], dtype=object)

In [151]# Replacing null values with 'Unknown City'
nyc_ser['City'].fillna('Unknown City',inplace=True)
nyc_ser['City'].isnull().sum()

In [152]0

In [153]# as now there are no null values

In [154]# grouping City and Complaint Type
nyc_ser.groupby(['City','Complaint Type']).size()

In [155]City Complaint Type
ARVERNE Animal Abuse 38
Blocked Driveway 35
Derelict Vehicle 2
Disorderly Youth 1
Drinking 1
Woodside Blocked Driveway 11
Derelict Vehicle 1
Illegal Parking 100
Noise - Commercial 1
Noise - Street/Sidewalk 5
Length: 782, dtype: int64

In [156]# Lets Figure out the major cities with higher complaints
per_City_Complaints = pd.crosstab(index=nyc_ser['City'], columns=nyc_ser['Complaint Type'])
print(per_City_Complaints.shape)

In [157](54, 24)

In [158]per_City_Complaints.head()

In [159]Complaint Type Agency Issues Animal Abuse Animal in a Park Bike/Roller/Skate Chronic Blocked Driveway Derelict Vehicle Disorderly Youth Drinking Ferry Complaint Graffiti ... Noise - House of Worship Noise - Park Noise - Street/Sidewalk Noise - Vehicle Panhandling Posting Advertisement Squeegee Traffic Urinating in Public Vending
ARVERNE 0 38 0 0 35 27 2 1 0 1 ... 11 2 29 7 1 0 0 0 1 1
ASTORIA 0 125 0 15 268 351 3 35 0 4 ... 19 61 386 204 1 1 0 47 9 54
Astoria 0 0 0 0 116 12 0 0 0 0 ... 0 0 114 0 0 0 0 0 0 0
BELLEROSE 0 37 0 0 377 198 1 1 0 3 ... 2 4 15 16 0 0 0 9 0 2
BELLEROSE 0 7 0 1 95 89 2 1 0 0 ... 1 1 13 10 1 1 0 7 1 0

In [160]# Making a stacked bar chart
per_City_Complaints.plot(kind='bar',stacked=True,figsize=(25,17))
plt.title('Complaints Per City',fontsize=28)

In [161]Text(0.5, 1.8, 'Complaints Per City')

In [162]Its observed that Brooklyn has the highest no. of complaints

In [163]# Lets analyze the major complaint types in 'Brooklyn'
nyc_ser_Brooklyn = nyc_ser[nyc_ser['City']=='BROOKLYN']
nyc_ser_Brooklyn.shape

In [164](9307, 54)

In [165]Brooklyn_values = nyc_ser_Brooklyn['Complaint Type'].value_counts()
plt.figure(figsize=(20,15))
Brooklyn_values.plot(kind='bar')
plt.title('Most Common Complaints in the city',fontsize=28)

In [166]Text(0.5, 1.8, 'Most Common Complaints in the city')

In [167]# Ordering the complaint types based on the average 'Request_Closing_Time', grouping them for different locations
nyc_ser['Request_Closing_Time'].isnull().sum()

In [168]2164

In [169]nyc_ser['Request_Closing_Time'].fillna(0,inplace = True)
nyc_ser['Request_Closing_Time'].isnull().sum()

In [170]0

In [171]nyc_avg_response_Time = nyc_ser.groupby(['City','Complaint Type']).Request_Closing_Time.mean()
nyc_avg_response_Time.head(15)

In [172]City Complaint Type
ARVERNE Animal Abuse 7753.852632
Blocked Driveway 8931.485734
Derelict Vehicle 16681.052361
Disorderly Youth 12928.580880
Drinking 8650.080900
Graffiti 5529.800880
Homeless Encampment 6523.250880
Illegal Parking 8335.13793
Noise - Commercial 8234.880880
Noise - House of Worship 5623.880880
Noise - Park 4520.680880
Noise - Street/Sidewalk 7375.620880
Noise - Vehicle 6895.571429
Panhandling 3720.880880
Urinating in Public 2491.880880
Name: Request_Closing_Time, dtype: float64

In [173]#Get the response Time across complaints
nyc_response_Time = nyc_ser.groupby(['Complaint Type']).Request_Closing_Time.mean().sort_values(ascending=True)
nyc_response_Time.head(20)

In [174]Ferry Complaint 0.080880
Posting Advertisement 7891.080154
Illegal Fireworks 1940.101139
Noise - Commercial 4217.680165
Noise - House of Worship 13471.378393
Noise - Park 22173.468162
Noise - Street/Sidewalk 12286.052292
Traffic 12489.731659
Disorderly Youth 12838.002080
Noise - Vehicle 12882.559790
Urinating in Public 13858.991554
Bike/Roller/Skate Chronic 13464.080880
Drinking 13846.291496
Vending 14523.502638
Squeegee 15639.418423
Homeless Encampment 15716.052336
Illegal Parking 14825.892519
Blocked Driveway 17615.421525
Animal Abuse 18743.545599
Name: Request_Closing_Time, dtype: float64

In [175]Perform a statistical test

In [176]#From the above data, its observed that the average response time across the complaint types are not equal.
#But, the following complaint types have response time which were too close.
#Panhandling - 15639.418423 & Homeless Encampment - 15716.052336
#Null hypothesis(H0) - Average response time across complaint types are equal
#Alternate hypothesis(H1) - Average response time across complaint types are not equal
Let's perform one way ANOVA for the above group of complaints

In [177]nyc_ser_Panhandling = nyc_ser[nyc_ser['Complaint Type']=='Panhandling']
nyc_ser_Homeless = nyc_ser[nyc_ser['Complaint Type']=='Homeless Encampment']
nyc_ser_Panhandling.head()

In [178]Request_Closing_Time
2076 137381.0
274 12730.0
346 616.0
7 810.0
3029 3297.0
9531 3297.0

In [179]nyc_ser_Homeless = nyc_ser[nyc_ser['Complaint Type']=='Homeless Encampment']
nyc_ser_Homeless = nyc_ser_Homeless.loc[:,['Request_Closing_Time']]
nyc_ser_Homeless.head()

In [180]Request_Closing_Time
382 2860.0
435 17110.0
459 7020.0
963 6206.0
966 11400.0

In [181]#value, pvalue = stats.f_oneway(nyc_ser_Panhandling,nyc_ser_Homeless)
pvalue

In [182]0.949815493

In [183]At p-value > 0.05, we will Accept Null hypothesis(H0) for Panhandling and Homeless Encampment

In [184]nyc_ser_city = pd.crosstab(nyc_ser['City'],nyc_ser['Complaint Type'])
nyc_ser_city.head()

In [185]Complaint Type Agency Issues Animal Abuse Animal in a Park Bike/Roller/Skate Chronic Blocked Driveway Derelict Vehicle Disorderly Youth Drinking Ferry Complaint Graffiti ... Noise - House of Worship Noise - Park Noise - Street/Sidewalk Noise - Vehicle Panhandling Posting Advertisement Squeegee Traffic Urinating in Public Vending
ARVERNE 0 38 0 0 35 27 2 1 0 1 ... 11 2 29 7 1 0 0 0 1 1
ASTORIA 0 125 0 15 268 351 3 35 0 4 ... 19 61 386 204 1 1 0 47 9 54
Astoria 0 0 0 0 116 12 0 0 0 0 ... 0 0 114 0 0 0 0 0 0 0
BELLEROSE 0 7 0 1 95 89 2 1 0 0 ... 1 1 13 10 1 1 0 7 1 0

In [186]5 rows x 24 columns

In [187]from scipy.stats import chi2_contingency

In [188]#contingency table
table = nyc_ser_city

In [189]#Get chi-square value , p-value, degrees of freedom, expected frequencies using the function chi2_contingency
stat, p, gsv, expected = chi2_contingency(table)

In [190]#chi2 significance value
alpha = 0.05

In [191]# Determine whether to reject or keep your null hypothesis
print('Significance=%.3f, p=%.3f' % (alpha, p))
if p < alpha:
print('Variables are associated (reject H0)')
else:
print('Variables are not associated(fail to reject H0)')
chi2_contingency(table)
Variables are associated (reject H0)

In []