

Table of contents

01

Introduction

02

Purpose

03

EDA

04

Data Cleaning

05

GroupBy, Crosstab 06

Visualization



Dataset Overview

The dataset comprises detailed information on various crime incidents reported within the city of Chicago from 2012 to 2017. Each record in the dataset represents a specific crime incident and includes attributes such as the unique identifier (ID) and case number assigned by the Chicago Police Department, the date and location of the incident, the type and description of the crime, arrest details, and geographic identifiers such as beat, district, ward, and community area.

Attributes

- ID: Unique identifier for each crime record.
- Case Number: Unique identifier assigned by the Chicago Police Department for each incident.
- Date: Date when the incident occurred, sometimes estimated.
- Block: Partially redacted address where the incident occurred, indicating the same block as the actual address.
- IUCR: Illinois Uniform Crime Reporting code linked to the Primary Type and Description
- Primary Type: Primary description of the crime.
- Description: Secondary description of the crimes, providing further details.
- Location Description: Description of where the incident occurred.
- · Arrest: Indicates whether an arrest was made in connection with the incident.

Attributes

- Domestic: Indicates if the incident was domestic-related according to Illinois Domestic Violence Act.
- · Beat: Indicates the smallest police geographic area where the incident occurred.
- District: Indicates the police district where the incident occurred.
- · Ward: City Council district where the incident occurred.
- · Community Area: Indicates the community area where the incident occurred.
- FBI Code: Federal Bureau of Investigation code, is a system used to classify different types of crimes based on their characteristics

Attributes

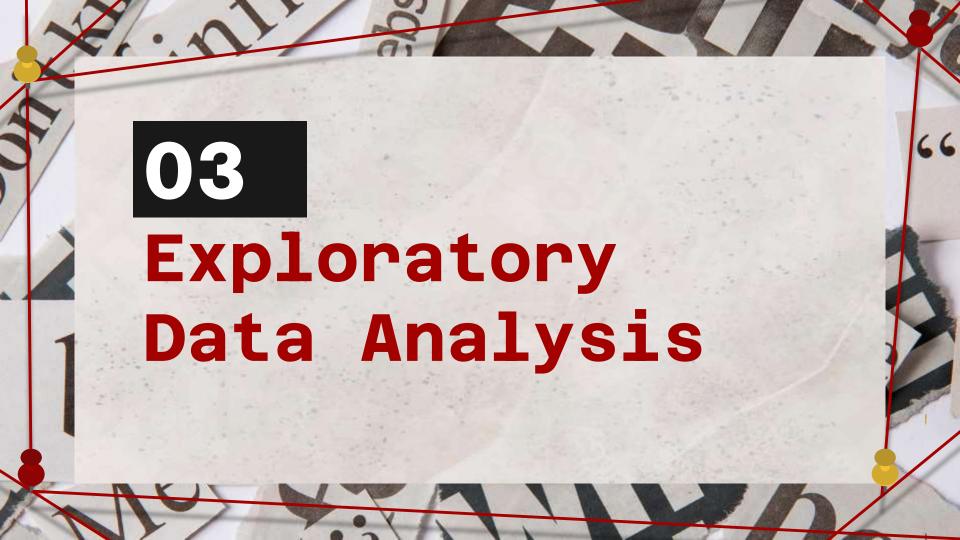
- X Coordinate: X-coordinate of the incident location in State Plane Illinois East NAD 1983 projection.
- Y Coordinate: Y-coordinate of the incident location in State Plane Illinois East NAD 1983 projection.
- · Year: Year when the incident occurred.
- Updated On: Date and time when the record was last updated.
- Latitude: Latitude of the incident location.
- · Longitude: Longitude of the incident location.
- Location: Location of the incident in a format suitable for geographic operations and mapping.



Purpose

Analyzing the Chicago Crime Dataset (2012-2017) helps us understand how crime changes over time and where it happens most often in the city. By looking at when and where crimes occur, we can figure out which areas need more attention from the police. We can also see how well the police are doing in catching criminals, especially in cases of domestic violence. This information helps city leaders make smarter decisions about how to keep people safe and reduce crime in Chicago.





- 1. Import the data
- 2. Displaying initial data

In [1]:	H	impo	ort pand ort numpy ort seab															
In [3]:	H	<pre>M df * pd.read_csv("CrimeData.csv") df.head()</pre>																
		ype	option (sha\AppDat on import : ead_csv("C	or set lo	w_memo			3344383926	0.py:1: Dty	peWarning:	Column	s (9,10)	have mixed	l type	s. Specif	/ di
Out[3]:		Unnamed: 0	ю	Case Number	Date	Block	IUCR	Primary Type	Description	Location Description	Arrest		Ward	Community Area	FBI Code	X Coordinate	Co
		0	3.0	10508693.0	HZ250496	05- 03- 2016 23:40	013XX S SAWYER AVE	486	BATTERY	DOMESTIC BATTERY SIMPLE	APARTMENT	True		24.0	29.0	08B	1154907.0	11
		1	89.0	10508695.0	HZ250409	05- 03- 2016 21:40	061XX S DREXEL AVE	486	BATTERY	DOMESTIC BATTERY SIMPLE	RESIDENCE	False	Jan-	20.0	42.0	088	1183066.0	.1
		2	197.0	10508697.0	HZ250503	05- 03- 2016 23:31	053XX W CHICAGO AVE	470	PUBLIC PEACE VIOLATION	RECKLESS CONDUCT	STREET	False		37.0	25.0	24	1140789.0	11
		3	673.0	10508698.0	HZ250424	05- 03- 2016 22-10	049XX W FULTON ST	460	BATTERY	SIMPLE	SIDEWALK	False	111	28.0	25.0	OBB	1143223.0	1
			911.0	10508699.0	HZ250455	05- 03-	003XX N LOTUS	820	THEFT	\$500 AND	RESIDENCE	False		28.0	25.0	6	1139890.0	11

3. Displaying last data

In [3]: N	<pre>df.dropna(inplace=True) df.tail()</pre>															
Out[3]:		Unnamed: 0	ID	Case Number	Date	Block	IUCR	Primary Type	Description	Location Description	Arrest	275	Ward	Community Area	FBI Code	Coord
	199994	2533684.0	8624881.0	HV298331	5/22/2012 17:30	016XX W 92ND PL	810	THEFT	OVER \$500	RESIDENCE- GARAGE	False	775	21.0	73.0	6	1166
	199995	2533685.0	8624882.0	HV298111	5/22/2012 15:40	133XX S LANGLEY AVE	460	BATTERY	SIMPLE	CHA PARKING LOT/GROUNDS	False	311	9.0	54.0	08B	1183
	199996	2533686.0	8624884.0	HV297965	5/22/2012 12:45	002XX W 110TH PL	820	THEFT	\$500 AND UNDER	RESIDENTIAL YARD (FRONT/BACK)	False	444	34.0	49.0	6	1176
	199997	2533687.0	8624885.0	HV298273	5/6/2012 12:00	106XX S HARDING AVE	486	BATTERY	DOMESTIC BATTERY SIMPLE	RESIDENCE- GARAGE	False	W	19.0	74.0	08B	1151
	199998	2533688.0	8624886.0	HV298416	5/22/2012 18:16	023XX N LOREL AVE	143A	WEAPONS VIOLATION	UNLAWFUL POSS OF HANDGUN	RESIDENTIAL YARD (FRONT/BACK)	True		37.0	19.0	15	1140
	5 rows ×	23 column	s													
	4		102	_	_											1

4. Statistical summary of numerical columns

Out[4]:		Unnamed: 0	ID	Beat	District	Ward	Community Area	X Coordinate	Y Coordinate	Year	Latitud
	count	1.924890e+05	1.924890e+05	192489.000000	192489.000000	192489.000000	192489.000000	1.924890e+05	1.924890e+05	192489.000000	192489.00000
	mean	1.873710e+06	9.104666e+06	1155.949857	11.213358	22.647481	37.888456	1.164457e+06	1.884956e+06	2013.196681	41,83991
	std	7.951918e+05	8,101823e+05	693.842145	6.898488	13.665015	21.405446	1.784708e+04	3.340824e+04	1.616702	0.09196
	min	3,000000e+00	2.085900e+04	111,000000	1.000000	1.000000	0.000000	0.000000e+00	0.000000e+00	2012.000000	36.61944
	25%	8.629250e+05	8.502760e+06	613.000000	6.000000	10.000000	23.000000	1.152440e+06	1.858253e+06	2012.000000	41.76628
	50%	2.436033e+06	8.581773e+06	1023.000000	10.000000	22.000000	32.000000	1.165915e+06	1.890190e+06	2012.000000	41.85439
	75%	2.485087e+06	9.953633e+06	1711.000000	16.000000	33.000000	58.000000	1.176388e+06	1,908404e+06	2015.000000	41.90440
	max	2.533688e+06	1.054914e+07	2535.000000	31.000000	50.000000	77.000000	1.205119e+06	1.951523e+06	2016.000000	42.02257

5. Displaying total number of rows and columns

6. Displaying total number of rows in each column

```
M df.count()
In [9]:
   Out[9]:
            Unnamed: 0
                                      192489
                                      192489
             Case Number
                                      192489
             Date
                                      192489
             Block
                                      192489
             IUCR
                                      192489
             Primary Type
                                      192489
             Description
                                      192489
             Location Description
                                      192489
             Arrest
                                      192489
             Domestic
                                      192489
             Beat
                                      192489
             District
                                      192489
            Ward
                                      192489
             Community Area
                                      192489
             FBI Code
                                      192489
             X Coordinate
                                      192489
             Y Coordinate
                                      192489
                                      192489
             Vear
             Updated On
                                      192489
             Latitude
                                      192489
             Longitude
                                      192489
             Location
                                      192489
             dtype: int64
```

7. Displaying Data types of each column and the number of non-null values

```
In [12]:
          M df.info()
             <class 'pandas.core.frame.DataFrame'>
             Int64Index: 192489 entries, 0 to 199998
             Data columns (total 23 columns):
                  Column
                                        Non-Null Count
                                                         Dtvpe
                  Unnamed: 0
                                        192489 non-null float64
                  ID
                                        192489 non-null float64
                  Case Number
                                        192489 non-null object
                                        192489 non-null object
                  Date
                  Block.
                                        192489 non-null object
                  IUCR
                                        192489 non-null object
                  Primary Type
                                        192489 non-null object
                  Description
                                        192489 non-null object
                  Location Description
                                       192489 non-null object
                                        192489 non-null object
                  Arrest
                  Domestic
                                        192489 non-null object
                  Beat
                                        192489 non-null float64
                  District
                                        192489 non-null float64
                                        192489 non-null float64
                  Community Area
                                        192489 non-null float64
                  FBI Code
                                        192489 non-null object
                  X Coordinate
                                        192489 non-null float64
                 Y Coordinate
                                        192489 non-null float64
                                        192489 non-null float64
                 Updated On
                                        192489 non-null object
                  Latitude
                                        192489 non-null float64
                 Longitude
                                        192489 non-null float64
              22 Location
                                        192489 non-null object
             dtypes: float64(11), object(12)
             memory usage: 35.2+ MB
```

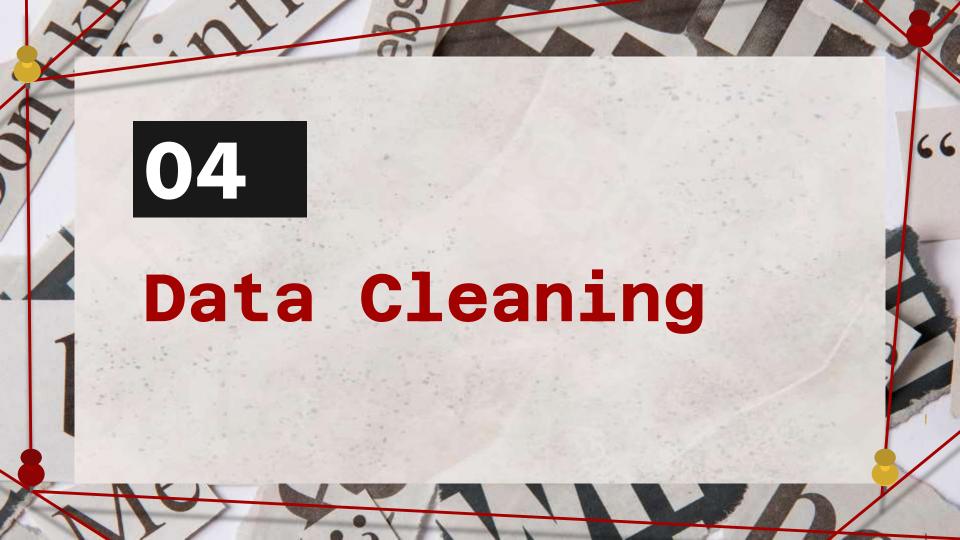
8. Displaying categorical and numerical columns

```
In [8]: M
    cat_col = [col for col in df.columns if df[col].dtype == 'object']
    print('Categorical columns:', cat_col)
    num_col = [col for col in df.columns if df[col].dtype != 'object']
    print('Numerical columns:', num_col)

Categorical columns: ['Case Number', 'Date', 'Block', 'IUCR', 'Primary Type', 'Description', 'Location Description', 'Arres t', 'Domestic', 'FBI Code', 'Updated On', 'Location']
    Numerical columns: ['Unnamed: 0', 'ID', 'Beat', 'District', 'Ward', 'Community Area', 'X Coordinate', 'Yea r', 'Latitude', 'Longitude']
```

9. Displaying number of unique values in each categorical columns

```
M df[cat_col].nunique()
Out[9]: Case Number
                                199993
                                 84770
        Date
        Block
                                  26157
        TUCR
                                   325
        Primary Type
                                    32
        Description
        Location Description
                                   107
        Arrest
        Domestic
        FBI Code
        Updated On
                                   246
        Location
                                110187
        dtype: int64
```



Data cleaning

1. Calculate the total number of missing values in each column

```
In [19]:
           M df.isnull().sum()
   Out[19]: Unnamed: 0
             Case Number
             Date
             Block
             IUCR
             Primary Type
             Description
             Location Description
             Arrest
             Domestic
             Beat
             District
             Ward
             Community Area
             FBI Code
             X Coordinate
             Y Coordinate
             Year
             Updated On
             Latitude
             Longitude
             Location
             dtype: int64
```

Data cleaning

2. Identify duplicate rows in a dataframe

```
M df.duplicated()
In [10]:
   Out[10]:
                       False
                       False
                       False
                       False
                       False
             199994
                       False
             199995
                       False
                       False
             199996
             199997
                       False
                       False
             199998
             Length: 192489, dtype: bool
          M df.duplicated().sum()
In [11]:
   Out[11]: 0
```

Data cleaning

3. Removing unwanted columns

Analysis of Categorical columns

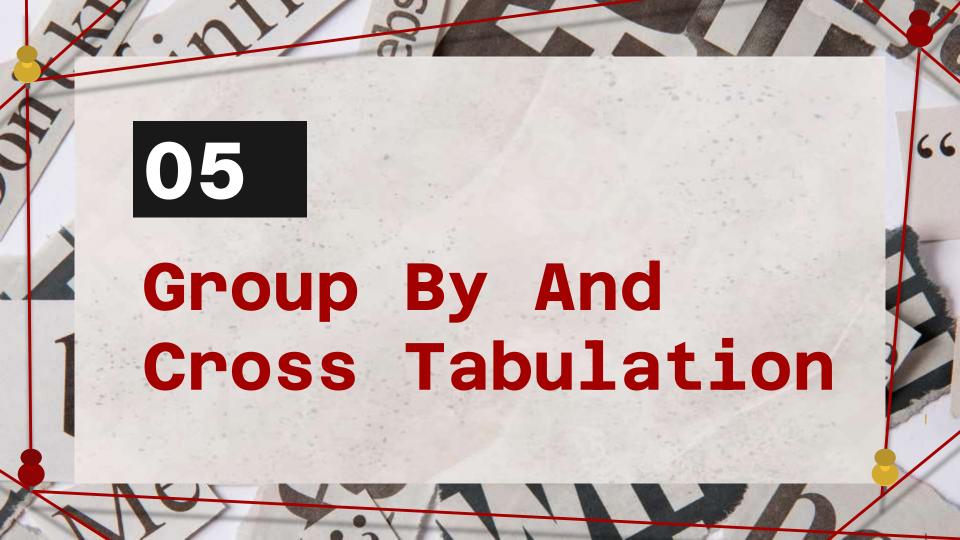
1. Count of each unique value in the 'Primary Type' column of the DataFrame, showing the frequency of Different types of crimes

```
M print(df['Primary Type'].value counts())
   Primary Type
   THEFT
                                          40390
   BATTERY
                                          35503
   NARCOTICS
                                          22137
   CRIMINAL DAMAGE
                                          20588
   ASSAULT
                                          12285
   OTHER OFFENSE
                                          11650
   BURGLARY
                                          10877
   MOTOR VEHICLE THEFT
                                           8849
                                           7953
   DECEPTIVE PRACTICE
   ROBBERY
                                           6876
   CRIMINAL TRESPASS
                                           5079
   WEAPONS VIOLATION
                                           2372
   PUBLIC PEACE VIOLATION
                                           1770
   OFFENSE INVOLVING CHILDREN
                                           1398
                                           1349
   PROSTITUTION
   INTERFERENCE WITH PUBLIC OFFICER
                                            825
   CRIM SEXUAL ASSAULT
                                            760
   SEX OFFENSE
                                            559
   LIQUOR LAW VIOLATION
                                            326
   ARSON
                                            288
   GAMBLING
                                            217
   KIDNAPPING
                                            129
   STALKING
                                            114
   INTIMIDATION
   HOMICIDE
   OBSCENITY
   PUBLIC INDECENCY
   NON-CRIMINAL
   CONCEALED CARRY LICENSE VIOLATION
   INDIAN TRAFFTENTING
```

Analysis of Categorical columns

1. Count of each unique value in the 'Location' column of the DataFrame, showing the frequency of Different types of locations where the crimes occurred

```
▶ print(df['Location Description'].value_counts())
  Location Description
  STREET
                               45363
  RESIDENCE
                               30508
  APARTMENT
                               24691
  SIDEWALK
                               21018
  OTHER
                                6878
  PORCH
  NEWSSTAND
  YARD
  PARKING LOT
   BARBER SHOP/BEAUTY SALON
  Name: count, Length: 107, dtype: int64
```



1. Displaying total number of arrets made for each type of crime

```
In [ ]:
             df.groupby('Primary Type')['Arrest'].sum()
  Out[61]:
            Primary Type
             ARSON
                                                       33
             ASSAULT
                                                     2956
             BATTERY
                                                     8213
             BURGLARY
                                                      617
             CONCEALED CARRY LICENSE VIOLATION
             CRIM SEXUAL ASSAULT
                                                      105
                                                     1463
             CRIMINAL DAMAGE
             CRIMINAL TRESPASS
                                                     3784
             DECEPTIVE PRACTICE
                                                     1263
             GAMBLING
                                                      217
             HOMICIDE
                                                       18
             HUMAN TRAFFICKING
                                                      783
             INTERFERENCE WITH PUBLIC OFFICER
             INTIMIDATION
                                                       25
             KIDNAPPING
                                                        9
             LIQUOR LAW VIOLATION
                                                      325
             MOTOR VEHICLE THEFT
                                                      567
             NARCOTICS
                                                    22093
             NON - CRIMINAL
                                                        0
             NON-CRIMINAL
                                                        0
             OBSCENITY
                                                       17
             OFFENSE INVOLVING CHILDREN
                                                      267
             OTHER NARCOTIC VIOLATION
                                                    False
             OTHER OFFENSE
                                                     2371
             PROSTITUTION
                                                     1344
             PUBLIC INDECENCY
             PUBLIC PEACE VIOLATION
                                                     1342
             ROBBERY
                                                      733
```

2. Displaying total number of arrets made in each police district

```
M df.groupby('District')['Arrest'].sum()
Out[63]:
         District
          1.0
                   2124
          2.0
                   1788
          3.0
                   2620
          4.0
                   3018
          5.0
                   2748
          6.0
                   3646
          7.0
                   3968
          8.0
                   3785
          9.0
                   2735
          10.0
                   2634
          11.0
                   6187
         12.0
                   1967
          13.0
                  False
          14.0
                   1492
          15.0
                   3703
          16.0
                   1192
          17.0
                   1213
          18.0
                   1839
          19.0
                   1686
          20.0
                    790
          22.0
                   1476
          24.0
                   1429
          25.0
                   3617
          31.0
          Name: Arrest, dtype: object
```

3. Displaying total number of domestic related incident in each police district

```
■ df.groupby('District')['Domestic'].sum()

In [ ]:
   Out[64]:
            District
             1.0
                      420
             2.0
                     1461
             3.0
                     2150
                     2239
             5.0
                     1686
             6.0
                     2312
             7.0
                     2289
             8.0
                     2038
             9.0
                     1434
             10.0
                     1575
                     2221
             11.0
                      963
             12.0
             13.0
                     True
             14.0
                      733
             15.0
                     1780
             16.0
                      803
             17.0
                      720
             18.0
                       311
             19.0
                      538
             20.0
                       327
             22.0
                     1038
             24.0
                      784
             25.0
                     1842
             31.0
             Name: Domestic, dtype: object
```

4. Displaying total number of arrest made in each year

5. Displaying total no of arrest made at each type of location

```
df.groupby('Location Description')['Arrest'].sum()
Location Description
ABANDONED BUILDING
                                                     298
AIRCRAFT
AIRPORT BUILDING NON-TERMINAL - NON-SECURE AREA
                                                      14
AIRPORT BUILDING NON-TERMINAL - SECURE AREA
AIRPORT EXTERIOR - NON-SECURE AREA
                                                      18
VEHICLE - OTHER RIDE SERVICE
VEHICLE NON-COMMERCIAL
                                                    1273
                                                      32
VEHICLE-COMMERCIAL
                                                      34
WAREHOUSE
YARD
Name: Arrest, Length: 107, dtype: object
```

6. Total no of arrest for each FBI crime code

```
df.groupby('FBI Code')['Arrest'].sum()
FBI Code
01A
           16
01B
            2
94A
          951
04B
          920
08A
         2024
         7293
08B
          117
10
11
         1107
12
13
           35
14
         1463
15
         1917
16
         1344
17
          165
18
        21043
          218
19
          118
20
          169
22
          325
24
         2041
26
         7412
          733
          617
         5024
          567
           33
Name: Arrest, dtype: object
```

Crosstab

1. Displaying the frequency of arrest and non-arrest for each year

```
pd.crosstab(df['Year'], df['Arrest'])

Arrest False True
Year

2012.0 86735 36083

2013.0 183 34

2014.0 487 46

2015.0 31824 14791

2016.0 17601 4704
```

Crosstab

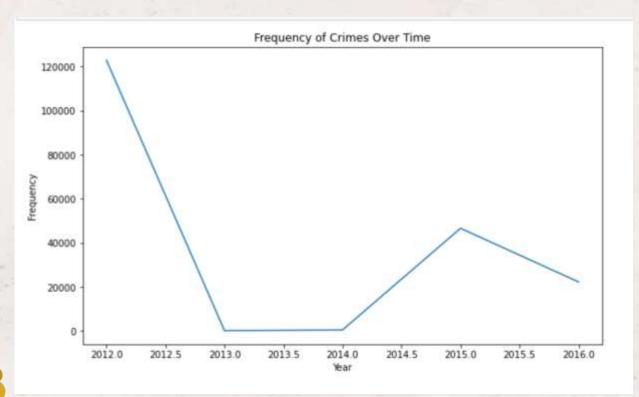
2. Displaying the count of Arrest and non-arrest for each type Location and then select the 10 rows with the largest counts of arrest.

```
pd.crosstab(df['Location Description'], df['Arrest']).nlargest(10, columns=True)
              Location Description
                        STREET 32225 13138
                      SIDEWALK
                                       11298
                    APARTMENT 20273
                     RESIDENCE
                                       4107
                         ALLEY
                                       2164
       SCHOOL, PUBLIC, BUILDING
                                        1602
             DEPARTMENT STORE
                                 1033
                                       1590
PARKING LOT/GARAGE(NON.RESID.)
                                        1523
       VEHICLE NON-COMMERCIAL
                                       1273
           GROCERY FOOD STORE
                                       1270
```



Visualization

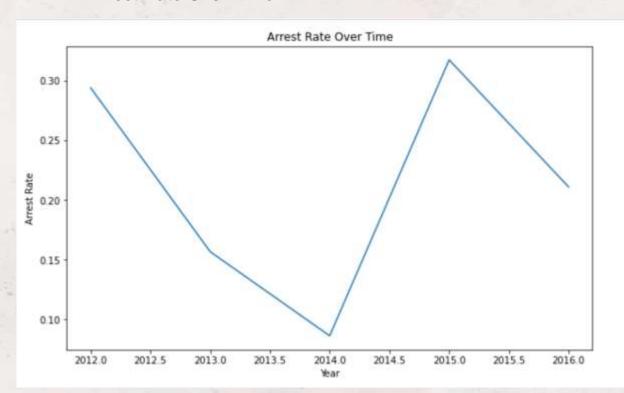
1. Frequency of Crimes Over Times:



The graph shows that crime was highest in 2012, but it dropped significantly in 2013 and 2014, reaching its lowest point during that time.

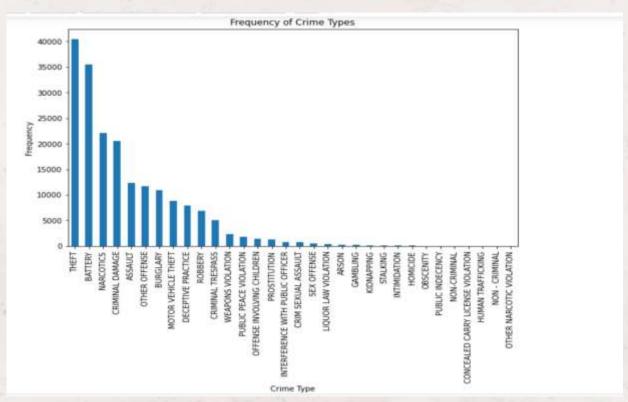
Visualization

2. Arrest Rate Over Time



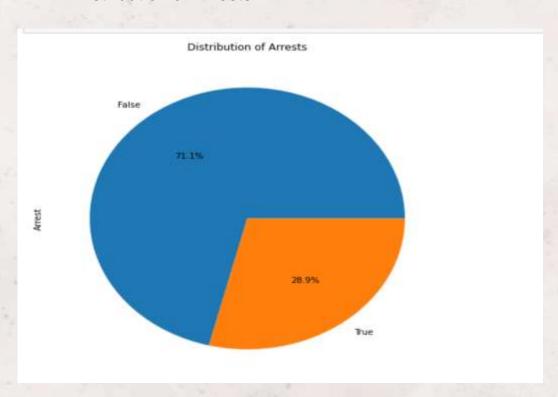
From this graph, we can see that the number of arrests was highest in the year 2015. On the other hand, the year 2014 had the fewest arrests compared to the other years shown in the graph.

3. Frequency of Crime types



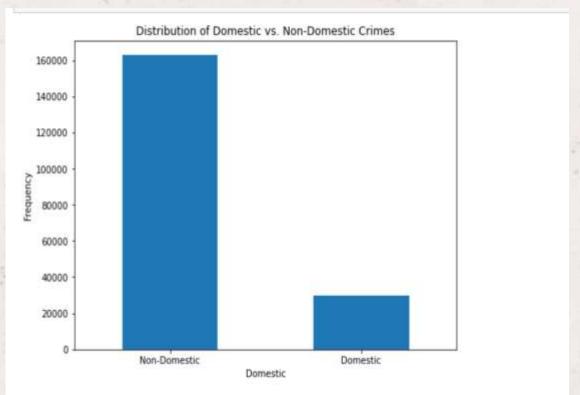
This graph tells us that theft was the most common type of crime that occurred.

4. Distribution of Arrests



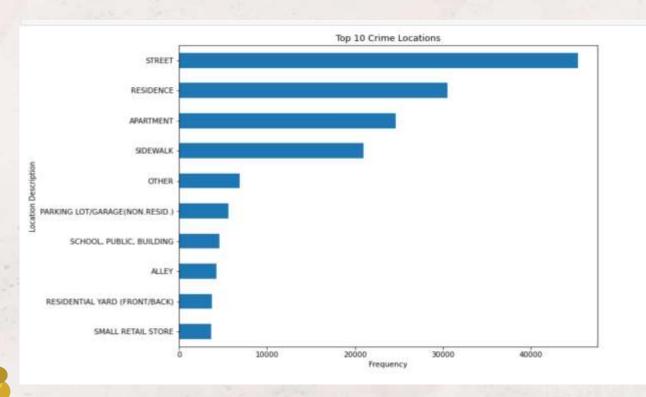
Looking at the pie chart, we can see that around 28.9% of the cases resulted in arrest, while the remaining 71.1% did not lead to arrest.

5. Distribution of Domestic vs Non-Domestic Crimes



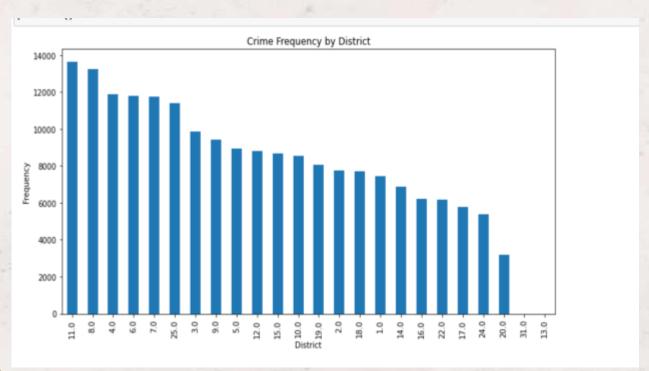
This graph clearly shows that there are more non-domestic crimes compared to domestic ones.

6. Top 10 Crime Locations



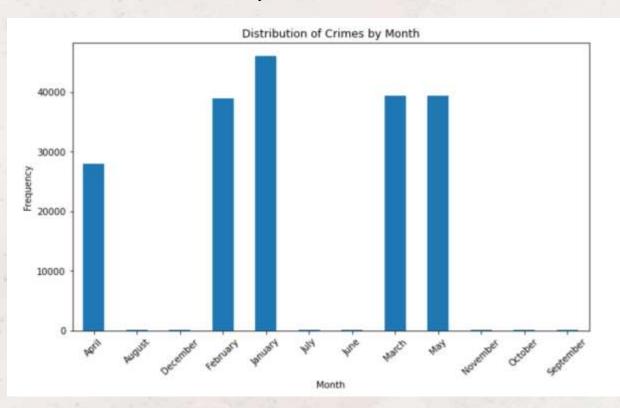
From this graph, it's evident that the most common location for incidents is on the street.

7. Crime Frequency by District



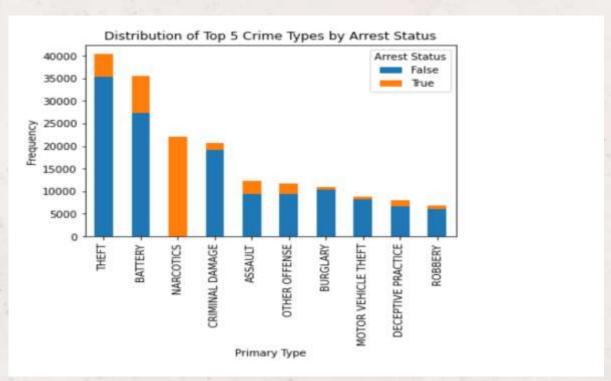
This graph indicates that District 11 experienced the highest occurrence of crimes compared to other districts.

8. Distribution of Crimes by Month



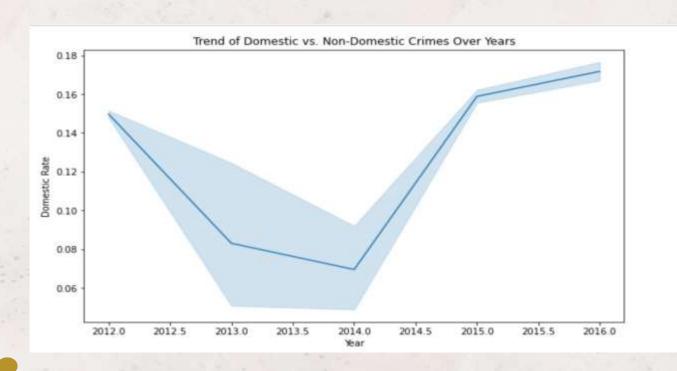
This graph indicates that January had the highest number of reported crimes compared to other months.

9. Distribution of Top 5 Crime Types by Arrest Status



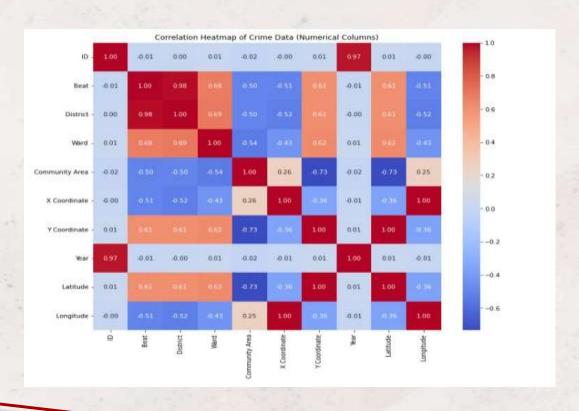
From this graph, we can analyse that for the crime type "Narcotics," all individuals involved were arrested.

10. Trend of Domestic vs. Non-Domestic Crimes over



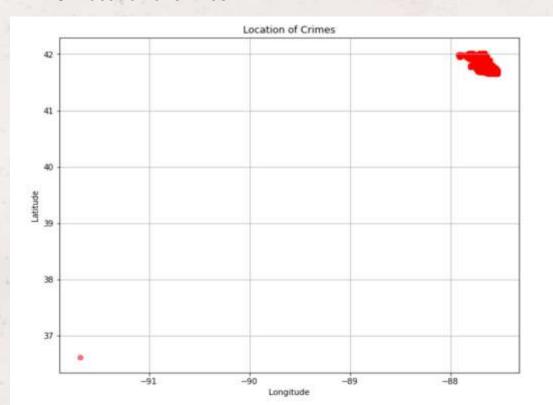
From this graph, we can analyze that the incidence of domestic crime peaked in the year 2016.

12. Correlation Heatmap of crime Data Numerical columns



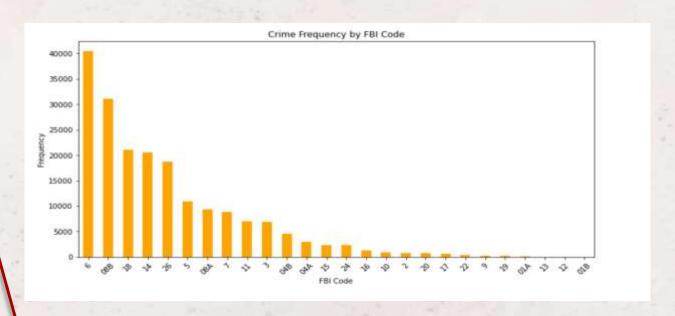
When the correlation coefficient is close to 1, it indicates a strong positive correlation the two attributes, between meaning that as one attribute increases, the other tends to increase as well. Converselv. when the correlation coefficient is close to -1, it signifies a strong negative correlation, indicating that as one attribute increases, the other tends to decrease. Α correlation coefficient near 0 suggests little linear relationship no between the attributes.

13. Location of crimes



The scatter plot displays the locations of crimes based on latitude and longitude. Each point represents a single crime incident. The dark shade color in the top-right corner indicates a higher concentration of crime incidents in that particular area.

14. Crime Frequency by FBI code



From this graph, it indicates that there are more occurrences of crimes classified under the FBI code 6 compared to other FBI codes.

Thank You