```python
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

dataset = pd.read_csv('crime_data_Proj1.csv')
columns = ['Unnamed: 0']
dataset.drop(columns, inplace = True, axis=1)
```

| | ID | Case Number | Date | Block | IUCR | Primary Type | Description | Location Description | Arrest | Domestic | ... | Ward | Co |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6407111 | HP485721 | 07/26/2008 02:30:00 PM | 085XX S MUSKEGON AVE | 1320 | CRIMINAL DAMAGE | TO VEHICLE | STREET | False | False | ... | 10.0 | |
| 1 | 11398199 | JB372830 | 07/31/2018 10:57:00 AM | 092XX S ELLIS AVE | 143C | WEAPONS VIOLATION | UNLAWFUL POSS AMMUNITION | POOL ROOM | True | False | ... | 8.0 | |
| 2 | 5488785 | HN308568 | 04/27/2007 10:30:00 AM | 062XX N TRIPP AVE | 0610 | BURGLARY | FORCIBLE ENTRY | RESIDENCE | True | False | ... | 39.0 | |
| 3 | 11389116 | JB361368 | 07/23/2018 08:55:00 AM | 0000X N KEELER AVE | 0560 | ASSAULT | SIMPLE | NURSING HOME/RETIREMENT HOME | False | False | ... | 28.0 | |
| 4 | 12420431 | JE297624 | 07/11/2021 06:40:00 AM | 016XX W HARRISON ST | 051A | ASSAULT | AGGRAVATED - HANDGUN | PARKING LOT / GARAGE (NON RESIDENTIAL) | False | False | ... | 27.0 | |
| 5 | 1699235 | G498287 | 08/21/2001 12:00:00 AM | 003XX W 28 PL | 0810 | THEFT | OVER $500 | STREET | False | False | ... | NaN | |
| 6 | 5061155 | HM660983 | 10/14/2006 10:00:00 PM | 006XX S CENTRAL AVE | 0320 | ROBBERY | STRONGARM - NO WEAPON | CTA PLATFORM | False | False | ... | 29.0 | |
| 7 | 9876456 | HX527438 | 12/02/2014 11:48:00 AM | 043XX W POTOMAC AVE | 1811 | NARCOTICS | POSS: CANNABIS 30GMS OR LESS | ALLEY | True | False | ... | 37.0 | |

8 rows × 22 columns

```
In [3]:    1  dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2278726 entries, 0 to 2278725
Data columns (total 22 columns):
 #   Column                Dtype
---  ------                -----
 0   ID                    int64
 1   Case Number           object
 2   Date                  object
 3   Block                 object
 4   IUCR                  object
 5   Primary Type          object
 6   Description           object
 7   Location Description  object
 8   Arrest                bool
 9   Domestic              bool
 10  Beat                  int64
 11  District              float64
 12  Ward                  float64
 13  Community Area        float64
 14  FBI Code              object
 15  X Coordinate          float64
 16  Y Coordinate          float64
 17  Year                  int64
 18  Updated On            object
 19  Latitude              float64
 20  Longitude             float64
 21  Location              object
dtypes: bool(2), float64(7), int64(3), object(10)
memory usage: 352.1+ MB
```

```
In [4]:   1  # idntifiying missing values
          2
          3  missing_val = dataset.isna().sum()
          4  missing_val
```

Out[4]:   ID                         0
          Case Number                1
          Date                       0
          Block                      0
          IUCR                       0
          Primary Type               0
          Description                0
          Location Description    2877
          Arrest                     0
          Domestic                   0
          Beat                       0
          District                  12
          Ward                  184695
          Community Area        184267
          FBI Code                   0
          X Coordinate           23985
          Y Coordinate           23985
          Year                       0
          Updated On                 0
          Latitude               23985
          Longitude              23985
          Location               23985
          dtype: int64

```
In [5]:  1  # check for percentage of the missing values if percentage less than < 5% then i will drop it.
         2  from pandas import read_csv
         3
         4  data = read_csv('crime_data_Proj1.csv', header= None, na_values= False )
         5
         6  for i in range(data.shape[1]):
         7      missing = data[[i]].isnull().sum()
         8      per = missing / data.shape[0] * 100
         9      print('> %d,  Missing: %d  (%.1f%%)' % (i, missing, per))
```

```
C:\Users\Shehu\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3444: DtypeWarning: Columns (1,11,
12,13,14,16,17,18,20,21) have mixed types.Specify dtype option on import or set low_memory=False.
  exec(code_obj, self.user_global_ns, self.user_ns)

> 0,  Missing: 2  (0.0%)
> 1,  Missing: 0  (0.0%)
> 2,  Missing: 1  (0.0%)
> 3,  Missing: 0  (0.0%)
> 4,  Missing: 0  (0.0%)
> 5,  Missing: 0  (0.0%)
> 6,  Missing: 0  (0.0%)
> 7,  Missing: 0  (0.0%)
> 8,  Missing: 2877  (0.1%)
> 9,  Missing: 1675252  (73.5%)
> 10,  Missing: 1965311  (86.2%)
> 11,  Missing: 0  (0.0%)
> 12,  Missing: 12  (0.0%)
> 13,  Missing: 184695  (8.1%)
> 14,  Missing: 184282  (8.1%)
> 15,  Missing: 0  (0.0%)
> 16,  Missing: 24013  (1.1%)
> 17,  Missing: 24013  (1.1%)
> 18,  Missing: 0  (0.0%)
> 19,  Missing: 0  (0.0%)
> 20,  Missing: 23985  (1.1%)
> 21,  Missing: 23985  (1.1%)
> 22,  Missing: 23985  (1.1%)
```

```
In [2]:   1  # treating ward and community Area features by replacing nan with 0
          2
          3  from numpy import nan
          4
          5  dataset['Ward'] = dataset['Ward'].replace(nan, 0)
          6  dataset['Community Area'] = dataset['Community Area'].replace(nan, 0)
```

```
In [3]:   1  # treating the missing value by droping nan
          2
          3  dataset.dropna(inplace = True)
```

```
In [5]:   1  dataset.isna().sum()
```

```
Out[5]:  ID                       0
         Case Number              0
         Date                     0
         Block                    0
         IUCR                     0
         Primary Type             0
         Description              0
         Location Description     0
         Arrest                   0
         Domestic                 0
         Beat                     0
         District                 0
         Ward                     0
         Community Area           0
         FBI Code                 0
         X Coordinate             0
         Y Coordinate             0
         Year                     0
         Updated On               0
         Latitude                 0
         Longitude                0
         Location                 0
         dtype: int64
```

```
In [9]:   1  print(round(2252860 / 2278726 * 100,2), "percentage of the data has been retained.")
```

```
98.86 percentage of the data has been retained.
```

```
In [10]:   1  dataset.describe(exclude='number').T.sort_values(by='unique')
```

Out[10]:

|  | count | unique | top | freq |
|---|---|---|---|---|
| **Arrest** | 2252860 | 2 | False | 1654990 |
| **Domestic** | 2252860 | 2 | False | 1941700 |
| **FBI Code** | 2252860 | 26 | 06 | 473222 |
| **Primary Type** | 2252860 | 35 | THEFT | 473222 |
| **Location Description** | 2252860 | 198 | STREET | 586596 |
| **IUCR** | 2252860 | 391 | 0820 | 183064 |
| **Description** | 2252860 | 513 | SIMPLE | 267014 |
| **Updated On** | 2252860 | 3620 | 02/10/2018 03:50:01 PM | 835879 |
| **Block** | 2252860 | 53378 | 100XX W OHARE ST | 4785 |
| **Location** | 2252860 | 545165 | (41.976290414, -87.905227221) | 4227 |
| **Date** | 2252860 | 1334502 | 01/01/2007 12:01:00 AM | 53 |
| **Case Number** | 2252860 | 2252813 | HK172551 | 3 |

```
In [11]:   1  # checking for duplicate in the dataset
           2
           3  # duplicate = dataset.duplicated()
           4  # duplicate.any()
           5
           6  # dataset = dataset.drop_duplicates()
           7
           8  duplicate = dataset.duplicated(keep = False).sum()
           9  duplicate
```

Out[11]: 0

```
In [6]:   1  dataset.shape
```

Out[6]: (2252860, 22)

```
In [15]:   1  # Number of distinct crimes in the city
           2
           3  crimes = dataset['Primary Type'].unique()
           4  print("The Number of distinct crimes are:", len(crimes))
           5  print()
           6  print("The Distinct Crimes are : \n", crimes)
```

The Number of distinct crimes are: 35

The Distinct Crimes are :
 ['CRIMINAL DAMAGE' 'WEAPONS VIOLATION' 'BURGLARY' 'ASSAULT' 'THEFT'
 'ROBBERY' 'NARCOTICS' 'MOTOR VEHICLE THEFT' 'BATTERY' 'OTHER OFFENSE'
 'PROSTITUTION' 'DECEPTIVE PRACTICE' 'INTIMIDATION'
 'INTERFERENCE WITH PUBLIC OFFICER' 'CRIMINAL TRESPASS' 'STALKING'
 'OFFENSE INVOLVING CHILDREN' 'PUBLIC PEACE VIOLATION' 'SEX OFFENSE'
 'CRIM SEXUAL ASSAULT' 'HOMICIDE' 'LIQUOR LAW VIOLATION'
 'CRIMINAL SEXUAL ASSAULT' 'KIDNAPPING' 'ARSON' 'GAMBLING'
 'CONCEALED CARRY LICENSE VIOLATION' 'PUBLIC INDECENCY' 'RITUALISM'
 'OBSCENITY' 'NON - CRIMINAL' 'OTHER NARCOTIC VIOLATION'
 'HUMAN TRAFFICKING' 'NON-CRIMINAL' 'NON-CRIMINAL (SUBJECT SPECIFIED)']

```
In [21]:   1  district_crimes = dataset['District'].unique()
           2  print(district_crimes)
           3
           4  # Filter out the Top 5 criminal districts
           5  top_5_district = dataset['District'].value_counts().sort_values(ascending=False).head()
           6  top_5_district
```

[ 4. 17. 11. 12.  2. 15. 25. 14.  8.  6. 19.  1.  7. 20. 18.  3.  9. 22.
  5. 16. 10. 24. 31.]

Out[21]:  8.0     152306
          11.0    144802
          6.0     131664
          7.0     131602
          4.0     128782
          Name: District, dtype: int64
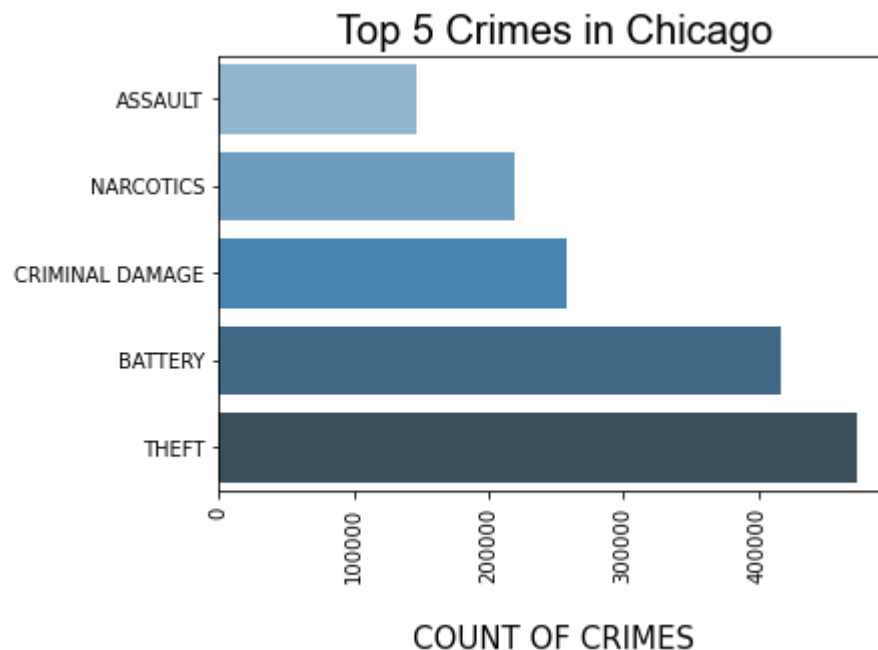
```
In [20]:   1  top_5_crimes = dataset['Primary Type'].value_counts().sort_values(ascending=False).head()
           2  top_5_crimes
```

```
Out[20]:  THEFT             473222
          BATTERY           416908
          CRIMINAL DAMAGE   258649
          NARCOTICS         219441
          ASSAULT           147046
          Name: Primary Type, dtype: int64
```

```
In [18]:   1  import seaborn as sns
           2  import matplotlib.pyplot as plt
           3
           4  crime = dataset.groupby('Primary Type', as_index = False).agg({"ID": "count"})
           5  crime_type = crime.sort_values(by = ['ID'], ascending = False).head()
           6  crime_type = crime_type.sort_values(by='ID', ascending = True)
           7  sns.barplot(x = 'ID', y = 'Primary Type', data = crime_type, palette = "Blues_d")
           8
           9  plt.title("Top 5 Crimes in Chicago", fontdict = {'fontsize': 20, 'fontname':'Arial', 'color': '#000000'})
          10  plt.xlabel("\nCOUNT OF CRIMES", fontdict = {'fontsize': 15})
          11  plt.ylabel("")
          12  plt.xticks(rotation=90)
          13  plt.show()
```



```
In [ ]:   1
```

```
In [ ]:   1  # feature engineering extracting the features month, period and day from the feature Date
```

```
In [4]:  1  from dateutil.parser import parse
         2  from datetime import datetime
         3  tCol = dataset.Date
         4
         5  List = [(datetime.ctime(parse(x[0:-3])), x[-2:]) for x in tCol]
         6  dayList = []
         7  monthList = []
         8  periodList = []
         9
        10  for row in List:
        11      day = row[0][0:4]
        12      month = row[0][4:7]
        13      if row[1] == 'AM':
        14          period = 'Morning'
        15      elif row[1] == 'PM' and int(row[0][11:13]) < 4:
        16          period = 'Afternoon'
        17      elif row[1] == 'PM' and int(row[0][11:13]) < 5:
        18          period = 'Evening'
        19      elif row[1] == 'PM' and int(row[0][11:13]) > 5:
        20
        21          period = 'Night'
        22      else:
        23          period = 'Unknown'
        24
        25      dayList.append(day)
        26      monthList.append(month)
        27      periodList.append(period)
        28
        29  print(len(dayList), len(monthList), len(periodList))
```

2252860 2252860 2252860

```
In [5]:  1  # create new features month, period and day into the dataset (extracting from feature Date)
         2
         3  dataset['month'] = monthList
         4  dataset['period'] = periodList
         5  dataset['day'] = dayList
```
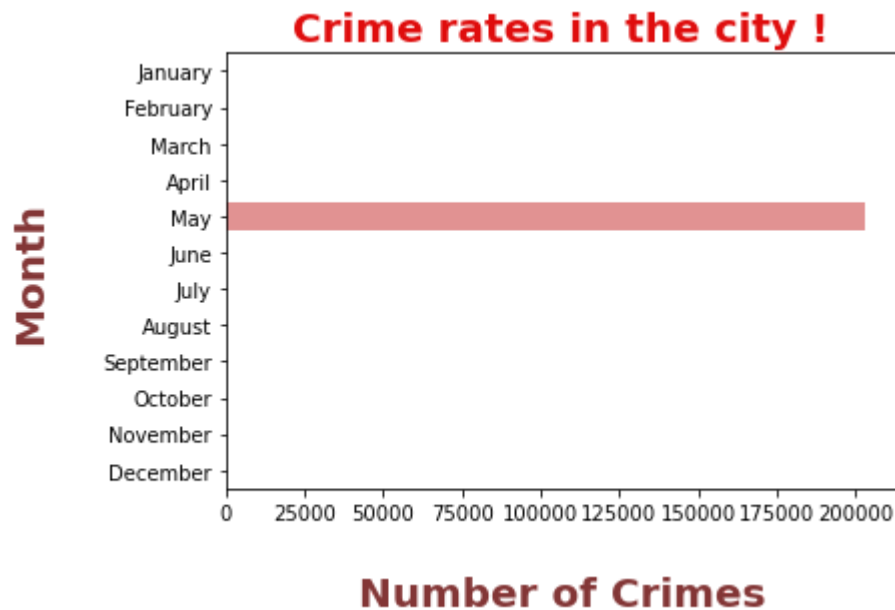
```
In [6]:    1  # dropping the following features to avoid duplicate
           2
           3  del dataset['Date']
           4  del dataset['Updated On']
```

```
In [7]:    1  dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2252860 entries, 0 to 2278725
Data columns (total 23 columns):
 #   Column                Dtype
---  ------                -----
 0   ID                    int64
 1   Case Number           object
 2   Block                 object
 3   IUCR                  object
 4   Primary Type          object
 5   Description           object
 6   Location Description  object
 7   Arrest                bool
 8   Domestic              bool
 9   Beat                  int64
 10  District              float64
 11  Ward                  float64
 12  Community Area        float64
 13  FBI Code              object
 14  X Coordinate          float64
 15  Y Coordinate          float64
 16  Year                  int64
 17  Latitude              float64
 18  Longitude             float64
 19  Location              object
 20  month                 object
 21  period                object
 22  day                   object
dtypes: bool(2), float64(7), int64(3), object(11)
memory usage: 382.4+ MB
```

```
1  sns.countplot(y = 'month', data = dataset, palette = ["#DF0D0D"], order = ['January', 'February', 'March',
2
3  plt.title("Crime rates in the city !", fontdict={'fontsize': 20, 'color': '#DF0D0D', 'fontname':'Agency FB'
4  plt.ylabel("Month\n", fontdict = {'fontsize': 20}, weight = "bold", color = "#833636")
5  plt.xlabel("\nNumber of Crimes", fontdict =  {'fontsize': 20}, weight = "bold", color = "#833636")
6
7  plt.xticks(fontsize = 10, color = 'black')
8  plt.yticks(fontsize = 10, color = 'black')
9  plt.show()
```

findfont: Font family ['Agency FB'] not found. Falling back to DejaVu Sans.



In [ ]:  `# The month of May have seen the most high crime rates in the city.`

```
In [26]:   1  district_crime = dataset.groupby(['District','month'], as_index =  False).agg({'Primary Type':"count"})
           2  district_crime.columns
           3
           4  district_crime = district_crime.pivot("District", "month", "Primary Type")
           5
           6
           7  plt.figure(figsize = (20,15))
           8  plt.title("District vs Month")
           9  with sns.axes_style("white"):
          10      sns.heatmap(district_crime, mask = district_crime.isnull(), cmap = "inferno", annot = True, fmt = "f")
```

District vs Month

| District | Apr | Aug | Dec | Feb | Jan | Jul | Jun | Mar | May | Nov | Oct | Sep |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1.0 | 7065.000000 | 8200.000000 | 6915.000000 | 6479.000000 | 7267.000000 | 8582.000000 | 7976.000000 | 7224.000000 | 7773.000000 | 7072.000000 | 7645.000000 | 7451.000000 |
| 2.0 | 8773.000000 | 9689.000000 | 7481.000000 | 7260.000000 | 8500.000000 | 10199.000000 | 9800.000000 | 8904.000000 | 9856.000000 | 8017.000000 | 9286.000000 | 9053.000000 |
| 3.0 | 9704.000000 | 10491.000000 | 8365.000000 | 7670.000000 | 8939.000000 | 10890.000000 | 10453.000000 | 9432.000000 | 10420.000000 | 8633.000000 | 9581.000000 | 9569.000000 |
| 4.0 | 10877.000000 | 11562.000000 | 9430.000000 | 8867.000000 | 10034.000000 | 12251.000000 | 11638.000000 | 10636.000000 | 11856.000000 | 9971.000000 | 11001.000000 | 10659.000000 |
| 5.0 | 8536.000000 | 8778.000000 | 7219.000000 | 6735.000000 | 7464.000000 | 9550.000000 | 9116.000000 | 8288.000000 | 9242.000000 | 7686.000000 | 8574.000000 | 8446.000000 |
| 6.0 | 10870.000000 | 11631.000000 | 9753.000000 | 9207.000000 | 10380.000000 | 12240.000000 | 11877.000000 | 11270.000000 | 12165.000000 | 10201.000000 | 11088.000000 | 10982.000000 |
| 7.0 | 10898.000000 | 12157.000000 | 9299.000000 | 8597.000000 | 10139.000000 | 12723.000000 | 12203.000000 | 11054.000000 | 12088.000000 | 9877.000000 | 11193.000000 | 11374.000000 |
| 8.0 | 12390.000000 | 13457.000000 | 11488.000000 | 10782.000000 | 12182.000000 | 13992.000000 | 13479.000000 | 12643.000000 | 13902.000000 | 12198.000000 | 12982.000000 | 12811.000000 |
| 9.0 | 9332.000000 | 9924.000000 | 7871.000000 | 7707.000000 | 8698.000000 | 10393.000000 | 10054.000000 | 9080.000000 | 9946.000000 | 8565.000000 | 9474.000000 | 9341.000000 |
| 10.0 | 8179.000000 | 8889.000000 | 6992.000000 | 6795.000000 | 7642.000000 | 9037.000000 | 8739.000000 | 8143.000000 | 8832.000000 | 7397.000000 | 8246.000000 | 8352.000000 |
| 11.0 | 12178.000000 | 13184.000000 | 10389.000000 | 10459.000000 | 11548.000000 | 13130.000000 | 12597.000000 | 12377.000000 | 13134.000000 | 11122.000000 | 12292.000000 | 12392.000000 |
| 12.0 | 8948.000000 | 9985.000000 | 8111.000000 | 7689.000000 | 8662.000000 | 10241.000000 | 9834.000000 | 9121.000000 | 9608.000000 | 8828.000000 | 9741.000000 | 9685.000000 |
| 14.0 | 7074.000000 | 8141.000000 | 6488.000000 | 6116.000000 | 6870.000000 | 8436.000000 | 7788.000000 | 7012.000000 | 7677.000000 | 6804.000000 | 7698.000000 | 7668.000000 |
| 15.0 | 8216.000000 | 8633.000000 | 7236.000000 | 6899.000000 | 7597.000000 | 9008.000000 | 8745.000000 | 8376.000000 | 9034.000000 | 7488.000000 | 8369.000000 | 8154.000000 |
| 16.0 | 6130.000000 | 6755.000000 | 5873.000000 | 5249.000000 | 6086.000000 | 7000.000000 | 6883.000000 | 6336.000000 | 6491.000000 | 5990.000000 | 6423.000000 | 6247.000000 |
| 17.0 | 5274.000000 | 5791.000000 | 4843.000000 | 4847.000000 | 5300.000000 | 5884.000000 | 5656.000000 | 5385.000000 | 5673.000000 | 5004.000000 | 5664.000000 | 5633.000000 |
| 18.0 | 7869.000000 | 9354.000000 | 7710.000000 | 6942.000000 | 7859.000000 | 9624.000000 | 9080.000000 | 7936.000000 | 8759.000000 | 7852.000000 | 8608.000000 | 8595.000000 |
| 19.0 | 7944.000000 | 9535.000000 | 7176.000000 | 6721.000000 | 7663.000000 | 9736.000000 | 9513.000000 | 8029.000000 | 8767.000000 | 7573.000000 | 8672.000000 | 8742.000000 |
| 20.0 | 3202.000000 | 3543.000000 | 2911.000000 | 2883.000000 | 3168.000000 | 3587.000000 | 3580.000000 | 3147.000000 | 3440.000000 | 3062.000000 | 3436.000000 | 3403.000000 |
| 22.0 | 5972.000000 | 6579.000000 | 5342.000000 | 4959.000000 | 5595.000000 | 6856.000000 | 6728.000000 | 6209.000000 | 6872.000000 | 5756.000000 | 6508.000000 | 6387.000000 |
| 24.0 | 5568.000000 | 6282.000000 | 5031.000000 | 4713.000000 | 5357.000000 | 6402.000000 | 6008.000000 | 5477.000000 | 6074.000000 | 5117.000000 | 5774.000000 | 5684.000000 |
| 25.0 | 10524.000000 | 11612.000000 | 9648.000000 | 9028.000000 | 10345.000000 | 11864.000000 | 11232.000000 | 10688.000000 | 11422.000000 | 10158.000000 | 11103.000000 | 10873.000000 |
| 31.0 | 3.000000 | 6.000000 | 7.000000 | 6.000000 | 9.000000 | 10.000000 | 6.000000 | 2.000000 | 2.000000 | 3.000000 | 5.000000 | 4.000000 |

```
In [ ]:  1  # district 8 and 11 are having the high crime rate.
```

```
In [285]:    1  # location attributes = ['Location Description','Beat','District', 'Ward', 'Community Area','X Coordinate',
             2  %matplotlib inline
             3  top_crime = dataset.groupby(['District', 'Primary Type']).size().reset_index(name='counts').groupby('Distri
             4  #print(topk)
             5
             6  # factor plot to make multiple plots
             7  graph = sns.catplot("Primary Type", y = 'counts', col = "District", col_wrap = 3,
             8                      data = top_crime, kind = 'bar')
             9  for ax in graph.axes:
            10      plt.setp(ax.get_xticklabels(), visible = True, rotation = 30, ha='right')
            11
            12  plt.subplots_adjust(hspace = 0.4)
```

C:\Users\Shehu\anaconda3\lib\site-packages\seaborn\_decorators.py:36: FutureWarning: Pass the following vari
able as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing
other arguments without an explicit keyword will result in an error or misinterpretation.
  warnings.warn(

```
In [28]:  1  crime_period = dataset[['Primary Type','period']]
          2  crime_period = crime_period.groupby(['Primary Type'])
          3  crime_period.head()
```

Out[28]:

|  | Primary Type | period |
| --- | --- | --- |
| 0 | CRIMINAL DAMAGE | Afternoon |
| 1 | WEAPONS VIOLATION | Morning |
| 2 | BURGLARY | Morning |
| 3 | ASSAULT | Morning |
| 4 | ASSAULT | Morning |
| ... | ... | ... |
| 668399 | RITUALISM | Evening |
| 867154 | NON - CRIMINAL | Morning |
| 880836 | RITUALISM | Night |
| 1031143 | NON-CRIMINAL (SUBJECT SPECIFIED) | Night |
| 1587959 | NON-CRIMINAL (SUBJECT SPECIFIED) | Unknown |

172 rows × 2 columns

```
In [29]:  1  fig, ax = plt.subplots(figsize=(10, 5))
          2  sns.countplot(x = 'period', data = dataset, palette = "viridis")
          3
          4  # Aesthetic appeal
          5  plt.title("Unsafest periods in Chicago", fontdict={'fontsize': 25, 'color': '#bb0e14','fontname':'Agency FB
          6  plt.xlabel("\nHour in the Day", fontdict={'fontsize': 10}, weight='bold')
          7  plt.ylabel("Number of Crimes\n", fontdict={'fontsize': 10}, weight="bold")
          8
          9  plt.text(3, 150000, 'Lowest Crime Rate', fontdict= {'fontsize': 12, 'color':"blue" }, weight='bold')
         10  plt.show()
```

findfont: Font family ['Agency FB'] not found. Falling back to DejaVu Sans.

```
In [30]:  1  day_crime = dataset.groupby(['Primary Type', 'day']).agg({'day':"count"})
          2  day_crime.columns = ["Count"]
          3  day_crime
```

Out[30]:

| Primary Type | day | Count |
|---|---|---|
| ARSON | Fri | 500 |
| | Mon | 566 |
| | Sat | 598 |
| | Sun | 652 |
| | Thu | 490 |
| ... | ... | ... |
| WEAPONS VIOLATION | Sat | 4558 |
| | Sun | 4225 |
| | Thu | 4239 |
| | Tue | 4284 |
| | Wed | 4109 |

235 rows × 1 columns

```
In [ ]:  1
```

```
In [ ]:  1  # Building Models
```

```
In [8]:    1  # treating the boolean features (Arresrt, Domestic) by replaced True and False by 1 and 0
           2
           3  # dataset = dataset.replace({False: 0, True: 1})
           4  Arrest = dataset['Arrest'].replace((True, False), (1, 0))
           5  Domestic = dataset['Domestic'].replace((True, False), (1, 0))
           6
           7  dataset['Arrest'] = Arrest
           8  dataset['Domestic'] = Domestic
```

```
In [9]:    1  import numpy as np
           2
           3  for categorical_column in ['Community Area', 'Ward', 'District']:
           4      dataset[categorical_column] = dataset[categorical_column].astype(np.int64)
           5
           6  for categorical_column in ['Primary Type', 'Description', 'Location Description', 'month']:
           7      dataset[categorical_column] = dataset[categorical_column].astype('category')
```

```
In [10]:   1  from sklearn.preprocessing import LabelEncoder
           2  le = LabelEncoder()
           3
           4  dataset['Primary Type'] = le.fit_transform(dataset['Primary Type']).astype('int64')
           5  dataset['Location Description'] = le.fit_transform(dataset['Location Description']).astype('int64')
           6  dataset['Description'] = le.fit_transform(dataset['Description']).astype('int64')
           7  dataset['month'] = le.fit_transform(dataset['month']).astype('int64')
           8  dataset['District'] = le.fit_transform(dataset['District']).astype('int64')
           9  dataset['Location'] = le.fit_transform(dataset['Location']).astype('int64')
          10  dataset['period'] = le.fit_transform(dataset['period']).astype('int64')
          11  dataset['day'] = le.fit_transform(dataset['day']).astype('int64')
```

```
In [11]:    1  dataset.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 2252860 entries, 0 to 2278725
Data columns (total 23 columns):
 #   Column                Dtype
---  ------                -----
 0   ID                    int64
 1   Case Number           object
 2   Block                 object
 3   IUCR                  object
 4   Primary Type          int64
 5   Description           int64
 6   Location Description  int64
 7   Arrest                int64
 8   Domestic              int64
 9   Beat                  int64
 10  District              int64
 11  Ward                  int64
 12  Community Area        int64
 13  FBI Code              object
 14  X Coordinate          float64
 15  Y Coordinate          float64
 16  Year                  int64
 17  Latitude              float64
 18  Longitude             float64
 19  Location              int64
 20  month                 int64
 21  period                int64
 22  day                   int64
dtypes: float64(4), int64(15), object(4)
memory usage: 412.5+ MB
```

```
In [ ]:     1
```

```
In [133]:   1  columns = ['Primary Type', 'Description', 'Location Description','Arrest','Domestic','Beat','District','War
            2            'Community Area','X Coordinate','Y Coordinate','Year','Latitude','Longitude','Location','month',
            3
            4  data = dataset[columns]
```

```
In [134]:   1  data.head()
```

Out[134]:

| | Primary Type | Description | Location Description | Arrest | Domestic | Beat | District | Ward | Community Area | X Coordinate | Y Coordinate | Year | Latitude | Long |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 6 | 462 | 171 | 0 | 0 | 423 | 3 | 10 | 46 | 1196638.0 | 1848800.0 | 2008 | 41.739980 | -87.55 |
| 1 | 34 | 468 | 139 | 1 | 0 | 413 | 3 | 8 | 47 | 1184499.0 | 1843935.0 | 2018 | 41.726922 | -87.59 |
| 2 | 3 | 216 | 145 | 1 | 0 | 1711 | 15 | 39 | 12 | 1146911.0 | 1941022.0 | 2007 | 41.994138 | -87.73 |
| 3 | 1 | 421 | 125 | 0 | 0 | 1115 | 10 | 28 | 26 | 1148388.0 | 1899882.0 | 2018 | 41.881217 | -87.73 |
| 4 | 1 | 25 | 134 | 0 | 0 | 1231 | 11 | 27 | 28 | 1165430.0 | 1897441.0 | 2021 | 41.874174 | -87.66 |

```
In [135]:   1  data.shape
```

Out[135]:  (2252860, 18)

```
In [138]:  1  # Checking for outliers
           2  plt.figure(figsize = (20, 15))
           3  for i in range (len(data.columns)):
           4      plt.subplot(4, 5, i+1)
           5      sns.boxplot(x = data.iloc[:, i])
           6      plt.xlabel(data.columns[i], size = 18)
```

```
In [139]:   1  # Removing outliers
            2  def Outliers(data, feature):
            3      IQ1 = data[feature].quantile(0.25)
            4      IQ3 = data[feature].quantile(0.75)
            5      IQR = IQ3 - IQ1
            6
            7      lower_bound = IQ1 - 1.5 * IQR
            8      upper_bound = IQ3 + 1.5 * IQR
            9
           10      index = data.index[ (data[feature] < lower_bound) | (data[feature] > upper_bound) ]
           11      return index
```

```
In [140]:   1  Outliers(data, 'period')
```

```
Out[140]: Int64Index([      0,      18,      24,      26,      28,      32,      44,
                       51,      78,      79,
                   ...
                  2278665, 2278667, 2278669, 2278687, 2278688, 2278690, 2278694,
                  2278699, 2278700, 2278721],
                 dtype='int64', length=341608)
```

```
In [144]:   1  # Getting index of all the outliers
            2
            3  index = []
            4
            5  for i in data.columns:
            6      index.extend(Outliers(data, i))
            7  index = set(index)
            8  print("Total number of outliers are {}".format(len(index)))
            9
           10  # Dropping all the outliers
           11  data.drop(index, inplace = True, axis = 0)
```

```
Total number of outliers are 0
```
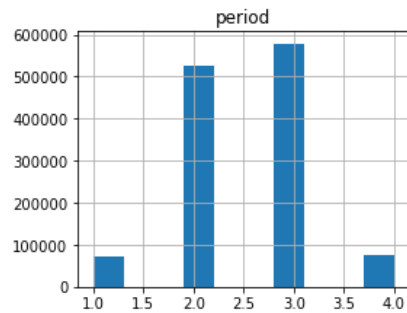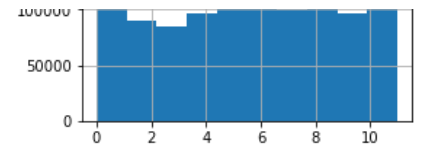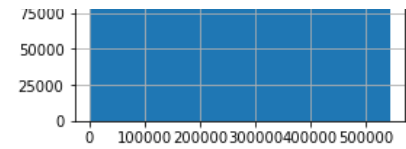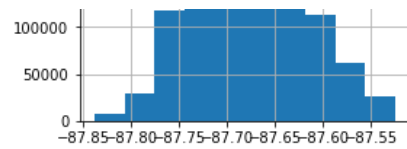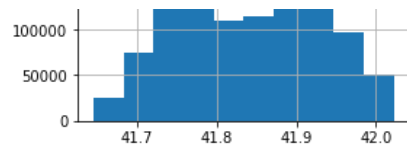
```
1  # checking for outliers in the dataseet
2
3  import matplotlib.pyplot as plt
4  import seaborn as sns
5
6  plt.figure(figsize = (20, 15))
7  for i in range(len(data.columns)):
8      plt.subplot(4, 5, i + 1)
9      sns.boxplot(x = data.iloc[:, i])
10     plt.xlabel(data.columns[i], size = 18)
```

```
1  # univariate
2
3  import matplotlib.pyplot as plt
4
5  data.hist()
6  plt.gcf().set_size_inches(20,20)
7  plt.show()
```

```
In [147]:    1  data.describe().T
```

Out[147]:

| | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| **Primary Type** | 1252041.0 | 1.675805e+01 | 12.223026 | 0.000000e+00 | 6.000000e+00 | 1.700000e+01 | 3.000000e+01 | 3.400000e+01 |
| **Description** | 1252041.0 | 2.888157e+02 | 157.418869 | 0.000000e+00 | 1.480000e+02 | 3.340000e+02 | 4.220000e+02 | 5.120000e+02 |
| **Location Description** | 1252041.0 | 1.589880e+02 | 15.787014 | 1.060000e+02 | 1.450000e+02 | 1.660000e+02 | 1.710000e+02 | 1.970000e+02 |
| **Arrest** | 1252041.0 | 2.684744e-01 | 0.443166 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 1.000000e+00 | 1.000000e+00 |
| **Domestic** | 1252041.0 | 0.000000e+00 | 0.000000 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 |
| **Beat** | 1252041.0 | 1.208138e+03 | 696.372231 | 1.110000e+02 | 6.310000e+02 | 1.113000e+03 | 1.732000e+03 | 2.535000e+03 |
| **District** | 1252041.0 | 9.933527e+00 | 6.090848 | 0.000000e+00 | 5.000000e+00 | 9.000000e+00 | 1.500000e+01 | 2.200000e+01 |
| **Ward** | 1252041.0 | 2.100356e+01 | 14.575182 | 0.000000e+00 | 8.000000e+00 | 2.100000e+01 | 3.200000e+01 | 5.000000e+01 |
| **Community Area** | 1252041.0 | 3.395031e+01 | 23.044400 | 0.000000e+00 | 1.700000e+01 | 2.900000e+01 | 5.300000e+01 | 7.700000e+01 |
| **X Coordinate** | 1252041.0 | 1.164286e+06 | 15664.644659 | 1.118919e+06 | 1.152698e+06 | 1.165365e+06 | 1.175652e+06 | 1.205119e+06 |
| **Y Coordinate** | 1252041.0 | 1.886036e+06 | 31498.042822 | 1.813894e+06 | 1.859274e+06 | 1.891464e+06 | 1.909992e+06 | 1.951622e+06 |
| **Year** | 1252041.0 | 2.009299e+03 | 5.884196 | 2.001000e+03 | 2.004000e+03 | 2.008000e+03 | 2.014000e+03 | 2.022000e+03 |
| **Latitude** | 1252041.0 | 4.184289e+01 | 0.086617 | 4.164459e+01 | 4.176921e+01 | 4.185780e+01 | 4.190884e+01 | 4.202291e+01 |
| **Longitude** | 1252041.0 | -8.767265e+01 | 0.056993 | -8.783807e+01 | -8.771471e+01 | -8.766852e+01 | -8.763123e+01 | -8.752453e+01 |
| **Location** | 1252041.0 | 2.667257e+05 | 147704.129911 | 1.000000e+00 | 1.405400e+05 | 2.764410e+05 | 3.797280e+05 | 5.451640e+05 |
| **month** | 1252041.0 | 5.586907e+00 | 3.457462 | 0.000000e+00 | 3.000000e+00 | 6.000000e+00 | 9.000000e+00 | 1.100000e+01 |
| **period** | 1252041.0 | 2.523199e+00 | 0.696131 | 1.000000e+00 | 2.000000e+00 | 3.000000e+00 | 3.000000e+00 | 4.000000e+00 |
| **day** | 1252041.0 | 2.981380e+00 | 2.018929 | 0.000000e+00 | 1.000000e+00 | 3.000000e+00 | 5.000000e+00 | 6.000000e+00 |

```
In [148]:    1  data.shape
```

Out[148]:  (1252041, 18)

```
In [201]:   1  # day_type as a target variable feature
            2  columns = ['Primary Type', 'Description', 'Location Description','Arrest','Beat','Ward',
            3            'Community Area','X Coordinate','Y Coordinate','Year','Location','month','day']
            4  day_data = data[columns]
            5
            6  # period as a target variable feature
            7  columns = ['Primary Type', 'Description', 'Location Description','Arrest','Beat','Ward',
            8            'Community Area','X Coordinate','Y Coordinate','Year','Location','month','period']
            9  period_data = data[columns]
           10
           11  # District as a target variable feature
           12  columns = ['Primary Type', 'Description', 'Location Description','Arrest','Beat','Ward',
           13            'Community Area','X Coordinate','Y Coordinate','Year','Location','month','District']
           14  district_data = data[columns]
```

```
In [151]:   1  day_data.columns
```

```
Out[151]: Index(['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
            'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year',
            'Location', 'month', 'day'],
          dtype='object')
```

```
In [152]:   1  day_data.head()
```

Out[152]:

| | Primary Type | Description | Location Description | Arrest | Beat | Ward | Community Area | X Coordinate | Y Coordinate | Year | Location | month | day |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 34 | 468 | 139 | 1 | 413 | 8 | 47 | 1184499.0 | 1843935.0 | 2018 | 59718 | 5 | 5 |
| 2 | 3 | 216 | 145 | 1 | 1711 | 39 | 12 | 1146911.0 | 1941022.0 | 2007 | 525966 | 0 | 0 |
| 3 | 1 | 421 | 125 | 0 | 1115 | 28 | 26 | 1148388.0 | 1899882.0 | 2018 | 316579 | 5 | 1 |
| 4 | 1 | 25 | 134 | 0 | 1231 | 27 | 28 | 1165430.0 | 1897441.0 | 2021 | 299524 | 5 | 3 |
| 5 | 33 | 334 | 171 | 0 | 2113 | 0 | 0 | 1174343.0 | 1885951.0 | 2001 | 256235 | 1 | 5 |

```python
In [283]: 1  # data transformation using scikit learn MinMaxScaler to Rescale
          2
          3  from sklearn.preprocessing import MinMaxScaler
          4
          5  array = day_data.values
          6  x = array[:,0:12]
          7  y = array[:,12]
          8
          9  model = MinMaxScaler(feature_range= (0, 1))
         10  rescaledx = model.fit_transform(x)
         11
         12  rescaledx_data = pd.DataFrame(rescaledx)
         13  rescaledx_data.columns = ['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
         14                            'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year', 'Location', 'mo
```

```python
In [284]: 1  rescaledx_data['day'] = day_data['day']
          2  rescaledx_data['day'] = rescaledx_data['day'].fillna(0)
```

```python
In [263]: 1  rescaledx_data.isna().sum()
```

```
Out[263]: Primary Type              0
          Description               0
          Location Description      0
          Arrest                    0
          Beat                      0
          Ward                      0
          Community Area            0
          X Coordinate              0
          Y Coordinate              0
          Year                      0
          Location                  0
          month                     0
          day                       0
          dtype: int64
```

```python
In [ ]: 1
```

```
In [264]:    1  rescaledx_data.shape
```

Out[264]:  (1252041, 13)

```
In [63]:    1  ndf = rescaledx_data.iloc[0:80000, 0:13]
```

```
In [357]:    1  ndf.shape
```

Out[357]:  (80000, 13)

```
In [157]:    1  rescaledx_data.describe().T
```

Out[157]:

|  | count | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|
| **Primary Type** | 1252041.0 | 0.492884 | 0.359501 | 0.0 | 0.176471 | 0.500000 | 0.882353 | 1.0 |
| **Description** | 1252041.0 | 0.564093 | 0.307459 | 0.0 | 0.289062 | 0.652344 | 0.824219 | 1.0 |
| **Location Description** | 1252041.0 | 0.582286 | 0.173484 | 0.0 | 0.428571 | 0.659341 | 0.714286 | 1.0 |
| **Arrest** | 1252041.0 | 0.268474 | 0.443166 | 0.0 | 0.000000 | 0.000000 | 1.000000 | 1.0 |
| **Beat** | 1252041.0 | 0.452615 | 0.287282 | 0.0 | 0.214521 | 0.413366 | 0.668729 | 1.0 |
| **Ward** | 1252041.0 | 0.420071 | 0.291504 | 0.0 | 0.160000 | 0.420000 | 0.640000 | 1.0 |
| **Community Area** | 1252041.0 | 0.440913 | 0.299278 | 0.0 | 0.220779 | 0.376623 | 0.688312 | 1.0 |
| **X Coordinate** | 1252041.0 | 0.526297 | 0.181724 | 0.0 | 0.391868 | 0.538817 | 0.658155 | 1.0 |
| **Y Coordinate** | 1252041.0 | 0.523799 | 0.228697 | 0.0 | 0.329490 | 0.563212 | 0.697738 | 1.0 |
| **Year** | 1252041.0 | 0.395188 | 0.280200 | 0.0 | 0.142857 | 0.333333 | 0.619048 | 1.0 |
| **Location** | 1252041.0 | 0.489257 | 0.270936 | 0.0 | 0.257793 | 0.507078 | 0.696538 | 1.0 |
| **month** | 1252041.0 | 0.507901 | 0.314315 | 0.0 | 0.272727 | 0.545455 | 0.818182 | 1.0 |
| **day** | 1252041.0 | 1.637712 | 2.107257 | 0.0 | 0.000000 | 0.000000 | 3.000000 | 6.0 |

```
In [358]:   1  # from pandas import read_csv
            2  # from sklearn.feature_selection import RFE
            3  # from sklearn.linear_model import LogisticRegression
            4
            5  # array = ndf.values
            6  # X = array[:, 0:12]
            7  # Y = array[:, 12]
            8
            9  # # feature extraction
           10  # model = LogisticRegression(solver= 'lbfgs', max_iter= 50000)
           11  # rfe = RFE(model, 8)
           12  # fit = rfe.fit(X, Y)
           13
           14  # # print("Num Features: {}".format(fit.n_features_))
           15  # # print("Selected Features: {}".format(fit.support_))
           16  # print("Feature Ranking: {}".format(fit.ranking_))
```

C:\Users\Shehu\anaconda3\lib\site-packages\sklearn\utils\validation.py:70: FutureWarning: Pass n_features_to_s
elect=8 as keyword args. From version 1.0 (renaming of 0.25) passing these as positional arguments will result
in an error
  warnings.warn(f"Pass {args_msg} as keyword args. From version "

Feature Ranking: [1 1 1 1 4 2 5 1 1 1 1 3]

```python
# Univeriate selection

from numpy import set_printoptions
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2

array = rescaledx_data.values
x = array[:, 0:12]
y = array[:, 12]

test = SelectKBest(score_func=chi2, k = 8)
fit = test.fit(x, y)

set_printoptions(precision = 3)
print(fit.scores_)
features = fit.transform(x)
print(features[0:2, :])
```

```
[2.038 1.567 0.622 5.64  0.238 1.084 0.676 0.029 0.267 1.4   0.436 2.22 ]
[[1.    0.914 0.363 1.    0.16  0.61  0.81  0.455]
 [0.088 0.422 0.429 1.    0.78  0.156 0.286 0.   ]]
```

```
In [275]:    1  dfeature = pd.DataFrame(features, columns=['a','b','c','d','e','f','g','h'])
             2  col_dict = dict()
             3  cols = list(rescaledx_data.columns)[0:-1]
             4
             5  for col in cols:
             6      dcol = rescaledx_data[col].values.tolist()
             7      col_dict[col] = dcol
             8
             9  count = 0
            10  col_len = len(dfeature.columns.tolist())
            11
            12  while count < col_len:
            13      col_label = ''
            14      if count == 0:
            15          col_label = 'a'
            16      elif count == 1:
            17          col_label = 'b'
            18      elif count == 2:
            19          col_label = 'c'
            20      elif count == 3:
            21          col_label = 'd'
            22      elif count == 4:
            23          col_label = 'e'
            24      elif count == 5:
            25          col_label = 'f'
            26      elif count == 6:
            27          col_label = 'g'
            28      elif count == 7:
            29          col_label = 'h'
            30      else:
            31          pass
            32
            33      coldata = dfeature.loc[:, col_label]
            34      coldatav = list(coldata.values)
            35
            36      for item in col_dict.items():
            37          if coldatav == item[1]:
            38              print(item[0])
            39
            40      count += 1
```

Primary Type
Description

```
Location Description
Arrest
Ward
Community Area
Year
month
```
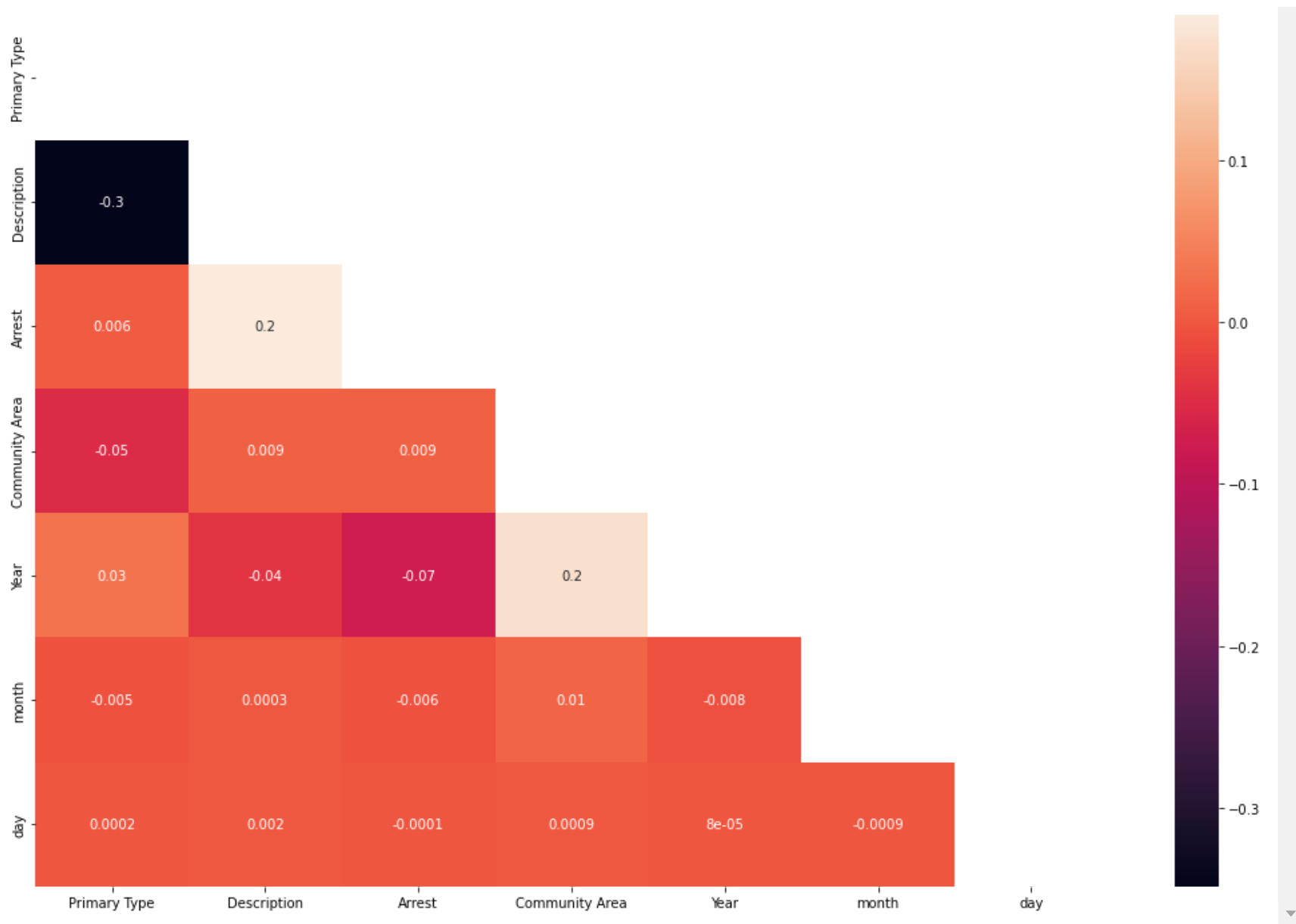
In [279]:
```python
dataframe = rescaledx_data.loc[:, ['Primary Type', 'Description','Arrest',
                                   'Community Area','Year', 'month', 'day']]
```

```
In [280]:   1  import numpy as np
            2
            3  matrix = np.triu(dataframe.corr())
            4  plt.figure(figsize = (18, 12))
            5  sns.heatmap(dataframe.corr(), annot = True, mask = matrix, fmt='.1g')
```

Out[280]:  <AxesSubplot:>

```
In [ ]:   1  import seaborn as sns
          2  import matplotlib.pyplot as plt
          3
          4  fig, ax = plt.subplots(figsize= (12,7))
          5  sns.countplot(dataframe['day'])
```

```
In [ ]:     1

In [271]:   1  # Evaluate using Cross Validation
            2  from sklearn.model_selection import KFold
            3  from sklearn.model_selection import cross_val_score
            4  from sklearn.linear_model import LogisticRegression
            5
            6  array = dataframe.values
            7  x = array[:, 0:6]
            8  y = array[:, 6]
            9
           10  num_folds = 10
           11  seed = 7
           12  kfold = KFold(n_splits=num_folds)
           13  model = LogisticRegression()
           14  results = cross_val_score(model, X, Y, cv=kfold)
           15
           16  print("Accuracy: {}.3f%%, {}.3f%%".format(results.mean()*100.0, results.std()*100.0))
```

Accuracy: 45.103954295633095.3f%%, 0.08728157932558331.3f%%

```
In [282]:   1  # Evaluate using a train and a test set
            2  from pandas import read_csv
            3  from sklearn.model_selection import train_test_split
            4  from sklearn.linear_model import LogisticRegression
            5
            6  array = dataframe.values
            7  x = array[:, 0:6]
            8  y = array[:, 6]
            9
           10  test_size = 0.33
           11  seed = 10
           12  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_size, random_state=seed)
           13
           14  model = LogisticRegression()
           15  model.fit(x_train, y_train)
           16  result = model.score(x_test, y_test)
           17
           18  print("Accuracy: {}".format(result * 100.0))
```

Accuracy: 52.46796749069399

```
In [ ]:     1

In [ ]:     1  # models

In [164]:   1  # CART Classification
            2  from pandas import read_csv
            3  from sklearn.model_selection import KFold
            4  from sklearn.model_selection import cross_val_score
            5  from sklearn.tree import DecisionTreeClassifier
            6
            7  array = dataframe.values
            8  x = array[:, 0:8]
            9  y = array[:, 8]
           10
           11  kfold = KFold(n_splits=10)
           12  model = DecisionTreeClassifier()
           13  results = cross_val_score(model, x, y, cv=kfold)
           14
           15  print(results.mean())
```

0.30925345087007694

```
In [198]:   1  # Linear Discriminant Aanaalysis
            2
            3  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
            4
            5  array = dataframe.values
            6  x = array[:, 0:8]
            7  y = array[:, 8]
            8
            9  num_folds = 10
           10  kfold = KFold(n_splits= num_folds, shuffle= True, random_state= 7)
           11  model = LinearDiscriminantAnalysis()
           12  results = cross_val_score(model, x, y, cv= kfold)
           13
           14  print(results.mean())
```

0.5339010489874335

```
In [73]:   1  # Gaussian Naive Bayes Classification
           2  from sklearn.model_selection import KFold
           3  from sklearn.model_selection import cross_val_score
           4  from sklearn.naive_bayes import GaussianNB
           5
           6  array = dataframe.values
           7  X = array[:, 0:7]
           8  Y = array[:, 7]
           9
          10  kfold = KFold(n_splits=10)
          11  model = GaussianNB()
          12  results = cross_val_score(model, x, y, cv = kfold)
          13
          14  print(results.mean())
```

0.5338858519667504

```
In [74]:  1  # LogissticRegression Classification
          2
          3  import numpy as np
          4  from sklearn.model_selection import KFold
          5  from sklearn.model_selection import train_test_split
          6  from sklearn.model_selection import cross_val_score
          7  from sklearn.model_selection import cross_val_predict
          8  from sklearn.linear_model import LogisticRegression
          9  from sklearn.metrics import confusion_matrix
         10  from sklearn.metrics import classification_report
         11
         12  array = dataframe.values
         13  x = array[:, 0:8]
         14  y = array[:, 8]
         15
         16  num_folds = 10
         17  test_size = 0.33
         18  seed = 7
         19
         20  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size= test_size, random_state= seed)
         21  kfold = KFold(n_splits= num_folds)
         22  model = LogisticRegression(max_iter= 500)
         23  results = cross_val_score(model, x, y, cv= kfold)
         24  score = np.mean(results)
         25  y_pred = cross_val_predict(model, x, y, cv= 3)
         26  con_m = confusion_matrix(y, y_pred)
         27  my_model = model.fit(x_train, y_train)
         28  predicted = my_model.predict(x_test)
         29  report = classification_report(y_test, predicted)
         30
         31  print(results)
         32  print()
         33  print('The mean score is', score * 100)
         34  print()
         35  print(con_m)
         36  print()
         37  print(report)
```

```
C:\Users\Shehu\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning:
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_di
vision` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

```
C:\Users\Shehu\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning:
Precision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_di
vision` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))

[0.534 0.535 0.533 0.533 0.533 0.534 0.535 0.534 0.532 0.535]

The mean score is 53.388585196675045

[[668471        0        0        0        0        0        0]
 [ 96561        0        0        0        0        0        0]
 [ 98846        0        0        0        0        0        0]
 [ 92614        0        0        0        0        0        0]
 [ 98470        0        0        0        0        0        0]
 [ 98095        0        0        0        0        0        0]
 [ 99029        0        0        0        0        0        0]]
              precision    recall  f1-score   support

         0.0       0.53      1.00      0.70    220414
         1.0       0.00      0.00      0.00     32067
         2.0       0.00      0.00      0.00     32450
         3.0       0.00      0.00      0.00     30577
         4.0       0.00      0.00      0.00     32648
         5.0       0.00      0.00      0.00     32215
         6.0       0.00      0.00      0.00     32818

    accuracy                           0.53    413189
   macro avg       0.08      0.14      0.10    413189
weighted avg       0.28      0.53      0.37    413189


C:\Users\Shehu\anaconda3\lib\site-packages\sklearn\metrics\_classification.py:1248: UndefinedMetricWarning: Pr
ecision and F-score are ill-defined and being set to 0.0 in labels with no predicted samples. Use `zero_divisi
on` parameter to control this behavior.
  _warn_prf(average, modifier, msg_start, len(result))
```

```
In [78]:   1  # cm = con_m
           2  # labls = ['N', 'Y']
           3  # sns.set(font_scale= 2)
           4  # cmn = cm.astype('float') / cm.sum(axis= 1)[:, np.newaxis]
           5  # fig, ax = plt.subplots(figsize= (20, 12))
           6  # sns.heatmap(cmn, annot= True, fmt='.2f', xticklabels= labls, yticklabels=labls)
           7  # plt.ylabel('Actual COF ', size= 20)
           8  # plt.xlabel('Predicted COF', size= 20)
           9  # plt.title('Confusion Matrix of Prediction', size= 24)
          10
          11  # # plt.figure(figsize= (16, 8))
          12  # plt.show(block=False)
```

```
In [ ]:    1
```

```
In [ ]:    1
```

```
In [ ]:    1  # period as a target variable
```

```
In [165]:  1  period_data.head()
```

Out[165]:

| | Primary Type | Description | Location Description | Arrest | Beat | Ward | Community Area | X Coordinate | Y Coordinate | Year | Location | month | period |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 34 | 468 | 139 | 1 | 413 | 8 | 47 | 1184499.0 | 1843935.0 | 2018 | 59718 | 5 | 2 |
| **2** | 3 | 216 | 145 | 1 | 1711 | 39 | 12 | 1146911.0 | 1941022.0 | 2007 | 525966 | 0 | 2 |
| **3** | 1 | 421 | 125 | 0 | 1115 | 28 | 26 | 1148388.0 | 1899882.0 | 2018 | 316579 | 5 | 2 |
| **4** | 1 | 25 | 134 | 0 | 1231 | 27 | 28 | 1165430.0 | 1897441.0 | 2021 | 299524 | 5 | 2 |
| **5** | 33 | 334 | 171 | 0 | 2113 | 0 | 0 | 1174343.0 | 1885951.0 | 2001 | 256235 | 1 | 2 |

```
In [166]:  1  period_data.columns
```

Out[166]:  Index(['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
               'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year',
               'Location', 'month', 'period'],
              dtype='object')

```
In [224]:   1  # data transformation using scikit learn MinMaxScaler to Rescale
            2
            3  from sklearn.preprocessing import MinMaxScaler
            4
            5  array = period_data.values
            6  x = array[:,0:12]
            7  y = array[:,12]
            8
            9  model = MinMaxScaler(feature_range= (0, 1))
           10  rescaledx = model.fit_transform(x)
           11
           12  rescaledx_data = pd.DataFrame(rescaledx)
           13  rescaledx_data.columns = ['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
           14                            'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year','Location', 'mon
```

```
In [225]:   1  rescaledx_data['period'] = period_data['period']
            2  rescaledx_data['period'] = rescaledx_data['period'].fillna(rescaledx_data['period'].median()).astype('int64
```

```
In [226]:   1  rescaledx_data.isna().sum()
```

```
Out[226]:  Primary Type            0
           Description             0
           Location Description    0
           Arrest                  0
           Beat                    0
           Ward                    0
           Community Area          0
           X Coordinate            0
           Y Coordinate            0
           Year                    0
           Location                0
           month                   0
           period                  0
           dtype: int64
```

```python
# Univeriate selection

from numpy import set_printoptions
from sklearn.feature_selection import SelectKBest
from sklearn.feature_selection import chi2

array = rescaledx_data.values
x = array[:, 0:12]
y = array[:, 12]

test = SelectKBest(score_func=chi2, k = 9)
fit = test.fit(x, y)

set_printoptions(precision = 3)
print(fit.scores_)
features = fit.transform(x)
print(features[0:2, :])
```

```
[0.786 0.458 0.057 1.238 0.442 1.123 0.137 0.2   0.245 1.022 0.383 0.219]
[[1.    0.914 1.    0.125 0.16  0.218 0.81  0.11  0.455]
 [0.088 0.422 1.    0.66  0.78  0.923 0.286 0.965 0.   ]]
```

```
In [228]:  1  dfeature = pd.DataFrame(features, columns=['a','b','c','d','e','f','g','h','i'])
           2  col_dict = dict()
           3  cols = list(rescaledx_data.columns)[0:-1]
           4
           5  for col in cols:
           6      dcol = rescaledx_data[col].values.tolist()
           7      col_dict[col] = dcol
           8
           9  count = 0
          10  col_len = len(dfeature.columns.tolist())
          11
          12  while count < col_len:
          13      col_label = ''
          14      if count == 0:
          15          col_label = 'a'
          16      elif count == 1:
          17          col_label = 'b'
          18      elif count == 2:
          19          col_label = 'c'
          20      elif count == 3:
          21          col_label = 'd'
          22      elif count == 4:
          23          col_label = 'e'
          24      elif count == 5:
          25          col_label = 'f'
          26      elif count == 6:
          27          col_label = 'g'
          28      elif count == 7:
          29          col_label = 'h'
          30      elif count == 8:
          31          col_label = 'i'
          32      else:
          33          pass
          34
          35      coldata = dfeature.loc[:, col_label]
          36      coldatav = list(coldata.values)
          37
          38      for item in col_dict.items():
          39          if coldatav == item[1]:
          40              print(item[0])
          41
          42      count += 1
```

Primary Type
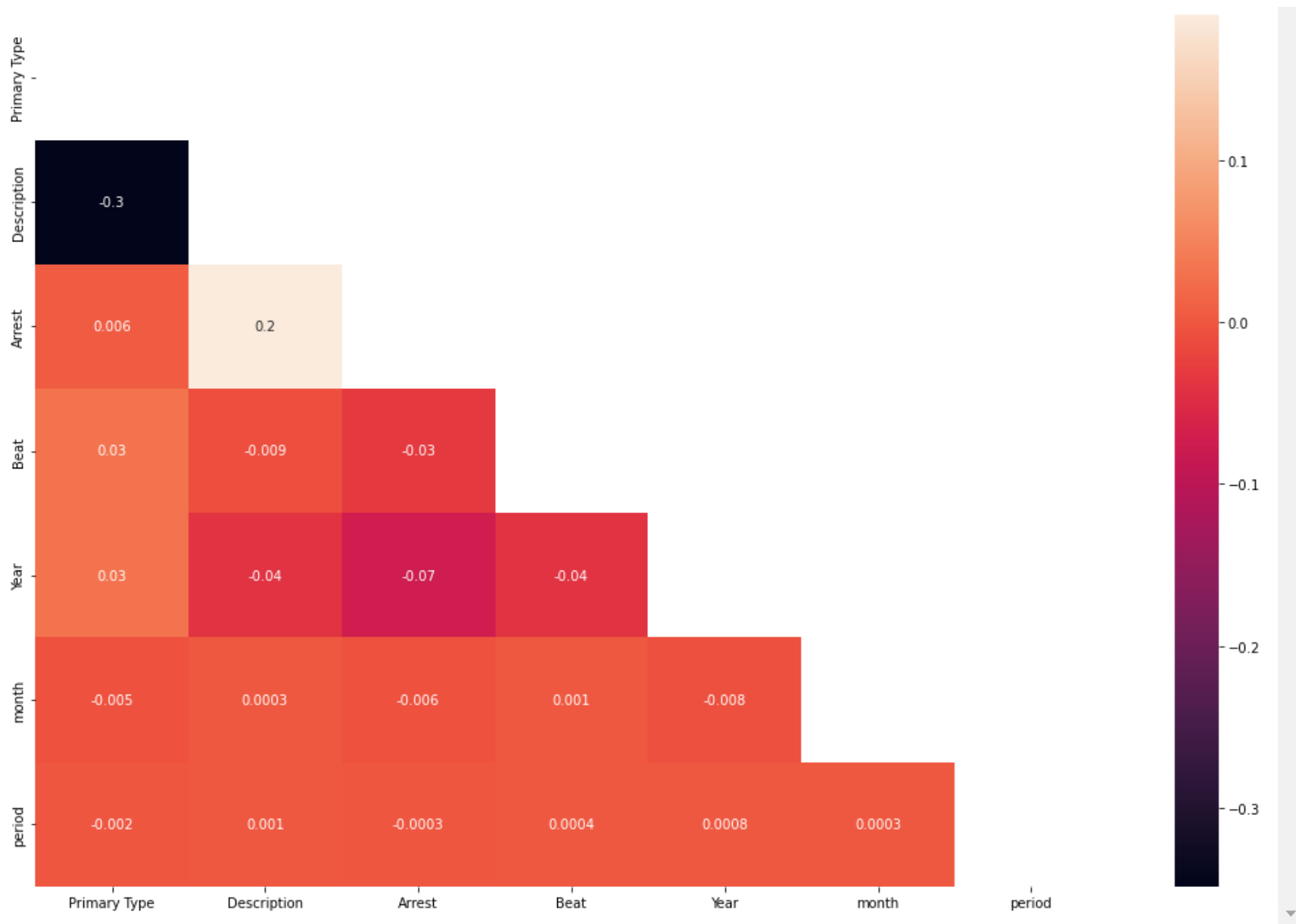Description
Arrest
Beat
Ward
Y Coordinate
Year
Location
month

In [257]:
```python
dataframe1 = rescaledx_data.loc[:, ['Primary Type','Description','Arrest','Beat',
                                     'Year', 'month', 'period']]
```

```
In [258]:   1  import numpy as np
            2
            3  matrix = np.triu(dataframe1.corr())
            4  plt.figure(figsize = (18, 12))
            5  sns.heatmap(dataframe1.corr(), annot = True, mask = matrix, fmt='.1g')
```

Out[258]:  <AxesSubplot:>

In [ ]: 1

```
In [259]:   1  # Evaluate using a train and a test set
            2  from pandas import read_csv
            3  from sklearn.model_selection import train_test_split
            4  from sklearn.linear_model import LogisticRegression
            5
            6  array = dataframe1.values
            7  x = array[:, 0:6]
            8  y = array[:, 6]
            9
           10  test_size = 0.33
           11  seed = 10
           12  x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_size, random_state=seed)
           13
           14  model = LogisticRegression()
           15  model.fit(x_train, y_train)
           16  result = model.score(x_test, y_test)
           17
           18  print("Accuracy: {}".format(result * 100.0))
```

Accuracy: 70.49257697725415

In [ ]:   1

```python
# Using Random Forest for classification
import sklearn.metrics as metrics
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report
# from sklearn.externals import joblib

array = dataframe1.values
x = array[:, 0:7]
y = array[:, 7]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 10 )

scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

classifier = RandomForestClassifier(n_estimators = 10, criterion = 'entropy', random_state = 10)
classifier.fit(X_train, y_train)

y_pred = classifier.predict(X_test)

print("Accuracy:",(metrics.accuracy_score(y_test, y_pred)*100),"\n")

cm = pd.crosstab(y_test, y_pred, rownames=['Actual Alarm'], colnames=['Predicted Alarm'])
print("\n----------Confusion Matrix----------------------------------")
print(cm)

print("\n----------Classification Report----------------------------------")
print(classification_report(y_test,y_pred))
```

Accuracy: 58.47398334243844


----------Confusion Matrix------------------------------------

| Predicted Alarm | 1.0 | 2.0 | 3.0 | 4.0 |
|---|---|---|---|---|
| Actual Alarm | | | | |
| 1.0 | 210 | 2013 | 7721 | 176 |
| 2.0 | 1494 | 14065 | 55403 | 1275 |
| 3.0 | 4606 | 43560 | 168568 | 3775 |

```
4.0                    215    1942    7801    187

----------Classification Report------------------------------------
              precision    recall   f1-score    support

        1.0       0.03      0.02       0.03      10120
        2.0       0.23      0.19       0.21      72237
        3.0       0.70      0.76       0.73     220509
        4.0       0.03      0.02       0.02      10145

    accuracy                           0.58     313011
   macro avg       0.25      0.25       0.25     313011
weighted avg       0.55      0.58       0.57     313011
```

In [245]:
```python
# Gaussian Naive Bayes Classification
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.naive_bayes import GaussianNB

array = dataframe1.values
x = array[:, 0:6]
y = array[:, 6]

kfold = KFold(n_splits=10)
model = GaussianNB()
results = cross_val_score(model, x, y, cv = kfold)

print(results.mean())
```

```
0.7047333113287193
```

In [189]:
```python
# CART Classification

from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier

array = dataframe1.values
X = array[:, 0:6]
Y = array[:, 6]

kfold = KFold(n_splits=10)
model = DecisionTreeClassifier()
results = cross_val_score(model, x, y, cv=kfold)

print(results.mean())
```

0.3196293095967449

In [190]:
```python
from sklearn.tree import DecisionTreeClassifier

array = dataframe1.values
X = array[:, 0:6]
Y = array[:, 6]

test_size = 0.33
seed = 10
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_size, random_state=seed)

model = DecisionTreeClassifier()
model.fit(x,y)

score = model.score(x_test, y_test)
pred = model.predict(x_test)

print(pred)
```

```
In [233]:  1  # Linear Discriminant Aanaalysis
           2
           3  from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
           4
           5  array = dataframe1.values
           6  x = array[:, 0:9]
           7  y = array[:, 9]
           8
           9  num_folds = 10
          10  kfold = KFold(n_splits= num_folds, shuffle= True, random_state= 7)
          11  model = LinearDiscriminantAnalysis()
          12  results = cross_val_score(model, x, y, cv= kfold)
          13
          14  print(results.mean())
```

0.7047333125407503

```
In [ ]:  1
```

```
In [ ]:  1  # district as a target variable
```

```
In [203]:  1  district_data.head()
```

Out[203]:

| | Primary Type | Description | Location Description | Arrest | Beat | Ward | Community Area | X Coordinate | Y Coordinate | Year | Location | month | District |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **1** | 34 | 468 | 139 | 1 | 413 | 8 | 47 | 1184499.0 | 1843935.0 | 2018 | 59718 | 5 | 3 |
| **2** | 3 | 216 | 145 | 1 | 1711 | 39 | 12 | 1146911.0 | 1941022.0 | 2007 | 525966 | 0 | 15 |
| **3** | 1 | 421 | 125 | 0 | 1115 | 28 | 26 | 1148388.0 | 1899882.0 | 2018 | 316579 | 5 | 10 |
| **4** | 1 | 25 | 134 | 0 | 1231 | 27 | 28 | 1165430.0 | 1897441.0 | 2021 | 299524 | 5 | 11 |
| **5** | 33 | 334 | 171 | 0 | 2113 | 0 | 0 | 1174343.0 | 1885951.0 | 2001 | 256235 | 1 | 1 |

```
In [202]:  1  district_data.columns
```

Out[202]: Index(['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
            'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year',
            'Location', 'month', 'District'],
           dtype='object')

```python
In [204]:   1   # data transformation using scikit learn MinMaxScaler to Rescale
            2
            3   from sklearn.preprocessing import MinMaxScaler
            4
            5   array = district_data.values
            6   x = array[:,0:12]
            7   y = array[:,12]
            8
            9   model = MinMaxScaler(feature_range= (0, 1))
           10   rescaledx = model.fit_transform(x)
           11
           12   rescaledx_data = pd.DataFrame(rescaledx)
           13   rescaledx_data.columns = ['Primary Type', 'Description', 'Location Description', 'Arrest', 'Beat',
           14                             'Ward', 'Community Area', 'X Coordinate', 'Y Coordinate', 'Year',
           15                             'Location', 'month']
```

```python
In [ ]:     1
```

```python
In [207]:   1   rescaledx_data['District'] = district_data['District']
            2   rescaledx_data['District'] = rescaledx_data['District'].fillna(0)
```

```python
In [ ]:     1   # Univeriate selection
            2
            3   from numpy import set_printoptions
            4   from sklearn.feature_selection import SelectKBest
            5   from sklearn.feature_selection import chi2
            6
            7   array = rescaledx_data.values
            8   x = array[:, 0:12]
            9   y = array[:, 12]
           10
           11   test = SelectKBest(score_func=chi2, k = 9)
           12   fit = test.fit(x, y)
           13
           14   set_printoptions(precision = 3)
           15   print(fit.scores_)
           16   features = fit.transform(x)
           17   print(features[0:2, :])
```

```python
dfeature = pd.DataFrame(features, columns=['a','b','c','d','e','f','g','h','i'])
col_dict = dict()
cols = list(rescaledx_data.columns)[0:-1]

for col in cols:
    dcol = rescaledx_data[col].values.tolist()
    col_dict[col] = dcol

count = 0
col_len = len(dfeature.columns.tolist())

while count < col_len:
    col_label = ''
    if count == 0:
        col_label = 'a'
    elif count == 1:
        col_label = 'b'
    elif count == 2:
        col_label = 'c'
    elif count == 3:
        col_label = 'd'
    elif count == 4:
        col_label = 'e'
    elif count == 5:
        col_label = 'f'
    elif count == 6:
        col_label = 'g'
    elif count == 7:
        col_label = 'h'
    elif count == 8:
        col_label = 'i'
    else:
        pass

    coldata = dfeature.loc[:, col_label]
    coldatav = list(coldata.values)

    for item in col_dict.items():
        if coldatav == item[1]:
            print(item[0])

    count += 1
```

```
dataframe2 = rescaledx_data.loc[:, ['Primary Type','Description','Arrest',
                                     'Year','Location', 'month', 'period']]
```

```
import numpy as np

matrix = np.triu(dataframe2.corr())
plt.figure(figsize = (18, 12))
sns.heatmap(dataframe1.corr(), annot = True, mask = matrix, fmt='.1g')
```

```
# Evaluate using a train and a test set
from pandas import read_csv
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression

array = dataframe2.values
x = array[:, 0:6]
y = array[:, 6]

test_size = 0.33
seed = 10
x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=test_size, random_state=seed)

model = LogisticRegression()
model.fit(x_train, y_train)
result = model.score(x_test, y_test)

print("Accuracy: {}".format(result * 100.0))
```

```
# CART Classification
from pandas import read_csv
from sklearn.model_selection import KFold
from sklearn.model_selection import cross_val_score
from sklearn.tree import DecisionTreeClassifier

array = dataframe1.values
X = array[:, 0:6]
Y = array[:, 6]

kfold = KFold(n_splits=10)
model = DecisionTreeClassifier()
results = cross_val_score(model, x, y, cv=kfold)

print(results.mean())
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [  ]:  1
```

```
In [17]:  1
```

```
In [  ]:  1
```