

Milestone 3: Model Implementation

For Milestones 1 and 2, we successfully completed several foundational steps that have set the stage for the progress in Milestone 3. We selected our project topic, developed a structured website with dedicated tabs to organize our findings, and clearly defined the aim and end goals of the project. These milestones also included performing comprehensive Exploratory Data Analysis (EDA) on the dataset, where we gained deeper insights into the variables, relationships, and patterns by creating and analyzing 10 detailed visualizations. Furthermore, we formulated 10 key questions that drive our analysis, providing a roadmap for answering critical aspects of the project. These efforts have equipped us with an in-depth understanding of our dataset and the project's scope, enabling us to approach Milestone 3 with a clear direction and actionable insights.

Website Link: <https://nikhil02jk.github.io/Crime-data-analysis/>

Screenshot of the data after the cleaning step for Milestone 2

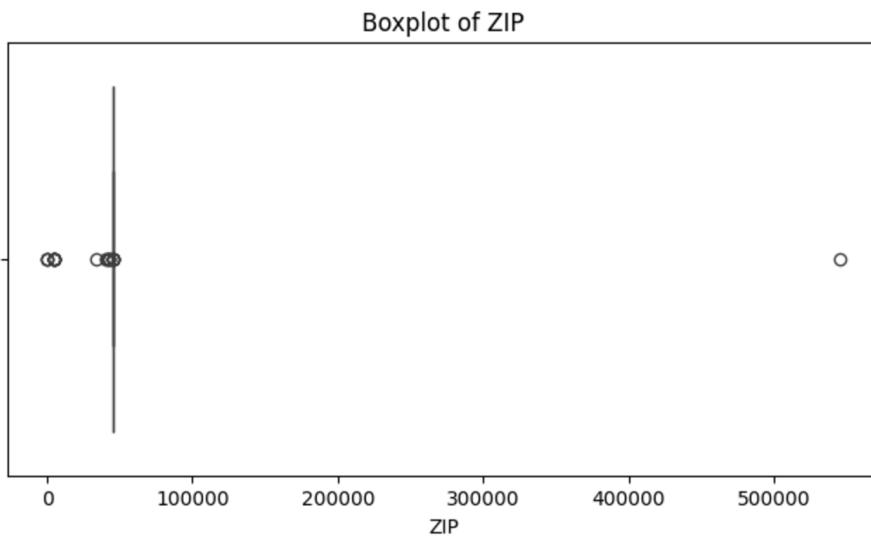
INCIDENT NO	CLOD	LOD	DFT	BERT	OFFENSE	LOCATION	THEFT_CODE	HIRE_BAS	DAYOFWEEK	CDP_NEIGHBORHOOD	WEAPONS	DATE_OF_CLEAVER	HOUR_FROM	HOUR_TO	ADDRESS	LONGITUDE	LATITUDE
0	Z20933241	F-CLEARED BY ARREST - ADULT	1400.0	3	CRIMINAL DAMAGE/DANGEROUS	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	TUESDAY	WESTWOOD	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/30/2023 12:00:00 AM	220.0	220.0	220X MONTANA AV	0.3591321199617140	0.50989	
1	Z20933250	Z-EARLY CLOSED	810.0	4	ASSAULT	02-MULTI FAMILY	N-NO BAS NOT APPLICABLE	TUESDAY	AVONDALE	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	12/15/2022 12:00:00 AM	30.0	330.0	99X GLENWOOD AV	0.657032196219850	0.54954	
2	Z20933251	D-VICTIM REFUSED TO COOPERATE	300.0	4	ASSAULTED ROBBERY	45-PARKING LOT	N-NO BAS NOT APPLICABLE	TUESDAY	CARTHAGE	12 - HANDGUN	01/15/2023 12:00:00 AM	190.0	193.0	19X SEMYOUNDR	0.87101991572220	0.83676	
3	Z20933252	D-VICTIM REFUSED TO COOPERATE	300.0	4	ASSAULTED ROBBERY	48-PARKING LOT	N-NO BAS NOT APPLICABLE	TUESDAY	CARTHAGE	12 - HANDGUN	01/15/2023 12:00:00 AM	190.0	193.0	19X SEMYOUNDR	0.870943157300	0.83660	
4	Z20933253	H-WARRANT ISSUED	600.0	5	DOMESTIC VIOLENCE	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	TUESDAY	EAST PRICE HILL	99 - NONE	10/15/2022 12:00:00 AM	184.0	185.0	100X COOLWOOD AV	0.424624584887790	0.27143	
5	Z20933254	H-WARRANT ISSUED	600.0	5	DOMESTIC VIOLENCE	29H-ALL OTHER LARCENY	N-NO BAS NOT APPLICABLE	TUESDAY	EAST PRICE HILL	99 - NONE	10/15/2022 12:00:00 AM	184.0	185.0	100X COOLWOOD AV	0.424624584887790	0.27143	
6	Z20933255	J-CLOSED	800.0	2	THEFT	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	TUESDAY	EVANSTON	99 - NONE	01/20/2023 12:00:00 AM	181.0	181.0	300X TRUMBULL AV	0.7160430277980	0.54177	
7	Z20933251	J-CLOSED	800.0	4	ASSAULTED MENACING	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	TUESDAY	AVONDALE	12 - HANDGUN	01/20/2023 12:00:00 AM	170.0	170.0	30X ROCKDALE AV	0.627074806725810	0.54991	
8	Z20934052	H-WARRANT ISSUED	701.0	1	UNAUTHORIZED USE OF MOTOR VEHICLE	02-GAMBLING FACILITY / CASINO / FACE TRACK	N-NO BAS NOT APPLICABLE	FRIDAY	C. S. O. / RIVERFRONT	99 - NONE	12/30/2022 12:00:00 AM	119.0	130.0	0			
9	Z20934053	H-WARRANT ISSUED	701.0	1	UNAUTHORIZED USE OF MOTOR VEHICLE	02-GAMBLING FACILITY / CASINO / FACE TRACK	N-NO BAS NOT APPLICABLE	FRIDAY	C. S. O. / RIVERFRONT	99 - NONE	12/30/2022 12:00:00 AM	119.0	130.0	0			
10	Z20934054	J-CLEARED BY ARREST - ADULT	270.0	1	DOMESTIC VIOLENCE	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	THURSDAY	CLINTON PARK/URBAN HEIGHTS	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	220.0	220.0	0			
11	Z20934046	J-CLEARED BY ARREST - ADULT	610.0	4	ASSAULT	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	THURSDAY	HARTWELL	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	211.0	230.0	19X GALTHER RD	0.8710497356490940	0.54901	
12	Z20934045	J-CLEARED BY ARREST - ADULT	802.0	1	ASSAULTED ROBBERY	01-MINIMAL HOME	N-NO BAS NOT APPLICABLE	THURSDAY	WEST END	U - UNKNOWN	02/16/2023 12:00:00 AM	190.0	201.0	40X ELIZABETH ST	0.5501035949240	0.30905	
13	Z20934041	H-WARRANT ISSUED	810.0	1	ASSAULT	02-MULTI FAMILY APARTMENT	N-NO BAS NOT APPLICABLE	THURSDAY	WEST END	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/24/2023 12:00:00 AM	200.0	203.0	20X DULLEY WALK	0.526294971248910	0.24743	
14	Z20934042	H-WARRANT ISSUED	810.0	1	ASSAULTED ROBBERY	01-MINIMAL HOME	N-NO BAS NOT APPLICABLE	THURSDAY	WEST END	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/24/2023 12:00:00 AM	184.0	185.0	20X DULLEY WALK	0.526294971248910	0.24743	
15	Z20934043	J-CLOSED	800.0	2	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	UNWOOD	99 - NONE	01/20/2023 12:00:00 AM	242.0	243.0	5400X ROOKSTER RD	0.9440204871467	0.46681	
16	Z20934044	J-CLOSED	800.0	2	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	317.0	320.0	6 V 17TH ST	0.59004171273330	0.30495	
17	Z20933775	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	UNWOOD	99 - NONE	01/20/2023 12:00:00 AM	110.0	103.0	300X 8TH ST	0.419204161630980	0.29869	
18	Z20933776	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
19	Z20933777	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	UNWOOD	99 - NONE	01/20/2023 12:00:00 AM	123.0	124.0	0			
20	Z20933778	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
21	Z20933780	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	UNWOOD	99 - NONE	01/20/2023 12:00:00 AM	123.0	124.0	0			
22	Z20933781	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
23	Z20933782	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
24	Z20933783	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
25	Z20933784	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
26	Z20933785	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
27	Z20933786	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
28	Z20933787	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
29	Z20933788	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
30	Z20933789	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
31	Z20933790	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
32	Z20933791	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
33	Z20933792	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
34	Z20933793	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
35	Z20933794	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
36	Z20933795	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
37	Z20933796	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
38	Z20933797	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
39	Z20933798	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
40	Z20933799	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
41	Z20933800	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
42	Z20933801	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
43	Z20933802	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
44	Z20933803	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
45	Z20933804	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
46	Z20933805	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
47	Z20933806	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
48	Z20933807	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
49	Z20933808	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
50	Z20933809	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
51	Z20933810	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
52	Z20933811	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
53	Z20933812	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
54	Z20933813	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
55	Z20933814	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PENNDOT	40 - PERSONAL WEAPONS (HANDS, FEET, TEETH, ETC)	01/20/2023 12:00:00 AM	123.0	124.0	0			
56	Z20933815	J-CLOSED	901.0	1	THEFT	04-PARKING LOT	N-NO BAS NOT APPLICABLE	THURSDAY	PEN								

Improvements on Data Cleaning from the Previous Review:

Outlier Detection

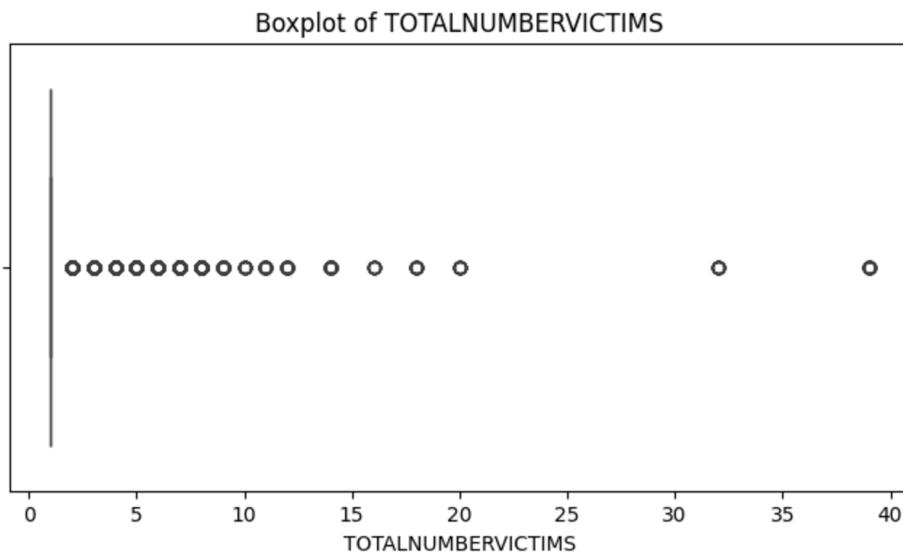
We have performed outlier detection on the numerical columns of our data, and these are the results

```
Column: ZIP
Number of outliers: 24
Lower bound: 45171.5
Upper bound: 45263.5
```



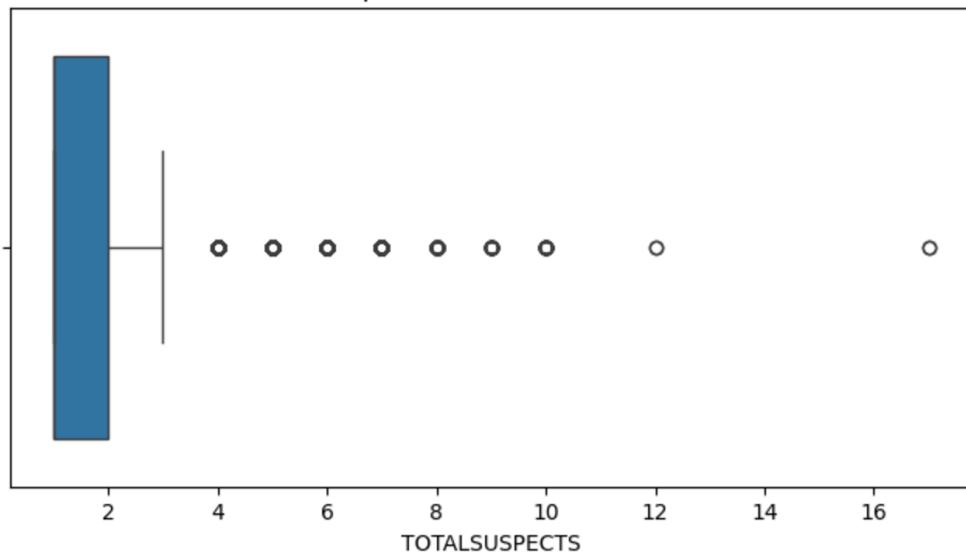
Skipping column HATE_BIAS because it contains only NaN values.

```
Column: TOTALNUMBERVICTIMS
Number of outliers: 28131
Lower bound: 1.0
Upper bound: 1.0
```



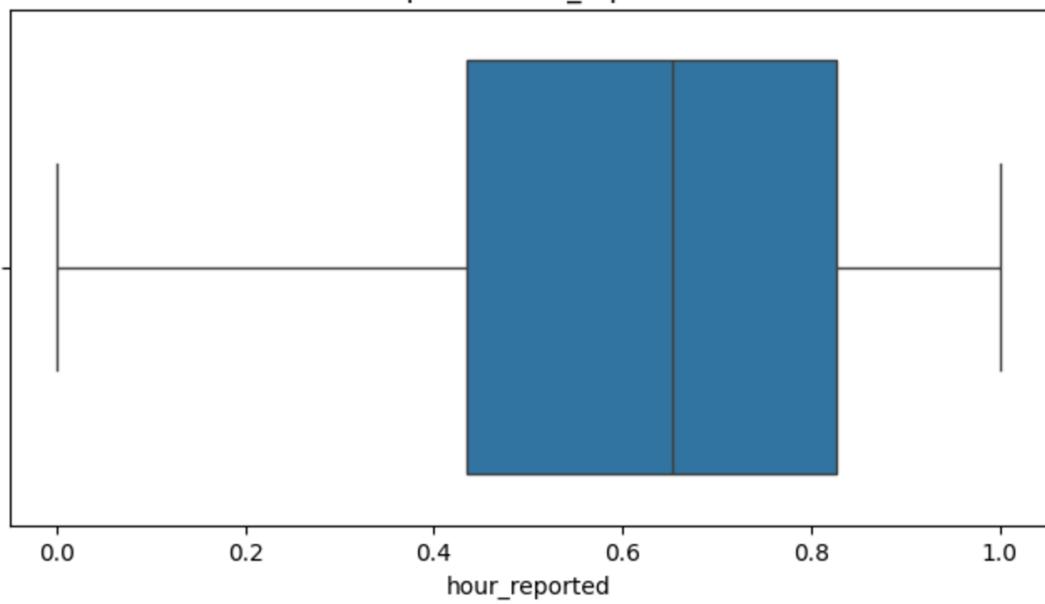
Column: TOTALSUSPECTS
Number of outliers: 6054
Lower bound: -0.5
Upper bound: 3.5

Boxplot of TOTALSUSPECTS



SKIPPING column OFFENSE because it contains only null values.
Column: hour_reported
Number of outliers: 0
Lower bound: -0.15217391304347833
Upper bound: 1.4130434782608696

Boxplot of hour_reported



Outlier Removal

Outlier detection revealed anomalies in the TOTALNUMBERVICTIMS and TOTALSUSPECTS columns. However, these values were retained, as such occurrences could represent valid cases involving a high number of suspects or victims, which are plausible in the context of criminal incidents. For the other columns used in modeling, outlier detection and handling were meticulously performed to ensure the accuracy and reliability of the models.

```
# Function to drop outliers using IQR method
def drop_outliers(df):
    numeric_cols = df.select_dtypes(include=['number']).columns # Select numeric columns
    for col in numeric_cols:
        if col != 'TOTALNUMBERVICTIMS' and col != 'TOTALSUSPECTS' and col != 'HATE_BIAS':
            Q1 = df[col].quantile(0.25) # First quartile
            Q3 = df[col].quantile(0.75) # Third quartile
            IQR = Q3 - Q1 # Interquartile range
            lower_bound = Q1 - 1.5 * IQR # Lower bound
            upper_bound = Q3 + 1.5 * IQR # Upper bound

            # Remove rows outside the bounds
            df = df[(df[col] >= lower_bound) & (df[col] <= upper_bound)]

    return df
```

Perform Null Value checks for the columns that we are using for modelling

```
[ ]: # Calculate and print the number of null values in each column
null_counts = data_subset.isnull().sum()

print("Number of null values in each column:")
print(null_counts)
```

```
Number of null values in each column:
LOCATION          0
SNA_NEIGHBORHOOD 0
ZIP              0
HATE_BIAS         0
TOTALNUMBERVICTIMS 0
TOTALSUSPECTS    0
UCR_GROUP         0
OFFENSE           0
hour_reported     0
dtype: int64
```

Screenshot of the data after cleaning

data_before_modelling									
LOCATION	SNA_NEIGHBORHOOD	ZIP	HATE_BIAS	TOTALNUMBERVICTIMS	TOTALSUSPECTS	UCR_GROUP	OFFENSE	time_reported	
02-MULTI FAMILY APARTMENT	WESTWOOD	45211.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	CRIMINAL DAMAGING/ENDANGERING	11:11:00 PM	
02-MULTI FAMILY	AVONDALE	45229.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	ASSAULT	08:55:20 PM	
48-PARKING LOT	CARTHAGE	45216.0	N-NO BIAS/NOT APPLICABLE	1.0	2.0	ROBBERY	AGGRAVATED ROBBERY	07:36:00 PM	
48-PARKING LOT	CARTHAGE	45216.0	N-NO BIAS/NOT APPLICABLE	1.0	2.0	ROBBERY	AGGRAVATED ROBBERY	07:36:00 PM	
02-MULTI FAMILY APARTMENT	EAST PRICE HILL	45205.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	AGGRAVATED MENACING	06:59:00 PM	
02-MULTI FAMILY APARTMENT	EAST PRICE HILL	45205.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	THEFT	THEFT	06:59:00 PM	
02-MULTI FAMILY APARTMENT	EVANSTON	45207.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	THEFT	THEFT	06:12:00 PM	
02-MULTI FAMILY APARTMENT	AVONDALE	45229.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	AGGRAVATED MENACING	06:00:00 PM	
02-GAMBLING FACILITY / CASINO / RACE TRACK		45202.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	UNAUTHORIZED USE	UNAUTHORIZED USE OF MOTOR VEHICLE	01:22:00 AM	
02-GAMBLING FACILITY / CASINO / RACE TRACK		45202.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	VIOLATE PROTECTION ORDER/CONSENT AGREEMENT	01:22:00 AM	
02-MULTI FAMILY APARTMENT	CUF	45220.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	DOMESTIC VIOLENCE	11:26:00 PM	
02-MULTI FAMILY APARTMENT	HARTWELL	45216.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	ASSAULT	11:05:00 PM	
01-SINGLE FAMILY HOME	WEST END	45203.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	AGGRAVATED MENACING	10:10:00 PM	
02-MULTI FAMILY APARTMENT	WEST END	45214.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	ASSAULT	08:38:00 PM	
01-SINGLE FAMILY HOME	EVANSTON	45207.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	AGGRAVATED MENACING	07:28:34 PM	
48-PARKING LOT	LINWOOD	45226.0	N-NO BIAS/NOT APPLICABLE	2.0	2.0	THEFT	THEFT	03:58:00 AM	
47-STREET	PENOLETON	45202.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	ROBBERY	ROBBERY	03:21:02 AM	
02-MULTI FAMILY APARTMENT	EAST PRICE HILL	45205.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	BURGLARY/BREAKING ENTERING	BURGLARY	02:07:00 AM	
48-PARKING LOT		45223.0	N-NO BIAS/NOT APPLICABLE	1.0	2.0	PART 2 MINOR	ASSAULT	01:23:00 AM	
48-PARKING LOT		45223.0	N-NO BIAS/NOT APPLICABLE	1.0	2.0	PART 2 MINOR	ASSAULT	01:23:00 AM	
20-BARBER/BEAUTY SHOP	WESTWOOD	45211.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	ASSAULT	11:34:00 PM	
02-MULTI FAMILY APARTMENT	MT. AIRY	45239.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	THEFT	THEFT	11:17:00 PM	
02-MULTI FAMILY APARTMENT	PENOLETON	45202.0	N-NO BIAS/NOT APPLICABLE	2.0	1.0	PART 2 MINOR	ASSAULT	07:36:46 PM	
02-MULTI FAMILY APARTMENT	PENOLETON	45202.0	N-NO BIAS/NOT APPLICABLE	2.0	1.0	PART 2 MINOR	ASSAULT	07:36:46 PM	
47-STREET	NORTH AVONDALE - PADDOCK HILLS	45229.0	N-NO BIAS/NOT APPLICABLE	2.0	1.0	ROBBERY	AGGRAVATED ROBBERY	07:29:32 PM	
47-STREET	NORTH AVONDALE - PADDOCK HILLS	45229.0	N-NO BIAS/NOT APPLICABLE	2.0	1.0	ROBBERY	AGGRAVATED ROBBERY	07:29:32 PM	
02-MULTI FAMILY 4 PLEX	MT. AIRY	45239.0	88--DOMESTIC VIOLENCE	1.0	1.0	PART 2 MINOR	DOMESTIC VIOLENCE	07:28:00 PM	
02-MULTI FAMILY		45211.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	AGGRAVATED ASSAULTS	FELONIOUS ASSAULT	06:30:00 PM	
02-MULTI FAMILY	WINTON HILLS	45232.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	PART 2 MINOR	CRIMINAL DAMAGING/ENDANGERING	06:14:00 PM	
02-MULTI FAMILY	WINTON HILLS	45232.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	PART 2 MINOR	IMPROPERLY DISCHARGING FIREARM AT/INTO HABITATION/SCHOOL	06:14:00 PM	
02-MULTI FAMILY	WINTON HILLS	45232.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	PART 2 MINOR	IMPROPERLY DISCHARGING FIREARM AT/INTO HABITATION/SCHOOL	06:14:00 PM	
02-MULTI FAMILY	WINTON HILLS	45232.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	PART 2 MINOR	CRIMINAL DAMAGING/ENDANGERING	06:14:00 PM	
32-OTHER	WEST PRICE HILL	45205.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	PART 2 MINOR	VIOLATE PROTECTION ORDER/CONSENT AGREEMENT	05:36:00 PM	
02-MULTI FAMILY APARTMENT	WESTWOOD	45211.0	N-NO BIAS/NOT APPLICABLE	1.0	1.0	THEFT	THEFT	05:31:38 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	AGGRAVATED ASSAULTS	FELONIOUS ASSAULT	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	HOMICIDE	MURDER	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	ROBBERY	AGGRAVATED ROBBERY	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	ROBBERY	AGGRAVATED ROBBERY	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	ROBBERY	AGGRAVATED ROBBERY	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	AGGRAVATED ASSAULTS	FELONIOUS ASSAULT	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	AGGRAVATED ASSAULTS	FELONIOUS ASSAULT	05:11:00 PM	
48-PARKING LOT	ROSELAWN	45237.0	N-NO BIAS/NOT APPLICABLE	3.0	1.0	HOMICIDE	MURDER	05:11:00 PM	

Screenshot of the data after cleaning and preprocessing and before building the model

[11]:	LOCATION	SNA_NEIGHBORHOOD	ZIP	HATE_BIAS	TOTALNUMBERVICTIMS	TOTALSUSPECTS	UCR_GROUP	OFFENSE	hour_reported
0	8	49	0.082921	16	0.0	0.0000	3	15	1.000000
1	6	0	0.082954	16	0.0	0.0000	3	10	0.869565
2	54	4	0.082930	16	0.0	0.0625	5	6	0.826087
3	54	4	0.082930	16	0.0	0.0625	5	6	0.826087
4	8	12	0.082910	16	0.0	0.0000	3	4	0.782609

These are the only few columns that we are using for model building in our data.we have converted categorical values to numerical values for model building.

Label Encoding

```
] categorical_columns = ['LOCATION', 'SNA_NEIGHBORHOOD', 'HATE_BIAS', 'UCR_GROUP', 'OFFENSE']
label_encoders = {col: LabelEncoder() for col in categorical_columns}
for col in categorical_columns:
    data_subset[col] = label_encoders[col].fit_transform(data_subset[col])

]: scaler = MinMaxScaler()
numeric_columns = ['ZIP', 'TOTALNUMBERVICTIMS', 'TOTALSUSPECTS', 'hour_reported']
data_subset[numeric_columns] = scaler.fit_transform(data_subset[numeric_columns])

]: data_subset.head()

]:   LOCATION SNA_NEIGHBORHOOD      ZIP HATE_BIAS TOTALNUMBERVICTIMS TOTALSUSPECTS UCR_GROUP OFFENSE hour_reported
0       8          49  0.082921        16           0.0      0.0000         3     15  1.000000
1       6          0  0.082954        16           0.0      0.0000         3     10  0.869565
2      54          4  0.082930        16           0.0      0.0625         5      6  0.826087
3      54          4  0.082930        16           0.0      0.0625         5      6  0.826087
4       8          12  0.082910        16           0.0      0.0000         3      4  0.782609
```

The provided code preprocesses the dataset by transforming categorical and numerical features to prepare it for machine learning. Categorical columns (LOCATION, SNA_NEIGHBORHOOD, etc.) are label-encoded, converting unique categories into integer values. Numerical columns (ZIP, TOTALNUMBERVICTIMS, etc.) are scaled using MinMaxScaler to normalize their range between 0 and 1. These steps ensure that categorical data is converted into a format suitable for modeling, and numerical features are on a consistent scale, improving model performance and interpretability. Finally, the processed dataset is displayed to confirm the transformations.

Splitting the data to testing and training

```
[ ]: X = data_subset.drop(columns=['OFFENSE'])
y = data_subset['OFFENSE']

[ ]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

[ ]: X_train = X_train.values
X_test = X_test.values

[ ]: X_train = X_train.reshape(X_train.shape[0], -1)
X_test = X_test.reshape(X_test.shape[0], -1)

[ ]: scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

[ ]: time_steps = 1
X_train = X_train.reshape(X_train.shape[0], X_train.shape[1], time_steps)
X_test = X_test.reshape(X_test.shape[0], X_test.shape[1], time_steps)

[ ]: print("Reshaped X_train shape:", X_train.shape)
print("Reshaped X_test shape:", X_test.shape)

Reshaped X_train shape: (94940, 8, 1)
Reshaped X_test shape: (23736, 8, 1)
```

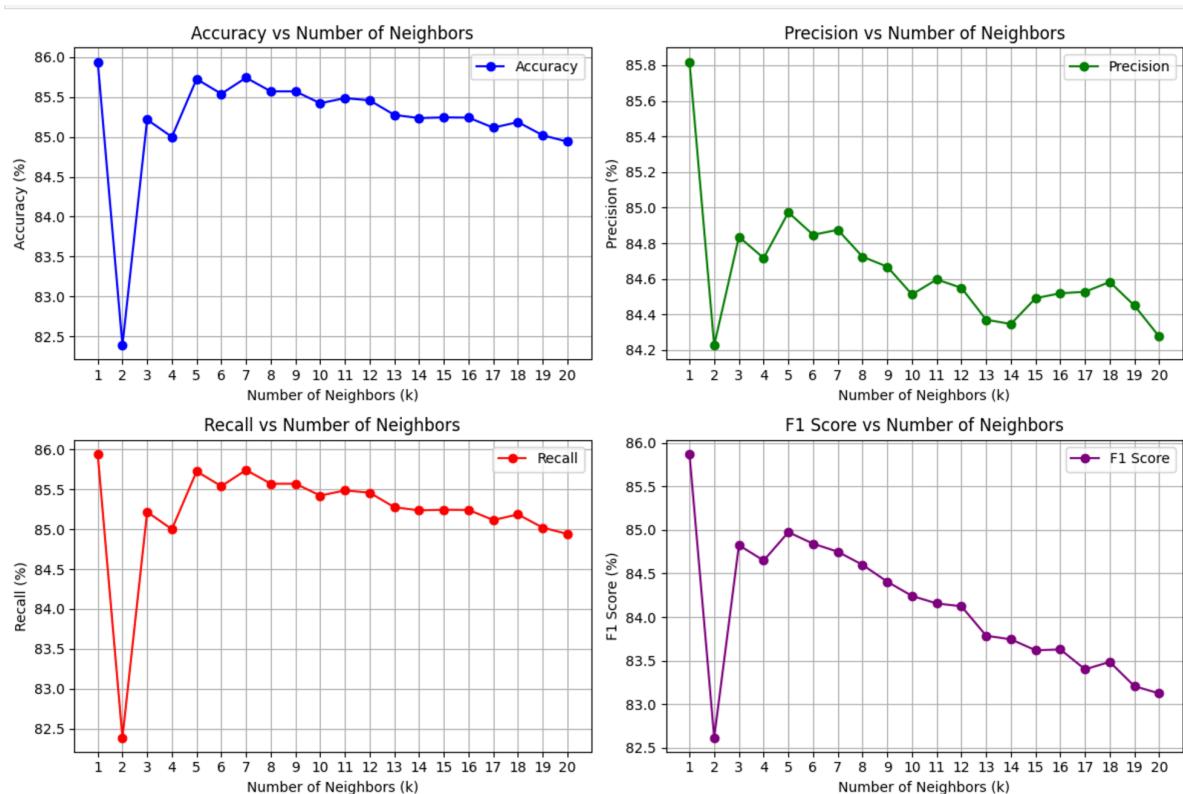
The code prepares the dataset for model training by separating features and the target variable. The feature set X is created by dropping the target column OFFENSE from the dataset, while the target variable y contains the OFFENSE values. The data is then split into training and testing subsets using an 80-20 split via train_test_split, ensuring reproducibility with a fixed random_state. The features in X_train and X_test are converted into NumPy arrays for further processing. A StandardScaler is applied to normalize the training and testing features, transforming them to have a mean of 0 and a standard deviation of 1. Next, the data is reshaped to ensure compatibility with models that require specific input dimensions, such as recurrent neural networks. The final shapes of X_train and X_test are printed to confirm the transformations, with the reshaped data now having dimensions suitable for time-series analysis or similar tasks.

DATA MODELLING

1.K NEAREST NEIGHBOUR

TARGET VARIABLE – UCR_GROUP

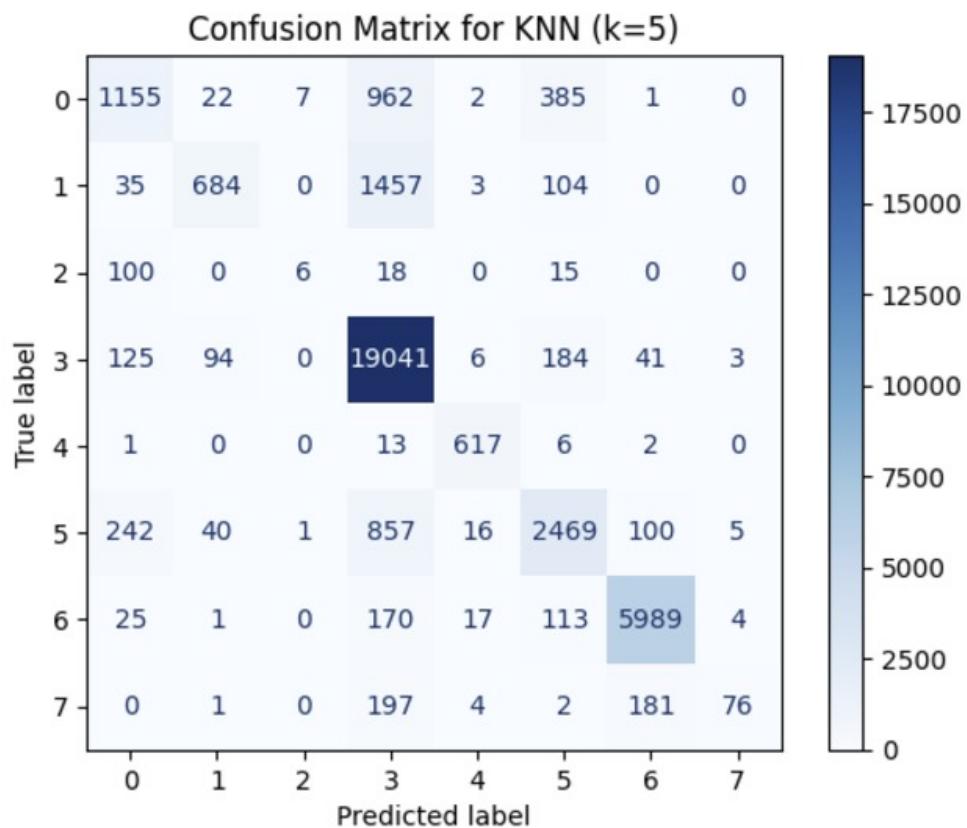
In our KNN modeling, we have used UCR (Uniform Crime Reporting) as our target variable because it provides a standardized way of categorizing crimes into broad classifications, such as violent crimes or property crimes. Using UCR as the target allows us to focus on understanding crime patterns and trends, which can help in resource allocation, law enforcement planning, and crime prevention efforts. Since UCR is a universally recognized crime reporting system, it makes our analysis more relevant and actionable, aligning it with real-world applications like policymaking and identifying crime hotspots. This ensures that our findings can be effectively used to address and understand criminal activity.



→ KNN Accuracy (using ADDRESS): 82.14%
KNN F1-Score (using ADDRESS): 79.96%

Classification Report:

	precision	recall	f1-score	support
0	0.69	0.46	0.55	2534
1	0.82	0.30	0.44	2283
2	0.46	0.04	0.08	139
3	0.85	0.95	0.90	19494
4	0.92	0.95	0.94	639
5	0.74	0.63	0.68	3730
6	0.80	0.92	0.86	6319
7	0.50	0.00	0.00	461
accuracy			0.82	35599
macro avg	0.72	0.53	0.56	35599
weighted avg	0.81	0.82	0.80	35599



The K-Nearest Neighbors (KNN) model achieved an accuracy of **82.14%** and an F1-score of **79.96%** on the test dataset. The classification report provides detailed metrics for each class, including precision, recall, and F1-score. The model performed exceptionally well for

class 3 and class 4, with F1-scores of **0.90** and **0.94**, respectively, indicating that these classes are predicted with high precision and recall. However, performance for class 2 was poor, with an F1-score of **0.08**, and class 7 had a score of **0.00**, indicating challenges in identifying these classes. The macro-average F1-score is **0.56**, reflecting the imbalance in performance across classes, while the weighted average F1-score is **0.80**, indicating the overall performance weighted by class support. These results suggest that while the model works well for majority classes, it struggles with underrepresented or minority classes, which might require further optimization or rebalancing of the dataset.

High Accuracy for Class 3: Class 3 has the highest correctly classified instances with 19,041, indicating excellent performance for this class.

Misclassifications: Class 0 shows significant misclassification, with 962 instances predicted as Class 3. Class 5 has 857 instances misclassified as Class 3 and 242 instances misclassified as Class 0. Other classes like 1, 4, 6, and 7 also show some degree of misclassification, though less severe.

Confusion for Neighboring Classes: The confusion is noticeable for classes that are "neighboring" or similar in data space, e.g., misclassification between Classes 5 and 3.

Sparse Misclassifications: Certain classes have very few non-zero values for off-diagonal elements, such as Class 7, where most predictions fall under Class 7 with minor misclassifications.

Dominant Class Sizes: Class 3 dominates the dataset in terms of size, as reflected by the large count of true and predicted values compared to other classes.

2. RANDOM FOREST

```
[ ]: from sklearn.ensemble import RandomForestClassifier
label_encoder = LabelEncoder()
data['VICTIM_ETHNICITY'] = label_encoder.fit_transform(data['VICTIM_ETHNICITY'].astype(str))

[ ]: data['UCR_GROUP'] = label_encoder.fit_transform(data['UCR_GROUP'])

X = data[['UCR_GROUP']]
y = data['VICTIM_ETHNICITY']

[ ]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

rf = RandomForestClassifier(random_state=42)
rf.fit(X_train, y_train)

129]: RandomForestClassifier(random_state=42)
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

[ ]: from sklearn.metrics import accuracy_score, classification_report, f1_score
y_pred = rf.predict(X_test)
accuracy = accuracy_score(y_test, y_pred) * 100
f1 = f1_score(y_test, y_pred, average='weighted') * 100

[ ]: print(f"Random Forest Accuracy (using ADDRESS): {accuracy:.2f}%")
print(f"Random Forest F1-Score (using ADDRESS): {f1:.2f}%")
print("\nClassification Report:")
print(classification_report(y_test, y_pred))

Random Forest Accuracy (using ADDRESS): 93.08%
Random Forest F1-Score (using ADDRESS): 89.74%
```

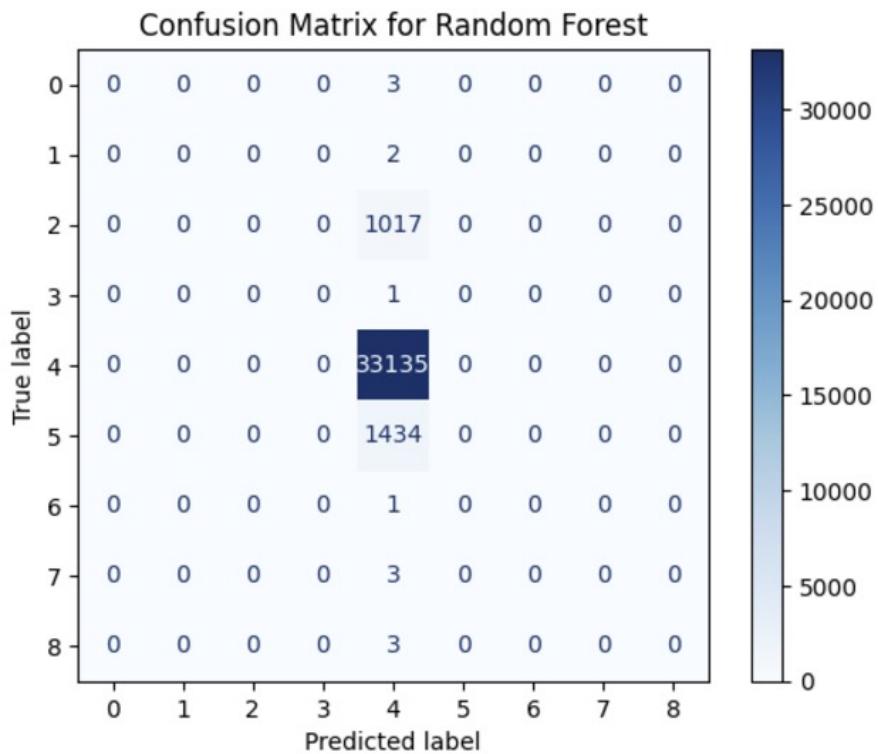
In this model, we are predicting the VICTIM_ETHNICITY using UCR_GROUP as the input feature. The significance of calculating this variable lies in its potential to provide insights into how different crime categories are associated with specific victim demographics. This can help law enforcement and policymakers identify trends, address potential biases, and allocate resources effectively to support vulnerable groups. By understanding these relationships, interventions can be tailored to improve community safety and inclusivity.

Random Forest Accuracy (using ADDRESS): 93.08%

Random Forest F1-Score (using ADDRESS): 89.74%

Classification Report:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	3
1	0.00	0.00	0.00	2
2	0.00	0.00	0.00	1017
3	0.00	0.00	0.00	1
4	0.93	1.00	0.96	33135
5	0.00	0.00	0.00	1434
6	0.00	0.00	0.00	1
7	0.00	0.00	0.00	3
8	0.00	0.00	0.00	3
accuracy				0.93
macro avg	0.10	0.11	0.11	35599
weighted avg	0.87	0.93	0.90	35599



The model's output shows an accuracy of **93.00%**, meaning that the predictions are highly accurate for the test dataset. The F1-score is **0.90**, which indicates that the model has a strong balance between precision and recall, ensuring reliable results across various

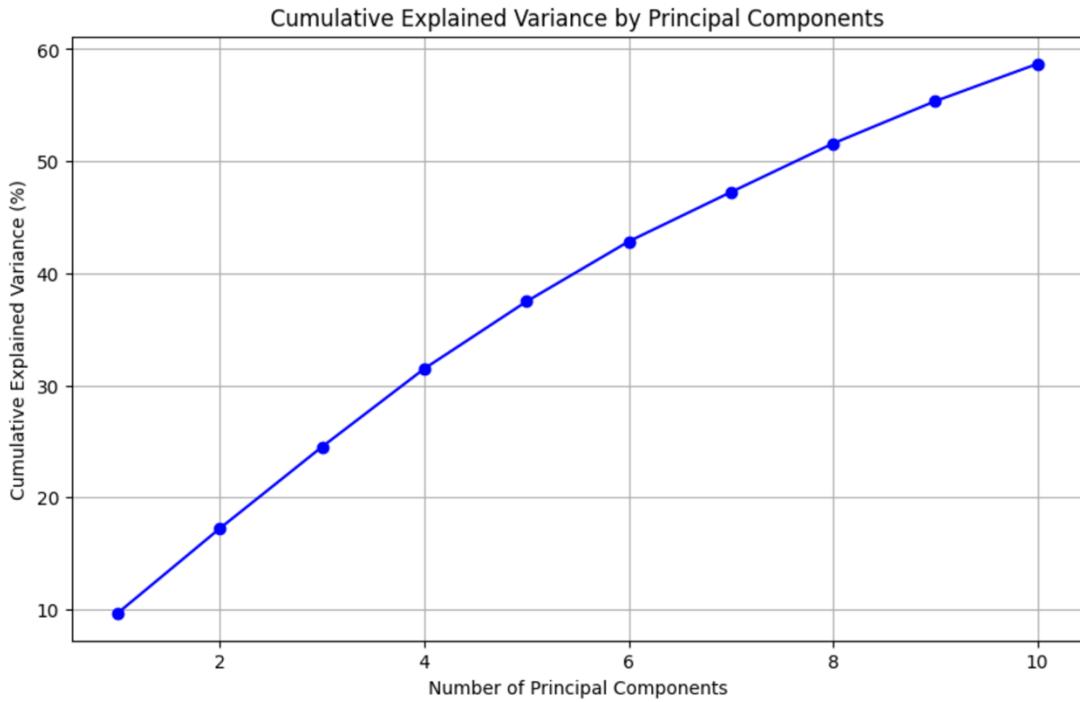
categories. Overall, the Random Forest classifier demonstrates its effectiveness in classifying victim ethnicity based on crime group data. However, further fine-tuning, feature engineering, or incorporating additional variables might enhance the model's performance and provide even deeper insights.

True labels are shown on the vertical axis, and predicted labels are on the horizontal axis. Most correct predictions are concentrated on the diagonal of the matrix, particularly for the class corresponding to True Label 4 (Predicted Label 4) with 33,135 correct predictions. Misclassifications are minimal but present: 1,017 instances of True Label 2 were predicted as Label 4. 1,434 instances of True Label 5 were predicted as Label 4. Classes with very low sample counts (e.g., Labels 0, 1, 3, 6, 7, and 8) show either no predictions or very minimal counts in both correct and incorrect predictions.

Overall, the model seems highly biased toward predicting Label 4, as most predictions fall in this category.

PRINCIPAL COMPONENT ANALYSIS

In the next step, we performed Principal Component Analysis (PCA) on the dataset to reduce dimensionality while retaining the maximum variance in the data. The graph above shows the cumulative explained variance as a function of the number of principal components. Each additional component contributes to the total variance captured, with the cumulative explained variance reaching around 60% after including 10 principal components. This indicates that these components collectively retain a significant portion of the data's variability, making PCA an effective preprocessing step. By using PCA, we aim to reduce computational complexity, eliminate redundant features, and enhance the performance of our models, especially when dealing with high-dimensional data. The curve also helps identify the optimal number of components to balance variance retention and dimensionality reduction.



3. PRINCIPAL COMPONENT ANALYSIS – RANDOM FOREST

```
[ ]: from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column].astype(str))
    label_encoders[column] = le

[ ]: data.fillna(data.mean(), inplace=True)

[ ]: y = data['THEFT_CODE']
X = data.drop(columns=['THEFT_CODE'])

[ ]: X_scaled = scaler.fit_transform(X)

[ ]: n_components = min(X.shape[1], 10)
pca = PCA(n_components=n_components)
X_pca = pca.fit_transform(X_scaled)

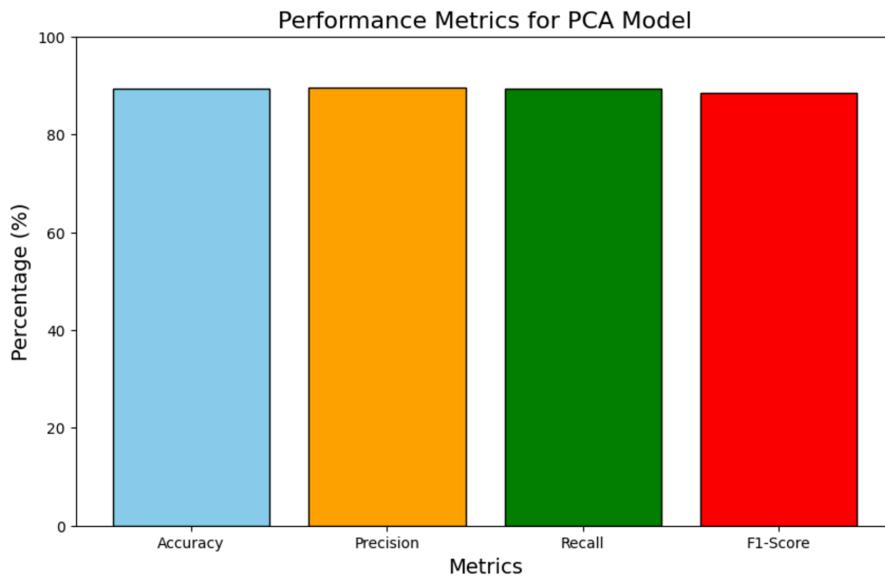
[ ]: X_train, X_test, y_train, y_test = train_test_split(X_pca, y, test_size=0.3, random_state=42)

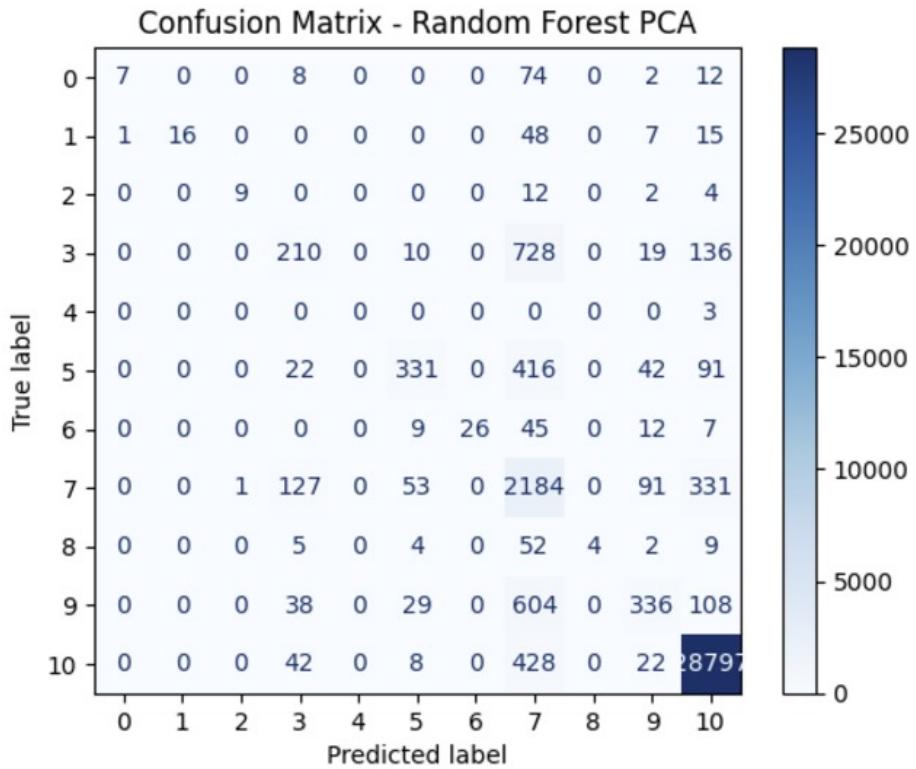
[ ]: from sklearn.ensemble import RandomForestClassifier
rf = RandomForestClassifier(random_state=42)
rf.fit(X_train, y_train)
```

The output shows the performance metrics for the Random Forest model trained on the PCA-transformed dataset. The model achieved an **accuracy of 89.37%**, indicating that nearly 9 out of 10 predictions were correct. The **precision of 89.57%** reflects the model's ability to minimize false positives and correctly identify true positives. The **recall of 89.37%** shows the model's capability to detect most of the actual positive cases, while the **F1-score of 88.43%** provides a balanced metric that considers both precision and recall.

it likely improved the computational efficiency and performance by reducing noise and redundancy in the features while retaining essential variance. This combination of PCA with Random Forest provides a robust approach for handling high-dimensional datasets effectively.

Classification Report:				
	precision	recall	f1-score	support
0	0.89	0.08	0.14	103
1	1.00	0.20	0.33	87
2	0.90	0.33	0.49	27
3	0.45	0.18	0.26	1103
4	0.00	0.00	0.00	3
5	0.71	0.36	0.48	902
6	0.96	0.23	0.37	99
7	0.48	0.73	0.58	2787
8	1.00	0.05	0.10	76
9	0.65	0.32	0.43	1115
10	0.97	0.98	0.98	29297
accuracy			0.89	35599
macro avg	0.73	0.32	0.38	35599
weighted avg	0.90	0.89	0.88	35599





Diagonal dominance:

Correct predictions (values on the diagonal) are most significant for True Label 10 (Predicted Label 10), with 28,797 correct predictions, indicating strong performance for this class.

Other classes have fewer correct predictions, such as: True Label 7: 2,184 correct predictions. True Label 6: 728 correct predictions.

Misclassifications: Many misclassifications occur across different labels: True Label 6 is often confused with Labels 7, 8, and 10. True Label 7 has significant misclassifications, with 331 instances predicted as Label 10.

True Label 5 shows notable misclassifications across Labels 6, 7, and 10. Some rows (e.g., True Labels 0, 1, 2) show very low correct prediction counts and significant misclassification spread.

Class bias: Predictions are skewed heavily towards Label 10, which has the largest count overall, indicating a potential class imbalance.

Low performance on rare classes: Classes with True Labels 0, 1, 2, and 4 have very few correct predictions, indicating poor performance on these less frequent labels.

4. SUPPORT VECTOR MACHINE

```
[ ]: from sklearn.svm import SVC
[ ]: from sklearn.linear_model import SGDClassifier
[ ]: from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

[ ]: label_encoders = {}
[ ]: for column in data.select_dtypes(include=['object']).columns:
[ ]:     le = LabelEncoder()
[ ]:     data[column] = le.fit_transform(data[column].astype(str))
[ ]:     label_encoders[column] = le

[ ]: X.fillna(X.median(), inplace=True) # Fill numerical NaNs with median
[ ]: categorical_cols = X.select_dtypes(include=['object']).columns
[ ]: for col in categorical_cols:
[ ]:     X[col].fillna("Unknown", inplace=True)

[ ]: y = data['TOTALSUSPECTS']
[ ]: X = data.drop(columns=['TOTALSUSPECTS'])

[ ]: data.fillna(data.mean(), inplace=True)

[ ]: scaler = StandardScaler()
[ ]: X_scaled = scaler.fit_transform(X)

[ ]: sgd_svc = SGDClassifier(loss='hinge', max_iter=1000, tol=1e-3, random_state=42)
[ ]: X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
[ ]: sgd_svc.fit(X_train, y_train)
[ ]: scores = cross_val_score(sgd_svc, X_scaled, y, cv=5)
```

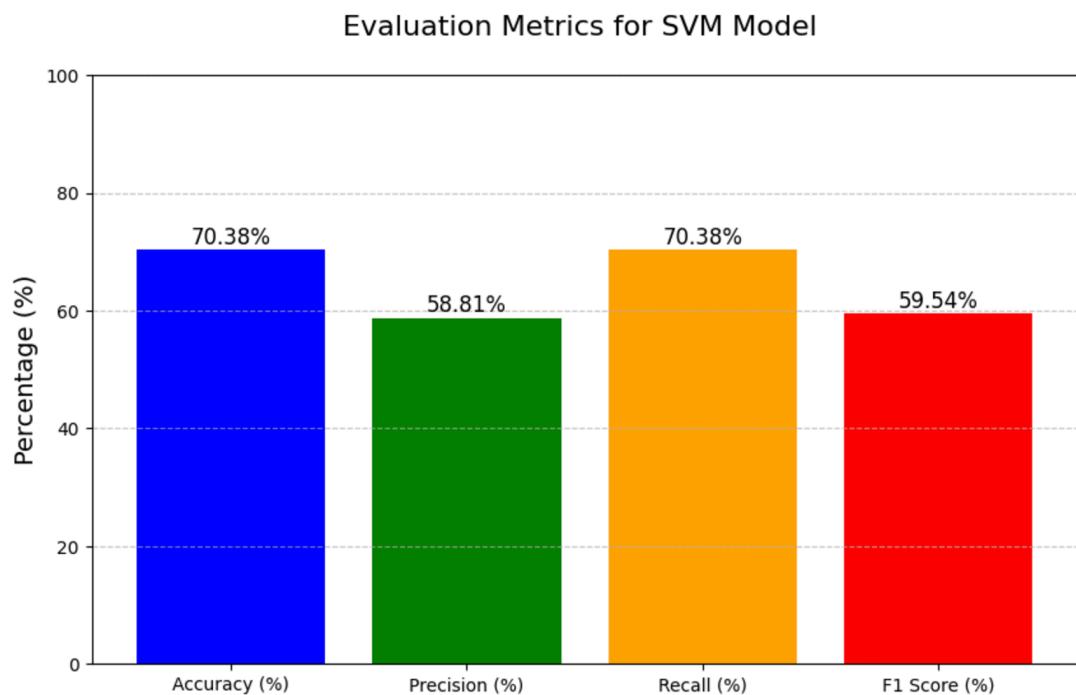
The code implements a Support Vector Machine (SVM) classifier using the SGDClassifier with a hinge loss function, which is commonly used for linear SVMs. The variable being predicted in this model is **TOTALSUSPECTS**, which represents the total number of suspects associated with each data point.

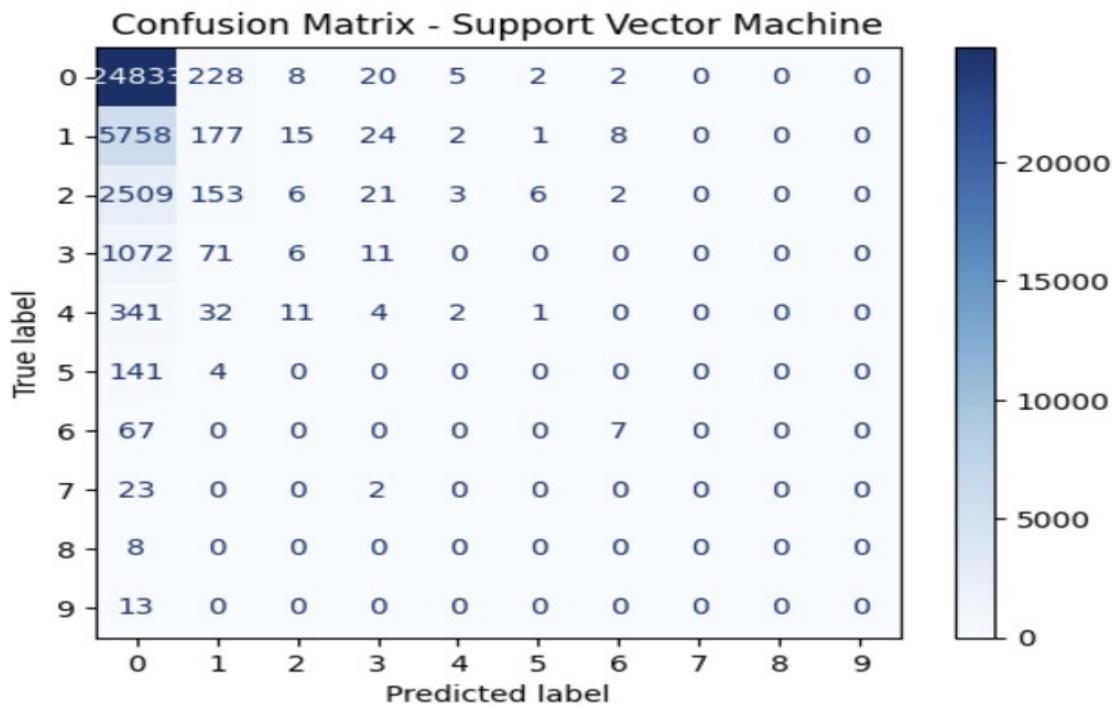
Predicting **TOTALSUSPECTS** is significant because it provides valuable insights into the patterns and characteristics associated with criminal incidents. By accurately forecasting the number of suspects involved in a case, law enforcement agencies can better allocate resources, prioritize investigations, and identify trends that could indicate larger issues, such as organized crime or recurring offenders. This information can also support data-driven decision-making in crime prevention strategies, helping policymakers address underlying factors contributing to higher numbers of suspects in certain areas or scenarios. Additionally, such predictions can enhance public safety by enabling more focused and effective interventions.

The preprocessing steps involve encoding categorical variables with LabelEncoder, filling missing numerical values with their median, and replacing missing categorical values with "Unknown." The target variable, TOTALSUSPECTS, is separated from the features, and all missing values in the dataset are replaced with the mean for numerical columns. The features are then standardized using StandardScaler to ensure that they are on the same

scale. Finally, the dataset is split into training and testing sets with a 70-30 split, and the SVM classifier is trained on the scaled training data.

The performance metrics show that the model achieved an **accuracy of 70.38%**, meaning it correctly predicted the total number of suspects in approximately 7 out of 10 cases. The **precision of 58.81%** reflects moderate success in minimizing false positives, while the **recall of 70.38%** indicates a good ability to identify most true cases. The **F1-score of 59.54%** provides a balanced measure of precision and recall. These results suggest reasonable performance but highlight opportunities for improvement through techniques like feature engineering or dataset balancing.





Correct Predictions (Diagonal Values): The largest number of correct predictions is for True Label 0, with 24,833 correctly classified instances. True Labels 1 and 2 also have relatively high correct predictions, with 5,758 and 2,509, respectively. Correct predictions for other labels drop significantly, with minimal correct classifications for labels like 3 through 9.

Misclassifications: Most misclassifications occur in the neighboring classes: For True Label 0, 228 instances are classified as Label 1, and 20 are classified as Label 3. For True Label 1, 177 instances are classified as Label 0, and 24 as Label 3. Misclassifications decrease significantly for labels with fewer samples (e.g., Labels 5, 6, 7, 8, and 9).

Class Imbalance or Bias: The model performs significantly better for dominant labels like 0, 1, and 2, suggesting potential class imbalance or insufficient separability for other labels. Rare or less-represented labels (5 through 9) have negligible correct predictions and are mostly misclassified as Label 0 or 1.

Model Performance: The overall trend shows that the model struggles to generalize for minor classes. The diagonal dominance is visible for major classes, but the confusion for less-represented classes highlights performance issues.

5.NAIVE BIAS

```
] from sklearn.naive_bayes import GaussianNB
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, classification_report
import matplotlib.pyplot as plt

]: label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    le = LabelEncoder()
    data[column] = le.fit_transform(data[column].astype(str))
    label_encoders[column] = le

]: data.fillna(data.mean(), inplace=True)

y = data['THEFT_CODE']
X = data.drop(columns=['THEFT_CODE'])

]: scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

]: X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)

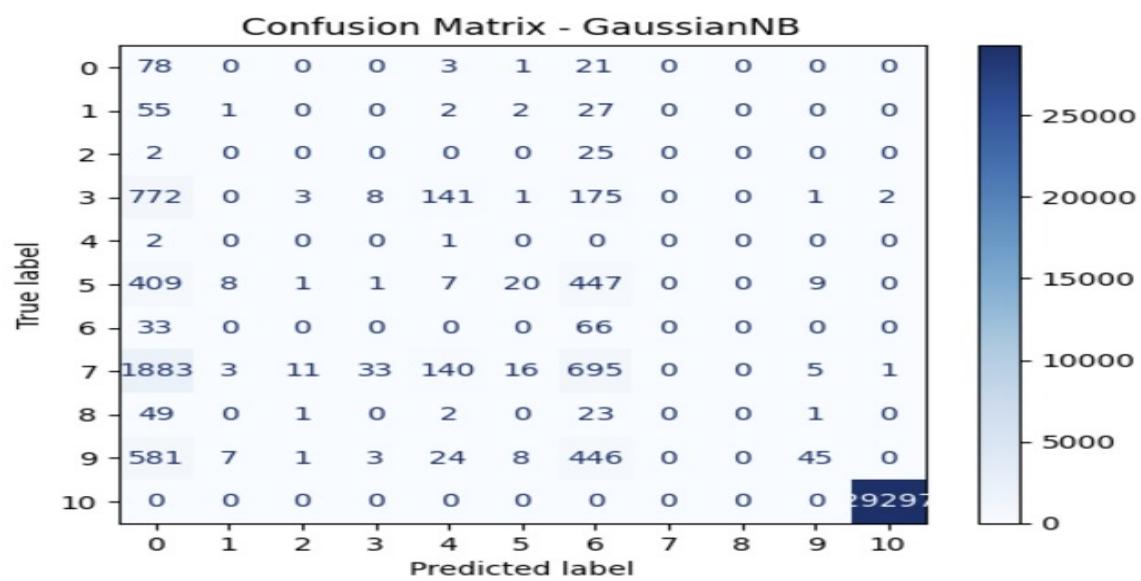
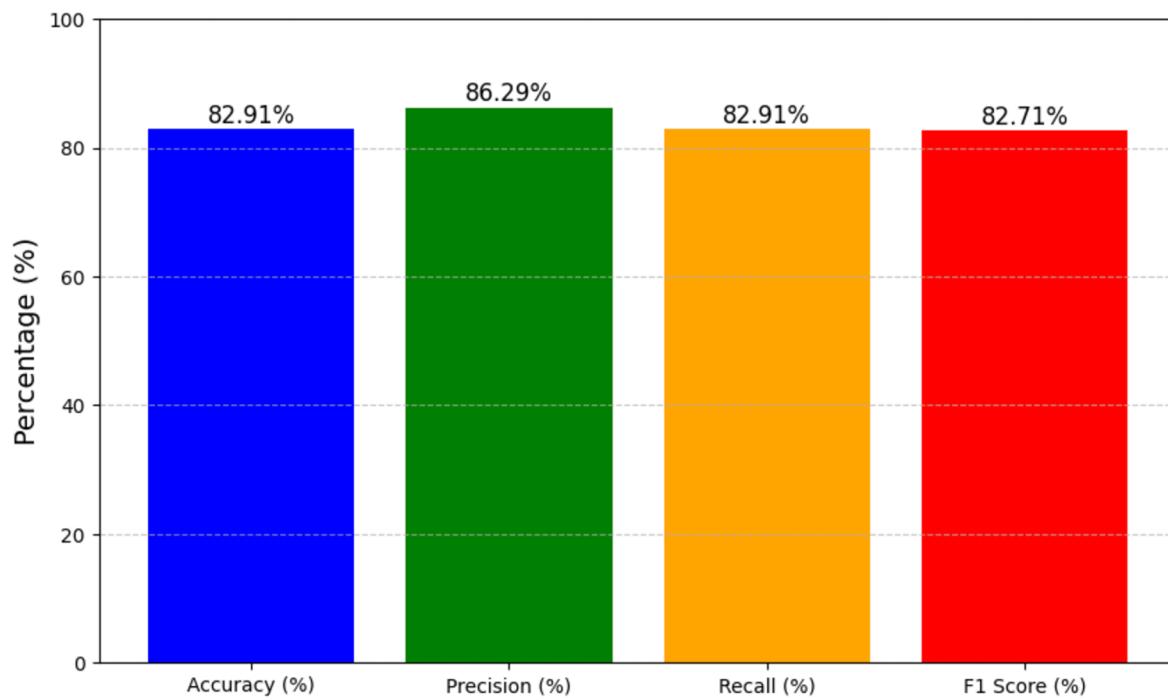
]: nb = GaussianNB()
nb.fit(X_train, y_train)
GaussianNB()
```

The variable being predicted by the Naive Bayes model is TOTALSUSPECTS, which represents the total number of suspects involved in each case. This prediction is significant because it provides valuable insights into crime patterns and allows law enforcement agencies to better understand the scale of incidents, allocate resources efficiently, and identify potential trends in criminal activity. By forecasting the number of suspects, this model can aid in strategic decision-making and help prioritize cases that may require more immediate attention.

Explanation

The Naive Bayes model achieved an accuracy of 82.91%, meaning it correctly predicted the number of suspects in approximately 83% of the test cases. The precision of 86.29% highlights the model's ability to avoid false positives, ensuring the predictions are specific and reliable. The recall of 82.91% indicates that the model successfully identified most true cases, capturing the majority of actual incidents involving suspects. The F1-score of 82.71% shows a balanced performance in terms of precision and recall, making the model robust for practical applications.

These results demonstrate that the Naive Bayes model effectively predicts TOTALSUSPECTS, providing a probabilistic approach to understanding and analyzing crime data. While the performance is strong, further optimization or additional features might improve the results even further, enhancing its utility for law enforcement and crime prevention efforts.



Strong Class 10 Performance: The model shows perfect classification for class 10 (29297 correct predictions), which appears to be the majority class based on the color intensity.

Major Misclassification Patterns: Class 7 has significant misclassifications: 1883 instances were incorrectly predicted as class 0, and 695 as class 6. Class 3 also shows notable confusion: 772 instances predicted as class 0, and 175 as class 6. Class 5 has 409 instances misclassified as class 0, and 447 as class 6.

Common Error Pattern: Classes 0 and 6 appear to be frequent incorrect predictions across multiple true classes, suggesting the model might be biased towards these classes.

Weak Performance Areas: Classes 7 and 8 show very few correct predictions. The model seems to rarely predict classes 7, 8, and 9, as indicated by the mostly zeros in these columns.

Diagonal Performance: The diagonal elements (representing correct predictions) are notably weaker for most classes compared to the misclassifications, indicating suboptimal overall performance except for class 10.

6. XGBOOST

```
[ ]: from xgboost import XGBRegressor
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, classification_report
import matplotlib.pyplot as plt

[ ]: label_encoders = {}
for column in data.select_dtypes(include=['object']).columns:
    if column != 'ZIP': # Don't encode the target
        le = LabelEncoder()
        data[column] = le.fit_transform(data[column].astype(str))
        label_encoders[column] = le

    # Handle missing values
data.fillna(data.mean(), inplace=True)

# Step 2: Select ZIP as the target variable
y = data['SNA_NEIGHBORHOOD'] # Use the 'ZIP' column as the target
X = data.drop(columns=['SNA_NEIGHBORHOOD']) # Drop the target column from features

# Step 3: Check if ZIP needs label encoding
if y.dtype == 'object' or isinstance(y[0], str):
    le = LabelEncoder()
    y = le.fit_transform(y) # Encode target as numeric if necessary

# Step 4: Standardize features
scaler = StandardScaler()
X_scaled = scaler.fit_transform(X)

# Step 5: Train-Test Split
X_train, X_test, y_train, y_test = train_test_split(X_scaled, y, test_size=0.3, random_state=42)
```

```

|: from sklearn.metrics import mean_absolute_error, mean_squared_error
xgb_regressor = XGBRegressor(random_state=42)
xgb_regressor.fit(X_train, y_train)

# Step 7: Evaluate the Model
y_pred = xgb_regressor.predict(X_test)

mae = mean_absolute_error(y_test, y_pred)
mse = mean_squared_error(y_test, y_pred)
rmse = mean_squared_error(y_test, y_pred, squared=False)

print(f"XGBoost Model Mean Absolute Error (MAE): {mae:.2f}")
print(f"XGBoost Model Mean Squared Error (MSE): {mse:.2f}")
print(f"XGBoost Model Root Mean Squared Error (RMSE): {rmse:.2f}")

XGBoost Model Mean Absolute Error (MAE): 0.88
XGBoost Model Mean Squared Error (MSE): 7.81
XGBoost Model Root Mean Squared Error (RMSE): 2.79

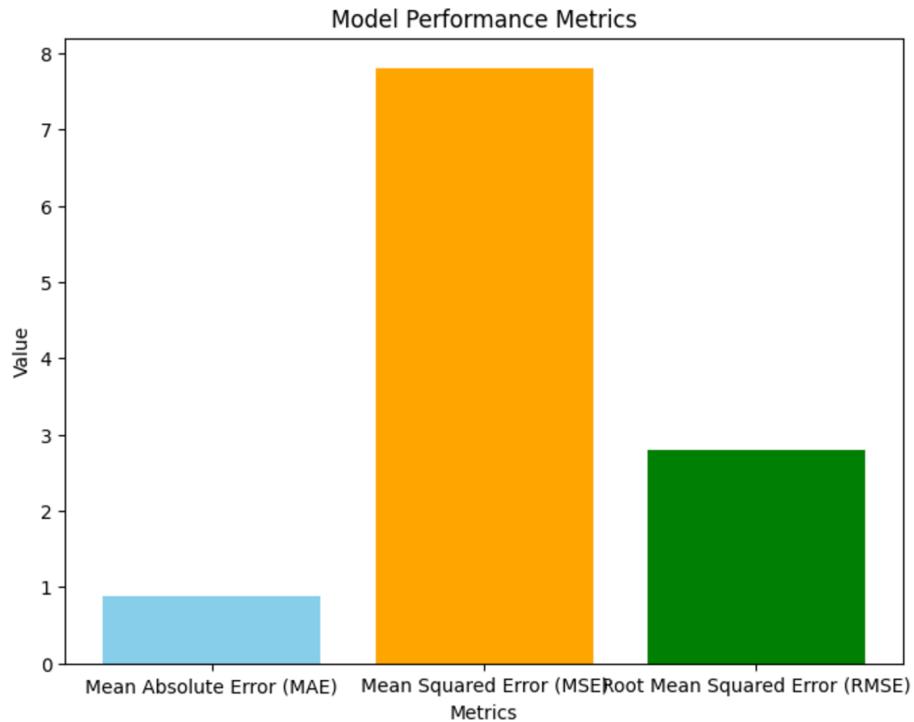
```

The code implements the XGBoost Regressor model to predict the target variable SNA_NEIGHBORHOOD, which likely represents a specific category or neighborhood-related feature. Predicting this variable is significant as it can help uncover localized patterns, trends, or relationships in the data. For example, understanding neighborhood-specific metrics could aid in better resource allocation, identifying hotspots for intervention, or gaining insights into regional variations that inform decision-making processes.

The data preprocessing involves several key steps to ensure the dataset is clean and ready for modeling. Categorical variables are label-encoded to transform them into numerical formats suitable for machine learning models. Missing numerical values are filled with their mean to maintain data completeness and prevent training issues. The SNA_NEIGHBORHOOD column is designated as the target variable, while the remaining columns serve as input features. Standardization is then applied to scale all numerical features, ensuring that the model processes features on a uniform scale, which is critical for gradient-boosting algorithms like XGBoost. The dataset is then split into training and testing sets, with 70% of the data used for training and 30% reserved for evaluation.

After training the XGBoost Regressor on the scaled training data, the model's performance is evaluated using metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), and Root Mean Squared Error (RMSE). The MAE of 0.88 indicates that the predictions, on average, deviate by less than one unit from the actual values. The MSE of 7.81 measures the squared average error, penalizing larger deviations more heavily, while the RMSE of 2.79 provides an interpretable measure of error in the same units as the target variable. These metrics demonstrate that the XGBoost model effectively captures patterns in the data with relatively low error.

Overall, the ability to predict SNA_NEIGHBORHOOD with this level of accuracy highlights the value of the model in providing actionable insights for neighborhood-specific analysis. The robustness of XGBoost, combined with effective preprocessing, ensures that the model can handle complex relationships and feature interactions.



7. LSTM

The code implements an LSTM (Long Short-Term Memory) model to predict the target variable HATE_BIAS, which categorizes hate-related biases. This prediction is highly significant as it provides insights into patterns of hate-related incidents, enabling law enforcement agencies and policymakers to identify and address targeted discrimination and violence effectively. By understanding these biases, proactive measures can be implemented to foster safer, more inclusive communities and improve support systems for those affected.

The dataset is preprocessed by separating the HATE_BIAS column as the target variable, while the remaining columns serve as features. The data is split into training and testing sets, with 80% used for training and 20% for testing. Since LSTM models require 3D input, the datasets are reshaped into the required format with dimensions of samples, time steps, and features. The target variable is also converted into categorical format, corresponding to the number of unique hate bias categories, to enable multi-class classification.

The LSTM model architecture includes an LSTM layer with 64 units and ReLU activation to capture temporal dependencies, followed by dropout layers to reduce overfitting. Dense layers handle the classification task, with the final output layer using softmax activation to predict probabilities for each hate bias category. The model is compiled with the adam optimizer and the categorical_crossentropy loss function and is trained for 10 epochs with a batch size of 64, using the test data for validation.

The model performed exceptionally well, achieving an accuracy of 96.35%, a precision of 95.98%, a recall of 96.35%, and an F1-score of 95.99%, indicating a strong balance between precision and recall. These metrics demonstrate that the model can effectively predict hate bias categories with high reliability. The use of LSTM ensures that temporal patterns and relationships in the data are captured, making it particularly suitable for sequential datasets. The results highlight the potential of this approach to provide actionable insights, supporting efforts to reduce hate-related incidents and improve community safety.

```
# Prepare features (X) and target (y)
X = data_subset.drop(columns=['HATE_BIAS'])
y = data_subset['HATE_BIAS']

# Split into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Reshape the data for LSTM input: (samples, time_steps, features)
time_steps = 1 # LSTM input requires 3D data
X_train_lstm = np.reshape(X_train.values, (X_train.shape[0], time_steps, X_train.shape[1]))
X_test_lstm = np.reshape(X_test.values, (X_test.shape[0], time_steps, X_test.shape[1]))

# Convert the target to categorical format for classification
num_classes = len(np.unique(y)) # Total number of hate bias categories
y_train_cat = to_categorical(y_train, num_classes)
y_test_cat = to_categorical(y_test, num_classes)

# Define the LSTM model architecture
model = Sequential([
    LSTM(64, input_shape=(time_steps, X_train.shape[1]), activation='relu', return_sequences=False),
    Dropout(0.2),
    Dense(32, activation='relu'),
    Dropout(0.2),
    Dense(num_classes, activation='softmax') # Output layer for classification
])

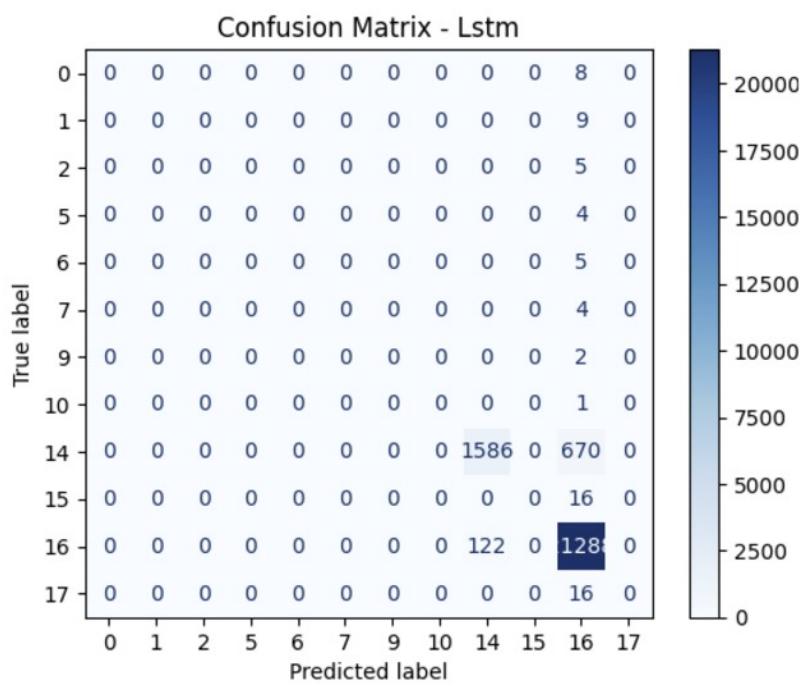
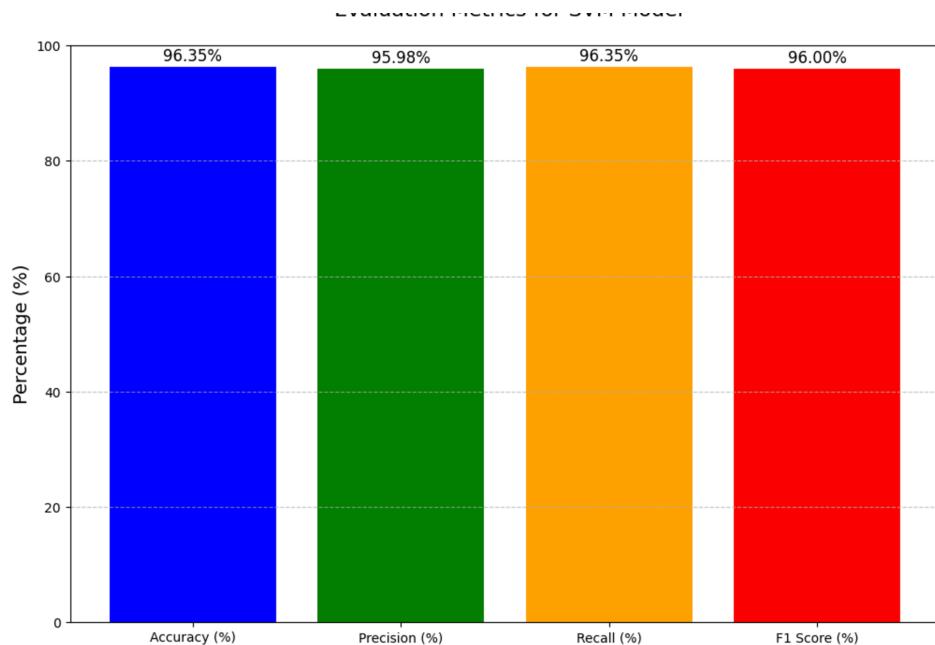
# Compile the model
model.compile(optimizer='adam', loss='categorical_crossentropy', metrics=['accuracy'])

# Train the model
history = model.fit(
    X_train_lstm, y_train_cat,
    epochs=10, # Adjust the epochs as needed
    batch_size=64,
    validation_data=(X_test_lstm, y_test_cat)
)
```

```

]: results = {
    "Accuracy (%)": accuracy,
    "Precision (%)": precision,
    "Recall (%)": recall,
    "F1 Score (%)": f1
}
print(results)
{'Accuracy (%)': 96.35153353555779, 'Precision (%)': 95.97671223840118, 'Recall (%)': 96.35153353555779, 'F1 Score (%)': 95.9982270634602
8}

```



Label Distribution: The matrix shows predictions across labels 0-17 (with some numbers missing in the sequence). The model appears to heavily favor predicting class 16, followed by class 14

Major Patterns: Class 14 shows mixed performance: 1586 correct predictions (class 14), but 670 instances misclassified as class 16. Class 16 has 1283 correct predictions, with 122 instances misclassified as class 14

Most other classes (0-10) are consistently misclassified as class 16, though with relatively small numbers (1-9 instances each)

Model Behavior: The model makes predictions primarily in classes 14 and 16. There's a clear binary confusion between classes 14 and 16. Many classes receive zero predictions, indicating the model isn't utilizing its full output space

8. GRADIENT BOOSTING

```
def build_and_train_model_with_metrics(data, target_column):
    data, target_encoder = preprocess_data(data, target_column)

    X = data.drop(columns=[target_column])
    y = data[target_column]

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

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

    model = GradientBoostingClassifier(random_state=42, n_estimators=50)
    model.fit(X_train, y_train)

    y_pred = model.predict(X_test)
    accuracy = accuracy_score(y_test, y_pred)
    cm = confusion_matrix(y_test, y_pred)

    print(f"Performance on {target_column}:")
    print(classification_report(y_test, y_pred))
    print(f"Accuracy: {accuracy}")

    plt.figure(figsize=(25, 17))
    sns.heatmap(cm, annot=True, fmt="d", cmap="Blues", xticklabels=target_encoder.classes_, yticklabels=target_encoder.classes_)
    plt.title(f"Confusion Matrix for {target_column}")
    plt.xlabel("Predicted")
    plt.ylabel("Actual")
    plt.show()

    return model, target_encoder
```

The code implements a Gradient Boosting Classifier to predict the target variable UCR_GROUP, which classifies crimes into predefined categories based on the Uniform Crime Reporting (UCR) system. This classification is vital for organizing and analyzing crime data, allowing law enforcement and policymakers to better understand and respond to various types of criminal activities.

The function `build_and_train_model_with_metrics` first preprocesses the dataset by separating features (X) and the target variable (y), in this case, UCR_GROUP. The data is split into training and testing sets (80-20 split) to train the model and evaluate its

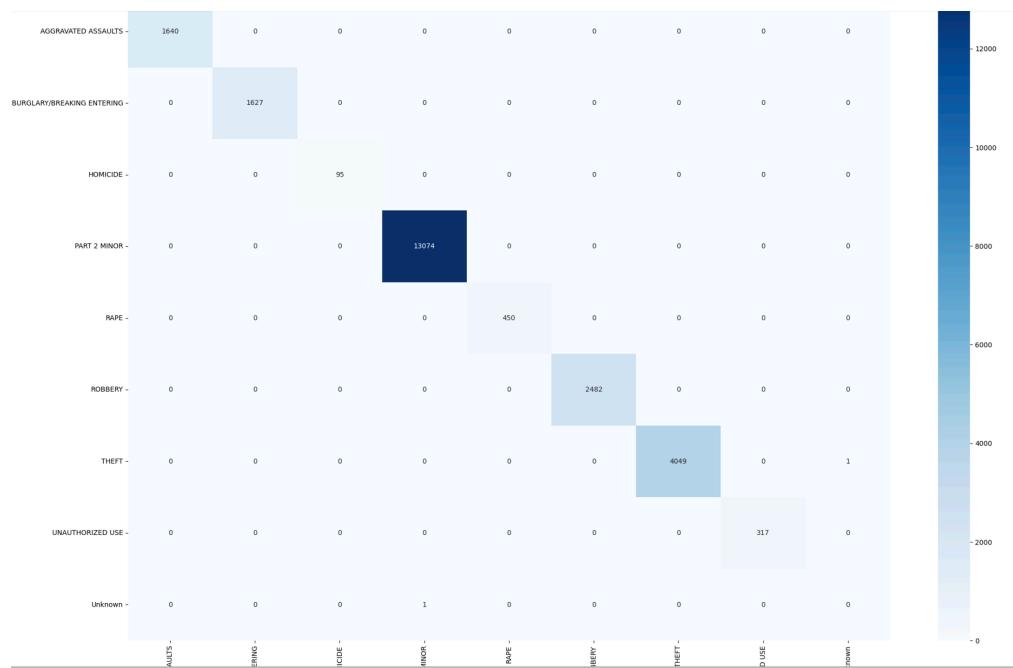
performance. Standard scaling is applied to the features to normalize them, ensuring consistent input ranges for the model.

The Gradient Boosting Classifier, configured with 50 estimators and a fixed random_state=42, is trained on the scaled training data. Gradient Boosting operates by building an ensemble of weak learners (decision trees) iteratively, improving classification performance with each step by focusing on misclassified samples.

Predicting UCR_GROUP...

Performance on UCR_GROUP:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	1640
1	1.00	1.00	1.00	1627
2	1.00	1.00	1.00	95
3	1.00	1.00	1.00	13074
4	1.00	1.00	1.00	450
5	1.00	1.00	1.00	2482
6	1.00	1.00	1.00	4050
7	1.00	1.00	1.00	317
8	0.00	0.00	0.00	1
accuracy			1.00	23736
macro avg	0.89	0.89	0.89	23736
weighted avg	1.00	1.00	1.00	23736



After training, the model makes predictions on the test set. The predicted results are compared with the actual labels using metrics such as accuracy, macro-average precision, recall, and F1-score. The confusion matrix, visualized as a heatmap, provides a detailed view of true positives, false positives, and false negatives for each crime category. The results demonstrate an accuracy of 99.99%, with a weighted average of 1.00 for precision, recall, and F1-score, indicating the model is almost perfect in its predictions.

Predicting UCR_GROUP using the Gradient Boosting Classifier has substantial implications. It automates the classification of crimes, saving time and reducing human error. The high accuracy of the model ensures reliable categorization, which is crucial for identifying trends, allocating resources, and devising crime prevention strategies. However, the near-perfect accuracy might indicate potential overfitting, suggesting that further evaluation on unseen data is necessary to confirm the model's generalizability. Nonetheless, this model provides a robust framework for enhancing data-driven decision-making in crime analysis and reporting.