

```
In [1]: import pandas as pd
import sklearn.preprocessing
import StandardScaler, MinMaxScaler
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from prophet import Prophet
import warnings
warnings.filterwarnings('ignore')
```

1. Data Acquisition: Download the dataset from the provided link and load it into your preferred data analysis tool

```
In [2]: #Task 1 - Download the dataset from the provided link and load it into your preferred data analysis tool
df_crime = pd.read_csv('Crime_Data_from_2020_to_Present.csv')
```

2. Data Inspection:

```
In [3]: #Task 2.1 - Display the first few rows of the dataset
df_crime.head()
```

	DR_NO	Date Rptd	DATE OCC	TIME OCC	AREA	AREA NAME	Rpt Dist No	Part 1-2	Crm Cd	Crm Cd Desc	...	Status	Status Desc	Crm Cd 1	Crm Cd 2	Crm Cd 3	Crm Cd 4	LOCATION	Cross Street
0	190326475	03/01/2020	03/01/2020	21:30	7	Wilshire	784	1	510	VEHICLE - STOLEN	...	AA	Adult Arrest	510.0	998.0	NaN	NaN	1900 S LONGWOOD AV	NaN
1	200106753	02/09/2020	02/08/2020	18:00	1	Central	182	1	330	BURGLARY FROM VEHICLE	...	IC	Invest Cont	330.0	998.0	NaN	NaN	1000 S FLOWER ST	NaN
2	200320258	11/10/2020	11/04/2020	17:00	3	Southwest	356	1	480	BIKE - STOLEN	...	IC	Invest Cont	480.0	NaN	NaN	NaN	1400 W 37TH ST	NaN
3	200907217	05/10/2023	03/10/2020	20:37	9	Van Nuys	964	1	343	SHOPLIFTING-GRAND THEFT (\$950.01 & OVER)	...	IC	Invest Cont	343.0	NaN	NaN	NaN	14000 RIVERSIDE DR	NaN
4	220614831	08/18/2022	08/17/2020	12:00	6	Hollywood	666	2	354	THEFT OF IDENTITY	...	IC	Invest Cont	354.0	NaN	NaN	NaN	1900 TRANSIENT	NaN

5 rows x 28 columns

```
In [4]: #Task 2.2 - Check the data types of each column
df_crime.dtypes
```

	DR_NO	Date Rptd	DATE OCC	TIME OCC	AREA	AREA NAME	Rpt Dist No	Part 1-2	Crm Cd	Crm Cd Desc	...	Status	Status Desc	Crm Cd 1	Crm Cd 2	Crm Cd 3	Crm Cd 4	LOCATION	Cross Street
count	9.82638e+05	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000	982638.000000
mean	1.2943964e+07	1338.954526	10.700277	1116.459887	1.404253	500.823555	29.079817	306.133008	363.840863	500.876661
min	8.170009e+02	1.000000	1.000000	101.000000	1.000000	110.000000	-4.000000	101.000000	101.000000	110.000000
25%	2.160609e+08	900.000000	5.000000	587.000000	1.000000	331.000000	0.000000	101.000000	311.000000	331.000000
50%	2.220816e+08	1420.000000	11.000000	1141.000000	1.000000	442.000000	30.000000	203.000000	442.000000	442.000000
75%	2.359153e+08	1900.000000	16.000000	1617.000000	2.000000	626.000000	44.000000	501.000000	400.000000	626.000000
max	2.499236e+08	2359.000000	21.000000	2199.000000	2.000000	956.000000	120.000000	976.000000	516.000000	956.000000

3.Data Cleaning:

```
In [6]: #Task 3.1 - Identify and handle missing data appropriately
df_crime.isnull().sum()
```

```
df_crime['Vict Sex'] = df_crime['Vict Sex'].replace({' ': 'X', '-': 'X', 'H': 'X'})
missing_values = ["n.a.", "NA", "n/a", "na", ""]
df_crime = pd.read_csv('Crime_Data_from_2020_to_Present.csv', na_values = missing_values)
df_crime.dropna(axis=0, inplace = True, how = 'all')
```

```
In [7]: #Task 3.2 - Check for and remove duplicate rows
df_crime = df_crime.drop_duplicates()
print(df_crime.duplicated().sum())
```

0

```
In [8]: #Task 3.3 - Convert data types if needed (e.g., dates to date format, numerical values to appropriate numeric types)
df_crime['Date Rptd'] = pd.to_datetime(df_crime['Date Rptd']).dt.date
df_crime['DATE OCC'] = pd.to_datetime(df_crime['DATE OCC']).dt.date
df_crime['AREA'] = df_crime['AREA'].astype(str)
df_crime['Rpt Dist No'] = df_crime['Rpt Dist No'].astype(int)
```

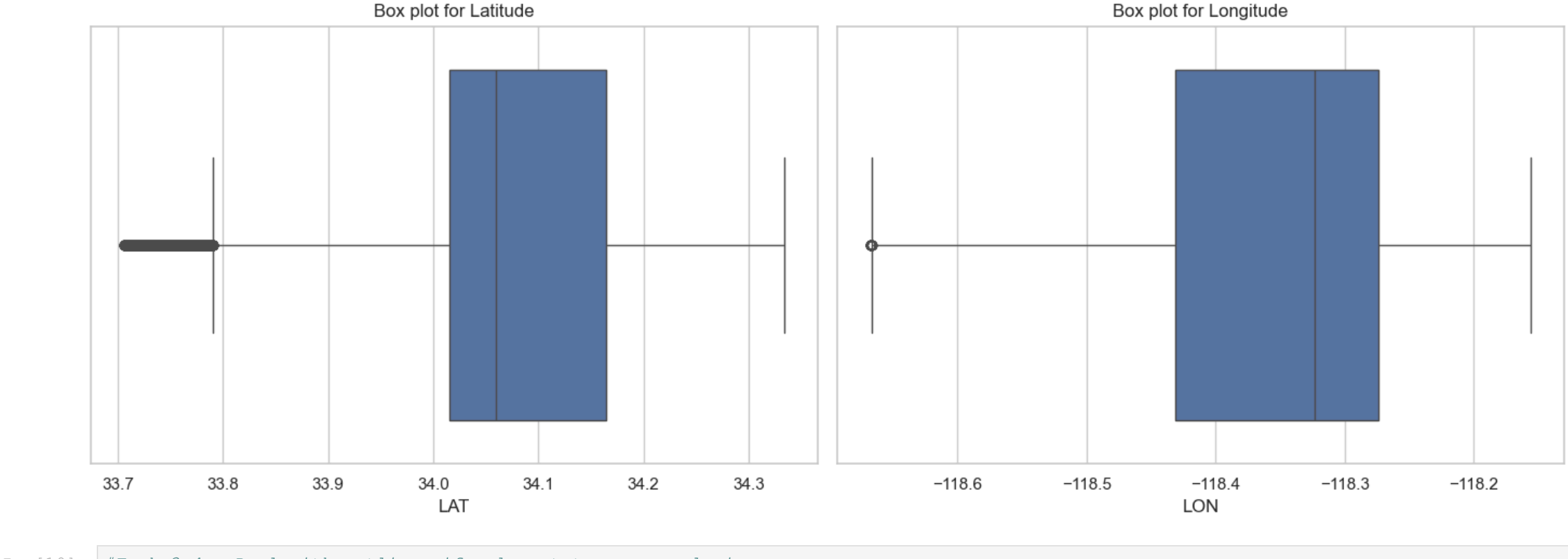
```
def convert_time(time):
    time = str(time).split(':')
    return f'{time[2]}:{time[1]}'
df_crime['TIME OCC'] = df_crime['TIME OCC'].apply(convert_time)
```

```
df_crime
```

	DR_NO	Date Rptd	DATE OCC	TIME OCC	AREA	AREA NAME	Rpt Dist No	Part 1-2	Crm Cd	Crm Cd Desc	...	Status	Status Desc	Crm Cd 1	Crm Cd 2	Crm Cd 3	Crm Cd 4	LOCATION	Cross Street
0	190326475	2020-03-01	2020-03-01	21:30	7	Wilshire	784	1	510	VEHICLE - STOLEN	...	AA	Adult Arrest	510.0	998.0	NaN	NaN	1900 S LONGWOOD AV	NaN
1	200106753	2020-02-09	2020-02-08	18:00	1	Central	182	1	330	BURGLARY FROM VEHICLE	...	IC	Invest Cont	330.0	998.0	NaN	NaN	1000 S FLOWER ST	NaN
2	200320258	2020-11-10	2020-11-04	17:00	3	Southwest	356	1	480	BIKE - STOLEN	...	IC	Invest Cont	480.0	NaN	NaN	NaN	1400 W 37TH ST	NaN
3	200907217	2023-05-10	2020-03-10	20:37	9	Van Nuys	964	1	343	SHOPLIFTING-GRAND THEFT (\$950.01 & OVER)	...	IC	Invest Cont	343.0	NaN	NaN	NaN	14000 RIVERSIDE DR	NaN
4	220614831	2022-08-18	2020-08-17	12:00	6	Hollywood	666	2	354	THEFT OF IDENTITY	...	IC	Invest Cont	354.0	NaN	NaN	NaN	1900 TRANSIENT	NaN
...
982633	24011172	2024-04-24	2024-08-20	23:00	20	Olympic	2033	1	341	THEFT-GRAND (\$950.01 & OVER)(EXCEPT GUNS,FOWLL...	...	IC	Invest Cont	341.0	NaN	NaN	NaN	WILSHIRE	NaN
982634	240710284	2024-07-24	2024-07-23	14:00	7	Wilshire	788	1	510	VEHICLE - STOLEN	...	IC	Invest Cont	510.0	NaN	NaN	NaN	400 23RD	NaN
982635	240104953	2024-01-15	2024-01-15	01:00	1	Central	101	2	745	MISDEMEANOR (\$395 OR UNDER)	...	IC	Invest Cont	745.0	NaN	NaN	NaN	13X SUNSET	NaN
982636	240309674	2024-04-24	2024-04-24	15:00	3	Southwest	358	1	230	ASSAULT WITH DEADLY WEAPON, AGGRAVATED ASSAULT	...	IC	Invest Cont	230.0	NaN	NaN	NaN	FLOWER	NaN
982637	240910892	2024-08-13	2024-08-12	23:00	9	Van Nuys	914	1	510	VEHICLE - STOLEN	...	IC	Invest Cont	510.0	NaN	NaN	NaN	VESPER	NaN

982638 rows x 28 columns

```
In [9]: #Task 3.4 Deal with outliers if relevant to your analysis.
sns.set(style='whitegrid')
fig, axes = plt.subplots(nrows=2, ncols=2, figsize=(14, 10))
sns.boxplot(x=df_crime['Vict Age'], axes=axes[0, 0])
sns.boxplot(x=df_crime['DR_NO'], axes=axes[0, 1])
sns.boxplot(x=df_crime['LAT'], axes=axes[1, 0])
sns.boxplot(x=df_crime['LON'], axes=axes[1, 1])
plt.tight_layout()
plt.show()
```



```
In [10]: #Task 3.4 - Deal with outliers if relevant to your analysis
# Replace 'Vict Age' less than or equal to 0 with NaN values to remove outliers
df_crime['Vict Age'] = np.where(df_crime['Vict Age'] <= 0, np.nan, df_crime['Vict Age'])
df_crime
```

```
df_crime
```

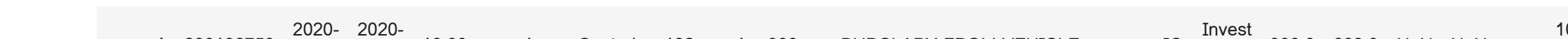
```
In [11]: #Task 3.5 - Standardize or normalize numerical data if necessary.
num_col = ['Vict Age', 'Rpt Dist No', 'Part 1-2', 'Premis Cd', 'LAT', 'LON']
# Standardization
scaler = StandardScaler()
df_crime[num_col] = scaler.fit_transform(df_crime[num_col])
# Normalization
scaler = MinMaxScaler()
df_crime[num_col] = scaler.fit_transform(df_crime[num_col])
df_crime
```

```
df_crime
```

```
In [12]: #Task 3.6 - Encode categorical data if required.
# There is no categorical data that required encoding.
```

4.Exploratory Data Analysis (EDA)

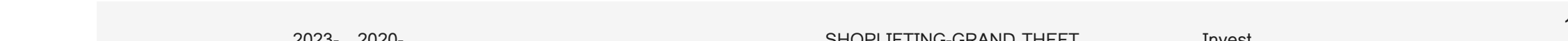
```
In [13]: #Task 4.1 - Visualize overall crime trends from 2020 to the present year
df_crime['year'] = pd.to_datetime(df_crime['DATE OCC']).dt.year
crime_per_year = df_crime.groupby('year').size()
```



```
In [14]: #Task 4.2 - Analyze and visualize seasonal patterns in crime data.
df_crime['month'] = pd.to_datetime(df_crime['DATE OCC']).dt.month
crime_by_month = df_crime.groupby('month').size()
```



```
In [15]: #Task 4.3 - Identify the most common type of crime and its trends over time.
most_common_crime = df_crime['Crm Cd Desc'].value_counts().head(10)
print(most_common_crime)
```

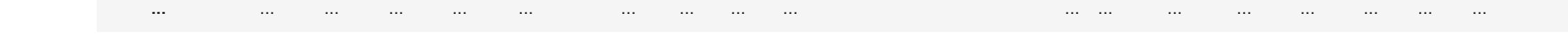


```
In [16]: #Task 4.3 - Identify the most common type of crime and its trends over time.
crime_type_counts = df_crime['Crm Cd Desc'].value_counts().head(10)
most_common_crime = crime_type_counts.idxmax()
print(f'Most common crime type: {most_common_crime}')
most_common_crime_data = df_crime[df_crime['Crm Cd Desc'] == most_common_crime]
```

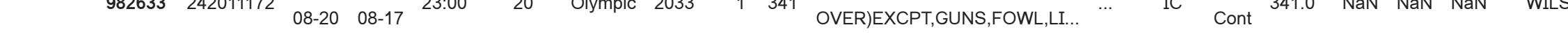
```
df_crime['year'] = pd.to_datetime(df_crime['DATE OCC']).dt.year
crime_trend_by_year = most_common_crime_data.groupby('year').size()
```



```
In [17]: #Task 4.4 - Investigate if there are any notable differences in crime rates between regions or cities
crime_by_area = df_crime.groupby('AREA NAME').size().reset_index(name='Count').sort_values('Count', ascending = False)
```

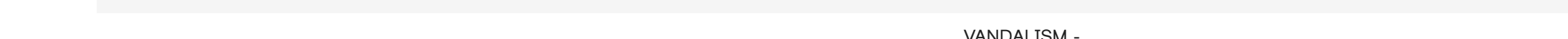


```
In [18]: #Task 4.4 - Investigate if there are any notable differences in crime rates between regions or cities
crime_by_loc = df_crime.groupby('LOCATION').size().reset_index(name='Count').sort_values('Count', ascending = False).head(10)
```



```
In [19]: #Task 4.5 - Explore correlations between economic factors (if available) and crime rates
# Cannot be visualized since economic data is not available.
```

```
In [20]: #Task 4.6 - Analyze the relationship between the day of the week and the frequency of certain types of crimes.
df_crime['Day of Week'] = pd.to_datetime(df_crime['DATE OCC']).dt.day_name()
top5_crime_data = df_crime[df_crime['Crm Cd Desc'].isin(top5_crimes)]
top5_crime_by_day_pivot = top5_crime_data.pivot(index='Day of Week', columns='Crm Cd Desc', values='Count')
```



```
In [21]: #Task 5 - a) Advanced Analysis
# Use predictive modeling techniques (e.g., time series forecasting) to predict future crime trends
df_crime['DATE OCC'] = pd.to_datetime(df_crime['DATE OCC'])
daily_crimes = df_crime.resample('D', on='DATE OCC').size().reset_index(name='Count')
```

