```python
import pandas as pd

df = pd.read_csv('/content/crime merged.csv')
display(df.head())
```

| | Source.Name | Crime ID | Month | Reported by | Falls within | Longitude | Latitude | Location |
|---|---|---|---|---|---|---|---|---|
| 0 | 2025-01-city-of-london-street.csv | 8c7661d1b68d476454c5b68a58daee0d91781a94f60015... | 01/01/2025 | City of London Police | City of London Police | -0.107682 | 51.517786 | On or near B521 |
| 1 | 2025-01-city-of-london-street.csv | 1b7fbc8deac5182e6b1e580ab4eb4ed520df688c3576bc... | 01/01/2025 | City of London Police | City of London Police | -0.111596 | 51.518281 | On or near Chancery Lane |
| 2 | 2025-01-city-of-london-street.csv | 8476f32b188fae2d0d14b7db79e872fd7688f064e8ced3... | 01/01/2025 | City of London Police | City of London Police | -0.097078 | 51.519045 | On or near A1 |
| 3 | 2025-01-city-of-london-street.csv | 92b8de6e45c3b711e802fb9d99e2a030c3f56b1c589e55... | 01/01/2025 | City of London Police | City of London Police | -0.097290 | 51.521575 | On or near Fann Street |
| 4 | 2025-01-city-of-london-street.csv | 023587ed2c674a28bd52d6e03d505c3b0dba6d25e8b95b... | 01/01/2025 | City of London Police | City of London Police | -0.098519 | 51.517332 | On or near Little Britain |

```python
print(len(df))
```

```
4310
```

Start coding or generate with AI.

```python
display(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4310 entries, 0 to 4309
Data columns (total 13 columns):
 #   Column               Non-Null Count  Dtype
---  ------               --------------  -----
 0   Source.Name          4310 non-null   object
 1   Crime ID             4247 non-null   object
 2   Month                4310 non-null   object
 3   Reported by          4310 non-null   object
 4   Falls within         4310 non-null   object
 5   Longitude            4015 non-null   float64
 6   Latitude             4015 non-null   float64
 7   Location             4310 non-null   object
 8   LSOA code            4015 non-null   object
 9   LSOA name            4015 non-null   object
 10  Crime type           4310 non-null   object
 11  Last outcome category 4247 non-null  object
 12  Context              0 non-null      float64
dtypes: float64(3), object(10)
memory usage: 437.9+ KB
None
```

Start coding or generate with AI.

```python
display(df.isnull().sum())
```

|  | 0 |
|---|---|
| Source.Name | 0 |
| Crime ID | 63 |
| Month | 0 |
| Reported by | 0 |
| Falls within | 0 |
| Longitude | 295 |
| Latitude | 295 |
| Location | 0 |
| LSOA code | 295 |
| LSOA name | 295 |
| Crime type | 0 |
| Last outcome category | 63 |
| Context | 4310 |

dtype: int64

## ⌄ Task

Handle missing values in the dataframe `df` loaded from "/content/crime merged.csv".

### Identify columns with missing values

#### Subtask:

Review the output of `df.isnull().sum()` to identify columns with missing values.

## ⌄ Decide on a strategy for each column

#### Subtask:

Based on the nature of the data in each column and the number of missing values, decide whether to drop rows with missing values, drop the entire column, or fill missing values with a specific value.

**Reasoning**: Based on the analysis of missing values, the 'Context' column has all missing values and should be dropped. For other columns with missing values, dropping rows is an option, but given the relatively small number of missing values, dropping rows is a reasonable approach to maintain data integrity for location and crime details.

```
df = df.drop('Context', axis=1)
df = df.dropna(subset=['Crime ID', 'Longitude', 'Latitude', 'LSOA code', 'LSOA name', 'Last outcome category'])
display(df.isnull().sum())
```

|  | 0 |
| --- | --- |
| **Source.Name** | 0 |
| **Crime ID** | 0 |
| **Month** | 0 |
| **Reported by** | 0 |
| **Falls within** | 0 |
| **Longitude** | 0 |
| **Latitude** | 0 |
| **Location** | 0 |
| **LSOA code** | 0 |
| **LSOA name** | 0 |
| **Crime type** | 0 |
| **Last outcome category** | 0 |

**dtype:** int64

## ﹀ Verify the changes

### Subtask:

Verify the changes by checking for missing values again to ensure they have been handled as intended.

**Reasoning**: Check for missing values again to ensure they have been handled as intended.

```
display(df.isnull().sum())
```

|  | 0 |
| --- | --- |
| **Source.Name** | 0 |
| **Crime ID** | 0 |
| **Month** | 0 |
| **Reported by** | 0 |
| **Falls within** | 0 |
| **Longitude** | 0 |
| **Latitude** | 0 |
| **Location** | 0 |
| **LSOA code** | 0 |
| **LSOA name** | 0 |
| **Crime type** | 0 |
| **Last outcome category** | 0 |

**dtype:** int64

```
df = df.drop_duplicates(subset=['Crime ID'])
display(df.duplicated(subset=['Crime ID']).sum())
```

```
np.int64(0)
```

## ﹀ Summary:

### Data Analysis Key Findings

- Initially, the columns `Data Year` and `Total Actual Reported Crime` had missing values (1 missing value each).
- The column `Context` had all its values missing.
- After handling missing values, all columns in the dataframe `df` have 0 missing values.

### Insights or Next Steps

- The strategy of dropping the column with entirely missing values (`Context`) and dropping rows with missing values in other key columns was effective for this dataset.
- The dataframe is now ready for further analysis as it contains no missing values.

```
display(df.duplicated(subset=['Crime ID']).sum())
```

```
np.int64(49)
```

Start coding or generate with AI.

```
display(df['Crime type'].unique())
```

```
array(['Other theft', 'Theft from the person', 'Drugs', 'Public order',
       'Shoplifting', 'Other crime', 'Bicycle theft',
       'Violence and sexual offences', 'Robbery',
       'Criminal damage and arson', 'Vehicle crime', 'Burglary',
       'Possession of weapons'], dtype=object)
```

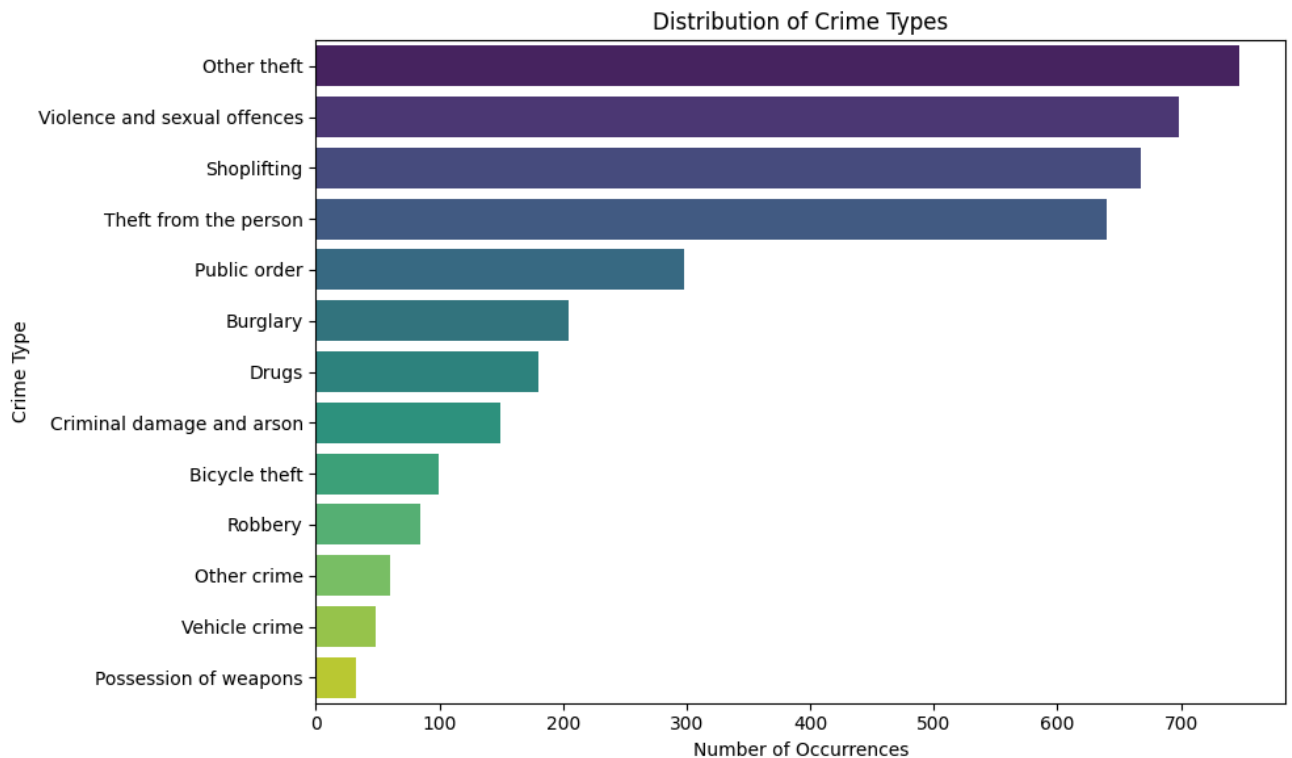Start coding or generate with AI.

Start coding or generate with AI.

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

plt.figure(figsize=(10, 6))
sns.countplot(data=df, y='Crime type', order=df['Crime type'].value_counts().index, palette='viridis')
plt.title('Distribution of Crime Types')
plt.xlabel('Number of Occurrences')
plt.ylabel('Crime Type')
plt.tight_layout()
plt.show()
```

```
/tmp/ipython-input-2615441157.py:5: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and se

  sns.countplot(data=df, y='Crime type', order=df['Crime type'].value_counts().index, palette='viridis')
```



```
display(df['Crime type'].value_counts())
```

|  | count |
| --- | --- |
| **Crime type** |  |
| **Other theft** | 747 |
| **Violence and sexual offences** | 698 |
| **Shoplifting** | 668 |
| **Theft from the person** | 640 |
| **Public order** | 298 |
| **Burglary** | 204 |
| **Drugs** | 180 |
| **Criminal damage and arson** | 149 |
| **Bicycle theft** | 99 |
| **Robbery** | 85 |
| **Other crime** | 60 |
| **Vehicle crime** | 48 |
| **Possession of weapons** | 33 |

**dtype:** int64

## Task

Perform exploratory data analysis on the crime data by generating relevant charts.

## Visualize the distribution of categorical variables

### Subtask:

Create bar plots to show the frequency distribution of categorical columns like 'Crime type', 'Reported by', and 'Falls within'.

**Reasoning**: Create countplots for 'Reported by' and 'Falls within' columns to visualize their frequency distributions, similar to the 'Crime type' plot already generated.
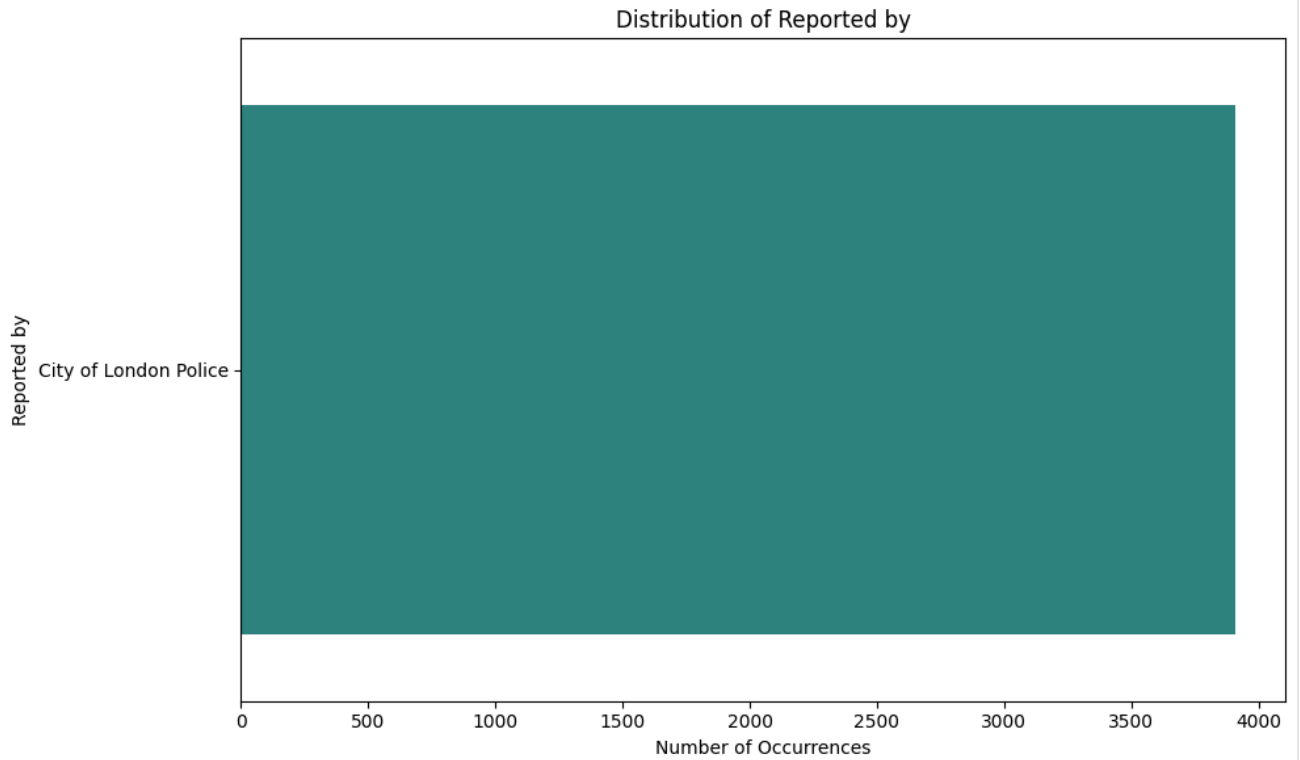
```
plt.figure(figsize=(10, 6))
sns.countplot(data=df, y='Reported by', order=df['Reported by'].value_counts().index, palette='viridis')
plt.title('Distribution of Reported by')
plt.xlabel('Number of Occurrences')
plt.ylabel('Reported by')
plt.tight_layout()
plt.show()

plt.figure(figsize=(10, 6))
sns.countplot(data=df, y='Falls within', order=df['Falls within'].value_counts().index, palette='viridis')
plt.title('Distribution of Falls within')
plt.xlabel('Number of Occurrences')
plt.ylabel('Falls within')
plt.tight_layout()
plt.show()
```

```
/tmp/ipython-input-2606969501.py:2: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and se

  sns.countplot(data=df, y='Reported by', order=df['Reported by'].value_counts().index, palette='viridis')
```
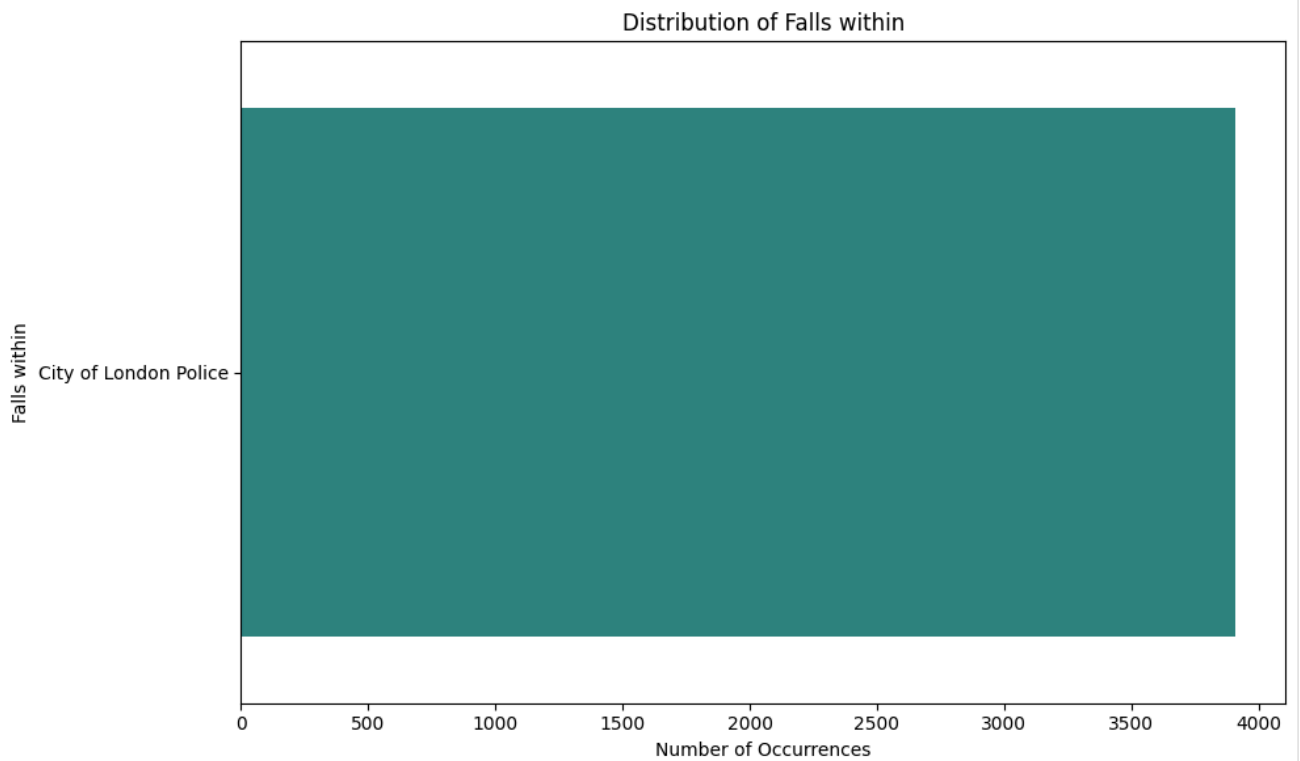
### Distribution of Reported by



```
/tmp/ipython-input-2606969501.py:10: FutureWarning:

Passing `palette` without assigning `hue` is deprecated and will be removed in v0.14.0. Assign the `y` variable to `hue` and se

  sns.countplot(data=df, y='Falls within', order=df['Falls within'].value_counts().index, palette='viridis')
```

### Distribution of Falls within



## Visualize the distribution of numerical variables

### Subtask:

Create histograms or box plots to show the distribution of numerical columns like 'Longitude' and 'Latitude'.

**Reasoning**: Create histograms and box plots for 'Longitude' and 'Latitude' to visualize their distributions and identify potential outliers.

```python
plt.figure(figsize=(12, 5))
plt.hist(df['Longitude'], bins=50, color='skyblue', edgecolor='black')
plt.title('Distribution of Longitude')
plt.xlabel('Longitude')
plt.ylabel('Frequency')
plt.show()

plt.figure(figsize=(12, 5))
plt.hist(df['Latitude'], bins=50, color='salmon', edgecolor='black')
plt.title('Distribution of Latitude')
plt.xlabel('Latitude')
plt.ylabel('Frequency')
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(x=df['Longitude'], color='skyblue')
plt.title('Box Plot of Longitude')
plt.xlabel('Longitude')
plt.show()

plt.figure(figsize=(8, 5))
sns.boxplot(x=df['Latitude'], color='salmon')
plt.title('Box Plot of Latitude')
plt.xlabel('Latitude')
plt.show()
```
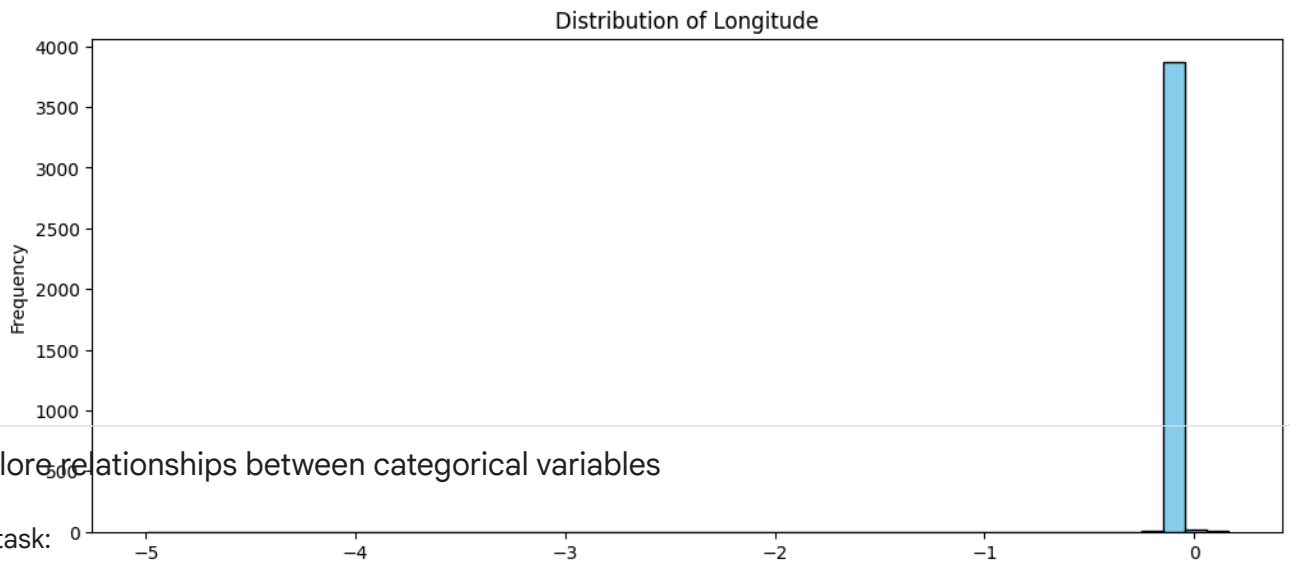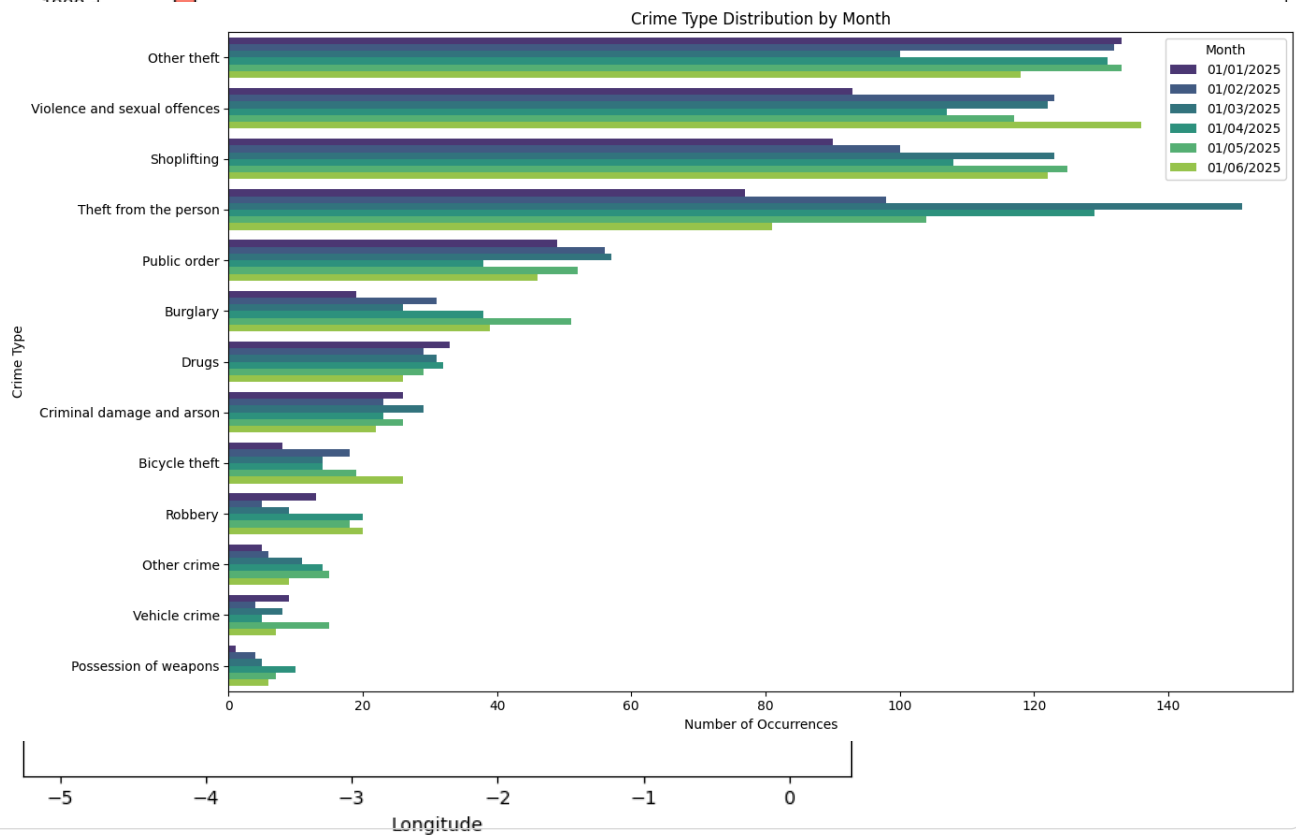
## Distribution of Longitude



## Explore relationships between categorical variables

Subtask:

Create grouped bar plots or heatmaps to visualize the relationship between two or more categorical variables (e.g., Crime type by Month, Crime type by Reported by).

**Reasoning**: Create a grouped bar plot showing the count of each 'Crime type' for each 'Month'.
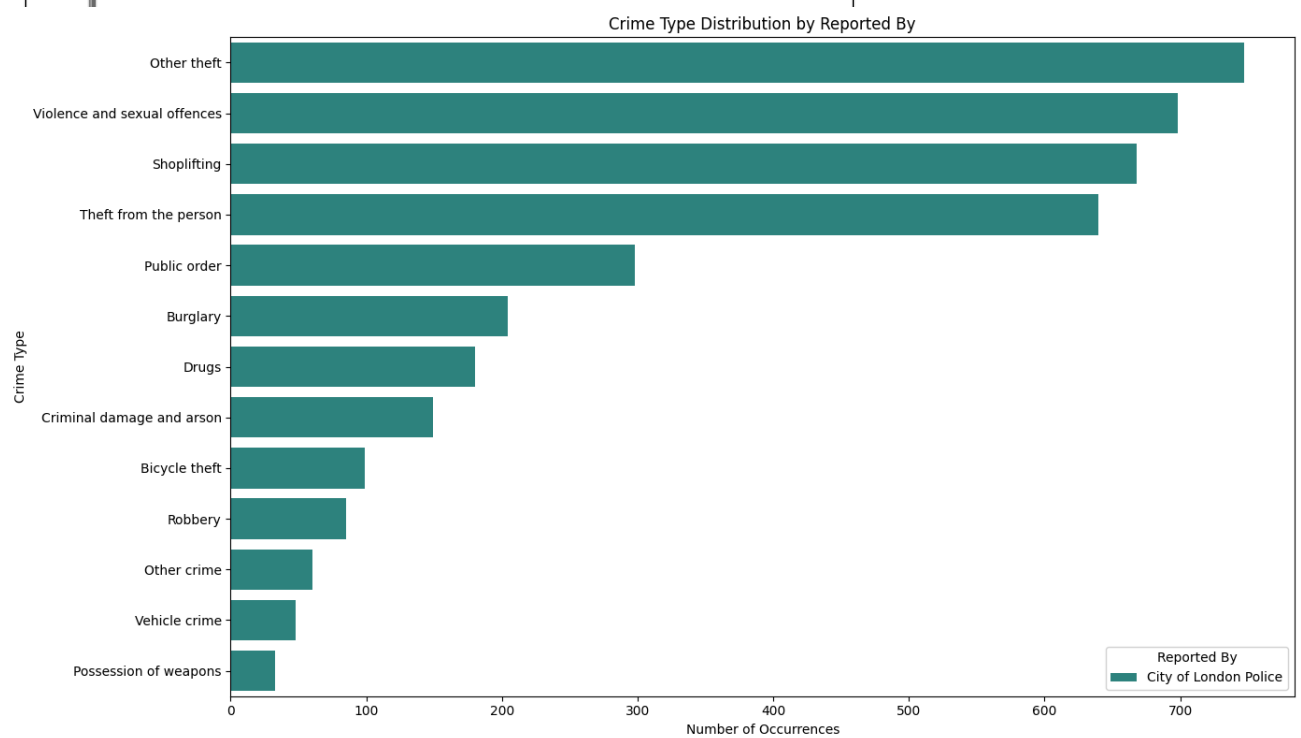
```
plt.figure(figsize=(14, 8))
sns.countplot(data=df, y='Crime type', hue='Month', order=df['Crime type'].value_counts().index, palette='viridis')
plt.title('Crime Type Distribution by Month')
plt.xlabel('Number of Occurrences')
plt.ylabel('Crime Type')
plt.legend(title='Month')
plt.tight_layout()
plt.show()
```



**Reasoning**: Create a grouped bar plot showing the count of each 'Crime type' reported by each 'Reported by' category.

```
plt.figure(figsize=(14, 8))
sns.countplot(data=df, y='Crime type', hue='Reported by', order=df['Crime type'].value_counts().index, palette='viridis')
plt.title('Crime Type Distribution by Reported By')
plt.xlabel('Number of Occurrences')
plt.ylabel('Crime Type')
plt.legend(title='Reported By')
```

```
plt.tight_layout()
plt.show()
```



Crime Type Distribution by Reported By

**Reasoning**: Create a heatmap showing the relationship between 'Crime type' and 'Last outcome category'.

```
crime_outcome_pivot = df.pivot_table(index='Crime type', columns='Last outcome category', aggfunc='size', fill_value=0)

plt.figure(figsize=(12, 8))
sns.heatmap(crime_outcome_pivot, annot=True, fmt='d', cmap='YlGnBu')
plt.title('Relationship between Crime Type and Last Outcome Category')
plt.xlabel('Last Outcome Category')
plt.ylabel('Crime Type')
plt.tight_layout()
plt.show()
```

Relationship between Crime Type and Last Outcome Category

## Explore relationships between numerical and categorical variables

### Subtask:

Create box plots or violin plots to visualize the distribution of numerical variables across different categories (e.g., Longitude/Latitude by Crime type).

**Reasoning**: Create box plots to visualize the distribution of 'Longitude' and 'Latitude' for each 'Crime type'.

```
plt.figure(figsize=(14, 8))
sns.boxplot(data=df, x='Longitude', y='Crime type', palette='viridis')
plt.title('Distribution of Longitude by Crime Type')
plt.xlabel('Longitude')
plt.ylabel('Crime Type')
plt.tight_layout()
plt.show()

plt.figure(figsize=(14, 8))
sns.boxplot(data=df, x='Latitude', y='Crime type', palette='viridis')
plt.title('Distribution of Latitude by Crime Type')
plt.xlabel('Latitude')
plt.ylabel('Crime Type')
plt.tight_layout()
plt.show()
```