

Introduction

This exploratory data analysis (EDA) investigates a dataset named "**Crime1.csv**", focusing on crime incidents. The dataset includes several categorical and numerical variables, with key attributes being **Category** (type of crime), **Description**, **DayOfWeek**, and **PdDistrict** (police district). The goal is to uncover patterns in the types, frequency, and distribution of crimes across time and geography. The analysis uses visual tools such as bar plots, count plots, and word clouds to interpret crime trends. The study uses libraries like Pandas, Seaborn, and WordCloud in Python to manipulate and visualize the data for meaningful insights.

Findings and Interpretation

The most frequent crime categories are identified using value counts, where crimes like **Larceny/Theft**, **Other Offenses**, and **Non-Criminal** dominate the dataset. These categories are visualized using a **horizontal bar plot** and a **word cloud**, offering both quantitative and visual frequency representation. Mathematically, let C_i denote the count of each crime type i , and the frequency distribution is given by $f(C_i) = \frac{C_i}{\sum C_i}$. The frequency plots exhibit a **right-skewed distribution**, meaning a few categories account for most crimes, following the **Pareto Principle** (80/20 rule).

Further analysis on the **Description** field mirrors this distribution. A word cloud constructed from unique crime descriptions illustrates frequent terms like "BATTERY", "ASSAULT", and "PETTY THEFT", indicating specific behaviors or subtypes under broader categories. In addition, the analysis explores the temporal distribution of crimes using the **DayOfWeek** variable. A bar plot visualizing counts per day shows a relatively even distribution with a slight increase mid-week. If D_j represents the crime count on day j , then the expectation $E[D]$ across the week suggests marginal deviation from uniform distribution, hinting at constant policing requirements throughout the week.

The spatial distribution is assessed using **PdDistrict** value counts. Central and dense districts such as **SOUTHERN**, **MISSION**, and **NORTHERN** have higher crime incidences. This is expected as urban centers generally have higher population density and socioeconomic diversity, increasing the chance of crime occurrences. These insights can be visualized through a bar chart where the y-axis is the district and the x-axis represents count, with peaks indicating hotspot regions.

Conclusion

In conclusion, the exploratory analysis of the crime dataset reveals meaningful patterns across crime categories, timing, and geographic distribution. The study finds that crimes are unevenly distributed, with a few types being disproportionately frequent. Temporal patterns show relatively stable crime occurrence throughout the week, while spatial patterns reveal district-specific crime hotspots. These findings suggest a need for focused law enforcement in high-incidence districts and ongoing crime-type monitoring to refine prevention strategies. The visualizations and quantitative summaries provide a strong basis for deeper predictive modeling or policy development in urban safety and policing efforts.

