# Crime Hotspot Identification in Los Angeles Using Machine Learning

## 1. Introduction

Urban crime poses significant challenges to public safety, resource allocation, and community well-being. Large metropolitan cities such as Los Angeles exhibit complex crime patterns influenced by spatial, temporal, and socio-environmental factors. Traditional crime analysis methods often rely on manual inspection and historical summaries, which may fail to capture hidden spatial structures and emerging risk zones.

The objective of this project is to develop a **data-driven, machine-learning-based framework to identify crime hotspots in Los Angeles** using historical LAPD crime data (2020–2024). The model aims to support law enforcement agencies, security planners, and policymakers by providing actionable insights into high-risk areas.

---

## 2. Dataset Description

The analysis is based on the LAPD Crime Data dataset covering incidents reported between 2020 and 2024. The dataset contains over **900,000 crime records** with attributes describing:

- Crime location (latitude and longitude)
- Time and date of occurrence
- Crime type and description
- Crime severity classification (Part 1 – serious crimes, Part 2 – less serious crimes)
- Area and reporting district information

The dataset was treated as confidential and was not uploaded or shared publicly.

---

## 3. Data Preprocessing

To ensure data quality and analytical reliability, a comprehensive preprocessing pipeline was implemented:

### 3.1 Data Cleaning

- Invalid geographic coordinates (e.g., latitude/longitude equal to 0 or outside Los Angeles boundaries) were removed.
- Date and time fields were parsed into proper datetime formats.
- Duplicate records were eliminated.

**3.2 Handling Missing Values**

- Critical features required for modeling (location, time, severity) were retained only where valid.
- Non-critical categorical attributes were imputed with descriptive placeholders (e.g., "Unknown").
- Numerical attributes were imputed using robust statistical measures where appropriate.

This ensured the dataset was both **visualization-ready and model-ready**.

---

# 4. Exploratory Data Analysis (EDA)

Exploratory analysis was conducted to understand crime behavior before modeling.

## 4.1 Temporal Analysis

- Crime frequency peaks during **midday and evening hours**.
- Weekday crimes exceed weekend crimes, indicating the influence of routine urban activity.

## 4.2 Crime Type Analysis

- Property and vehicle-related crimes dominate the dataset.
- Violent crimes occur less frequently but pose higher public safety risk.

## 4.3 Spatial Analysis

- Crime incidents are **not uniformly distributed**.
- Clear geographic clusters are visible, particularly in central and densely populated regions.

**EDA Conclusion:** Crime in Los Angeles shows strong spatial concentration and time dependency, validating the need for hotspot-based modeling rather than uniform city-wide analysis.

---

# 5. Feature Engineering

Based on EDA insights, the following features were engineered and selected for modeling:

- **Latitude & Longitude:** Core spatial indicators for hotspot detection
- **Hour of Crime:** Captures time-of-day patterns
- **Day of Week / Weekend Indicator:** Reflects behavioral differences
- **Crime Severity (Part 1-2):** Enables risk-aware interpretation

Feature relevance was justified through exploratory analysis rather than algorithmic coefficients.

---

# 6. Modeling Approach

## 6.1 Learning Paradigm

This project follows an **unsupervised learning approach** because no labeled ground truth for crime hotspots exists.

## 6.2 Algorithm Selection

**DBSCAN (Density-Based Spatial Clustering of Applications with Noise)** was selected as the primary model due to its ability to: - Detect arbitrarily shaped clusters - Automatically identify noise (low-risk isolated incidents) - Avoid forcing all points into clusters

K-Means clustering was also explored as a comparative baseline.

---

# 7. DBSCAN Hotspot Detection

## 7.1 Model Training

DBSCAN was applied on scaled spatial features (latitude and longitude). Model parameters were iteratively tuned to balance hotspot granularity and interpretability.

## 7.2 Clustering Results

- A dominant large cluster representing continuous dense crime regions was identified.
- Multiple medium-sized clusters represented **localized, actionable hotspots**.
- A significant number of incidents were labeled as noise, corresponding to isolated crimes.

Small clusters were filtered to retain only **operationally meaningful hotspots**.

---

# 8. Model Evaluation

## 8.1 Accuracy Considerations

Traditional accuracy metrics are **not applicable** due to the unsupervised nature of the problem.

## 8.2 Evaluation Metrics Used

- **Silhouette Score (Sampled):**
- A sampled silhouette score of approximately **−0.3** was obtained.

- Negative values are expected for large, continuous spatial datasets with overlapping regions.

- **Spatial Coherence:**

• Visual inspection confirmed that identified clusters align with known high-crime areas.

• **Domain Validation:**

• Hotspots matched insights from EDA and known urban crime patterns.

**Evaluation Conclusion:** For density-based hotspot detection, spatial interpretability and domain consistency are more meaningful than silhouette optimization.

---

## 9. Severity Analysis

Analysis of crime severity within and outside hotspots revealed:

• Hotspots contain a high concentration of both serious and non-serious crimes.
• Noise regions still include a substantial proportion of serious crimes, indicating that not all high-severity incidents occur within persistent hotspots.

This highlights the importance of combining **hotspot-based prevention** with **city-wide rapid response strategies**.

---

## 10. Business and Security Implications

• Security resources should be geographically prioritized toward identified hotspots.
• Medium-sized hotspots offer the highest potential return for targeted interventions.
• Noise incidents require responsive policing rather than long-term preventive deployment.
• Data-driven hotspot identification improves collaboration between law enforcement and urban planners.

---

## 11. Limitations

• Lack of labeled ground truth prevents direct accuracy computation.
• DBSCAN parameter sensitivity may affect cluster granularity.
• The model identifies existing hotspots but does not forecast future crime.

---

## 12. Future Scope

• Incorporation of supervised learning using labeled hotspot data
• Severity-weighted clustering
• Temporal forecasting of hotspot evolution
• Integration of socio-economic and environmental data

---

## 13. Conclusion

This project demonstrates the effectiveness of **density-based unsupervised learning** for crime hotspot identification in Los Angeles. By combining rigorous preprocessing, exploratory analysis, and DBSCAN clustering, the model delivers interpretable and actionable insights into spatial crime concentration. The approach supports evidence-based decision-making for public safety initiatives and lays the foundation for future predictive crime analytics.