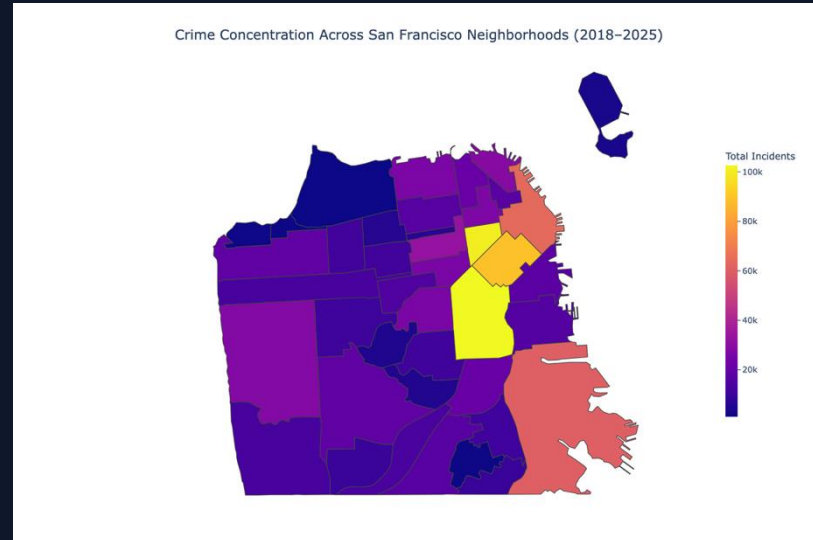# San Francisco Crime Trend Analysis and Forecasting (2018 - 2025)

From Data to Decision Support

Sileshi Hirpa | Data Analyst

# The Problem: Data Exists, Decisions Lag

- San Francisco publishes extensive crime data, but it is not decision-ready.

- Incident records are noisy, seasonal, and difficult to interpret.

- Short-term spikes are often mistaken for long-term change.

- Result: decisions become reactive rather than strategic.



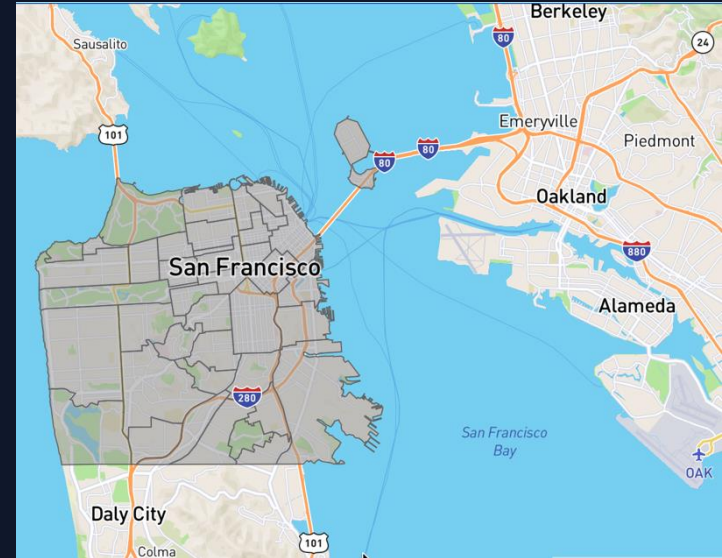Crime Concentration Across San Francisco Neighborhoods (2018–2025)

# Project Goal

- Transform raw crime data into a decision-support framework.

- Explain where crime concentrates across neighborhoods.

- Clarify when incidents occur by time and season.

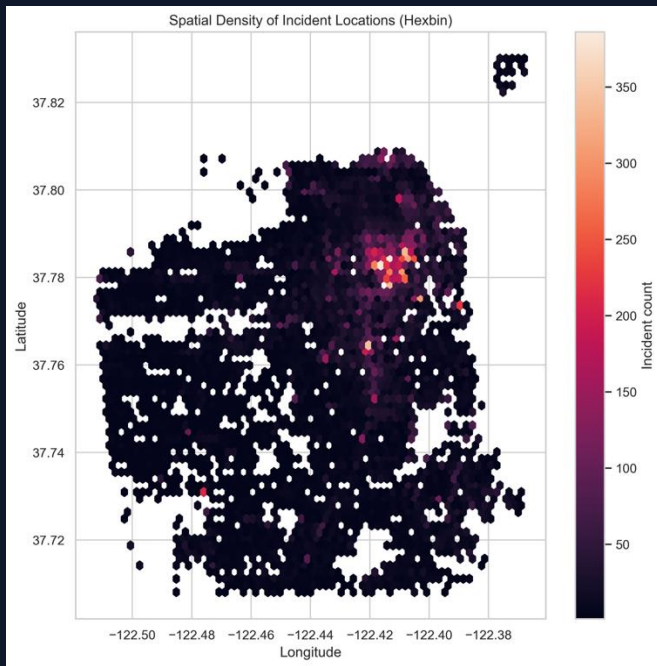- Provide a validated outlook for near-term planning.

# Key Questions This Project Answers

- **Where** are incidents most concentrated across neighborhoods?

- **How** does crime vary by hour, weekday, and season?

- **How** has crime changed from 2018 through 2025?

- **Did crime rebound** after COVID or structurally decline?

- **What** does the 2026 outlook suggest?

# Dataset Scope & Integrity

- **994,600 SFPD incident records.**
- **Full coverage** from Jan 1, 2018 to Dec 31, 2025.
- **41 official** Analysis Neighborhoods.
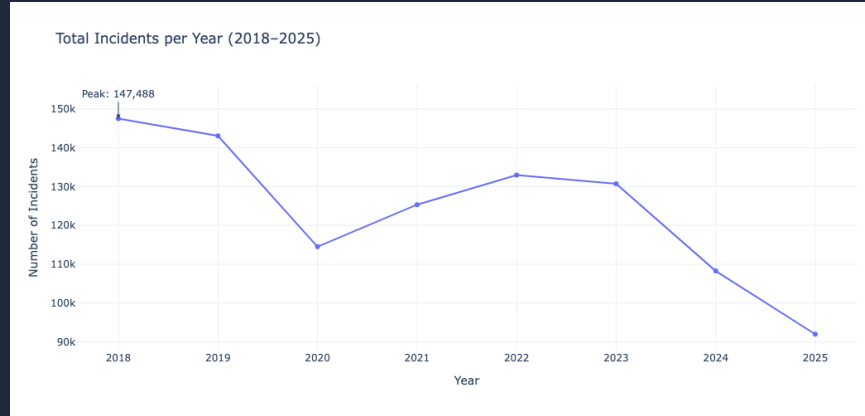- **Zero invalid timestamps** after validation.



Spatial Density of Incident Locations (Hexbin)

# Long-Term Change: Crime Volume (2018–2025)

**Total Reported Incidents by Year (2018 - 2025)**

↓ **37.6%**

$$\frac{(92001 - 147488)}{(147488)} \times 100\%$$



Total Incidents per Year (2018–2025)

Key Insight: The decline reflects a structural shift rather than short-term volatility.
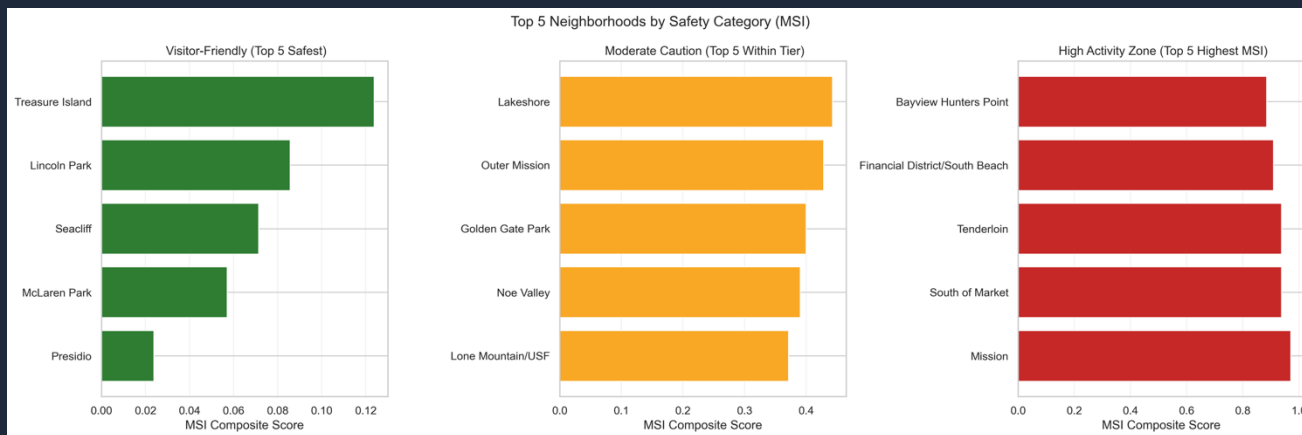
# Executive Summary: What the Data Says

- **Citywide crime declined 37.6%** from 2018 to 2025.

- **Crime is highly concentrated** in a small number of neighborhoods.

- **Time of day explains more variation** than day of week.

- **Post-COVID crime levels stabilized** rather than rebounded.

- **2026 outlook** shows continuity, not escalation.

# How I Approach Data Problems

- **Establish data integrity first.**
- **Identify structural patterns** before modeling.
- **Use relative metrics** for fair comparison.
- **Validate results** before forecasting.
- **Translate findings** into decision-relevant insight.

# Mobility Safety Index (MSI)
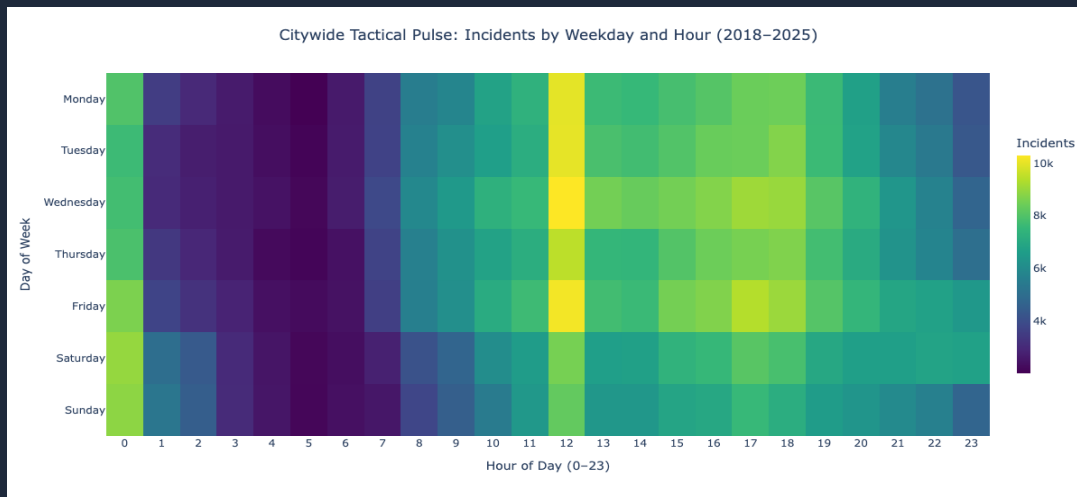
**Top 5 neighborhoods ranked by MSI score**



Top 5 Neighborhoods by Safety Category (MSI)

**MSI combines** incident volume, visitor-impact exposure, and recent momentum into a single score.

**On the 0 – 1 scale,** values closer to 0 indicate lower exposure, while values closer to 1 indicate higher and more persistent activity.

**Purpose:** comparison, not labeling.

# When Crime Happens



Hour (rows) × Day of Week (columns)

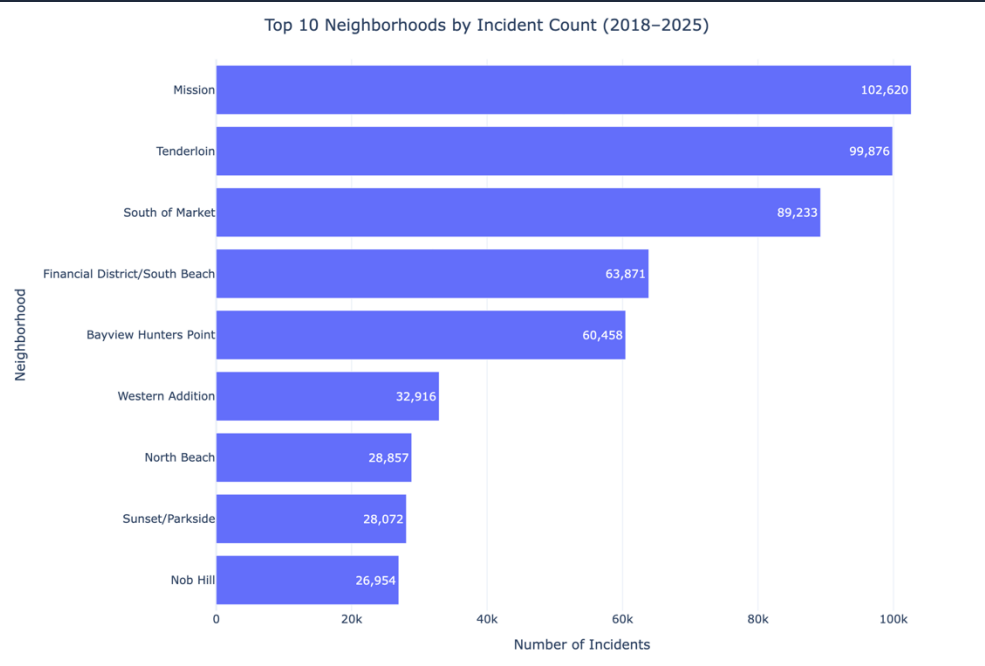Citywide Tactical Pulse: Incidents by Weekday and Hour (2018–2025)

- Incidents peak around midday on weekdays, when city activity is highest.
- Weekday incident activity reaches its highest levels around midday, reflecting peak urban movement.

**Key Insight:** Time of day matters more than day of week

# Where Crime Concentrates

### Incidents by neighborhood

- **Total incidents citywide (2018–2025): 939,401**

- **Top 3 neighborhoods: 291,729 incidents**
  - ➤ **31.1% of all incidents**

- **Top 5 neighborhoods: 416,058 incidents**
  - ➤ **44.3% of all incidents**

Top 10 Neighborhoods by Incident Count (2018–2025)

| Neighborhood | Number of Incidents |
|---|---|
| Mission | 102,620 |
| Tenderloin | 99,876 |
| South of Market | 89,233 |
| Financial District/South Beach | 63,871 |
| Bayview Hunters Point | 60,458 |
| Western Addition | 32,916 |
| North Beach | 28,857 |
| Sunset/Parkside | 28,072 |
| Nob Hill | 26,954 |

**Key Insight:** A small number of neighborhoods account for a disproportionate share of total incidents.
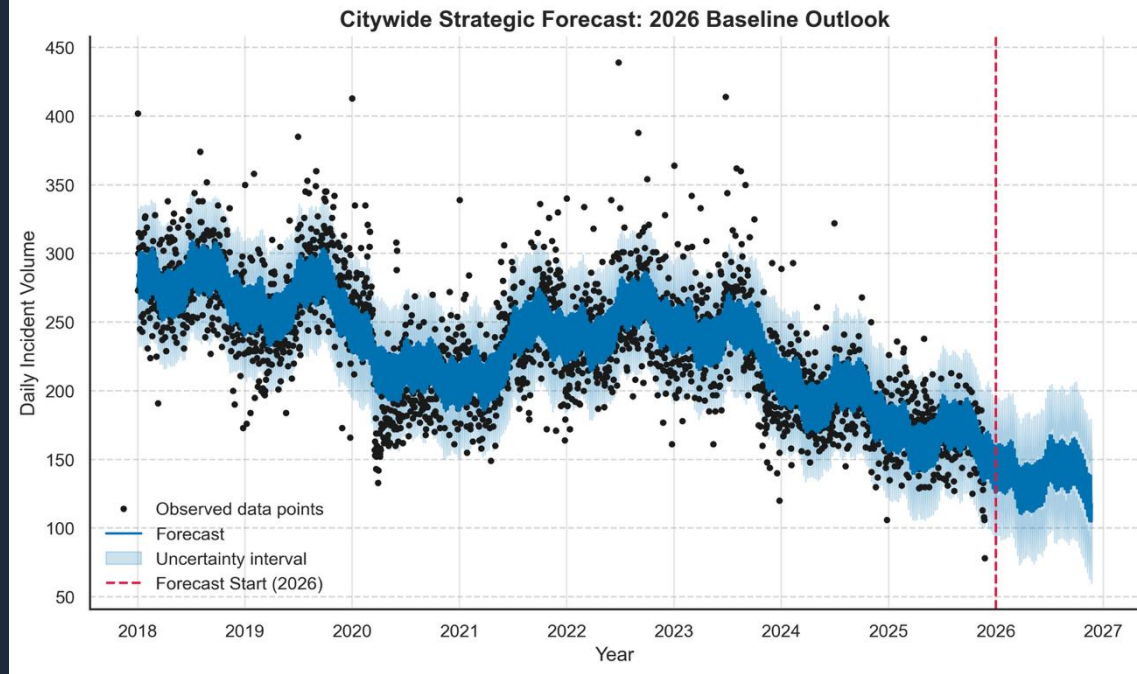
# Model Selection as a Business Decision

- **Multiple models evaluated** to avoid bias.
- **Criteria:** stability, interpretability, seasonal reliability.
- **Prophet selected** for balance of accuracy and usability.
- **Chosen** to support stakeholder decision-making.

# Citywide Baseline Forecast

## Actual vs Predicted (2025)

- The model captures long-term decline and recurring seasonal patterns observed since 2018.
- The 2026 forecast extends this downward trend, assuming no major structural shocks.
- Uncertainty increases into 2026, reflecting greater variability as projections move forward.



Citywide Strategic Forecast: 2026 Baseline Outlook

**Key Insight:** The forecast extends an existing downward trend, with uncertainty widening into 2026.

# Forecast Performance & Validation (2025 Holdout)

- **Citywide Prophet:** MAE ~127 incidents/month, RMSE ~156 (explains 84% variance)

- **Cross-validation:** 5-fold CV on 2018-2024 training window; consistent performance

- **2025 validation:** Model tested on blind holdout year before 2026 projection

- **Structural assumptions:** No reversion to 2018 peaks; stabilization at 2023-2024 baseline (~8.5-9.5K/month citywide)
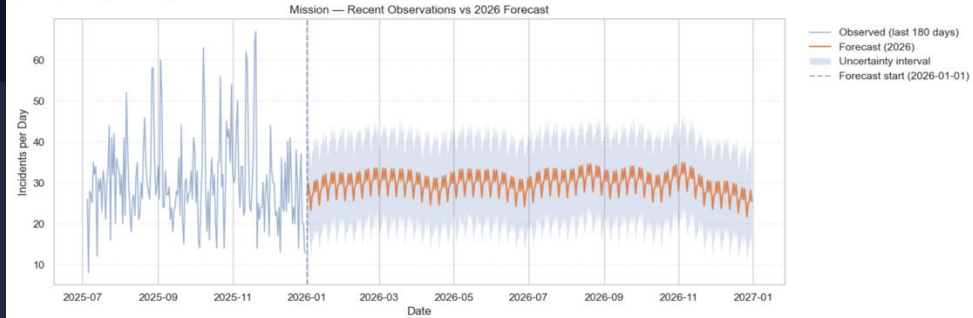
# Top 3 Hotspot Forecasts: 2026

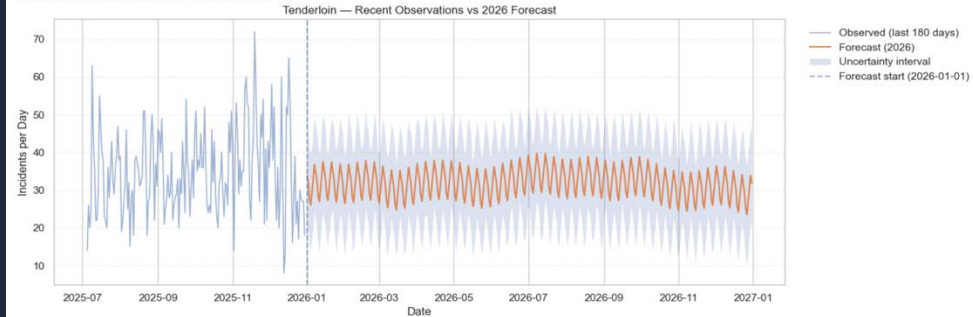## Forecast Line with confidence intervals

- All three hotspots remain elevated into 2026, even as citywide levels decline.
- Forecasts smooth daily noise to reveal stable underlying patterns.
- Uncertainty widens over time, reflecting sensitivity to external conditions.

**Interpretation:** Citywide improvement does not eliminate persistent hotspot risk.
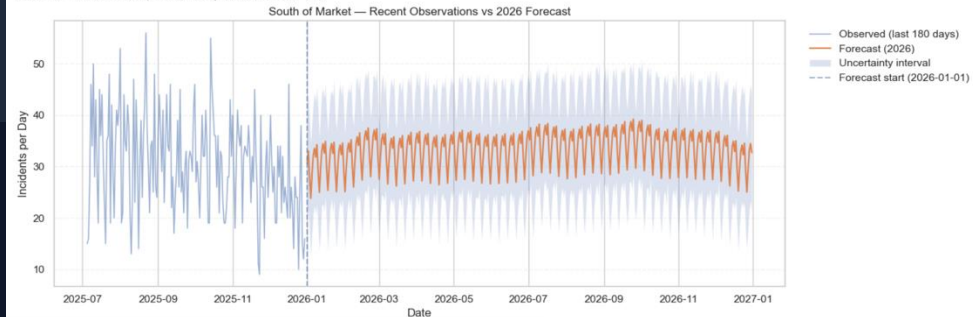
# What I Would Recommend as an Analyst

- **Align staffing** to time-of-day patterns.
- **Focus resources** on top 3-4 neighborhoods.
- **Use MSI** as a monitoring signal, not a label.
- **Review trends quarterly** to detect drift.

# Limitations & Responsible Use

- Reported incidents only.

- No socioeconomic causal modeling.

- Not a surveillance escalation tool.

- Insights should guide prevention and investment.

# From Data to Decisions

- **Start with data integrity.**

- **Identify structure and trends.**

- **Validate before forecasting.**

- **Communicate insights for decisions.**
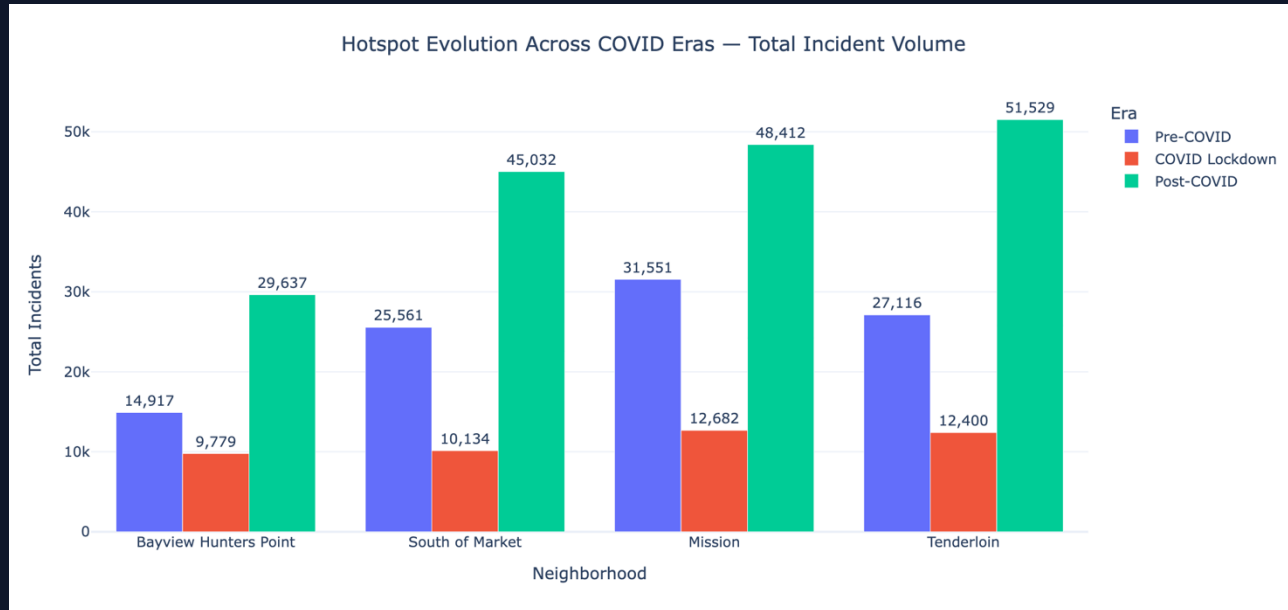
# Live Dashboard & Reproducible Artifacts

🔗 **LIVE Dashboard (Streamlit)**   https://project2sfcrimeportfolio-g6jhgqizljzqexcb3ss7wd.streamlit.app/

- ✓ Interactive filters by year range, neighborhood, and incident category

- ✓ Monthly incident trends from precomputed parquet aggregates

- ✓ Top neighborhoods and categories based on filtered views

- ✓ Exploratory 2026 outlook aligned with offline forecasting results

- Jupyter: Full reproducible pipeline (cleaning → EDA → MSI construction → forecasting → validation)

- GitHub repository: Clean, modular codebase with requirements.txt and documented methodology

- Data Artifacts: Parquet aggregates (~ 995K records), CSV exports, neighborhood tables

# Strategic Insights for Decision Makers

• **Concentrated risk:** Nearly half of incidents (44%) are concentrated in the top five neighborhoods, supporting focused interventions.

• **Structural stabilization:** Post-2020 incident levels remain ~15–25% below the 2018–2019 baseline, with no automatic rebound evident.

• **Category matters:** Larceny (29%) dominates; assault (6.5%) concentrated in nightlife; drug offenses structural in Tenderloin

• **Mobility integration:** High-activity areas align with transit hubs, tourist corridors, and nightlife density.



Hotspot Evolution Across COVID Eras — Total Incident Volume

# Strategic Applications & Stakeholder Impact

## Law Enforcement

Predictive patrol optimization. Resource allocation to top 3 neighborhoods (40% of incidents). Shift staffing aligned to hourly patterns.

## City Planning

Environmental design interventions in high-larceny zones. Lighting & transit improvements in assault hotspots. Community safety programs (seasonal).

## Community/Tourism

Risk-aware travel guidance. Neighborhood safety scorecards. Visitor information systems. Predictive event security planning.

# Limitations & Recommended Extensions

- **Data scope:** Reported crimes only (excludes unreported, victimization bias); no socioeconomic covariates

- **Forecast horizon:** 2026 only; longer horizons require external features (economic, policy, events)

- **Next steps:** Integrate foot traffic data, economic indicators, weather, major events; build micro-location grids (500m²)

- **Prescriptive layer:** Evolve from predictive → optimization (optimal patrol allocation, intervention timing)

# Conclusion: From Data to Decision Support

- Unified framework: MSI + Hierarchical forecasting = data-driven resource planning
- Rigorous validation: 2025 holdout year, cross-validation, model benchmarking ensure confidence in 2026 outlook
- Actionable insight: Specific recommendations for police, city planning, community stakeholders
- Reproducibility: Full pipeline archived in GitHub; dashboard scales to new data (monthly updates)