# Detailed Analysis of the SafeRouteMLModel: A Machine Learning-Based Route Safety Evaluation System

August 31, 2025

## Contents

# 1    Introduction

The SafeRouteMLModel is an advanced machine learning framework designed to assess and rank travel routes based on crime risk factors. By integrating geospatial clustering, nearest-neighbor searches, and ensemble regression, the model quantifies safety using a score from 10 to 100. This report expands on the original analysis by incorporating mathematical formulations, such as the geodesic distance formula and regression equations, to provide a deeper understanding suitable for academic scrutiny.

In urban environments like Delhi, where the dataset is sourced, route safety is paramount. The model uses historical crime data to predict risks, considering variables like crime severity and time of day. Mathematically, safety can be modeled as a function $s = f(\mathbf{x}, \theta)$, where $\mathbf{x}$ is a feature vector and $\theta$ are learned parameters.

The system's hybrid approachcombining rule-based penalties with ML predictionsensures interpretability and accuracy. Reported custom accuracy reaches 98%, defined as predictions within $\pm 5$ points of ground truth, calculated as:

$$\text{Accuracy} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{I}(|y_i - \hat{y}_i| \leq 5) \times 100\%$$

where $\mathbb{I}$ is the indicator function, $y_i$ the true score, and $\hat{y}_i$ the predicted score.

# 2    System Architecture and Dependencies

## 2.1    Dependencies and Library Analysis

The architecture revolves around the `SafeRouteMLModel` class, with Flask for API deployment. Key dependencies include:

- **Flask and CORS**: Enable web serving and cross-origin access.

- **Pandas**: For DataFrame operations, e.g., `df = pd.read_csv(file_path)`.

- **Geopy**: Computes distances via `geodesic(point1, point2).km`.

- **SciPy KDTree**: Spatial indexing for $O(\log n)$ queries.

- **NumPy**: Array manipulations, e.g., `np.mean(cluster_points, axis=0)`.

- **Scikit-learn**: Provides RandomForestRegressor, LabelEncoder, DBSCAN, and train_test_split.

- **Pickle and Logging**: Model persistence and debugging.

Example code from initialization:

```
self.crime_points = df[['latitude', 'longitude']].values
self.kd_tree = KDTree(self.crime_points)
```

## 2.2 Mathematical Foundations of Key Libraries

KDTree uses balanced tree partitioning in $k$-dimensional space (here $k = 2$ for lat/lon). Query time is $O(\log n)$ on average. DBSCAN groups points with density reachability, detailed in Section 4. Random Forest aggregates decision trees, with prediction:

$$\hat{y} = \frac{1}{T} \sum_{t=1}^{T} h_t(\mathbf{x})$$

where $h_t$ is the $t$-th tree's output, $T = 100$.

# 3 Data Loading and Preprocessing

## 3.1 CSV Parsing and Validation

Loading begins with `pd.read_csv`, checking columns like 'Latitude'. If missing:

```
if not all(col in df.columns for col in required_columns):
    raise ValueError("Missing required columns in dataset")
```

## 3.2 Spatial Data Conversion and Indexing

Conversions: `df['latitude'] = df['Latitude'].astype(float)`. Points array: $\mathbf{P} = [(lat_i, lon_i)]_{i=1}^{m}$, $m$ records. KDTree construction balances points for queries.

## 3.3 Mathematical Representation of Geospatial Points

Points are in spherical coordinates but treated as Euclidean for small areas. Actual distance uses the Haversine formula in geodesic:

$$d = 2r \arcsin\left(\sqrt{\sin^2\left(\frac{\Delta\phi}{2}\right) + \cos\phi_1 \cos\phi_2 \sin^2\left(\frac{\Delta\lambda}{2}\right)}\right)$$

where $\phi, \lambda$ are latitude/longitude, $r = 6371$ km. Maximum severity: $s_{\max} = \max(s_i)$, used in normalizations.

# 4 Crime Data Clustering with DBSCAN

## 4.1 DBSCAN Algorithm Overview and Parameters

DBSCAN identifies clusters as dense regions separated by noise. Parameters: $\epsilon = 0.001°$ ($\approx 111$m), `min_samples=5`. Fit:

```
clustering = DBSCAN(eps=0.001, min_samples=5).fit(self.
    crime_points)
```

## 4.2 Centroid and Severity Calculations

For each cluster $c$:

$$\mathbf{centroid}_c = \frac{1}{|c|} \sum_{\mathbf{p} \in c} \mathbf{p}, \quad avg_severity_c = \frac{1}{|c|} \sum_{i \in c} s_i$$

## 4.3 Mathematical Formulation of Clustering

A point $\mathbf{p}$ is core if $|N_\epsilon(\mathbf{p})| \geq \text{min\_samples}$, where $N_\epsilon$ is the $\epsilon$-neighborhood. A cluster is a maximal set of density-connected points.

# 5 Feature Engineering and Extraction

## 5.1 Nearby Crimes Querying with KDTree

For route points $R = [r_1, \ldots, r_k]$, query: `indices = kd_tree.query_ball_point(`$r_i$`, r/111)`, $r = 0.1°$. Deduplicate with set of CrimeIDs.

## 5.2 Distance Calculations Using Geodesic Formula

Total distance:

$$d = \sum_{i=1}^{k-1} \text{geodesic}(r_i, r_{i+1})$$

## 5.3 Aggregate Feature Computations

Total severity: $\sum s_j$ for nearby $j$. Average: $\bar{s} = \sum s_j / n$, $n = \text{total\_crimes}$. High severity count: $\sum \mathbb{I}(s_j \geq 0.6 s_{\max})$. Crime types: $|\text{unique categories}|$.

## 5.4 Penalty and Safety Score Formulas

Penalties:

$$severity_penalty = \frac{\sum s_j}{max_crimes \cdot s_{\max}}, \quad high_severity_penalty = \frac{high_count}{max_crimes} \cdot 0.5, \quad hotspot_penalty = hig$$

Safety score:

$$s = \max(10, 100 \cdot (1 - p_{sev} - p_{high} - p_{hot}) \cdot m_t)$$

where $m_t$ is the time multiplier (e.g., 0.7 for Night).

## 5.5 Time Encoding and Multipliers

LabelEncoder maps to integers 0–3. Multipliers model temporal risk. Features: $[n, \bar{s}, \max s, high_count, d$

# 6 Synthetic Data Generation for Training

## 6.1 Random Route Simulation

Loop 2000: Base lat/lon $\sim U(28.4, 28.8)$, $U(77.0, 77.4)$. Coords: 15 points $\pm 0.05°$. Time random choice.

## 6.2 Statistical Distribution of Synthetic Data

Routes approximate normal distribution around Delhi center, ensuring uniform feature space sampling.

# 7 Model Training with Random Forest Regressor

## 7.1 Data Splitting and Random Forest Basics

Split: 80/20, seed 42. RF: Bootstrap aggregated trees.

## 7.2 Mathematical Model of Random Forest Regression

Each tree trained on bootstrap sample, random features. Prediction as above. Variance reduction:

$$Var(\hat{y}) = \rho\sigma^2 + \frac{1-\rho}{T}\sigma^2$$

where $\rho$ is correlation.

## 7.3 Training Process and Serialization

Fit on $X_{\text{train}}$, $y_{\text{train}}$. Serialize with Pickle.

# 8 Model Evaluation and Custom Accuracy Metrics

## 8.1 Prediction and Error Analysis

$y_{\text{pred}} = \text{model.predict}(X_{\text{test}})$. Errors: $\text{abs}(y_{\text{test}} - y_{\text{pred}})$.

## 8.2 Custom Tolerance-Based Accuracy

Tolerance $\tau = 4$:

$$acc = \frac{\sum \mathbb{I}(|e_i| \leq \tau)}{n} \times 100 = 92\%$$

## 8.3 Statistical Interpretation of Performance

Mean Absolute Error implied $< 5$ for $92\%$. Bootstrap resampling could add confidence intervals.

# 9 Route Generation Using OSRM API

## 9.1 API Query Construction

URL with `alternatives=3`, geojson. Parse coordinates.

## 9.2 Route Deduplication and Filtering

Use tuple keys for uniqueness. Filter $d \leq 1.5 \cdot \min d$, top 7.

## 9.3 Mathematical Optimization in Route Selection

Minimizes distance while maximizing alternatives, akin to multi-objective optimization.

# 10 Route Evaluation and Safety Scoring

## 10.1 Dynamic Parameter Adjustment

Set `max_crimes` over routes, `max_crimes_per_route`.

## 10.2 Hybrid Scoring Mechanism

Final $= 0.5 \cdot \text{raw} + 0.5 \cdot \text{pred}$, clamped [10,100].

## 10.3 Ranking and Output Serialization

Sort by score descending.

# 11 Serialization and API Integration

## 11.1 Recursive Serialization Function

Converts NumPy types to Python.

## 11.2 Flask Endpoints and Data Flow

`/evaluate_routes`: POST source/dest/time $\rightarrow$ ranked JSON. `/model_performance`: GET metrics.

# 12 Performance Analysis and Insights

## 12.1 Accuracy Breakdown and Visualizations

98% indicates excellent fit. Example plot:

```python
import matplotlib.pyplot as plt
plt.scatter(y_test, y_pred)
plt.show()
```

## 12.2 Computational Complexity Analysis

KDTree query: $O(k \log m)$, $k$ route points, $m$ crimes. RF predict: $O(T \cdot \text{depth}) \approx O(100 \cdot 10)$.

## 12.3 Real-World Applicability and Case Studies

Example: Route A vs B, scores 85 vs 60 due to hotspots. Night multipliers reduce scores by 30%.

# 13  Mathematical Extensions and Derivations

## 13.1  Derivation of Safety Score Formula

Start with base 100, subtract normalized risks. Penalty as expected loss: $p_{\text{sev}} \approx P(\text{crime} \mid \text{route}) \approx \text{density} \cdot \text{severity}$.

## 13.2  Probabilistic Interpretation of Features

Total crimes $\sim \text{Poisson}(\lambda)$, $\lambda = \text{density} \cdot \text{length}$.

## 13.3  Optimization Formulations

Route selection as $\arg\max_s$ routes, subject to $d \leq$ threshold.

# 14  Limitations and Potential Improvements

## 14.1  Current Constraints

Synthetic data bias; no temporal decay in crimes.

## 14.2  Proposed Enhancements with Mathematical Backing

Add decay: $s_i \cdot \exp(-t/\tau)$, $\tau$ half-life. Use NN: loss $= \text{MSE}(y, \hat{y})$.

# 15  Conclusion

This detailed report, with mathematical integrations, underscores the model's sophistication. Convert to PDF using a LaTeX compiler like Overleaf or `latexmk`.

# 16  Appendices

## 16.1  Code Snippets

Full `train_model` code analyzed.

## 16.2  Example Calculations

For a route with 10 crimes, avg $s = 3$, max=5: penalties calculated.

## 16.3  Glossary of Mathematical Terms

- **Geodesic**: Shortest path on sphere.

- **Indicator** $\mathbb{I}$: 1 if true, 0 else.