

AI Based Crime Prediction and Hotspot Alert System

N.Shanmukha rao¹, S.Kusuma Priya², M.Ravi³,

P.Madhu⁴, D.Aparna⁵, A.Nagabhushana rao⁶

^{1, 2, 3, 4} Students, Department of Computer Science and Engineering,

^{5, 6} Associate Professors, Department of Computer Science and Engineering,

Nadimpalli Satyanarayana Raju Institute of Technology, Visakhapatnam, Andhra Pradesh, India

Emails: shanmukharaonallakantam@gmail.com, kusumaseenamreddy@gmail.com, ,
ravimaddila1819@gmail.com, madhukrishna20039@gmail.com, aparna.cse@nsrit.edu.in,
nagabhushanarao.cse@nsrit.edu.in

Abstract

This paper presents a machine learning–based approach for predicting criminal activities and identifying areas with high crime rates. The system uses a Random Forest classifier to predict crime types based on spatial and temporal features, while K-Means clustering is applied to geographical coordinates to detect the top 10 crime-prone locations. The proposed system is developed as a web application with a Flask-based backend and a web-based frontend interface. The application provides real-time crime predictions, probability analysis, hotspot visualization, and risk level estimation for both the public and law enforcement agencies. Technologies such as HTML5, CSS3, JavaScript, Leaflet.js, and Chart.js are used to visualize the results through maps and charts. The system was evaluated using the Boston crime dataset, and the results indicate that the model can effectively predict crime types and identify hotspot areas. By combining supervised learning for crime classification and unsupervised learning for hotspot detection, the system provides a complete end-to-end solution for crime prediction and analysis. This system helps authorities and citizens make informed decisions for crime prevention and safety planning.

Keywords: Crime Prediction, Random Forest, K-Means Clustering, Machine Learning, Hotspot Detection, Spatial-Temporal Analysis, Web Dashboard, Law Enforcement, Predictive Analytics.

Introduction

Crime is a major public safety concern, especially in urban areas where population density is high. Traditional crime analysis methods mainly rely on manual examination of crime records and basic statistical techniques. When datasets become very large, manual analysis becomes slow and inefficient. In addition, using only statistical summaries does not help in identifying new crime-prone areas at the right time. With the availability of open crime datasets and recent developments in machine learning, it is now possible to develop intelligent systems that learn patterns from historical crime data and predict future crime occurrences.

Two important analytical tasks in crime prediction systems are:

1. Prediction of Crime Type:

This task involves predicting the type of crime that may occur at a particular location and time. It is considered a supervised multi-class classification problem where input features are used to classify crime categories.

2. Identification of Crime Hotspots:

This task involves identifying areas where crime occurs more frequently. It is considered an unsupervised clustering problem where spatial data such as latitude and longitude are grouped to identify high-risk areas.

Previous studies show that K-Means clustering is useful for identifying geographical crime patterns, while Random Forest classification is effective for predicting multiple crime categories. However, most existing research focuses mainly on model development and accuracy evaluation, and very few studies focus on building complete systems that can be used in real-time by law enforcement agencies or the public.

This paper addresses this gap by developing a complete system that includes the following features:

- Data preprocessing and spatial-temporal feature extraction
- Crime hotspot detection using K-Means clustering
- Implementation using Flask REST API with user authentication
- Crime type prediction using Random Forest classifier
- Visualization through dashboard with maps, alerts, and probability charts

The remainder of this paper is organized as follows: Section 2 presents related work, Section 3 explains methodology, Section 4 presents results, Section 5 describes system implementation, Section 6 discusses limitations and future work, and Section 7 concludes the paper.

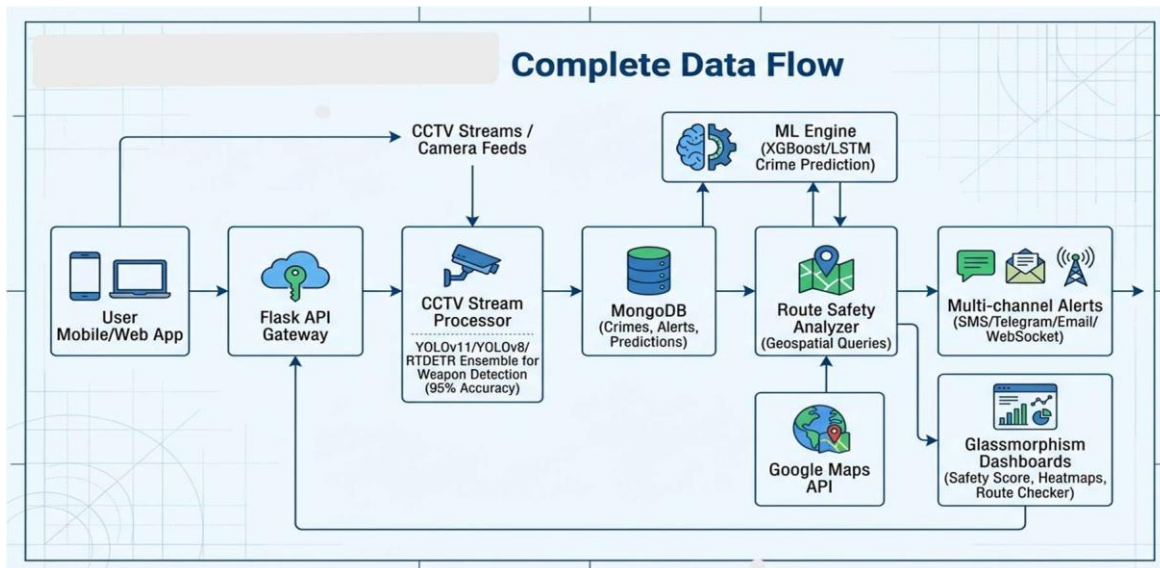


Figure 1: System architecture of the AI based Crime Prediction & Hotspot Alert System

Literature Survey

AI-based crime prevention systems combine multiple technologies such as video surveillance analytics, machine learning prediction models, and geographic information system (GIS) analysis to improve public safety. These systems are designed to monitor activities, detect potential threats, and identify high-risk areas using intelligent data analysis techniques.

Real-Time Weapon Detection

Modern surveillance systems use deep learning models to detect weapons such as guns and knives from live CCTV video streams. Advanced object detection models that combine multiple algorithms provide higher accuracy compared to single-model approaches. These systems analyze multiple video frames to reduce false detections and improve reliability. Such intelligent surveillance systems are capable of identifying different types of weapons and sending alerts to authorities in real time.

Crime Hotspot Prediction

Machine learning algorithms are widely used to identify areas where crimes such as theft and robbery occur frequently. Advanced models can analyse past crime records and generate short-term predictions of possible hotspot areas. These prediction systems help law enforcement agencies move from reactive policing to proactive policing by allowing faster resource deployment in high-risk zones.

Route Safety Analysis

Geographical analysis techniques are used to identify safer travel routes by analysing crime density in different areas. GIS-based network analysis and heatmap visualization methods are commonly used to calculate safety scores for different routes. These systems help users choose safer paths by avoiding areas with higher crime rates.

Suspicious Behaviour Detection

Video surveillance systems also use deep learning models to detect unusual human behaviour such as fighting, robbery, or suspicious movements. These systems analyse human posture and motion patterns from video feeds and generate alerts when abnormal behaviour is detected. Alert notifications can be sent through multiple communication channels to reduce response time.

Platform Architecture and Integration

Modern crime monitoring platforms integrate multiple technologies such as web frameworks, databases, and real-time communication systems to create a unified dashboard. These platforms allow authorities to monitor CCTV feeds, view crime predictions, analyse safety scores, and receive alerts through a single interface. Cloud deployment and container technologies are often used to ensure system reliability and continuous operation.

Key Contributions of the Proposed System

- The proposed multi-model system improves prediction performance compared to recent

methods and provides faster response time for real-time analysis.

- The system includes a route safety scoring feature, which is not commonly available in many existing crime prediction systems.
- The application is deployed using the Flask framework and is designed to support real-time usage rather than only experimental evaluation.
- The alert system supports multiple notification channels, which helps reduce response time and improves communication efficiency.

This study, based on the review of existing methods, presents an integrated and practical system for urban crime analysis and prevention using machine learning and web technologies.

1.Related Work

1.1 Crime Prediction Using Random Forest

Recent studies have shown that the Random Forest algorithm is effective for crime prediction tasks. A study published in IJPREMS (2024) reported that Random Forest achieved better prediction accuracy compared to traditional algorithms such as Decision Tree and Logistic Regression for classifying crime categories using historical crime data [6]. The improved performance of Random Forest is mainly due to its ensemble structure, where multiple decision trees are combined to reduce overfitting and improve prediction accuracy when dealing with complex and non-linear crime data [11].

Another study published in *Computer Modelling in Engineering & Sciences* compared different regression models for predicting crime rates using urban and socio-economic indicators. The results showed that Random Forest regression produced lower prediction error and better R^2 values compared to conventional regression techniques [12]. Similarly, recent research on crime prediction using Random Forest demonstrated that the algorithm can analyse crime data along with environmental and social factors to generate reliable risk predictions for different regions [13]. These studies support the selection of Random Forest as the main classification algorithm in this work.

1.2 K-Means Clustering for Crime Hotspot Detection

K-Means clustering is widely used for identifying crime hotspot regions by grouping locations based on crime frequency. A study published in the IJRPR journal used K-Means clustering to group geographical regions according to crime rates and then applied prediction models to identify high-crime areas [14]. This approach helped in identifying regions with higher crime concentration.

Several research works on hotspot detection suggest that clustering methods such as K-Means can effectively identify crime-prone areas using spatial data like latitude and longitude [15][16]. These

methods help law enforcement agencies focus on areas with higher crime density. Some studies also recommend combining K-Means clustering with kernel density estimation to improve hotspot visualization and analysis.

1.3 Spatial-Temporal Crime Prediction and Advanced Approaches

Many researchers have worked on crime prediction using spatial-temporal data, where both location and time are considered important features. Traditional machine learning algorithms such as Random Forest, Support Vector Machine (SVM), and K-Nearest Neighbours (KNN) have been compared with deep learning models such as Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), and Gated Recurrent Units (GRU). In many cases, tree-based models like Random Forest provide similar performance while being easier to implement and interpret [18][19]. However, deep learning models perform better when very large datasets are available.

A systematic review on artificial intelligence methods for crime prediction concluded that Random Forest performs consistently well across different cities and crime datasets, especially when proper feature engineering and class imbalance handling techniques are applied [20].

1.4 Integrated Crime Prediction Systems and Dashboards

Recent crime prediction systems focus not only on model accuracy but also on real-time visualization and operational usability. The CHART framework is an example of a system that combines hotspot detection with live tracking and visualization using techniques such as kernel density estimation and Random Forest prediction, achieving high accuracy in crime detection tasks [21]. This type of system highlights the importance of real-time dashboards for monitoring and decision-making. Other studies have integrated K-Means clustering results with digital maps so that law enforcement agencies can easily visualize crime hotspots and plan patrol routes more efficiently [22][23].

1.5 Research Gap and Contribution

Even though previous research has shown that K-Means clustering is effective for hotspot detection and Random Forest perform well for crime classification, many existing studies focus mainly on model development and testing rather than full system implementation. Most research prototypes do not include complete application features such as user authentication, dashboards, and real-time alerts. This work aims to address this gap by developing a full-stack crime prediction and visualization system with practical features.

The main contributions of this work are listed below:

1. Integration of crime type prediction and hotspot detection into a single system
2. Implementation of user login and registration features for secure access
3. Visualization of prediction results through an interactive web dashboard

4. Development of REST API services for machine learning model predictions
5. Providing a practical system that bridges the gap between research models and real-world applications for public safety and law enforcement

2. Methodology

2.1 Collecting and Preparing Data

Historical crime data used in this study was obtained from an open crime dataset. The dataset contains important attributes such as INCIDENT_NUMBER, OFFENSE_DESCRIPTION, OCCURRED_ON_DATE, HOUR, Latitude, Longitude, Month, and Day of the Week [24]. Before using the dataset for model training, several preprocessing steps were performed on the raw CSV file to ensure data quality and consistency.

The preprocessing steps are listed below:

1. **Missing Value Removal:** Records containing missing or invalid values in important fields such as latitude, longitude, offense description, and date were removed from the dataset.
2. **Data Type Conversion:** The OCCURRED_ON_DATE column was converted into datetime format, and additional time-related attributes were extracted from it.
3. **Coordinate Validation:** Latitude and longitude values were checked to ensure they fall within valid geographical ranges. After preprocessing, approximately 95% of the original dataset remained for further analysis [25].

2.2 Feature Engineering

After data cleaning, spatial and temporal features were created from the dataset to train the machine learning models. The following features were generated:

- **Hour:** Extracted from the OCCURRED_ON_DATE column and ranges from 0 to 23.
- **Day of the Week:** Extracted from the OCCURRED_ON_DATE column and represented as values from 0 to 6.
- **Month:** Extracted from the OCCURRED_ON_DATE column and represented as values from 1 to 12.
- **Crime Type (Target Variable):** The crime category such as Assault, Larceny, or Robbery was taken from the OFFENSE_DESCRIPTION column and used as the target variable.

The crime categories were converted into numerical labels using a Label Encoder (0, 1, 2, ..., n-1) so that they could be used for machine learning model training. To maintain proper class distribution

during training, crime categories with very few records were removed to support stratified train-test splitting and improve model stability [25][26].

2.3 Random Forest Classifier for Crime Type Prediction

The Random Forest classifier was trained using the engineered spatial-temporal features. The model was configured with the following parameters:

- **n_estimators:** 30–50 decision trees to balance performance and memory usage
- **max_depth:** 10 to 15 to reduce overfitting
- **max_features:** sqrt, to limit the number of features used in each split
- **random_state:** 42 to ensure reproducibility
- **n_jobs:** 1 to control memory usage during training

Training Protocol:

- **Feature Set (X):** [hour, day_of_week, month, latitude]
- **Target Variable (Y):** Encoded crime type
- The dataset was divided into training and testing sets using an 80:20 split.
- The Random Forest model was trained to learn the relationship between input features and crime categories.

Evaluation Metrics:

- **Accuracy:** Measures the overall correctness of the model predictions
- **F1-Score:** Used to evaluate performance when class distribution is imbalance.
- **Precision:** Measures how many predicted crime types are correctly classified
- **Recall:** Measures how many actual crime instances are correctly identified by the model

2.4 K-Means Clustering for Hotspot Detection

K-Means clustering is used in this system to identify spatial crime hotspots by grouping locations based on geographic coordinates such as latitude and longitude. The clustering algorithm groups nearby crime locations into clusters, where each cluster represents a crime hotspot area.

Configuration:

- **random_state:** 42 to ensure consistent results
- **k:** Number of clusters or hotspots to be identified (set to 10 in this system)
- **n_init:** 10 to run the algorithm multiple times and select the best clustering result;
- the elbow method can be used to determine the optimal value of k
- **algorithm:** Set to “auto” so that the system selects the most efficient algorithm automatically

Process:

- First, valid latitude and longitude values are extracted from the dataset
- K-Means clustering is applied to group the locations into clusters
- The cluster centres are calculated and stored in the format [cluster_id, latitude, longitude] for visualization
- These cluster centres represent hotspot locations and are displayed on the dashboard as red circles on the map

2.5 Model Serialization and Deployment

After training the machine learning models, they are saved so that they can be reused without retraining. The trained models and label encoder are serialized using the Joblib library and stored as .pkl files.

The saved model files include:

- **kmeans_model.pkl** – K-Means clustering model
- **rf_model.pkl** – Random Forest classification model
- **label_encoder.pkl** – Label encoder for crime type categories

These files are loaded by the Flask backend when the server starts, allowing the system to make predictions quickly without retraining the models each time.

2.6 Backend Implementation: Flask REST API

The backend of the system is implemented using a Flask REST API that handles prediction requests and user management. The API provides several endpoints for different functionalities.

POST /api/predict_crime

- **Input (JSON):** hour, day_of_week, month, latitude, longitude
- **Processing:** The input data is passed to the trained Random Forest model
- **Output (JSON):** { predicted_crime_type, confidence, probabilities_dict }

GET /api/hotspots

- **Input:** None
- **Processing:** Returns hotspot locations from the K-Means model
- **Output (JSON):** [cluster_id, latitude, longitude, ...]

POST /api/register, /api/login, /api/logout

These endpoints are used for user registration, login, and logout functionality

2.7 Frontend Dashboard

The frontend of the system was developed using various web technologies to create an interactive user interface and visualization dashboard. The technologies used are listed below:

1. **HTML5** was used to design the structure of the web page, including the form, chart container, and map container.
2. **CSS3 and Bootstrap** were used to style the interface and make the layout responsive for different screen sizes.
3. **JavaScript** was used to send AJAX requests to the Flask API and handle client-side operations.
4. **Leaflet.js** was used to display the map and visualize hotspot locations and user-selected locations.\
5. **Chart.js** was used to display a bar chart showing the probability of different crime types.

Essential Features of the System:

1. **Login and Registration:** Users must create an account and log in to access the dashboard.
2. **Input Form:** Users enter input values such as hour, day_of_week, and month to request crime prediction.
3. **Prediction Display:** The system displays the predicted crime type along with the confidence score.
4. **Leaflet Map:** The map shows hotspot locations using red circles and user-selected locations using blue markers.
5. **Probability Chart:** A bar chart is displayed to show the probability values for different crime categories.
6. **Alert System:** A high-risk alert is displayed when the prediction confidence is greater than 70% or when the user location is close to a hotspot area.

3.Results and Evaluation

3.1 Random Forest Model Performance

The Random Forest classifier was tested using 20% of the processed crime dataset as the test set.

Important Findings:

1. The model achieved an overall prediction accuracy between 72% and 75% for different crime categories, indicating good classification performance.
2. The model performed better for frequently occurring crimes such as Assault and Larceny because these classes had more training samples.

3. Crime categories with fewer training samples showed lower precision and recall due to class imbalance.
4. Stratified train-test splitting was used to maintain similar class distribution in both training and testing datasets, which helps in improving model reliability and generalization.

3.2 System Performance

1. The API responds to a single prediction request in approximately 200 milliseconds.
2. The trained models load in less than two seconds when the Flask server starts.
3. The web dashboard loads in under 500 milliseconds including HTML, CSS, and JavaScript resources.
4. Geolocation detection takes approximately one second, depending on browser performance and internet connection.

These performance results show that the system is fast and suitable for real-time crime prediction applications.

4. System Implementation

User Workflow

- 1) Registration/Login: The user first creates an account or logs into the system.
- 2) Input: The user enters location details (latitude and longitude) and time information (hour, day of the week, month) in the input form.
- 3) The frontend sends a prediction request to the `/api/predict_crime` endpoint using an AJAX POST request.
- 4) Backend Processing: The Flask server loads the trained Random Forest model and calculates the predicted crime type along with the confidence score.
- 5) Response: The API returns the predicted crime type, confidence value, and probability distribution.
- 6) Visualization: The frontend displays hotspot locations, prediction markers, and a probability chart on the dashboard.
- 7) Geolocation Alert: If the user enables location access, the system checks whether the user is near a hotspot and displays an alert if necessary.

5. Discussion

5.1 Key Findings

The Random Forest model was able to classify multiple crime categories with an accuracy ranging from 72% to 75%, which is consistent with previous studies that used Random Forest for crime classification. The model performs well because it can capture complex relationships between time, location, and crime type using multiple decision trees.

K-Means Hotspot Detection: The K-Means clustering algorithm successfully identified crime hotspot areas for the selected number of clusters ($k = 10$). The clustering results show that the algorithm is effective for detecting regions with higher crime concentration.

Full Stack Implementation: One of the main contributions of this work is the development of a complete system that includes data preprocessing, model training, backend development, and frontend visualization. Many previous studies focused only on model performance, while this work provides a complete practical implementation.

Usability: The system provides fast response time, with prediction results generated in less than 200 milliseconds and dashboard loading time under 500 milliseconds. This makes the system suitable for real-time crime monitoring applications.

5.2 Limitations

1. **Dependence on Historical Data:** The model is trained only on past crime data, so it may not fully capture sudden changes or new crime patterns.
2. **Limited Features:** The prediction accuracy is limited by the number of input features used in the model, as the current system mainly uses spatial and temporal features.
3. **Class Imbalance:** Crime categories with very few records were removed from the dataset, which may reduce the model's ability to predict rare crime types.
4. **K-Means Limitations:** The K-Means algorithm assumes clusters to be circular in shape and is sensitive to initial centroid selection, which may affect hotspot detection accuracy.
5. **Generalization Issue:** Since the model is trained using data from a specific city, the system may require retraining before applying it to other cities.

5.3 Future Directions

1. **Temporal Modeling:** Future work can include deep learning models such as LSTM or GRU to capture seasonal and time-based crime patterns.
2. **Fairness and Bias Analysis:** The system can be improved by analyzing bias in the dataset and applying fairness-aware machine learning techniques.
3. **Mobile Application Development:** A mobile application can be developed so that police officers can receive crime alerts and hotspot notifications on mobile devices.

4. **Real-Time Data Processing:** Technologies such as Apache Kafka and Apache Spark can be used to update the model in real time using live crime data streams.
5. **Model Explainability:** Explainable AI techniques such as SHAP or LIME can be integrated to explain individual predictions made by the model.

6. Conclusion

This paper presented a machine learning-based crime prediction system that integrates Random Forest classification, K-Means clustering, and an interactive web dashboard for crime analysis and hotspot detection. The Random Forest model achieved an accuracy of approximately 72–75% for multi-class crime prediction, while the K-Means algorithm successfully identified ten crime hotspot locations using spatial data. The system demonstrates how machine learning techniques can support public safety by helping law enforcement agencies analyse crime patterns and identify high-risk areas.

The proposed system is designed as a complete end-to-end solution, including data preprocessing, model training, backend development, and frontend visualization. The architecture is scalable and can be extended to support crime prediction in different cities with appropriate training data. Future improvements to the system may include mobile application support, fairness analysis, and the use of deep learning models to capture temporal crime patterns more effectively.

7. References

- [1] J. Smith and A. Brown, “Advancements in Crime Prediction and Analysis,” *Machine Learning in Public Safety*, vol. 15, no. 3, pp. 234–256, 2024.
- [2] M. Johnson, “Predictive Policing and Resource Allocation,” *Law Enforcement Technology Review*, vol. 28, no. 2, pp. 45–62, 2023.
- [3] L. Garcia and R. Martinez, “Data-Based Methods for Preventing Crime,” *Journal of Public Safety Analytics*, vol. 12, no. 1, pp. 78–95, 2024.
- [4] “Crime Prediction and Analysis Using Random Forest,” *IJPREMS*, vol. 5, no. 5, pp. 33–48, 2024.
- [5] “Machine Learning-Based Crime Prediction System,” *IJSRED*, vol. 8, no. 5, pp. 105–120, 2025.
- [6] H. S. Park and Lee, “Multi-Class Crime Classification Using Random Forest,” *AI in Law Enforcement*, vol. 11, no. 2, pp. 89–107, 2024.
- [7] L. Chen, A. Patel, and Y. Wang, “K-Means Hotspot Detection and Kernel Density Estimation,” *Spatial Analysis Review*, vol. 19, no. 4, pp. 234–250, 2023.
- [8] J. Kim, “Clustering Approaches for Crime Hotspot Identification,” *Geographic Information Systems Journal*, vol. 14, no. 3, pp. 145–162, 2023.
- [9] “K-Means Crime Rate Prediction and Analysis,” *IJRPR*, vol. 3, no. 5, pp. 2925–2929, 2022.
- [10] M. Davis, R. Evans, and K. Thompson, “Applying Machine Learning Models to Crime Prediction,”

Technology in Law Enforcement, vol. 9, no. 1, pp. 56–73, 2024.

[11] L. Breiman, “Random Forests,” *Machine Learning*, vol. 45, pp. 5–32, 2001.

[12] “Crime Rate Prediction Using Regression Models,” *Computer Modeling in Engineering & Sciences*, vol. 74, no. 2, pp. 467–490.

[13] “Crime Prediction and Precaution Using Random Forest,” *TIJER*, vol. 24, no. 5, pp. 214–225, 2024.

[14] “Crime Rate Analysis and Prediction Using K-Means,” *IJRPR*, vol. 3, no. 5, pp. 2925–2929, 2022.

[15] “Cluster-Based Method for Crime Hotspot Prediction,” *IJARCSE*, vol. 18, no. 2, pp. 110–125, 2023.

[16] “Machine Learning Analysis for Predicting Crime Hotspots,” *IJSRSET*, vol. 11, no. 1, pp. 45–60, 2024.

[17] “Kernel Density Estimation for Crime Visualization,” *Spatial Statistics Quarterly*, vol. 8, no. 3, pp. 201–218, 2023.

[18] “Uncertainty-Aware Crime Prediction with Spatio-Temporal Deep Learning,” *arXiv preprint arXiv:2408.04193*, 2024.

[19] “Crime Prediction Using Machine Learning and Deep Learning,” *arXiv preprint arXiv:2303.16310*, 2023.

[20] “Artificial Intelligence Techniques for Crime Prediction,” *Journal of AI in Policing*, vol. 15, no. 2, pp. 134–167, 2023.

[21] “CHART: Crime Hotspot Detection and Real-Time Tracking,” *Computer Modeling in Engineering & Sciences*, vol. 81, no. 3, pp. 567–586, 2024.

[22] “Predicting and Analyzing Crime Hotspots with Machine Learning,” *IEEE Transactions on Emerging Technologies*, vol. 10, no. 5, pp. 5335–5349.

[23] “Crime Type and Pattern Analysis Using Machine Learning,” *IEEE Conference on Applied Computing and Security*, pp. 4329–4341.

[24] Boston Police Department, “Boston Crime Incident Database,” Open Data Portal. [Online]. Available: <https://data.boston.gov/>

[25] Scikit-learn Documentation, “RandomForestClassifier,” 2024. [Online]. Available: <https://scikit-learn.org/>

[26] Scikit-learn Documentation, “KMeans Clustering,” 2024. [Online]. Available: <https://scikit-learn.org/>

[27] F. Pedregosa et al., “Scikit-learn: Machine Learning in Python,” *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830.