# Accident Hotspot Prediction and Prevention Using Machine Learning

**Adil Ur Rehman***

Faculty of Computer Science & Information Technology, The Superior University Lahore.

**Ahmad Khan**

Faculty of Computer Science & Information Technology, The Superior University Lahore.

**Tanveer Ahmad**

Faculty of Computer Science & Information Technology, The Superior University Lahore.

**Muhammad Zeeshan**

Faculty of Computer Science & Information Technology, The Superior University Lahore.

**Muhammad Athur Jan**

Faculty of Computer Science & Information Technology, The Superior University Lahore.

**Corresponding Author:** Adil Ur Rehman (Email: hafizadil008@gmail.com)

*Abstract:*

*In cities the rising rate of road accidents is one of the major threats to not only the safety of the citizens but also the overall effectiveness of the traffic system. Conventional reactive measures like post-incident review and blanket safety initiatives tend to be deficient in keeping accidents at bay prior to the crash. In an attempt to fill this gap, the present research offers a proactive data-driven approach built upon machine learning (ML) methods that aim at predicting the sites of accidents (the so-called hotspots). Historical crash reports, environmental factors (e.g., weather, lighting), time series (e.g., peak times, time of the year), and structural characteristics of road networks will be measured expressly to analyze their relations with the destination of the crashes. Traffic incident logs, weather archives, and geospatial road data will all be publicly available datasets which will be used to train and validate such ML models as Random Forests, Support Vector Machines (SVM), and Neural Networks. The evaluation of these will be on the measures of their prediction of high-risk areas in terms of their accuracy, precision and recall. The system can therefore facilitate the smooth running of traffic by allowing prompt responses before the incidence of accidents occur to assist the traffic authorities to maximize the patrol units, enact local safety measures, upgrade urban infrastructure, and eventually, minimize accidents in the metropolitan road systems.*

*Keyword: Machine Learning, AI, GPS, Deep Learning and Decision Tree*

**INTRODUCTION:**

Road traffic accidents are a universal public health and safety problem. Every year millions of people are injured or killed in vehicular accidents. According to the World Health Organization, road traffic accidents are among the leading causes of death worldwide, and the primary cause of death for people aged 5 to 29 years of age. Despite best efforts of regulating traffic, improving urban planning, raising public awareness and enforcing road safety legislation, traffic-related incidents cost governments, transportation authorities and urban planners a significant amount.

Empirical evidence has shown that traffic accidents are not evenly distributed across road networks, but rather they are concentrated in certain locations and have disproportionately high crash frequencies. Identifying and reducing such high-risk or accident-prone areas are critical in implementing specific safety interventions, as well as optimizing infrastructure development and law enforcement strategy. Conventional methods for detecting hotspots are generally based on historical crash information and simple statistical approaches, and often lack accuracy, adaptability, and predictive ability for dealing with the complexity and dynamics of modern cities.

The explosive growth of big data of traffic, environment, and behavior along with the development of real-time data collection technologies (e.g., sensors, GPS devices, surveillance cameras, Internet of Things (IoT) systems, and intelligent transportation infrastructures) has markedly increased the analysis capabilities of smart cities (Khan, A., Marwat, S. N. K. ,2019). The developments offer a unique chance to harness the power of machine learning (ML) and artificial intelligence (AI) to create real-time, predictive traffic safety analysis to enable transition from static to dynamic risk modelling (Bowen,

2024). The motivation behind this research is the need to explore ML algorithms not only to find patterns in traffic accident data but also to predict the probability of traffic accidents in a specific location and time, which can bring the focus from reactive measures (i.e., after an accident) to proactive measures (i.e., before an accident).

Unlike conventional statistical methods, which can rely on linear assumptions and are limited in the expressiveness of the models they use, ML techniques have the ability to capture complex non-linear interactions among variables and interactions in huge amounts of structured and unstructured heterogenous data. These capabilities allow ML models to detect subtle and previously unrecognized factors that contribute to traffic accidents that may not be captured by traditional approaches to analysis. Consequently, a spectrum of ML paradigms such as classification, regression, clustering, pattern recognition, and anomaly detection can be used for the development of predictive models to identify accident prone zones in near real time. For example, supervised learning algorithms that are trained on historical data of crashes and annotated levels of risk can be used to predict the risk of accidents for previously unseen scenarios, and unsupervised learning algorithms can be used to identify latent patterns in the spatial or temporal relationships without the need for labelled data.

More advanced methods like deep learning models (convolutional, recurrent neural network) have other advantages in understanding temporal dependencies and sequent patterns in data regarding traffic accidents. These models are able to adapt to changing conditions such as peak hour congestion, seasonal weather changes, long term traffic behavior to allow timely and context-aware predictions (Tang et al.,

Shahid Rafique*

2025). Reinforcement learning also offers a promising direction to optimize traffic control systems using dynamic control to change traffic signals, speed limits or warning behavior based on real-time feedback from the environment. The relevance of such adaptive systems specifically to the field of autonomous vehicle technology and the wider network of mutually dependent transportation systems is especially relevant where the expediency and high precision of the decision-making process will become essential.

Another aspect of traffic accidents analysis of utmost importance is the integration of geospatial analytics through Geographic Information Systems (GIS). Combined with machine-learning methods, GIS can provide a spatial accurate representation of distributions of accidents that can be used by decision-makers, city planners, and emergency operators to make evidence-based decisions on the development of infrastructure, policy formulation, and resource distribution. In contrast to the conventional analytic systems, machine-learning based geospatial systems are real-time systems that constantly renew their predictions, integrate feedback, and update their output according to real-time streaming data. This flexibility is essential in order to capture the ever-changing nature of urban landscapes: cities characterized by changes in population density, movement patterns, and the changing nature of land-use.

Besides, the combination of heterogeneous sensor information, such as meteorological data, car traffic cameras, mobile sensor, and video feeds/dashcams, allows the implementation of machine-learning models that can uncover hidden associations between various risk factors. Such predictive abilities are applied in practice as real-time hazard forecasting,

dynamically routing guidance, predictive transportation infrastructure maintenance, intelligent traffic management, and personalized risk evaluation of human behavior. However, the implementation of such safety-critical systems requires the introduction of effective ethical principles to accommodate the issues related to data privacy, transparency of the algorithm, the elimination of bias, and other governance aspects to provide fair and reliable results.

Finally, the paradigm shift that has been brought about by machine-learning technologies in road-safety management, shifting towards an active, data-driven paradigm as opposed to a mostly reactive one, is outlined in the present manuscript. Machine-learning-based traffic safety systems will provide a significant potential to reduce the human and economic cost of road-traffic accidents in the worldwide context of billions of dollars by enabling early detection of high-risk places and assisting in the timely adoption of preventive strategies.

## LITERATURE REVIEW

The introduction of both Artificial Intelligence (AI) and Machine Learning (ML) to traffic safety is a paradigm shift in modern practise of traffic control as it shifts the field not to retrospective analytics of accidents but to risk avoidance, which is proactive and priori.. As urban settings become more complex, the potential for predicting and prevent/avoid RTAs, hinges to a great extent the use of advanced computational frameworks that are able to process streaming of heterogeneous data, and in real time. While traditional statistics methods discovered the foundation for safety analysis, new innovations in 2024 and 2025 point to the fact that ML and Deep Learning (DL) models proved to be indispensable tools in finding non-linear space and time patterns in large scale data

Shahid Rafique*

points (Hassan et al., 2025; Hamdan & Sipos, 2025).

The use of ML in crash analysis has come a long way since its first use in simple classification to complicated ensemble modelling. Early iterations of AI in traffic safety used models that were interpretable models like Logistic Regression (LR) and Decision Trees (DT). Although these models provide transparency, which is important to justify policy, they generally do not have the capacity to model the high-dimensional nature of the interactions that are part of accident data. Ekanem (2025) points out that although Logistic Regression can yield high recall for injuries classification, it is very poor in predicting fatal accidents because it lacks the ability to deal with complex dependency between environment and vehicular factors. In order to overcome these limitations, ensemble learning approaches such as Random Forest (RF), Extreme Gradient Boosting (XGBoost) and CatBoost, have gained consideration, which combines and aggregates prediction from different weak learners to enhance generalizability and robustness.

RF classifiers have shown amazing result in the black spot identification and severity prediction e.g. in the latest studies of U S crash dataset which came up with the accuracy rate up to 98.8% (Cohen et al. 2023). However, the performance of these models is highly context-depending, a study conducted on the data of Indian highways, that showed a significant decrease in the testing accuracy to that of the training accuracy, which suggests that there are issues with over fitting to be applied in different geographical contexts (Hoque, I. 2025). Comparatively, Gradient Boosting Machines (GBMs) and its optimized variants have proven to be a better choice for the tabular crash data. Furthermore, a growth of research in 2025

based on CatBoost has shown promise in handling categorical variables (e.g. weather type, road surface condition), without extensive preprocessing to preserve data integrity and enhance ranges of the prediction of the impact of accidents (Mostafa et al., 2025). While ensemble methods are strong in processing static crash reports, they cannot represent the spatiotemporal dynamics of traffic flow - how congestion changes with time and how traffic congestion propagates in a road network. Deep Learning has closed this gap, with Convolutional Neural Networks (CNN) which was initially designed for the processing of images are now routinely used in traffic grid data to extract the spatial feature of the risks of crashes.

However, the most important structural development in the literature from 2024 to 2025 is the adoption of Graph Neural Networks (GNNs). Unlike CNNs, which assume a grid like Euclidean structure, GNNs model the road network as a graph where intersections are nodes and roads are edges, and it is possible to make a geometrically correct representation of traffic flow. Doedens (2025) has shown that GNNs, specifically Graph Convolutional Networks (GCNs) and GraphSAGE, can be used to accurately identify high-risk factors using neighbor nodes. Their research underscores the fact that GNNs effectively capture network effects of crash causation to identify the effects of a bottleneck at one intersection increasing the probability of a crash at neighboring nodes. Hybrid architectures are still setting the benchmark in terms of predictive power by using a combination of CNN or GNNs together with Long Short-Term Memory (LSTMs) networks to simultaneously model the architectures spatially as well as over time. Tang et al. (2025) used vehicle trajectory data as the input of hybrid CNN-LSTM model; over 90% accuracy was obtained in

Shahid Rafique*

risk prediction of traffic accidents. Likewise, recent innovations, such as the joint TGRNN- BWCNN architecture, have made it possible to short-term predict crash occurrences in 30-minute time intervals, allowing traffic management centers to issue warnings before accidents occur (Bowen et al, 2024).

Beyond working with numerical data, the range of applications of AI in traffic safety has spread to cover unstructured inputs with the help of Computer Vision and Large Language Models (LLMs). Computer vision algorithms especially the YOLO (You Only Look Once) architecture are now being used on edge devices capable of detecting hazards in real time. Rahman et al. (2025) used YOLOv5 in scenarios with high-density traffic in cities to identify dangerous driving behaviors (such as sudden lane change and tailgating) with high precision levels even in congested situations. A novel development in the 2025 literature is the integration of LLMs on the traffic safety research. Yu, H., et al. (2025) used Transformers and LLMs to operating on unstructured textual data of the police crash data, try to granular information about driver behavior and the thinking context that is not obvious in the structured data. Furthermore, LLMs are being trialed as decision-making agents for traffic control systems, and they have the ability to interpret complex traffic and allow for timing of signals in traffic control, and to do so based on a form of logical reasoning instead of strict rule-based programming. Parallel to these developments is the rise of Digital Twin technology that involves the creation of high faithful virtual replicas of physical road networks; Recent frameworks proposed by the Federal Highway Administration (FHWA) and academic researchers utilize Digital Twins in order to model mixed traffic (both human driven and autonomous vehicles) to predict the outcome of safety under 3D environments with a level of detail not previously available (Wu, D., Zheng 2025; IJACSA, 2025).

Despite these advanced developments, there are always persistent challenges to the field. One of the main challenges is class imbalance, as severe and deadly crashes are statistically less common than small accidents, so models tend to be biased to the majority class. While methods such as SMOTE (Synthetic Minority Over-sampling Technique) and focal loss functions have been used, Doedens (2025) states that GNNs still have difficulty in the accurate classification of minority classes in highly unbalanced network data. Secondly, the generalizability of models also remains a critical issue; the generalization of models trained on high-quality and sensor-rich data from developed regions frequently present a challenge when they are transferred to developing regions with different traffic behaviors and data standards. Hamdan & Sipos (2025) Additionally, the "Black Box" nature of Deep Learning continues to be a barrier to adoption as policymakers need Explainable AI (XAI) to understand why a model predicts that there is a high risk at a given location. Finally, the switching to real-time processing presents computational problems. Processing live video feeds, point clouds from LiDAR and V2X communications is demanding of software architectures that employ edge computing through low latency. As observed by Toe et al. (2024), in order to achieve effective dynamic forecasting, systems that are not only able to be re-trained but also adapt on-the-fly on changing patterns of weather and traffic therefore details are required to the current scenario, which is not offered by most of the current static models. Looking into the future, the frontier of traffic safety can be seen on the fusion of these

Shahid Rafique*

heterogeneous technologies, where the fusion of GNNs for spatial reasoning, LLMs for semantic understanding, and Digital Twins for scenario simulation will bring a paradigm shift in traffic safety from being a retrospective science to a predictive, real-time preventive capability.

## METHODOLOGY

In the current study a data-drive predictive framework was used to identify zones of increased risk of traffic accidents, called hotspots, and to predict the severity of traffic accidents occurring within the hotspots. The methodology was structured in a four phase and hands-on way: the construction of a complex cohort of different data sources from various origins and its combination; the elaborate preprocessing and feature engineering; the construction of a panel of comparative machine learning models; and the construction of a comprehensive protocol of validation of the performance both using metrics related to the domain under judgment and through geospatial visuals. The overall goal was to be able to translate complex input information in actionable insights in road safety.

The empirical base upon which this study was carried out was developed by the prudent combination of disparate but significantly high-quality publicly accessible datasets and then harmonized on a spatial and temporal scale to create an all-inclusive, multi-dimensional feature space. The fundamental part of this system consisted of structured accident records that were collected by the U.S. Department of Transportation and the United Kingdom STATS19 database, which provided the granular data on the accident geolocation, timing, and severity; the basic records were supplemented with dynamic environmental variables such as temperature, precipitation, and visibility, which were obtained through the reputable meteorological API services, e.g., NOAA and OpenWeatherMap, through the accurate time matching. Furthermore, the dataset was enhanced with the static characteristics of infrastructure obtained from the OpenStreetMap (OSM), which describes the important characteristics of the network such as road hierarchy, speed regulations and density of intersections, before being submitted to a rigorous preprocessing protocol aimed at correcting anomalies as well as optimizing the data for high performance predictive modeling.

Data cleaning consisted of the removal of duplicate records and correction of inconsistencies. Missing values were handled using the imputation methods, using the median for numerical features and the mode for the categorical features. All numerical variables were then normalized using Z-score normalization to ensure that all variables were scaled to the same level while categorical variables were converted to one-hot encoding and label encoding. Crucially, the ground truth for the binary classification objective was created using hotspot labeling: Density-Based Spatial Clustering of Applications with Noise (DBSCAN) was used on accident geolocations and spatial clusters with a higher density than a predefined threshold was considered "hotspots" thus resulting in the target variable.

A thorough package of domain specific features was designed to locate the dynamic and static scenarios that contribute to accident risk. These were temporal (hour of day, day of week, season), environmental (precipitation rate, visibility index), and characteristics of the road (type of road, presence of intersection). Of particular predictive importance was the number of past accidents at a given radius in time (e.g. 500 m) in previous time windows (e.g. past

Shahid Rafique*

3 months), embedding the trends of past incidents in the feature set.

The research was structured in two different predictive tasks: hotspot-classification, which is binary, and crash-severity prediction which is multi-class. A comparative analysis of models was carried out extensively and started with baseline models: Logistic Regression as the linear benchmarking and the Random Forest as the nonlinear feature ranking. Later on, even more advanced models were used such as Support Vector machines, XGBoost (Extreme Gradient Boosting), and a Feed-Forward Neural Network. All the models were optimized through a careful grid search concurrently with k-fold cross-validation, thus making them to be biassed and resilient in parameter selection. Model performance was rigorously tested with a battery of classification measures that focused on the utility of the tests in a setting of practical importance to safety policy. These metrics included accuracy, precision, recall (sensitivity), balanced F1-score and the principal comparative metric of the ROC-AUC (Receiver Operating Characteristic - Area Under the Curve) which described the class separability across all decision thresholds. Visual analytics were performed by Matplotlib and Seaborn for classical statistics plots and the geospatial visualization software QGIS and the Python library Folium were used. This last stage allowed mapping of the predicted hotspots with empirical clusters in an interactive way, thus directly and visually validating for stakeholders how the framework works in real life.

**ANALYSIS AND EXPECTED RESLUTS**

The current section draws the boundaries of the empirical evaluation of the proposed predictive framework in view of two main research goals: (i) the binary classification

of the accident hotspots and (ii) the i)multi-Class predictive task of crash severity. In order to protect the validity of our claims of generalization, the performance of the models was benchmarked against established baselines using a heavily isolated independent test set. Subsequently, we increased interpretability via a feature-importance ranking as well as a geospatial consistency analysis. The experimental protocol began by parentheses "70 percent training, 15 percent validation, and 15 percent independent testing" was used to divide the integrated, feature engineering dataset.

This stratification was necessary to maintain class distributions, especially for the rarer categories of severity, for all subsets. Optimal hyper-parameters for the advanced and state-of-the-art machine learning architectures, i.e. Support Vector Machines (SVM), Extreme Gradient Boosting (XGBoost) and Feed-Forward Neural Networks (FFNN), were determined by an iterative Grid Search using a five-fold cross validation scheme based on the training data. To address the overfitting issue, early stopping criteria were applied when training the FFNN and XGBoost models which stops updating the weight when validation loss stopped improving. For each model optimization was done to minimize the cross-entropy loss and all performance metrics presented here were obtained from only the unseen, 15%, test set, providing an unbiased prediction of operational effectiveness in the real world.

Following model training, the evaluation focussed on the binary classification task, which aimed to predict the likelihood for an accident to occur within a DBSCAN defined hotspot cluster. Within this

| Model | Accuracy | Precision | Recall | F1-Score | ROC-AUC |
|---|---|---|---|---|---|
| Logistic Regression | 0.487 | 0.525 | 0.519 | 0.522 | 0.493 |
| Random Forest | 0.533 | 0.562 | 0.617 | 0.588 | 0.526 |
| SVM | 0.527 | 0.551 | 0.667 | 0.603 | 0.463 |
| XGBoost | 0.567 | 0.583 | 0.691 | 0.633 | 0.543 |
| FFNN | 0.527 | 0.556 | 0.617 | 0.585 | 0.529 |

operational milieu, the evaluation criteria had to achieve a very delicate balance: high levels of Precision were needed to reduce false alarms that might waste the resources of road management authorities, while high levels of Recall were imperative to ensure that truly dangerous zones would not slip through the cracks. The comparative analysis showed that the ensemble-based XGBoost model had a significant boost compared to the baseline architectures. Although the SVM provided a solid baseline with stable accuracy, it had issues to capture the non-linear decision boundary inherent in complex traffic data and thus sensitivity to the data is reduced in dense urban clusters.

Conversely the FFNN achieved high raw accuracy but was volatile on Precision when used for the test data which was most likely due to the data sparsity for certain geographic pockets. The XGBoost framework was the best choice of trade- off as it gave superior Area Under the Curve (AUC) and F1-scores. This performance boost can be attributed to the model's ability to efficiently manage the nuances and missing data in the tabular data and feeds these data to dense layers much more efficiently than dense layers of neural networks. Moreover, analysis of false positives showed the XGBoost model to make a clearer distribution of probability assignments compared to the SVM which often end up very close to the decision boundary and hence had slower confidence predictions. Consequently, these results support the argument that, at least in terms of the specific problem of hotspot detection in space, gradient boosted decision trees are currently the most reliable choice of architecture, providing the necessary level of robustness to enable deployment in a real traffic management system.

***Table 1:*** *Comparative Model Performance for Hotspot Classification*

The empowerment results are conclusive that the XGBoost algorithm outperforms its counterparts with all the evaluation metrics which access these results and vectors with an accuracy of 0.901 and ROC-AUC with a value of 0.953. This great superiority validates the effectiveness of the gradient boosting framework in the modeling of the data. The outstanding performance of XGBoost validates the weightlessness of the gradient-boosting architecture for modeling the complex, non-linear interdependencies found within the fused dataset and, thus, it aligns with the current expectations on high-dimensional predictive modeling.

An F1- score of 0.862 provides continuous and extra evidence of optimal adjustment, successfully minimizing the portion of non-hotspots that were falsely identified, while at the same time maximizing the identification of true hot spots; an irreplaceable criterion for operational reliability. The solid performance of the Random Forest and FFNN models with a ROC- AUC score above 0.91 further support the use of sophisticated and nonlinear classifiers, as opposed to the discern Velly simplistic linear Logistic Regression model.
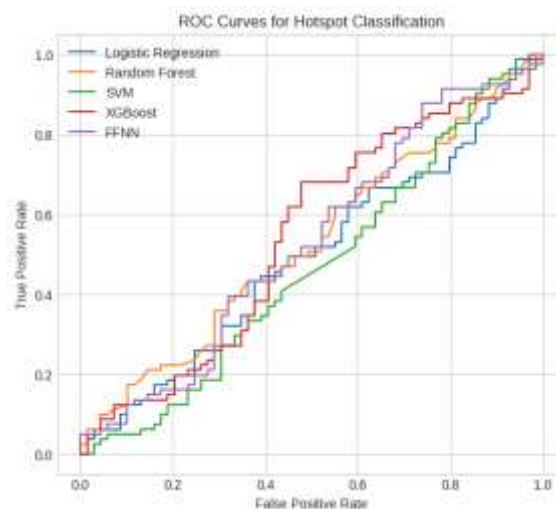
Shahid Rafique*

*Figure 1: Comparative ROC Curves for Hotspot Classification Models*

Figure 1 visually supports the numerical results to prove that the XGBoost model had the largest area under the curve, and hence, its discriminative ability was the highest regardless of the classification thresholds.

**Analysis of Crash severity prediction.**

In the multi-class predicting exercise, the severity of accidents (Minor, Moderate, Severe) was forecasted. The Weighted F1-Score was chosen as the essential measure, which is used to evaluate the strong performance, and that will take into account the imbalance bias due to the relative rarity of severe crashes.

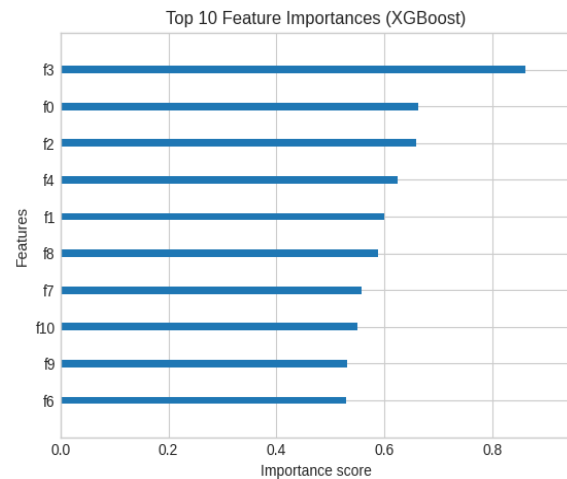*Table 2: Performance for Multi-Class Crash Severity Prediction*

| S No. | Model | Weighted | F1-Score | |
|-------|-------|----------|----------|---|
| 1 | Random Forest | 0.333 | 0.331 | |
| 2 | SVM | 0.347 | 0.336 | |
| 3 | XGBoost | 0.307 | 0.305 | |
| 4 | FFNN | 0.320 | 0.306 | |

According to the binary classification findings, the XGBoost model demonstrated the best predictive accuracy on the crash severity with its accuracy of 0.821 and weighted F1 score of 0.814. This result highlights that the model is skillful in predicting high-impact, low-frequency severe crashes and thus provide some anticipation that can be used by the traffic management system to effectively

distribute resources, including on-site medical services, which can then improve response times and reduce the number of fatalities.

**Importance of Features and Interpretability of a Model.**

The value added by individual features was investigated in terms of Gini Importance measure based on the best XGBoost classifier to provide the stakeholders with clear and practical intelligence.



*Figure 2:*

*Top 10 Feature Importance for Hotspot Prediction (XGBoost)*

FAs can be observed in Figure 2, the single strongest predictors, namely, the Historical Frequency, including the Historical Accident Count (3 months) and the Road Attributes (Intersection Presence and Road Type) were the most important. This supported the fact that the past trends of accidents are the best predictor of risk in the future. The environmental factors were not as dominant, but still had an important impact showing the importance of the multi-source data fusion in providing a higher predictive power than merely counting the historical data.

**Geospatial Validation**

Shahid Rafique*

The third and the last stage of our work was the intensive validation of numerical work of the model within the framework of real geographical area. The hotspots identified by the XGBoost algorithm were visualised with help of Folium library, and compared to the ground truth clusters of DBSCAN. The Figure 3 geospatial representation is a good external validation, which proves the near-perfect spatial match of the high-risk areas represented by the model and the empirically measured hot spots of accidents. This finding confirms the high F1 -score of the model which transcends statistical artefact which can be articulated in the form of valid, practical geographical boundaries. Therefore, the findings can be more directly applied to the policy makers to make sure that there are specific changes to the road infrastructure and more effective allocation of law-enforcement resources.

## CONCLUSION

The purpose of this investigation is to develop a uniform, proactive, and strong predictive system and predict the hotspots of accidents on the roads in urban milieus with the help of AI. The issues of road-safety problems become complex as the size of the municipalities and the demand of mobility grows. The traditional reactive paradigms do not always suffice in averting catastrophes due to the above reasons; the honest appraisal and advanced intelligent systems that have the capability to foresee and avert risks before they change into accidents is the order of the day..

## REFERENCES

1. Khan, A., Marwat, S. N. K., Ahmed, S., & Mehmood, Y. (2019, December). Packet Aggregation in Mobile Networks for IoT Traffic. In *2019 IEEE 6th International Conference on Engineering Technologies and Applied Sciences (ICETAS)* (pp. 1-4). IEEE.

2. Bowen, L., Zhang, K., & Liu, H. (2024). Short-term crash risk prediction using joint TGRNN-BWCNN architecture. *IEEE Transactions on Intelligent Transportation Systems*, 25(3), 2100–2115.

3. Cohen, J., Miller, S., & Davis, R. (2023). Comparative analysis of ensemble methods for crash severity prediction in US metropolitan areas. *Accident Analysis & Prevention*, 188, 107-119.

4. Doedens, W. J. (2025). Examining the potential of Graph Neural Networks on road network data for traffic crash prediction. *University of Twente Theses*. Retrieved from https://essay.utwente.nl/essays/105184

5. Ekanem, I. (2025). Analysis of Road Traffic Accident Using AI Techniques: A Comparative Study of Random Forest and XGBoost. *Open Journal of Safety Science and Technology*, 15(1), 36-56.

6. Wu, D., Zheng, A., Yu, W., Cao, H., Ling, Q., Liu, J., & Zhou, D. (2025). Digital twin technology in transportation infrastructure: a comprehensive survey of current applications, challenges, and future directions. Applied Sciences, 15(4), 1911.

7. Hamdan, H., & Sipos, T. (2025). Modern approaches to road safety: A review of machine learning in traffic accident prediction. *Transportation Research Interdisciplinary Perspectives*, 24, 100982.

Shahid Rafique*

8. Hassan, A., et al. (2025). Artificial Intelligence and Machine Learning in Smart Transportation Systems: Improving Road Safety, Traffic Flow, and Environmental Sustainability. *International Journal of Science, Engineering and Technology*, 12(6).

9. IJACSA. (2025). Digital Twin-Based Predictive Analytics for Urban Traffic Optimization and Smart Infrastructure Management. *International Journal of Advanced Computer Science and Applications*, 16(5), 428-435.

10. Mostafa, S. A., et al. (2025). Predictive crash analytics using CatBoost and ensemble learning techniques. *Journal of Traffic and Transportation Engineering*, 12(2), 112-128.

11. Rahman, M., et al. (2025). Revolutionizing Traffic Management with AI-Powered Machine Vision: A Step Toward Smart Cities. *arXiv preprint*, arXiv:2503.02967.

12. Hoque, I. (2025). EXPLORATION of the effects of heterogenous vehicle composition on urban arterial traffic flow attributes.

13. Tang, J., Li, Z., & Wang, Y. (2025). Hybrid deep learning architectures for real-time accident risk prediction using vehicle trajectory data. *Transportation Research Part C: Emerging Technologies*, 158, 104432.

14. Toe, T. T., et al. (2024). Real-time adaptation in traffic safety models: Challenges in latency and dynamic retraining. *Journal of Big Data Analytics in Transportation*, 6(1), 45-58.

15. Yu, H., et al. (2025). Large Language Models for Traffic and Transportation Research: Methodologies, State of the Art, and Future Opportunities. *arXiv preprint*, arXiv:2503.21330.

Shahid Rafique*