



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** V **Month of publication:** May 2026

DOI: <https://doi.org/10.22214/ijraset.2026.83148>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Crash Severity Prediction Using Machine Learning Models and Explainable Artificial Intelligence Techniques

Olusola Theophilus Faboya

Department of Computing and Information Science, Bamidele Olumilua University of Education, Science and Technology, Ikere-Ekiti, Nigeria

Abstract: Road crashes remain one of the leading causes of mortality globally, with fatal incidents representing a relatively small proportion of the total crashes that occur, yet remain a critically important class in accident datasets. This study developed a machine learning-based framework that incorporated explainable artificial intelligence techniques for crash prediction. A road accident dataset that contained 10,000 accident records and 25 features related to driver, vehicle and environmental conditions was used. Three machine learning models, namely Random Forest, XGBoost and Artificial Neural Networks, were developed and evaluated using accuracy, precision, recall, and F1-score metrics. The results from the experiments indicated that all models achieved strong predictive performance, with XGBoost having the best result, with 99.20% accuracy and 99.10% F1-score value. Shapley Additive explanation analysis was applied to improve model interpretability and to identify the influential predictors of crash severity. The outcome showed that vehicle-related features contributed most significantly to prediction results, while feature interaction effects remained relatively weak. Confusion matrix analysis of the three models indicated poor minority-class identification due to class imbalance within the dataset. This study shows that integrating machine learning with explainable artificial intelligence techniques provides both accuracy and interpretability in crash severity prediction models, which are capable of supporting policymakers with robust road safety management and policy development.

Keywords: Crash Severity, XAI, Machine Learning

I. INTRODUCTION

Road crashes constitute a major public health and safety concern globally, resulting in large human and economic losses each year (Amiri et al., 2025). The international road safety reports confirm that millions of people are injured annually, while road crash fatalities account for a large proportion of preventable deaths (WHO, 2018). In developing countries, including Nigeria, poor road infrastructure, human error, and inadequate enforcement of traffic regulations further increase the concerns for the problem.

The increasing complexity of traffic dynamics and its interdependence variables requires a robust predictive framework to predict explainable solutions to fatalities and optimise intervention strategies (Cicek et al., 2023; Benfaress et al., 2025). Traditional statistical models, including Poisson regression and logistic regression, have been employed in analysing crash severity. But these approaches have often been found to struggle in capturing the complexity and non-linear relationships inherent in multifaceted accident datasets. Machine learning models such as Random Forest, XGBoost, and Artificial Neural Networks have proven to be better with superior predictive capability due to their ability to model nonlinear and high-dimensional datasets (Somvanshi et al., 2026). However, despite the high predictive performance of ML models, limitations come from their inherent opacity in a black-box nature, which limits interpretability and trustworthiness, and hence their practical policy-making utility and real-world safety management (Madushani et al., 2023). To address this gap, integrating explainable artificial intelligence (XAI) techniques such as SHapley Additive exPlanations (SHAP) becomes necessary to interpret model decisions, ensuring that prediction insights are transparent and adequate in explaining and identifying critical contributing factors in crash severity analysis (Aboulola, 2024; Somvanshi et al., 2026).

This study proposes the development of a framework for accident severity prediction that integrates XAI techniques with machine learning models for predictive accuracy and interpretability. The proposed framework will compare the predictive performance of Random Forest, XGBoost, and ANN models, and apply Explainable Artificial Intelligence techniques for model interpretation for the provision of adequate insights for road safety management and policy formulation.

The remaining part of the paper is structured as follows: Section II presents the review of related work, while the study's methodology is discussed in Section III, the results and discussion from the experiments are discussed in Section IV, and Section V provides the conclusion to the study.

II. RELATED WORK

Several previous studies have applied machine learning techniques such as Logistic Regression, Decision Trees, and Neural Networks to predict crash severity. Obasi & Benson (2023) conducted a study that involved the use of Naïve Bayes, Random Forest, Logistic Regression and Artificial Neural Networks to predict accident severity. Xiao et al. (2023) in their study evaluated the efficacy of a Rare Events Logistic Model (RELM) compared to the Logit Model (LM) in enhancing the precision of fatal crash estimations. Mostafa et al. (2025) developed an AI-driven machine learning framework for traffic severity prediction. The study used clustering techniques that include K-means and HDBSCAN, with oversampling methods-SOMTE and ADASYN due to class imbalance. The study showed that the Extra Trees (ET Classifier) ensemble model demonstrated superior performance, achieving 96.19% accuracy and an F1-score of 95.28%, which ensures a well-balanced prediction. Assi et al. (2020) developed machine learning models to predict crash injury severity using fuzzy c-means, which enhanced the predictive capability. The study showed that the SVM-FCM model outperformed the other developed models in terms of accuracy and F1 score, and concluded that the FCM clustering algorithm enhanced the prediction power of FNN and SVM models.

Although previous studies delivered promising prediction accuracy, but the gaps remain that there is a need for transparent and interpretable ML frameworks capable of supporting practical safety decisions and context-specific insights. Many advanced AI models operate as "black boxes", producing predictions without transparent reasoning processes (Dong et al., 2022; Sajid et al., 2024). This lack of interpretability raises concerns regarding trust, clarity and ethical compliance. As a result, Explainable Artificial Intelligence (XAI) emerged as a research field aimed at improving the transparency and interpretability of AI systems (Dong et al., 2022). It is methods and techniques that allow human users to comprehend and trust the outcomes of a machine learning algorithm. The intrinsic interpretable models, such as Linear Regression, Logistic Regression, and Decision Trees, are easy to understand without external explanation tools, while complex models such as Deep Neural Networks (DNNs), Random Forest, XGBoost, and Support Vector Machines often achieve high accuracy but lack transparency in their interpretation.

For instance, a neural network can be represented as:

$$a_j^l = f\left(\sum_i w_{ij}^l a_i^{(l-1)} + b_j^{(l)}\right)$$

Where:

a_j^l = activation of neuron

w_{ij} = weights

b_j = bias

f = activation function

Due to millions of parameters, understanding internal reasoning becomes difficult, and hence the need for a tool to provide explanations for hidden facts. There are many kinds of these tools with varying capabilities. The model-agnostic methods, such as Local Interpretable Model-Agnostic Explanations (LIME), Shapley Additive explanation (SHAP) and partial Dependence Plots (PDP), can explain any machine learning model regardless of architecture. While LIME explains individual predictions by approximating the black-box model locally with a simpler interpretable model, the SHAP is based on cooperative game theory and calculates feature contributions using Shapley values to provide both local and global explanations (Wang, 2024).

The SHAP is a theoretically grounded XAI technique that is represented as:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} [f(S \cup \{i\}) - f(S)]$$

Where:

ϕ_i = condition of feature i

F = feature set

S = subset of features

XAI has become a critical area of modern AI research due to the increasing deployment of intelligent systems in high stake environment, such as crash severity, to improve transparency, fairness, and by making AI systems understandable to humans.

Therefore, recent studies on severity prediction have increasingly focused on applying explainable machine learning methods to transportation safety analysis (Dong et al., 2022; Sajid et al., 2024). Wei et al. (2023), applied an explainable machine learning technique for rural interstate crash severity modelling and observed that XGBoost outperformed other ML algorithms on imbalance crash datasets. Somvanshi et al. (2026) Integrated Feedforward Neural Networks with SHAP analysis for crash severity prediction using naturalistic driving data, and demonstrated the effectiveness of XAI in understanding contributing crash factors. Also, Xiao & Duan (2025) combined SHAP-based feature importance analysis with deep neural networks to propose a framework for multisource crash severity prediction.

III. METHODOLOGY

A. Dataset Description

The dataset obtained consists of road crash records containing 10,000 samples and 25 features, including driver characteristics, vehicle type, road condition, and crash severity labels.

B. Data Preprocessing and Cleaning

The inconsistencies in the datasets were identified and replaced with standardised text data to ensure categories are grouped correctly for counting and visualisation. The affected variables include Vehicle make, Gender, Vehicle type, etc. This is to ensure that categorical groups are correctly identified for counting and visualisation. Each of the categorical columns is checked for unexpected categories. Also, the missing values in the numeric columns were handled using a combination of mean imputation for numerical values and mode imputation for categorical variables. Missing values were found in nine out of the twenty-five features of the datasets. Categorical columns with more than 30% missing values were replaced with a new category, “unknown”; we decided not to remove them to prevent biased data.

Duplicate values in the dataset were removed to avoid double-counting, and outliers were identified using the Interquartile Range (IQR) method and capped to the 5th and 95th percentiles to reduce their influence on model training.

Label encoding was applied to the ordinal features, such as “Time_of_day”, “Day_of_Week”; and the binary features, including “Transmission_Type”, “Crash_Location”, to convert them into simple integers. One-hot encoding was used to map the nominal to numbers. This is to avoid the mistakes of ranking, i.e., vehicle type “Tata” that was marked a “3” might mistakenly be classified as greater than “Honda” labelled as “1”. The resulting shape of the data after encoding has forty-one features.

To prepare for modelling, the dataset was divided into two sets as 80% training data and 20% testing before scaling to prevent data leakage. The numerical features were normalised using Min-Max scaling to ensure uniform feature contribution during model training. Then, MinMaxScaler was used on the data to squeeze all values between 0 and 1.

IV. RESULT AND DISCUSSION

Crash Severity and Drivers' Age

The ages range from 18 to 83 years, indicating that drivers from wide spectrum of age groups were represented in the study. The absence of extreme outliers indicates that the dataset does not represent a particular age group.

The distribution appears to be uniform across all categories; however, slight peaks can be observed around the age ranges of 28-32, 50-55 and 70-75, which suggest that these groups recorded a comparatively higher number of drivers. The Kernel Density Estimation (KDE) curve also supports this observation by showing mild fluctuations rather than sharp spikes.

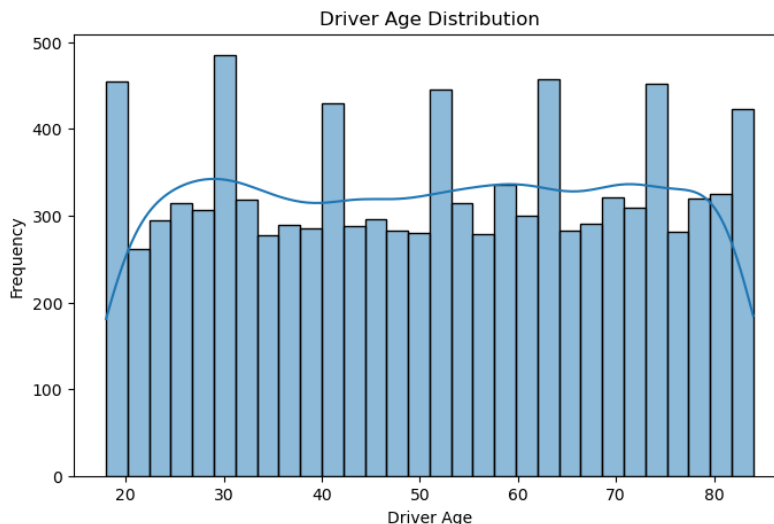


Figure 1: Age Distribution and Crash Frequency

From the perspective of this study, the balanced age distribution is beneficial to the models as it reduces age-related bias that allows the models to learn patterns across diverse driver age groups, hence, improving the generalisation and robustness of the predictive analysis.

Table 1: Machine Learning Model Comparison Table

Model	Accuracy	Precision	Recall	F1-Score
Random Forest	0.9840	0.982045	0.9840	0.982915
XGBoost	0.9920	0.990034	0.9920	0.991014
ANN	0.9845	0.983500	0.9845	0.983999

Table 1 presents the performance comparison of the three classification models on their level of Accuracy, Precision, Recall, and F1-Score metrics. The results from the metrics indicate that all models achieved high predictive performance with accuracy values exceeding 98%, demonstrating the robustness of the proposed classification framework

XGBoost achieves the best result of the three models by recording 99.20% accuracy, a precision of 99.000%, a recall of 99.20% and 99.10% F1-score. The result suggests that XGBoost provided the most balanced and reliable classification capability across all evaluation metrics with the datasets used in this study. The ANN model produced the second-best performance, with an accuracy of 98.45%, a precision of 98.35%, a recall of 98.45%, and an F1-score of 98.40%. It shows that ANN successfully captured complex nonlinear relationships within the dataset. Although ANN has slightly lower performance than XGBoost, the table indicates that ANN still demonstrates excellent classification capability and strong generalisation performance. The Random Forest also demonstrated strong performance, achieving an accuracy of 98.40%, precision of 98.20%, recall of 98.40% and F1-score of 98.29%. The model maintained a good balance between sensitivity and precision; it also exhibited reliable classification ability. When compared with XGBoost and ANN, Random Forest showed slightly lower performance across all evaluation metrics. This could be due to the limitation of bagging-based ensemble learning in handling highly complex feature interactions when compared to boosting techniques such as XGBoost.

The comparative analysis of the three models indicates that all the models are highly effective in predicting crash severity with the dataset used in this study; however, one of the aims of this study is to gain insight into the complex relationships among variables and the interrelated nonlinear features that exist among them. The confusion matrices in Figures 2 to 4, the SHAP interaction summary plot in Figure 5, as well as the SHAP Feature Importance in Figure 6 provide further explanation of the hidden relationships.

A. Models Confusion Matrix Comparison

Figure 2 presents the confusion matrix for the Random Forest, and it summarises the

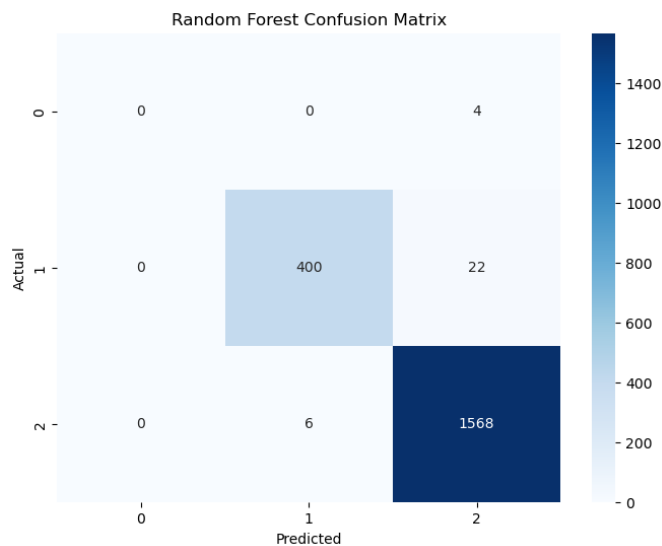


Figure 2: Random Forest Confusion Matrix

classification performance for the Random Forest model across class 0, class 1 and class 2. In the confusion matrix, the rows represent the actual classes, while the columns represent the predicted classes. The correct classification appears along the diagonal, while off-diagonal values indicate miscalculations. The confusion matrix in Figure 2 revealed that the Random Forest classifier performs well for the dominant classes while exhibiting poor sensitivity toward the minority class. It achieves excellent prediction for class 2 with correct identification of 1568 out of 1574 samples, and only 6 instances are misclassified as class 1. It also classified correctly 400 out of 422 instances for class 1, and 22 samples were incorrectly assigned to class 2. However, the model failed to correctly classify any samples belonging to class 0. All 4 class 0 instances were classified as class 2. It shows that the model was unable to learn the distinguishing characteristics of the minority class. This could be due to imbalanced data conditions.

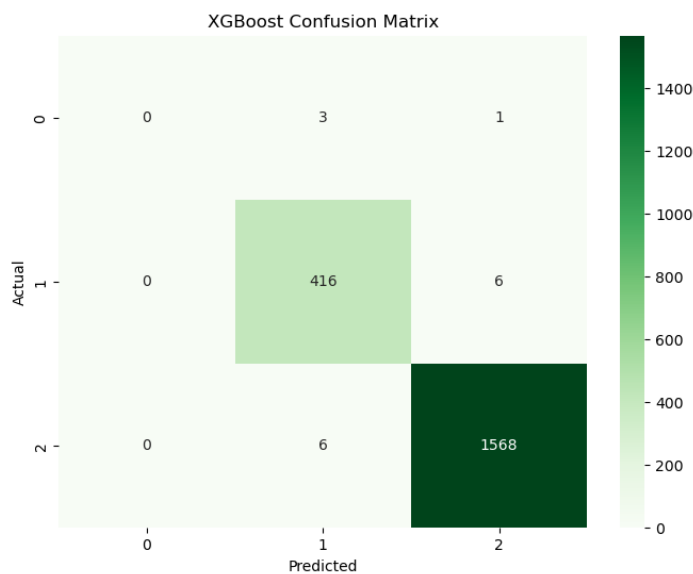


Figure 3: XGBoost Confusion Matrix

The confusion matrix in Figure 3 presents the classification performance of the XGBoost model across three classes (0,1,2). It demonstrated strong predictive performance for class 1 and 2 in particular, with most observations were correctly classified, as evidenced by high values along the diagonal of the confusion matrix. It demonstrated that 416 samples of class 1 and 1568 samples of class 2 were accurately predicted, with only minimal inter-class confusion. However, the model failed to correctly identify any samples belonging to Class 0, suggesting substantial class imbalance within the dataset. Despite an overall accuracy of 99.2% the classifier exhibited poor minority-class sensitivity, which indicates that accuracy alone may not adequately reflect model robustness. This suggests that there is a need for further investigation and the need to employ the use of data imbalance handling strategies.

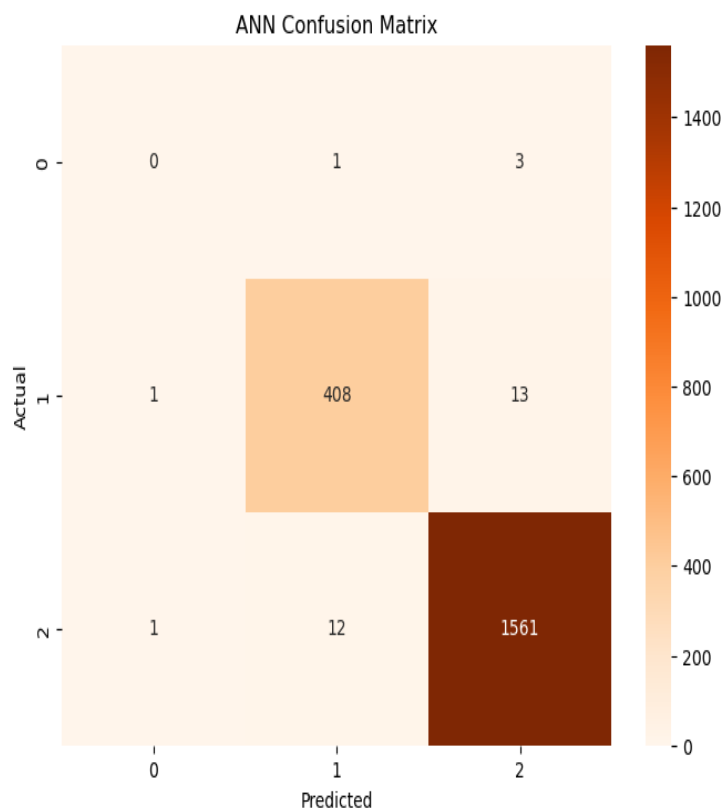


Figure 4: The Confusion Matrix for ANN

The confusion matrix in Figure 4 indicates that the ANN achieved strong predictive capability for classes 1 and 2, while class 0 exhibited weak recognition performance. The model correctly predicted 408 and 1561 instances for classes 1 and 2, respectively. However, the classifier failed to correctly identify samples belonging to Class 0, indicating poor minority-class recognition. The misclassification was observed between Classes 1 and 2, although the error rate remained relatively low. The findings suggest that the ANN effectively learned dominant class patterns but failed in handling imbalanced data distributions. Future improvements may require cost-sensitive learning approaches to enhance minority-class prediction performance.

B. SHAP Interaction Summary Plot

The interaction Summary plot breaks down how features interact with each other to impact the model’s prediction. Figure 5 presents the pairwise interaction effect among three explanatory variables, namely Engine_displacement, Number_of_cylinders, and Vehicle_year. The SHAP interaction analysis revealed that pairwise feature interactions contributed minimally to the predictive behaviour of the model. As observed, the SHAP values for Engine_displacement, Number_of_cylinders, and Vehicle_year were tightly centred around zero, indicating that the model predominantly relied on additive feature effects rather than nonlinear interdependencies.

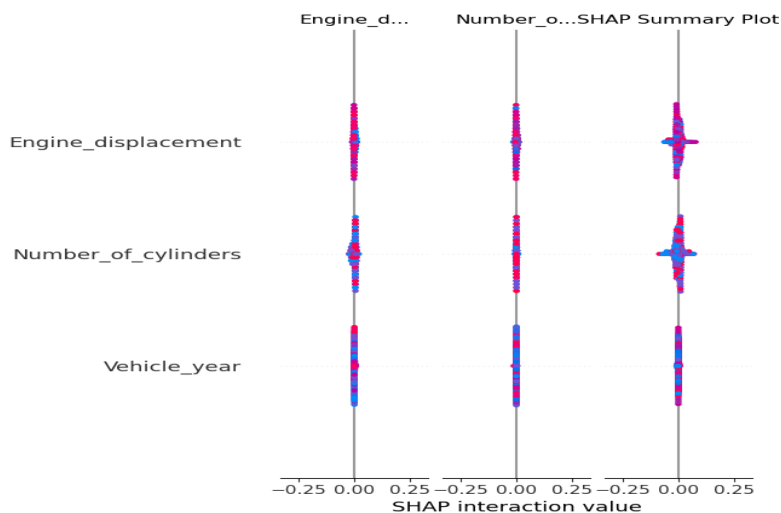


Figure 5: SHAP Interaction Summary Plot

Among the evaluated variables, Number_of_cylinders exhibited marginally greater interaction dispersion, suggesting limited conditional dependence with engine displacement. Also, Vehicle_year displayed negligible interaction influence, implying temporal effects were incorporated independently within the prediction structure. The results suggest that the modelled systems are governed by direct feature contributions, but there are weak pairwise interactions overall. The SHAP interaction further reveals dominance of additive feature contributions and limited nonlinear dependence among predictors. It is therefore interpretable, structurally stable and only mildly dependent on higher-order feature interaction.

C. Shapley Additive exPlanations Feature Importance Plot

The SHAP plot in Figure 6 illustrates the contributions of each variable to the crash severity prediction model. The horizontal bar represents the mean absolute SHAP value, which measures the average magnitude of influence each feature has on the model's predictions. The analysis reveals that the model relies primarily on vehicle-related attributes, especially vehicle manufacturer and engine type, to predict crash severity as indicated by larger SHAP values.

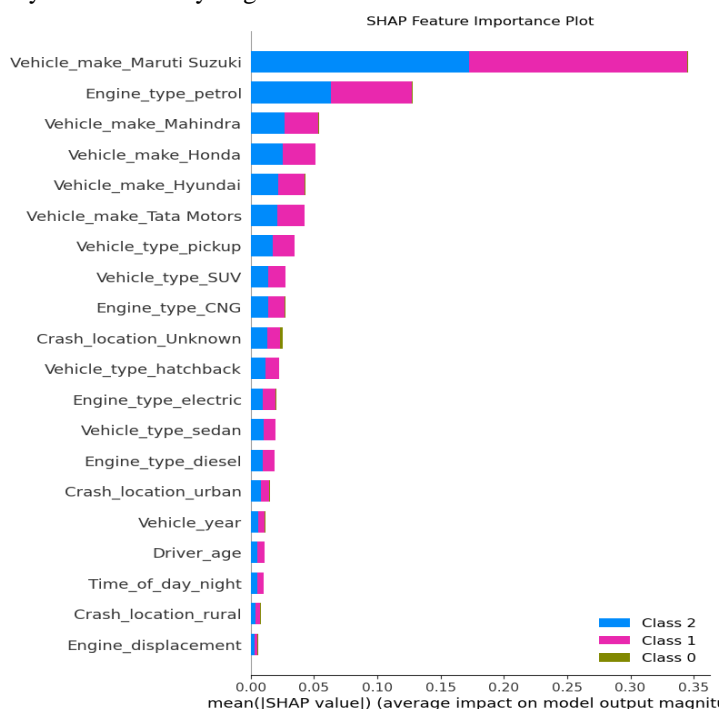


Figure 6: The SHAP Feature Importance

The results in Figure 6 show that Vehicle_make_Maruti Suzuki has the most influential feature in the model behaviour, with a higher SHAP value than any other variable, while environmental, temporal, and driver-related variables contribute less significantly. The dominance of a few features suggests that the crash severity model is highly sensitive to specific vehicle characteristics. This finding could be a useful insight for transportation safety policy, vehicle regulation, and targeted accident prevention strategies.

V. DISCUSSION

The findings from this study demonstrated the performance of the three machine learning models evaluated on crash severity prediction. All the models achieved high predictive capability, with accuracy values exceeding 98.0%, indicating the suitability of the machine learning approaches for crash severity analysis. However, despite the excellent overall evaluation metrics, the confusion matrix analysis revealed an important limitation associated with all the models regarding the imbalanced nature of the dataset. All three models (Random Forest, XGBoost and ANN) failed to correctly classify the minority class (class 0), with most minority instances being classified into the dominant classes. This points to the fact that overall accuracy alone may not be enough to assess a model's robustness when dealing with highly imbalanced datasets. The inability of the models to adequately recognise the minority class suggests the need for advanced imbalance-handling strategies such as Synthetic Minority Oversampling Technique (SMOTE) and Adaptive Synthetic Sampling (ADASTNY), and cost-effective learning or focal loss optimisation.

The application of SHAP-based explainability analysis in this study provided additional insights into the internal behaviour of the developed models. The SHAP feature importance plot revealed that vehicle-related variables contributed significantly to crash severity prediction. Variables such as vehicle manufacturer and engine type exhibited the highest SHAP values, suggesting that vehicle characteristics influence crash severity outcomes. These findings support previous research, such as Dong et al. (2022) and Sajid et al.(2024)that identified vehicle-related attributes as critical determinants of accident severity. In contrast, environmental and temporal factors contributed insignificantly to the prediction, which implies that the predictive framework relied more on intrinsic vehicle characteristics than on external driving conditions. Furthermore, the SHAP interaction analysis demonstrated that pairwise feature interactions contributed minimally to the predictive structure of the model.

The SHAP interaction values were concentrated around zero for most variables, which indicates that the models primarily relied on additive features' contributions rather than strong nonlinear feature dependencies. This finding suggests that the framework developed in this work possesses structural stability and interpretability that can improve transparency in model decision-making. The findings also demonstrate that combining advanced machine learning algorithms with explainable AI techniques can significantly improve both predictive accuracy and interpretability in crash severity modelling.

VI. CONCLUSION

The framework developed in this study evaluated three machine learning models, namely Random Forest, XGBoost, and Artificial Neural Networks, integrated with explainable Artificial Intelligence (XAI) techniques. The three models were trained and evaluated using a road accident dataset containing 10,000 crash records and multiple explanatory variables related to driver, vehicle and environmental conditions. The experiments showed that the three models achieved high predictive performance, with XGBoost performing best with an accuracy of 99.20% and the highest precision, recall, and F1-score values. The findings confirm the effectiveness of boosting-based ensemble learning techniques in modelling complex and nonlinear relationships. ANN and Random Forest also maintained reliable classification performance across evaluation metrics. However, the integration of SHAP-based explainability techniques enhanced the transparency and interpretability of the predictive framework by identifying the most influential variables contributing to crash severity outcomes. This addresses one of the major limitations of traditional black-box machine learning models by enabling interpretable prediction analysis for transportation safety applications.

This study demonstrates that integrating machine learning models with explainable artificial intelligence techniques provides an effective and interpretable framework for crash severity prediction. And offer valuable insights for transportation safety. However, future studies may improve the framework by incorporating real-time traffic data and addressing the limitation of imbalanced datasets with relevant machine learning techniques.

REFERENCES

- [1] Aboulola, O. I. (2024). Improving traffic accident severity prediction using MobileNet transfer learning model and SHAP XAI technique. *PLoS ONE*, 19(4 April), 1–18. <https://doi.org/10.1371/journal.pone.0300640>
- [2] Amiri, M. A., Afshari, S., & Soltani, A. (2025). Machine learning approaches to traffic accident severity prediction: Addressing class imbalance. *Machine Learning with Applications*, 22(November), 100792. <https://doi.org/10.1016/j.mlwa.2025.100792>



- [3] Assi, K., Rahman, S. M., Mansoor, U., & Ratrou, N. (2020). Predicting crash injury severity with machine learning algorithm synergized with clustering technique: A promising protocol. *International Journal of Environmental Research and Public Health*, 17(15), 1–17. <https://doi.org/10.3390/ijerph17155497>
- [4] Benfaress, I., Bouhoute, A., & Zinedine, A. (2025). Improving Intelligent Systems with Explainable AI for Early Prediction and Analysis of Traffic Accidents. 8th International Conference on Networking, Intelligent Systems & Security (NISS), 85–92. <https://doi.org/https://doi.org/10.1109/NISS66502.2025.00021>
- [5] Cicek, E., Akin, M., Uysal, F., & Topcu Aytas, R. M. (2023). Comparison of traffic accident injury severity prediction models with explainable machine learning. *Transportation Letters*, 15(9), 1043–1054. <https://doi.org/10.1080/19427867.2023.2214758>
- [6] Dong, S., Khattak, A., Ullah, I., & Zhou, J. (2022). Predicting and Analyzing Road Traffic Injury Severity Using Boosting-Based Ensemble Learning Models with SHAPley Additive exPlanations.
- [7] Madushani, J. P. S. S., Sandamal, R. M. K., Meddage, D. P. P., Pasindu, H. R., & Gomes, P. I. A. (2023). Evaluating expressway traffic crash severity by using logistic regression and explainable & supervised machine learning classifiers. *Transportation Engineering*, 13(June), 100190. <https://doi.org/10.1016/j.treng.2023.100190>
- [8] Mostafa, A. M., Aldughayfiq, B., Tarek, M., Alaerjan, A. S., Allahem, H., Elbashir, M. K., Ezz, M., & Hamouda, E. (2025). AI-based prediction of traffic crash severity for improving road safety and transportation efficiency. *Scientific Reports*, 15(1), 1–24. <https://doi.org/10.1038/s41598-025-10970-7>
- [9] Obasi, I. C., & Benson, C. (2023). Evaluating the effectiveness of machine learning techniques in forecasting the severity of traffic accidents. *Heliyon*, 9(8), e18812. <https://doi.org/10.1016/j.heliyon.2023.e18812>
- [10] Sajid, A., Jalayer, M., Das, S., & Bin, A. (2024). International Journal of Transportation Application of machine learning models and SHAP to examine crashes involving young drivers in New Jersey. *International Journal of Transportation Science and Technology*, 14, 156–170. <https://doi.org/10.1016/j.ijtst.2023.04.005>
- [11] Somvanshi, S., Liu, J., Chakraborty, R., Tamakloe, R., & Das, S. (2026). Predicting Crash Severity using Naturalistic Driving Data and Neural Networks. *International Journal of Intelligent Transportation Systems Research*, 2023. <https://doi.org/10.1007/s13177-025-00624-3>
- [12] Wang, Y. (2024). A Comparative Analysis of Model Agnostic Techniques for Explainable Artificial Intelligence. *Research Reports on Computer Science*, 3(2), 25–33. <https://doi.org/10.37256/rrcs.3220244750>
- [13] Wei, Z., Zhang, Y., & Das, S. (2023). Applying Explainable Machine Learning Techniques in Daily Crash Occurrence and Severity Modelling for Rural Interstates. *Transportation Research Record: Journal of the Transportation Research Board*, 2677(5), 611–628.
- [14] WHO. (2018). *Global Status Report on Road Safety 2018 (Issue 1)*. <https://www.who.int/publications/i/item/9789241565684>
- [15] Xiao, Y., & Duan, Z. (2025). An explainable multi-task deep learning framework for crash severity prediction using multi-source data. *Scientific Reports*, 15(1), 1–20. <https://doi.org/10.1038/s41598-025-09226-1>
- [16] Xiao, Y., Lin, L., Zhou, H., Tan, Q., Wang, J., Yang, Y., & Xu, Z. (2023). Fatal crashes and rare events logistic regression: an exploratory empirical study. *Frontiers in Public Health*, 11. <https://doi.org/10.3389/fpubh.2023.1294338>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)