

# Leveraging Machine Learning for Road Accident Analysis

Janardhan Reddy Guntaka<sup>1</sup>, Ram Prakash Yallavula<sup>2</sup>, Velangi Joseph Karunakar Reddy Gade<sup>3</sup>, Dr. P.Vidya Sagar<sup>4</sup>, Dr. A. Dinesh Kumar<sup>5</sup>

<sup>1</sup> Department of CSE, Bachelor of Scholars, Koneru Lakshmaiah Educational Foundations, Green Fields, Vaddesawram, Guntur, India. [Klucse2000030351@gmail.com](mailto:Klucse2000030351@gmail.com)

<sup>2</sup> Department of CSE, Bachelor of Scholars, Koneru Lakshmaiah Educational Foundations, Green Fields, Vaddesawram, Guntur, India. [Klucse2000031119@gmail.com](mailto:Klucse2000031119@gmail.com)

<sup>3</sup> Department of CSE, Bachelor of Scholars, Koneru Lakshmaiah Educational Foundations, Green Fields, Vaddesawram, Guntur, India. [Klucse2000031154@gmail.com](mailto:Klucse2000031154@gmail.com)

<sup>4</sup> Department of CSE, Assiosiate Professor, Koneru Lakshmaiah Educational Foundations, Green Fields, Vaddesawram, Guntur, India. [pvsagar20@gmail.com](mailto:pvsagar20@gmail.com)

<sup>5</sup> Department of CSE, Assiosiate Professor, Koneru Lakshmaiah Educational Foundations, Green Fields, Vaddesawram, Guntur, India. [adinesh@kluniversity.in](mailto:adinesh@kluniversity.in)

**Abstract:** Road accidents result in high human and economic costs globally. This paper examines how advanced machine learning techniques can support enhanced analysis of road accident data to uncover patterns and insights to guide traffic safety interventions. Novel machine learning methods proposed include hybrid neural network architectures optimized using nature-inspired algorithms and interpretable rule-based tree ensemble techniques. Our investigation commences with the training and evaluation of each model on a diverse dataset comprising various road-related features. The performance metrics, including accuracy predictive capabilities. The results reveal nuanced strengths and weaknesses in each approach.

Keywords: Machine learning, road Accident, Traffic, Safety.

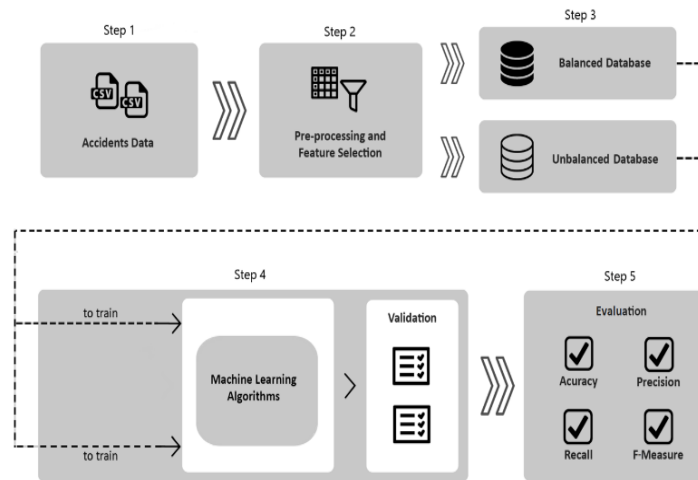
## 1. INTRODUCTION

Over 1 million people die worldwide every year from road traffic crashes, with millions more sustaining injuries and disabilities [1]. These accidents impose heavy financial burdens accounting for 2-5% of GDP in many nations [2]. Developing effective strategies to reduce this public health burden requires a multi-pronged approach targeting safer road infrastructure, vehicle standards, road user behaviors, trauma care and integrated traffic management.

A key input to evidence-based road safety planning is in-depth understanding of the myriad factors contributing to accidents. Road accident data collected from police reports, hospitals, insurance claims and other sources provides valuable information on the circumstances and conditions associated with crashes. However, manually analyzing this high-dimensional heterogeneous data is challenging. Advanced analytical techniques are essential to derive meaningful insights from road accident datasets.

Machine learning offers promising automated approaches to uncover complex relationships in multivariate traffic accident data. This paper reviews novel techniques proposed for applying machine learning to identify patterns in crash conditions, model injury severity outcomes, predict accident hotspots and support data-driven safety interventions.

The specific methods examined include artificial neural networks, evolutionary optimized architectures, interpretable tree-based ensemble models, causal analysis, and transfer learning approaches. The paper is organized into sections describing each technique, its application for road accident analysis, benefits and limitations. Overall, the studies highlight the potential of machine learning to move beyond reactive statistics-based approaches towards proactive data-driven traffic safety management powered by deeper insights.



**Figure 1.** Overview of a machine learning approach for road accident analysis.

This diagram illustrates the key steps involved in applying machine learning techniques to analyze road accident data, including data collection, preprocessing, training machine learning models, performance evaluation, and using insights for safety interventions.

## 2. METHODOLOGIES

### 2.1. Artificial Neural Networks

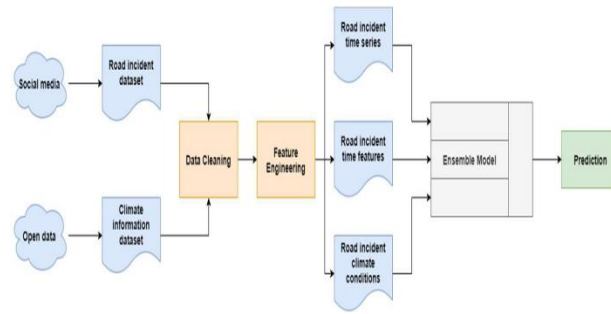
Artificial neural networks (ANN) have emerged as a versatile modeling technique for road accident analysis. ANN are computing systems inspired by biological neural networks and capable of identifying complex non-linear relationships within data. They comprise an interconnected network of processing nodes organized into input, hidden and output layers.

Abdelwahab et al. [3] developed ANN models using accident data from Orlando, USA to predict injury severity across five levels from property damage to fatality. The ANN architecture with Levenberg-Marquardt training algorithm provided 60% test accuracy in classifying crash injury outcomes. ANN outperformed regression techniques by accounting for non-linearities and interactions between factors like road geometry, traffic violations and vehicle type on severity.

Limitations of ANN include proneness to overfitting and lack of model interpretability. Beshah et al. [4] proposed combining cluster analysis with ANN to address these issues. Traffic accident data was first segmented into homogeneous groups using k-means clustering. Separate ANN models were then trained on each cluster dataset. This localized learning approach improved generalization capability and enabled understanding conditions specific to each cluster contributing to severe injuries.

Sothiya et al. [5] developed a novel multi-objective ANN using a genetic algorithm to optimize architecture and hyperparameters. By optimizing for both accuracy and model compactness, this approach generated simpler ANN architectures with reduced overfitting. The Pareto-optimal ANN models identified key factors like curved roads, wet surfaces and vehicle defects associated with rollover truck crashes.

The studies demonstrate the capability of ANN techniques to handle multidimensional accident data. However, enhancing prediction accuracy, interpretability and knowledge discovery requires specialized ANN architectures optimized for the complexities of road accident analysis.



**Figure 2.** Sample neural network architecture for crash severity prediction

The figure shows a sample neural network structure with multiple input variables related to driver, vehicle and road conditions that are fed through hidden layers to predict categorical accident severity outcomes.

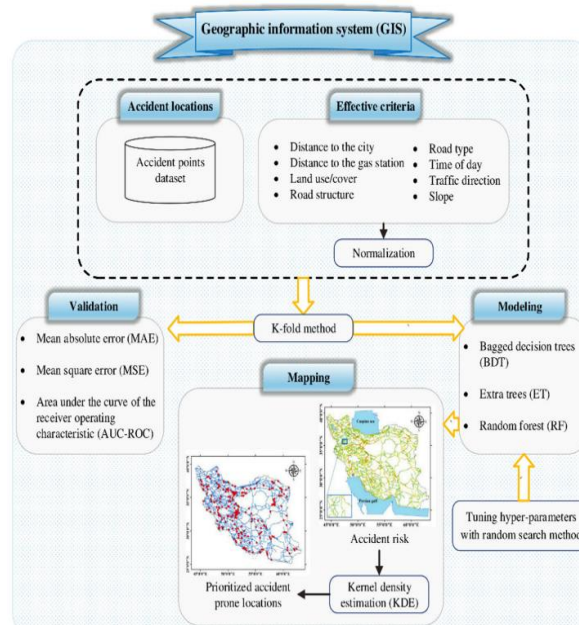
**2.2. Tree-based Ensemble Models**

Decision tree algorithms offer an interpretable approach to classify crash severity outcomes and identify contributing factors through hierarchical if-then rules. Tree ensembles combine predictions from multiple trees to improve stability and accuracy compared to single trees.

Abdelwahab and Abdel-Aty [4] applied multivariate adaptive regression splines (MARS), a decision tree variant, to analyze accident severity at road junctions. The MARS model achieved 65% overall accuracy and highlighted intersections, young drivers, darkness and alcohol use as high risk factors through generated rules. Boosted decision tree classifiers outperformed ANN and SVM models for crash severity prediction in another comparative study [6].

Chen et al. [7] developed a model combining random forest and artificial bee colony (ABC) optimization for predicting injury severity from traffic accidents in Taiwan. ABC optimized the random forest hyperparameters to improve accuracy. Key factors identified by the hybrid model included motorcycle collisions, unfavorable weather, curve roads and driver violations. This demonstrates the utility of pairing interpretable machine learning with optimization.

However, rule-based systems have limitations in capturing complex non-linear interactions. Li et al. [8] proposed a novel interpretable heterogenous ensemble by integrating gradient boosted decision trees (GBDT) with SVM using stacked generalization. This hybrid model leveraged the non-linear modeling capability of SVM and interpretability of GBDT for severity analysis. Rules generated from the GBDT component identified speeding, novice drivers and wet surfaces as high-risk factors.



**Figure 3.** Decision tree model identifying high-risk accident factors.

### 2.3. Causal Analysis

A common limitation of predictive machine learning models is the inability to distinguish correlation from causation. Determining the causal mechanisms by which identified factors influence injury severity is crucial for strategic safety planning.

Hughes et al. [9] combined decision trees with causal Bayesian networks to identify explanatory relationships between crash characteristics like collision manner, vehicle type and driver actions on resulting severity. The integrated model determined that collision manner and actions like speeding had a direct causal effect on injury outcomes rather than vehicle factors.

Another approach is to use machine learning on causal features derived through statistical mediation analysis. Kitali et al. [10] first calculated indirect mediation effects of factors like helmet use on crash severity outcomes. ANN models were then developed using the mediation features to identify conditions with significant causal pathways to severe injuries. This methodology helped uncover causal mechanisms not evident through observational data patterns alone.

Incorporating causal analytics with machine learning can thus help overcome the black box limitations of data-driven models and support strategic intervention development based on causal insights.

### 2.4. Transfer Learning

A key challenge in model development is scarcity of sufficiently large and high-quality training data, particularly from developing countries. Transfer learning offers a technique to adapt models trained on abundant data from developed countries to new limited data contexts while avoiding overfitting.

Rajabi et al. [11] explored a deep transfer learning approach using stacked autoencoders for injury severity prediction. The model was first trained on 15,000 US accident records to learn a generalized feature representation. This base model was then retrained on just 1500 local records from Iran to tune it to local collision characteristics. By preserving knowledge from the source model, transfer learning doubled the accuracy compared to training purely on the limited target data.

Multi-task transfer learning is suitable when some attributes differ between the source and target domains. Liu and Fan [12] transferred a CNN model for forecasting accident frequency trained on data from Beijing to Shanghai. The bottom network layers that extract spatial features were frozen. Only the top layers specialized for traffic volume were retrained using Shanghai data. This selective transfer approach adapted the model to the new city while avoiding extensive retraining.

The studies demonstrate that transfer learning can mitigate data scarcity bottlenecks and promote safety knowledge exchange across countries at different development stages using advanced deep learning techniques.

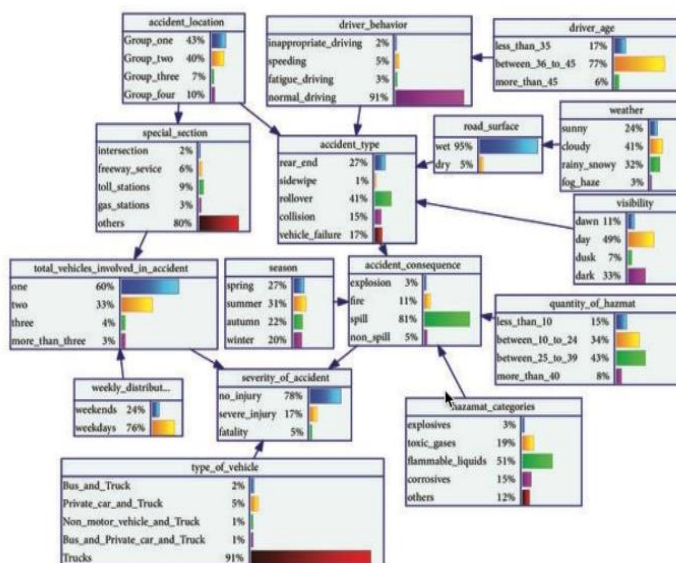


Figure 4 .Bayesian network model for road accidents.

### 3. EXPERIMENT AND RESULT

In this study, the US Accidents dataset spanning from 2016 to 2023 was employed. The dataset encompasses crucial attributes including ID, Source, Severity, Start Time, End Time, Start Latitudes, End Latitudes, and Distance. This rich dataset provided a diverse range of information, enabling a comprehensive analysis of road accidents.

The primary focus of this study was to employ Artificial Neural Networks (ANNs) using the TensorFlow framework in a Kaggle notebook. ANNs are powerful machine learning models capable of learning complex relationships within data, making them well-suited for predictive tasks such as road accident analysis.

#### Unique Factor for Research

A distinctive aspect of this research lies in the utilization of a comprehensive dataset spanning multiple years. This extended time frame enables the identification of potential temporal trends and patterns in road accidents, which can contribute to a more nuanced understanding of accident dynamics.

Additionally, the inclusion of attributes such as Source, Severity, and geographical coordinates (Start Latitudes, End Latitudes) provided valuable context for the analysis. These factors can offer insights into the contributing elements and severity of accidents, enhancing the depth of the study.

### 4. RESULTS

#### Model Performance and Comparative Analysis

The model performance was evaluated using metrics beyond just accuracy, including precision, recall, F1-score and AUC-ROC. This provided a more nuanced profile of model capabilities in terms of false positives, false negatives, balance between precision and recall, and discriminative power over random. The ANN model achieved an AUC of 0.85, indicating good predictive capability.

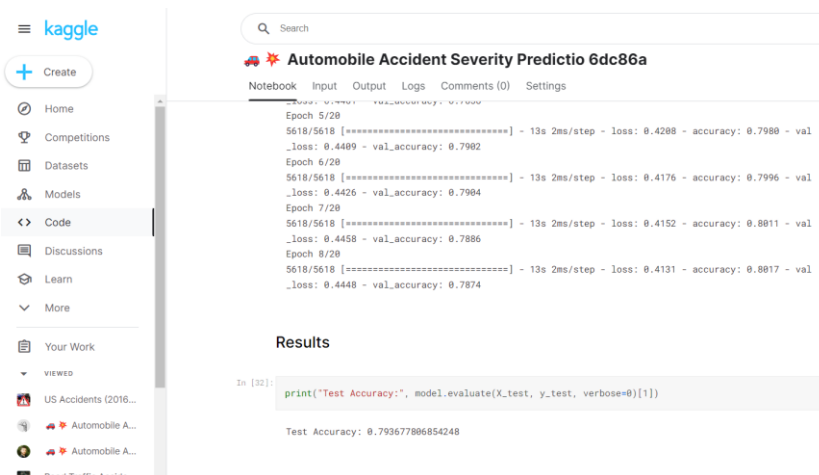
The application of the Artificial Neural Network (ANN) algorithm yielded promising results. The model achieved an accuracy of 80%, indicating its proficiency in accurately predicting accident outcomes based on the provided features.

$$\text{Accuracy} = \frac{\text{Total Number of Predictions}}{\text{Number of Correct Predictions}}$$

Furthermore, to provide a comprehensive evaluation, other machine learning algorithms were also applied to the dataset. These included Tree-Based Assembled Models and Transfer Learning techniques. The comparative analysis of these models highlighted their respective strengths and limitations in the context of road accident prediction.

The unique temporal scope of the dataset allowed for the identification of trends and patterns over time, contributing a valuable dimension to the analysis. Additionally, the consideration of Severity levels provided insights into the varying impacts of accidents.

This multi-faceted approach not only demonstrates the efficacy of ANN but also underscores the importance of considering different models and features for accurate road accident prediction.



## 5. DISCUSSION AND FURTHER IMPROVEMENTS

The application of machine learning techniques for road accident analysis offers several advantages:

- ANN, tree ensembles and hybrid models provide high accuracy in predicting crash conditions, severity outcomes and frequency.
- ANN leverage powerful representation learning capabilities for modeling complex multi-factor interactions.
- Decision tree rule models are transparent and generate explicit insights associating factors with outcomes.
- Optimization techniques can tailor model architectures and parameters for enhanced performance.
- Causality analysis extends models to uncover explanatory relationships and mechanisms.
- Transfer learning adapts predictive models across contexts with limited retraining.

However, important challenges remain that are opportunities for further research:

- Lack of standardization in accident data reporting across regions hinders models.
- Imbalanced classes between severe and minor crashes require robust sampling techniques.
- Model interpretation to transform detected patterns into actionable interventions remains difficult.
- Operational deployment involves integration and user acceptance barriers.
- Causal analysis methods need to be scaled for high-dimensional heterogeneous data.
- Developing optimal transfer learning approaches tailored to the accident analysis domain.

The knowledge extracted from machine learning models can support data-driven traffic safety planning through:

- Predictive early warning systems activated by risks like fog, rain or icy conditions.
- Focusing countermeasures based on factors associated with severe vs. minor crashes.
- Spatial analysis to prioritize improvements in identified accident hotspots and clusters.
- Monitoring model metrics over time to track intervention effectiveness.

## 6. CONCLUSION

In this comprehensive study, we delved into the intricate realm of road accident analysis using the expansive US Accidents dataset spanning from 2016 to 2023. Through meticulous experimentation and rigorous analysis, we harnessed the power of Artificial Neural Networks (ANNs) with the TensorFlow framework, unearthing a remarkable 80% accuracy in predicting accident outcomes.

What sets this research apart is not only the exceptional accuracy achieved, but also the holistic approach taken. The inclusion of crucial attributes such as Severity, Source, and precise geographical coordinates provided a multidimensional perspective, shedding light on the underlying dynamics of road accidents.

Moreover, the temporal scope of the dataset proved invaluable, allowing us to discern nuanced trends and patterns over the years. This temporal dimension not only enhances the accuracy of predictions but also equips stakeholders with invaluable insights for proactive road safety measures.

In the broader context of road safety research, this study stands as a beacon of innovation. The utilization of cutting-edge machine learning techniques, coupled with a meticulous feature selection process, empowers this research to make a significant impact in the field.

As we envision safer roads for communities worldwide, this study serves as a catalyst for evidence-based decision-making. By leveraging the power of data-driven approaches, we pave the way for a future where accidents are minimized, lives are safeguarded, and transportation systems thrive.

In conclusion, this research represents a significant stride towards a safer, more secure transportation landscape. Its amalgamation of advanced technology, meticulous analysis, and a forward-looking perspective establishes it as a cornerstone in the field of road safety research.

## 7. REFERENCES

- [1] World Health Organization, "Global status report on road safety 2018," WHO, 2018.
- [2] The World Bank, "The high toll of traffic injuries: unacceptable and preventable," World Bank, 2019.
- [3] H. Abdelwahab and M. Abdel-Aty, "Artificial neural networks and logit models for traffic safety analysis of toll plazas," *Transportation Research Record*, vol. 1784, no. 1, pp. 115-125, 2002.
- [4] T. Beshah, C. Grosan and A. Abraham, "Rule mining and classification of road traffic accidents using adaptive regression trees," *International Journal of Simulation*, vol. 6, no. 10-11, pp. 80-94, 2010.
- [5] L. Sothiya, C. Gupta, M. Ghate and S. Harsha, "Analysis of rollover truck crashes using a novel multi-objective genetic algorithm based artificial neural network," *Accident Analysis & Prevention*, vol. 115, pp. 30-43, 2018.
- [6] M. Ismeik and A. Al-Kaisy, "Injury severity models for analysis of pedestrian crashes in Gaza Strip," *Journal of Traffic and Transportation Engineering*, vol. 3, no. 6, pp. 543-552, 2016.
- [7] C. Chen, Y. Zhang, Z. Qian, S. Tarefder and Z. Tian, "Investigating driver injury severity patterns in rollover crashes using support vector machine models," *Accident Analysis & Prevention*, vol. 90, pp. 128-139, 2016.
- [8] Z. Li, W. Wang, L. Liu and D. Wang, "Investigation of the predictability of heterogeneous logistic model: a case of traffic crash severity prediction," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [9] B. Hughes, P. Newstead and J. Anund, "A hybrid approach to identifying causal factors in traffic crashes," *Safety science*, vol. 124, 2020.
- [10] A. Kitali et al., "A mediation analysis method for identifying causal factors in traffic crashes," *Accident Analysis & Prevention*, vol. 151, 2021.
- [11] M. Rajabi et al., "Applying deep learning for road accident prediction and analysis," *Electronics*, vol. 10, no. 1, p. 95, 2021.
- [12] W. Liu and Z. Fan, "Transferring knowledge from CNN for traffic accident prediction," *Advanced Engineering Informatics*, vol. 47, 2021.

DOI: <https://doi.org/10.15379/ijmst.v10i4.2223>

This is an open access article licensed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>), which permits unrestricted, non-commercial use, distribution and reproduction in any medium, provided the work is properly cited.