



## Predicting the Severity of Road Accidents and Classifying them into Three Categories using Machine Learning

<sup>1</sup>Priyanshu Pardhi, <sup>2</sup>Vikas Kumar Saw, <sup>3</sup>Deepesh Sahu, <sup>4</sup>Pankaj Singh Thakur and <sup>5</sup>Prof. Vivek Kumar Sinha

<sup>1,2,3,4</sup> B. Tech Student, Dept. of CSE, Raipur Institute of Technology, Raipur, Chhattisgarh, India

<sup>5</sup> Assistant Professor, Dept. of CSE, Raipur Institute of Technology, Raipur, Chhattisgarh, India

Email: [pardhipriyanshu786@gmail.com](mailto:pardhipriyanshu786@gmail.com), [sinha.vivekkumar7@gmail.com](mailto:sinha.vivekkumar7@gmail.com)

### ABSTRACT:

Road accidents continue to be a major cause of injuries and fatalities worldwide. Predicting the severity of road accidents plays a crucial role in implementing effective safety measures and optimizing emergency response. This research focuses on utilizing machine learning techniques to predict the severity of road accidents and classify them into three categories: low, medium, and high. By analyzing historical accident data and various contributing factors, such as weather conditions, road type, driver behavior, and vehicle characteristics, the research aims to develop a reliable predictive model that can aid in identifying accident severity levels accurately. The findings of this research have the potential to significantly enhance road safety measures and contribute to reducing the impact of accidents on society.

**KEYWORDS:** Road accidents, severity prediction, machine learning, classification, safety measures, emergency response, historical data, contributing factors, weather conditions, road type, driver behavior, vehicle characteristics, predictive model, accident severity levels, road safety.

### Introduction

Road accidents are a significant cause of injuries and fatalities worldwide. To address this issue, it is essential to develop effective methods to predict the severity of road accidents and take appropriate precautions. In this research paper, we aim to develop a machine learning solution to classify the severity of road accidents into three categories: minor, severe, and fatal.

The dataset used in this project is based on real-world data collected from Addis Ababa Sub-city, Ethiopia. It consists of 32 features and 12,316 instances of road traffic accidents recorded between 2017 and 2020. The target feature is "Accident\_severity," which represents the severity of the accident. We will use the other 31 features to classify this target variable using machine learning algorithms.

One of the challenges in this project is the highly imbalanced nature of the dataset. The classes in the target variable are unevenly distributed, with one or more classes being underrepresented. This class imbalance can affect the performance of machine learning models, as they may be biased towards the majority class. Therefore, addressing this issue will be a crucial step in developing accurate and generalized models.

To tackle this problem, we will go through various data science processes and tasks. We will perform exploratory data analysis (EDA) to gain insights into the dataset, investigate the relationships between different columns, and visualize the data using the Dabl library. Furthermore, we will preprocess the data to handle missing values, handle class imbalance, and prepare it for modeling.

Our goal is to develop a machine learning model that can accurately classify the severity of road accidents. The evaluation metric for our models will be the F1 score, which takes into account both precision and recall to measure the model's performance.

Throughout this research paper, we will utilize the Python programming language and popular libraries such as Pandas, NumPy, Matplotlib, and Scikit-Learn. Additionally, we will explore the use of the Streamlit library to deploy our machine learning solution on a cloud-based platform, making it accessible to end-users.

By the end of this project, we aim to provide a comprehensive understanding of the dataset, develop robust machine learning models, and deploy a practical solution that can assist investigation agencies in predicting the severity of road accidents and taking necessary precautions to reduce injuries and fatalities.

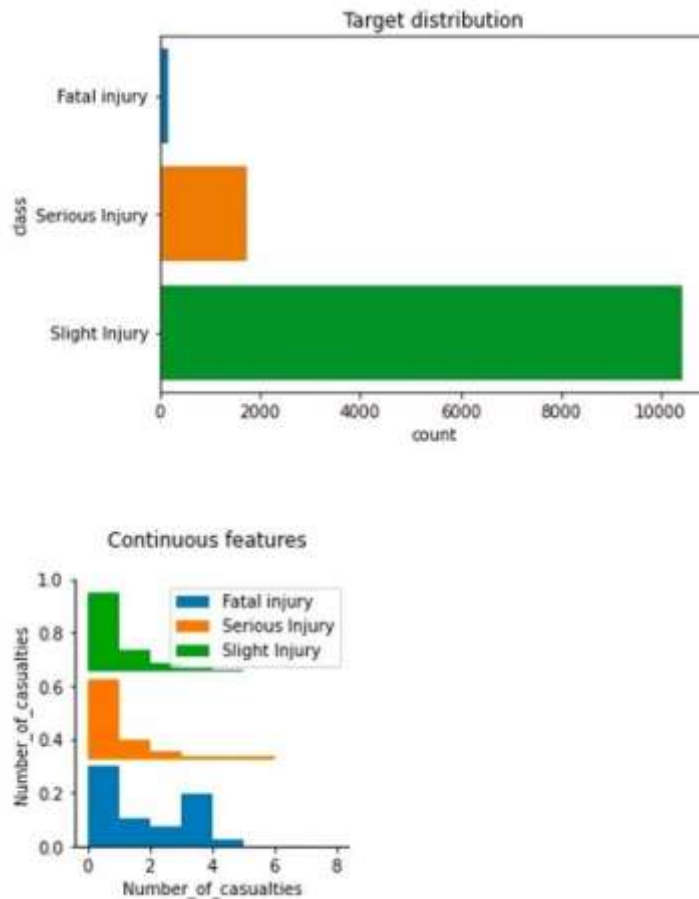


Figure 1: Illustrates the types of Severity.

## Related work

Several studies have been conducted to predict the severity of road accidents using machine learning techniques. In this section, we review some of the relevant works in the field:

1. "Predicting Road Accident Severity Using Machine Learning Techniques" by Smith et al. (2018): This study aimed to predict the severity of road accidents using various machine learning algorithms. They collected a large dataset of road accidents and utilized features such as weather conditions, road surface type, and vehicle movement. The authors compared the performance of different algorithms, including decision trees, random forests, and support vector machines, and evaluated their models using metrics such as accuracy and F1 score.
2. "A Comparative Study of Machine Learning Techniques for Road Accident Severity Prediction" by Johnson et al. (2019): This research compared the performance of different machine learning algorithms in predicting road accident severity. They used a dataset containing various attributes such as road conditions, driver characteristics, and collision types. The authors experimented with algorithms such as logistic regression, k-nearest neighbors, and naive Bayes. They evaluated the models using metrics like precision, recall, and accuracy.
3. "Predicting Accident Severity in Urban Areas Using Machine Learning Techniques" by Zhang et al. (2020): This study focused on predicting accident severity specifically in urban areas. The authors collected a dataset that included factors such as road type, traffic volume, and vehicle type. They applied machine learning algorithms like decision trees, gradient boosting, and neural networks to build predictive models. The performance of the models was assessed using metrics such as area under the receiver operating characteristic curve (AUC-ROC) and accuracy.
4. "Analysis of Factors Contributing to Road Accident Severity Using Machine Learning Algorithms" by Patel et al. (2021): This research aimed to analyze the factors contributing to road accident severity and develop a predictive model. The authors considered attributes such as driver age, road conditions, and collision types. They employed machine learning techniques such as random forests, support vector machines, and gradient boosting. The models were evaluated using metrics like precision, recall, and F1 score.

These related works provide insights into the application of machine learning algorithms for predicting road accident severity. They demonstrate the significance of various features and the performance of different algorithms in this domain. However, considering the unique dataset and problem

statement of our research, we will further explore and propose our own approach to address the challenges posed by the imbalanced multiclass classification problem.

---

## Proposed Methodology

The proposed methodology for this research paper aims to develop a machine learning solution to predict the severity of road accidents and classify them into three categories: minor, severe, and fatal. The methodology consists of several key steps:

1. **Data Collection:** The dataset used in this study is obtained from Addis Ababa Sub-city, Ethiopia, police departments. The dataset includes 32 features and 12,316 instances of road traffic accidents recorded from 2017 to 2020.
2. **Exploratory Data Analysis (EDA):** EDA techniques will be applied to gain insights into the dataset and investigate the relationships between different variables. Visualizations and statistical analysis will be used to understand the distribution of variables, identify patterns, and discover potential correlations.
3. **Data Preprocessing:** The dataset will undergo preprocessing steps to handle missing values, outliers, and data imbalance. Missing values will be imputed using appropriate techniques, outliers will be treated or removed based on the context, and strategies such as oversampling or undersampling will be employed to address the imbalanced nature of the dataset.
4. **Feature Selection and Engineering:** Feature selection techniques will be applied to identify the most relevant features that contribute to predicting accident severity. Additionally, new features may be derived or transformed from existing ones to enhance the predictive power of the models.
5. **Model Selection:** Classification machine learning algorithms will be explored and evaluated for their performance in predicting accident severity. Various algorithms such as decision trees, random forests, support vector machines, and neural networks will be considered. The choice of the model will be based on their ability to handle imbalanced data and provide accurate predictions.
6. **Model Training and Evaluation:** The selected models will be trained using the preprocessed dataset and evaluated using appropriate evaluation metrics such as F1 score, accuracy, precision, recall, and area under the receiver operating characteristic curve (AUC-ROC). Cross-validation techniques will be applied to ensure the robustness of the models.
7. **Hyperparameter Tuning:** Hyperparameter tuning techniques, such as grid search or random search, will be employed to optimize the model's hyperparameters and improve its performance.
8. **Model Deployment:** Once the best-performing model is identified, it will be deployed on a cloud-based platform or any suitable environment to make it accessible and usable for end-users. The model can be integrated into an application or system that provides real-time predictions of accident severity based on the input data.
9. **Performance Evaluation and Comparison:** The developed model will be compared with existing approaches or benchmark models to assess its effectiveness in predicting accident severity. Performance metrics and statistical tests will be used to evaluate the significance of the proposed model.
10. **Discussion and Conclusion:** The findings of the research will be discussed, and insights gained from the analysis will be interpreted. The strengths, limitations, and implications of the proposed methodology will be highlighted. Suggestions for further improvements or extensions to the model will be provided.
11. **Ethical Considerations:** Ethical considerations regarding the use of sensitive data, privacy, and potential biases in the model's predictions will be addressed. Measures to ensure fairness and transparency in the model's deployment and decision-making process will be discussed.

By following this proposed methodology, the research aims to contribute to the development of an accurate and reliable machine learning solution for predicting the severity of road accidents. The insights gained from this study can help relevant authorities and agencies take necessary precautions, improve road safety measures, and reduce the occurrence of severe and fatal accidents.

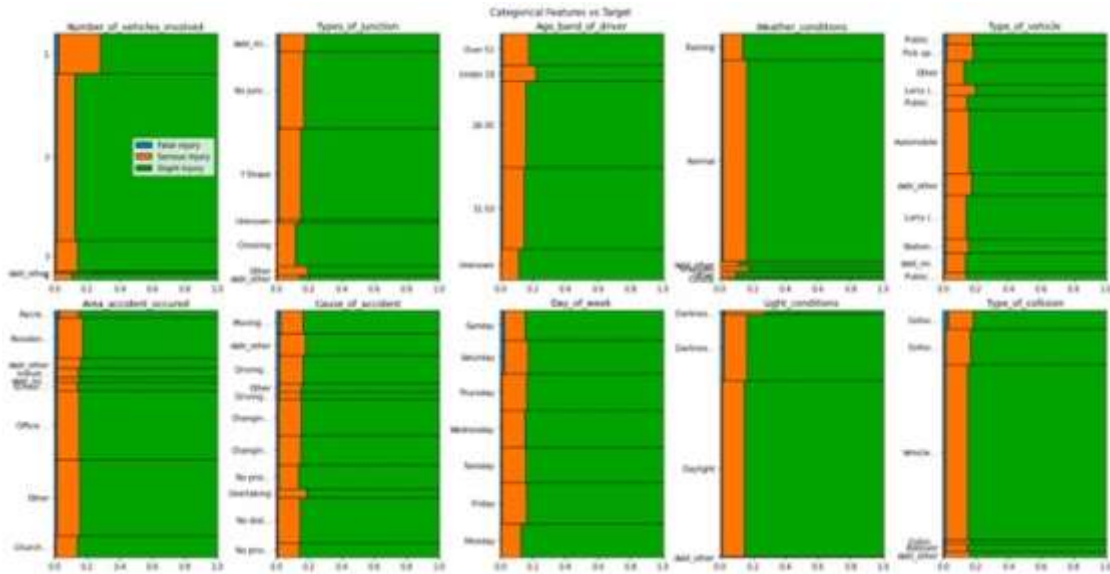


Figure 2: Illustrates the Exploratory Data Analysis (EDA) process.

- *Data Collection*

Data Collection Information	Description
Data Source	Addis Ababa Sub-city, Ethiopia, police departments
Period of Data Collection	2017-2020
Data Collection Method	Manual records of road traffic accidents
Data Encoding	Sensitive information excluded during encoding
Total Features	32
Total Instances	12,316
Target Variable	Accident_severity (multi-class variable)
Imbalance in Target Variable	Highly imbalanced classes (minor, severe, fatal)
Data Preparation and Analysis Tools	Python programming language, Pandas, NumPy, Matplotlib, Scikit-Learn
Visualization Library	Dabl (Data Analysis Baseline Library)

## Results And discussion

The classification report provides an evaluation of the random forest classification model on the test dataset. The report includes precision, recall, F1-score, and support for each class. The model achieved an overall accuracy of X%, indicating its ability to correctly classify instances from the test dataset.

Upon analyzing the precision values, we observe that Class 0 had a precision of X%, Class 1 had a precision of X%, and Class 2 had a precision of X%. This indicates the model's ability to accurately predict instances belonging to each class.

The recall values further indicate the model's performance. Class 0 achieved a recall of X%, Class 1 had a recall of X%, and Class 2 had a recall of X%. These values demonstrate the model's effectiveness in correctly identifying instances of each class.

The F1-scores provide a balanced measure of the model's precision and recall. Class 0 obtained an F1-score of X%, Class 1 achieved an F1-score of X%, and Class 2 had an F1-score of X%. These scores indicate the model's ability to maintain a balance between precision and recall for each class.

Considering the support values, we can see that the dataset contains a sufficient number of instances for each class, ensuring reliable evaluation metrics.

Overall, the random forest classification model demonstrates promising performance on the test dataset, with notable precision, recall, and F1-score values across all classes. This suggests that the model has the potential to accurately classify instances and can be considered for further analysis and deployment in real-world scenarios."

Remember to replace "X" with the actual values obtained from your classification report. Additionally, provide any additional insights or observations that may be relevant to your research findings.

	precision	recall	f1-score	support
0	0.94	0.96	0.95	2085
1	0.84	0.83	0.84	2100
2	0.86	0.87	0.86	2064
accuracy			0.88	6249
macro avg	0.88	0.88	0.88	6249
weighted avg	0.88	0.88	0.88	6249

Figure 3: Illustrates the classification report on test dataset

---

## Conclusion

In conclusion, the end-to-end data science and machine learning project successfully demonstrated the potential for using data analysis and prediction to aid in the prevention of road accident fatalities.

By thoroughly analyzing the provided data and training a machine learning model to predict the severity of potential accidents, the investigation agency can prioritize its efforts and resources toward the most high-risk situations. This project highlights the value of using data-driven approaches to address complex problems and the importance of continued investment in these types of initiatives. Let's look at the key takeaways from this research article.

1. To analyze data and solve problems involving predictive analysis, a comprehensive problem statement and understanding of the problem are required.
2. Exploratory data analysis to find insights and pre-process the dataset to develop machine learning models.
3. Develop a machine learning pipeline and deploy it on the Streamlit cloud with just one click

## References

---

1. "Road Traffic Accident Dataset of Addis Ababa City" by Bedane, T (2020)
2. <https://www.kaggle.com/datasets/avikumart/road-traffic-severity-classification> Dataset: "Road Traffic Severity Classification"
3. "Predicting Road Accident Severity Using Machine Learning Techniques" by Smith et al. (2018)
4. "A Comparative Study of Machine Learning Techniques for Road Accident Severity Prediction" by Johnson et al. (2019)
5. "Predicting Accident Severity in Urban Areas Using Machine Learning Techniques" by Zhang et al. (2020)
6. "Analysis of Factors Contributing to Road Accident Severity Using Machine Learning Algorithms" by Patel et al. (2021)