

Background and Research Questions

Traffic accidents are a widespread problem, resulting in many injuries and deaths worldwide. In 2022, car crashes resulted in 42,795 deaths, incurring costs of \$100 billion.

Through a crash data analysis and identify the trend of leading factors the following research questions are identified:

- What factors contributed to severe injury or fatality in recent years?
- What factors contributed to minor or no injury in recent years?
- How do changes in leading factors influence the injury severity over the years?
- Which modeling scheme can more accurately represent the statistical property of the data and precisely predict the injury severity?

Methodology

Based on the research questions, the following steps are considered:

- 1. Data Collection: Recent accident data from the CRSS from 2019-2021.
- 2. Data Preprocessing: Data cleaning and data engineering to prepare the data set for modeling.
- 3. Labeling: Categorized the severity of accidents into two groups of No injury/Minor injury (label 0) and Major injury/Fatal as label 1.
- 4. Feature Selection: Identified various factors including the year, geographic region, month, speed, alcohol involvement, weather conditions, age group, use of seatbelt, and driver distraction.
- 5. Model Training: Performed correlation analysis to gain insights into the factors and how they interact with each other and the severity.
- 6. Model Evaluation: Used binary classification methods of Logistic Regression and Random Forest, model analysis, and visualization of the data.

Data Engineering

Figure 1. No. of crash injuries per month over year by injury severity

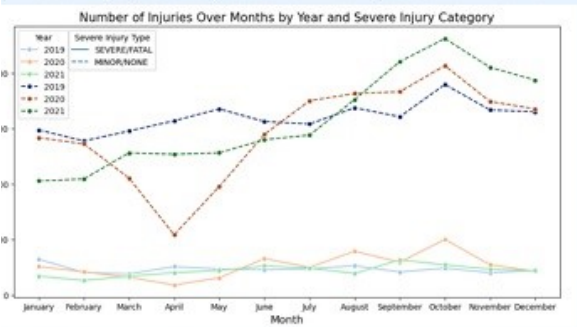


Figure 1. No. of crash injuries per month over year by injury severity

Findings toward modeling

The data is imbalanced: Only 13.3% of the data represent fatal/serious (label 1). A simple theoretical model classifying all cases as Minor/no injury could achieve 86.7% accuracy. Advanced techniques for dealing with imbalanced data such as sampling are required. Oversampling with SMOTE resulted in computational challenges. Undersampling was identified as the remedy for imbalanced data. In 2020, there were 120 recorded the lowest number of crashes in each category. The South region experienced the highest frequency of crashes. The year 2019 had the most major/fatal injuries. The year 2021 had the most minor/no injuries.

Factor Distribution-Alcohol

Alcohol impaired crashes by age group:

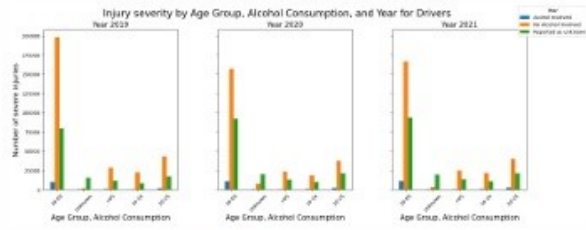


Figure 2. Variable distribution

Alcohol impaired crashes by region:

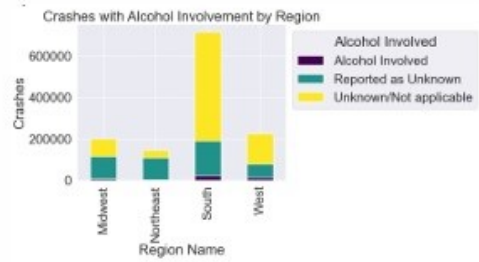


Figure 3. Variable distribution

Correlation Analysis

Correlation analysis to decipher the relationships between independent variables:

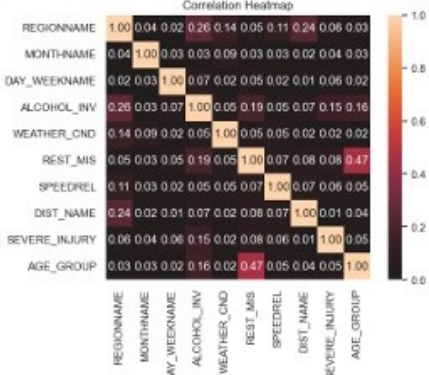
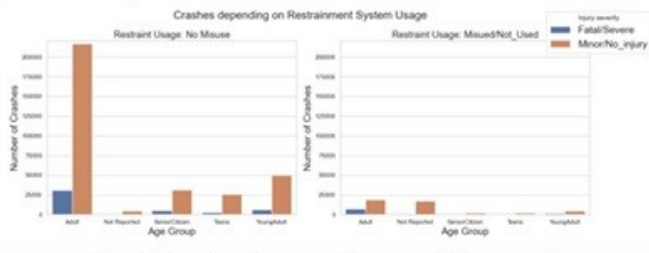


Figure 4. Correlation Matrix of Variables

- restraint showed the strongest correlation with the age group.
- adults and young adults tend to have a higher incidence of restraint misuse or non-use compared to instances where restraints were not misused.
- More crashes reported when restraint was used for senior citizens.



Modeling

Modeling schemes

1)Logistic Regression (LR) 2)Random Forest (RF)

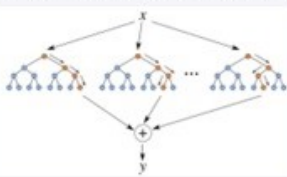
LR models the relationship between predictors and a categorical response variable as follows:

Logit(P) = β0 + β · X + ε

Where {x1,x2,...,xn} represent the independent variables used in the modeling. Hence we can derive:

P(Y = 1) = 1 / (e^-(β0+βX) + 1)

RF uses labeled data to learn and classify unlabeled data, with architecture:



Modeling Challenge

Downsampling can enhance the precision and recall of the model, however:

- Downsampling can lead to a loss of information.
- Reducing the size of the majority class affects the overall accuracy of the model.
- Comparison of sampling rate vs accuracy is critical to decide the best model.

Results and Conclusions

Trade-off between downsampling and evaluation metrics of modeling:

Modeling	Injury Severity Class	Precision	Recall	F1
LR-100% sampling-rate	0	0.59	0.71	0.64
	1	0.64	0.52	0.58
LR-75% sampling-rate	0	0.63	0.86	0.73
	1	0.64	0.34	0.42
LR-50% sampling-rate	0	0.70	0.94	0.81
	1	0.63	0.21	0.33
LR-25% sampling-rate	0	0.81	0.99	0.90
	1	0.56	0.06	0.10
RF-100% sampling-rate	0	0.79	0.78	0.78
	1	0.78	0.79	0.78
RF-75% sampling-rate	0	0.81	0.84	0.82
	1	0.77	0.73	0.75
RF-50% sampling-rate	0	0.84	0.90	0.87
	1	0.76	0.65	0.70
RF-25% sampling-rate	0	0.89	0.96	0.92
	1	0.76	0.51	0.61

Table 1. Modeling performance in presence of different sampling rates

Key Findings

- RF exceeds the modeling performance of LR.
- RF with 75% sampling-rate resulted in the best performance (with an accuracy of 0.92).
- LR with 100% sampling-rate performs the best among LR cases (with an accuracy of 0.64).
- This analysis showed that the South region experienced the highest number of alcohol-involved crashes.
- West region reported the second highest.
- Distraction by Outside/Others is the main category of known distraction contributing to both minor and major/fatal crashes.
- This category caused the highest fatalities in 2020.
- Speeding Contributed to the increasing trend in crashes over the recent years.