
Car Crash Prevention

Hotspot Identification and Prevention

Kevin Cha¹ Jonathan Lam¹ Emily Zhou¹

Abstract

Emily Zhou

Our research question for this project is how to predict and prevent mobile vehicle crashes in Virginia, especially with crashes being a leading cause of death in the United States. By identifying hotspots for car crashes and the common causes for car crashes, we can work with state officials to better identify how to promote safer driving and reduce fatalities from car crashes.

1. Data

To answer our research question, we have used the VDOT's dataset on crashes within Virginia. <https://dashboard.virginiadot.org/pages/safety/crashes.aspx>

1.1. Dataset Jonathan Lam

This dataset consists of motor vehicle crash records in the state of Virginia since 2016. Each record corresponds to a different crash and provides detailed information on roadway conditions, driver factors, environmental conditions, location details, and affected individuals, along with the outcome of the incident. The raw dataset contains 69 variables and over 1 million observations. This dataset was chosen because it provided a comprehensive and reliable overview of all the different competing factors involved in a car crash. As such, we can better analyze how they influence the frequency and severity of collisions to build machine learning models that can assist government officials with public safety and maintenance decisions.

^{*}Equal contribution ¹Department of Computer Science, University of Virginia, Charlottesville, VA, USA. Correspondence to: Kevin Cha <hpb2gv@virginia.edu>, Jonathan Lam <rgm6yp@virginia.edu>, Emily Zhou <csz6wd@virginia.edu>.

1.2. Key Variables Jonathan Lam

In order to build an effective model, we need to identify key variables to serve as features. Because we wish to discern the patterns behind car crashes, variables describing the severity and outcome of the accident, vehicles involved, environmental and weather conditions, road conditions, crash dynamics, traffic control, driver conditions, special populations, and whether it was work-zone-related are all included. Variable groups are listed instead of each individual variable due to the vast number of variables. Although driver conditions may not classify as an external contributing factor to a collision, context on whether a given driver was under the influence, distracted, drowsy, unbuckled, or speeding provides a more holistic understanding behind an accident's severity. Remaining variables containing metadata or identifiers are also dropped in the data cleaning pipeline.

These variables groups are key because they represent the broad domain of contributing factors in a collision. Variables pertaining to the crash severity and outcome serve as a baseline for predicting crashes, while variables about the vehicles involved (trucks, motorcycles, etc.) may also influence a crash's severity. Environmental and weather conditions have a direct effect on driving performance, and, consequently, the likelihood of an accident; heavy rain or snow limit visibility and influence reaction time. Road conditions play an equally impactful role, where potholes or icy roads could result in serious injuries. Crash dynamics and traffic control provide additional knowledge on an accident's severity. As aforementioned, driver conditions directly relate to incident severity, and special populations and work zone accidents highlight special at-risk groups and conditions.

1.3. Data Wrangling Challenges Emily Zhou

The first challenge to handle with this dataset was the sheer size of the data from VDOT, as it had data from 2016 and had around 1 million entries. We wanted relevant data, considering our research question of what can be done to prevent crashes in Virginia. Take for example the case of roadway defects such as potholes, since they could have been fixed or new ones have emerged over time and use of the roads. The next step for data cleaning was to determine

which variables were still useful to our model, since while there were key variables, there were other variables such as OBJECTID which would not support our model so we dropped those columns completely.

The next big steps for cleaning the data were handling both categorical and numerical data, where we had to decide how to handle null values and also turn our categorical data into meaningful numerical data for our model to train on. For null values, it was determined that there were a few variables that only had a few null values so those rows were just dropped but for the following variables: Max Speed Diff, Functional Class, and Facility Type, had more than a few thousand null values and required further analysis. Max Speed Diff was only showing data where the difference between the speed of the car at the incident compared to the speed limit where the crash happened, so it was fair to fill out those null values with 0's. Then for Functional Class, it was a variable that determined where the crash occurred and since 10% of the values were NaN's, imputing them as "Unknown" sufficed. With the last variable with a significant number of null values being Facility Type, it also explained the format of the road where the crash occurred, so the remaining null values were imputed as "Unknown" as well.

With our categorical data, we had some issues discerning what the variables mean since there was no specific data dictionary for this dataset but after some thorough search the most confusing variables, K_people, A_people, B_people, C_people were determined to be on a scale referred to as KABCO. This scale indicated the damage to human lives, where K meant fatality, A meant serious injury, B meant minor injury, C meant possible injury, and O meant property damage only. The best way to handle this in a meaningful manner was to assign weight mappings to the variables, since K does have a higher severity due to lives being lost.

2. Methods and Results Kevin Cha

2.1. Method Overview

The task at hand is a classification task to predict the crash severity of an incident given various factors. The important features discussed previously are used with an 80:20 train/test split. The models used will be Random Forest, XGBoost, and Neural Networks. 5-fold cross validation will also be used to ensure that the model is not overfitting, and to optimally tune hyperparameters. The data will then be transformed by one-hot encoding the categorical features, and scaling the numerical features for stronger Neural Network performance.

2.2. Feature Filtering, Engineering, and Manipulation

After data wrangling, a few features were filtered, engineered, and manipulated to improve the model performance and avoid feeding junk features. Starting with the filtered features - any features with extremely low variance were dropped, since they are essentially constant and thus would not contribute meaningfully to the model's prediction aside from adding noise.

Existing features were also manipulated and engineered to produce features with less bins to reduce sparsity in certain categoricals. The following features were manipulated: coordinates, speed difference from limit, weather conditions, and light conditions/nighttime. Coordinates were modified by creating clusters using a KMeans cluster approach to create 100 area patches. The remaining features had some of their values binned and condensed - for example, for weather features, "rain", "drizzle", and "shower" were condensed to just "rain". This was done in cases where a feature may have had too many highly similar categories. A new feature was created by identifying cases where it was rainy and night, an expected combination of high crash risk.

These features were fed into a pipeline. This pipeline imputed any missing values. It also scaled the numericals with scikit-learn's StandardScaler, and one hot encoded the categoricals.

After these transformations were made to the data, further feature refinement was done by selecting the K best features with a mutual info classifier from scikit-learn.

Finally, the data was stratified, shuffled, and split via scikit-learn's StratifiedShuffleSplit. This was done to preserve sample ratios for the different classes when creating the train and test set.

2.3. Model Selection & Justification

Random Forest was selected as the first model choice for its robustness stemming from being an ensemble of Decision Trees. It typically performs well in classification tasks due to its capability to capture more complex relationships, which is important due to our large amount of features. Another strong suit of Random Forest is their interpretability, which allows for further analysis and understanding of what factors may influence crash severity.

XGBoost was selected as the second model choice due to its historically strong performance on classification tasks. Boosting as an ensemble learning method aims to optimize both bias and variance, so it may outperform Random Forest (a bagging ensemble learning method). Similar to Random Forest, there is strong interpretability from results as feature importance can be extracted.

A Neural Network is selected as the third model choice due

to its capability to handle very complex relationships thanks to the many nodes that can be added in hidden layers. Due to how many different features are present in the cleaned dataset, Neural Networks may outperform Random Forest or XGBoost given a more complex architecture capable of handling the features.

2.4. Model Training/Hyperparameters

For the Random Forest model, the following hyperparameters were considered and tuned: `n_estimators`, `max_depth`, and `min_samples_split`. `n_estimators` is an important hyperparameter because it tunes the number of trees to fit, which would lead to a more robust model for voting. `max_depth` is important to limit the depth of the tree and minimize overfitting. `min_samples_split` requires more samples for a category to justify a split, which is important to prevent overfitting on smaller less important categories in features.

For the XGBoost model, the following hyperparameters were considered and tuned: `n_estimators`, `max_depth`, `subsample`, and `colsample_bytree`. The `n_estimators` and `max_depth` parameters are tuned for similar reasons to the RandomForest model given that XGBoost is also a tree-based ensemble learning method and is equally important. The `subsample` and `colsample_bytree` forces subsampling of the training data for the trees, allowing for less overfitting and more robustness overall.

For the Neural Network, the final architecture contained two hidden layers, each with a dropout of 0.2 to prevent overfitting. It was optimized with the Adam optimizer. Training was completed with 60 epochs and a batch size of 256 at a time - early stop was also implemented with a patience of 10 epochs. The class weights were calculated to handle class imbalance.

2.5. Performance Metrics and Loss Function

The performance metric used for training and testing was F1 score. This is because in calculating crash severity the model should optimally maximize both precision and recall. Precision would measure how many crashes predicted as a severity rating were correct, while recall would measure how many crashes of a particular severity rating were correct. The model should perform well in both aspects. The loss function used is cross-entropy loss, as it is standard for classification tasks.

2.6. Preliminary Results

The models have not been run yet because more testing needing to be completed with respect to hyperparameter tuning and because the models take a lot of time. However, they will be run through Google Colab and results will be prepared by the next milestone.