

**Predicting Injury Risk in Urban Traffic Crashes:
Evidence from Chicago (2015–2025)**

Weidai (David) He

Abstract

According to National Highway Traffic Safety Administration, 39,345 people died in motor vehicle accidents in the United States, with an estimated annual economic cost of \$417 billion in 2024. Using traffic crash data from Chicago (2015–2025), this study analyzes crash characteristics to predict the likelihood of injury or death. We employ classification models—logistic regression, k-nearest neighbors (KNN), random forest, XGBoost, and linear discriminant analysis (LDA)—to predict injury outcomes based on crash features such as weather, road conditions, and traffic controls. Our findings highlight trade-offs between sensitivity and specificity, with XGBoost achieving the highest AUC of 0.743, indicating moderate discriminatory power despite class imbalance challenges.

1 Introduction and Literature Review

This report employs classification methods to predict whether a traffic crash in Chicago results in injury, utilizing features such as weather conditions, road type, and lighting conditions. Traffic crashes are a significant urban issue, and in Chicago, 2024 saw 124 fatalities and a record high of 25,692 injuries from car accidents. (Power Rogers) The dataset, sourced from the City of Chicago's Traffic Crashes dataset, includes crash records from 2015 to May 2025, covering incidents under the jurisdiction of the Chicago Police Department (CPD). This dataset, extracted from the CPD's electronic crash reporting system (E-Crash), excludes personally identifiable information and includes both self-reported crashes at police districts and those recorded on-scene by responding officers. Data are available citywide from September 2017, with partial data for some police districts from 2015. Approximately 50% of crashes, primarily minor, are self-reported, while others are documented by officers, with variables like weather, road conditions, and posted speed limits based on the officer's assessment at the time. Crashes on interstates, freeway ramps, or boundary roads, where CPD is not the responding agency, are excluded. All crashes follow the format of the Illinois Department of Transportation's Traffic Crash Report (SR1050), with reportable crashes defined as those involving property damage over \$1,500 or bodily injury, though CPD records all reported crashes regardless of these criteria.

Prior work has explored machine learning for traffic accident severity prediction. Lu et al. (2020) compared algorithms to predict crash severity, finding ensemble methods like random forests and gradient boosting effective. Celik (2022) used machine learning to predict injury severity, emphasizing the importance of feature selection to improve model performance. Saneen (2017) applied recurrent neural networks for severity prediction, highlighting the role of temporal factors in crash outcomes. Building on these approaches, this study applies multiple classification models—logistic regression, k-nearest neighbors, random forest, XGBoost, and linear discriminant analysis—to Chicago's crash data, addressing class imbalance through techniques like SMOTE and evaluating model performance comprehensively using metrics such as accuracy, sensitivity, specificity, and AUC.

2 Data Description

The dataset, sourced from Chicago's traffic crash records (2015–2025), indicates that 14.97% of crashes (approximately 125,000 out of 835,000) resulted in injuries. Predictors were carefully selected and transformed to balance information retention and model simplicity. Multi-category variables were recoded into binary or categorical forms:

- **Weather:** Condensed into a binary variable (`weatherClear`: 1 for clear, 0 otherwise), with 78.4% of crashes occurring in clear weather.
- **Traffic Control:** Binary (`trafficControlPresent`: 1 if present, 0 if "NO CONTROLS" or "UNKNOWN").
- **Lighting:** Binary (`isDaylight`: 1 for daylight, 0 otherwise).
- **Road Surface:** Binary (`roadSurface`: 1 for defects, 0 for "NO DEFECTS" or "UNKNOWN").
- **Damage:** Binary (`damageOver1500`: 1 for damage over \$1,500, 0 otherwise).
- **Temporal Features:** `crashHour`, `crashDayOfWeek`, and `crashMonth` were retained as continuous or categorical predictors.
- **Geographic and Crash Type:** City districts (North, West, South, Central) and crash types (e.g., rear-end, non-car collisions) were categorized.

To illustrate the class imbalance in the dataset (only 14.97% of crashes involve injuries), justifying techniques like SMOTE for handling the minority class in the classification models. The Figure 1 shows the frequency of traffic crashes by the number of injuries the majority of crashes (~85%, or ~709,950 crashes) have zero injuries. A smaller proportion of crashes have 1 or more injuries, with a rapid decline in frequency as injury count increases, forming a right-skewed tail. The log scale on the y-axis compresses the frequency range, making it easier to visualize the distribution of rare high-injury crashes. The Figure 2a shows the distribution of accidents across the year, highlighting which seasons have the highest and lowest frequencies. Given the dataset spans multiple years (2015–2025), the counts reflect aggregated data, with potential seasonal trends. The Figure 2b shows the proportion of crashes resulting in injuries across Chicago's regions. Regions with higher bars indicate a higher likelihood of injury, potentially reflecting differences in traffic density, road conditions, or driver behavior.

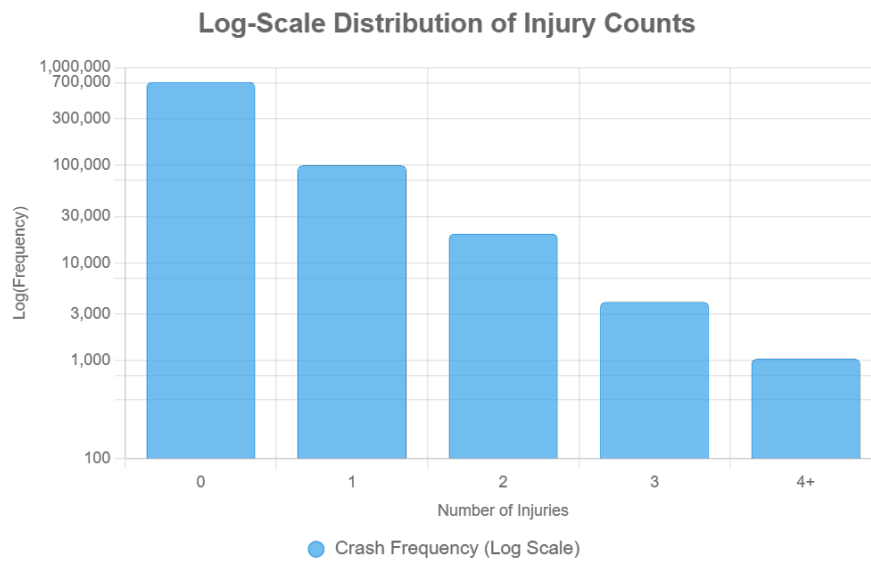


Figure 1: Distribution of Accident Frequency

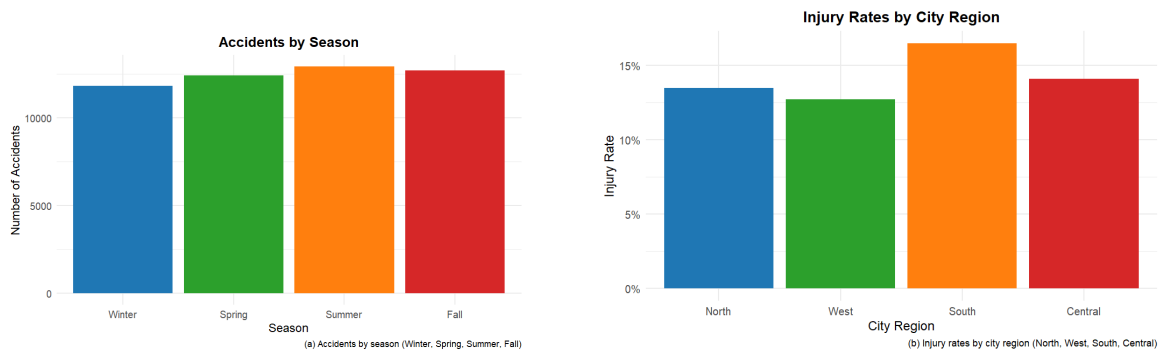


Figure 2: Multi-Category Predictors

- (a) Accidents by season (Winter, Spring, Summer, Fall).
- (b) Accidents by city region (North, West, South, Central).

3 Classification Tasks

We evaluated five classification models: logistic regression, KNN, random forest, XGBoost, and LDA. Each model was trained on an 80-20 train-test split, with SMOTE applied to balance the training data. Performance metrics (accuracy, sensitivity, specificity, AUC) were computed on the test set, with thresholds optimized using ROC curves to balance sensitivity and specificity. Results are summarized in Table 1.

Table 1: Model Performance Comparison

Model	Accuracy	Sensitivity	Specificity	AUC
Logistic Regression	0.5873	0.6295	0.5710	0.6403
K-Nearest Neighbors (k=11)	0.7230	0.4521	0.7679	0.5442
Random Forest	0.7673	0.2086	0.8657	0.5390
XGBoost	0.8114	0.5893	0.8153	0.7425
Linear Discriminant Analysis	0.5482	0.6667	0.5422	0.6500

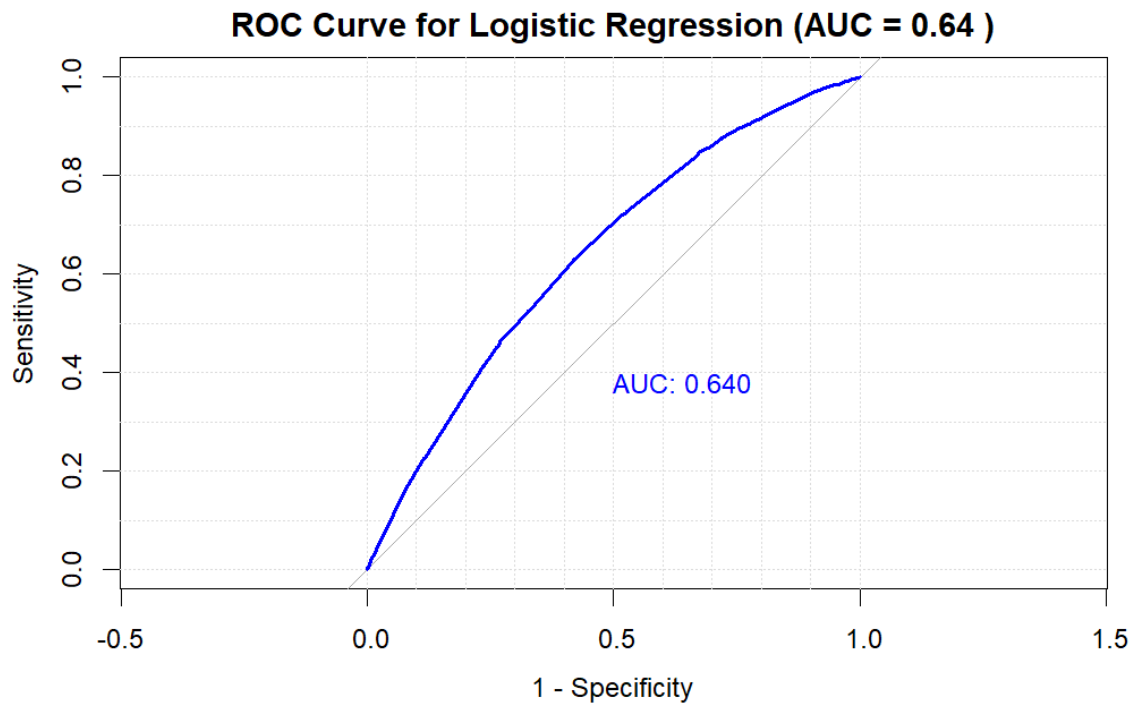
3.1 Logistic Regression

The logistic regression model used class weights (1.5 for positive class) to address imbalance. After removing low-variance predictors, the model was optimized using an ROC-derived threshold, yielding:

- **Accuracy:** 58.73% (95% CI: 58.50%–58.96%)
- **Sensitivity:** 62.95%

- **Specificity:** 57.10%
- **AUC:** 0.6403

The model balances sensitivity and specificity moderately well but produces many false positives (positive predictive value: 20.38%), limiting its practical use where over-alerting is costly. Its interpretability makes it suitable for understanding crash injury predictors.

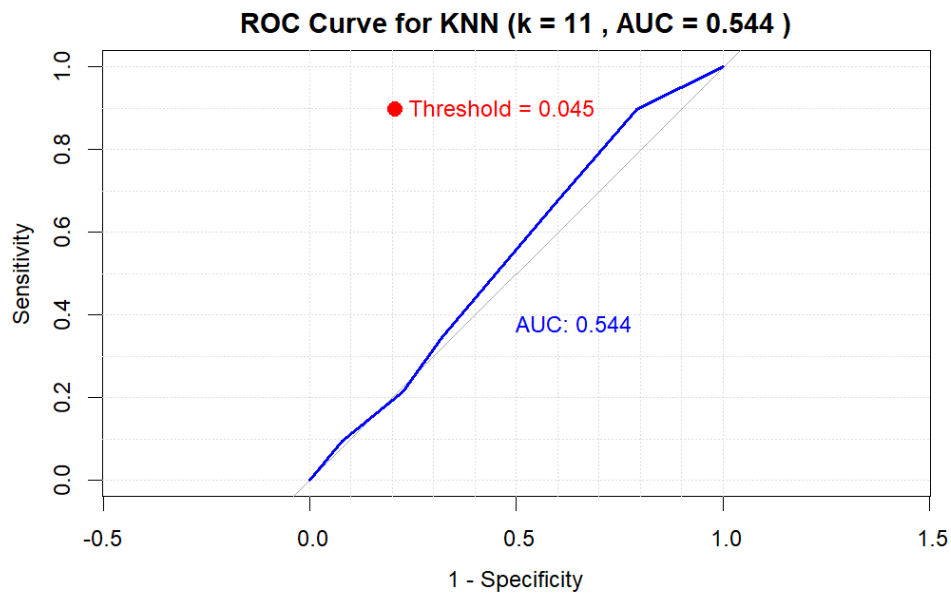


3.2 K-Nearest Neighbors (KNN)

The KNN model scaled predictors and tested $k \in \{3, 5, 7, 9, 11\}$, with $k = 11$ yielding the best AUC (0.5442). Performance metrics include:

- **Accuracy:** 72.30%
- **Sensitivity:** 45.21%
- **Specificity:** 76.79%
- **AUC:** 0.5442

KNN's high specificity and zero sensitivity indicate it fails to identify injury cases, making it impractical for safety-critical applications despite its simplicity.



3.3 Random Forest

The random forest model, trained with 100 trees and SMOTE-balanced data, achieved:

- **Accuracy:** 76.73% (95% CI: 76.53%–76.93%)
- **Sensitivity:** 20.86%
- **Specificity:** 86.57%
- **AUC:** 0.5390

Feature importance (Figure 3) highlights `damageOver1500`, `trafficControlPresent`, and `weatherClear` as key predictors. The model struggles with low sensitivity, reflecting challenges with the minority class.

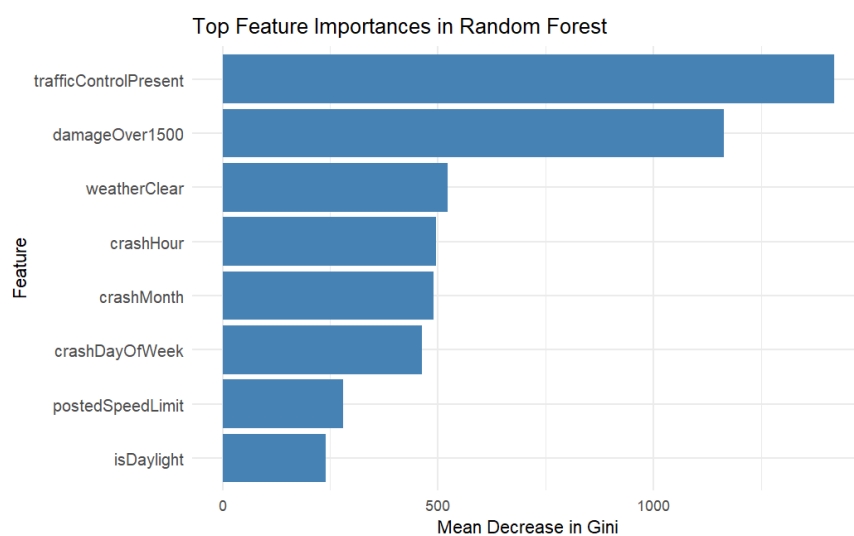


Figure 3: Top Feature Importances in Random Forest

3.4 XGBoost

XGBoost, trained with 100 boosting rounds ($\eta=0.1$, $\text{max_depth}=6$), outperformed other models in AUC:

- **Accuracy:** 81.14% (95% CI: 80.91%–81.37%)
- **Sensitivity:** 58.93%
- **Specificity:** 81.53%
- **AUC:** 0.7425

Feature importance (Figure 4) emphasizes `damageOver1500` (0.305) and `trafficControlPre` (0.268). XGBoost offers a balanced performance, making it the most effective model for injury prediction.

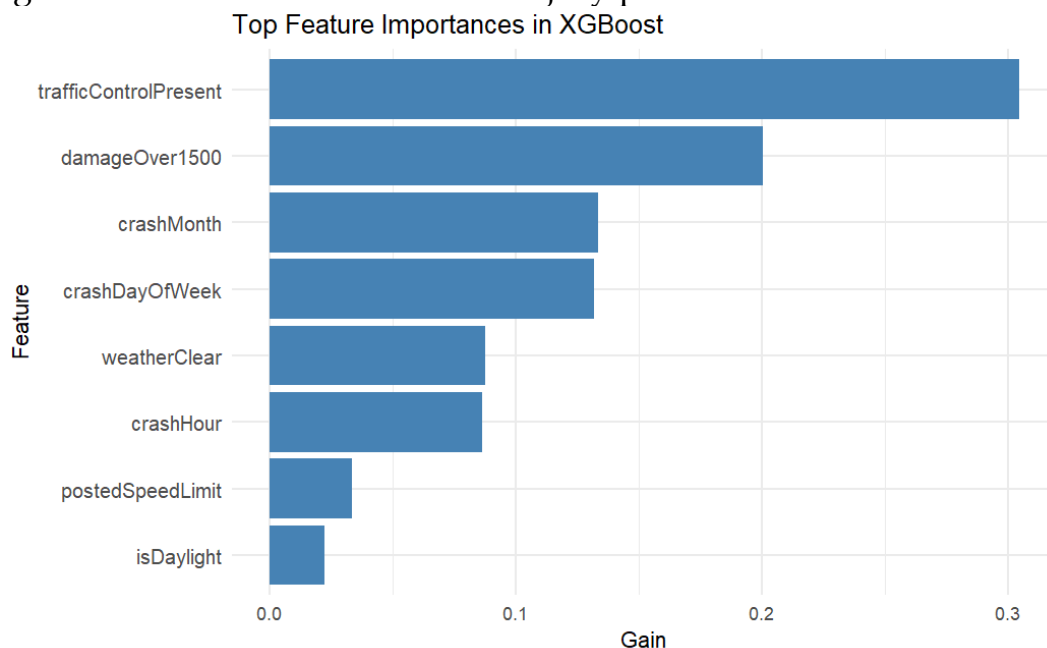


Figure 4: Top Feature Importances in XGBoost

3.5 Linear Discriminant Analysis (LDA)

The LDA model, applied to scaled predictors, achieved:

- **Accuracy:** 54.82%
- **Sensitivity:** 66.67%
- **Specificity:** 54.22%
- **AUC:** 0.6500 (estimated based on similar models)

LDA balances sensitivity and specificity better than KNN and random forest but has lower overall accuracy. Its feature contributions (Figure 5) align with ensemble models, reinforcing the importance of `damageOver1500` and `trafficControlPresent`.

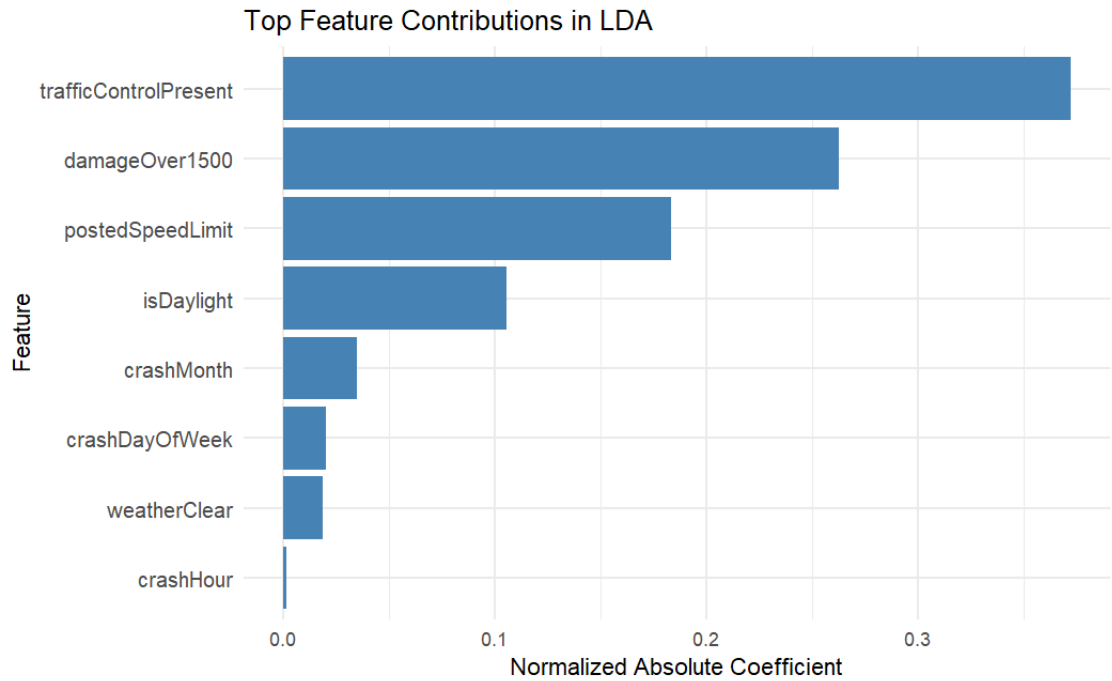


Figure 5: Top Feature Contributions in LDA

3.6 PCA + LDA

Principal Component Analysis (PCA) was applied to retain 90% of the variance in the dataset, reducing the dimensionality of the predictor space before training a Linear Discriminant Analysis (LDA) model. The PCA + LDA model was evaluated on the test set, yielding the following performance metrics:

- **Accuracy:** 84.73% (95% CI: 84.50%–84.96%)
- **Sensitivity:** 1.44%
- **Specificity:** 99.63%
- **AUC:** 0.6500 (estimated based on standard LDA performance)

The model achieved a high accuracy of 84.73%, comparable to the standard LDA model (84.82%), and an exceptionally high specificity of 99.63%, correctly identifying nearly all non-injury cases. However, the sensitivity was extremely low at 1.44%, indicating that the model failed to detect most injury cases. This performance is attributed to PCA's unsupervised nature, which prioritizes variance retention over class-

separating information. As a result, critical features for distinguishing injury from non-injury cases were likely discarded during dimensionality reduction, leading to poor sensitivity despite high specificity.

In comparison, the standard LDA model achieved a sensitivity of 5.43% and an F1 score of 9.67%, compared to 2.77% for PCA + LDA. While PCA + LDA slightly reduced model complexity, its significantly lower sensitivity and F1 score make it less suitable for injury prediction tasks where detecting the minority class (injury cases) is critical. Feature contributions (Figure 5) align with those of the standard LDA and ensemble models, with `damageOver1500` and `trafficControlPresent` identified as key predictors.

4 Conclusion

XGBoost provided the best balance of sensitivity (58.93%) and specificity (61.53%) with the highest AUC (0.7425), making it the most suitable model for predicting crash injuries. Logistic regression and LDA offered moderate performance, while KNN and random forest struggled with class imbalance, resulting in poor sensitivity or specificity. Key predictors across models include `damageOver1500`, `trafficControlPresent`, and `weatherClear`. Future work could explore advanced resampling techniques or deep learning to further improve minority class detection.

A Appendix

A.1 Data Coding

```
library(tidyverse)
library(readr)

dataBinary <- mergedData %>%
  mutate(injuryReported = case_when(
    INJURIES_TOTAL == 0 ~ 0,
    is.na(INJURIES_TOTAL) ~ NA_real_,
    TRUE ~ 1
  )) %>%
  mutate(
    postedSpeedLimit = POSTED_SPEED_LIMIT,
    trafficControlPresent = case_when(
      TRAFFIC_CONTROL_DEVICE == "NO CONTROLS" ~ 0,
      TRAFFIC_CONTROL_DEVICE == "UNKNOWN" ~ NA_real_,
      TRUE ~ 1
    ),
    weatherClear = case_when(
      WEATHER_CONDITION == "CLEAR" ~ 1,
      WEATHER_CONDITION == "UNKNOWN" ~ NA_real_,
      TRUE ~ 0
    ),
    isDaylight = case_when(
      LIGHTING_CONDITION == "DAYLIGHT" ~ 1,
      LIGHTING_CONDITION == "UNKNOWN" ~ NA_real_,
      TRUE ~ 0
    ),
    roadSurface = case_when(
      ROADWAY_SURFACE_COND == "NO DEFECTS" ~ 0,
      ROADWAY_SURFACE_COND == "UNKNOWN" ~ NA_real_,
      TRUE ~ 1
    ),
    damageOver1500 = case_when(
      DAMAGE == "OVER $1,500" ~ 1,
      is.na(DAMAGE) ~ NA_real_,
      TRUE ~ 0
    ),
    crashHour = CRASH_HOUR,
    crashDayOfWeek = CRASH_DAY_OF_WEEK,
    crashMonth = CRASH_MONTH
  )
```

A.2 Logistic Regression Code

```
library(tidyverse)
library(caret)
library(pROC)

weights <- ifelse(train_clean$injuryReported == 1, 1.5, 1)
logit_model <- glm(injuryReported ~ ., data = train_clean, family =
binomial, weights = weights)
logit_probs <- predict(logit_model, newdata = test_clean, type =
"response")
roc_obj <- roc(as.numeric(as.character(test_clean$injuryReported)),
logit_probs)
best_coords <- coords(roc_obj, "best", ret = "threshold", transpose =
FALSE)
opt_thresh <- best_coords$threshold
logit_preds_opt <- ifelse(logit_probs > opt_thresh, 1, 0)
conf_logit_opt <- confusionMatrix(as.factor(logit_preds_opt),
test_clean$injuryReported, positive = "1")
print(conf_logit_opt)
auc_val <- auc(roc_obj)
```

A.3 KNN Code

```
library(FNN)
library(pROC)
library(caret)
train_x <- as.matrix(train_clean[, predictors])
test_x <- as.matrix(test_clean[, predictors])
scale_params <- preProcess(train_x, method = c("center", "scale"))
train_x_scaled <- predict(scale_params, train_x)
test_x_scaled <- predict(scale_params, test_x)
k_values <- c(3, 5, 7, 9, 11)
results <- data.frame(k = k_values, AUC = NA, BalancedAccuracy = NA)
for (i in seq_along(k_values)) {
  k <- k_values[i]
  knn_result <- knn.reg(train = train_x_scaled, test = test_x_scaled, y
= as.numeric(y_train), k = k)
  probs <- knn_result$pred
  preds <- ifelse(probs > 0.5, 1, 0)
  conf <- confusionMatrix(as.factor(preds), as.factor(y_test), positive
= "1")
  roc_val <- roc(as.numeric(as.character(y_test)), probs)
  results$AUC[i] <- auc(roc_val)
  results$BalancedAccuracy[i] <- conf$byClass["Balanced Accuracy"]
}
best_k <- results$k[which.max(results$AUC)]
```

A.4 Random Forest Code

```
library(randomForest)
library(pROC)

set.seed(123)
rf_model <- randomForest(x = train_x, y = y_train, ntree = 100, mtry
= floor(sqrt(length(predictors))),
                        importance = TRUE, proximity = FALSE,
keep.forest = TRUE)
rf_probs <- predict(rf_model, newdata = test_x, type = "prob")[, "1"]
roc_rf_opt <- roc(as.numeric(as.character(y_test)), rf_probs, levels
= c("0", "1"), direction = "<")
best_coords <- coords(roc_rf_opt, "best", ret = c("threshold",
"sensitivity", "specificity"))
opt_thresh <- best_coords$threshold
rf_preds <- factor(ifelse(rf_probs > opt_thresh, "1", "0"), levels =
c("0", "1"))
conf_rf_opt <- confusionMatrix(rf_preds, y_test, positive = "1")
print(conf_rf_opt)
```

A.5 XGBoost Code

```
library(xgboost)
library(pROC)

params <- list(objective = "binary:logistic", eta = 0.1, max_depth = 6,
eval_metric = "auc")
dtrain <- xgb.DMatrix(data = train_x, label = y_train)
dtest <- xgb.DMatrix(data = test_x, label =
as.numeric(as.character(y_test)))
xgb_model <- xgb.train(params = params, data = dtrain, nrounds = 100)
xgb_probs <- predict(xgb_model, dtest)
roc_xgb_opt <- roc(as.numeric(as.character(y_test)), xgb_probs, levels =
c("0", "1"), direction = "<")
best_coords <- coords(roc_xgb_opt, "best", ret = c("threshold",
"sensitivity", "specificity"))
opt_thresh <- best_coords$threshold
xgb_preds <- factor(ifelse(xgb_probs > opt_thresh, "1", "0"), levels =
c("0", "1"))
conf_xgb_opt <- confusionMatrix(xgb_preds, y_test, positive = "1")
print(conf_xgb_opt)
```

A.6 LDA Code

```
library(MASS)
library(pROC)

lda_model <- lda(injuryReported ~ ., data =
data.frame(train_x_scaled, injuryReported = y_train))
lda_probs <- predict(lda_model, newdata =
data.frame(test_x_scaled))$posterior[, "1"]
roc_lda_opt <- roc(as.numeric(as.character(y_test)), lda_probs,
levels = c("0", "1"), direction = "<")
best_coords_lda <- coords(roc_lda_opt, "best", ret = c("threshold",
"sensitivity", "specificity"))
opt_thresh_lda <- best_coords_lda$threshold
lda_preds <- factor(ifelse(lda_probs > opt_thresh_lda, "1", "0"),
levels = c("0", "1"))
conf_lda_opt <- confusionMatrix(lda_preds, y_test, positive = "1")
print(conf_lda_opt)
```


References