# Predicting Crash Severity of Chicago Streets

Hawraa Salami

SPRINGBOARD Capstone Project 1

# 2018 Chicago Crash Facts

| Total Crashes | 98,859 |
|---|---|
| Total Injuries | 24,400 |
| Total Incapacitating Injuries | 2,609 |
| Total Deaths | 132 |

- What conditions lead to severe crashes?

- Data-driven models aim to understand the severity of crashes.

# How can the models help?

**City of Chicago**

**Car Manufacturer**

**Insurance Companies**

Develop better traffic control policies

Incorporate more safety features

Perform better risk assessment

# Chicago Crashes Dataset

- The data is available from Chicago online portal*.
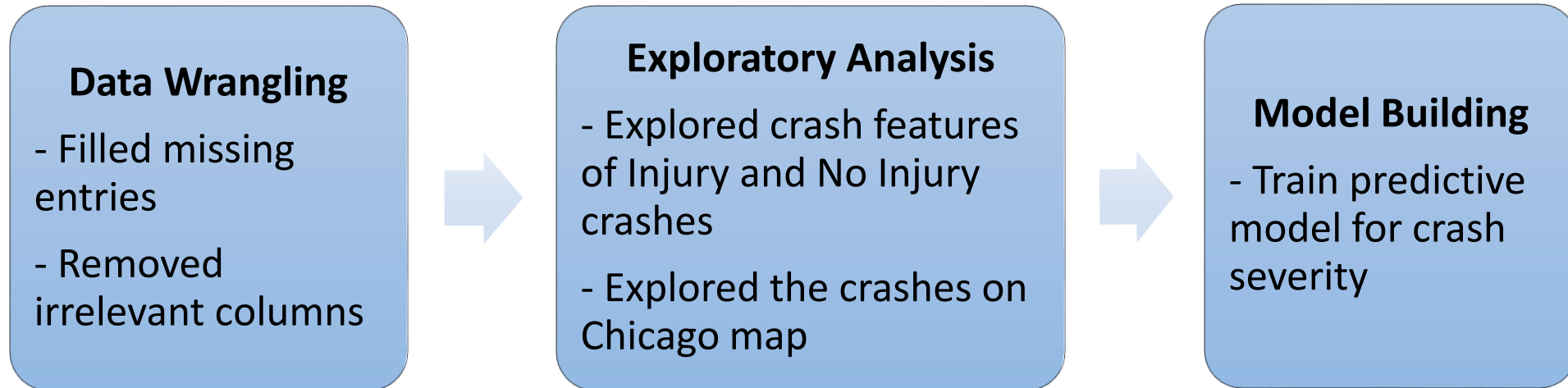- It contains information of Chicago crashes from 2015 to present.

| Crash Location and Time | External Conditions | Crash Cause and Description |
|---|---|---|
| Crash hour, day, month | Weather & Lighting Conditions | Primary Cause (driving behavior) |
| Crash address | Road alignment, type and surface | Type of Collision |
| | Speed Limit and Control Device | Type of crash: Injury or No injury |

* https://data.cityofchicago.org/Transportation/Traffic-Crashes-Crashes/85ca-t3if

# Data Analysis Steps

- **Goal:** Build a model that predicts the type of crash (Injury or No Injury)

- **Steps:**

**Data Wrangling**

- Filled missing entries

- Removed irrelevant columns

**Exploratory Analysis**

- Explored crash features of Injury and No Injury crashes

- Explored the crashes on Chicago map
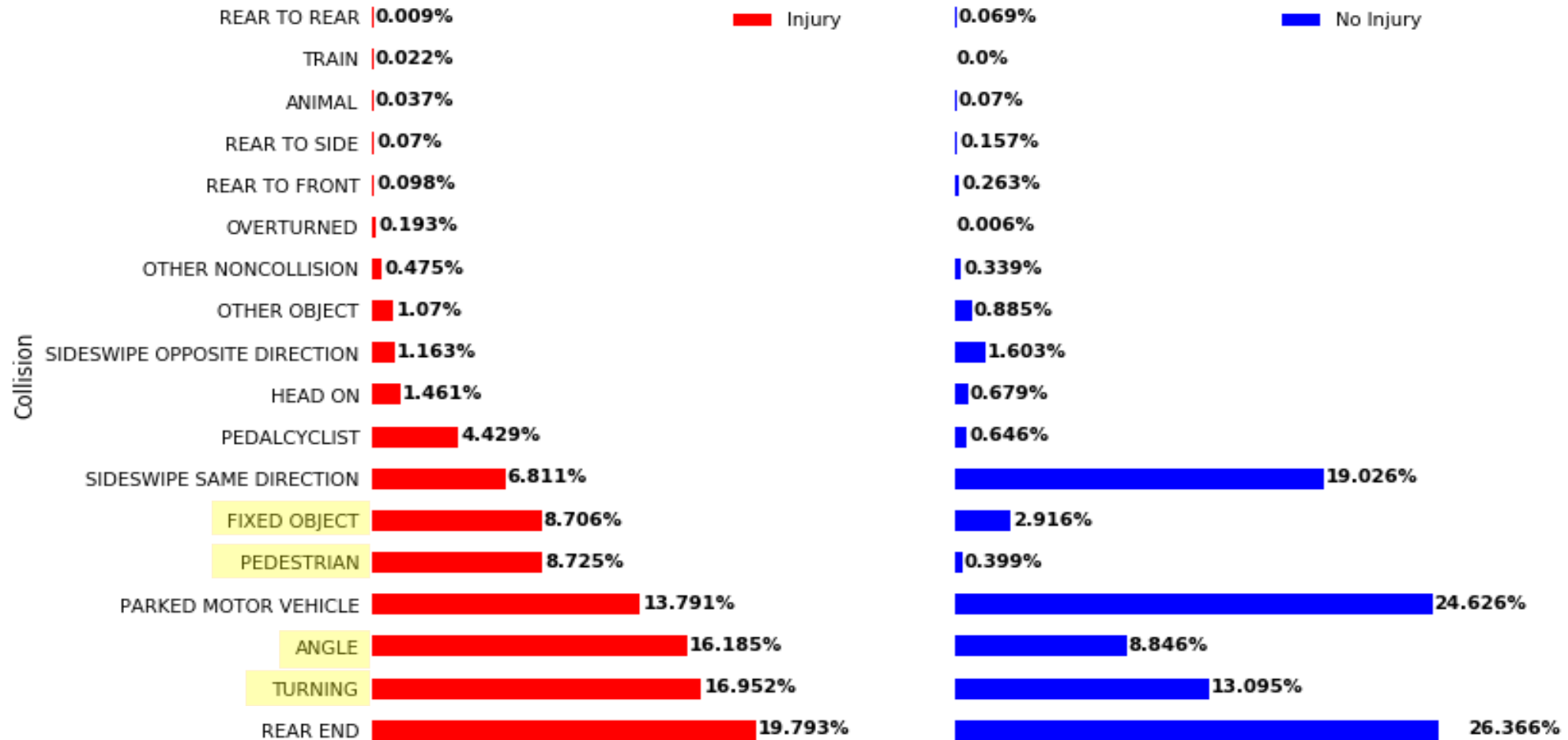
**Model Building**
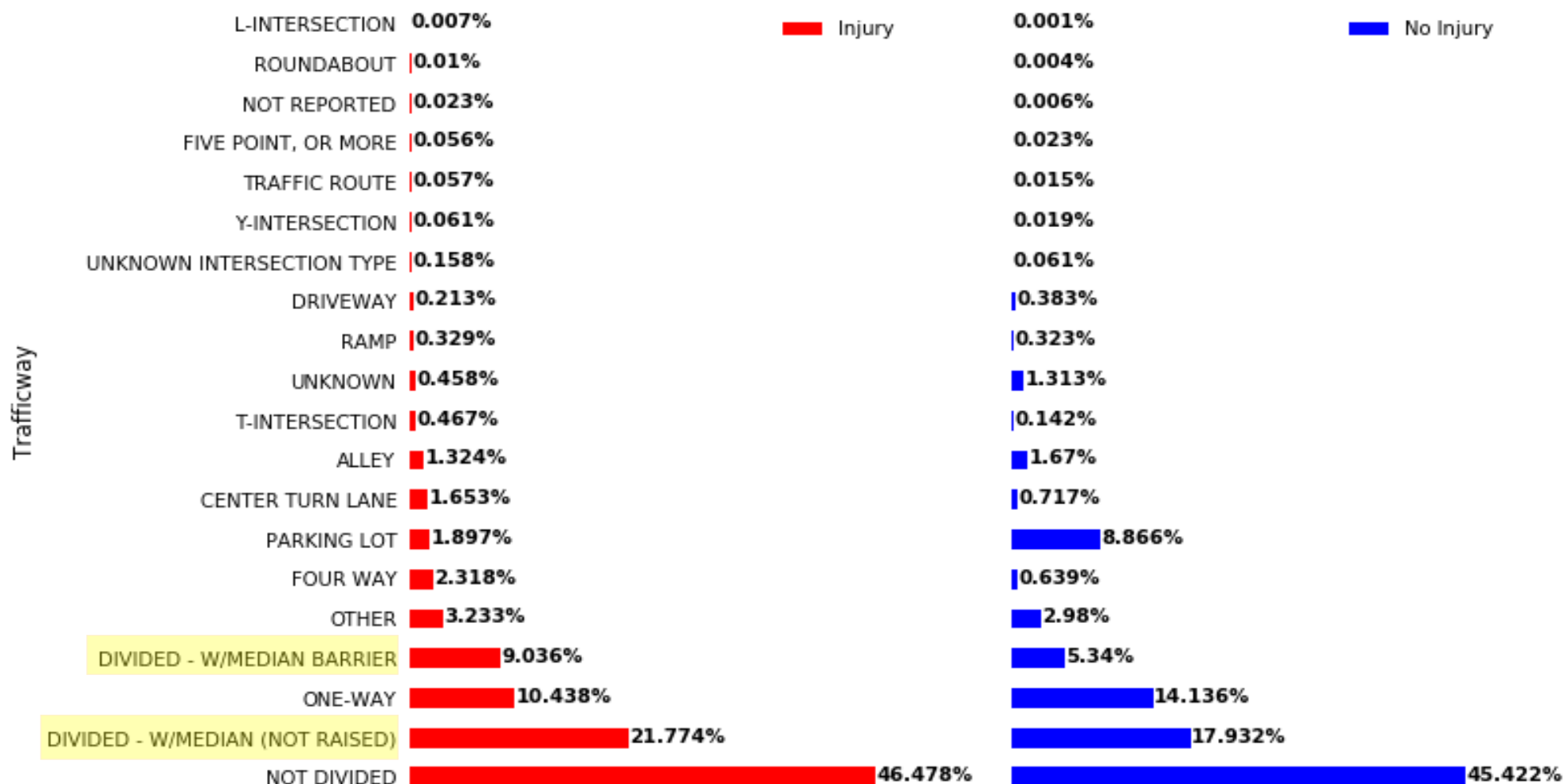
- Train predictive model for crash severity

# Exploratory Analysis

Data Visualization & Statistical Analysis

# Collision Distribution for Each Crash Type

# Trafficway Distribution for Each Crash Type

| Trafficway | Injury | No Injury |
|---|---|---|
| L-INTERSECTION | 0.007% | 0.001% |
| ROUNDABOUT | 0.01% | 0.004% |
| NOT REPORTED | 0.023% | 0.006% |
| FIVE POINT, OR MORE | 0.056% | 0.023% |
| TRAFFIC ROUTE | 0.057% | 0.015% |
| Y-INTERSECTION | 0.061% | 0.019% |
| UNKNOWN INTERSECTION TYPE | 0.158% | 0.061% |
| DRIVEWAY | 0.213% | 0.383% |
| RAMP | 0.329% | 0.323% |
| UNKNOWN | 0.458% | 1.313% |
| T-INTERSECTION | 0.467% | 0.142% |
| ALLEY | 1.324% | 1.67% |
| CENTER TURN LANE | 1.653% | 0.717% |
| PARKING LOT | 1.897% | 8.866% |
| FOUR WAY | 2.318% | 0.639% |
| OTHER | 3.233% | 2.98% |
| DIVIDED - W/MEDIAN BARRIER | 9.036% | 5.34% |
| ONE-WAY | 10.438% | 14.136% |
| DIVIDED - W/MEDIAN (NOT RAISED) | 21.774% | 17.932% |
| NOT DIVIDED | 46.478% | 45.422% |

# Driving Behavior for Each Crash Type



| Cause | Injury | No Injury |
|---|---|---|
| OBSTRUCTED CROSSWALKS | 0.006% | 0.002% |
| MOTORCYCLE ADVANCING LEGALLY ON RED LIGHT | 0.006% | 0.004% |
| PASSING STOPPED SCHOOL BUS | 0.01% | 0.015% |
| RELATED TO BUS STOP | 0.013% | 0.012% |
| BICYCLE ADVANCING LEGALLY ON RED LIGHT | 0.017% | 0.017% |
| DISREGARDING YIELD SIGN | 0.075% | 0.023% |
| TEXTING | 0.077% | 0.036% |
| DISTRACTION - OTHER ELECTRONIC DEVICE (NAVIGATION DEVICE, DVD PLAYER, ETC.) | 0.079% | 0.034% |
| TURNING RIGHT ON RED | 0.095% | 0.057% |
| ANIMAL | 0.119% | 0.073% |
| HAD BEEN DRINKING (USE WHEN ARREST IS NOT MADE) | 0.176% | 0.086% |
| DISREGARDING ROAD MARKINGS | 0.201% | 0.128% |
| CELL PHONE USE OTHER THAN TEXTING | 0.219% | 0.113% |
| ROAD ENGINEERING/SURFACE/MARKING DEFECTS | 0.31% | 0.332% |
| ROAD CONSTRUCTION/MAINTENANCE | 0.31% | 0.26% |
| EVASIVE ACTION DUE TO ANIMAL, OBJECT, NONMOTORIST | 0.337% | 0.137% |
| DISREGARDING OTHER TRAFFIC SIGNS | 0.449% | 0.128% |
| DISTRACTION - FROM OUTSIDE VEHICLE | 0.782% | 0.342% |
| DRIVING ON WRONG SIDE/WRONG WAY | 0.917% | 0.285% |
| EXCEEDING SAFE SPEED FOR CONDITIONS | 0.997% | 0.346% |
| VISION OBSCURED (SIGNS, TREE LIMBS, BUILDINGS, ETC.) | 1.035% | 0.445% |
| IMPROPER BACKING | 1.099% | 5.595% |
| DISTRACTION - FROM INSIDE VEHICLE | 1.211% | 0.59% |
| EXCEEDING AUTHORIZED SPEED LIMIT | 1.386% | 0.343% |
| EQUIPMENT - VEHICLE CONDITION | 1.438% | 0.315% |
| PHYSICAL CONDITION OF DRIVER | 1.606% | 0.238% |
| UNDER THE INFLUENCE OF ALCOHOL/DRUGS (USE WHEN ARREST IS EFFECTED) | 1.65% | 0.136% |
| WEATHER | 2.175% | 1.441% |
| OPERATING VEHICLE IN ERRATIC, RECKLESS, CARELESS, NEGLIGENT OR AGGRESSIVE MANNER | 2.219% | 0.868% |
| IMPROPER OVERTAKING/PASSING | 2.309% | 5.637% |
| DISREGARDING STOP SIGN | 2.395% | 0.637% |
| DRIVING SKILLS/KNOWLEDGE/EXPERIENCE | 2.597% | 3.317% |
| IMPROPER LANE USAGE | 2.745% | 4.459% |
| IMPROPER TURNING/NO SIGNAL | 3.671% | 3.307% |
| NOT APPLICABLE | 4.231% | 5.786% |
| DISREGARDING TRAFFIC SIGNALS | 4.59% | 0.77% |
| FAILING TO REDUCE SPEED TO AVOID CRASH | 6.955% | 3.148% |
| FOLLOWING TOO CLOSELY | 8.551% | 12.022% |
| FAILING TO YIELD RIGHT-OF-WAY | 17.285% | 9.476% |
| UNABLE TO DETERMINE | 25.636% | 39.022% |

# Control Device for Each Crash Type

| Device | Injury | No Injury |
|---|---|---|
| NO PASSING | 0.005% | 0.002% |
| BICYCLE CROSSING SIGN | 0.006% | 0.001% |
| FLASHING CONTROL SIGNAL | 0.011% | 0.008% |
| PEDESTRIAN CROSSING SIGN | 0.035% | 0.013% |
| POLICE/FLAGMAN | 0.035% | 0.041% |
| DELINEATORS | 0.04% | 0.019% |
| SCHOOL ZONE | 0.064% | 0.033% |
| RAILROAD CROSSING SIGN | 0.1% | 0.094% |
| OTHER REG. SIGN | 0.12% | 0.08% |
| OTHER WARNING SIGN | 0.159% | 0.062% |
| YIELD | 0.217% | 0.116% |
| LANE USE MARKING | 0.577% | 0.294% |
| OTHER | 0.645% | 0.565% |
| UNKNOWN | 1.937% | 3.437% |
| STOP SIGN/FLASHER | 14.419% | 8.269% |
| TRAFFIC SIGNAL | 32.431% | 26.608% |
| NO CONTROLS | 49.19% | 60.35% |

Legend: ■ Injury (red)  ■ No Injury (blue)

# Lighting Conditions for Each Crash Type

# Speed Limit Distribution for Each Crash Type

# Weather Condition for Each Crash Type

■ Injury    ■ No Injury

| Weather | Injury | No Injury |
|---|---|---|
| BLOWING SNOW | 0.001% | 0.0% |
| SEVERE CROSS WIND GATE | 0.023% | 0.021% |
| FREEZING RAIN/DRIZZLE | 0.026% | 0.016% |
| FOG/SMOKE/HAZE | 0.269% | 0.163% |
| SLEET/HAIL | 0.28% | 0.128% |
| CLOUDY/OVERCAST | 3.89% | 2.773% |
| SNOW | 3.983% | 3.585% |
| RAIN | 11.747% | 8.755% |
| CLEAR | 79.776% | 84.555% |

# Number of Injury Crashes per Hour

# Number of No Injury Crashes per Hour

# Association Measure between Crash Features and Crash Type

| Crash Feature | Cramer's V Coefficient |
|---|---|
| Collision | 0.354896 |
| Primary Cause (Driving Behavior) | 0.305708 |
| Trafficway Type | 0.166856 |
| Crash Hour | 0.140647 |
| Lighting Conditions | 0.132704 |
| Posted Speed Limit | 0.127492 |
| Control Device | 0.121500 |
| Road Surface Condition | 0.064132 |
| Weather Condition | 0.057022 |
| Road Alignment | 0.056438 |
| Crash Day | 0.039466 |
| Crash Month | 0.022364 |

# Kernel Density Estimation of Crashes' Location

Crashes with injury locations

Crashes with no injury locations

# Crashes' Distribution Across the Community Areas of Chicago



Distribution of Injury Crashes Across Chicago Areas

Distribution of No Injury Crashes Across Chicago Areas

# Key Findings

The features that have the strongest association with crash severity are: *driving behavior* and *type of collision*.

- Example of such behaviors: disregarding traffic signal and stop sign, failing to reduce speed, and not giving the right-of-way.

- Example of collisions: angle or turning collision, and collisions with pedestrians.

# Key Findings

- Crashes with Injury do not only happen at rush hours during weekdays, but also during early morning hours of the weekend.

- Injury Crashes are not only located in the central part of Chicago but also in the west side of the city.

# Building the Predictive Model

In-Depth Analysis

# Additional Preprocessing

*Additional Preprocessing Steps*:

- Added one column "Area" to designate the Chicago area of the crash

- Converted nominal categorical features into numerical entries:
  Tried different encoding schemes (one hot, binary, leave-one-out)

# Data Imbalances

- No Injury crashes: 78% ,    Injury crashes: 22%

- To address imbalances:
  - Metrics used: precision, recall and F1-score
  - Considered under-sampling of the majority class

# Steps of Model Building

- Split the data into training and testing sets:

| 80% | 20% |
|:---:|:---:|

- Performed 5-fold cross validation on the training set:

**5-fold cross validation**

   - train different combination of: encoding scheme and training model, with/out under-sampling
   - select the final model using F1-score

- Tested the final chosen model on the testing set:

# Training Models

We tried various models and compared their performance:

- Naïve Bayes

- Logistic Regression

- Linear SVM

- Random Forest

- Ada Boost

- Gradient Boosting

- Balanced Random Forest

# Final Model and its Performance

- Final Model Selected (trained after under-sampling):

  **Gradient Boosting** (parameters: *n_estimators=600, max_depth=4*)

  + **Leave-one-out Encoding**

- Performance on the testing set:

| | |
|---|---|
| **Accuracy** | **0.747** |
| **Precision** | **0.697** |
| **Recall** | **0.457** |
| **F1-Score** | **0.552** |

# Features' Importance

| Crash Feature | Importance |
|---|---|
| Collision | 0.418 |
| Primary Cause (Driving Behavior) | 0.287 |
| Trafficway Type | 0.071 |
| Area | 0.054 |
| Crash Hour | 0.042 |
| Posted Speed Limit | 0.041 |
| Lighting Conditions | 0.033 |
| Control Device | 0.016 |
| Crash Month | 0.0086 |
| Crash Day | 0.0066 |
| Road Alignment | 0.0062 |
| Weather Condition | 0.0059 |
| Road Surface Condition | 0.0058 |

# Possible Future Works

*Additional Work*:

- Incorporate more features related to driver and vehicle's information
- Consider stacking of the models
- Account for location in terms of zip code instead of the code area
- Perform streets segments analysis
- Focus on crashes with injury and analyze the conditions of possible types of injuries (fatal, incapacitating and non-incapacitating)

# Recommendations

**Driving behavior**: important feature in predicting severity of crashes especially at *intersections.*

Efforts should be focused on:

- Pushing drivers to drive cautiously and carefully
- Keeping on educating drivers of defensive driving techniques
- Helping drivers staying alert while driving