

# Machine Learning for Predicting Severity of Collisions

## 1 Introduction

### 1.1 Background

Would it not be great if you could tell in advance how severe an accident will be if it occurs?! In this project, we will use a car-accident data of Seattle city in the Washington State, US from 2004 to 2013. The data can be downloaded from [1].

The data contains many useful features such as weather and road conditions. The severity of the accident is classified into several classes ranging from property damage to fatality. This problem is interesting for the Departments of Transportation especially in the Seattle City. It can help the department linking weather, road and light conditions, certain locations and certain streets to very severe accidents.

### 1.2 Problem Statement

Car accidents are one of the leading causes of death in the US. In this project, we set out to predict the severity of car accidents using machine learning techniques and real-world data.

It is expected that the machine learning classifier will be able to classify the severity of an accident based on its attributes. Furthermore, the exploratory data analysis will shed light on different hidden patterns in the data that can be useful for decision makers.

## 2 Data

### 2.1 Data Understanding

The data contains information about the car accidents occurred in the Seattle city in the US. The data contains 221144 records and 40 attributes (features). The values of the target label (the severity of accidents) are described in Table 1.

Value	Description
0	Unknown
1	Property damage
2	Injury
2b	Serious injury
3	Fatality

Table 1: Target label.

## 2.2 Data Cleaning

1. Features with 50% or more missing values are dropped since imputing their missing values will alter the real distribution of the data significantly and therefore, can be misleading to the machine learning algorithm.
2. Missing values are replaced using a roulette-wheel replacement method that preserve the distribution of the data.

## 2.3 Data Wrangling & Feature Engineering

The following techniques are used to generate new, hopefully more informative, features out of the original features.

1. **Binning** is used to group values that are invariant to the target label to reduced the dimensions of the data.
2. **One-hot encoding** is used to convert categorical variables into binary columns.
3. **Binary encoding** is used when one-hot encoding is not practical (such as when there are too many levels for a categorical variable).

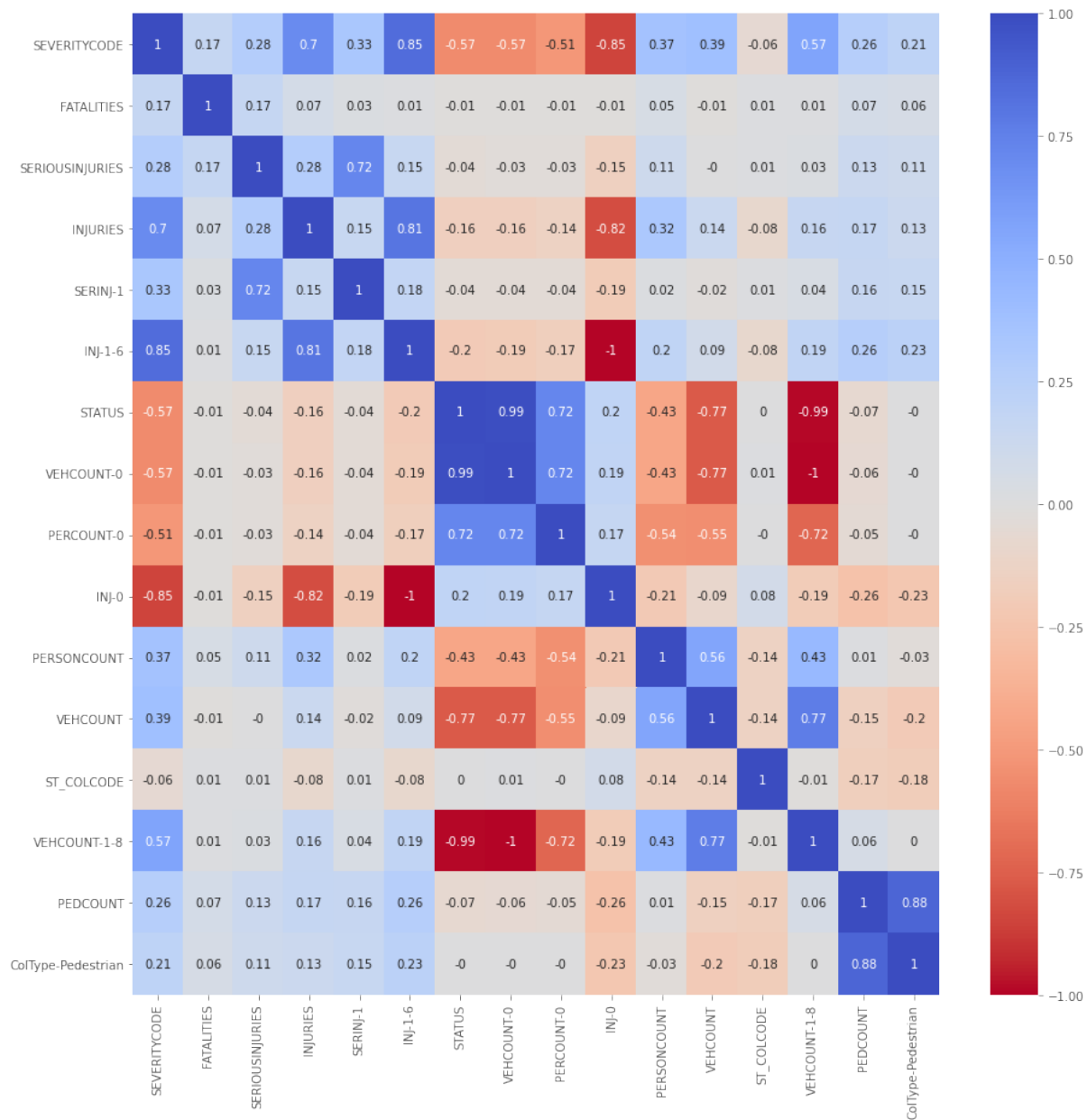
# 3 Methodology

## 3.1 Feature Selection

After the feature engineering step, we end up with 81 features that can be used in modeling. Feeding all these features to the model is likely to degrade the performance of the model since it will incur a very high number of parameters to be learned. Therefore, we perform feature selection as follows:

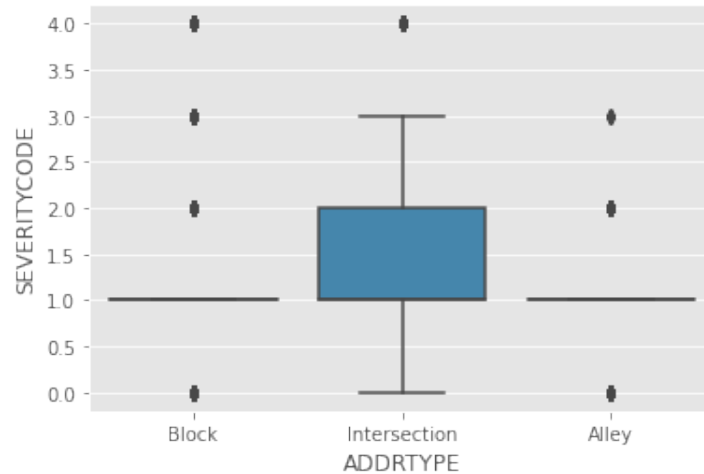
1. We use the **chi test** together with the correlation matrix to select the top features.
2. The feature selection is performed by considering the top positively and negatively correlated features using a correlation matrix.

The correlation matrix of the 15 features that are used in the model.

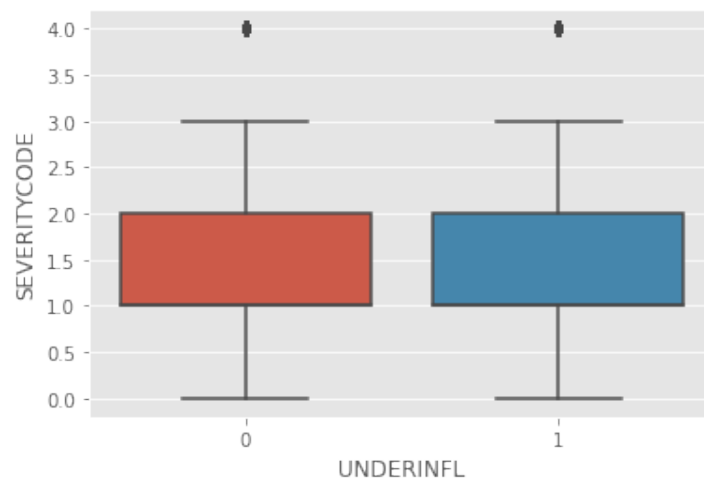


### 3.2 Exploratory Data Analysis

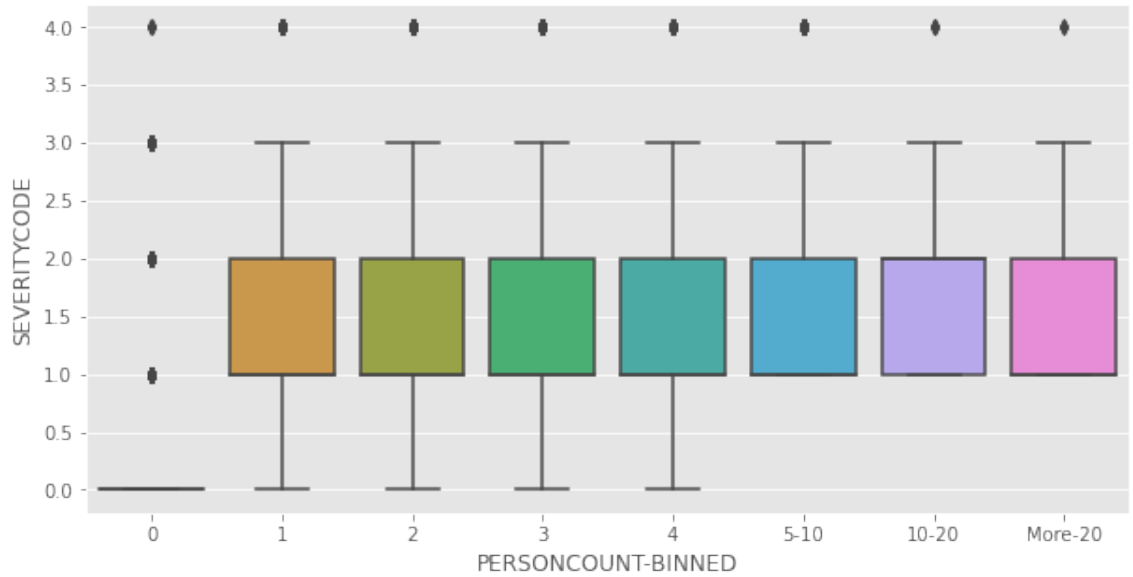
Accidents occurring at intersections are more severe.



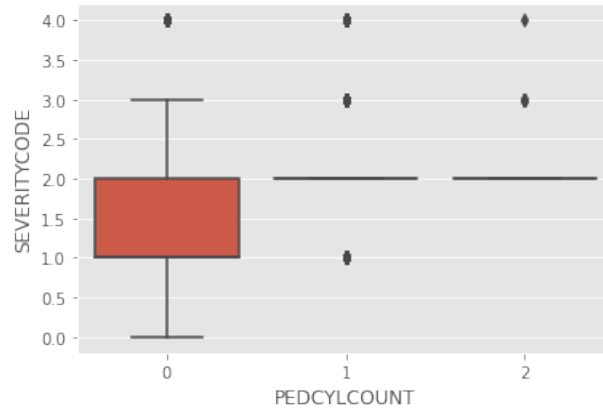
Whether the driver was under the influence of alcohol or not, the severity of the accident remains the same which counter our intuition. This could possibly be as a result of trying to drive more carefully when under alcohol influence. This is true for those who are not totally drunk.



An accident involving pedestrians is more severe. If 5 or more pedestrians are involved, the accident tends to be even more severe.



Accidents involving bicyclers are more sever.



### 3.3 Modeling

The *random forest classifier* is used due its ability to deal with imbalanced classes without the need to *up-sampling* or *down-sampling* the original data. This is the case since the decision trees have hierarchical structures that enables them to effectively distinguished different classes regardless of the actual number of classes.

## 4 Result & Discussion

The predictive models that we built achieves an accuracy of 99% which is very high and is able to predict all five classes despite the fact that the problem is highly imbalanced. This ensures that the high accuracy of the model is not due to being very good at predicting the majority class. On contrary, the model

high predictive power is due to being able to distinguish all five classes. Random forest classifiers is used since its hierarchical structures enables decision trees to effectively distinguish different classes even if the data is imbalanced.

**Important Observations:**

1. Surprisingly, the geographic data tends to be insignificant for this problem. The target label is very weakly correlated with the locations at which the collisions occur.
2. The under-influence-of-alcohol state does not influence the severity of the accident which counters our intuition. In particular, the data reveals that the distribution of the severity of the accident remains invariant with the under-influence-of-alcohol state.
3. The poor weather and road conditions do not cause sever accidents. This could possibly be explained by the extra care exercised by the drivers under poor weather or road conditions.
4. Accidents at the intersection are more sever than those in alleys and blocks.

**Recommendations:**

1. The most sever accident involved pedestrians and bicyclers. Therefore, better policies to ensure the safety of pedestrians and bicyclers are needed.
2. Accidents at intersection are more sever than those in alleys and blocks. This needs to be investigated physically to see whether these intersections have working traffic lights. If not, perhaps, installing new ones will reduce the number of accidents.
3. Though very important, car speeding is mostly missing! This could be due to the difficulty of measuring the exact speed just before the accidents. However, we believe it is a very important predictive feature that should be better recorded.

## 5 Conclusion

The project aims at predicting the severity of car accident using data acquired from the Seattle City, US. The data is rich with 40 features and 221144 records.

The problem we dealt with is imbalanced as the number of records across classes. Therefore, we carefully chose the random forest classifiers owing to its power in distinguishing different imbalanced classes as this study showed. We ignored the resampling technique to balance the classes since the up-sampling will unacceptably increase the size of the data and the down-sampling will reduce it significantly.

The model is highly accurate (achieved an accuracy of 99%) and is able to predict all five classes which ensured that the high accuracy is not solely produced by predicting the majority class.

## Reference

[1] Data Source: <https://data.seattle.gov>