

# IBM CAPSTONE PROJECT “CAR ACCIDENT SEVERITY” REPORT

## Introduction

The costs of fatalities and injuries due to traffic accidents have a great impact on the society. In recent years, researchers have paid increasing attention to determining factors that significantly affect severity of driver injuries caused by traffic accidents.

Annual Global Road Crash Statistics:

- Approximately 1.35 million people die in road crashes annually, on average 3,700 people lose their lives every day on the roads.
- An additional 20-50 million suffer non-fatal injuries, often resulting in long-term disabilities. •
- More than half of all road traffic deaths occur among vulnerable road users—pedestrians, cyclists, and motorcyclists. •
- Road traffic injuries are the leading cause of death among young people aged 5-29.
- Young adults aged 15-44 account for more than half of all road deaths. •
- On average, road crashes cost countries 3% of their gross domestic product.

## Business Understanding

For problem understanding let's see the situation: You are driving to another city for work or to visit some friends. It is rainy and windy, and on the way, you come across a terrible traffic jam on the other side of the highway. Long lines of cars barely moving. As you keep driving, police car start appearing from afar shutting down the highway. Oh, it is an accident and there's a helicopter transporting the ones involved in the crash to the nearest hospital. They must be in critical condition for all of this to be happening. Now, wouldn't it be great if there is something in place that could warn you, given the weather and the road conditions about the possibility of you getting into a car accident and how severe it would be, so that you would drive more carefully or even change your travel if you are able to.

The case study is to predict the severity of an accident.

To reduce the frequency of car collisions we have to develop a model to predict the severity of an accident given the current weather, road and visibility conditions. When conditions are bad, this model will say to be more careful.

The machine learning model should be able to predict accident "severity".

## Data Understanding

The dataset contains 194673 observations (rows) and 37 attributes (columns). The machine learning model should be able to predict accident "severity". The target of prediction is 'SEVERITYCODE' (it is used to measure the severity of an accident). In the dataset there are only 2 variants (1 - prop damage and 2 - injury). But the target or label columns should be accident "severity" in terms of human fatality, traffic delay, property damage, or any other type of accident bad impact. The attributes we can use to predict the severity of an accident are 'WEATHER', 'ROADCOND' and 'LIGHTCOND'.

## Data Preparation:

This dataset is not fit for analysis perfectly. We should not use all attributes for our model.

Most of the attributes are text-type, so we should convert them to a numerical type. We should use label encoding to convert the features.

The target variable SEVERITYCODE is only 42% balanced. the quantity of severitycode in class 1 is 136485 and the class 2 is 58188. We can fix this by downsampling the class 1.

There are a lot of empty fields - fill them by zero.

Features selected (X):

Feature	Description	Reason for Dropping
X	Latitude	Can't be modelled in classification
Y	Longitude	Can't be modelled in classification
OBJECTID	ESRI unique identifier	ID not relevant
INCKEY	Secondary key for the incident	ID not relevant
COLDKEY	Identifying key	ID not relevant
LOCATION	Description of Location	ADDRTYPE captures this
REPORTNO	Report Number	ID not relevant
STATUS	Matched/Unmatched	ID not relevant
INTKEY	Intersection key for collision	ID not relevant
EXCEPTRSNCODE	Blank	No data
EXCEPTRSNDESC	Blank	No data
SEVERITYCODE	Label	Label to be predicted
SEVERITYDESC	Description of Severity	Label to be predicted
INCDATE	The date of the incident.	INCDTTM captures this
SDOT_COLCODE	Collision code	Collision type captures this
SDOT_COLDESC	A description of the collision corresponding to the collision code.	Collision type captures this
SDOTCOLNUM	A number given to the collision by SDOT.	Collision type captures this
SEGLANEKEY	A key for the lane segment in which the collision occurred.	ID not relevant
CROSSWALKKEY	A key for the crosswalk at which the collision occurred.	ID not relevant

Features dropped:

Feature	Description	Reason for Selecting
ADDRTYPE	Collision at Alley, Block, Intersection	Gives the likelihood of collision at these places
PERSONCOUNT	Number of people involved in the collision	Gives an indication of severity
PEDCOUNT	Number of pedestrians involved in the accident	Gives an indication of severity
PEDCYLCOUNT	Number of cyclists involved in the accident	Gives an indication of severity
VEHCOUNT	Number of vehicles involved in the accident	Gives an indication of severity
INCDTTM	The date and time of the incident	Time of accident: midnight/ day time
INATTENTIONIND	Whether the person was not paying attention	Not paying attention can result in accident
UNDERINFL	Whether the person was driving under influence	DUI can cause accidents
WEATHER	Weather conditions	Bad weather can cause accidents
ROADCOND	Road conditions	Wet roads can cause skidding
LIGHTCOND	Light conditions	Light conditions affect visibility
PEDROWNOTGRNT	Pedestrian right of way was granted or not	
SPEEDING	Whether speeding or not	Speeding causes accidents
COLLISIONTYPE	Collision Type	Type of collision gives severity of accident
HITPARKEDCAR	Whether or not the collision involved hitting a parked car.	Hitting a parked car causes property damage

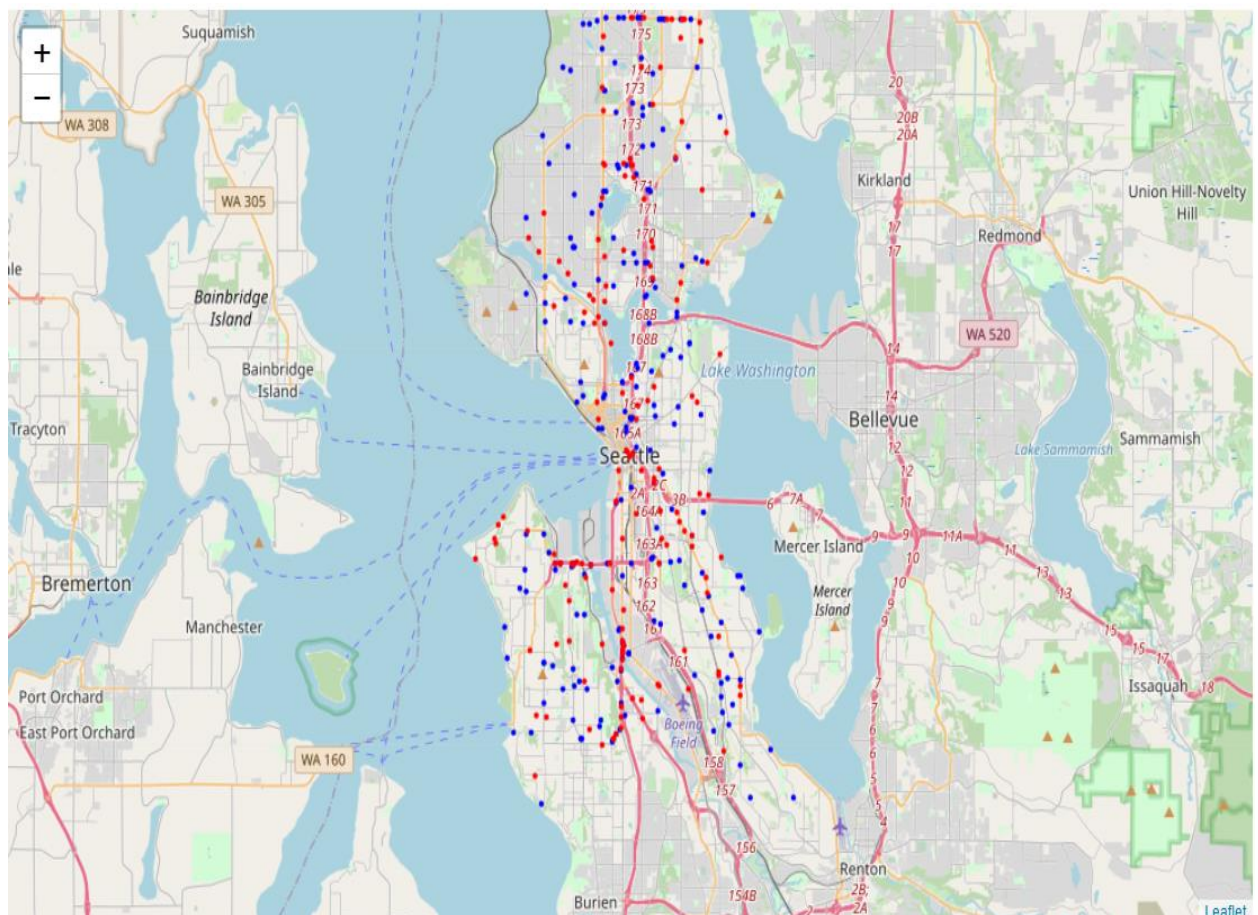
# Data Analysis

Plotting factors on the map to get density of areas where accidents were caused by the features in question:

1. Speeding:
2. Under Influence (DUI)
3. Inattention
4. Hitting a parked car

## Plotting density of accidents sorted by Severity caused by Speed

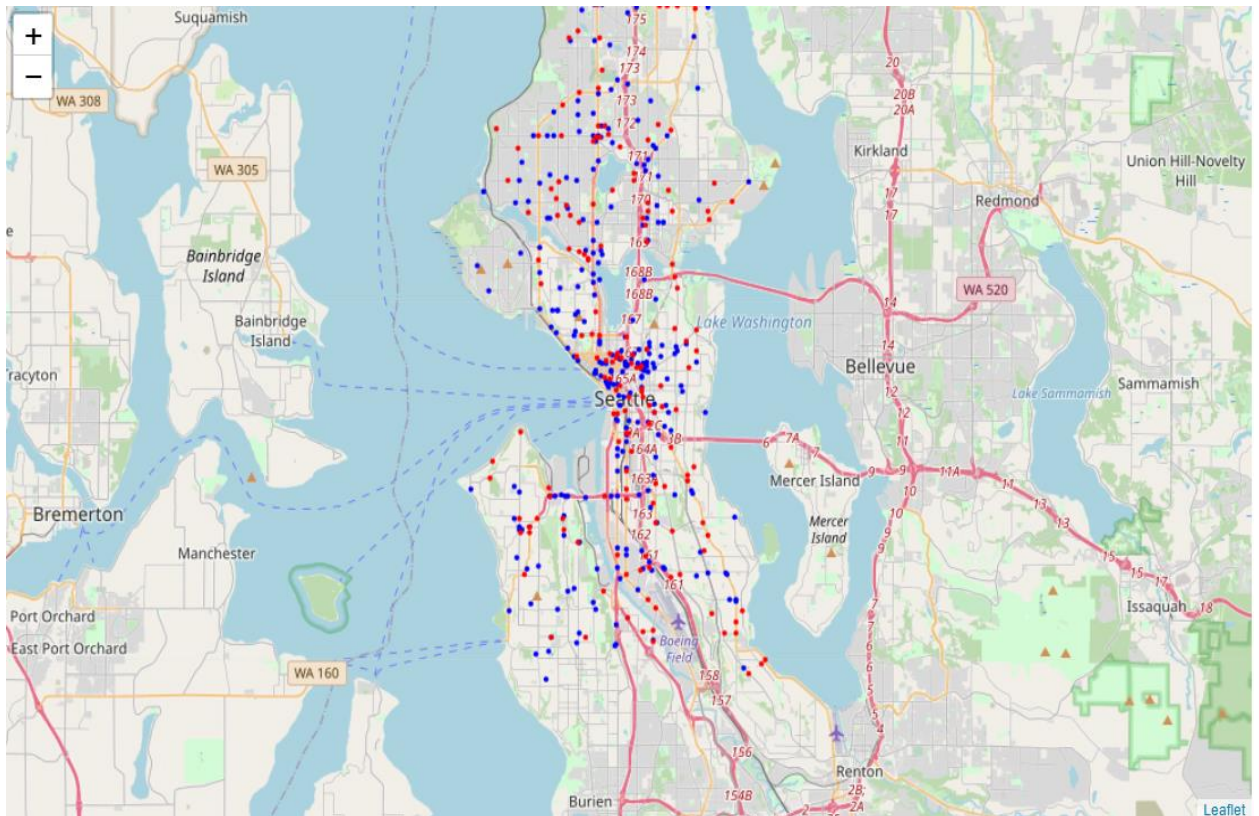
Certain roads have a lot of accidents which occur on them due to speeding. The government of Seattle can introduce proper traffic management in the form of speed restricting interventions (e.g. speed bumps). This can cause reduction in accidents due to speeding.



## Plotting density of accidents sorted by Severity caused by Driving Under Influence

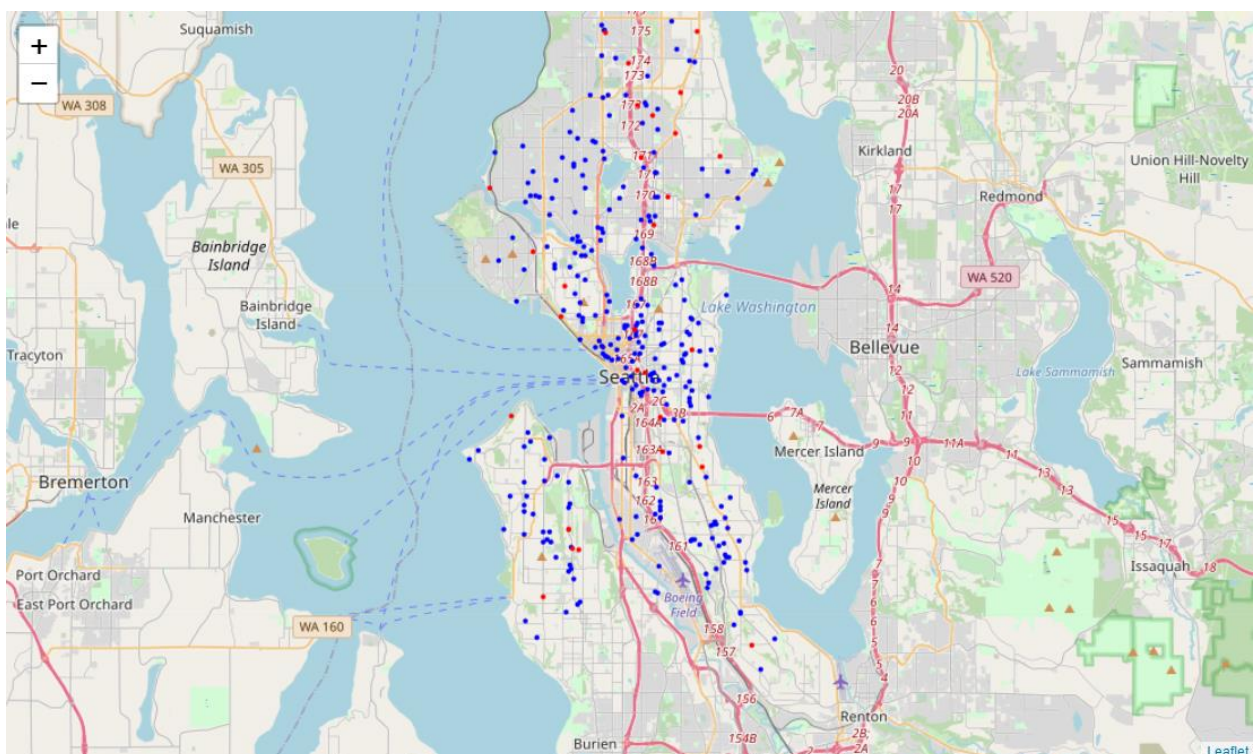
The above map shows, the points where accidents are caused due to DUI. The Seattle government can introduce Police check-ups on vehicles which are entering nodes where one has high density of accidents caused by DUI. This can reduce potential accidents before they happen.





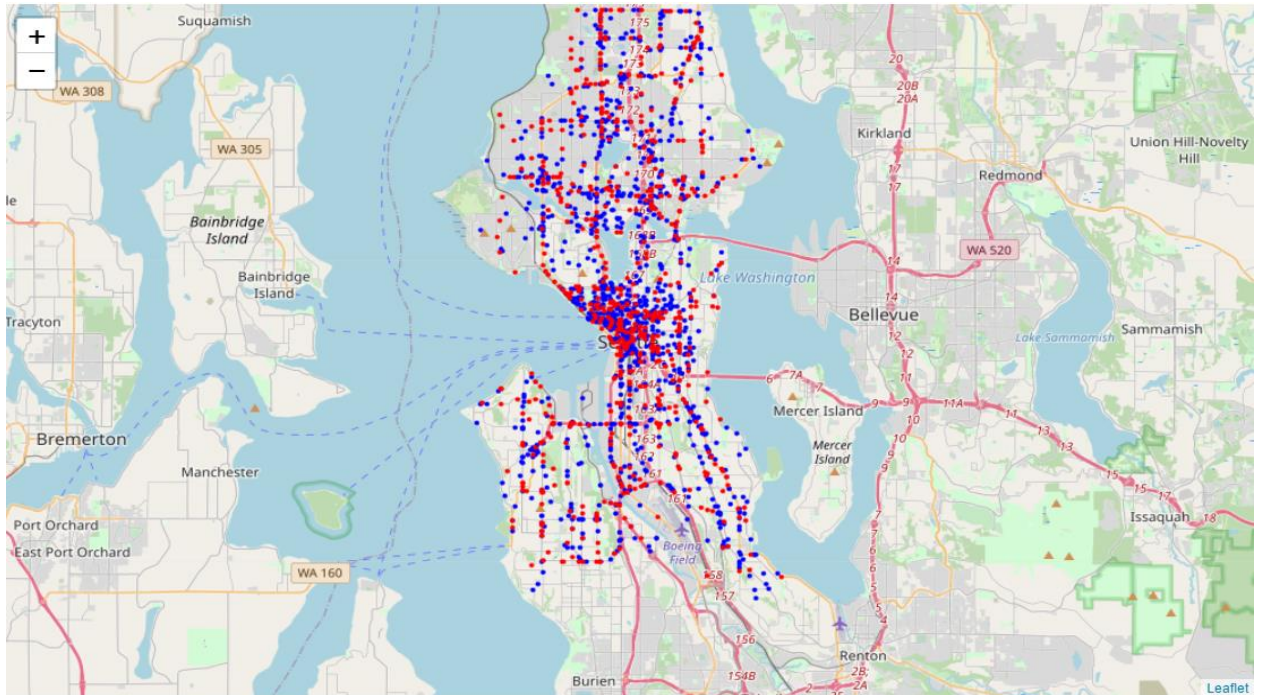
## Plotting density of accidents sorted by Severity caused by Hitting Parked Cars

We can see places where accidents in which parked cars were hit. The areas can be used by Insurance companies to tweak their car insurance premiums for individuals living in those areas.



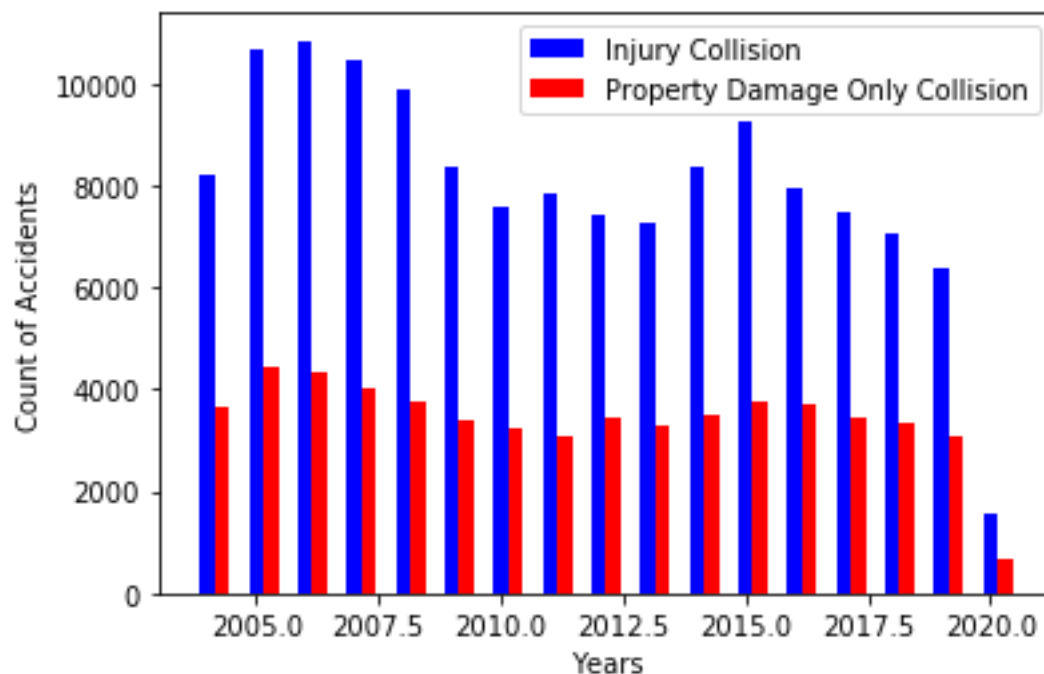
## Plotting density of accidents sorted by Severity caused by Inattention

It shows that a huge majority of accidents are caused by inattention. Perhaps a product monitoring the attention of drivers can be developed/marketed citing this data.



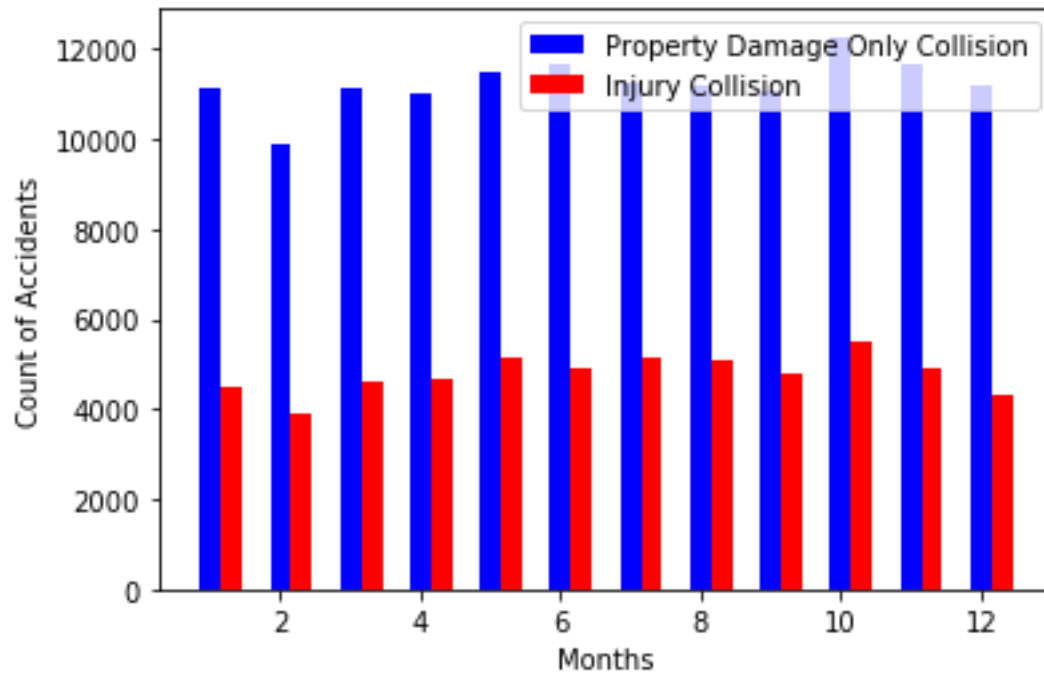
## Plotting count of accidents based on the following factors:

### 1. Year:



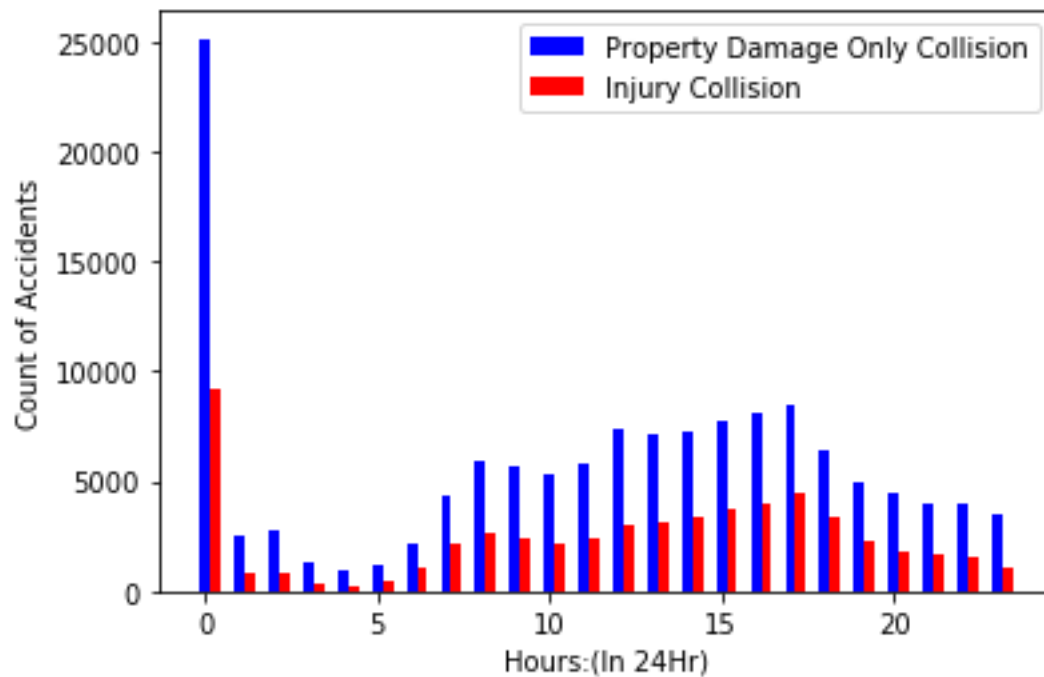
1. The Number of accidents in both the Severity classes have been decreasing over the years.
2. The drastic drop in 2020 is due to there being data for part of the year.

## 2. Months



October has the most number of crashes in the year. Number of crashed further decreases in the months of November and then December.

## 3. Hours in 24-hour format



The highest number of accidents by far happen at midnight from 12AM to 1 AM.



# Methodology

Our data is now ready to be fed into machine learning models.

We will use the following models:

**K-Nearest Neighbor (KNN)** KNN will help us predict the severity code of an outcome by finding the most similar to data point within k distance.

**Decision Tree** A decision tree model gives us a layout of all possible outcomes so we can fully analyze the consequences of a decision. In context, the decision tree observes all possible outcomes of different weather conditions.

**Logistic Regression** Because our dataset only provides us with two severity code outcomes, our model will only predict one of those two classes. This makes our data binary, which is perfect to use with logistic regression.

Also we will use GradientBoosting, XGBClassifier, RandomForest and Support Vector Machine.

Exploratory data analysis was performed on the relevant categorical variables: address type, collision type, person count, pedestrian count, cyclist count, vehicle count, junction type, SDOT type, under the influence, weather, road conditions, light conditions, lane key, crosswalk key, and if a parked car was hit. The amount of categories in each variable ranged from a few to over a dozen. As a result, categories with less than a dozen variables were turned into countplots to better visualize the data. Statistical testing was not performed because the data revolved around categorical variables, not numerical ones. Unfortunately, key variables such as pedestrian right of way, inattentive drivers, and whether the car was speeding had a majority of null values. Therefore, they were dropped and not part of the analysis. However, it is likely that these variables play a key factor in vehicle accidents.

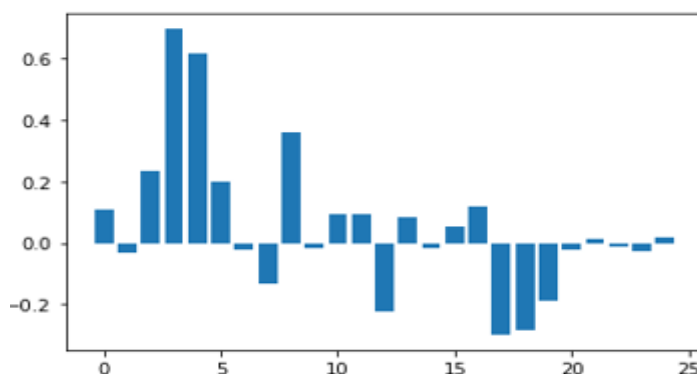
# Modeling

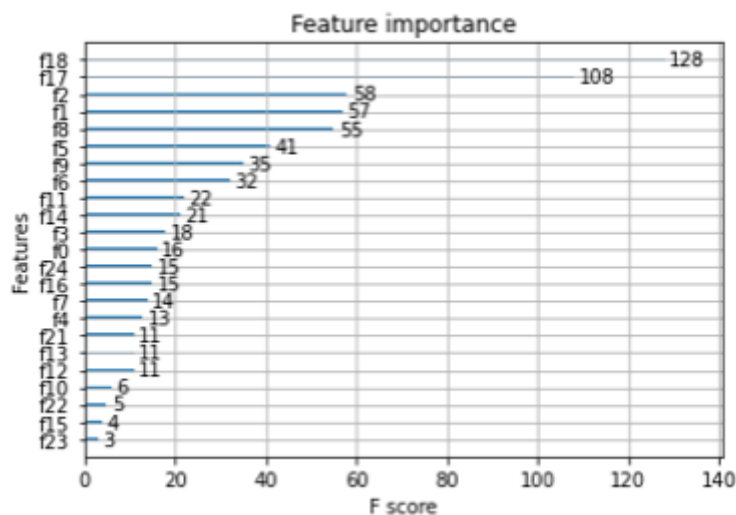
Classification models (Decision tree, K neighbors and logarithmic regression) are supervised learning algorithms that seek to classify data into a set of discrete value classes. This type of algorithm learns from the relationship between a set of categorical variables and a variable of interest, in which the latter is a categorical variable with discrete values.

Logistic regression aims to classify data from a set considering the input variables in order to predict a target variable which must be binary. It is important to note that the independent variables must be continuous.

The K-Nearest Neighbors algorithm is a classification method which groups diverse points with common characteristics to learn from the relationships between them to label unknown information. Those data that are close to each other are called neighbors.

## Feature Importance Analysis





## Results

Interestingly, the data shows that most accidents occur during the day with normal drivers and conditions. Most of the cases involved property damage. As expected, most occur at a block or intersection. Most vehicle accidents occur during the best driving times when it is clear, during the day, and the roads are dry.

*The Results of the classification are as follows:*

Model	Accuracy
K Nearest Neighbours	0.7453191216129447
Logistic Regression Accuracy	0.756876846025427
DecisionTrees	0.6912032875304995
<b>GradientBoosting</b>	<b>0.7662771285475793</b>
XGBoost	0.7648131501219982
Random Forest's Accuracy	0.7630152818800565
Support Vector Machine	0.7629639142160011

Very good results are from **GradientBoosting (best result)**, XGBoost, Random Forest and Support Vector Machine.

## Discussion

The first striking observation has to do with the dependent variable. It seems highly unlikely that over the date range of the data no serious injuries or fatalities occurred. This may be a warning sign that the severity codes were somehow altered when the data set was being created, or that the sample data is incomplete and missing those reports. The main recommendation has to do with key variables such as pedestrian right of way, inattentive drivers, and if the car was speeding. In many of the records, these values were null. However, this data should be collected in order to draw new insights or create better prediction models. It was determined that most accidents occur during normal weather and road conditions. However, further data is needed to analyze this trend. It may be that these types of days constitute the highest number of days in the year. Therefore, further data on the weather needs to be analyzed. For example, it may be sunny for 100 days and then snow on 1 day. If you look at data for accidents, there may be 1000 accidents occur during sunny days



and only 20 on snowy days. Really, the average number of accidents per weather type is much higher on snowy days. Though, because there are few days like that during the year the total number of accidents appears low.

Once we analyzed and cleaned the data, it was then fed through three ML models; K-Nearest Neighbor, Decision Tree and Logistic Regression. Although the first two are ideal for this project, logistic regression made most sense because of its binary nature. Evaluation metrics used to test the accuracy of our models were jaccard index and f-1 score. Choosing different k, max depth and hyperparameter C values helped to improve our accuracy to be the best possible.

## **Conclusion**

The data-set has been used to classify the severity of the accidents based on certain select features. The exploratory data analysis shows density of accidents based on geography based on Speeding, Driving Under Influence, Inattention and Hitting Parked Cars. From a machine learning standpoint. The most important features were: Collision Type, Person Count, Vehicle Count and Address Type. The Gradient Boost algorithm performed the best.

The data showed that most vehicle accidents occur during good conditions with normal drivers. This means it will be harder for the Seattle transportation department to mitigate accidents. However, as most accidents only involve property damage or minor injuries, there is not a serious problem that needs to be dealt with right away. This shows that infrastructures are being designed and operating properly. Therefore, the focus should be an emphasis on drivers being more careful.

Based on historical data from weather conditions pointing to certain classes, we can conclude that particular weather conditions have a somewhat impact on whether or not travel could result in property damage (class 1) or injury (class 2).