

# Report: Predicting accident severity in Seattle

## IBM Capstone Project – Applied Data Science

*C. Cotton*

*October 4, 2020*

### Table of Contents

<b>IBM CAPSTONE PROJECT – APPLIED DATA SCIENCE .....</b>	<b>1</b>
INTRODUCTION:.....	2
BUSINESS UNDERSTANDING: .....	2
DATA SOURCE & DESCRIPTION: .....	3
METHODOLOGY:.....	3
<i>1. Data evaluation .....</i>	<i>3</i>
<i>2. Data Exploration.....</i>	<i>4</i>
<i>3. Data Pre-processing.....</i>	<i>6</i>
<i>4. Model Development and Evaluation .....</i>	<i>6</i>
CONCLUSIONS: .....	7

### Introduction:

The number of accidents has been rising globally due to increases in population and transportation. Traffic accidents are a daily source of death, injury and property damage on roads and highways resulting in huge losses at economic and social levels.

According to the World Health Organisation (WHO), in 2018<sup>1</sup> around 1.35 million died as a result of a traffic accident.

If we look more precisely at the figures for Seattle, in 2017<sup>2</sup>, there were 10,959 police reported collisions and a further 1516 self-reported ones.

Fatalities were 19, and serious injuries 168.

### Business Understanding:

As the demand for vehicles rises, the number of vehicles on the road and traffic jams increase, especially during rush hours.

The local government of Seattle is trying to implement measures to alert vehicle users, police, traffic and health systems about critical situations – to reduce the number of accidents on the road. I will attempt to build a model to predict the severity of an accident given the weather, road conditions and location.

By doing this, we should be able to predict the severity of an accident.

This will be useful to the traffic department also to implement measures to decrease the risk of an accident.

### Sources:

1. WHO Global status report on road safety 2018: <https://apps.who.int/iris/bitstream/handle/10665/277370/WHO-NMH-NVI-18.20-eng.pdf?ua=1>

2. Seattle Department of Transport 2018 Traffic Report: [https://www.seattle.gov/Documents/Departments/SDOT/About/DocumentLibrary/Reports/2018\\_Traffic\\_Report.pdf](https://www.seattle.gov/Documents/Departments/SDOT/About/DocumentLibrary/Reports/2018_Traffic_Report.pdf)

## Data source & description:

The data used in this study was provided by Coursera, and had records on all types of collisions from 2004 through to present day.

The data is regarding the **severity of each accident** along with the time and conditions under which each accident occurred.

The model aims to predict the severity of an accident, considering that, the variable of Severity Code was in the form of 1 (Property Damage Only) and 2 (Physical Injury) which were encoded to the form of 0 (Property Damage Only) and 1 (Physical Injury). Following that, 0 was assigned to the element of each variable which can be the least probable cause of severe accident whereas a high number represented adverse condition which can lead to a higher accident severity.

As there are null values in some records, the data needs to be pre-processed before any further processing.

The data will be analysed to predict the likelihood of a collision using 3 machine learning modules and it's severity given a number of variables such as weather conditions, road conditions, time of day, period of year, etc.

These predictions will help the Department of Transport of Seattle to implement preventive measures as required.

A quick look at the map of Seattle, shows that the majority of accidents occur downtown.

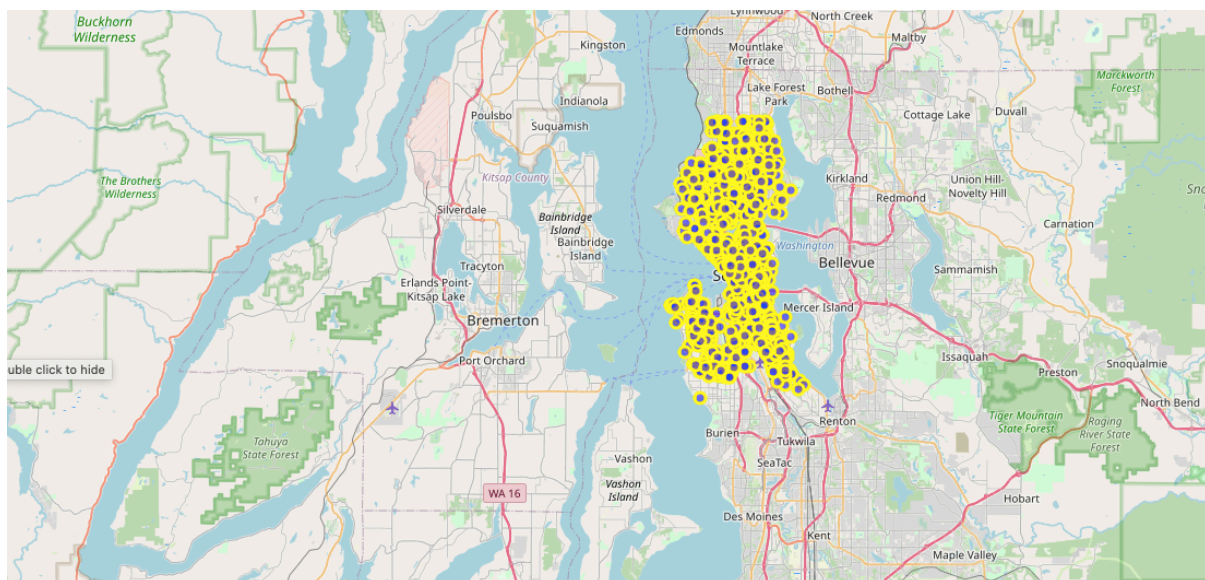


Figure 1:

Map showing the 1st 1000 collisions

## Methodology:

### 1. Data evaluation

The original data set had 37 columns and 194.673 rows of entries. The columns are (with the count of null values):

	<b>NaN Count</b>	<b>EXCEPTRSNCODE</b>	109862	<b>SDOT_COLCODE</b>	0		
<b>SEVERITYCODE</b>	0	<b>EXCEPTRSNDESC</b>	189035	<b>SDOT_COLDESC</b>	0		
<b>X</b>	5334	<b>SEVERITYCODE.1</b>	0	<b>INATTENTIONIND</b>	164868		
<b>Y</b>	5334	<b>SEVERITYDESC</b>	0	<b>UNDERINFL</b>	4884		
<b>OBJECTID</b>	0	<b>COLLISIONTYPE</b>	4904	<b>WEATHER</b>	5081		
<b>INCKEY</b>	0	<b>PERSONCOUNT</b>	0	<b>ROADCOND</b>	5012		
<b>COLDKEY</b>	0	<b>PEDCOUNT</b>	0	<b>LIGHTCOND</b>	5170		
<b>REPORTNO</b>	0	<b>PEDCYLCOUNT</b>	0	<b>PEDROWNOTGRNT</b>	190006		
<b>STATUS</b>	0	<b>VEHCOUNT</b>	0	<b>SDOTCOLNUM</b>	79737		
<b>ADDRTYPE</b>	1926	<b>INCDATE</b>	0	<b>SPEEDING</b>	185340	<b>SEGLANEKEY</b>	0
<b>INTKEY</b>	129603	<b>INCDTTM</b>	0	<b>ST_COLCODE</b>	18	<b>CROSSWALKKEY</b>	0
<b>LOCATION</b>	2677	<b>JUNCTIONTYPE</b>	6329	<b>ST_COLDESC</b>	4904	<b>HITPARKEDCAR</b>	0

We can drop those columns that have more than 3% of null values as well as those that are not relevant to our study:

- INTKEY
- EXCEPTRSNCODE
- EXCEPTRSNDESC
- INATTENTIONIND
- PEDROWNOTGRNT
- SDOTCOLNUM
- SPEEDING
- INCDATE

## 2. Data Exploration

As we want to clean the data, we can start by looking at the “critical” variables:

- weather
- light conditions
- road conditions
- location of accident

# Number of collisions

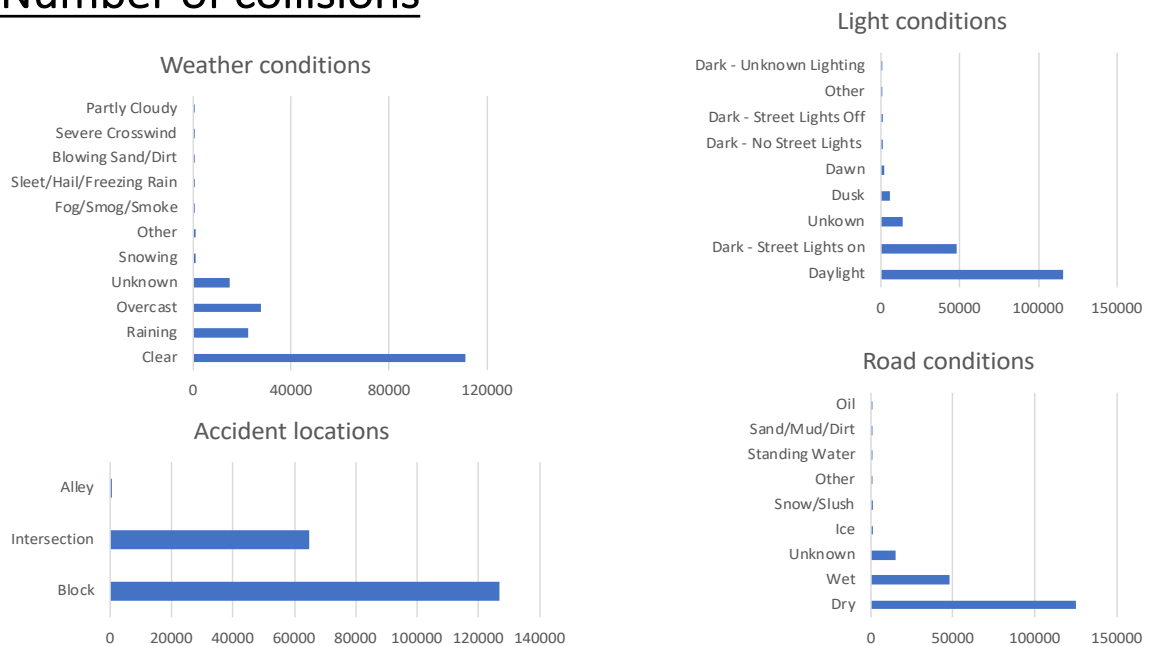


Figure 2: Accidents analysed with different variables

The majority of accidents did not involve injuries, but there were some that did.

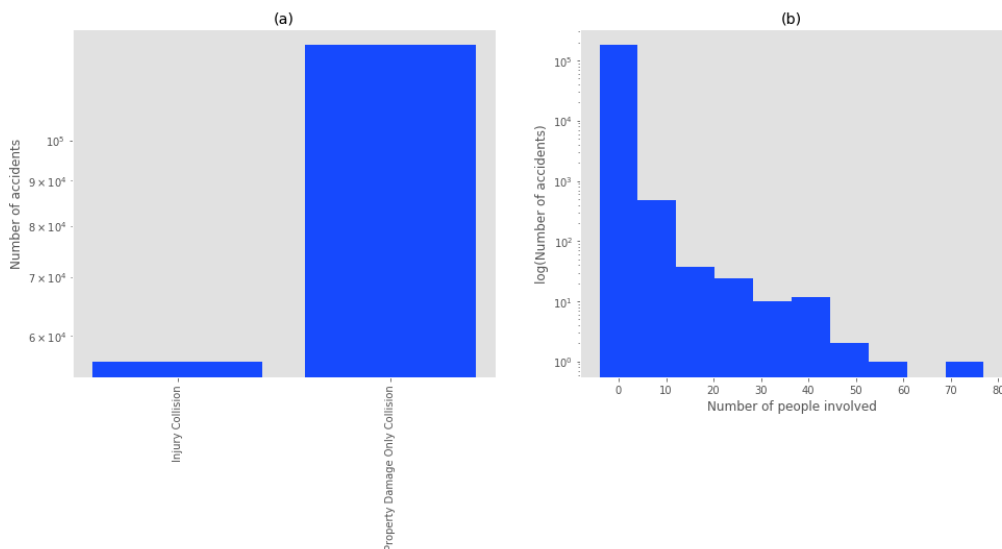


Figure 3: Type of collision & number of people involved

We can see that Fridays were the day with the most accidents, and October was the month.

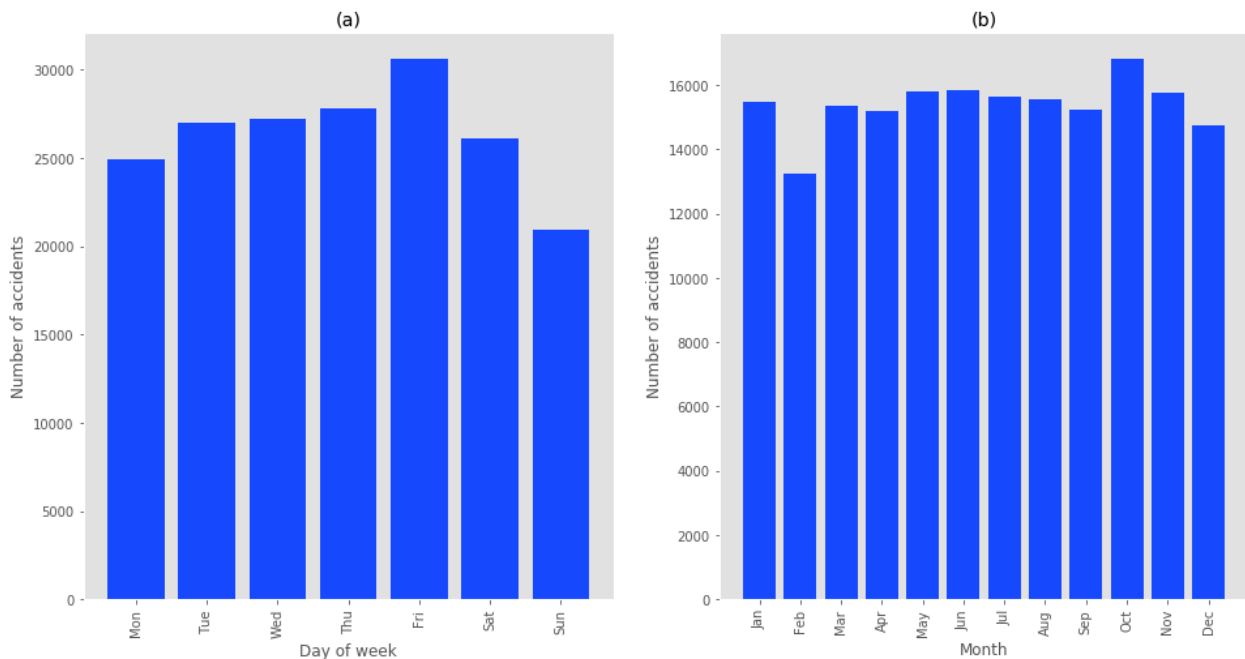


Figure 4: When the collisions took place

### 3. Data Pre-processing

The dataset providing the Seattle collisions information from 2004 – present day are not suitable for quantitative data analysis in its original form, as:

- Data is incomplete. Around 15% of recorded accidents are missing one or more key features, including in some cases the target variable (accident severity code). These needed to be removed from the dataset.
- Dataset includes redundant columns. These also were removed.
- Data are imbalanced and non-standardised.

### 4. Model Development and Evaluation

To develop a model for predicting accident severity, the re-sampled and cleaned dataset was split into testing and training sets (containing 30% and 70% of the samples). By using the scikit learn “train\_test\_split” method, 3 models were trained and evaluated.

- Decision Tree model
- K-Nearest Neighbors (kNN) model
- Logistic Regression model

It was found that the Decision Tree produced the highest accuracy.

	Algorithm	Jaccard	F1-score	Precision
0	KNN	0.73	0.68	0.71
1	Logistic Regression	0.7	0.58	0.67
2	Decision Tree	0.75	0.69	0.77

### Conclusions:

This work highlights that machine learning techniques can be applied to historical data to be able to make reliable predictions about the outcome of road traffic accidents.

The model can be extended to include new features so that city planners can gain insight into the conditions associated with high accident severity and use this to improve road or traffic flow design.