# Predict the Severity of a Traffic Accident

## Monica
LinkedIn | GitHub

9 October 2020

## Contents

# 1   Introduction

## 1.1 Background

Every year car accidents cause hundreds of thousands of deaths worldwide. According to a research conducted by the World Health Organization (WHO) there were 1.35 million road traffic deaths globally in 2016, with millions more sustaining serious injuries and living with long-term adverse health consequences. Globally, road traffic crashes are a leading cause of death among young people, and the main cause of death among those aged 15–29 years. Road traffic injuries are currently estimated to be the eighth leading cause of death across all age groups globally and are predicted to become the seventh leading cause of death by 2030[1]. Leveraging the tools and all the information nowadays available, an extensive analysis to predict traffic accidents and its severity would make a difference to the death toll. Analysing a significant range of factors, including weather conditions, locality, type of road and lighting among others, an accurate prediction of the severity of the accidents can be performed. Thus, trends that commonly lead to severe traffic incidents can help identifying the highly severe accidents. This kind of information could be used by emergency services, to send the exact required staff and equipment to the place of the accident, leaving more resources available for accidents occurring simultaneously. Moreover, this severe accident situation can be warned to nearby hospitals which can have all the equipment ready for a severe intervention in advance. Consequently, road safety should be a prior interest for governments, local authorities and private companies investing in technologies that can help reduce accidents and improve overall driver safety.

## 1.2 Problem

Data that might contribute to determining the likeliness of a potential accident occurring might include information on previous accidents such as road conditions, weather conditions, exact time and place of the accident, type of vehicles involved in the accident, information on the users involved in the accident and off course the severity of the accident. This project aims to forecast the severity of accidents with previous information that could be given by a witness informing the emergency services.

## 1.3 Interest

Governments should be highly interested in accurate predictions of the severity of an accident, to reduce the time of arrival and to make a more efficient use of the resources, and thus save a significant amount of people each year. Others interested could be private companies investing in technologies aiming to improve road safeness.

# 2  Data

## 2.1 Data Source

The dataset is a part of an example dataset in IBM Data Science Professional Certificate. It can be downloaded from my GitHub repository https://github.com/monica110394/Coursera_Capstone. The name of the file is Data-Collisions.csv.

## 2.2 Variable Selection

After reading the data in pandas, analysis is performed to find the number of samples present in the dataset, also the number of attributes. The dataset has 194,673 records and 38 attributes including the target variable. The columns present in the dataset are:

```
'SEVERITYCODE', 'X', 'Y', 'OBJECTID', 'INCKEY', 'COLDETKEY', 'REPORTNO,
'STATUS', 'ADDRTYPE', 'INTKEY', 'LOCATION', 'EXCEPTRSNCODE',
'EXCEPTRSNDESC', 'SEVERITYCODE.1', 'SEVERITYDESC', 'COLLISIONTYPE',
'PERSONCOUNT, 'PEDCOUNT', 'PEDCYLCOUNT', 'VEHCOUNT', 'INCDATE',
'INCDTTM', 'JUNCTIONTYPE', 'SDOT_COLCODE', 'SDOT_COLDESC',
'INATTENTIONIND', 'UNDERINFL', 'WEATHER', 'ROADCOND', 'LIGHTCOND',
'PEDROWNOTGRNT', 'SDOTCOLNUM', 'SPEEDING', 'ST_COLCODE', 'ST_COLDESC',
'SEGLANEKEY', 'CROSSWALKKEY', 'HITPARKEDCAR'
```

Going through the description of every attribute, which can also be found at my GitHub repository with the file name Metadata.pdf, intuitionally it can be seen that many of the attributes will not contribute to the prediction. The attributes that have been ignored for this prediction are:

```
'X', 'Y', 'OBJECTID', 'INCKEY', 'COLDETKEY', 'REPORTNO, 'STATUS', 'INTKEY', 'LOCAT
ION', 'EXCEPTRSNCODE', 'EXCEPTRSNDESC', 'SEVERITYCODE.1',
'SEVERITYDESC', 'INCDATE', 'INCDTTM', 'SDOT_COLDESC', 'PEDROWNOTGRNT',
'SDOTCOLNUM', 'SPEEDING', 'ST_COLDESC', 'SEGLANEKEY', 'CROSSWALKKEY'
```

The attributes that are kept for the prediction are (17 variable including target variable):

```
'SEVERITYCODE', 'ADDRTYPE', 'COLLISIONTYPE', 'PERSONCOUNT', 'PEDCOUNT',
'PEDCYLCOUNT', 'VEHCOUNT', 'JUNCTIONTYPE', 'SDOT_COLCODE',
'INATTENTIONIND', 'UNDERINFL', 'WEATHER', 'ROADCOND', 'LIGHTCOND',
'SPEEDING', 'ST_COLCODE', 'HITPARKEDCAR'
```
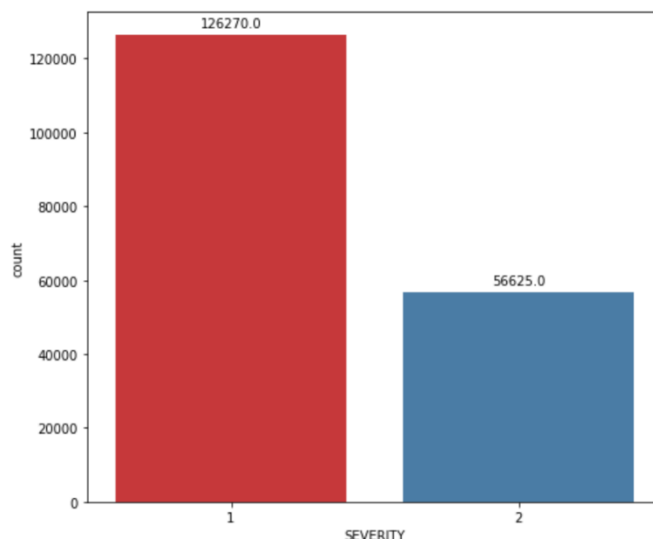
Target variable: `'SEVERITYCODE'`

## 2.3 Feature Engineering

After knowing about the dataset, feature engineering is performed such as:

- Null Values are detected
  - **SEVERITYCODE          0**
  - ADDRTYPE          1926
  - COLLISIONTYPE     4904
  - PERSONCOUNT          0
  - PEDCOUNT             0
  - PEDCYLCOUNT          0
  - VEHCOUNT             0
  - JUNCTIONTYPE      6329
  - SDOT_COLCODE         0
  - INATTENTIONIND  164868
  - UNDERINFL         4884
  - WEATHER           5081

```
o   ROADCOND              5012
o   LIGHTCOND             5170
o   SPEEDING            185340
o   ST_COLCODE             18
o   HITPARKEDCAR            0
```
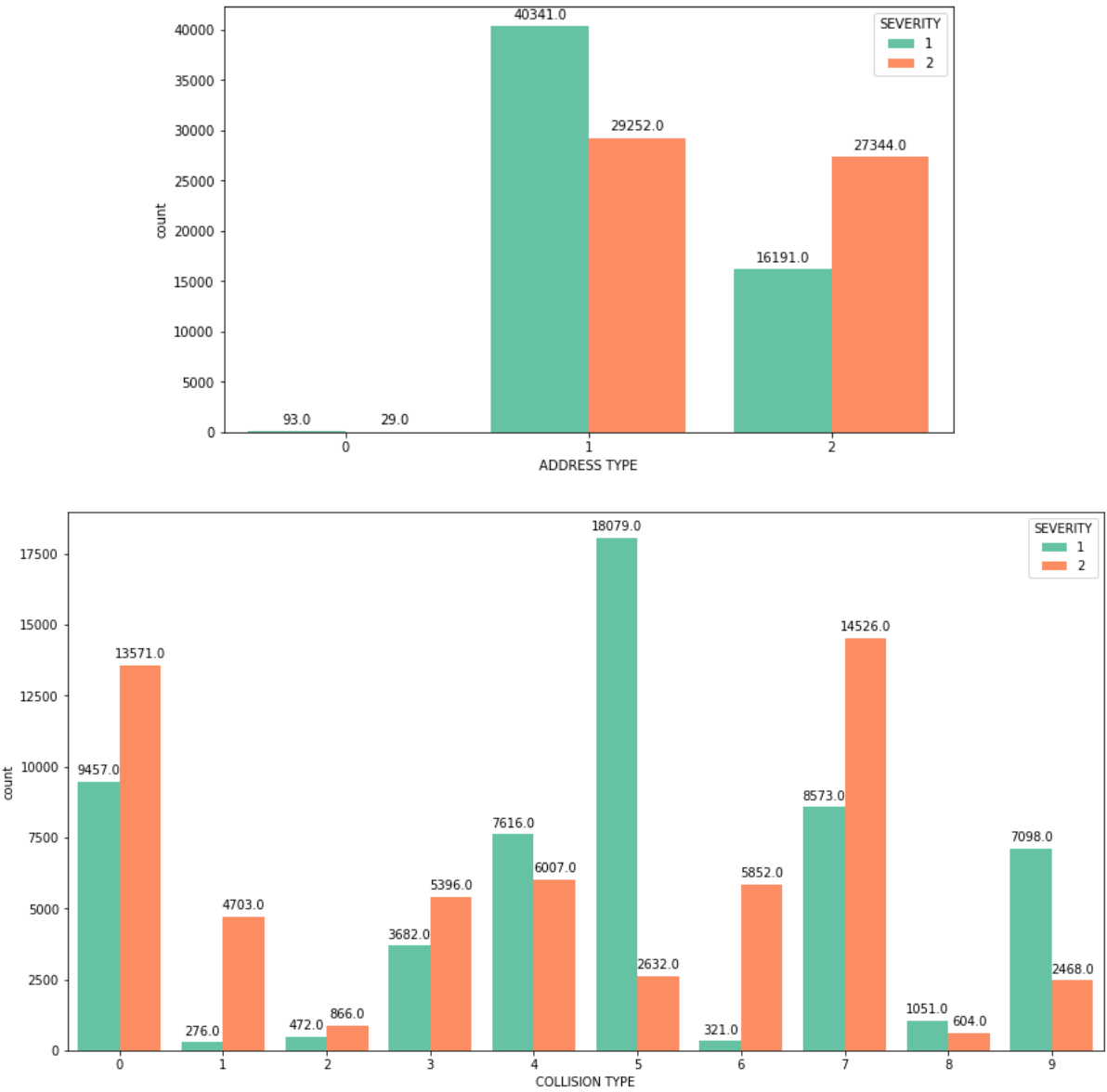
- Attributes that have most of the values, more than 70% values, as null, are dropped
  ```
  o   INATTENTIONIND   164868
  o   SPEEDING         185340
  ```

- Current number of rows, columns and attributes
  ```
  o   Number of rows: 194673
  o   Number of Columns 15
  o   'SEVERITYCODE', 'ADDRTYPE', 'COLLISIONTYPE', 'PERSONCOUNT',
      'PEDCOUNT', 'PEDCYLCOUNT', 'VEHCOUNT', 'JUNCTIONTYPE',
      'SDOT_COLCODE', 'UNDERINFL', 'WEATHER', 'ROADCOND',
      'LIGHTCOND', 'ST_COLCODE', 'HITPARKEDCAR'
  ```

- Records or rows that have missing values are dropped as well.
  ```
  o   Number of rows: 182895
  o   Number of Columns 15
  ```

- Label Encoded all the categorical columns converting the 'string' labels into numbers and then changed their data types from 'integer' to 'category'.
  ```
  o   'ADDRTYPE'
  o   'COLLISIONTYPE'
  o   'JUNCTIONTYPE'
  o   'SDOT_COLCODE'
  o   'UNDERINFL'
  o   'WEATHER'
  o   'ROADCOND'
  o   'LIGHTCOND'
  o   'ST_COLCODE'
  o   'HITPARKEDCAR'
  ```

- Numeric columns
  ```
  o   'PERSONCOUNT'
  o   'PEDCOUNT'
  o   'PEDCYLCOUNT'
  o   'VEHCOUNT'
  ```

- Renaming columns to reasonable names

- On plotting a count plot of the target variable, **'SEVERITYCODE'**, it is found that the data set is highly unbalanced, which can cause the prediction to be skewed. Hence, the dataset is balanced by down sampling the category which as greater number of samples, in this case, when **'SEVERITYCODE'** is 1.
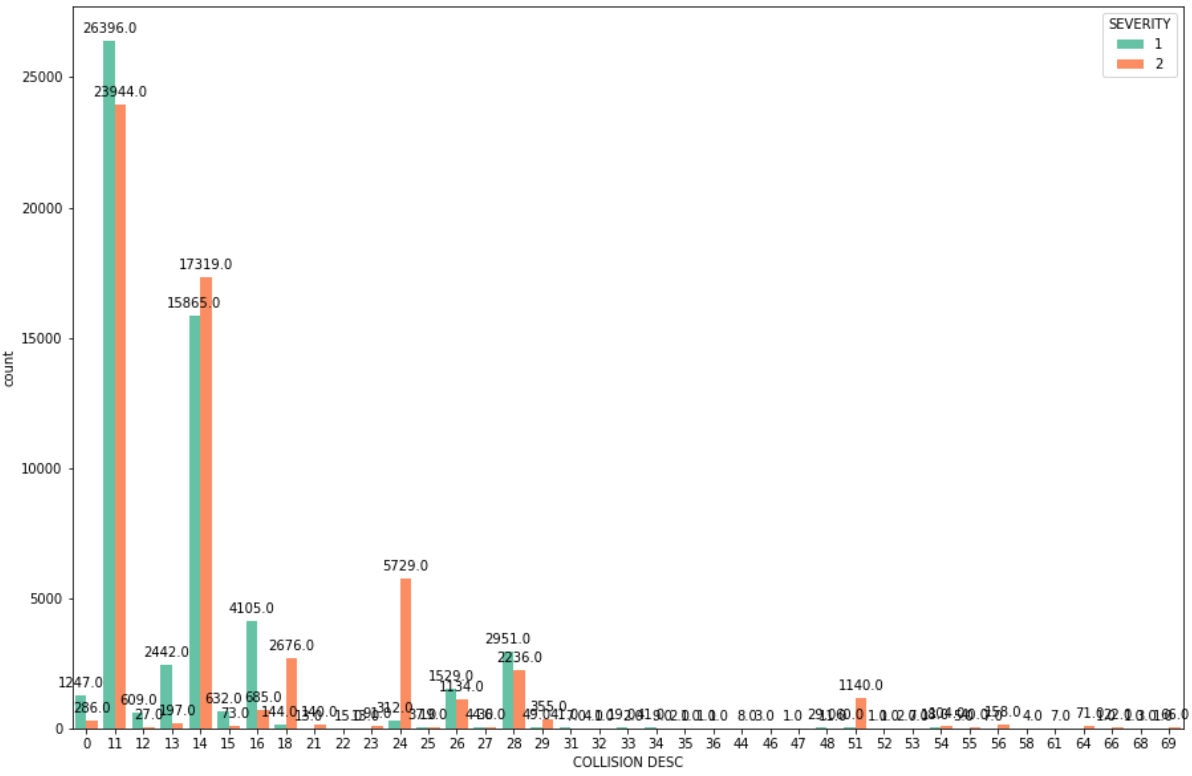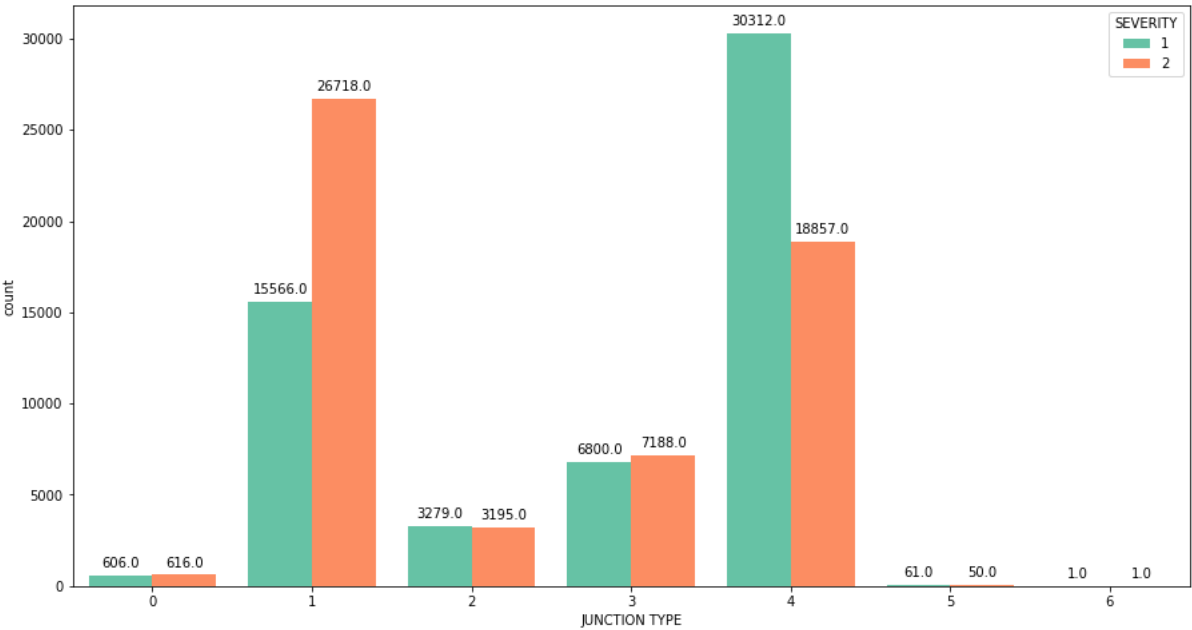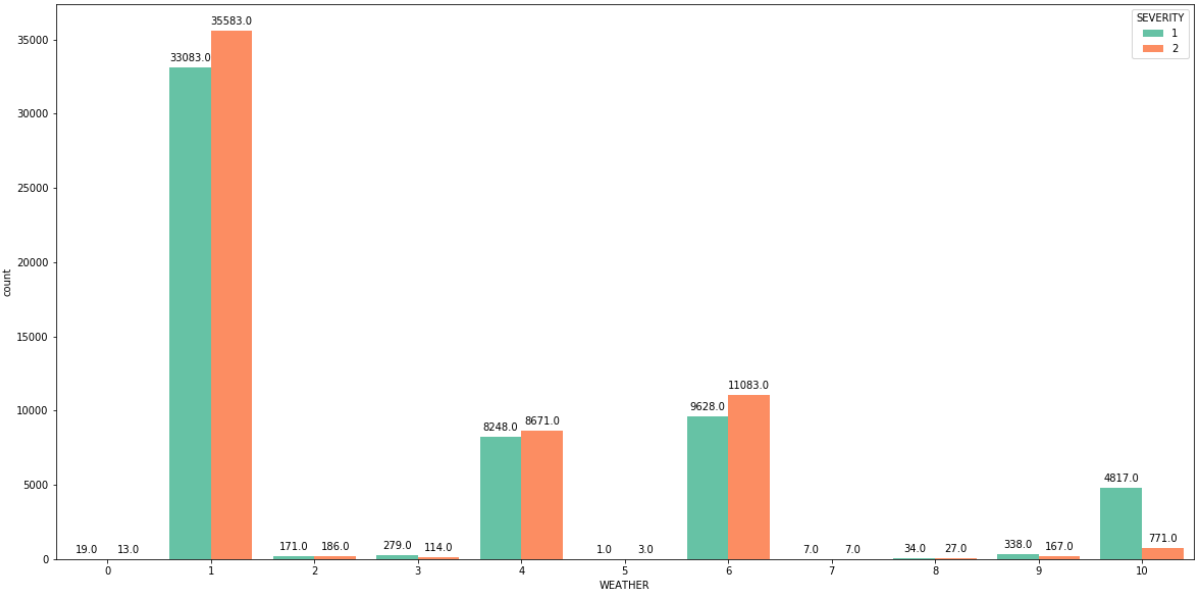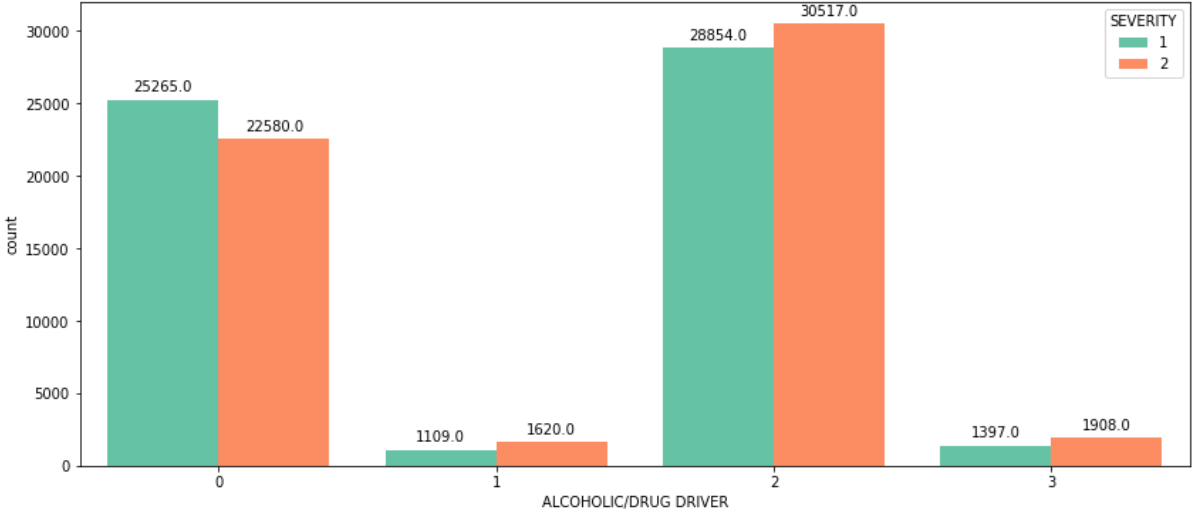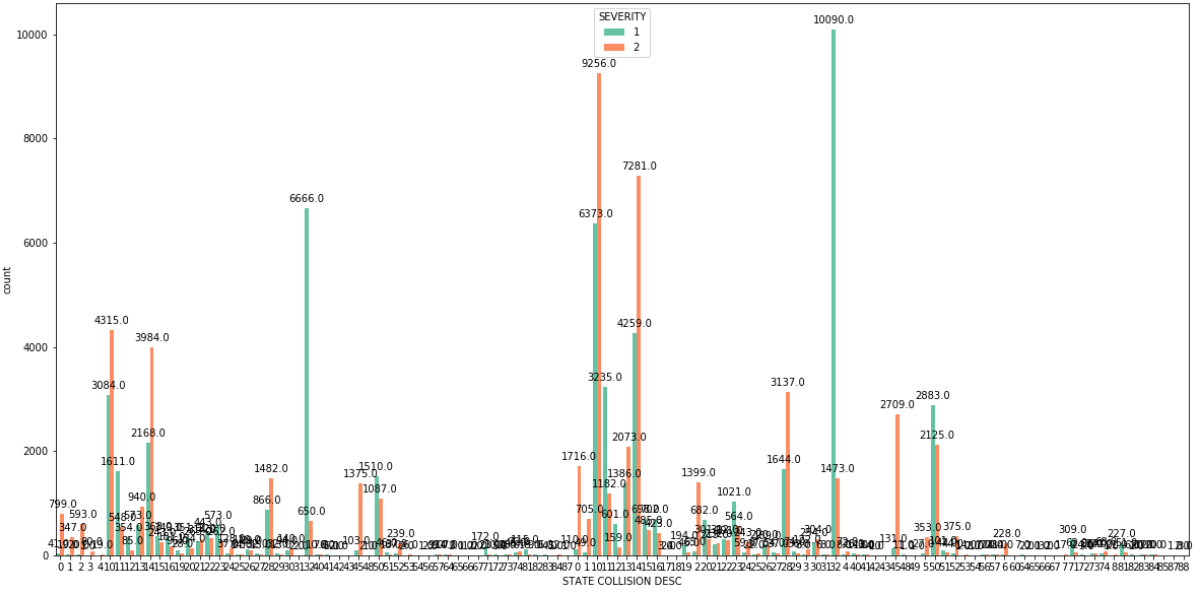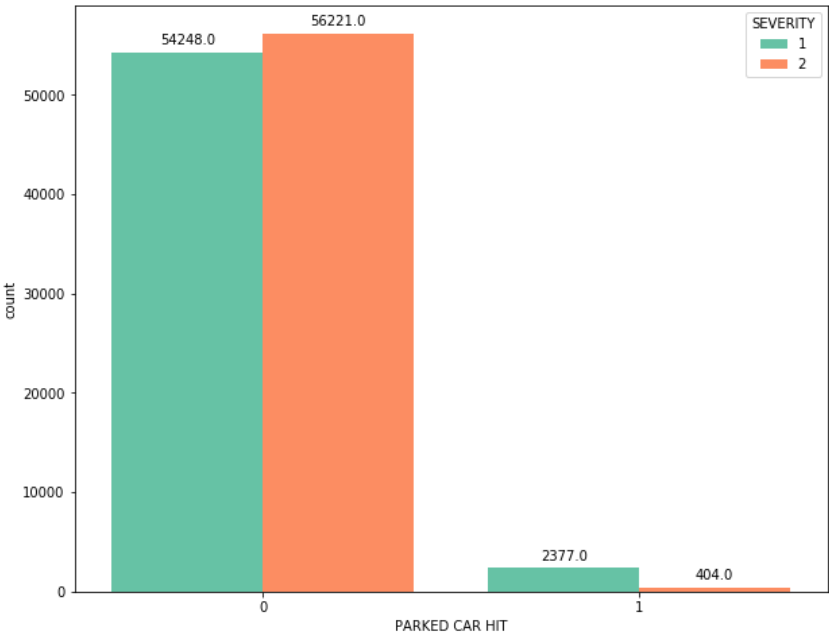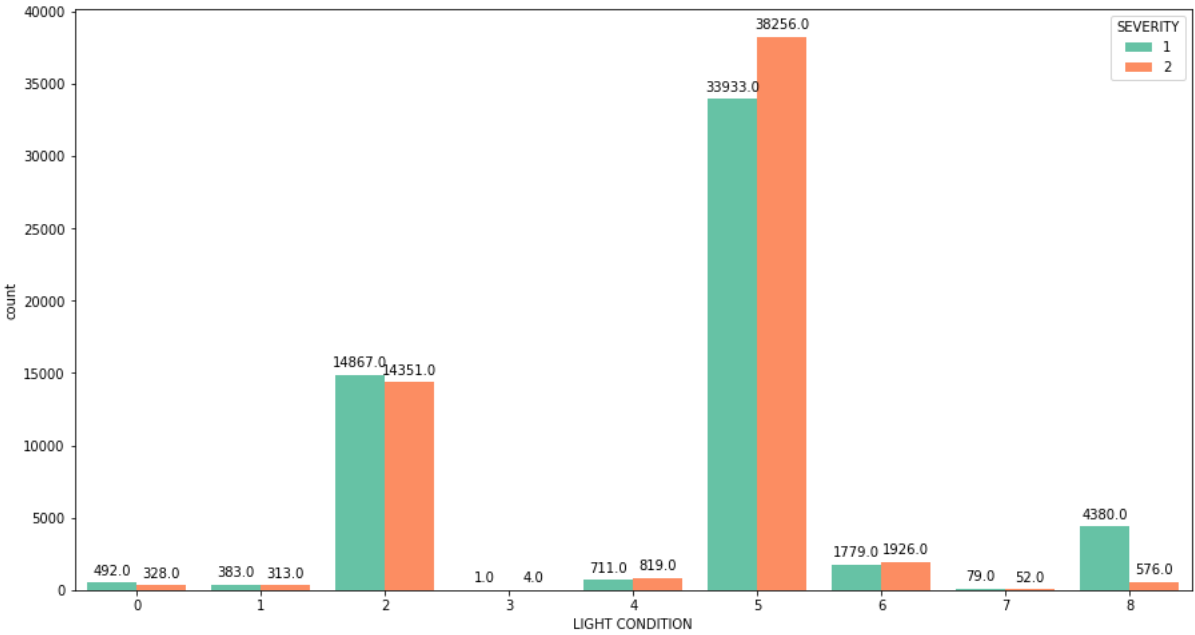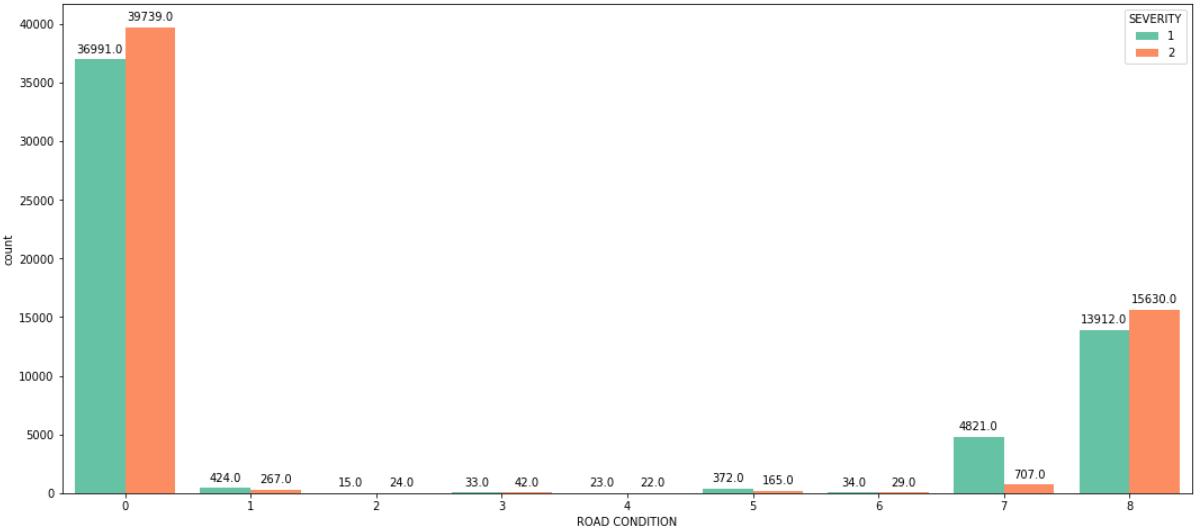
# 3   Exploratory Data Analysis

## 3.1 Categorical Variables

The best way to analyse categorical variables is using a count plot against the target variable `'SEVERITYCODE'`.
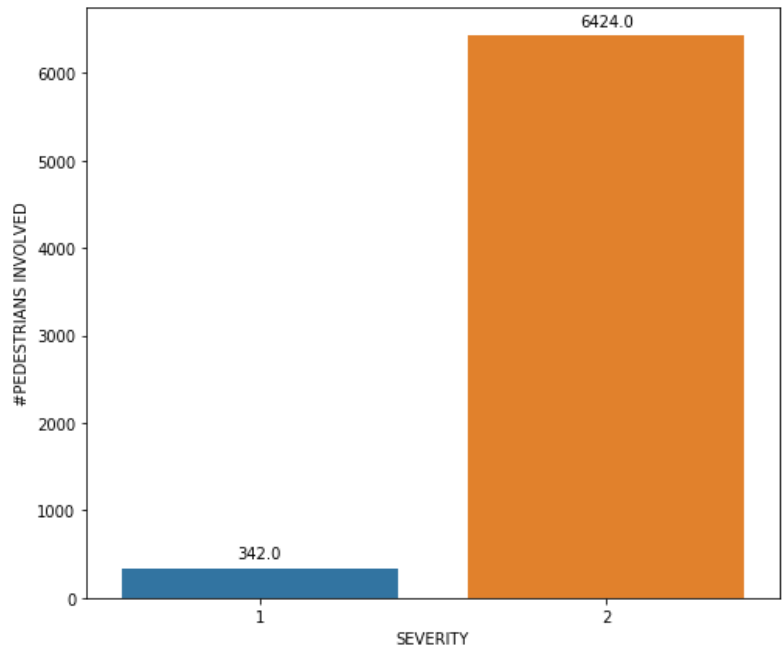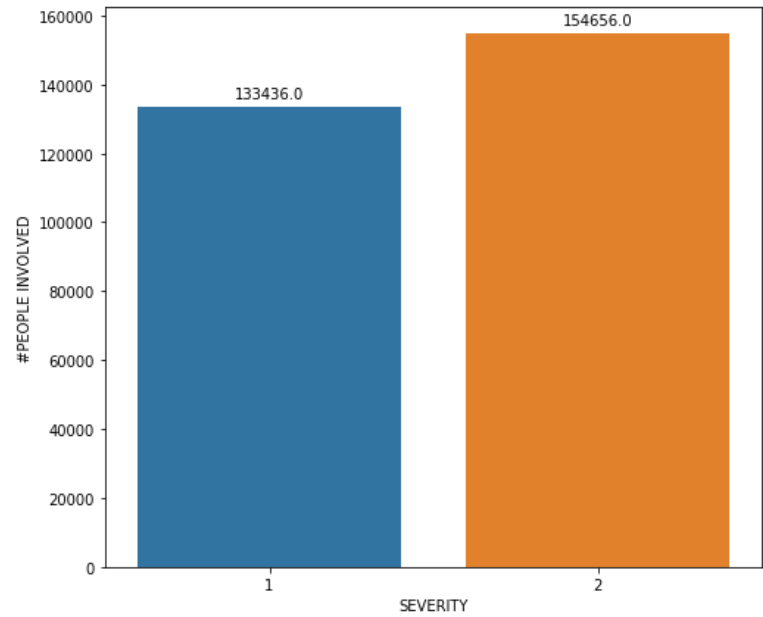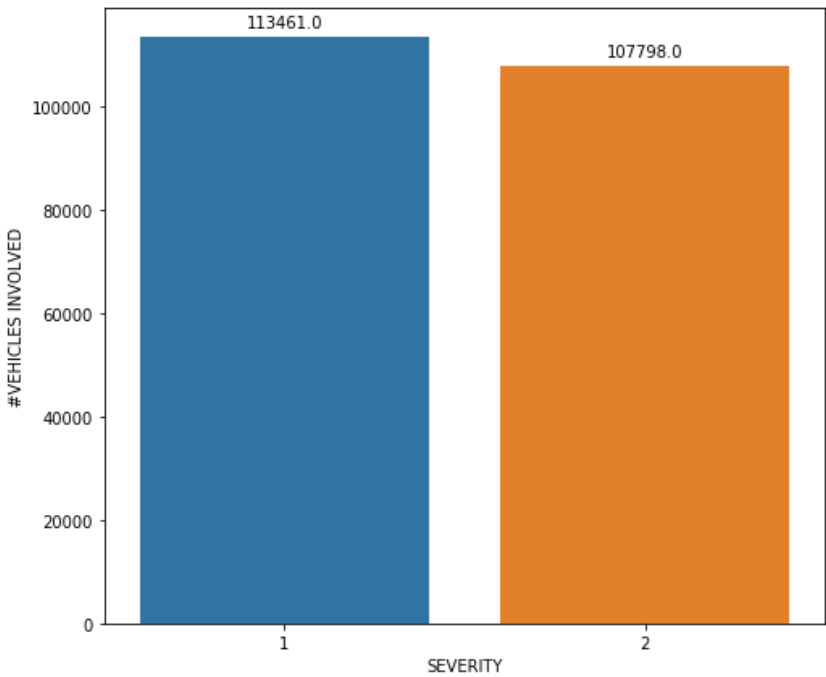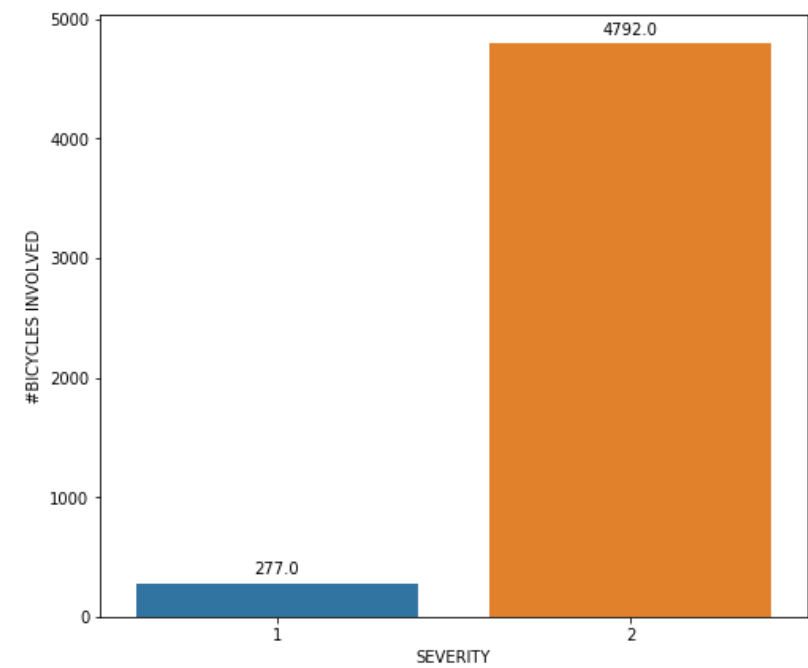
## 3.2 Numerical Variables

For numerical variables, they are grouped and aggregated as sum over both the categories of the target variable `'SEVERITYCODE'.` After which, bar plots is used to analyse the variables against the target variables.

| | SEVERITY | #PEOPLE INVOLVED | #PEDESTRIANS INVOLVED | #BICYCLES INVOLVED | #VEHICLES INVOLVED |
|---|---|---|---|---|---|
| **0** | 1 | 133436 | 342 | 277 | 113461 |
| **1** | 2 | 154656 | 6424 | 4792 | 107798 |

# 4   Predictive Modeling

## 4.1 Data Preparation

Firstly, the features and the target variables are saved into two different data frames. Next, the features' dataset is normalized using 'StandardScaler()', to bring the dataset to zero mean and unit variance.

```
from sklearn import preprocessing
X = preprocessing.StandardScaler().fit(features).transform(features)
```

After the normalization process the whole dataset it divided into train and test datasets with a test size of 20% of the whole dataset.

```
Train set: (90600, 14) (90600,)
Test set: (22650, 14) (22650,)
```
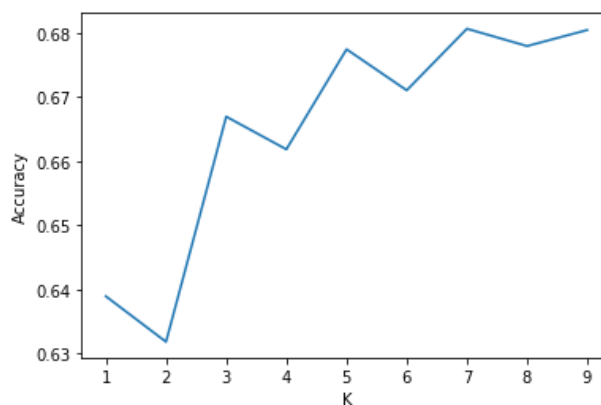
## 4.2 Data Modeling

To predict the severity of a traffic accident, 3 different classification methods are used. Since most of the features are categorical, appropriate classification methods are used. All the classification methods are implemented using the scikit-learn package that is available as a python library. For every classification, the best hyperparameter value is computed using grid search to maximum accuracy. The algorithms that were used are:

- K Nearest Neighbours
- Decision Tree
- Logistic Regression

### 4.2.1   KNN

The best accuracy score achieved was 0.6806 with k=7.
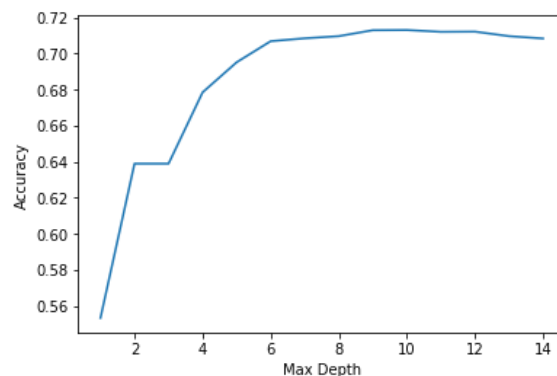


```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
                     metric_params=None, n_jobs=None, n_neighbors=7, p=2,
                     weights='uniform')
```

```
Train set Accuracy:  0.7195
```

```
Test set Accuracy:  0.6806
```

### 4.2.2   Decision Tree

The best accuracy score achieved was 0.6806 with max_depth = 10
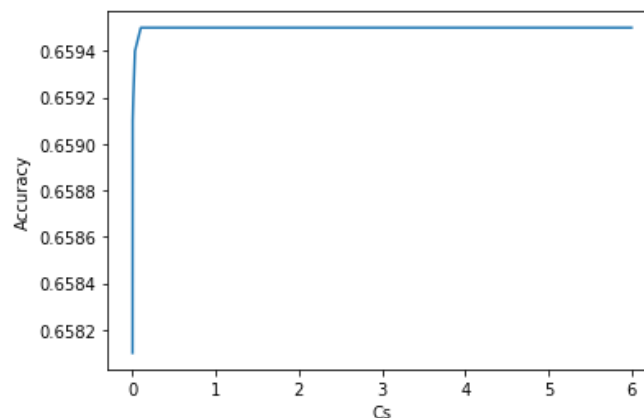


```
DecisionTreeClassifier(ccp_alpha=0.0, class_weight=None, criterion='entropy',
                max_depth=10, max_features=None, max_leaf_nodes=None,
                min_impurity_decrease=0.0, min_impurity_split=None,
                min_samples_leaf=1, min_samples_split=2,
                min_weight_fraction_leaf=0.0, presort='deprecated',
                random_state=None, splitter='best')
```

```
Train set Accuracy:  0.7183
```

```
Test set Accuracy:  0.7128
```

### 4.2.3   Logistic Regression

The best accuracy score achieved was 0.6806 with C=0.1



```
LogisticRegression(C=0.1, class_weight=None, dual=False, fit_intercept=True,
                intercept_scaling=1, l1_ratio=None, max_iter=100,
                multi_class='auto', n_jobs=None, penalty='l2',
                random_state=None, solver='liblinear', tol=0.0001, verbose=0,
                warm_start=False)
```

```
Train set Accuracy:  0.6591
```
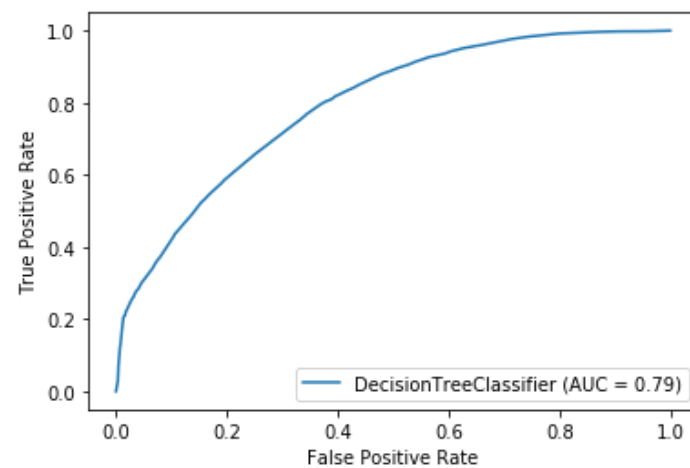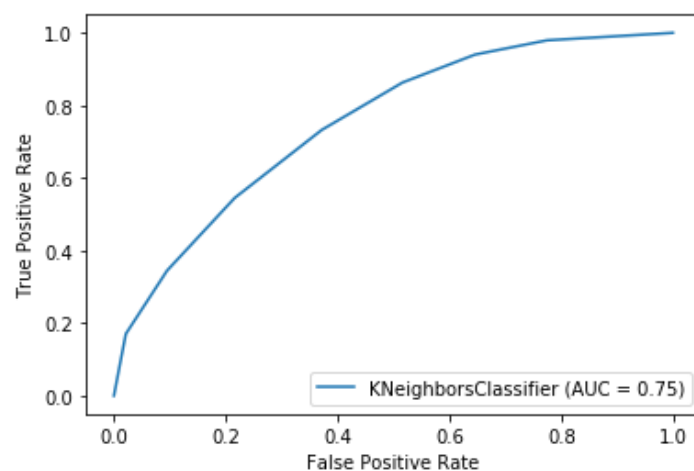
```
Test set Accuracy:  0.6595
```

# 5  Results

All the classification models are evaluated against each other using evaluations metrics provided by the scikit learn library. The metrics are:

- Jaccard Index
- F1 Score
- Precision Score
- Recall Score
- ROC
  In this problem, lower false positive rate is less important than higher true positive rate. In other words, it is more important to properly predict the high-severity accident properly, if there is room for doubt it is better to prevent.
- Log Loss (Logistic Regression only)

| Algorithm | Jaccard | F1 Score | Precision | Recall | AUC (ROC) | Log Loss |
|---|---|---|---|---|---|---|
| KNN | 0.6806 | 0.6797 | 0.6991 | 0.6284 | 0.75 | NA |
| Decision Tree | 0.7128 | 0.7104 | 0.7582 | 0.6207 | 0.79 | NA |
| Logistic Regression | 0.6595 | 0.6585 | 0.6417 | 0.7144 | 0.74 | 0.5899 |

# 6  Conclusion

In the analysis done above and looking at the table, it can be said that Decision Tree is the best classifier out of the three with the best Jaccard Index 0.71, best F1 Score 0.75 and best AUC. Although, there is a possibility to achieve a better accuracy by a Random Forest Classifier, or with a reduced number of features, or adding relevant features to the dataset.

These models can have multiple application in real life. For instance, imagine that emergency services have an application with some default features such as date, time and department/municipality and then with the information given by the witness calling to inform on the accident they could predict the severity of the accident before getting there and so alert nearby hospitals and prepare with the necessary equipment and staff. Also, by identifying the features that favour the most the gravity of an accident, these could be tackled by improving road conditions or increasing the awareness of the population.

# 7  Discussion

For this study, only three classifiers were used, that is, KNN, Decision Tree and Logistic Regression but other classifiers can also be applied to the data such as Naïve Bayes, SVM and Random Forest, which might or might not produce a better accuracy metric.

As per the data, for this problem, i.e. to predict the severity of an accident, I feel a lot more relevant features could be added such as number of accidents per year which could be helpful to determine the frequency of accidents happening and also, the age of the person involved in the accident. This dataset had lots of missing values in two of the most important features of the dataset, SPEEDING (speeding or not) and INATTENTIONID (collision due to inattention), which could have increased the accuracy of the classification models.

The next step on this problem could be to add an accident prediction model able to not just predict the accuracy but also the critical time and spots where potential accidents can occur in advance.

# 8   References

Josep-at-work/Coursera_Capstone. (2020). Retrieved 9 October 2020, from https://github.com/Josep-at-work/Coursera_Capstone/blob/master/Predicting_Traffic_Accident_Severity_.pdf