

Predict the Severity of a Traffic Accident

IBM Data Science Professional Certificate

Monica

[LinkedIn](#) | [GitHub](#)

9 October 2020

Data Source

- Part of an example dataset in IBM Data Science Professional Certificate
- Can be downloaded from my GitHub repository https://github.com/monica110394/Coursera_Capstone

Variable Selection

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

Feature Engineering

- Null Values are detected
- Attributes are removed that have most of the values as null
- Records are dropped that had missing values
 - Number of rows: 182895
 - Number of Columns: 15
 - 'SEVERITYCODE', 'ADDRTYPE', 'COLLISIONTYPE', 'PERSONCOUNT', 'PEDCOUNT', 'PEDCYLCOUNT', 'VEHCOUNT', 'JUNCTIONTYPE', 'SDOT_COLCODE', 'UNDERINFL', 'WEATHER', 'ROADCOND', 'LIGHTCOND', 'ST_COLCODE', 'HITPARKEDCAR'

Label Encoding (categorical variables)

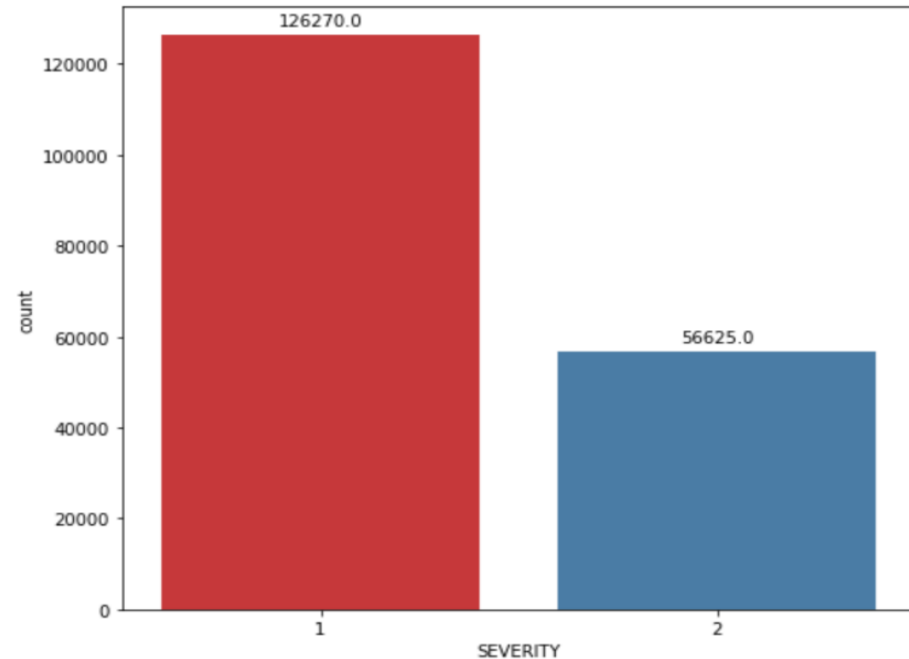
- 'ADDRTYPE'
- 'COLLISIONTYPE'
- 'JUNCTIONTYPE'
- 'SDOT_COLCODE'
- 'UNDERINFL'
- 'WEATHER'
- 'ROADCOND'
- 'LIGHTCOND'
- 'ST_COLCODE'
- 'HITPARKEDCAR'

Numerical variables

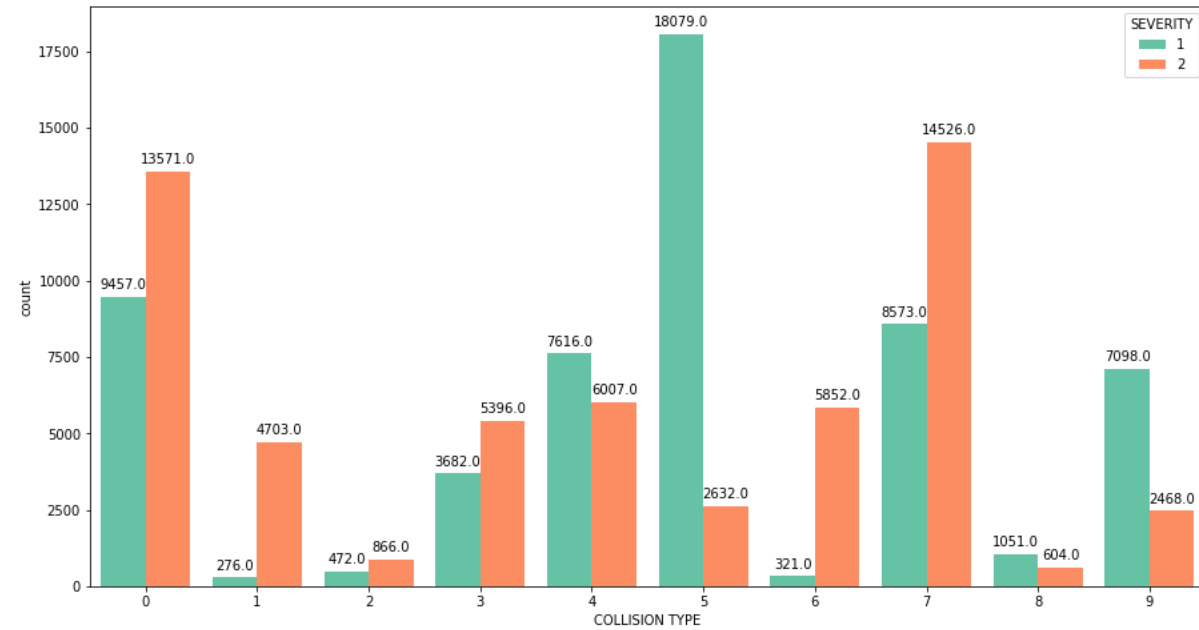
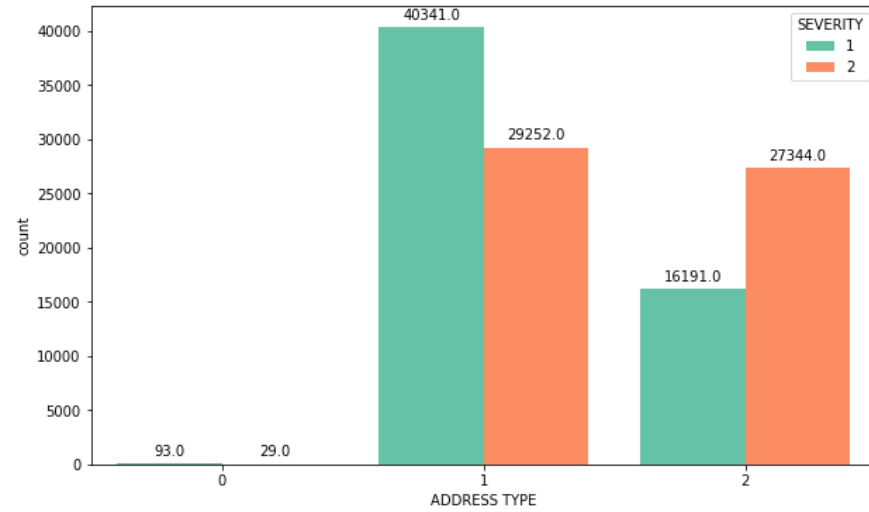
- 'PERSONCOUNT'
- 'PEDCOUNT'
- 'PEDCYLCOUNT'
- 'VEHCOUNT'

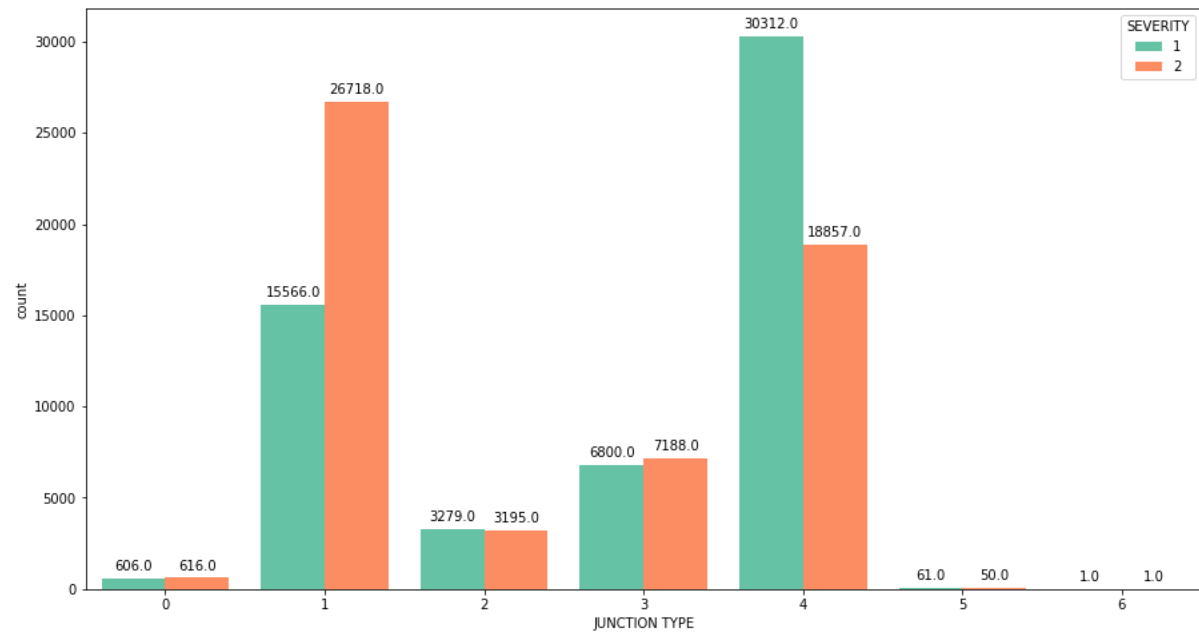
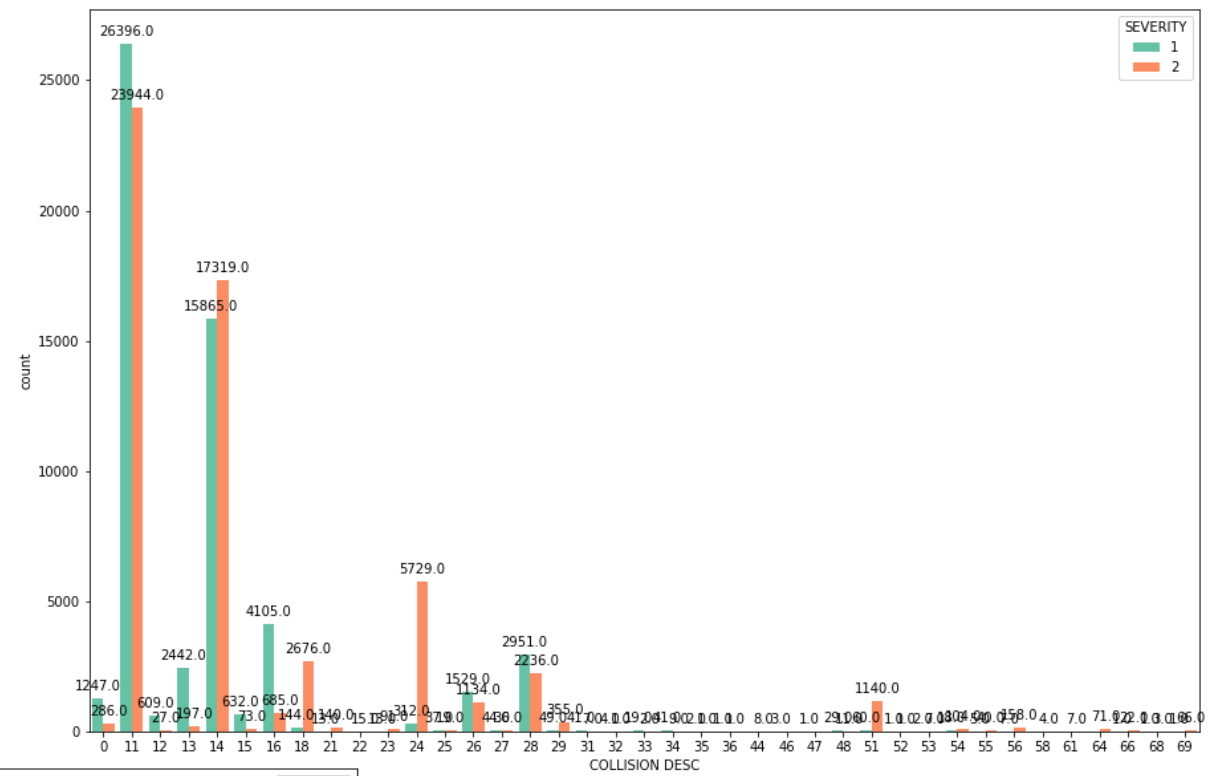
Balancing dataset:

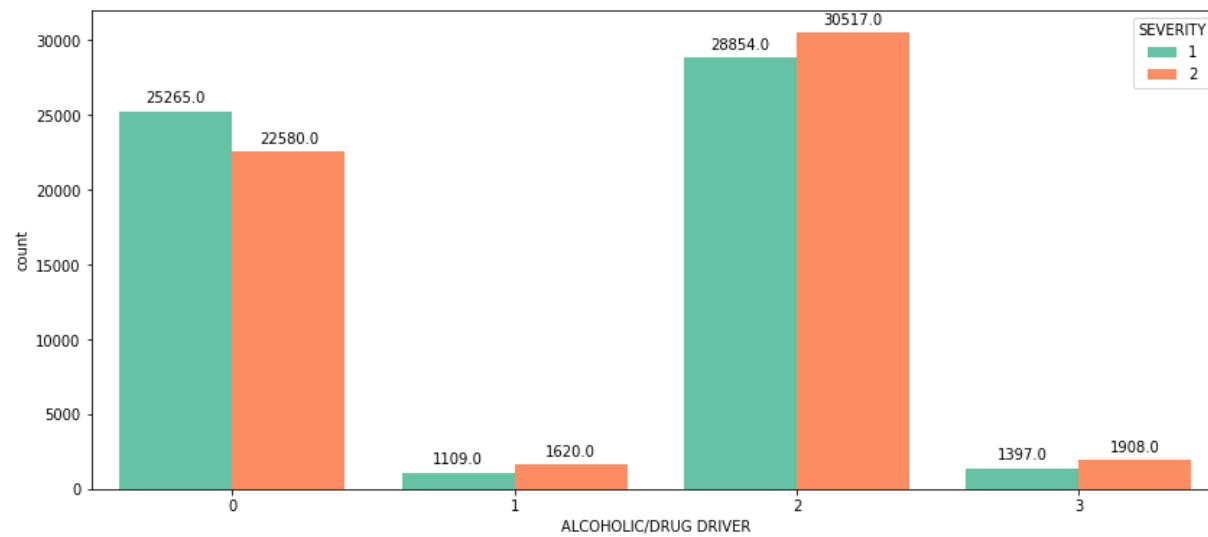
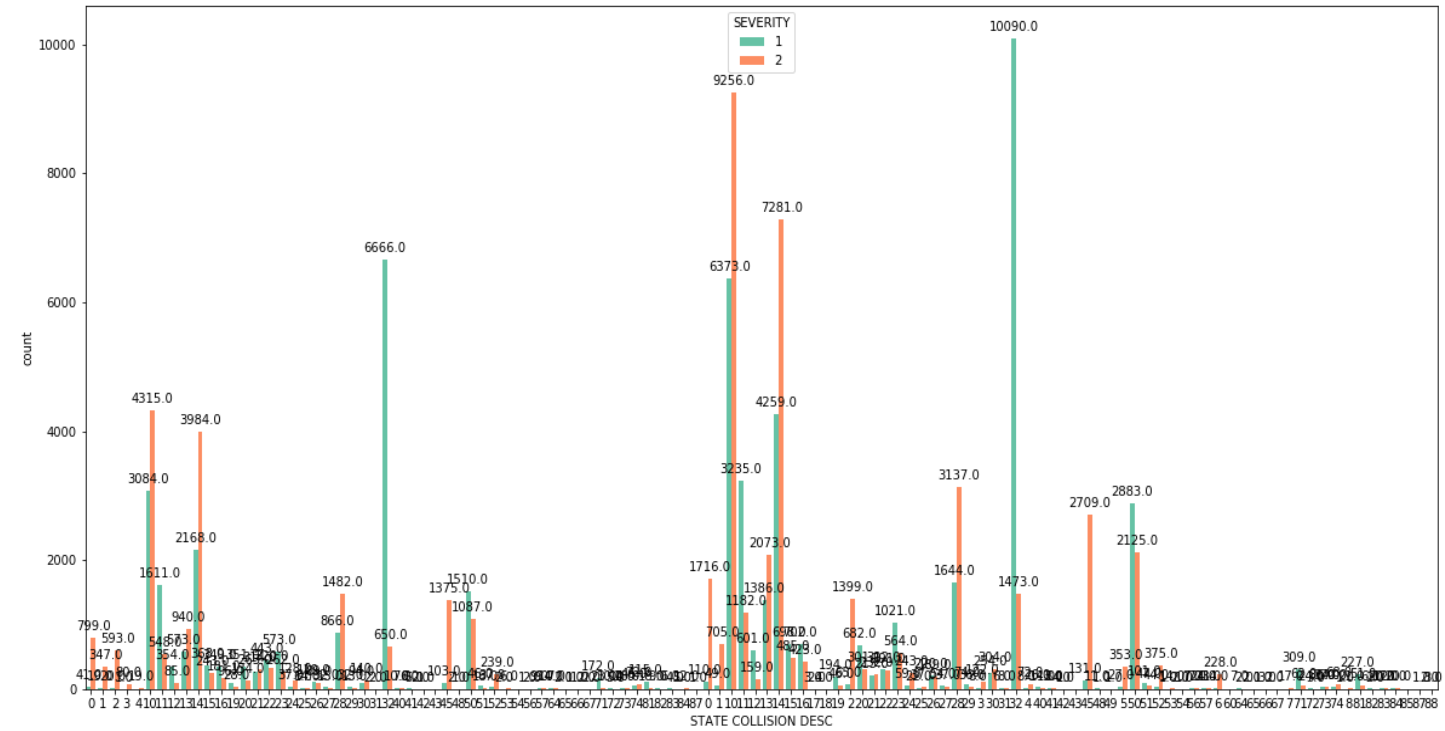
- Unbalanced dataset can cause the prediction to be skewed
- Dataset is balanced by down sampling the category that has greater number of samples
- In this case it is 'SEVERITYCODE' = 1

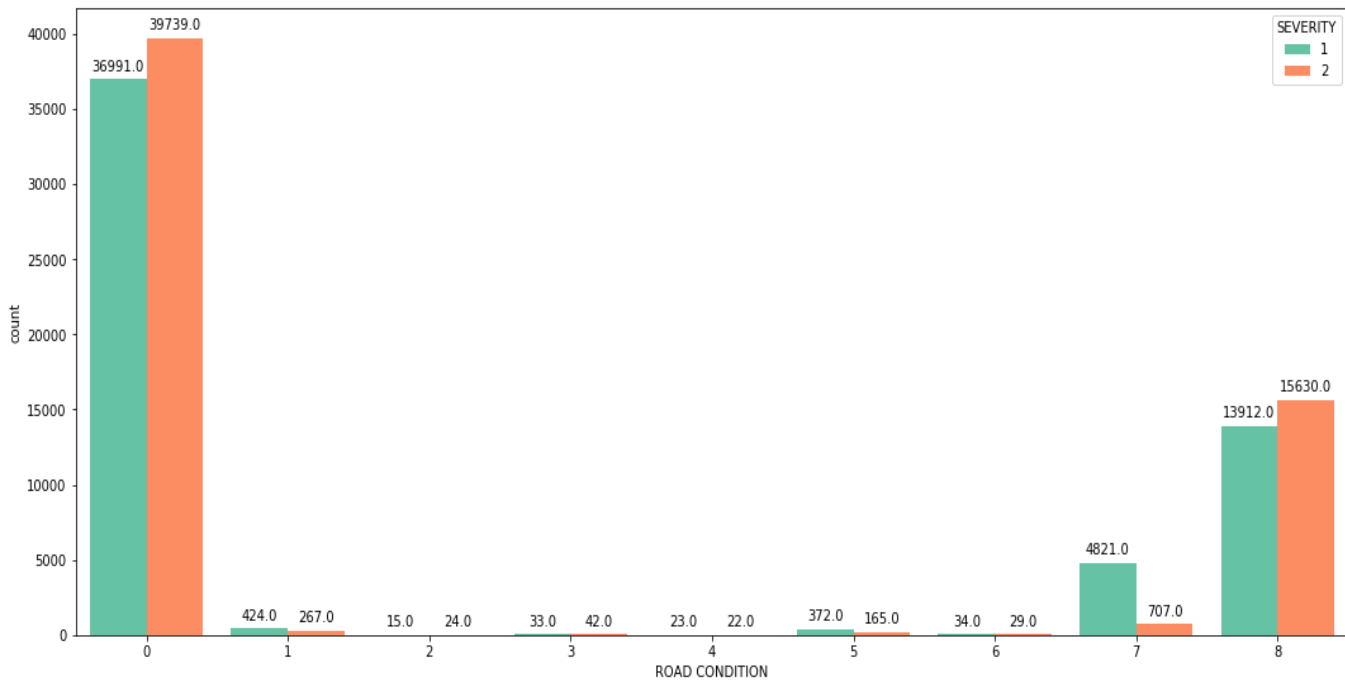
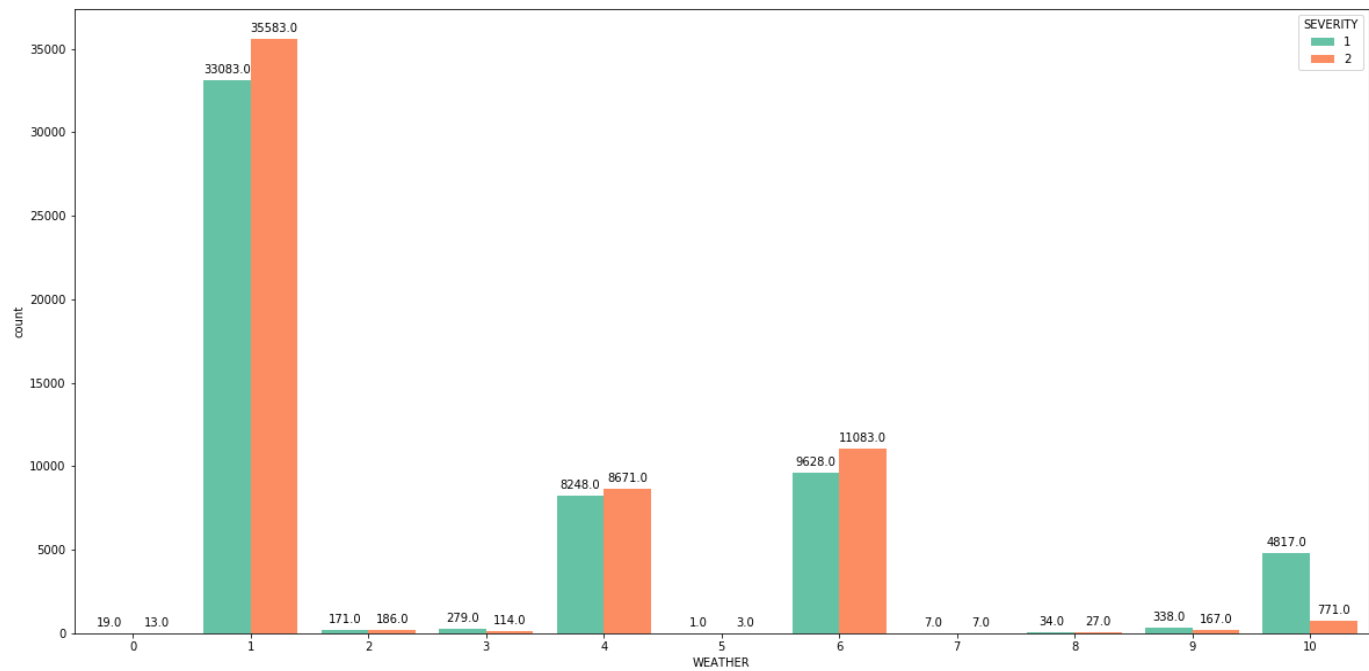


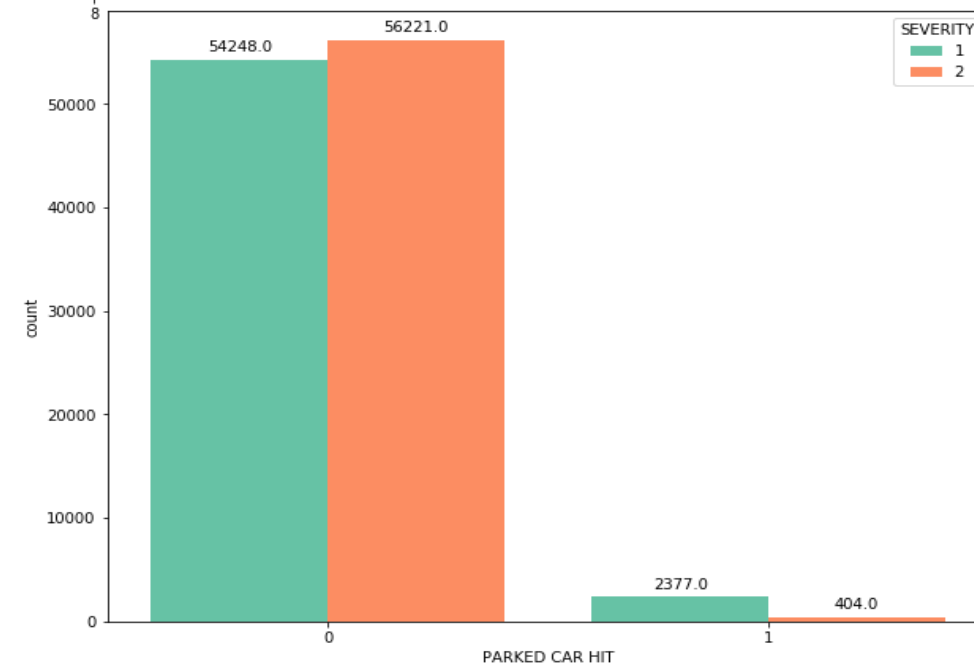
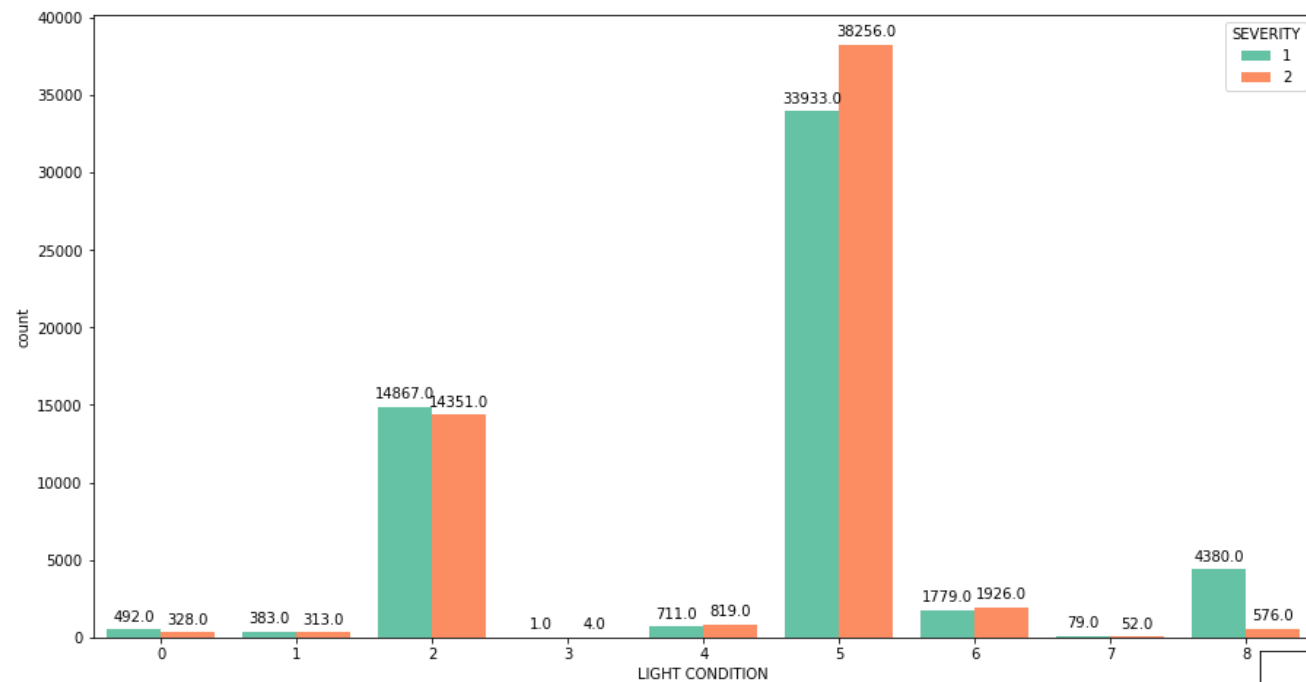
Exploratory Data Analysis (categorical variables)



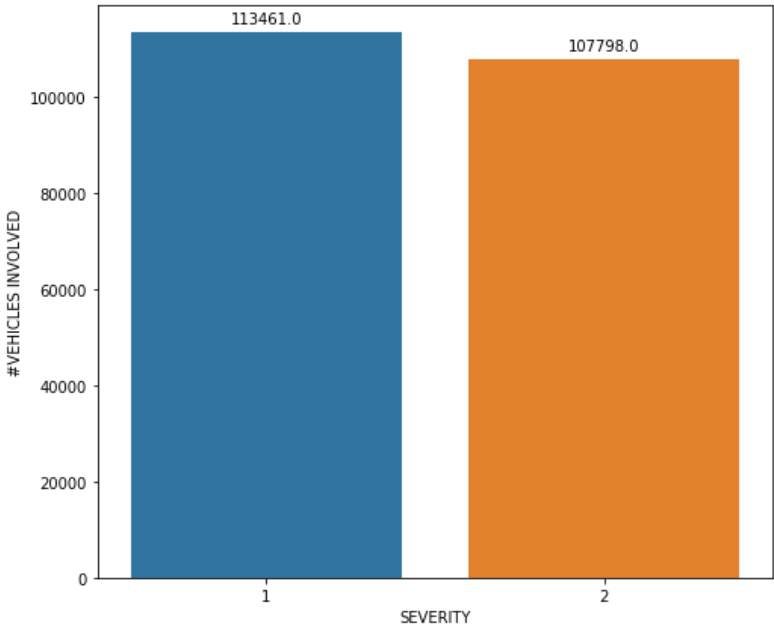
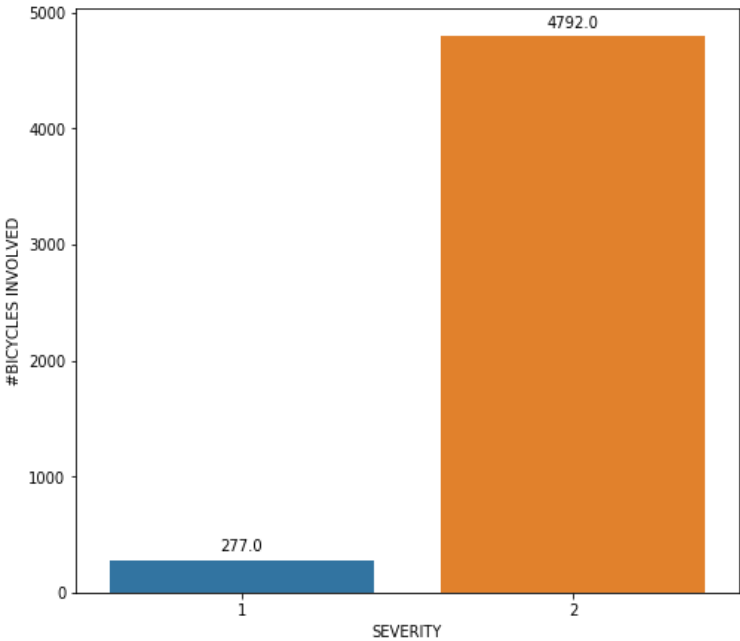
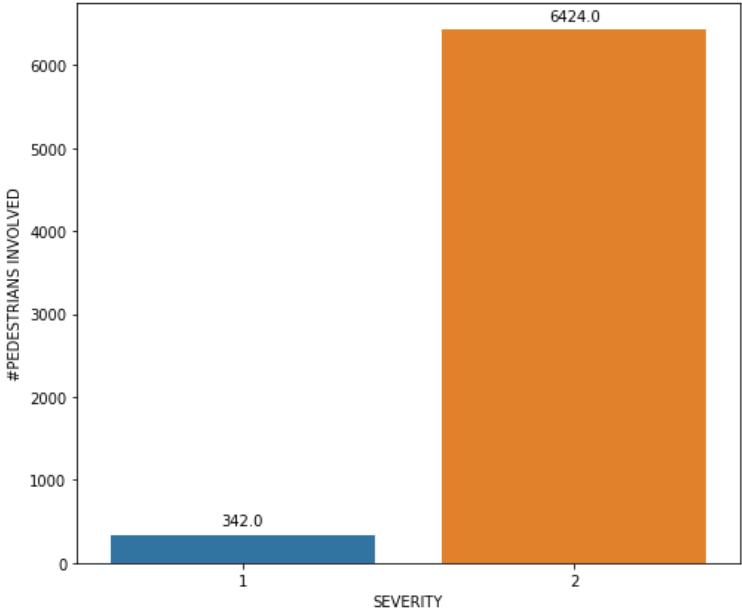
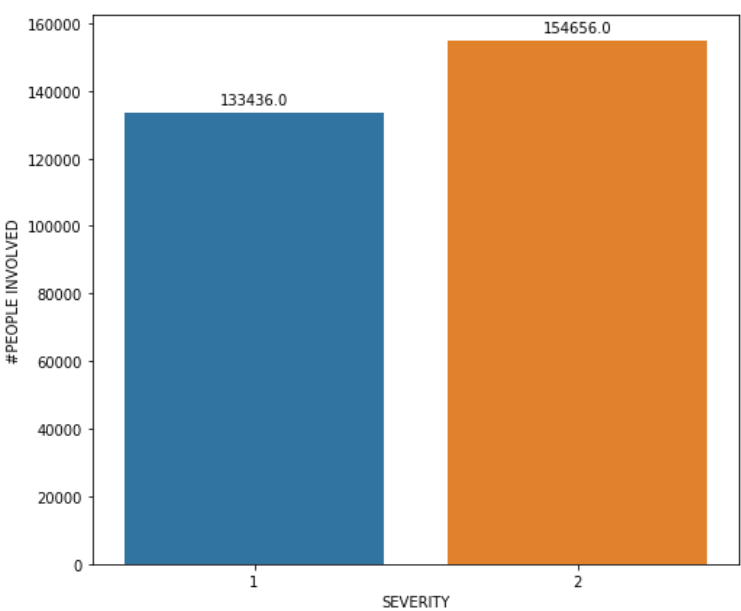








Exploratory Data Analysis (numerical variables)



Data Preparation

- Normalize the feature set

```
from sklearn import preprocessing
```

```
X = preprocessing.StandardScaler().fit(features).transform(features)
```

- Dataset split into train and test sets

```
Train set: (90600, 14) (90600,)
```

```
Test set: (22650, 14) (22650,)
```

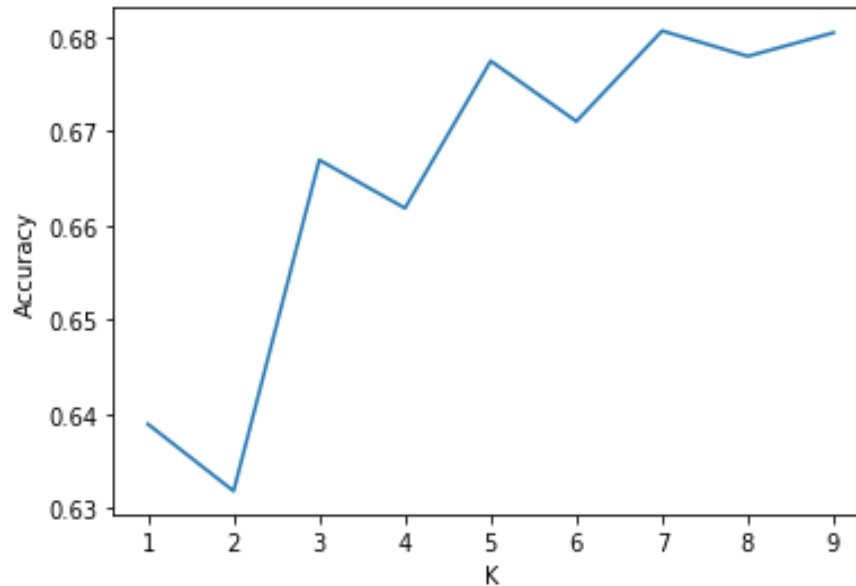
K Nearest Neighbours

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



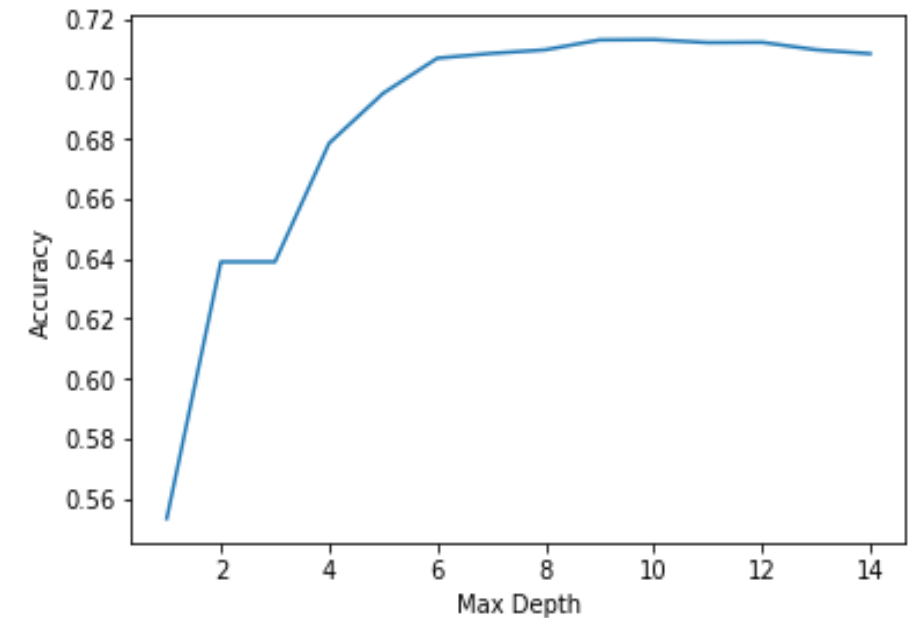
Decision Tree

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



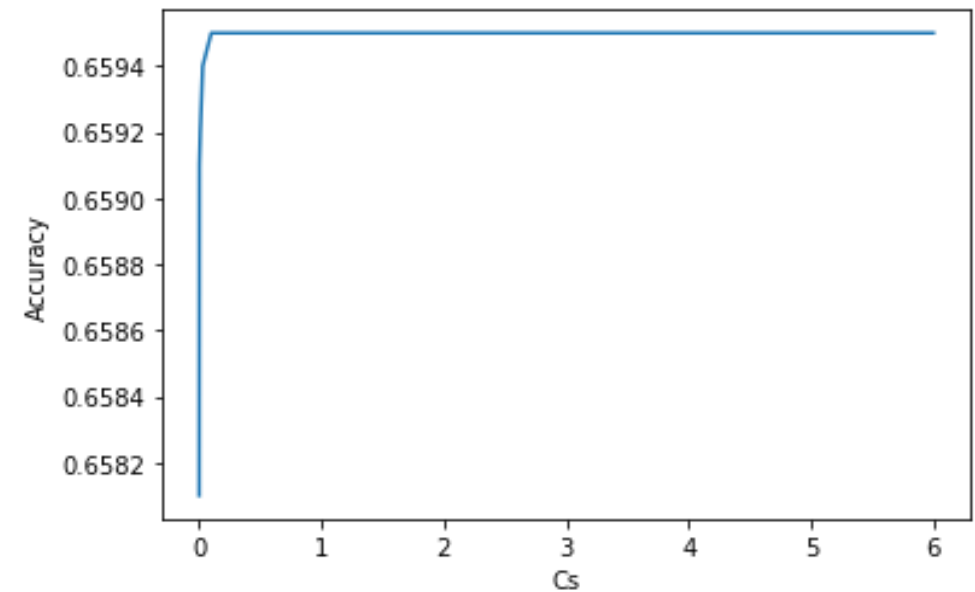
Logistic Regression

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

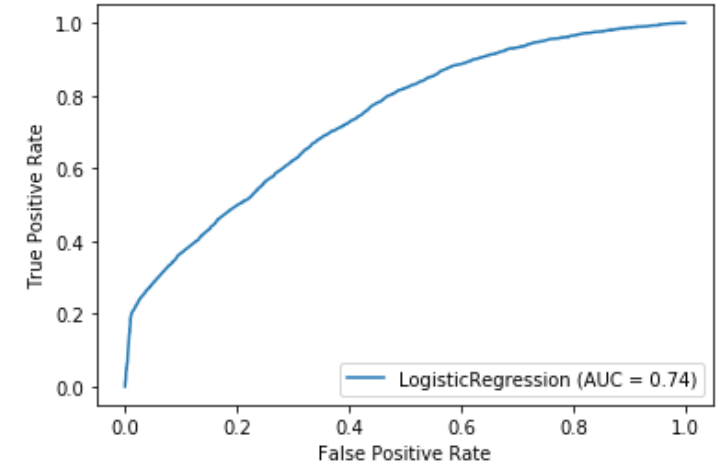
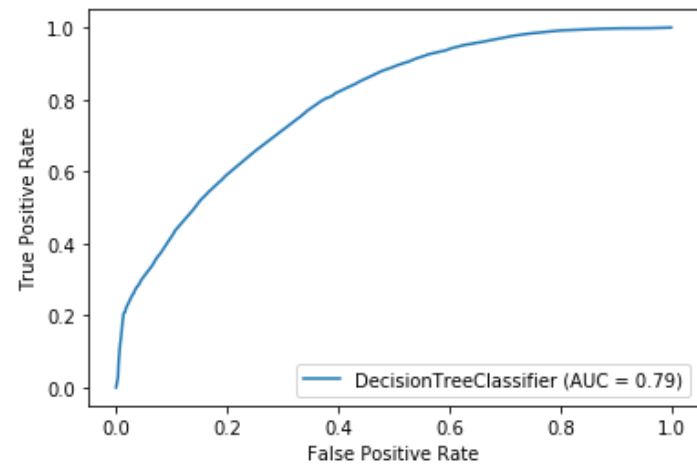
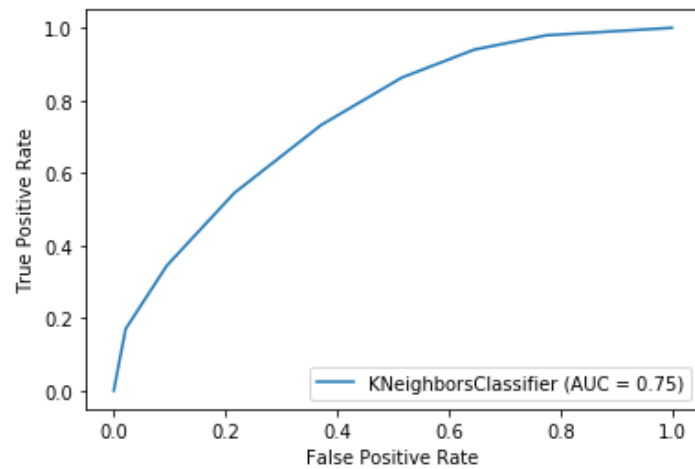


Results

- 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 |

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 0.79