

PREDICTING ROAD ACCIDENT SEVERITY

New Ru Wee

3 September 2020

1. Introduction

1.1 Background

Road accidents is a major problem in the modern world – it is said more people are killed by road accidents than plane crashes. The objective of this analysis to investigate which are the factors driving the severity of a collisions and provide insights to the community on possible to reduce road accidents.

The dataset is taken from SDOT Traffic Management Division, Traffic Records Group from 2004 to 2020. The dataset has 38 columns consisting of 194673 observations without any duplicates. This dataset is a hybrid combination of numerical and categorical data types.

1.2 Problem

The first problem would be to identify which are the features which are driving the road accidents severity. Each column in the dataset represent a feature. In process of doing so, the features of importance are also being identified.

1.3 Interest

This analysis will be of interest to the authorities, insurance company and community when trying to reduce severe road accidents. Road accidents not only cause loss of lives but also cause injuries and monetary loss to individual, insurance companies and community.

2. Data Cleansing

The dataset has many columns, several of which are repetitive – for example some columns are represented twice, once with text description and another time with numerical representation, both of which carries the same information. Therefore, the first step taken was to select the feature of interest based on correlation from Exploratory Data Analysis (EDA) – this will be further discussed in the EDA section of this report.

On top of that, there are null values (indicated by “NaN” in pandas) in some of the records which needs to be removed as all machine learning algorithm in Python would not be able to work if there are null values in the dataset.

The dataset has many columns, several of which are repetitive – for example some columns are represented twice, once with text description and another time with numerical representation, both of which carries the same information. Therefore, the first step taken was to select the feature of interest based on correlation from Exploratory Data Analysis (EDA) – this will be further discussed in the EDA section of this report.

On top of that, there are null values (indicated by “NaN” in pandas) in some of the records which needs to be removed as all machine learning algorithm in Python would not be able to work if there are null values in the dataset.

In order to use the sklearn library in Python to run the machine learning methods, all the columns have to be converted to numpy array.

Categorical data has to be converted to numerical representation so that it could fit into the machine learning algorithm. This was achieved using LabelEncoder from sklearn library. Next the data has to be normalized using MinMaxScaler also from sklearn library.

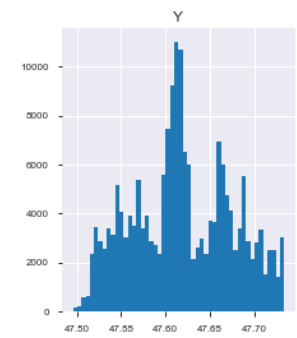
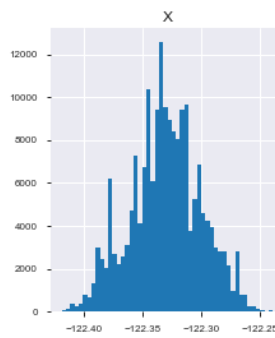
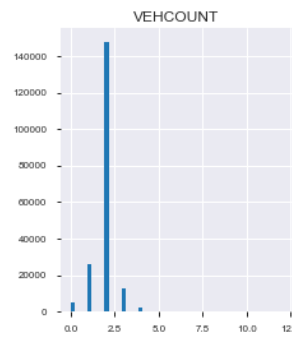
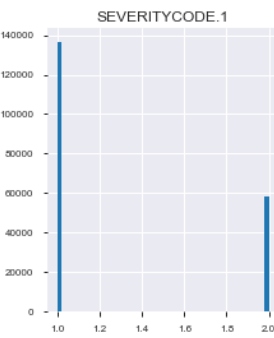
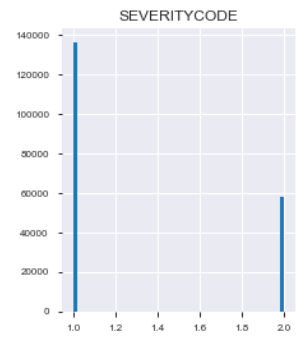
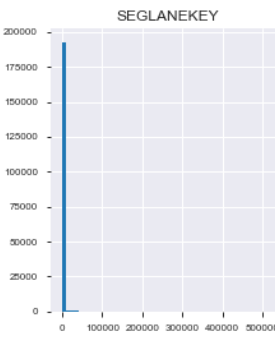
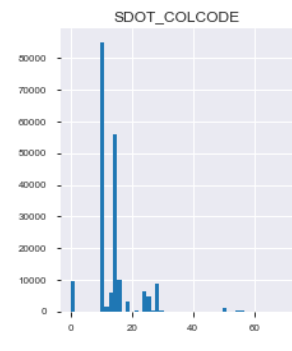
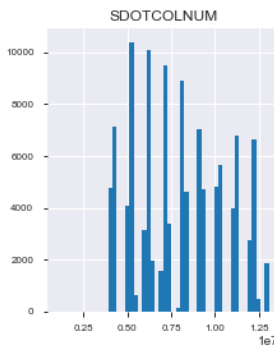
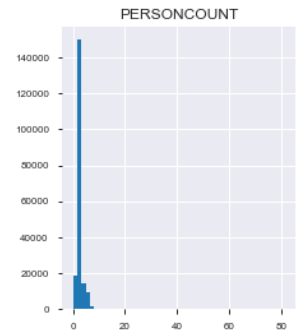
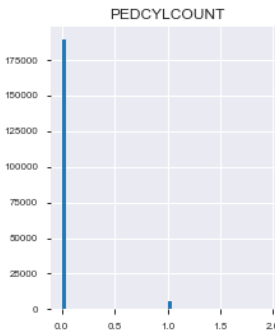
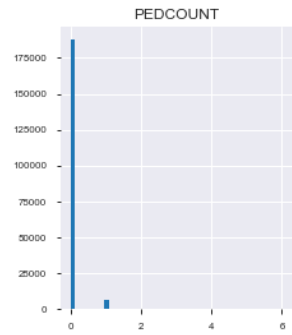
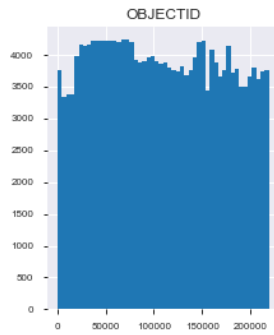
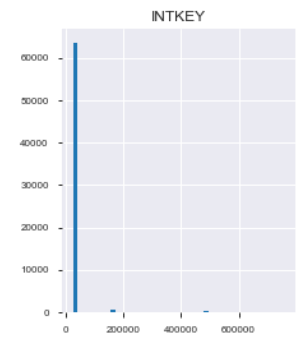
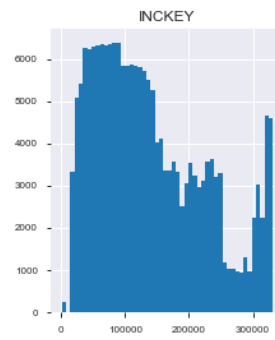
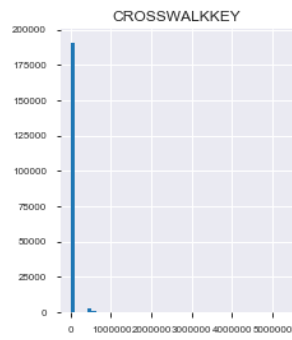
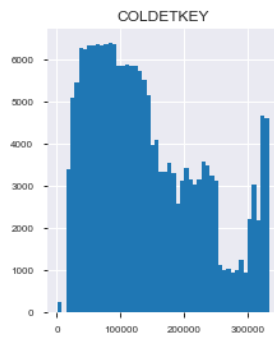
As the dataset is imbalanced, we need to perform upsampling to deal with imbalanced dataset.



3. EDA

The EDA was first done using pandas dataframe’s describe function to get the Central Frequency of each feature. Next matplotlib’s graphing function was used to visualize the EDA separately for both numerical and categorical features.

Last but not least, I also run pandas profiling to get the EDA – this is where I got the correlation between the features.





4. Classification

6 types of classification methods available in sklearn namely Logistic Regression, Decision Tree, KNN, Linear Discriminant Analysis, Gaussian Naïve Bayes and Support Vector Machine have been used to model the data.

5. Conclusion

From the accuracy obtained, it seems like Decision Tree is the best model for this dataset.

6. Recommendation