# Predicting Accident Severity in Seattle using Weather, Road conditions and Light conditions.

Elijah Chipato

September 6, 2020

## 1 Data Section

The example dataset for Seattle city is used in this study. Figure 1 shows the first 5 rows of the original



Figure 1: First 5 rows of the dataset and a couple of columns.

dataset which has about 38 features/columns. However, in this study only a couple of features will be used for the modelling. These features are shown on Figure 2 below. The severity code feature will be used as



Figure 2: First 5 rows of the selected feature to be used in this study.

the label of the dataset to train the models. In Figure 2 there are columns of incident date (INCDATE) and speeding (SPEEDING), these will be used during feature engineering to assess whether the day of the week is important in predicting accident severity. If observed that it has no impact these two will be dropped. To assess the quality of the data a heat-map is used to reveal the location or concentration of null values in our dataset. Figure 3 shows a heatmap which reveals that the speeding column has the most null values in the selected features for modelling. WEATHER, ROADCOND and LIGHTCOND also have similar number of null values. For the speeding column the attributes are Y for yes and NaN for No meaning the vehicle was not speeding. The NaN values are replaced by N to represent No. The null values for the rest of the features namely WEATHER, ROADCOND and LIGHTCOND are dropped. Figure 4 shows the new heatmap after dropping all null values and it shows that the datset is now clean and free of null values.
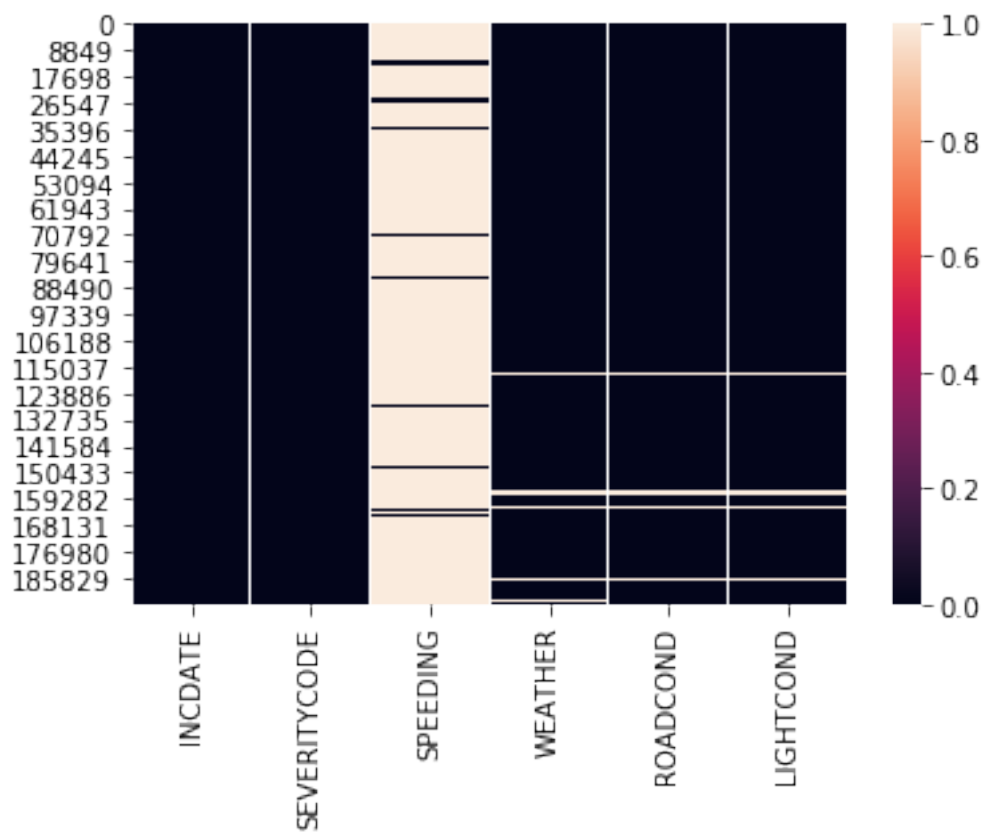
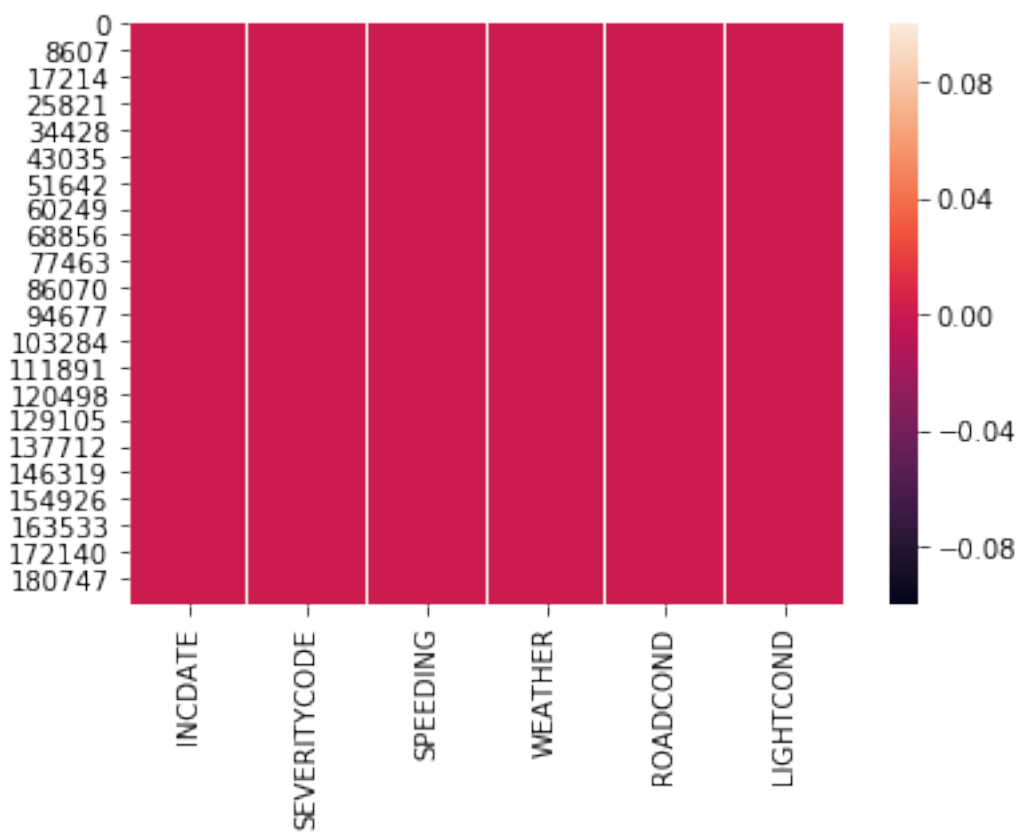Figure 3: Heatmap showing null values in the dataset.



Figure 4: Heatmap showing absence of null values in the dataset.