

CAPSTONE PROJECT - CAR ACCIDENT SEVERITY

1. INTRODUCTION

For the final capstone project in the IBM certificate course, we want to analyse the accident “severity” in terms of human fatality, traffic delay, property damage, or any other type of accident bad impact. The data was collected by Seattle SPOT Traffic Management Division and provided by Coursera via a link. This dataset is updated weekly and is from 2004 to present. It contains information such as severity code, address type, location, collision type, weather, road condition, speeding, among others.

The target audiences of this study are those people who really care about the traffic records, especially in the transportation department. Also, we want to figure out the reason for collisions and help to reduce accidents in the future.

2. DATA

There are 194,673 observations and 38 variables in this data set. Since we would like to identify the factors that cause the accident and the level of severity, we will use SEVERITYCODE as our dependent variable Y, and try different combinations of independent variables X to get the result.

Since the observations are quite large, we may need to filter out the missing value and delete the unrelated columns first. Then we can select the factor which may have more impact on the accidents, such as address type, weather, road condition, and light condition. The target Data to be predicted under (SEVERITYCODE 1-prop damage 2-injury) label.

Other important variables include:

- ADDRTYPE: Collision address type: Alley, Block, Intersection
- LOCATION: Description of the general location of the collision
- PERSONCOUNT: The total number of people involved in the collision helps identify severity involved
- PEDCOUNT: The number of pedestrians involved in the collision helps identify severity involved
- PEDCYLCOUNT: The number of bicycles involved in the collision helps identify severity involved
- VEHCOUNT: The number of vehicles involved in the collision identify severity involved
- JUNCTIONTYPE: Category of junction at which collision took place helps identify where most collisions occur

- **WEATHER:** A description of the weather conditions during the time of the collision
- **ROADCOND:** The condition of the road during the collision
- **LIGHTCOND:** The light conditions during the collision
- **SPEEDING:** Whether or not speeding was a factor in the collision (Y/N)
- **SEGLANEKEY:** A key for the lane segment in which the collision occurred
- **CROSSWALKKEY:** A key for the crosswalk at which the collision occurred
- **HITPARKEDCAR:** Whether or not the collision involved hitting a parked car

3. METHODOLOGY

We used Jupyter Notebook to do the data analysis. To generate the table and graph for the dataset, we imported Python libraries (Pandas, Numpy, Matplotlib, and Seaborn).

First, we imported the data through `pd.read_csv`. We noticed that it had 194,673 rows and 38 columns. Therefore, we narrowed it down to 8 columns ('Severity', 'X', 'Y', 'Location', 'Vehcount', 'Weather', 'Roadcond', 'Lighdcond') and delete the missing values, which made the final dataset with 184,167 observations and 8 variables.

We have to select the most important features to weigh the severity of accidents in Seattle. Among all the features, the following features have the most influence in the accuracy of the predictions - The '**WEATHER**', '**ROADCOND**' and '**LIGHTCOND**' attributes.

The dataset in the original form is not ready for data analysis. In order to prepare the data, first, we need to drop the non-relevant columns. In addition, most of the features are of object data types that need to be converted into numerical data types. After analysing the data set, I have decided to focus on only four features, severity, weather conditions, road conditions, and light conditions, among others.

To get a good understanding of the dataset, I have checked different values in the features. The results show, the target feature is imbalance, so we use a simple statistical technique to balance it.

I have run a value count on road ('ROADCOND') and weather condition ('WEATHER') to get ideas of the different road and weather conditions. I also have run a value count on light condition ('LIGHTCOND'), to see the breakdowns of accidents occurring during the different light conditions.

After balancing SEVERITYCODE feature, and standardizing the input feature, the data has been ready for building machine learning models.

I have employed three machine learning models:

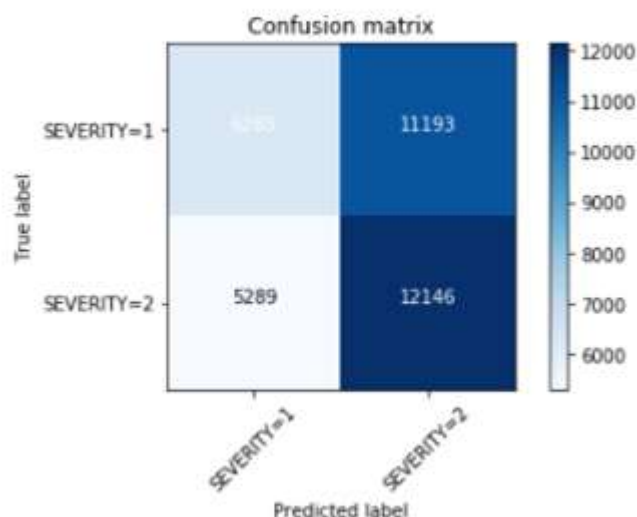
- K Nearest Neighbour (KNN)
- Decision Tree
- Logistic Regression

After importing necessary packages and splitting pre-processed data into test and train sets, for each machine learning model, I have built and evaluated the model.

4. RESULTS

Confusion Matrix:

```
Confusion matrix, without normalization
[[ 6285 11193]
 [ 5289 12146]]
```



Classification Report:

```
print (classification_report(y_test, LRpred))
```

	precision	recall	f1-score	support
1	0.54	0.36	0.43	17478
2	0.52	0.70	0.60	17435
micro avg	0.53	0.53	0.53	34913
macro avg	0.53	0.53	0.51	34913
weighted avg	0.53	0.53	0.51	34913

Comparison:

	F1-score	Jaccard-score	Log Loss
KNN	0.51	0.52	NA
Decision Tree	0.54	0.56	NA
Logistic Regression	0.52	0.53	0.68

5. OBSERVATIONS

Decision Tree is the best algorithm to be used in this case.

6. CONCLUSION

Based on historical data from the collision in Seattle, we can conclude that particular weather, road and light conditions have an impact on whether or not the car ride could result in one of the two classes property damage (class 1) or injury (class 2).