Capstone Project: Car Accident Severity

Applied Data Science Capstone by IBM/Coursera

Introduction: Business Problem

Average number of car accidents in the U.S. every year is 6 million.

3 million people in the U.S. are injured every year in car accidents.

Analysing the conditions that contribute to these accidents would lead to the prevention of significant loss of life and financial resources.

The project is aimed at predicting the severity of a car accident given the location, weather, road and visibility conditions in order to reduce the frequency of car collisions in a community based on a dataset provided by Seattle PD.

Consequently, this analysis would aid drivers to exercise more caution while driving or even choose an alternative route or time for their travel if possible.

It could also potentially help the local government, the police department and car insurance providers to gain deeper insight into road accidents.

_

Data

The detailed dataset of all road collisions (since 2004 to present) can be found here. This data was provided by the Seattle Police Department and recorded by the Traffic Records Department. The dataset consists of 37 independent fields and 194673 records, which includes both numerical and categorical data. The dependent field or label for the data set is SEVERITY CODE, which describes the fatality of an accident. The values under this label are categorised into fatality (3), serious injury (2b), injury (2), prop damage (1) and unknown (o).

Methodology

First, the necessary packages are imported.

```
#import necessary packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import seaborn as sns
from scipy import stats
from pylab import rcParams
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import jaccard_score, log_loss
from sklearn.metrics import precision_score, recall_score
from sklearn.metrics import f1_score
from sklearn.model_selection import train_test_split
```

Then the dataset is read from the .csv file and assigned to a DataFrame df.

DATA WRANGLING

The fields necessary for the analysis are identified and a new dataset df_acc is created with these selected fields. These fields are -

SEVERITYCODE,

ADDRTYPE,

LOCATION,

JUNCTIONTYPE,

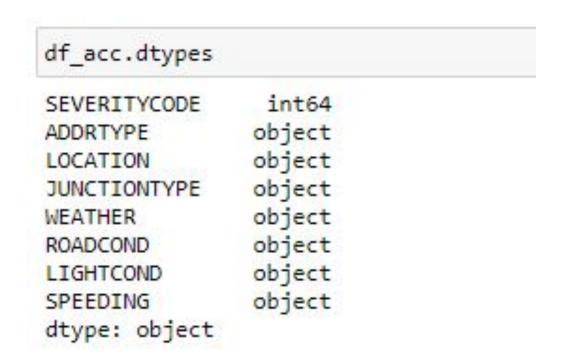
WEATHER,

ROADCOND,

LIGHTCOND,

SPEEDING

The datatype of each of the included columns is examined and the following output is observed.



Null values in the DataFrame are then identified using the following code.

```
missing_data = df_acc.isnull()
for column in missing_data.columns.values.tolist():
    print(column)|
    print (missing_data[column].value_counts())
    print("")
```

Subsequently, null values in the fields 'JUNCTIONTYPE', 'WEATHER', 'ROADCOND', 'LIGHTCOND', and 'SPEEDING' are replaced with suitable values.

UNDERSAMPLING

On performing a count of values on the DataFrame, it is observed that the dependent variable is an unbalanced label.

```
df_acc.value_counts("SEVERITYCODE")

SEVERITYCODE
1 136485
2 58188
dtype: int64
```

Hence, undersampling is performed on df_acc to ensure that the prediction data obtained after passing the sample data through the various Machine Learning models is not biased.

```
#balance label SEVERITYCODE"

target="SEVERITYCODE"
minority_class_len = len(df_acc[df_acc[target] ==2])
majority_class_indices = df_acc[df_acc[target] ==1].index
random_majority_indices = np.random.choice(majority_class_indices,minority_class_len, replace = False)
minority_class_indices = df_acc[df_acc[target] ==2].index

under_sample_indices = np.concatenate([minority_class_indices, random_majority_indices])
df_acc = df_acc.loc[under_sample_indices]
df acc.value counts("SEVERITYCODE")
```

```
SEVERITYCODE
2 58188
```

1 58188

dtype: int64

EXPLORATORY DATA ANALYSIS

The DataFrame is examined to obtain the following output.

df_acc.describe(include="all")

	SEVERITYCODE	ADDRTYPE	LOCATION	JUNCTIONTYPE	WEATHER	ROADCOND	LIGHTCOND	SPEEDING
count	116376.000000	115473	115136	116376	116376	116376	116376	116376
unique	NaN	3	19829	7	11	9	9	2
top	NaN	Block	AURORA AVE N BETWEEN N 117TH PL AND N 125TH ST	Mid-Block (not related to intersection)	Clear	Dry	Daylight	N
freq	NaN	71587	177	49590	68073	76182	71623	110371
mean	1.500000	NaN	NaN	NaN	NaN	NaN	NaN	NaN
std	0.500002	NaN	NaN	NaN	NaN	NaN	NaN	NaN
min	1.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN
25%	1.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN
50%	1.500000	NaN	NaN	NaN	NaN	NaN	NaN	NaN
75%	2.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN
max	2.000000	NaN	NaN	NaN	NaN	NaN	NaN	NaN

Further, the distribution of SEVERITY CODE by different columns is explored through the *value_counts*() as shown below.

SEVERITYCODE JUNCTIONTYPE	1	2	% of Total
Mid-Block (not related to intersection)	30186	19404	42.61
At Intersection (intersection related)	15036	27174	36.27
Mid-Block (but intersection related)	6613	7297	11.95
Driveway Junction	3126	3234	5.47
Unknown	2543	402	2.53
At Intersection (but not related to intersection)	637	623	1.08
Ramp Junction	47	54	0.09

% of Tota	2	1		SEVERITYC WEAT	Total	2	1	SEVERITYCODE
58.4	35840	2233	Clear	(177	120	57	AURORA AVE N BETWEEN N 117TH PL AND N 125TH ST
17.6	11176	9353	ining	Rai		1,000	10000	TO SERVICE EXPRESSION OF THE SERVICE AND ADDRESS
14.3	8745	8006	ercast	Over	166	94		N NORTHGATE WAY BETWEEN MERIDIAN AVE N AND CORLISS AVE N
8.3	1900	7780	nown	Unkn	166	86	80	BATTERY ST TUNNEL SB BETWEEN AURORA AVE N AND ALASKAN WY VI SB
0.4	171	322	owing	Sno	164	107	57	6TH AVE AND JAMES ST
0.3	116	270	Other	C	164	78	86	BATTERY ST TUNNEL NB BETWEEN ALASKAN WY VI NB AND AURORA AVE N
0.3	187	172	moke	Fog/Smog/Sn	157	88	69	AURORA AVE N BETWEEN N 130TH ST AND N 135TH ST
0.0	28	32	Rain	Sleet/Hail/Freezing	156	94	62	RAINIER AVE S BETWEEN S BAYVIEW ST AND S MCCLELLAN ST
0.0	15	15	d/Dirt	Blowing Sand	134	78	56	ST SEATTLE BR EB BETWEEN ALASKAN WY VI NB ON RP AND DELRIDGE-W SEATTLE BR EB ON RP
0.0	7	4	swind	Severe Cross	132	63	69	ALASKAN WY VI NB BETWEEN S ROYAL BROUGHAM WAY ON RP AND SENECA ST OFF RP
0.0	3	1	loudy	Partly Clo	122	68	54	ALASKAN WY VI SB BETWEEN COLUMBIA ST ON RP AND ALASKAN WY VI SB EFR OFF RP
of Total	2 %			SEVERITYCODE SPEEDING	I)	of Tota	%	SEVERITYCODE 1 2 % of Total SEVERITYCODE 1 2 ROADCOND LIGHTCOND
94.84	57	546	5571	N	1	61.54		Dry 36118 40064 65.46 Daylight 33079 38544
								AND

SEVERITYCODE ROADCOND	1	2	% of Total	SEVERITYCODE LIGHTCOND	1	2	% of Total	SEVERITYCODE SPEEDING	1	2	% of Tota
Dry	36118	40064	65.46	Daylight	33079	38544	61.54	N	55714	54657	94.84
Wet	13463	15755	25.11	Dark - Street Lights On	14549	14475	24.94	Y	2474	3531	5.16
Unknown	7742	1809	8.21	Unknown	7270	1695	7.70				
Ice	414	273	0.59	Dusk	1624	1944	3.07				
Snow/Slush	340	167	0.44	Dawn	696	824	1.31				

MODEL EVALUATION

The values in the current dataset are not in a suitable format to be used as testing and training data for the prediction models. Thus, label encoding is carried out to convert the categorical data to numerical data ranging from digits 0 to 9. The data is then normalized using the following code.

The normalized data is then split into 70% training and 30% testing data.

```
#Splitting the data as 70 % for training and 30 % for testing

x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.3, random_state=0)
print("Train set:", x_train.shape, y_train.shape)
print("Test set:", x_test.shape, y_test.shape)

Train set: (81463, 5) (81463,)
Test set: (34913, 5) (34913,)
```

MODELING AND PREDICTIONS

To determine the best fit for predicting the SEVERITY CODE, three models are considered.

K-Nearest Neighbors

F1 Score: 0.6148349280182468

```
# Training the Model
from sklearn.neighbors import KNeighborsClassifier
k=25

kneigh = KNeighborsClassifier(n_neighbors = k).fit(x_train, y_train)
k_y_pred = kneigh.predict(x_test)
k_y_pred[0:5]

array([2, 1, 1, 1, 1], dtype=int64)

#Model Evaluation
j1=jaccard_score(y_test,k_y_pred)
f1=f1_score(y_test,k_y_pred, average = 'macro')
print("Jaccard Score: ",j1)
print("F1 Score: ",f1)

Jaccard Score: 0.4510596594389318
```

Decision Tree

```
# Training the Model
from sklearn.tree import DecisionTreeClassifier
dt = DecisionTreeClassifier (criterion = 'entropy', max_depth = 7)

dt.fit(x_train, y_train)
dt_y_pred = dt.predict(x_test)
dt_y_pred[0:5]

array([2, 1, 1, 1, 1], dtype=int64)
```

```
#Model Evaluation
j2=jaccard_score(y_test,dt_y_pred)
f2=f1_score(y_test,dt_y_pred, average = 'macro')
print("Jaccard Score: ",j2)
print("F1 Score: ",f2)
```

Jaccard Score: 0.43849031227672003 F1 Score: 0.6171862904671328

Linear Regression

```
# Training the Model
from sklearn.linear model import LogisticRegression
from sklearn.metrics import confusion matrix
lr = LogisticRegression(C = 6, solver = 'liblinear').fit(x train, y train)
lr y pred = lr.predict(x test)
lr y prob = lr.predict proba(x test)
lr y prob
array([[0.39458097, 0.60541903],
       [0.66094284, 0.33905716],
       [0.64601659, 0.35398341],
       ...,
       [0.53628072, 0.46371928],
       [0.58710319, 0.41289681],
       [0.51892818, 0.48107182]])
#Model Evaluation
j3=jaccard score(y test,lr y pred)
f3=f1 score(y test, lr y pred, average = 'macro')
print("Jaccard Score: ",j3)
print("F1 Score: ",f3)
print("Log Loss: ",log loss(y test,lr y prob))
Jaccard Score: 0.45384
F1 Score: 0.6082536169528481
Log Loss: 0.6561502998049485
```

Results and Discussion

- -Evaluation metrics used to test the accuracy of our models were the Jaccard index and the f-1 score.
- -Choosing different k, max depth, and hyperparameter C values helped to improve the accuracy of our models.
- -Through our exploratory data analysis, we have also observed that approximately 43% of all recorded accidents occur at the Blocks but the severity of the accident tends to be higher (SEVERITY CODE= 2 Injury) at the intersection as compared to the blocks.
- -From this exercise, we can draw the result that our model can be trained to determine the severity of an accident to a certain extent using the given dataset and specific parameters such as WEATHER, ROADCOND, LIGHTCOND, SPEEDING, and JUNCTION TYPE.

Conclusion

Based on the evaluation of the above models, it can be concluded that Linear Regression is the most ideal model for this case.