

Analysis of Car Accidents In USA

By Lior Ben David 208778225
and Sahar Ziv 312160369

Car Accidents link to code

Contents

Introduction	3
Correcting the data	3
Challenges and bias in the data -	5
Model Training	6
Results and conclusions.....	6
Appendices.....	8
Appendix 1 –	8
Appendix 2 –	9
Appendix 3 –	10
Appendix 4 –	11

Introduction

In our analysis we focused on car accidents data from USA.

In the US, a fatal accident occurs every 15 minutes on average.

The issue is of interest to the insurance companies and the transportation and government ministries.

In our analysis we will try to classify the severity of the injury of the drivers involved in road accidents with the help of personal data, vehicle data and the weather.

The data was taken from the [Kaggle](#) website.

Our data was collected in 2018 and includes three files:

Number of features	Number of observations	Datasets
51	48443	Acc_18
54	120230	Pers_18
87	86105	Veh_18

The feature that links all the tables is 'CASENUM' which represents the number of accidents.

For each of the data sets, they tried to complete missing values (Imputation), therefore there are identical features, and the features that have already been completed we would like to delete to remain only with the original data.

Our label feature is 'INJ_SEV' which describes the severity of the injury to a person in the crash.

Correcting the data

At first, because our research question is about drivers, we only filtered out the drivers' data from the person's dataset. Where the feature 'PER_TYP' == '1' it means that it is the observation of the driver. So now we have 85,916 observations of drivers in our data.

The second step, we merged all data sets to one data set by the feature 'CASENUM' and remove duplicate columns.

After that, we checked if our data included null values, there was not and it seems weird to us, so we checked the distribution of values for each feature.

We noticed that the missing values are displayed in the data as 9/8, 98/99, 999, etc., there for we have summarized the data on the features in a table so that we can decide whether to delete them and whether we need to fill in missing values and if so, with what digit they are represented. Then we converted all the values that represent missing values to null.

Now our data is ready to imputation. We decided to use KNN imputation because it is simple and gets very good results for imputation. (We saved our data after imputation just because it takes a little bit time to do the KNN and we didn't want to do it every time we work on our analysis)

Our label contains 7 levels after imputation:

- 0 No Apparent Injury
- 1 Possible Injury
- 2 Suspected Minor Injury
- 3 Suspected Serious Injury
- 4 Fatal Injury
- 5 Injured, Severity Unknown
- 6 Died Prior to Crash

We decided to combine between the levels of the label to 2 levels:

- 0 - no injury to the driver or minor
- 1 - severe injury to the driver

After that we checked the correlations between the features and dropped variables that have a high correlation between them.

Then we perform several actions:

1. We paired our variables into numeric and categorical features.
2. We split our data to train data (70%) and test data (30%)
3. We scaled the numeric values.
4. We did One-hot Encoding on the categorical features.
5. We noticed that our data is imbalanced, there were more observations with label "0" than label "1" so we fixed this by oversampling with SMOTE so the number of observations from each category would be the same.

6. We checked the correlation of the numeric features before and after over-sampling and there was not a big difference which is good.
(Appendix 1)
7. We performed feature selection with LASSO.
8. We calculated the correlation between the features and the label to know if there is feasibility for our project. We have accepted that the three most influencing variables on the label are: (appendix 2)
 1. REST_USE_3 which means that shoulder and lap belt used by this person at the time of the crash.
 2. RELJCT2_3 which means that the accident was intersection Related
 3. TRAV_SP - the speed the vehicle was traveling prior to the occurrence of the crash
9. We checked whether we have outliers in the data by fitting logistic regression model. We noticed that the residuals are normally distributed with the mean 0. We saw that we had one outlier observation, so we dropped it. (Appendix 3)
10. Again we saved our data only to save time working on the project.

Challenges and bias in the data -

We presented the number of accidents per maker, and we saw that the number of accidents is higher to American manufacturers than non-American manufacturers.

It makes sense because our data is sampled from the USA, and it makes a choosing bias and can cause wrong conclusions.

Therefore, we decided to remove the features MAKE and MODEL.

We also saw that there are states that weren't sampled, which means that we have choosing bias in our data. We didn't have something to do with it, so we use it as is and refer to it in the conclusion.

One challenge that we had was unbalanced Data. To overcome it we over-sampled the data with SMOTE to make the label balanced.

Model Training

We trained 4 models.

1. **Logistic Regression** with L1 penalty ("Lasso") – We chose this model because it can easily predict a binary label.
2. **K-Nearest Neighbors** - We chose this model because it is simple and can predict well and does not require any prior preparation of the data.
3. **Classification Tree** - We chose a classification tree to check if it classifies the data well without scaling
4. **Random Forest** - We know that this model has good classification results because it uses classification from many trees and thus reduces its error.

All the model's tuning parameters were estimated in 5-fold Cross-Validation and we will show the results of the best estimators.

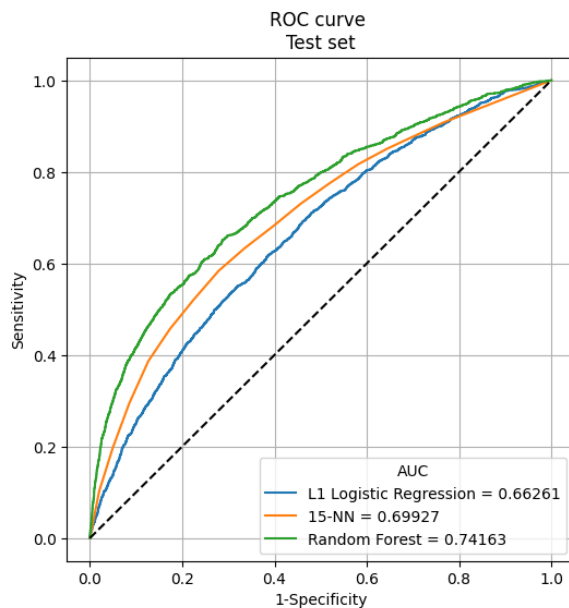
We will compare the method by Accuracy, Sensitivity, Specificity and AUC of ROC curve.

Results and conclusions

Results –

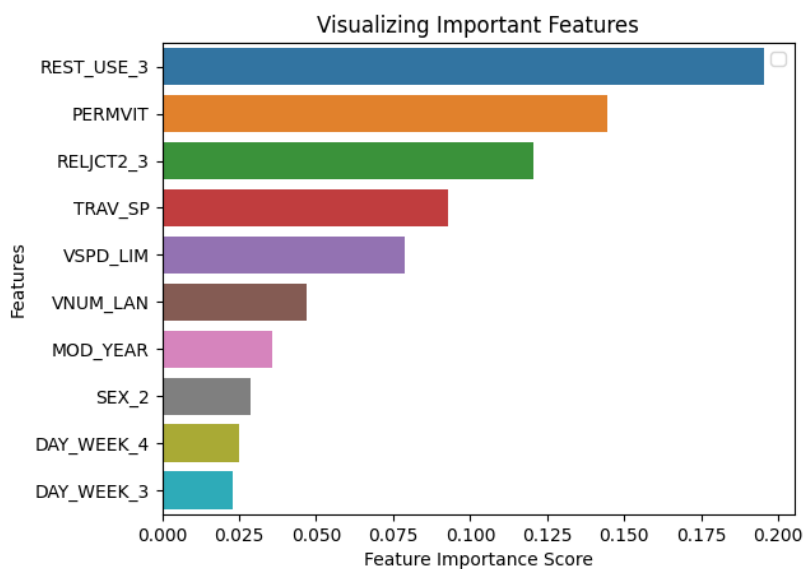
Model	Accuracy	Sensitivity	Specificity
Logistic regression	0.85998	0.2604	0.89677
KNN	0.66262	0.63356	0.6644
Random Forest	0.83371	0.49329	0.8546
Classification tree	0.942		

Although the accuracy of the classification tree is very high, we saw that it classified almost all observations into a certain group, so it is not a good model for our data. (appendix 4)



We can see that the Random Forest model is the best model because its AUC is the biggest.

This are the most important features:



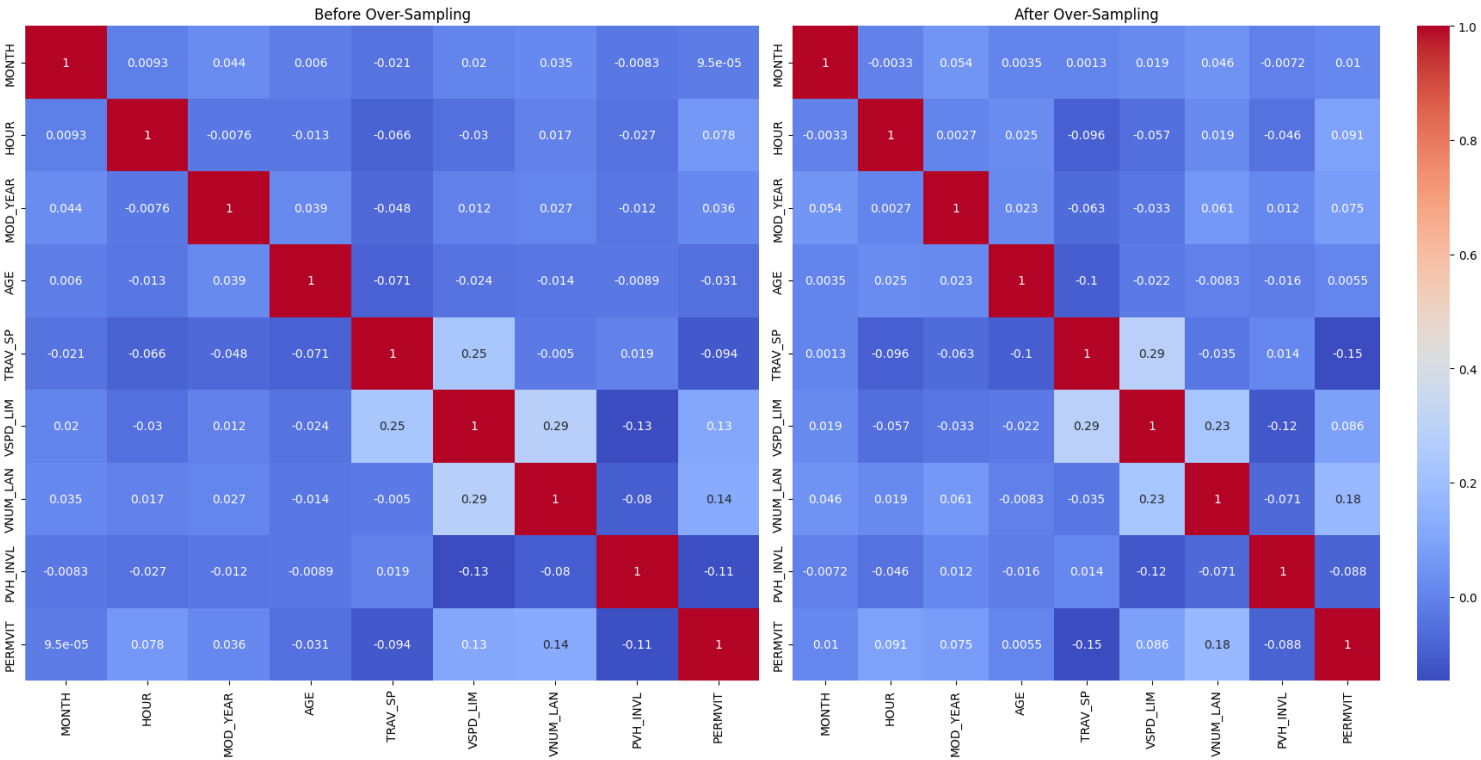
1. REST_USE_3 which means that shoulder and lap belt used by this person at the time of the crash.
2. PERMVIT – which is the number of persons in motor vehicles.
3. RELJCT2_3 which means that the accident was intersection Related.
4. TRAV_SP - the speed the vehicle was traveling prior to the occurrence of the crash.

Work to be continued –

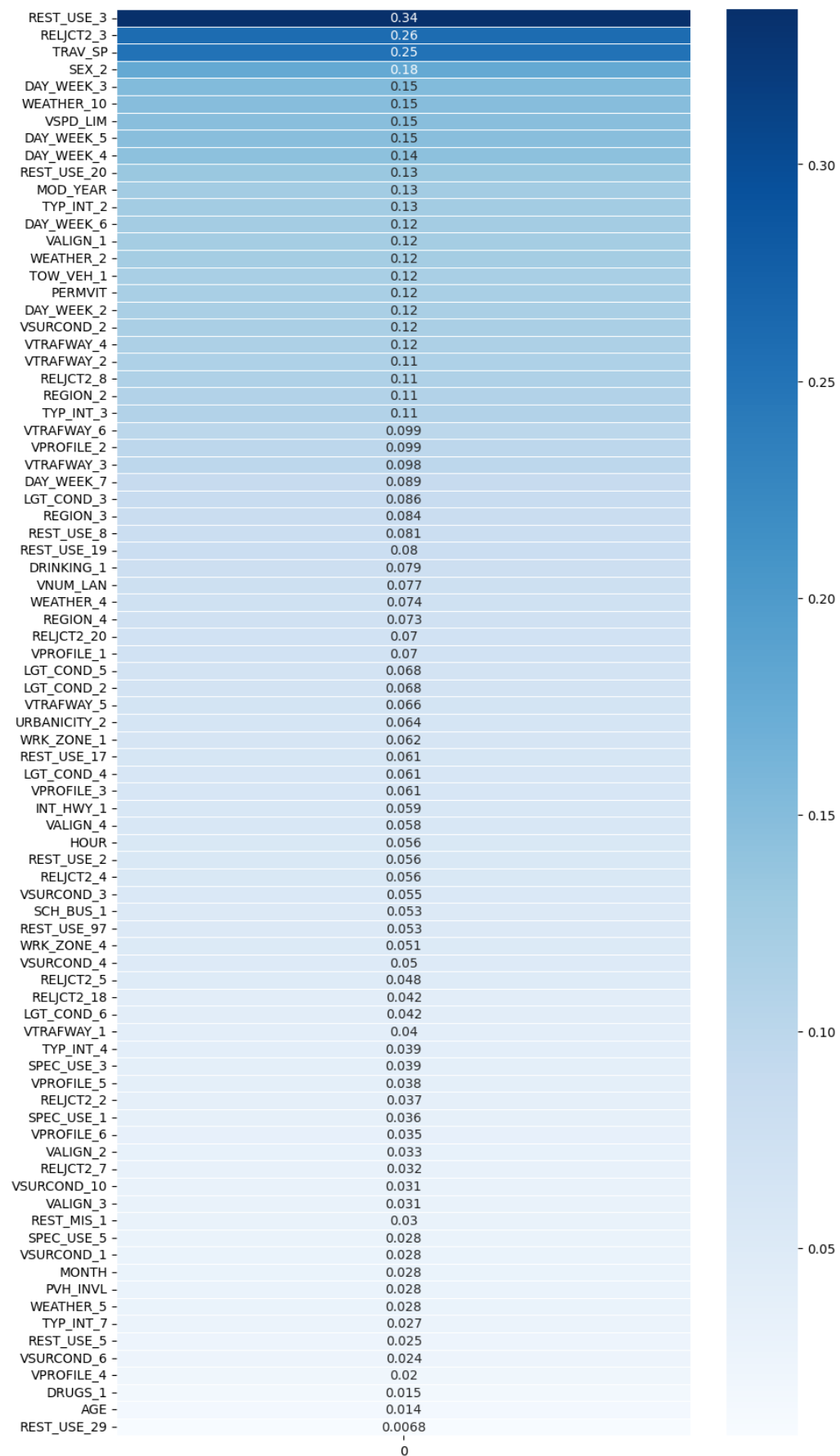
In the future, the performance of these models can also be evaluated on data from other years or other places.

Appendices

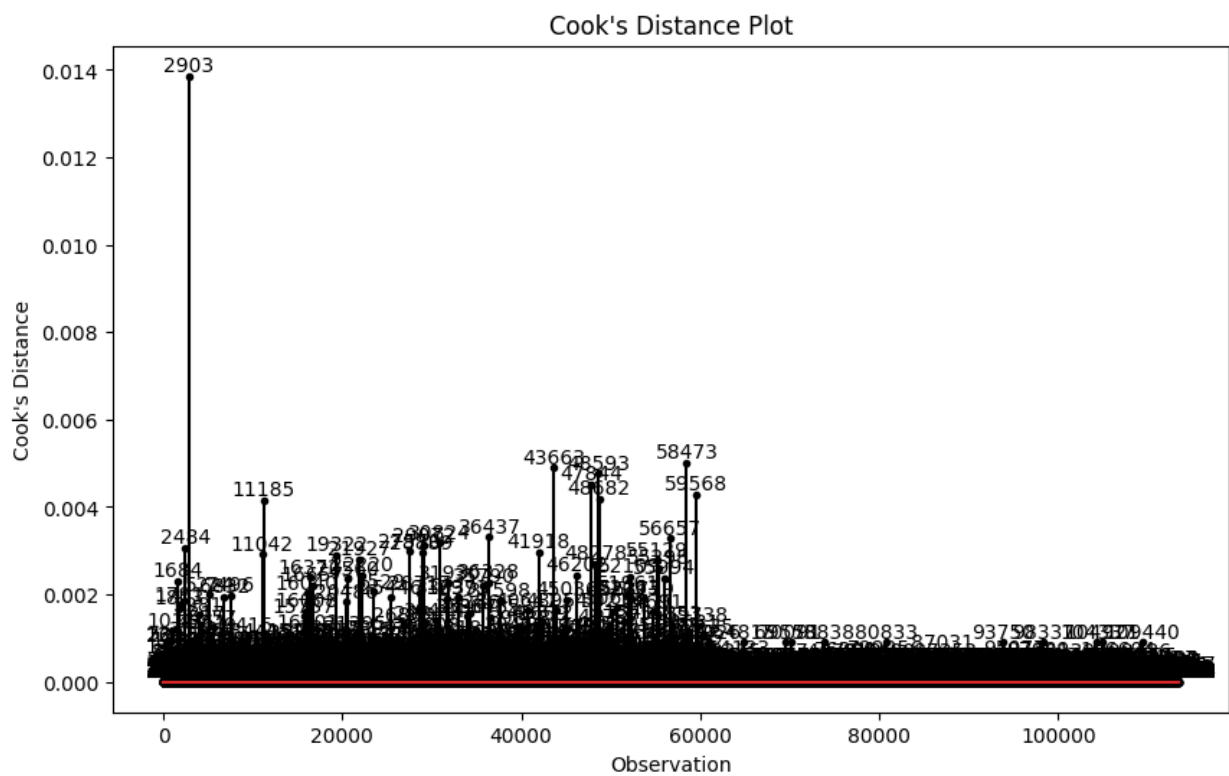
Appendix 1 –



Appendix 2 –



Appendix 3 –



Appendix 4 –

