

Introduction

It is important to analyse and predict traffic accidents in order to limit their occurrence. This project will focus on the analysis of UK road accident data given by the Department of Transport. The UK government publishes data studies on traffic accidents every year. As of 2019, there were 1,752 recorded road traffic fatalities and 153,158 reported casualties of all types, with 25,945 of them suffering significant injuries.

There are two objectives for this project; Conduct an exploratory analysis of 2019 road accident datasets to get insight into the current condition of road accidents in the UK; to investigate how well machine learning algorithms predict the occurrence and severity of traffic accidents in the United Kingdom to enhance road safety.

In this project, mining of road traffic data will be done using an association rule-based algorithm. The Apriori Algorithm is used to establish a relationship between features in the dataset, and then the rules are investigated with high lift and high support. K-nearest Neighbors and Naive Bayes classifiers are used in classifying the independent of attributes. To create clusters and evaluate them based on characteristics, K-means is utilized. The positive findings will surely help the government make better decisions and prevent traffic accidents.

Dataset

The dataset in this project is separated into four categories as collected from the UK government data portal.

- Accident Datasets
- Vehicle Datasets
- Casualty Datasets
- Adjustments Datasets
- Variable Lookup Dataset

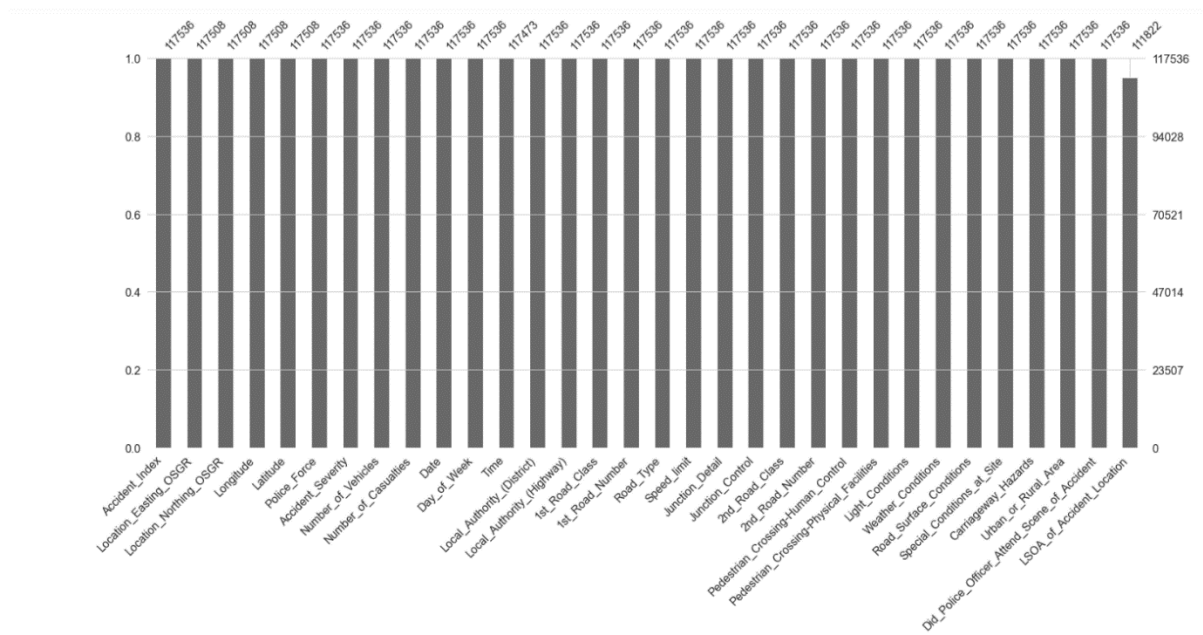
At most granular level, the Accidents dataset may be regarded as a superclass of the vehicle and casualty dataset linked via a shared 'accident index' identifier.

Data Pre-processing

The process of preparing data for analysis is known as data pre-processing. Cleaning the data, engineering new features through data transformation, and determining which portions of the data to keep and which to discard were all part of the process(Alkheder et al., 2017). This section explains how the dataset was prepared for analysis, modeling, and predictive analysis using pre-processing techniques.

Data Cleaning

Missing value



The source datasets' overall strategy was to assign unknown data items a value of (-1), which is a standard method in data reporting. However, there were times when data was completely missing. To comply with the basic strategy, the entire dataset was analysed to detect this (null,missing n/a, None). The initial thought was to replace all missing values with -1 to comply with the report's general approach, but a closer examination of some columns ('Location Northing OSGR', 'Location Easting OSGR','Longitude','Latitude', 'Time','LSOA of Accident Location') reveals that most of the data in them have unique values, and the percentage of missing values are insignificant. Missing values in these columns were dropped entirely while that of other columns were replaced with -1.

After standardising the value for all missing data, columns with less than 5% missing values (-1) were replaced with the most frequent (modal) value and the rows of the ones with more missing values (especially in the vehicle and casualty dataset) were dropped from the record. This was an appropriate strategy to utilize and not have a detrimental influence on the model's performance since a greater percentage of the dataset was available for analysis.

Feature Selection

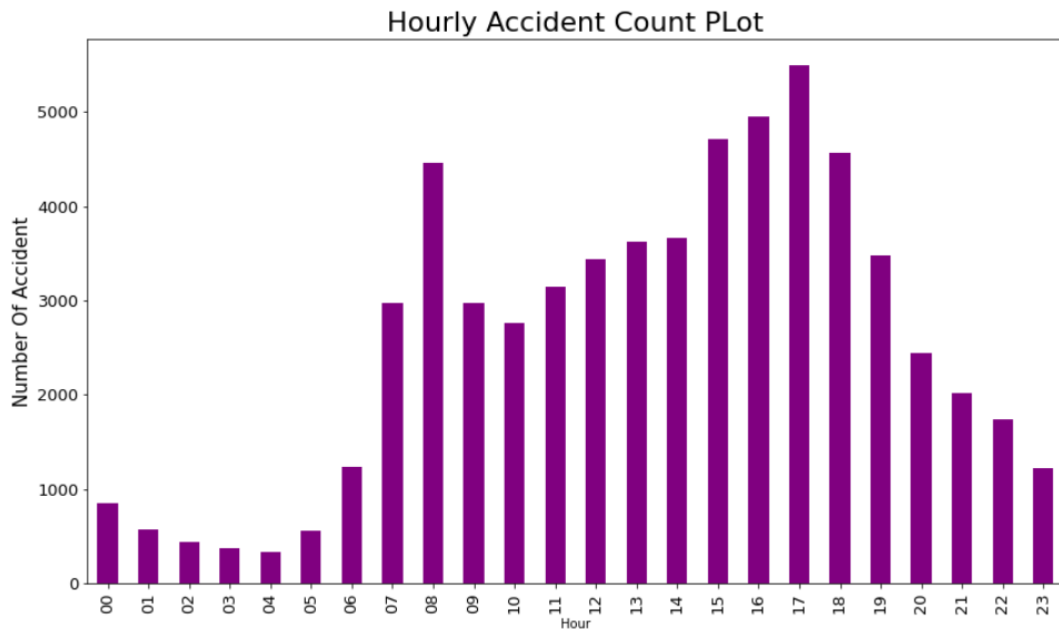
This project used an iterative feature selection method as different analyses and models needed different data transformations and feature engineering with the goal of enhancing the analysis quality. An initial examination of each of the three datasets was conducted to ensure that the structure of the datasets was consistent

Exploratory Data Analysis

An exploratory analysis of the dataset was done to better grasp the information and acquire a better knowledge of the attributes in the data. Focus points were considered, and certain EDA are made to guide this study.

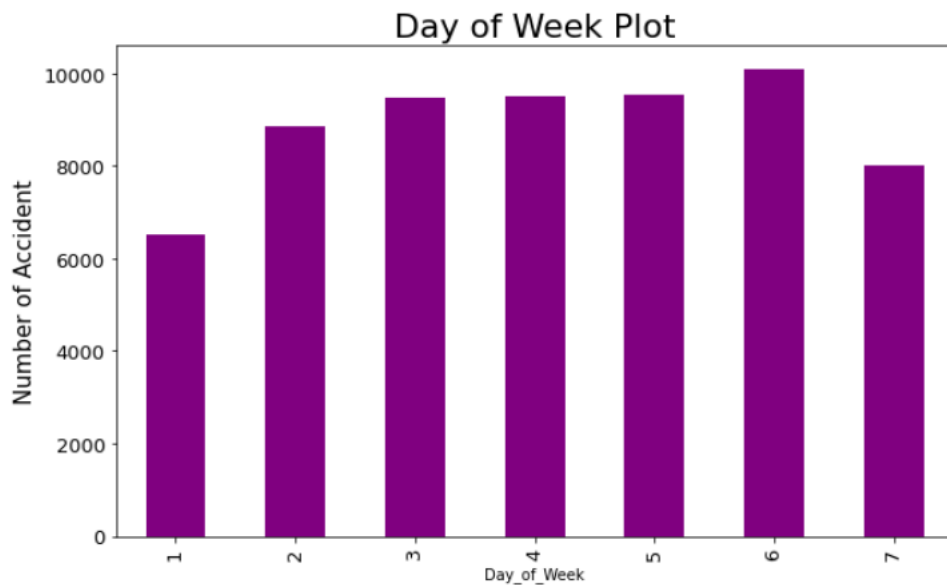
A. Are there significant hours of the day, and days of the week, on which accidents occur?

HOURS OF THE WEEK



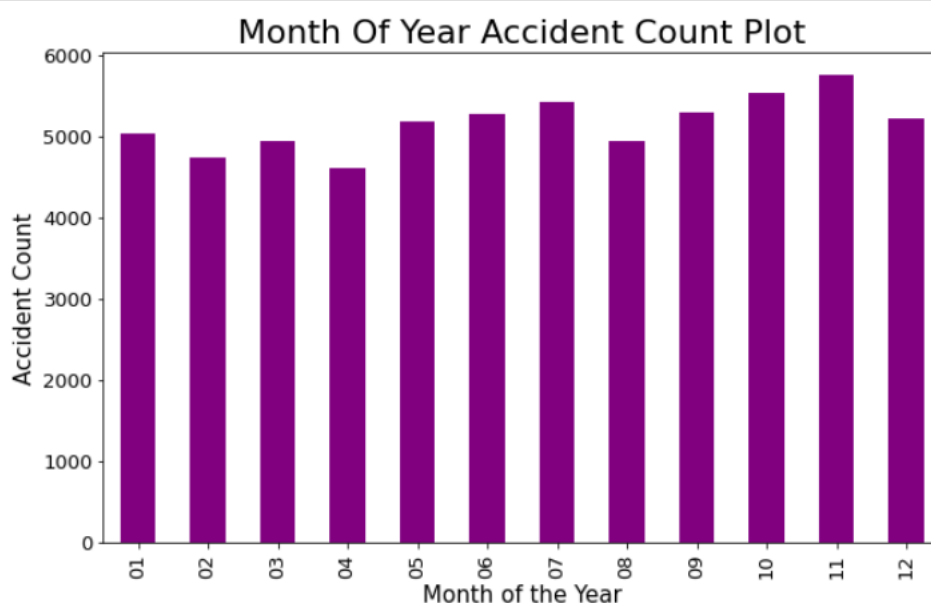
Accidents are most prevalent in the evening, as one might expect because people are going home after a hard day at work, and the plot reveals that more accidents occur at this time, with 5pm having the largest number of accidents. From 3 to 6 p.m., the number of accidents tends to increase. By 8 a.m., when people are hurrying to work, the graph reveals a substantial number of accidents during this time. The safest time to be on the road is between 12 to 5 a.m., when the roads are normally clear.

DAYS OF THE WEEK



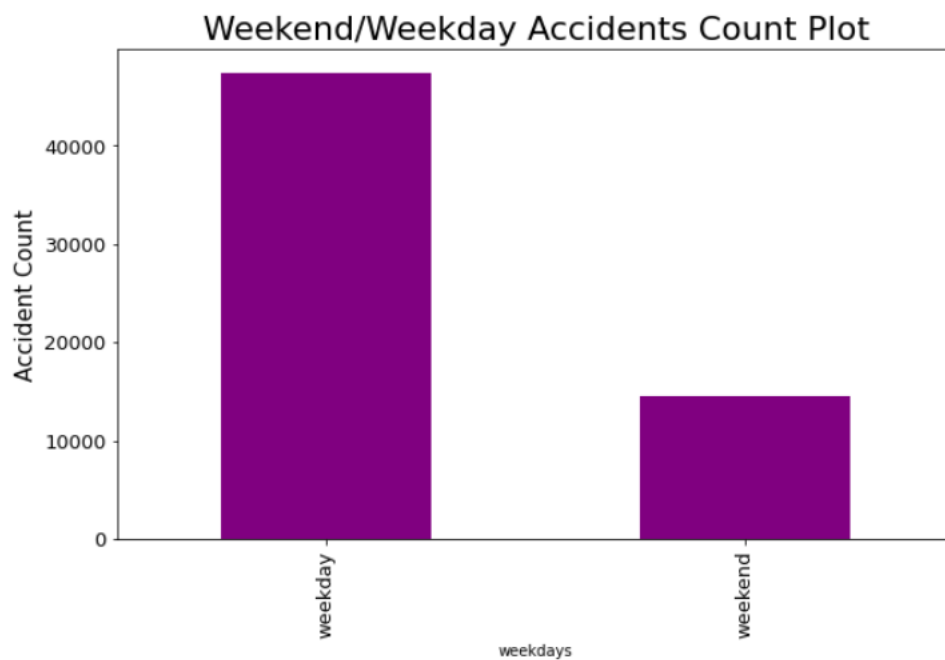
The graph shows Monday having the least accidents while Saturday has the highest. The high number of accidents on Saturday might be due to a multitude of causes. There are usually a lot more vehicles on the road on Saturday. People want to spend their days off work by going out which entails driving to numerous places. Second, drivers may be less alert or aware on Saturdays than they are during the week, which may hinder them from being as vigilant or aware as they are during the week.

MONTH OF THE YEAR



The month of November has the highest number of traffic accidents. While one may anticipate the winter to have the most number of accidents because of the poor weather and long nights, the summer and spring months constitute 3 of the top 6 in terms of the highest number of accidents. This might be because there are more vehicles on the road as people enjoy the sunlight.

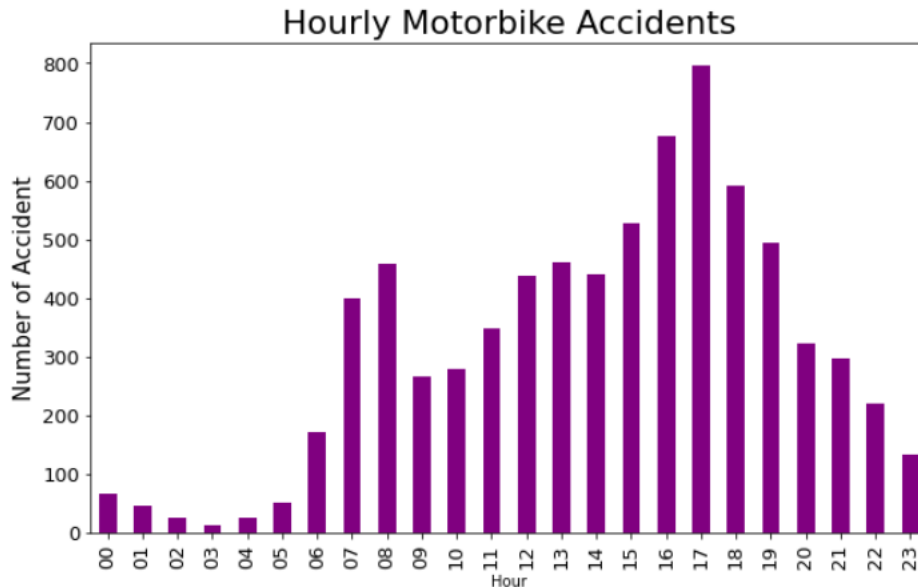
ACCIDENTS ON WEEKDAY OR WEEKEND



The previous graph reveals Saturday has more accidents than any day of the week.

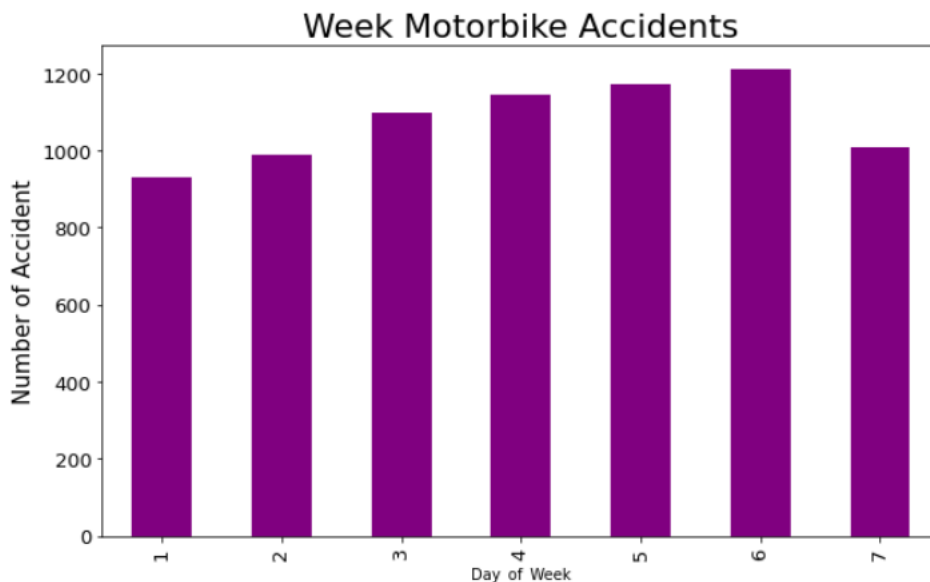
B. For motorbikes, are there significant hours of the day, and days of the week, on which accidents occur?

MOTORBIKES – HOURS OF THE DAY



According to the data, most motorbike accidents occur in the evening, between 4:00pm and 6:00 p.m., with a peak around 5 p.m. Between the hours of 12:00 a.m. and 5:00 a.m., there are the fewest motorbike accidents.

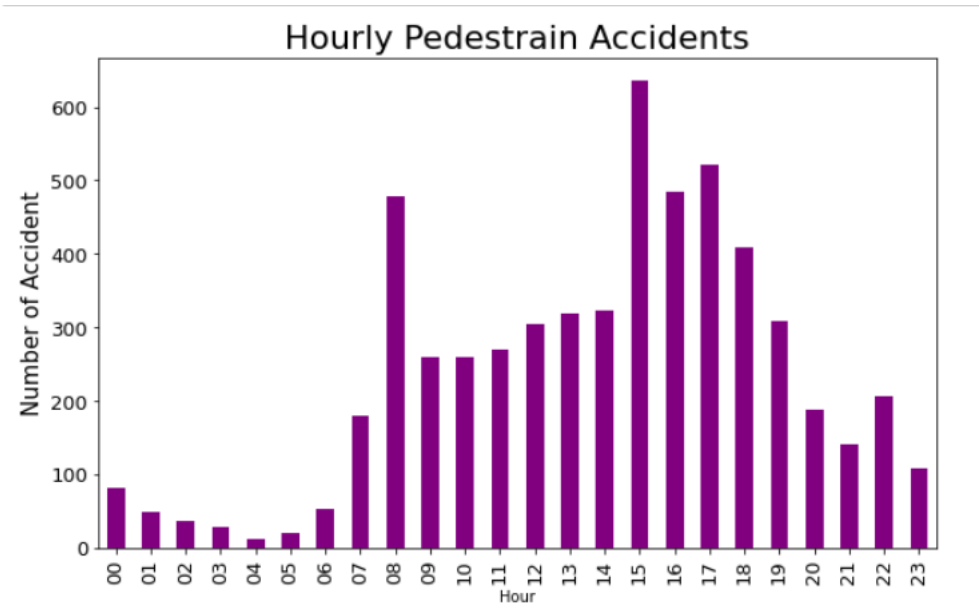
MOTORBIKES – DAYS OF THE WEEK



Just like the general accident, motorbike accidents tend to be more on Saturdays and the least on Monday.

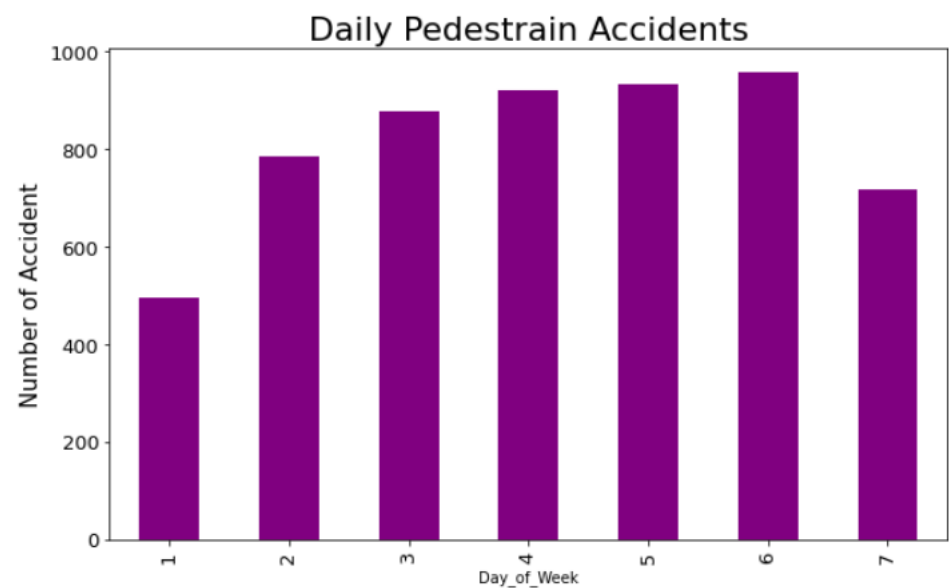
C. For pedestrians involved in accidents, are there significant hours of the day, and days of the week, on which they are more likely to be involved?

PEDESTRIANS – HOURS OF THE DAY



According to the graph, walking by 8 a.m. and in the evening hours is particularly perilous, since more pedestrian accidents occur at these times, possibly due to increased traffic as people attempt to get home.

PEDESTRIANS – HOURS OF THE DAY



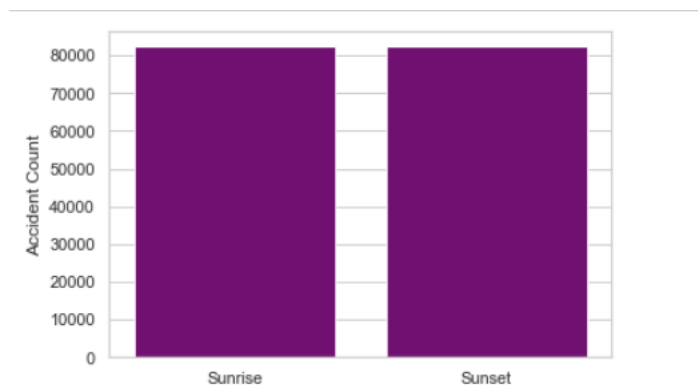
Just like the general accident, pedestrian accidents tend to be more on Saturdays and the least on Monday.

D. What impact, if any, does daylight savings have on road traffic accidents in the week after it starts and stops?

TIME	NUMBER OF ACCIDENTS	DIFFERENCE
Week before Start of DST	1916	293
Week after DST	1623	
Week before End of DST	1584	25
Week after After end of DST	1559	

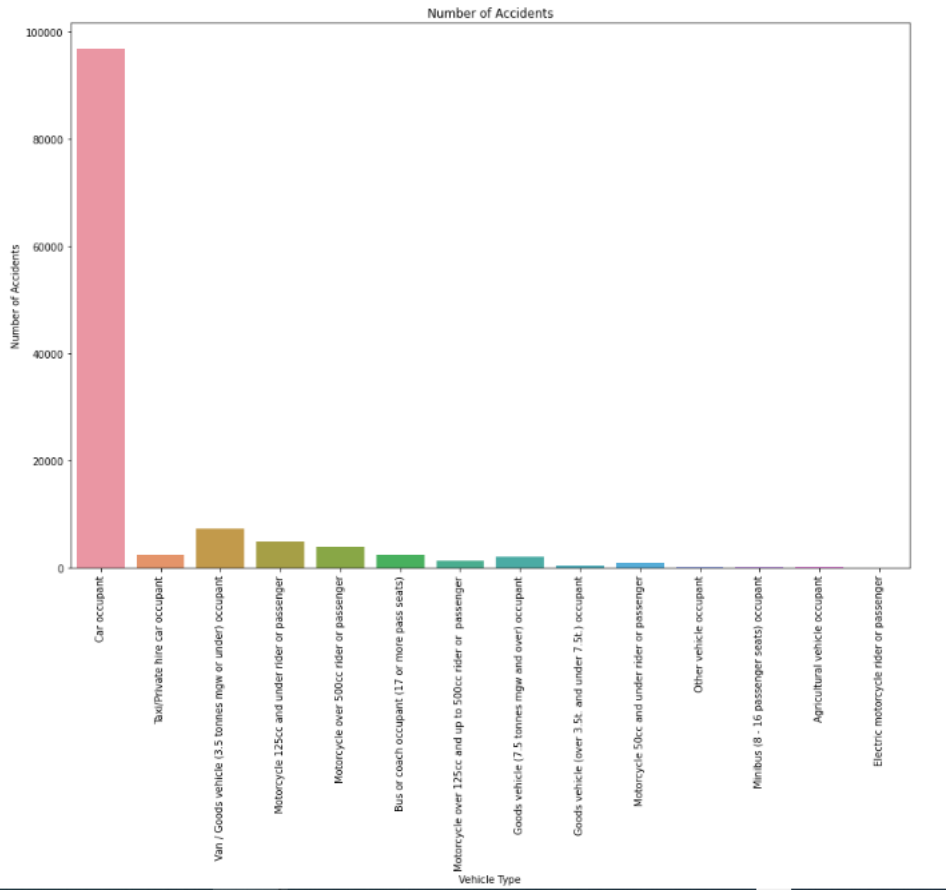
Changes to and from DST may have an impact on the frequency of the days after the shift because of sleep loss(Lahti et al., 2010). In 2019, we looked at the number of road accidents a week before and after the switch. Changes in DST did not increase rather it decreased the number of road accidents, according to our findings.

E. What impact, if any, do sunrise and sunset times have on road traffic accidents?



The light sparkles straight into the eyes of drivers at the sunrise, causing glare drivers. At sunset, drivers' eyes battle to acclimate to the shifting degree of brightness, making it harder to identify potential hazards(Gomes-Franco et al., 2020). The report does not show any impact of sunset and sunrise on road accidents in 2019.

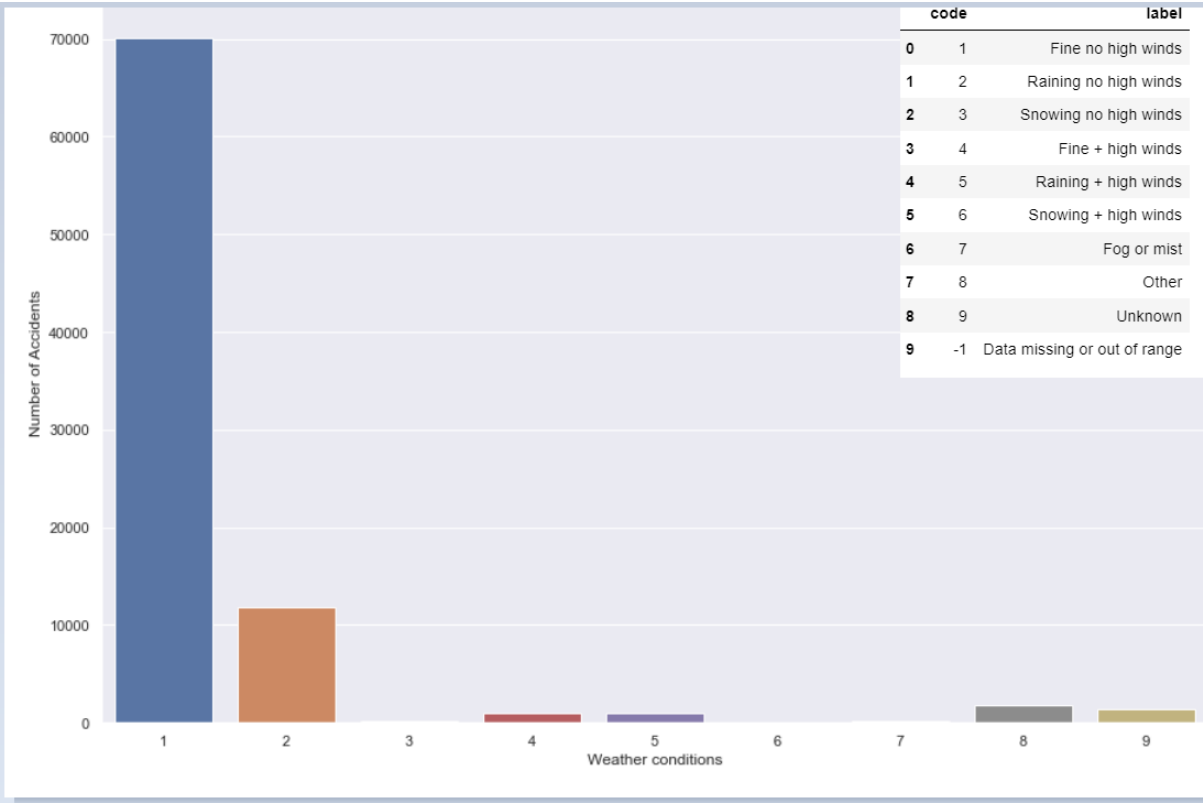
F. Are there particular types of vehicles (engine capacity, age of vehicle, etc.) that are more frequently involved in road traffic accidents?



From the graph, cars are the major culprits of road traffic accidents consisting of more than 90% of all road traffic accidents. The modal age of vehicles involved in an accident is 3 years while the cars with Engine capacity(cc) 1598 tend to be more prone to accidents.

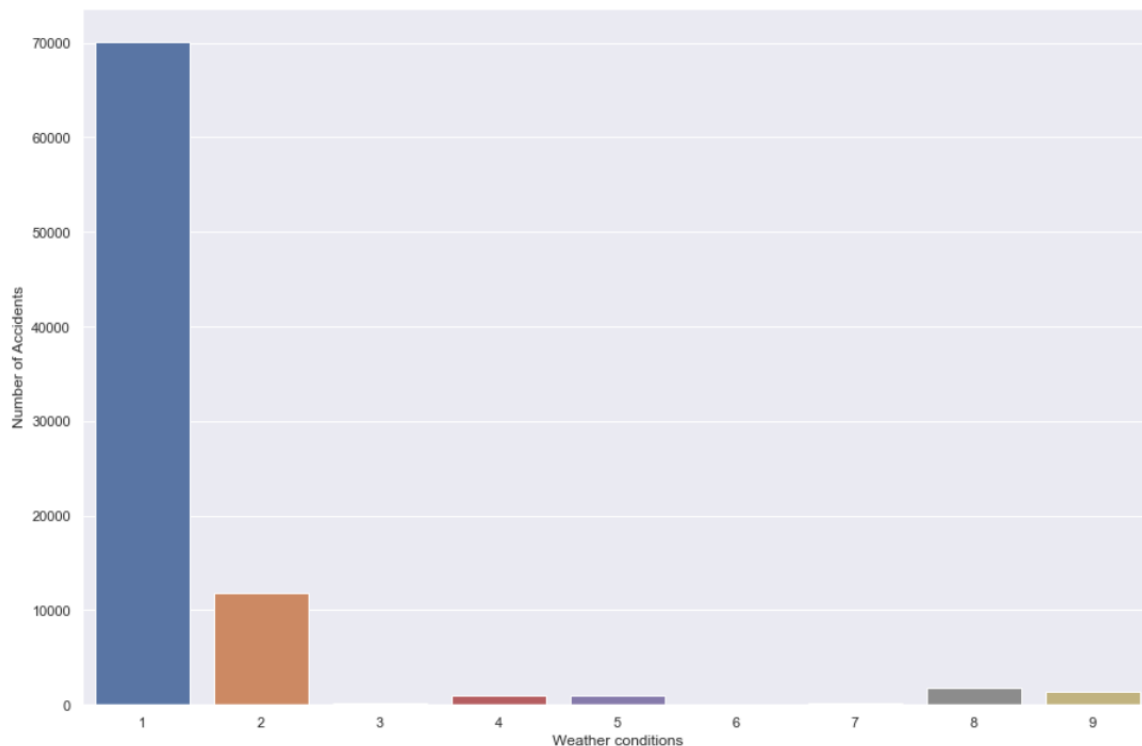
G. Are there particular conditions (weather, geographic location, situations) that generate more road traffic accidents?

Impact of Weather on Road Traffic Accidents



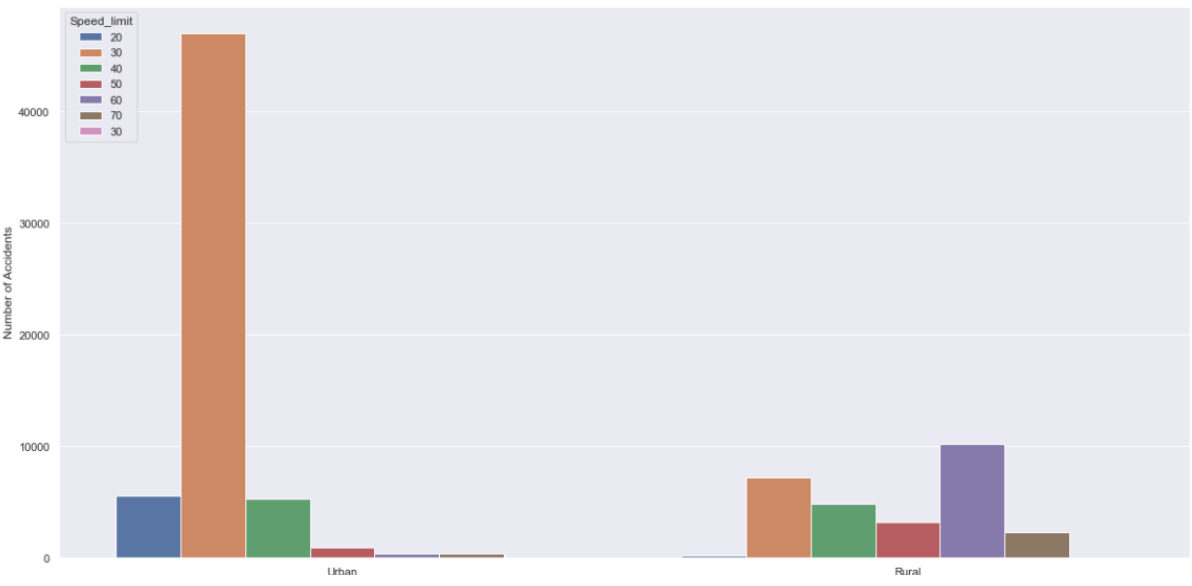
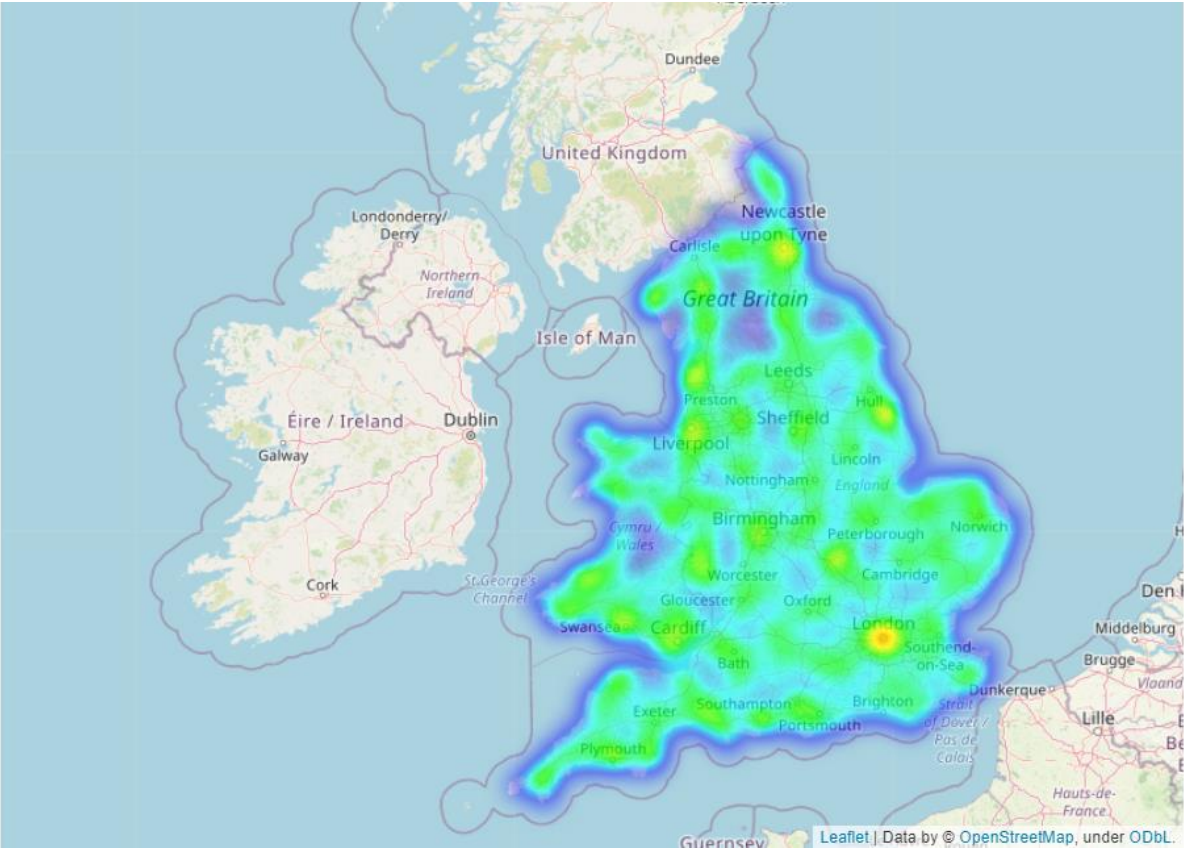
Most accidents happened while the weather was fine. During the rain, a significant number of accidents occurred. There were very few incidents while it was snowing, fog and mist. It is worthy of note that more cars tend to be on the road when the weather is fine than during any other weather conditions.

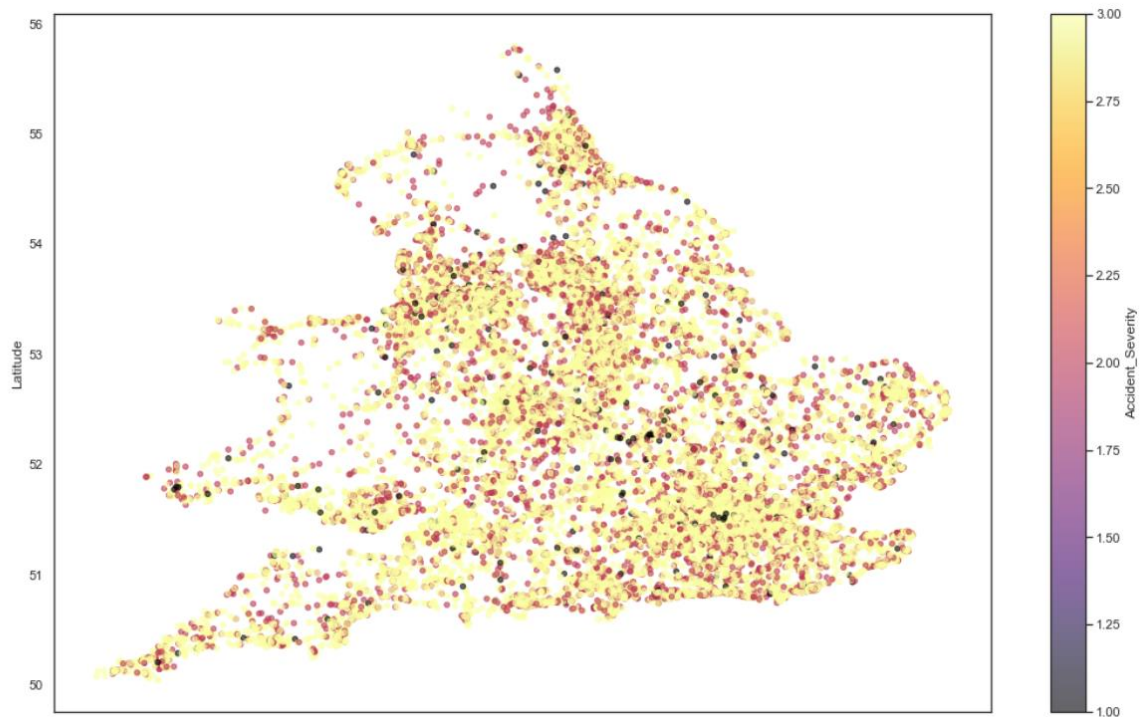
Impact of Light Conditions on Road Traffic Accidents



More accidents tend to happen during daylight and a significant amount happened during darkness. Just because there are more accidents during the day does not indicate that driving at night is safer. According to the graph, more than 35% of all accidents occurred in the dark. This is a startling result given that there is 60percent lesser traffic. That implies your chances of being in an accident are substantially higher at night.

IMPACT OF SPEED LIMIT AND GEOGRAPHIC LOCATIONS ON ACCIDENT

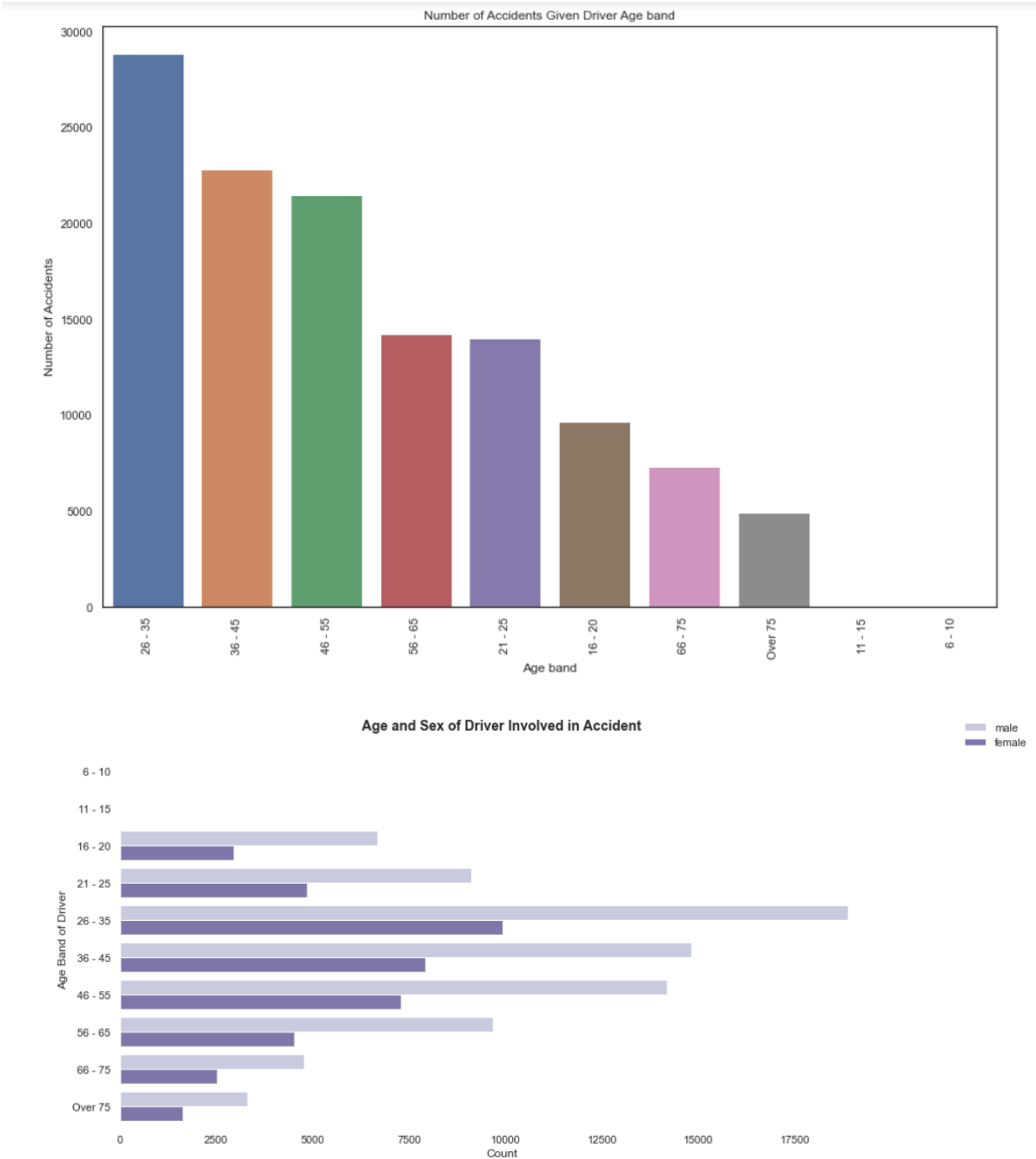




The graph and maps shows that more accidents happened in the urban cities than in rural areas but interestingly, it reveals that the severity of most of these accidents is slight as the speed is just 30mph whereas most accidents in the rural areas are fatal with a speed limit above 30mph.

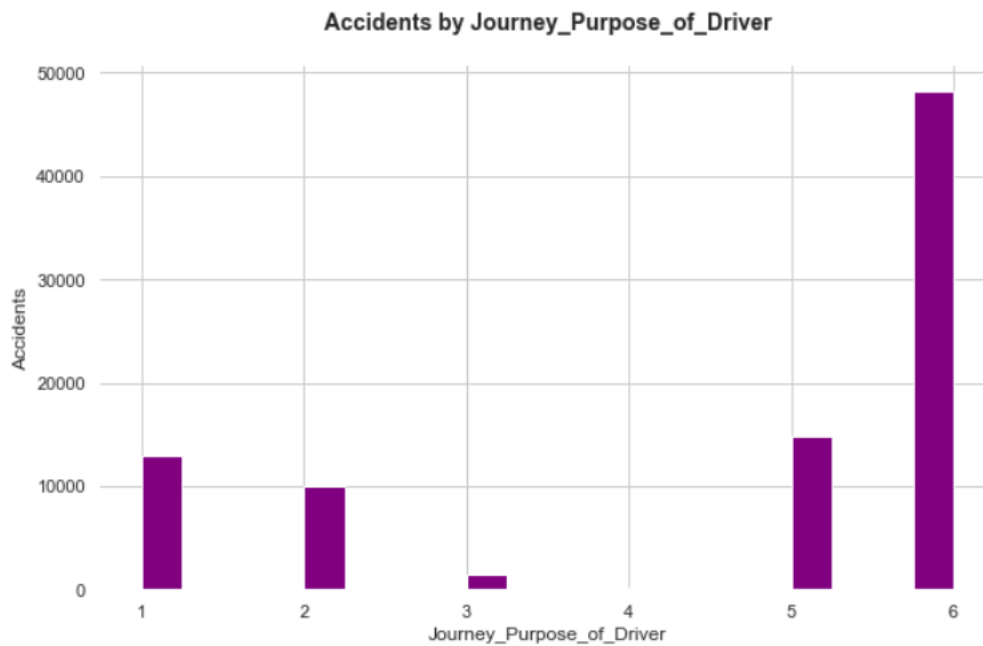
H. How does driver related variables affect the outcome (e.g., age of the driver, and the purpose of the journey)?

AGE OF THE DRIVER



For a long time, it has been known that younger and older drivers are more likely to be involved in an accident(Gomes-Franco et al., 2020). In the report, the age band of 26-35 are prone to accidents while ages over 75 are least involved in an accident. Also, more males of different age bands are involved.

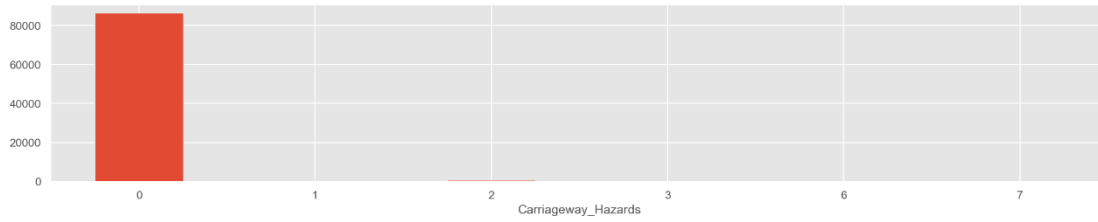
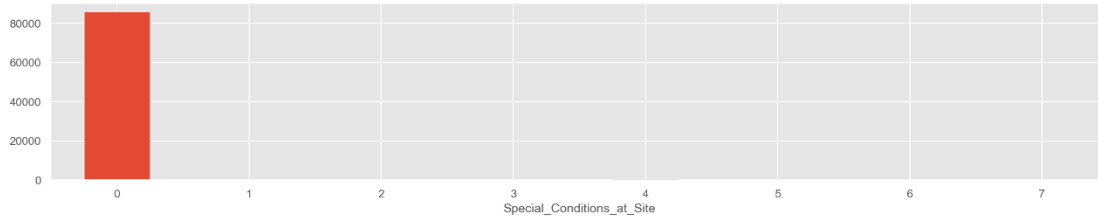
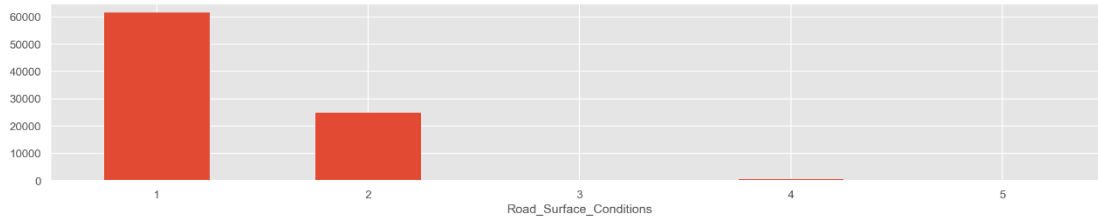
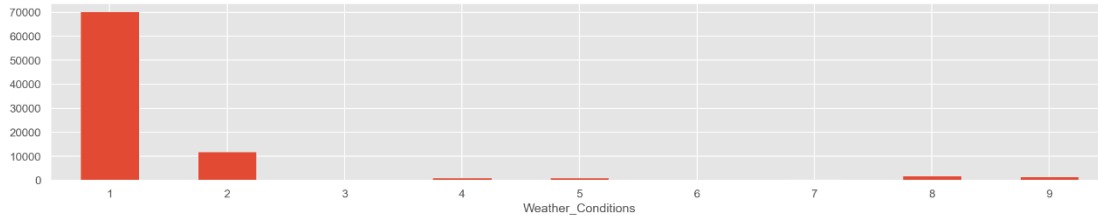
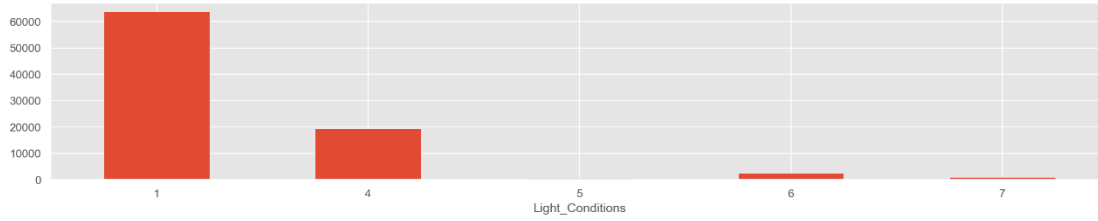
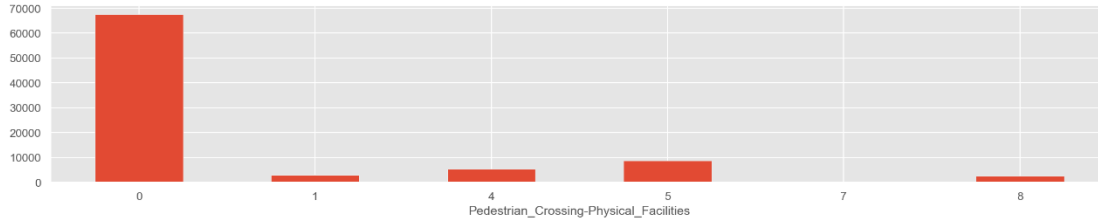
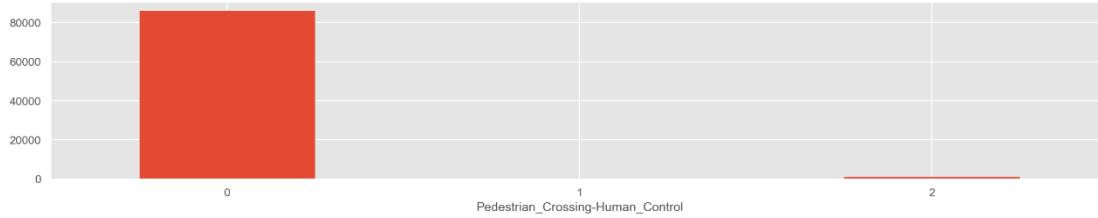
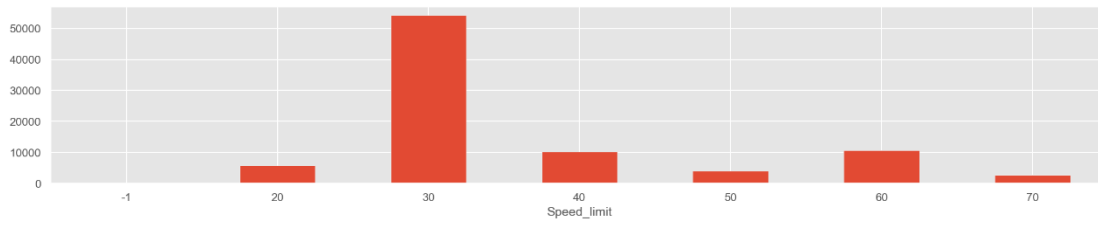
PURPOSE OF JOURNEY

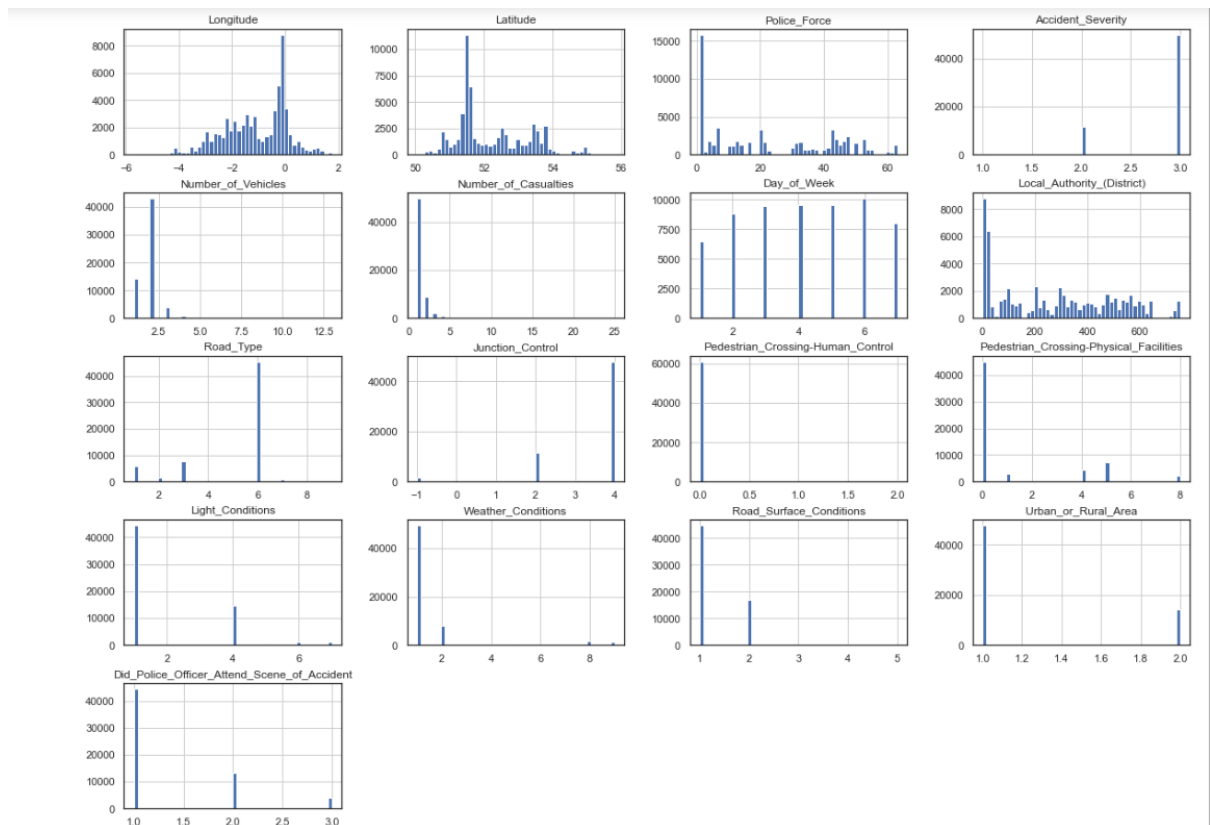
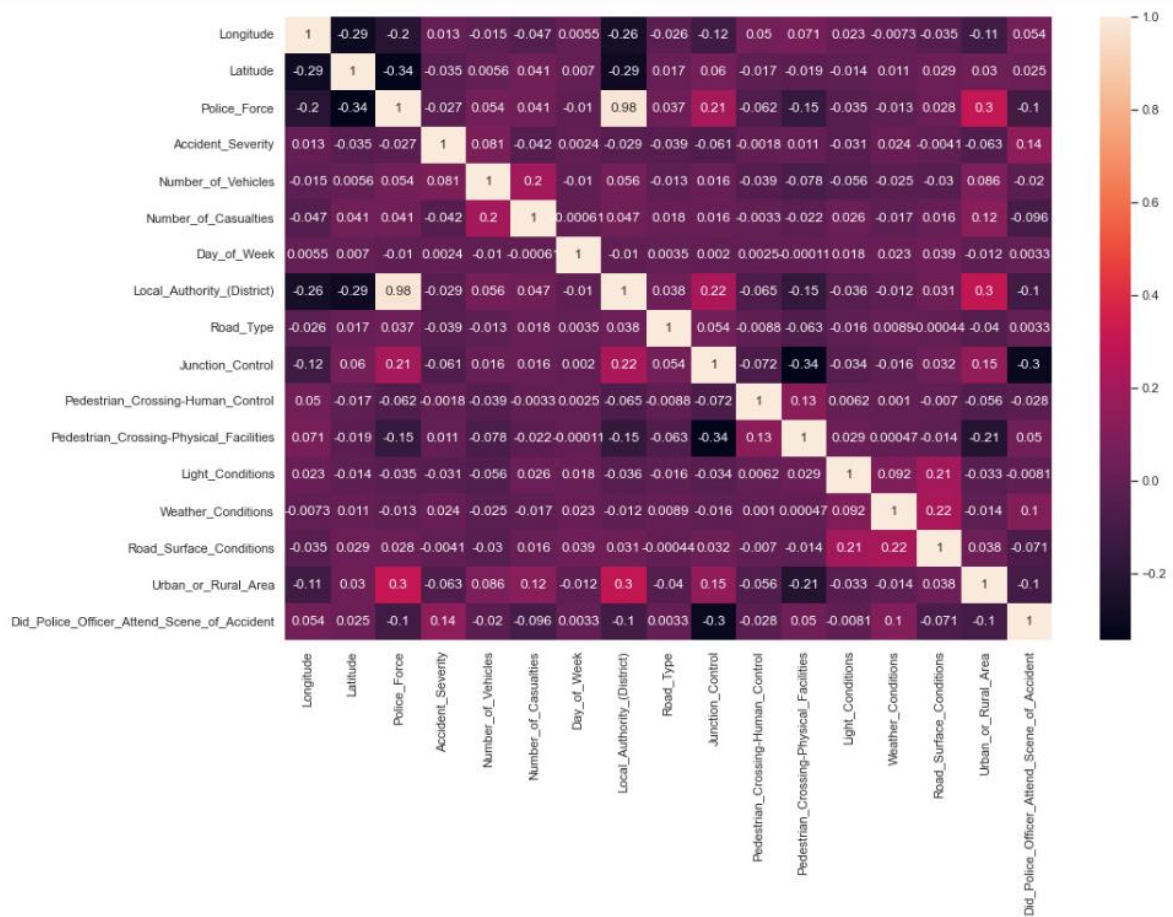


According to the data, the purpose of the drivers' journey in most of the incidents is unknown, but for all known purposes, the two most involved in an accident are 'commuting to/from work' and 'journey as part of work.'

I. Can we make predictions about when and where accidents will occur, and the severity of the injuries sustained from the data supplied to improve road safety? How well do our models compare to government models?

From analysis and as seen in the graphs below, most accidents happen when speed limit = 30, Ped Cross - Human = None within 50 metres, Ped Cross - Physical = No physical crossing facilities within 50 metres, light condition = daylight, weather condition = Fine no high winds, road surface condition = dry . The Heatmap below shows some correlations between attributes.



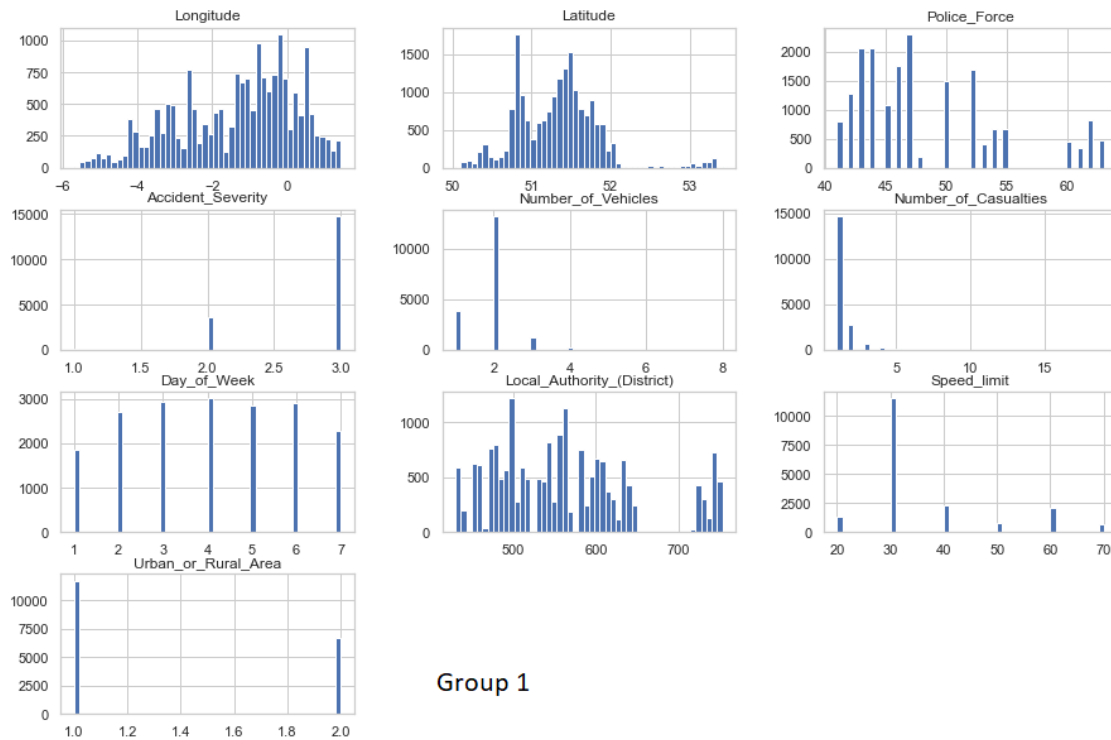


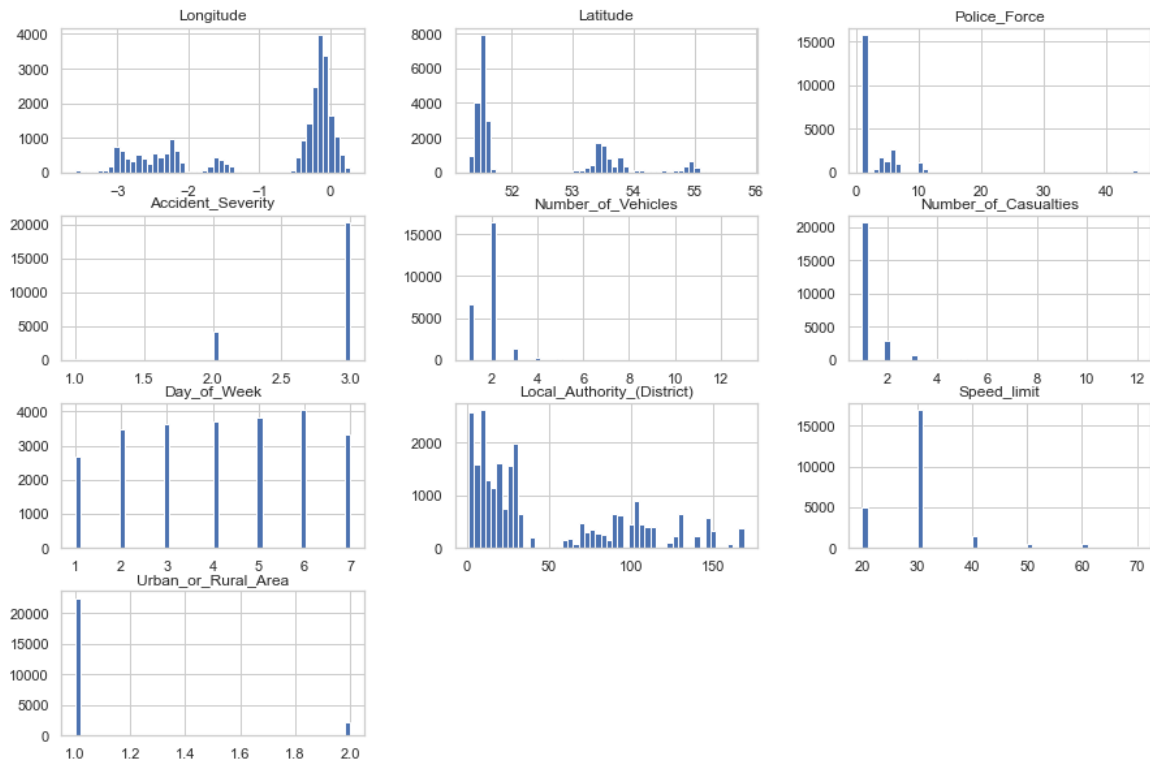
PREDICTION

A prediction algorithm is a type of ML approach that analyses historical and current data to predict future behaviour.

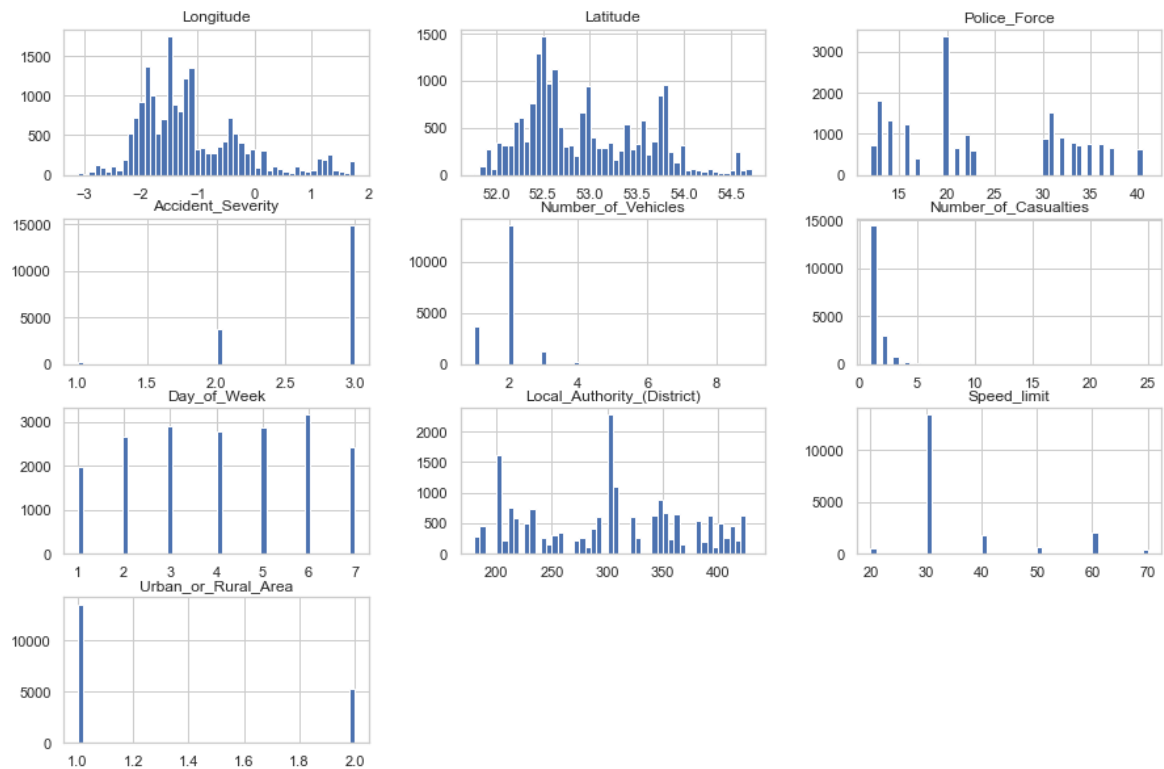
K-Means clustering

Heterogeneity in the data might lead to incorrect modelling and prediction(Kyriakopoulou & Kalamboukis, 2008). We are utilizing clustering to help us solve this challenge. In the database, we used cluster analysis to check the closest data to each other, allowing us to detect a pattern in the causes of the recorded accidents and perform trend analysis.





GROUP 2



GROUP 3

Apriori Algorithm

We used the Apriori algorithm to generate useful rules from the data. For usage in the algorithm, the attributes in the dataset are inserted into an array. The amount of rules created is controlled by three parameters: support, confidence, and lift. The minimum support and confidence are set to 0.1 and 0.5, respectively, while the minimum lift and maximum lift are set to 2. We can examine how the values of various attributes are related and the rule for the occurrence of an accident in a specific circumstance. As a result, it becomes a complement to the K-mean. The table below summarizes some of the rules.

Rule: Saturday -> weekend -> Urban	Support: 0.10001935546307945 Confidence: 0.7731920199501248 Lift: 4.330294385986254
Rule: Saturday, 0 -> High -> weekend	Support: 0.10148714474660472 Confidence: 0.7859519291587602 Lift: 4.301929012639248
Rule: High, weekend -> Saturday -> 0	Support: 0.10148714474660472 Confidence: 0.5485136430999913 Lift: 4.301929012639248
Rule: Saturday, 0 -> High -> 1 -> weekend	Support: 0.100212910093874 Confidence: 0.7859582542694497 Lift: 4.298883092033289
Rule: High, 1, weekend -> Saturday -> 0	Support: 0.100212910093874 Confidence: 0.5481252756947508 Lift: 4.298883092033289
Rule: Saturday -> High -> 0 -> weekend	Support: 0.10148714474660472 Confidence: 0.7845386533665836

	Confidence: 0.5557812914053528 Lift: 4.296425000567216
Rule: Saturday -> High -> weekend	Support: 0.1028097680570341 Confidence: 0.7947630922693267 Lift: 4.295503634775845
Rule: High, weekend -> Saturday	Support: 0.1028097680570341 Confidence: 0.5556621044372766 Lift: 4.295503634775845

Association rules with high support

End of document ■

Preventive measures can be adopted, guided by the rules, to make faster choices and limit the likelihood of accidents as against the government model. For example, the rules (weather condition = Raining no high winds, light Condition = Darkness - lights lit, weekday = Saturday) Accident Severity = High suggest that fatal accident is more likely to occur on rainy night when the road lights are turned on. These associations are powerful rules that tell us where accidents will occur and how severe they will be.

KNN and Naïve Bayes Classifiers

Finally, we use predictive models (K-Nearest Neighbors and Naive Bayes) to train our database to forecast if a particular condition would result in a serious accident. The Bayes theorem is used to create a probabilistic classifier called Naive Bayes(Budiawan et al., 2019). It presupposes that variables are unrelated to one another. In KNN new data point are assigned a value based on how close it fits the features in the training set(Lv et al., 2009). We can see that the Naive Bayes method did better in categorization, and we can also infer that it is obviously possible to forecast the severity of an accident based on the scenario even with 79.41 percent accuracy.

Recommendation

It is vital for governments throughout the world to understand the key causes of road traffic accidents and the circumstances in which they occur to take action to reduce the number of fatalities because of these occurrences. The following guidelines can reduce the likelihood and severity of road accidents.

Control speed limit to keep it at a safe level.

Implement existing Laws and Regulations

Increase public awareness and education on road safety.

Accelerate emergency response and care

Design smarter Roads

References

Alkheder, S., Taamneh, M. & Taamneh, S. (2017) Severity prediction of traffic accident using an artificial neural network. *Journal of Forecasting*, 36 (1), 100-108.

Budiawan, W., Saptadi, S., Tjioe, C. & Phommachak, T. (2019) Traffic accident severity prediction using naive bayes algorithm-a case study of semarang toll road. *IOP Conference Series: Materials Science and Engineering*. IOP Publishing.

Gomes-Franco, K., Rivera-Izquierdo, M., Martín-delosReyes, L. M., Jiménez-Mejías, E. & Martínez-Ruiz, V. (2020) Explaining the association between driver's age and the risk of causing a road crash through mediation analysis. *International Journal of Environmental Research and Public Health*, 17 (23), 9041.

Kyriakopoulou, A. & Kalamboukis, T. (2008) Combining clustering with classification for spam detection in social bookmarking systems. *ECML PKDD*.

Lahti, T., Nysten, E., Haukka, J., Sulander, P. & Partonen, T. (2010) Daylight saving time transitions and road traffic accidents. *Journal of Environmental and Public Health*, 2010 .

Lv, Y., Tang, S. & Zhao, H. (2009) Real-time highway traffic accident prediction based on the k-nearest neighbor method. *2009 international conference on measuring technology and mechatronics automation*. IEEE.