OM 420 LEC B02 - Winter 2024

Project Presentation

Explanatory Data Analysis & Prediction models



Sheikh Safwan Hossain & Mohammad Shahriar Hossain

Table of contents



O1 Exploratory Data Analysis (EDA)

Investigate the dataset to uncover patterns and insights through statistical summaries and visualizations.

O2 Variables and Data Preparation

Detail the selection of relevant variables and the steps taken to clean and prepare the data for modeling.

03 Prediction Model Development

Describe the process of building various predictive models and selecting the target variable for prediction.

04 Model Comparison and Evaluation

Compare the performance of different models using evaluation metrics to determine the most effective one.



01 Exploratory Data Analysis (EDA)

Table: Counts of Vehicle Types

Vehicle Type	Count
::	::
PCycle	17437
MC50	1489
MC125	8053
MC500	2119
MC500+	5228
Taxi	4696
Car	152686
Minibus	405
Bus	3902
R.Horse	65
Agri	423
Tram	26
Van	12579
GV7.5	982
GV7.5+	3357
MScooter	250
EMCycle	65
Other	1005
MC-Unknown	415
GV-Unknown	890
Missing	309
1::	l : : l



- Breakdown of vehicle types in road accidents.
- Cars most frequently involved; followed by taxis, buses.
- Horse riders, trams least involved lower counts.
- Data informs road safety measures and interventions.

	2	Severity			
	Туре	Fatal	Serious	Slight	
	PCycle	117	3945	13375	
Severity of road	MC50	9	380	1100	
accidents by vehicle	MC125	56	2125	5872	
type.	MC500	40	679	1400	
Cars most involved in	MC500+	279	2293	2656	
fatal, serious, and slight	Taxi	26	680	3990	
injuries.	Car	1833	25692	125161	
 High-capacity motorcycles show 	Minibus	5	71	329	
increased fatal accident	Bus	71	696	3135	
rates.	R.Horse	0	22	43	
 Highlights need for 	Agri	21	137	265	
vehicle-specific safety	Tram	2	8	16	
interventions.	Van	197	2316	10066	
	GV7.5	27	182	773	
	GV7.5+	223	695	2439	
	MScooter	12	49	189	
	EMCycle	0	18	47	
_	Other	28	285	692	
	MC-Unknowr	n 5	156	254	
	GV-Unknowr	n 18	159	713	
	Missing	0	50	259	

....



This world map displays the geographic distribution of road accident locations, with a concentrated cluster of accidents depicted in a specific region. The visualization highlights spatial patterns in accident occurrences, which could be crucial for regional safety assessments and the development of location-specific traffic safety strategies.





- Detailed map of UK road accident locations.
- Dense clusters indicate accident hotspots.
- Can inform local traffic safety and planning.
- Aids in addressing highrisk areas.







The boxplot displays age distribution of road accident casualties by vehicle type, revealing the median, range, and age outliers. Cars show a wide casualty age range, while motorcycles predominantly affect younger individuals. This aids in directing safety campaigns towards the most vulnerable age groups per vehicle type.



Dark - ng lights unlit
Dark - no light
Dark - unknown light
missing

 Stacked bar chart of road casualties by light conditions.

...

- Majority of accidents in daylight, yearround.
- Notable accidents in dark with artificial lighting.
- Fewer accidents in complete darkness, a persistent issue.

Ö





Q

Bar chart shows road casualties peak in summer, particularly for cars and motorcycles, suggesting seasonal influences on travel behavior.





Chart reveals more males are casualties in road accidents, especially ages 20 to 44, highlighting a key demographic for safety interventions.



Line graph indicates a seasonal trend in road casualties, with higher numbers in summer and lower in winter, reflecting potential changes in driving habits and conditions.

Some interesting findings!







Prevalence in Passenger Vehicles

Cars are the leading contributors to road accidents and casualties, indicating a significant area of focus for traffic safety improvements and driver education programs.

Illumination and Incidents

Daylight sees the highest occurrence of road accidents, suggesting that despite good visibility, other factors such as increased traffic volume during daylight hours may influence accident rates.

Grading the Gravitas

The majority of accident outcomes are classified as slight, with a smaller yet noteworthy proportion resulting in serious injuries, and fatal outcomes being the least common yet most critical concern.



02 Variables and Data Preparation

 $\bullet \bullet \bullet$

....

Data Preparation for Predictive Modeling



Ö

1 Target Variable Creation

A binary 'Severity_Target' was created to distinguish between accidents of high severity (fatal or serious) and others, facilitating a more focused analysis on severe accidents. Streamlining Data:

2 Streamlining Data

Unnecessary columns were removed to streamline datasets, ensuring that only the most relevant information is retained for model building.

3 Data Cleansing

Rows with missing values (NA) or placeholder values (-1) were removed to enhance data quality, which is crucial for the accuracy of predictive models.

4 Data Integration

The cleaned accident, vehicle, and casualty datasets were merged using 'Accident_Index' as a key, creating a comprehensive dataset for analysis.



5 Handling NA Values Post-Join

NA rows that emerged after joining datasets were omitted to maintain the integrity of the dataset.

6 Categorization of Variables

Selected columns were converted into factors to accurately represent categorical data within the modeling process.

7 Reproducibility and Sampling

A random seed was set for reproducibility, and the dataset was split into training, validation, and testing subsets to evaluate model performance.

8 Dataset Division for Model Assessment

The data was proportionately split into 70% training, 15% validation, and 15% testing sets. This separation is crucial for not only building the model but also for fine-tuning and assessing its generalizability to new, unseen data.



03 Prediction Model Development





K-Nearest Neighbors (KNN)

Applied KNN using two sets of 9 predictors from 18 total, to leverage diverse aspects of the data in predicting accident severity.

Linear Discriminant Analysis (LDA)

Conducted LDA with different combinations of 9 predictors, aiming to find the linear combination that best separates high severity accidents from others.



Decision Trees

Developed decision tree models using distinct sets of 9 predictors, providing a transparent and interpretable approach to classify the severity of road accidents.

Ö

04 Model Comparison and Evaluation

















Optimal Model Selection: Efficacy Proven by Error Rate

- Rigorous error rate comparisons guide our model selection process.
- Both LDA models outperform their KNN and CART counterparts, with LDA Model 2 showing the lowest error rate overall.
- LDA Model 1 also presents a lower first error rate compared to the corresponding KNN and CART models.
- Based on the compelling low error rates, LDA is chosen for its superior accuracy in predicting accident severity.



Thanks!

Do you have any questions?

