# Nationwide Traffic Insights: Visualizing Road Collisions Across the UK

Sahan Chowdhury

**Abstract**— This study utilises visual analytics approaches to gain insight into road collisions within the UK from 2018 to 2022. Despite the UK being amongst the safest in terms of road safety, the persistent high rate of accident and casualties requires a comprehensive analysis. This spatiotemporal analysis is conducted on a dataset collated by STATS19 with over 500,000 records of incidents attended by police officers. With the aid of time analysis through bar plots for 'year','month', 'day of the week' and 'hour' and spatial analysis through colour intensity chloropleth maps and k-means clustered scatter plots, the study aims to uncover patterns and trends amongst different variables with a focus on fatal accidents.

---  ◆  ---

## 1   PROBLEM STATEMENT

Roads serve as an essential means of transportation, enabling easy communication and travel from one point to another. While the UK is renowned for being amongst the safest in the world for road safety, every 16 minutes someone is killed or seriously injured on UK roads and road deaths are up by 10% since 2021 [1].

The UK government had injected a further £47.5 million for safer roads funding in April 2023 [2]. Visual analytics methods can be used to discern patterns and better understand the dynamics of road safety, while also contributing valuable insights towards the ongoing efforts to improve road safety across the UK.

This analysis utilises temporal and spatial patterns as well as clustering to unravel critical insights into road collisions in the UK from 2018 -2022. It aims to address the following questions:

- How has total road collisions fluctuated from 2018 to 2022?
- Are there specific months, days, or hours more prone to collisions?
- Which areas experienced the highest road collisions and is there a pattern?
- Which regions were fatal accidents most frequent?
- What specific factors contribute to road collisions?

The dataset used was obtained from government open data source [3]. It compromises of 538,462 rows of entries which includes geographical information (Local_authority_ons_district, urban_or_rural) temporal information (Year, Day, Time), collision factors (weather_type, road_surface_conditions, light_conditions, speed_limit, road_type) and accident details (accident_severity, number_of_casualties). The dataset offers extensive coverage for both spatial and temporal analysis, while also aligning with the analytical objectives.

## 2   STATE OF THE ART

Road safety is an essential aspect of transportation and societal well-being, by ensuring protection of human life, facilitating the smooth flow of traffic and even beyond that preserving the ecosystem. The UK government had always addressed the efforts to enhance safety measures, with more and more funding going into the transport sector each year. With increasing cars on the roads, it becomes imperative to understand factors contributing to accidents and devising effective strategies for a safer and sustainable future.

A recent study was conducted using the same dataset over an extensive timeframe (2009 -2019) by J. K. Choudhary, N. Rayala, A. E. Kiasari, and F. Jafari [4]. Which looked at the temporal trends, spatial trends, factors influence and injury severity. Their findings offer valuable insights guiding safety measures and enhancing the overall road safety strategy. With the use of bar plots against time and count heatmaps, it seemed to be very insightful, a similar approach will be adapted in the instance for this study. However, the study appeared to lack in depth analysis, such as comparing different severities and drawing conclusive spatial insights.

A second study which was looked at was that of C. Sinclair and S. Das [5]. This study used K-means clustering to map UK road collisions based on clustering. Although the dataset was not the same, the context of the data and research was very similar. The study used three clusters and it seemed to be effective, however it didn't use statistical justification for the chosen number of clusters. In this study k-means clustering will be applied as it has been proven to be successful, while also substantiating the number of clusters with the calculation of silhouette scores. Silhouette score measures how well separated clusters are, and this can be used to justify the number of clusters used.

The final study looked at, was that of P. Michalaki, M. Quddus, D. Pitfield, and A. Huetson [6]. This study focused on the time-series analysis of motorway collision in England. This study was relevant, as it looked at the time analysis, and compared variables such as speed and road type over time. The aim of this study was to identify which of the factors had a greater influence in motorway collisions, with an implication on main carriageway and hard shoulder roads. The study employed various line graphs against time to confirm whether a trend exists. This is relevant to the current study as correlation of factors will be looked at, and time trends will also be analyzed. Although the study was comprehensive, it was found that there were not enough

variables used to conclude a result, and the dataset was relatively small, as for this study a relatively large dataset was used, and a range of variable types will be included as part of the analysis.

The research papers provide a valuable insight, into shaping the analysis for this study.

## 3  PROPERTIES OF THE DATA

The data was collected through the STATS19 dataset and covers 5 years of data (2018-2022). STATS19 is a repository for reported road traffic collisions, a form completed by police officers who attended the scene of the accident. The dataset encompasses numerous variables and contains 538,461 rows and 36 columns. Several of the variables were encoded in accordance with STATS20 document [7].

Parameters include.
- *Accident_year*
- *Accident severity:* Encoded as, 1- fatal, 2 – severe 3- slight.
- *Number of casualties*
- *Date*
- *day of the week*: Encoded 1-7, where 1 indicated Monday and 7 - Sunday,
- *time*: time at which accident was recorded.
- *local_authority_ons_district*: corresponding district code.
- *road_type*: encoded as 1 – Roundabout, 2- One way street, 3 - Dual carriageway, 6 -Single carriageway, 7 - Slip Road, 9 – Unknown.
- *speed limit*
- *light_conditions:* encoded as 1-Daylight: streetlights present, 2 -Daylight: no street lighting, 3 - Daylight: street lighting unknown, 4 - Darkness: streetlights present and lit, 5 - Darkness: streetlights present but unlit, 6 -Darkness: no street lighting, 7 -Darkness: street lighting unknown.
- *weather_type*: Encoded as; 1- Fine without high winds, 2 - Raining without high winds, 3 -Snowing without high winds, 4 - Fine with high winds, 5 - Raining with high winds, 6 - Snowing with high winds, 7 - Fog or mist, 8- Other, 9 -Unknown.
- *Road Surface:* Encoded as, 1- Dry, 2 - Wet/Damp, 3- Snow, 4 - Frost/Ice, 5 - Flood (surface water over 3cm deep).
- *Urban_or_rural*

A geoJson file obtained from government geoportal website [8] was used alongside the dataset, for spatial analysis. It is also noteworthy that district names for Northern Ireland were missing from the dataset, resulting in exclusion of the analysis.

The dataset was then checked and addressed for missing values. Missing values in this dataset was denoted by -1 and NaN. All columns containing 'NaN' and not deemed necessary for the analysis were removed. This left 2067 cells with value '-1'. All '-1' values for 'weather_type' column

was categorized under label 'unknown' (encoded as [9]). Remaining missing entries were removed due to its insignificance, compromising 0.38% of the entire dataset. This ensured missing values do not impact the integrity of analytical outcomes. Lastly the 'date' column was converted into datetime format for temporal analysis. Figure [1] shows the data types.

```
accident_year                    int64
accident_severity                int64
number_of_vehicles               int64
number_of_casualties             int64
date                     datetime64[ns]
day_of_week                      int64
time                            object
local_authority_ons_district    object
road_type                        int64
speed_limit                      int64
light_conditions                 int64
weather_conditions               int64
road_surface_conditions          int64
urban_or_rural_area              int64
district_name                   object
hour                             int32
month                            int32
day                             object
year                             int32
dtype: object
```
*Figure [1]*

Outlier detection involved plotting a line graph of counts against time, see figure [2]. Quantiles were used to distinguish extremes. Upper quantile was set to 99% and lower quantile was set 0.05%. 19/01/2018 saw the highest counts of 504 and 03/29/2020 was the lowest 46. While the accidents count where extreme, the corresponding total casualties provide a crucial perspective. 19/01/2018 there was a total of 631 casualties and for 03/29/2020 there was 57. This indicates that number of accidents correlates with magnitude of casualties, and hence reinforcing the rationale for not designating them as outliers. The final modified data frame consisted of 536437 rows and 14 columns.
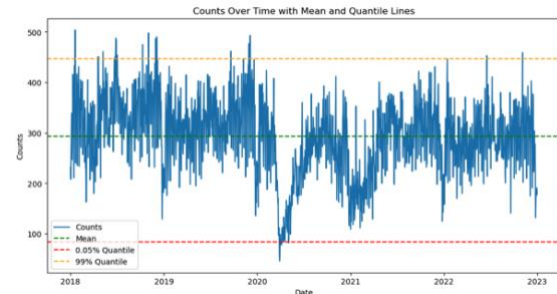

*Figure [2]*

## 4  ANALYSIS

### a.  Analysis Approach

Throughout this analysis Python will be used, due to its versatility, readability and extensive libraries making it suitable for this study. Figure [3] shows the analysis workflow plan.

Task 1:
The initial stage of analysis is data pre-processing. This is a vital step as it ensures better computational accuracy and allows better human reasoning. In the instance of data pre-processing human reasoning helps select the relevant columns and computational methods identify and processes missing values

Task 2:

Preliminary exploration of data will be carried out using histograms. This will aid in understanding the distribution of key variables.

Task 3:
Temporal analysis will use aggregate counts. To plot against time as well as counts. The date column will be extracted using computational methods and a new column for 'day' and 'month' will be created, which will be plotted against number of counts. Pearsons correlation co-efficient will be used to plot heatmaps to visualise correlation between time and day. Human judgement will then be required to interpret trends and identify potential causes for observed patterns.

Task 2:
Spatial analysis involves plotting aggregate counts on a colour intensity choropleth map for each district using geoJson file. The geoJson file will correspond district names to the 'local_aurhority_ons_district'. K-means clustering will group district based upon counts and geographical location will be performed, this will be validated using silhouette scores, where human judgement will be required to compare different numbers of clusters. Furthermore, these clusters will be filtered for fatal accidents (accident_severity = 1). To determine which region has the highest fatal accidents. Clusters percentage counts will be computed to obtain which clusters have the highest fatal accidents. Human judgement will be needed to interpret the choropleth map and identify the spatial distribution.

Task 4:
Building upon temporal analysis and spatial analysis. All variables are analysed using line plots, alongside accident severity. Human judgement is required to interpret which variables and at what conditions fatal accidents occur more often.
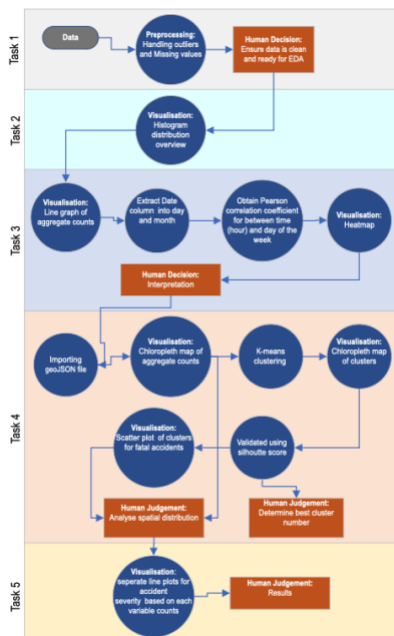


*Figure [3]*

## b. Analysis Process

Figure [4] displays a line graph of total collision counts every year. The trend reveals that number of vehicle collisions in the UK has fallen since 2018. However, this decrease only lasted till 2020, where since it's been on the rise. 2018 in total had 122,635 counts, compared to 2022 of 106,004. Notably between 2019 and 2020 there was a sharp decrease in total counts, this could be explained by the lockdown implemented during COVID 19, restricting movement and the necessity for the public to be on the road, in fact road traffic fell by 21% overall between 2019-2020, almost equating that off 1993 [9]. A similar figure of 22% was obtained in this analysis by calculation of percentage difference between the two dates. The net percentage change since 2018 was observed to be minus 11%.
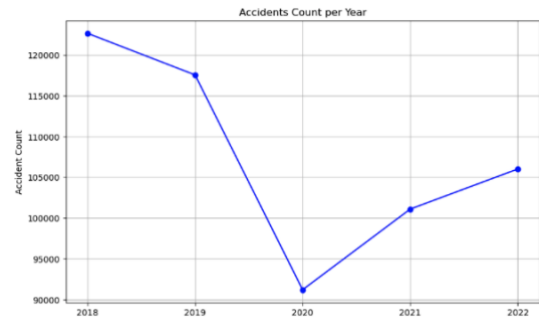


*Figure [4]*

To enhance understanding of the temporal patterns, the 'date' column was extracted into 'day' and 'months', and from the 'time' column 'hour' was extracted. The dataset already included a 'day of the week' column. The figure below displays the counts graph for 'month', 'day of the week' and 'hour'.
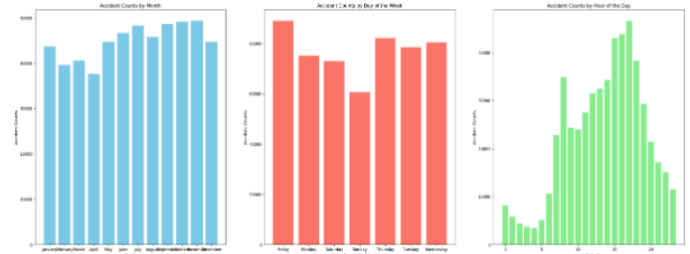


*Figure [5]*

The bars contain aggregate of the 5 years. Upon examining the 'month' graph, it's clear to see that count seemed relatively consistent. November had counts of greater than 45,000 whereas April witnessed the lowest number of vehicle collisions of 37,510 counts; however, this is very likely due to the first national lockdown being announced around this month (26/03/2020). A study which investigated the effect of a national lockdown on road collisions found the largest fall between 2019 and 2020 to be the month April, when accidents were 65% lower than in 2019 [10]. The national lockdown surely had an impact on the data for April, resulting in the lowest number of vehicle collisions in the year.

From the 'day of the week' graph it was quite clear to see that Sunday had the lowest counts of 60,350 accidents and Friday being the highest with 88,586 counts, there doesn't seem to be any clear trend other than the fact the weekends (Saturday and Sunday) displayed the lowest. Accidents are expected to be substantially higher during peak times, due to the vast traffic volume. The hourly counts graph shown in figure [5], reveals heightened accident occurrences during the early peak hour which is 8am, and the evening peak hours which is 3pm to 6pm, followed by a gradual decline. A heatmap was created to visually see the correlation between 'Hour' and 'Day of the week' as shown in figure [6].
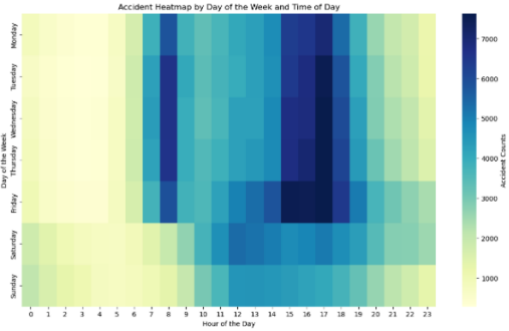


*Figure [6]*

The figure above solidifies the finding that collisions peak weekdays between 8am to 9am and 3pm to 6 pm. The heatmap was created using Pearson correlation coefficient. The lowest correlation coefficient noted was -0.97 which was identified to be Friday 5pm, having accounts of 7628 vehicle collisions. Historically 17% of all recorded accidents took place on a Friday, and 5pm time stamp contributes to 9% of these incidents [11]. On the other hand, Tuesday 3am, which had a correlation co-efficient of 0.99 had the lowest counts of 291.

Next spatial analysis was applied to observe the spatial distribution of collisions. A GeoJSon file containing UK district boundaries was used to compute choropleth maps of the UK (Northern Ireland excluded). The column '*local_authority_ons_district*' was used to obtain the corresponding district name from the GeoJson file. Subsequently, a new column was created 'District_name' containing the district name for each entry. The total counts for each district were aggregated over the 5 years of data and outputted in a bar plot (figure not shown). Birmingham, Leeds, Westminster were amongst the highest in terms of total counts, with Birmingham holding top spot with total count of 11,559. Total counts were then used to create a choropleth map using the GeoJson file as shown in figure [7].
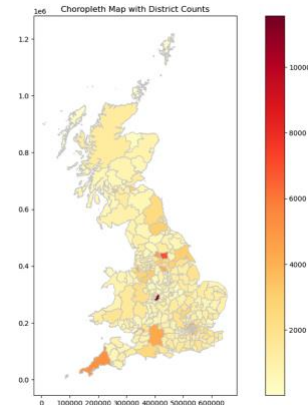


*Figure [7]*

Darker regions indicate higher counts of road collision. Birmingham, Cornwall, Wiltshire, London, and Leeds were prominently dark, while less visible districts within London such as Westminster, Tower Hamlets, Lambeth averaged more 5000 counts each. Although there was no spatial trend identified till this point, it was evident England had far greater road collisions, potentially due to England having the greater population and consequently meaning a greater number of vehicles. England also apportioned of 77% of total Roads in the UK [10]. As expected, a trend between number of accidents and number of casualties was observed, Birmingham which had the highest number of collisions also had also had highest number of casualties of 15,525, compared to Isles of Scilly which had only 5 counts of accidents and 5 casualties over the 5 years.

K-means clustering analysis was also conducted. The exploration of different cluster numbers revealed that 4-clusters yielded the highest silhouette score of 0.49. Human judgement was required to compare different clusters numbers to find the optimal number of clusters. The clustered choropleth map segmented the UK into 4 sections, Southeast Southwest, Midlands, and North.
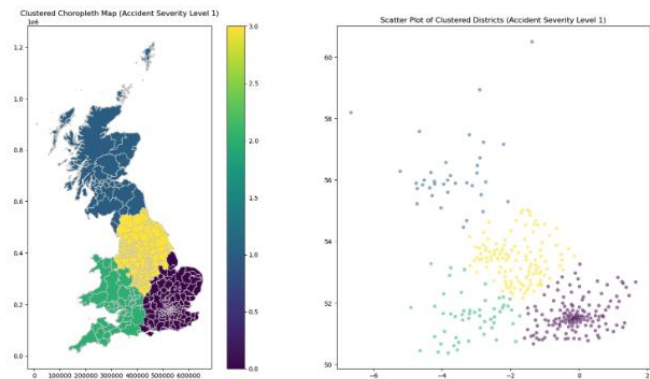


*Figure [8]*

This was then reflected on a scatter plot figure [8] to visualise in which region most of the fatal accidents occurred. A fatal collision is described to be the most dangerous form of collision and lead to a death. This was encoded as 'accident_severity = 1'. Total cluster count percentages were

calculated to pinpoint regions with the highest fatal accidents. Of the total number of fatal crashes between 2018 and 2022 which was 7796, majority accidents occurred in the Southeast (denoted by purple) which equated to 42% of all fatal accidents within the UK, followed by midlands (denoted as yellow) which equated 33% of fatal accidents. The lowest was 10% which was North of the UK (denoted as blue). This finding is also supported by a study which revealed that Southeast UK had almost 3,000 fatal reports in the last decade. [12]

Proceeding into looking into collision factors. Initially line graphs were plotted for all variables thought to exert an indirect influence on collisions. Recent study by UCL found that the situation of the environment such as poor lighting, road surface and road layout increases likelihood of a vehicle collision [13] in fact in 2022, amongst the leading causes of accidents it was found that slippery road due to weather condition and exceeding speed limit accounted for 15% of crashes [14]. This highlights the need to investigate external factors. The variables which were looked at was Road types, light conditions, weather conditions, Road surface conditions, Urban or Rural, and Speed limit. as shown in figure [9]. Initial line graphs for the three different severities (fatal, severe, slight) were plotted, however due to the inherit differences it was difficult to observe patterns. Therefore, log graphs were taken to better observe the patterns more effectively.

Initial observations reveal accident severities are very strongly correlated, an increase in one coincides with the rise of others. 'Accident severity =3' which is slight collisions exceeds in every graph and 'accident severity = 1' which is fatal collisions is the lowest. This aligns with expectations, given the substantial number of slight counts (421745 observations).

For the first subplot Road type vs Counts. There were 9 different road types. Most collisions (388744) were reported to be in road type 6 which is single carriageway. Whereas one-way streets (road type 2) and slip roads (road type 7) observed the lowest 12,646 and 8680 respectively, this disparity might be due to the specific traffic flow characteristics associated with these road types. Single carriage way has four and half times more than any other type of road. [15]

Light conditions are extremely important when it comes to driving. Lighting conditions were encoded from 1 to 7. Surprisingly when streetlight was present and lit for both daytime and darkness there were more accidents, this was also denoted by a study by the Royal society for the prevention of accidents, who found that there was no evidence of an association between reduced lighting and night-time collision. [16].

The weather condition line graph against counts saw significant fluctuations across different weather types. Most notably the weather type with the highest collisions observed was fine without high winds followed by raining without high winds. Conversely the lowest was snowing without high winds. Similarly with the road surface graph, dry Road

surface saw the highest number of accidents followed by wet damp.

For the most part vehicle collisions took place in urban regions, this is also seen in the line graph. However, for fatal accidents rural areas was higher.

The speed limit which observed the highest collisions was 30mph. For fatal accidents both 30mph and 60mph were notably significant.
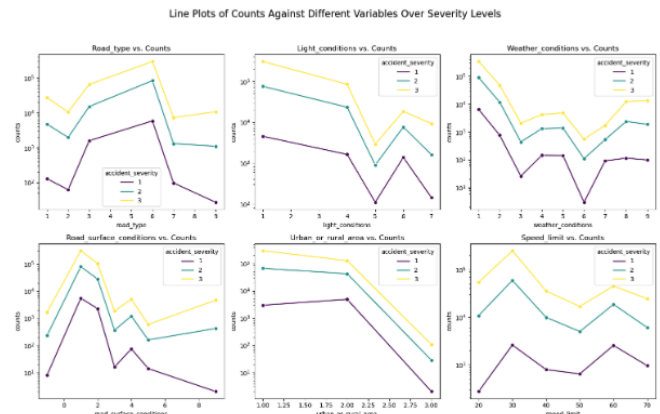


*Figure [9]*

### c.   Analysis Results

The analysis revealed a subtle understanding of UK road collisions. Accidents declined from 2018 to 2020, largely due to the COVID-19 Lockdowns, but surged afterwards. Weekday collisions peaked during rush hours. Birmingham, Leeds, Westminster, and Tower hamlets topped collision counts and this also corresponded to the number of casualties. K-means clustering identified regions with high fatal accidents, with the denoted Southeast region having highest. Looking at collision factors single carriageways sees the most collisions, especially during daytime, a cause for this may be because accidents occur more at peak times in urban regions, which are more well-lit. Urban roads also have a lower speed limits average of 30mph and this is also the speed limit at which most accidents occurred, as for rural most accidents occurred at 60mph. Weather also plays a key role on road surface conditions this is further emphasised on figure [10] showing a positive correlation between the 2 variables of 0.38, dry roads and dry weather saw the highest number of accidents , and this could be explained as majority of accidents occurred during the drier months in the UK.
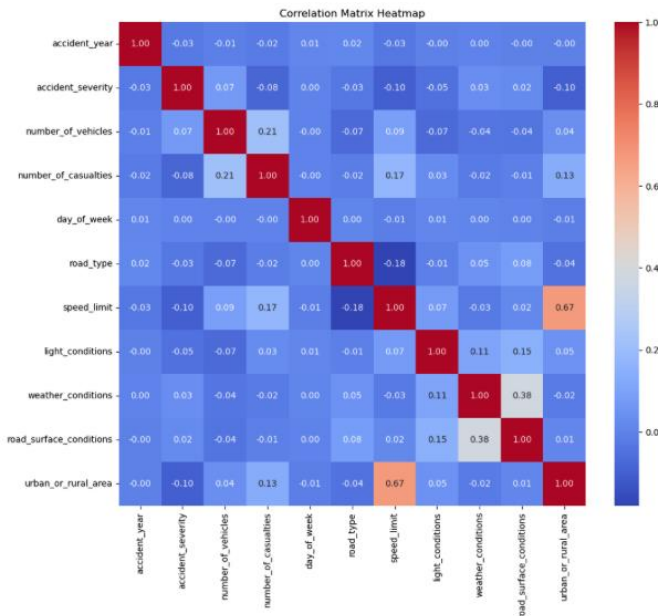
Correlation Matrix Heatmap

*Figure [10]*

## 5 CRITICAL REFLECTION

The analysis of UK road collisions from 2018 to 2022 provided valuable insights. Temporal analysis revealed time trends, with collisions initially decreasing until 2020 due to COVID-19 lockdowns, followed by a subsequent increase. However, its crucial to note the decrease started from 2018 suggesting the causation was beyond the pandemic's influence. Observing specific 'day of the week', 'hour' and 'month' patterns unveiled November's Fridays at 5pm as notable standout on an aggregate level. However, this ignores subtle differences between different days and months, therefore it is important to consider occasions like national holidays such as Christmas day and events that may have an impact on collision frequencies. Future temporal analysis may benefit to break down the years into various periods.

Spatial analysis was implemented into this study. It was found that Birmingham, Leeds, Westminster, and Tower Hamlets had the highest total aggregated counts, in the 5-year period. The use of choropleth maps visually presented the distribution of accidents using colour intensity. K-means clustering was used to distinguish regions by fatal accidents, and silhouette score was used to justify the use of the number of clusters, although exploring alternative clustering methods as suggested by the study conducted by C. Sinclair and S. Das [5], would be more effective. For future work analysing spatial distribution over time to observe how accident hotspots have evolved would provide a more comprehensive understanding. The exclusion of Northern Ireland should also be considered for.

Factor analysis explored variables influencing collisions. Line plots emphasized patterns related to road type, light conditions, weather, road surface, urban/rural classification, and speed limits and the severity of accidents. Unexpected associations, such as higher collisions in well-lit conditions,

underscored the need for critical interpretation. For future work incorporating a wider range of variables such as driver demographics, vehicle information, district population or population density would also provide some good context as to why accidents are denser in certain areas and try identifying more tailored safety measures needed to be implemented.

## 6 TABLE OF WORD COUNTS

| Problem Statement | 243 |
|---|---|
| State of the art | 481 |
| Properties of the data | 495 |
| Analysis Approach | 323 |
| Analysis Process | 1425 |
| Analysis Results | 184 |
| Critical reflection | 328 |
| Total | 3479 |

## 7 REFERENCES

[1] "UK collision and casualty statistics," Brake. https://www.brake.org.uk/get-involved/take-action/mybrake/knowledge-centre/uk-road-safety#:~:text=Every%20death%20and%20serious%20injury

[2] "Government re-launch THINK! campaign in continued drive to improve road safety," GOV.UK. https://www.gov.uk/government/news/government-re-launch-think-campaign-in-continued-drive-to-improve-road-safety#:~:text=re%2Dlaunch%20THINK (accessed Jan. 07, 2024).

[3] "Road Safety - Collisions last 5 years - data.gov.uk," www.data.gov.uk. https://www.data.gov.uk/dataset/cb7ae6f0-4be6-4935-9277-47e5ce24a11f/road-safety-data/datafile/77c26682-93e7-4e77-aa4f-5681510ca57c/preview

[4] J. K. Choudhary, N. Rayala, A. E. Kiasari, and F. Jafari, "Road Safety in Great Britain: An Exploratory Data Analysis," International Journal of Transport and Vehicle Engineering, vol. 17, no. 7, pp. 273–287, 2023, Available: https://repository.uel.ac.uk/item/8wz6y

[5] C. Sinclair and S. Das, "Traffic Accidents Analytics in UK Urban Areas using k-means Clustering for Geospatial Mapping." Available:https://ore.exeter.ac.uk/repository/bitstream/handle/10871/125145/CameraReady_PID237.pdf?sequence=1&isAllowed=y

[6] P. Michalaki, M. A. Quddus, D. Pitfield, and A. Huetson, "Exploring the factors affecting motorway accident severity in England using the generalised ordered logistic regression model," Journal of Safety Research, vol. 55, pp. 89–97, Dec. 2015, doi: https://doi.org/10.1016/j.jsr.2015.09.004.

[7] "STATS 20 Department for Transport Instructions for the Completion of Road Accident Reports," 2004. Accessed: Dec. 07, 2023. [Online]. Available: https://assets.publishing.service.gov.uk/media/60d0cd62d3bf7f4bd323e29f/stats20-2005.pdf

[8] "Open Geography Portal," geoportal.statistics.gov.uk. https://geoportal.statistics.gov.uk/datasets/196d1a072aaa4882a50be333679d4f63/explore?location=54.780722%2C-4.374937%2C9.96 (accessed Jan. 07, 2024).

[9] C. Baker, "How were road accidents and traffic affected by lockdowns in 2020?," commonslibrary.parliament.uk, Dec. 2021, Available: https://commonslibrary.parliament.uk/how-were-road-accidents-and-traffic-affected-by-lockdowns-in-2020/

[10] "Friday afternoon is peak time for RTAs | Simpson Millar," www.simpsonmillar.co.uk, Nov. 01, 2022. https://www.simpsonmillar.co.uk/media/road-traffic-accidents/friday-afternoon-is-peak-time-for-rtas/

[11] "Road lengths in Great Britain: 2021," GOV.UK. https://www.gov.uk/government/statistics/road-lengths-in-great-britain-2021/road-lengths-in-great-britain-2021#:~:text=2001%20to%202021-

[12] J. Ruppert, "What counties have the highest accident rate in the UK?," Free Car Mag, Nov. 19, 2020. https://www.freecarmag.com/what-counties-have-the-highest-accident-rate-in-the-uk#:~:text=Region%20by%20region (accessed Jan. 07, 2024).

[13] UCL, "Half of London car crashes take place in 5% of the city's junctions," UCL News, Aug. 09, 2018. https://www.ucl.ac.uk/news/2018/aug/half-london-car-crashes-take-place-5-citys-junctions#:~:text=This%20suggests%20that%20the%20environment (accessed Jan. 07, 2024).

[14] E. Yurday, "Leading Causes of Car Accidents UK 2022," www.nimblefins.co.uk, Apr. 14, 2020. https://www.nimblefins.co.uk/cheap-car-insurance/top-causes-car-accidents-uk

[15] N. Zapolski, "The Most Dangerous Roads in the UK," ChooseMyCar - Find The Best Deal on a Cheap Car Loan, Apr. 17, 2020. https://choosemycar.com/resources/choosing-the-right-car/the-most-dangerous-roads-in-the-uk#:~:text=The%20Types%20of%20Road%20with%20the%20Most%20Accident&text=This%20is%20likely%20because%20these (accessed Jan. 07, 2024).

[16] "The Royal Society for the Prevention of Accidents Street Lighting and Road Safety," 2020. Available: https://www.rospa.com/media/documents/road-safety/factsheets/street-lighting-and-road-safety.pdf