# Interactive UK Road Safety tool

## Jheronimus Academy of Data Science: Data Visualization

Camila Matoba*
TiU Student number (SNR): 2067717
Tu/e Student number: 1728938

Sietske Wijffels†
TiU Student number (SNR): 2051929
Tu/e Student number: 0940371

### ABSTRACT

The interactive and exploratory tool was envisioned to corroborate actionable data driven decisions qua government decision makers to improve road safety in the UK. This tool empowers the comprehensive amount of data collected over the years by the British government to pinpoint a district location that needs improvement to avoid casualties. The main principle of the tool is that it requires as little cognitive effort as possible, according to the guidelines of data visualisation best practices, not to the detriment of completeness. For that reason it is divided in two pages: **Road Casualties Map** and **Attributes correlation analysis**. The former is intended to give an overview of the situation and the latter allows more in depth investigation on correlation analysis between attributes.

Changes from the interim report to this final version are the addition of this abstract, sections 6 and 7. Minor changes in other sections in are highlighted in **bold** and sections 4 and 5 completely changed, to comply with the teacher's feedback and the team's progress.

## 1 INTRODUCTION

The growing information society enables professionals to learn more efficiently about their domain every day. The transformation from available data to useful insights is different for every specific field and for every goal. Within this assignment we aim to create insightful visualizations for the British government and municipalities to enhance safety on the British roads. More specifically, the main goal is to emphasize roads that can be considered as outliers when it comes to accidents in order to take data driven actions by finding the root causes. This tool is focused on analyzing and monitoring traffic accidents and highlights when specific roads are involved with accidents too frequently. Government and municipalities have the responsibility to enhance a safe traffic situation and our tool could provide them with the initial information suggesting for intervention in form of road reconstruction or safety measures. To achieve this goal, we are provided with traffic accident data form Great Britain from 2016 to 2020 [5]. The available data includes information about the accident, the vehicle and the casualty.

For this visualization tool we are mostly interested in the accident data to look for road related causes. Wegman (2017) describes the importance of data driven identification of conditions and circumstances under which crashes occur. For instance the time of day/night, alcohol involvement and road type [12]. A visualization dashboard is the perfect tool for municipalities to get better insights in the underlying causes of traffic accidents, however the data that is used for these visualizations is limited to only reported accidents and does not include any statements of people involved in the accidents. For this reason the information retrieved from the visualizations can be taken into account when considering road adjustments, but

---

*e-mail: Camila@Matoba.com.br
†e-mail: s.wijffels@tilburguniversity.edu

this information should be additive to the decision making process. This report follows the 'four nested levels of vis design' framework described by Tamara Munzner in Visualization Analysis & Design [9].

## 2 RELATED WORK

The official government authorities have published a Power BI dashboard visualizing many statistics about traffic accidents using the same dataset. The visualization tool [2], developed on PowerBI, consists of 5 pages, each of them focusing on different aspects of the accident. The first page formally informs the user about the origin of the data and the setup of the tool. The second page provides statistics of different features plotted against the sex of the casualty. The third page draws a timeline of the different accidents. The last two pages visualize the amount of accidents per region in form of a color map of Great Britain. The data is mostly aggregated, which means that you can only draw global conclusions from the visualizations. Different filter options are available to zoom in more in specified topics and look at statistics for only one specific group of interest. The above described implementation of the data serves as inspiration for the foundation of our project but is not developed for the same goal. Where the Power BI dashboard shows a map that highlights which areas are severely involved with traffic accidents, one cannot zoom in on the map on municipalities or neighborhoods. The Dashboard does not provide the option to inspect separate traffic accidents nor does it show any environmental conditions that could have been of influence for the accident. Part of our goal is to embed previously described shortcomings into our visualizations.

A tool that does show the separate accidents on a map is crashmap.co.uk [1]. When zooming in on a road, separate traffic accidents pop up in a color related to the severity of the injury. However, the free version of the tool does not give more information on the accident excepts for the date, severity and number of casualties and vehicles involved. This tool could be used out of curiosity, but will not provide you with any statistics on possible causes of the traffic accidents.

## 3 DATA ANALYSIS

### 3.1 Domain Data Specification

The road safety data descended from police forces and is also known as STATS19. For this project, only the years 2016 untill 2020 are used, being described as 'last-5-year' datasets which are three different datasets including information about the casualty, accident and vehicle. An additional dataset named 'Road-Safety-Open-Dataset-Data-Guide' holds all categorical descriptions and is used to decrypt the integer values of the 'last-5-year' datasets. The three datasets are combined using outer join and missing values are imputed with -1. Instances with data errors are deleted. The combined dataset holds information of 1213544 traffic accidents spread across the continent of Great Britain. The casualty attributes relate to the scope of our project since they can be used as filter options for the user of the application. The accident attributes include information about the external circumstances of the accident, which are the elements that will give an indication of related causes. The vehicle attributes are interesting, since the distribution of vehicle types is changing over

the years and telling something about the road security for specific vehicle types. Later in the project, a comparison between electric vehicles and all other vehicles will be made.

**Some attributes have an extensive amount of unknown or missing information. For full data transparency, it is possible to choose between 'All data' or 'Exclude' the data that is declared as 'Unknown', 'Not known', 'Other/Not known', 'Other','Undefined', 'Data missing or out of range', 'Data missing', 'Unclassified', 'Unallocated', 'unknown (self reported)', 'Unknown vehicle type (self rep only)', 'Other vehicle', -1**

**Furthermore, the dataset used [5] contains data for the years 2016 to 2020, even though there is data available for the years 1979 to 2020. A shorter and more recent timeframe is chosen to ensure faster calculations/loading times while encountering enough data on electric vehicles.**

## 3.2 Data Abstraction

The format of the data is tabular and describes all characteristics of the accidents in one row. The different datasets are connected using an outer join on 'accident index', 'accident year' and 'accident reference', which are the key attributes for the merged dataset. **The attributes used are summarized in table 1 and most of the attributes are categorical.**

| Attribute summary | | |
|---|---|---|
| Attribute name | Missing values | Attribute type |
| accident year | 0 | quantitative |
| age of driver | 0 | quantitative |
| age of vehicle | 0 | quantitative |
| age of casualty | 0 | quantitative |
| number of casualties | 0 | quantitative |
| speed limit | 71 | quantitative |
| day of week | 0 | ordinal, cyclic |
| year-month | 0 | ordinal, cyclic |
| time | 0 | quantitative, cyclic |
| age band of casualty | 445512 | categorical, ordinal |
| sex of driver | 62 | categorical |
| propulsion code | 0 | categorical |
| vehicle manoeuvre | 2375 | categorical |
| accident severity | 0 | categorical, ordinal |
| sex of casualty | 433422 | categorical |
| casualty home area type | 530780 | categorical |
| road type | 18832 | categorical |
| junction detail | 18 | categorical |
| urban or rural area | 5 | categorical |
| did police officer attend scene of accident | 5 | categorical |
| weather conditions | 59435 | categorical |
| latitude | 271 | quantitative, geometric |
| longitude | 271 | quantitative, geometric |
| special conditions at site | 0 | categorical |
| carriageway hazards | 0 | categorical |
| road surface conditions | 0 | categorical |

Table 1: Attribute summary

## 4 TASK ANALYSIS

4 tasks are intended to be accomplished with this project. After presenting them in this section, they will be referred by their numbers.

### 4.1 Domain Specific tasks

The visualization tool is thought out so that it can answer four questions and their follow ups. The extensive Statistical Release by the British Department for Transport [10] main task is to present a commented overview of the 2019 reported road casualties to the general public. However, it does not provide information detailed enough to be acted upon. It also states that the "total value of prevention of unreported injury accidents are around £17bn a year", which is a great financial motivation to explore which are the actions the government can do to improve safety. In order to have compelling material as evidence in an actionable plan in urban road planning, it is important to recognize where have the most accidents happened over the years, so that preventive maintenance can be performed. Furthermore, one of the points mentioned in it is that "There is no single underlying factor that drives road casualties. Instead, there are a number of influences. These include: (...)

1. The mix of transport modes used;
2. The mix of groups of people using the road (e.g. changes in the number of newly qualified or older drivers);
3. External effects such as the weather, which can influence behaviour (e.g. encouraging/discouraging travel, or closing roads) or change in the risk on roads (by making the road surface more slippery)".

In spite of this statement, there is no data visualization accompanying it nor any further explanation. The first three tasks derive from these aspects that the mentioned report could be complemented upon. Thus, our first and main task is:

**Task 1:** Where are the exact roads that hoarded the highest amount of casualties in each district and their characteristics?

**Task 1.1:** How do these specific road characteristics compare to those of the whole district?

To investigate what could be the underlying factors that lead them to be the worst roads, a further attribute analysis is crucial. The attributes chosen to be further investigated are: road type, speed limit, junction detail, vehicle manoeuvre, since they are more complete and informational than the others. Detailed information about the completeness of the data is present on table 1

**Task 2: Is there a correlation to the selected attributes and number of accidents?**

The last task related to the report is:
**Task 3:** To what extent do the factors in the British statistical release vehicle type, age, weather and risk on roads influence fatality?

The fourth task is related to the rise of electrical vehicles (EVs) in the United Kingdom. In the analysed period, there were 31.889 Battery Electric Vehicles (BEVs) and 54.981 Plug-in Hybrid Electric Vehicles (PHEVs) in 2016. In the final analysed year of 2020, there was a rise of a little over five times in the amount of EVs, which per category was 207.051 BEVs and 232.492 PHEVs [11]. This increase is expected to continue, since there is a global effort in tackling transport related emissions. The Prime Minister himself has announced that the sale of new petrol and diesel cars and vans by 2030 will end, with all new cars and vans being fully zero emission from 2035 [4]. For that reason, it is important to compare how the increase of EVs have influenced traffic behavior.

**Task 4:** How do the same factors explored in task 1.1 fare for EVs?

## 4.2 Task Abstraction

From the tasks defined in the previous section, the task abstraction is respectively:

**Task 1:** Task 1 can be subdivided into two separate tasks. *Locate outliers:* targets are outliers geographically located on a map. The target is known, but the location is unknown. This could be visualized using a map, showing the locations of the targets. *Present features:* visualize the details of the target in an interpretable manner. One implementation could be a hovermode on the map, introducing additional information. Another option is to provide the information in a tableform, since outliers are by definition not with high amounts compared to the total dataset of interest.

**Task 1.1:** *Compare distributions:* compare the distribution of one category with the distribution of the same category for a different dataset. This can be visualized using stacked or layered barcharts for categorical variables, but could also be visualised in a piechart for comparison with the total.

**Task 2:** *Discover correlations:* to see if a specific set of values for different attributes correlates with the target variable, a parallel coordinate plot can be used. When the attributes are categorical, a parallel category plot is preferred.

**Task 3:** task 3 can be subdivided into two separate tasks. *present distribution:* presenting the distribution of attributes for the extremes of the target attribute. This can be achieved using a histogram for quantitative attributes and bar charts for categorical attributes. *discover correlation:* see if there are specific values for attributes that are of influence for the target attribute. This can be achieved using a bar chart for categorical attributes or a line chart for quantitative attributes.

**Task 4:** *Compare distributions:* compare the distribution of one category with the distribution of the same category for a different dataset. This can be visualized using stacked or layered barcharts for categorical variables, but could also be visualised in a piechart for comparison with the total.

## 5 THE INTERACTIVE UK ROAD SAFETY TOOL

The main goal of this project is making data visualisations that can lead to precise data driven decisions. On the Road Casualties Map, first there is an overview of the full UK following the rule of thumb "Overview First, Zoom and Filter, Detail on Demand" [9]. Another reason for this is because plotting all the precise locations of casualties made the map very cluttered, losing a clear view of the bigger picture. This also greatly impacted the responsiveness of the visualisation, since less dots needed to be plotted. As an interactive tool, it is possible to choose map aggregation between:

- **General:** Overview of in which districts the most road casualties happen in the UK. By clicking on a district you can see a table underneath with the top 10 addresses with the most casualties. By clicking on a district, it is possible to enter the Precise mode.

- **Precise:** Pinpoint the exact roads where most accidents happen in the chosen district from the Dropdown menu. Additionally, hovering over a point allows more information about it, such as the sum of number of casualties, the mean of the speed limit, and the mode for road type, junction detail, urban or rural area, and vehicle manoeuvre. This fulfills **Task 1**. Lastly, underneath a table presents the top 10 streets where the most casualties happened and their amount.

for **Task 1.1** there is an accompanying attribute chart of the total number of casualties that is synched with the district chosen for the map. The idiom choice is used according to the attribute type by Munzer.

For **Quantitative attributes**, such as accident year, age of driver, age of vehicle, age of casualty, number of vehicles, number of casualties, and speed limit the idiom chosen was line charts, as length is one of the most effective magnitude channels for quantitative values. The marks used for this are lines, because they make it easier to find trends, for example. The x-axis is the value of the attribute. For ordinal attributes, such as day of week, year month , time (rounded to the hour), age band of driver, and age band of casualty, in a similar manner as the quantitative attributes, the x-axis presents the attribute in the correct order, with it's respective name.

For **Categorical**, as most of the attributes on the dataset are, they can be subdivided for easier understating as:

- The driver and their vehicle: sex of driver, vehicle type, propulsion code, ev, vehicle manoeuvre, car passenger, accident severity

- The casualty: casualty severity, casualty class, casualty type, sex of casualty, casualty home area type

- The road: road type, first road class, second road class, junction location, junction control, junction detail, urban or rural area, did police officer attend scene of accident, special conditions at site, road surface conditions, carriageway hazards, weather conditions

The idiom chosen for these are horizontal bar charts. Bars are used instead of lines because of the expressiveness principle, and they are horizontal for legibility of the lengthy categories. They have been ordered by quantity for easier comparison between categories as well as to be easier browse for the highest in casualties for that category. Since all attributes sum up to the total amount of casualties, in case there are less than six categories, the idiom chosen is a pie chart for faster part-to-whole judgement. In this case the marks used are separate colored areas with the channels color for categorical attribute; and angle for quantitative attribute. A more concrete example of how this visualization can be used to complete task 1 is presented on the section 6.

The tasks 3 and 4 can also be accomplished by this chart by utilising the filters that work also for the map. For the checkboxes, if no option is selected, the filter is the same as selecting all categories. Each of the attributes mentioned by the report of **Task 3** can be individually selected on the attributes drop down menu. The attributes are casualty class (driver, passenger, etc), age, weather and risk on roads. To analyse only fatalities, Accident severity level is selected as 'Fatal' and Slight, Serious can be deselected. These categories are according to the Severity adjustment figure guidance 2020 [3].

Each of the elements mentioned on task 3 were unexpectedly debunked. The assumptions seemed reasonable enough, but the data proves otherwise. By filtering out only fatalities and excluding missing data for the full UK dataset:

1. Transport modes used: vehicle type is 71% cars, indicating very little mix of transports.
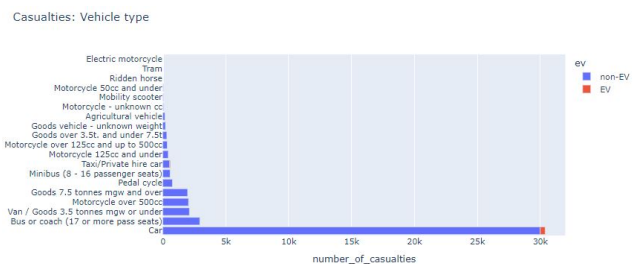


Figure 1: Sum of casualties per vehicle type

2. The highest amount of casualties was between 26-45 years old (37% of all cases). This is counter intuitive to the belief that young drivers are the main age band for accidents.
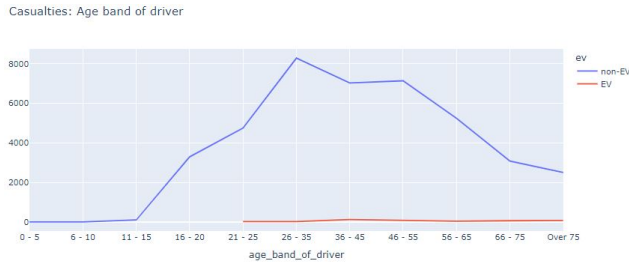


Figure 2: Sum of casualties per age band

3. External effects: Weather was fine with no winds for 85% of fatalities, including a Dry road surface for 67% of the cases. Furthermore, There were no special conditions at site nor carriageway hazards for over 97% of the cases.
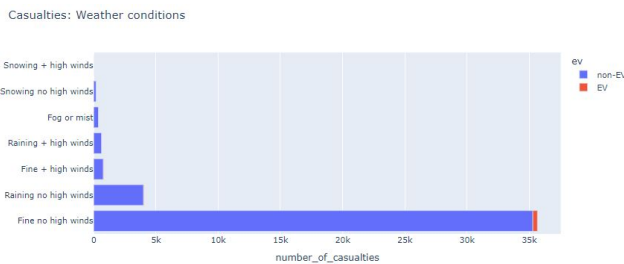


Figure 3: Sum of casualties per Weather conditions

For **Task 4** there is the checkbox between EV and non-EV vehicles. EV vehicles are the ones with propulsion code in 'Electric', 'Hybrid electric', 'Electric diesel', 'New fuel technology'. However, the EV presence in the UK market is still very small, reaching (1,71%) of the full dataset. That might be the reason why some strange trends appeared in comparison with the non-EV vehicles. For example, when looking up at which hours do most accidents happen, non-EVs have all levels of severity around rush hour (18:00).



Figure 4: Sum of casualties per rounded time slot

However, for electric vehicles, the time where most fatalities occur is much later, at 22:00. Since the sample of EVs is so small it is difficult to assure that this is a trend.

Lastly, the categorical bar charts are stacked to highlight the differences between the different propulsions. But since their quantities are highly different, EVs can be better visualised when filtered on their own.
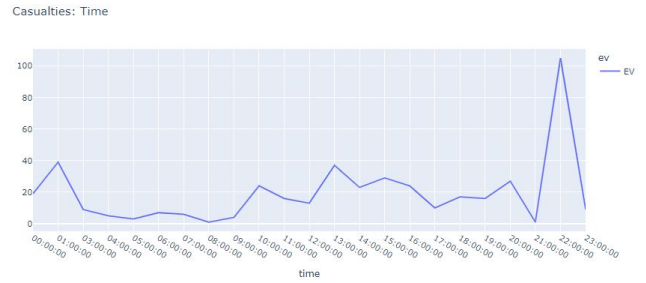


Figure 5: Sum of casualties per rounded time slot for EVs

For **Task 2**, a parallel category plot and heat map are shown on the Attributes correlation analysis page of the tool. These visualizations reveal some of the characteristics of traffic accidents and enable the user to discover relationships between Vehicle Manoeuvre, Road Type, Speed Limit and Junction Location. The parallel category plot provides a view on the relations between the attributes and shows which combinations occur most frequently. For each local authority district, a custom parallel coordinate plot and heat map are generated. It provides municipalities with useful information about what specific combination of circumstances result in the most accidents within their district. This information can be considered when constructing new roads or could function as a guide to reconstruct old roads. The parallel category plot constructed with plotly express does not yet have an extensive functionality for hovering, which required the use of an external legend. To provide the user with a different perspective on the same data, the heat map is added to the page. The advantage of the heat map in comparison with the parallel category plot is that it shows the labels of the categories, which makes it easier to interpret directly. The disadvantage, on the other hand, is that only two attributes can be compared at the same time.

### 5.1 Implementation

The code and how to run it is available at GitHub: `https://github.com/octokami/uk_road_safety`. The app was fully implemented in Python with Plotly and Dash [7]. The main challenge was that neither of the group members worked with dash before, so that the beginning was very slow. Dash is, however, an extremely practical and powerful tool to use for data visualisation. Other python packages used were:

- **geopandas [8]**: To find the centroids for the Genral view

- **geopy.geocoders's Nominatim [6]**: to get addresses from the coordinates on the table

Other functions needed for the functionality were developed by the team members and they are in the beginning of the Jupyter Notebook.

## 6 USE CASES

Our main intended end user are the people in charge of planning road maintenance in the UK. This person would click on their city (or select from the dropdown box) and locate where the most dangerous roads are located, as well as the list of top 10 most dangerous roads. In the example in Appendix A for Epping Forest there is a clear discrepancy in 'Linkside'. The amount of fatalities is 5 times the number of the second road. Exploring further, even though in the district casualties are usually not at junctions as it can be seen on the bar chart, this crossroad is definitely an outlier. This information about the point can be seen by hovering, along with other data such as speed of 30mph and that most drivers had accidents there when moving off (calculated by mode).

Another use case would be detecting the city's accident 'rush hour', to implement measures to incentive drives to change hours of driving, for example. In London, this hours are very well defined as 9:00 and 18:00, when the average of the UK starts at 16:00 (from figure 4), indicating that this peak could be more spread and hopefully lowered.
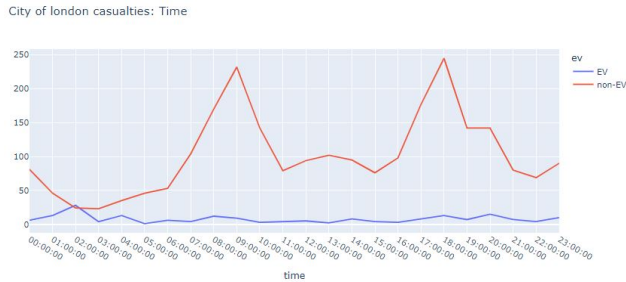


Figure 6: Hours with most accidents in London

Finally, a simple use case is to enjoy, since the tool allows a lot of interaction for exploration. Some interesting trends are for example that more accident happen in urban areas (60%), but 75% of all fatalities are rural. Another intriguing trend is the steep reduce in casualties during the beginning of COVID-19.
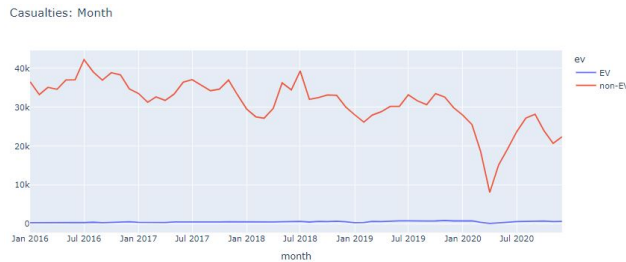


Figure 7: Steep reduce in casualties during the beginning of COVID-19

## 7 CONCLUSION AND FUTURE WORK

This report shows the progress of a first version of a supportive system for road (re)construction development. The tool envisions the main outliers in traffic accidents and provides complementary information about the conditions of these accidents. These conditions can be compared to the general attribute trends within the local authority district, with a separate functionality of filtering on electric vehicles and severity of the accident. Additionally, combinations of circumstances can be compared to find traffic situations with a higher risk for accidents per local authority district. Multiple views are provided, in form of a zoom map with hovering details, charts and tables.

The fourth task about EVs was completed, but as there is still very little amount of them in the UK, some districts even have zero casualties for EVs. This does not imply that EV drivers there have had zero casualties, but could be that there are actually none or very little vehicles of this category. Some interesting trends could be found as in figure 5, but as the sample size is small this could be only considered as a possible indication of what might happen when more EVs are in circulation and not as a definite fact.

In general we can conclude that we are content with the achievements we made considering our initial project goal. The development of the tool was highly challenging considering the new

software, data cleaning and data aggregation, but in the end remains a feeling of excitement when navigating through your own tool.

### 7.1 Future work

For future work it would be interesting to develop a better integration of the two dashboard pages. It is likely that the user is only interested in one municipality, which makes it recommendable to keep the location settings of the map page the same for the attribute page. Moreover, The heat map serves as a different view on the information envisioned in the parallel category plot, but would become redundant when the category labels could be envisioned while hovering over the parallel category plot. When more time would be provided, it would be interesting to develop the integration between the accidents per address and the combination of traffic situation specific attributes. If the top 10 most dangerous traffic situations all share the same sort of attributes, one can conclude that the underlying cause lies (partly) within the combination of those factors.

Some disadvantages of the current implementation of the tool could not be solved within the given time span. One difficulty is the loading speed of the charts on the Road Casualties Map page. Due to the calculation of the extra hover information, updating the graphs takes more time. This was a trade-off, where we prioritized the hover information over the loading speed of the charts. Furthermore, extra hover information could not be implemented for the parallel category plot, which made it obligatory to add a separate legend to the graph. Unfortunately this took the position of the heat map that need to divert to the bottom of the page, which makes the Attribute page a scrollable page. It would be more preferable to have both plots on the same screen to give a better overview and provide the possibility to compare both graphs.

To give the tool better grounds, it would also be highly valuable to talk to experts from the field of traffic safety. Within the scope of this project, we limited ourselves to literature, however, the problem of road safety is a very pragmatic one and might be different for any municipality/district. Talking to experts would give a different perspective on the problem and might reveal other factors that could be interesting to envision within the tool.

### REFERENCES

[1] crashmap. *CrashMap Data: Great Britain 1999 - 2021 (verified) - 2021 Provisional data to June*, 2022.

[2] Department for Transport. *Reported road casualty statistics in Great Britain: interactive dashboard*, 2017.

[3] Department for Transport. *Severity adjustment figure guidance 2020*, 2020.

[4] Department for Transport. *Outcome and response to ending the sale of new petrol, diesel and hybrid cars and vans*, July 2021.

[5] Department for Transport. *Road Safety Data [Data set]*, February 2022.

[6] S. Hoffman. Nominatim.

[7] P. T. Inc. Collaborative data science, 2015.

[8] K. Jordahl, J. V. den Bossche, M. Fleischmann, J. Wasserman, J. McBride, J. Gerard, J. Tratner, M. Perry, A. G. Badaracco, C. Farmer, G. A. Hjelle, A. D. Snow, M. Cochran, S. Gillies, L. Culbertson, M. Bartos, N. Eubank, maxalbert, A. Bilogur, S. Rey, C. Ren, D. Arribas-Bel, L. Wasser, L. J. Wolf, M. Journois, J. Wilson, A. Greenhall, C. Holdgraf, Filipe, and F. Leblanc. geopandas/geopandas: v0.8.1, July 2020. doi: 10.5281/zenodo.3946761

[9] T. Munzner. *Visualization analysis and design*. CRC press, 2014.

[10] A. Murphy. Reported road casualties in great britain: 2019 annual report. Department for Transport Statistical Release, 2020.

[11] RAC. *The road to electric - in charts and data*, 2022.

[12] F. Wegman. The future of road safety: A worldwide perspective. *IATSS research*, 40(2):66–71, 2017.
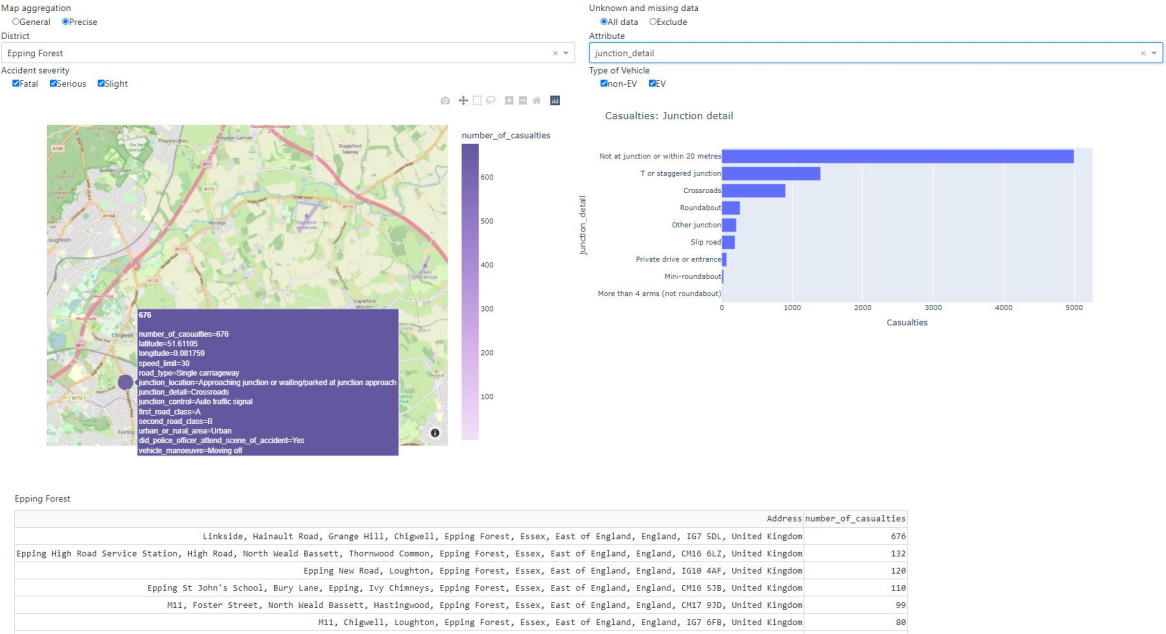
## 8 APPENDIX

## A ROAD CASUALTIES MAP



Figure 8: Use Case for Epping Forest

## B ATTRIBUTES CORRELATION ANALYSIS
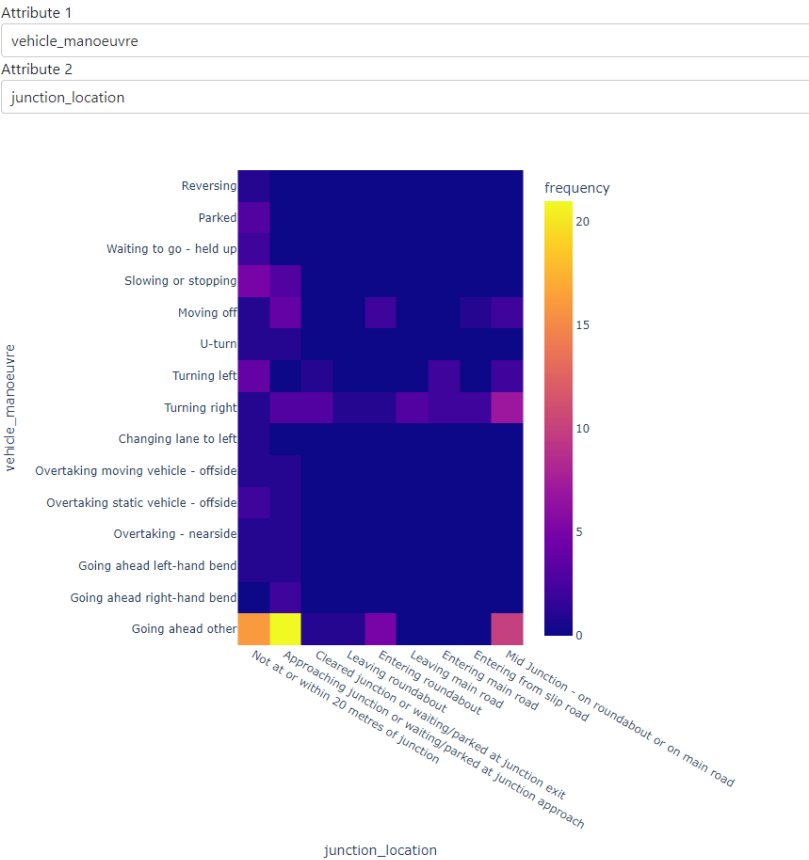


Figure 9: Parallel category plot

## C  HEAT MAP

Attribute 1

vehicle_manoeuvre

Attribute 2

junction_location



Figure 10: Heat map