

Vehicle Road Crash Analysis in Victoria

Contents

Data checking and cleaning	2
First Irregularities	2
Second Irregularities	2
Third Irregularities	3
Data exploration	4
Question no 1	4
Question no 2	5

Data checking and cleaning

After loading the dataset into Tableau, it will look like the picture below:

Accident No	Accident Date	Accident Time	Day Of Week	Day Week Desc	Accident
T20120000009	1/01/2012	30/12/1899 2:25:00 am	1	Sunday	
T20120000012	1/01/2012	30/12/1899 2:00:00 am	1	Sunday	
T20120000013	1/01/2012	30/12/1899 3:35:00 am	1	Sunday	
T20120000018	1/01/2012	30/12/1899 5:15:00 am	1	Sunday	
T20120000021	1/01/2012	30/12/1899 7:30:00 am	1	Sunday	
T20120000028	1/01/2012	30/12/1899 4:00:00 am	1	Sunday	
T20120000032	1/01/2012	30/12/1899 12:55:00 am	1	Sunday	
T20120000043	1/01/2012	30/12/1899 12:45:00 am	1	Sunday	
T20120000044	1/01/2012	30/12/1899 4:25:00 pm	1	Sunday	
T20120000046	1/01/2012	30/12/1899 4:25:00 pm	1	Sunday	
T20120000050	1/01/2012	30/12/1899 3:00:00 pm	1	Sunday	
T20120000053	1/01/2012	30/12/1899 4:50:00 pm	1	Sunday	
T20120000054	1/01/2012	30/12/1899 6:20:00 pm	1	Sunday	
T20120000056	1/01/2012	30/12/1899 4:15:00 pm	1	Sunday	
T20120000058	1/01/2012	30/12/1899 6:40:00 am	1	Sunday	

Figure 1 Dataset Overview

First Irregularities

The initial irregularities I encountered pertained to duplicate values for the accident number (ACCIDENT_NO). To solve this issue, I decided to remove the redundant rows using Python. This approach was chosen as duplicate entries offer no additional value to the dataset and may skew analysis results if left unattended.

Before	After																																				
<p>Duplicate Rows</p> <table> <tr> <th>Accident No</th><th>Count of Acci..</th></tr> <tr> <td>T20170014004</td><td>2</td></tr> <tr><td>T20240002923</td><td>1</td></tr> <tr><td>T20240002413</td><td>1</td></tr> <tr><td>T20240001940</td><td>1</td></tr> <tr><td>T20240001799</td><td>1</td></tr> <tr><td>T20240001324</td><td>1</td></tr> <tr><td>T20240001125</td><td>1</td></tr> <tr><td>T20240000743</td><td>1</td></tr> <tr><td>T20240000535</td><td>1</td></tr> <tr><td>T20240000350</td><td>1</td></tr> <tr><td>T20240000296</td><td>1</td></tr> <tr><td>T20230030840</td><td>1</td></tr> <tr><td>T20230030781</td><td>1</td></tr> <tr><td>T20230030758</td><td>1</td></tr> <tr><td>T20230030418</td><td>1</td></tr> </table>	Accident No	Count of Acci..	T20170014004	2	T20240002923	1	T20240002413	1	T20240001940	1	T20240001799	1	T20240001324	1	T20240001125	1	T20240000743	1	T20240000535	1	T20240000350	1	T20240000296	1	T20230030840	1	T20230030781	1	T20230030758	1	T20230030418	1	<p>After Removing Duplicate Rows</p> <table> <tr> <th>Accident No</th><th>Count..</th></tr> <tr> <td>T20170014004</td><td>1</td></tr> </table>	Accident No	Count..	T20170014004	1
Accident No	Count of Acci..																																				
T20170014004	2																																				
T20240002923	1																																				
T20240002413	1																																				
T20240001940	1																																				
T20240001799	1																																				
T20240001324	1																																				
T20240001125	1																																				
T20240000743	1																																				
T20240000535	1																																				
T20240000350	1																																				
T20240000296	1																																				
T20230030840	1																																				
T20230030781	1																																				
T20230030758	1																																				
T20230030418	1																																				
Accident No	Count..																																				
T20170014004	1																																				

Figure 2 Duplicate Values (Before and After)

Second Irregularities

I identified an issue related to the SPEED_ZONE column, which had an outlier. Outliers can have a significant impact on statistical analysis and model performance. To address this, I utilized Python to impute the median value for the outlier. By replacing the outlier with the median, we ensure that we have a more representative measure of central tendency and preserve the overall distribution of the data (Horsch, 2021).

The image's left side depicts the dataset before preprocessing, where duplicate values, outliers, and inconsistent formatting are evident. The right side shows the cleaned dataset after addressing these issues, resulting in a more organised and reliable dataset for analysis.

Data exploration

Question no 1

Compare and contrast when and what type of accidents occur over time (2012-2024). Consider different time periods: hour of day, day of week, month of year, year in the dataset. What do you notice about the pattern over time? What do you know about that time period in order to provide some explanation of increases, decreases, and similar numbers of accidents?

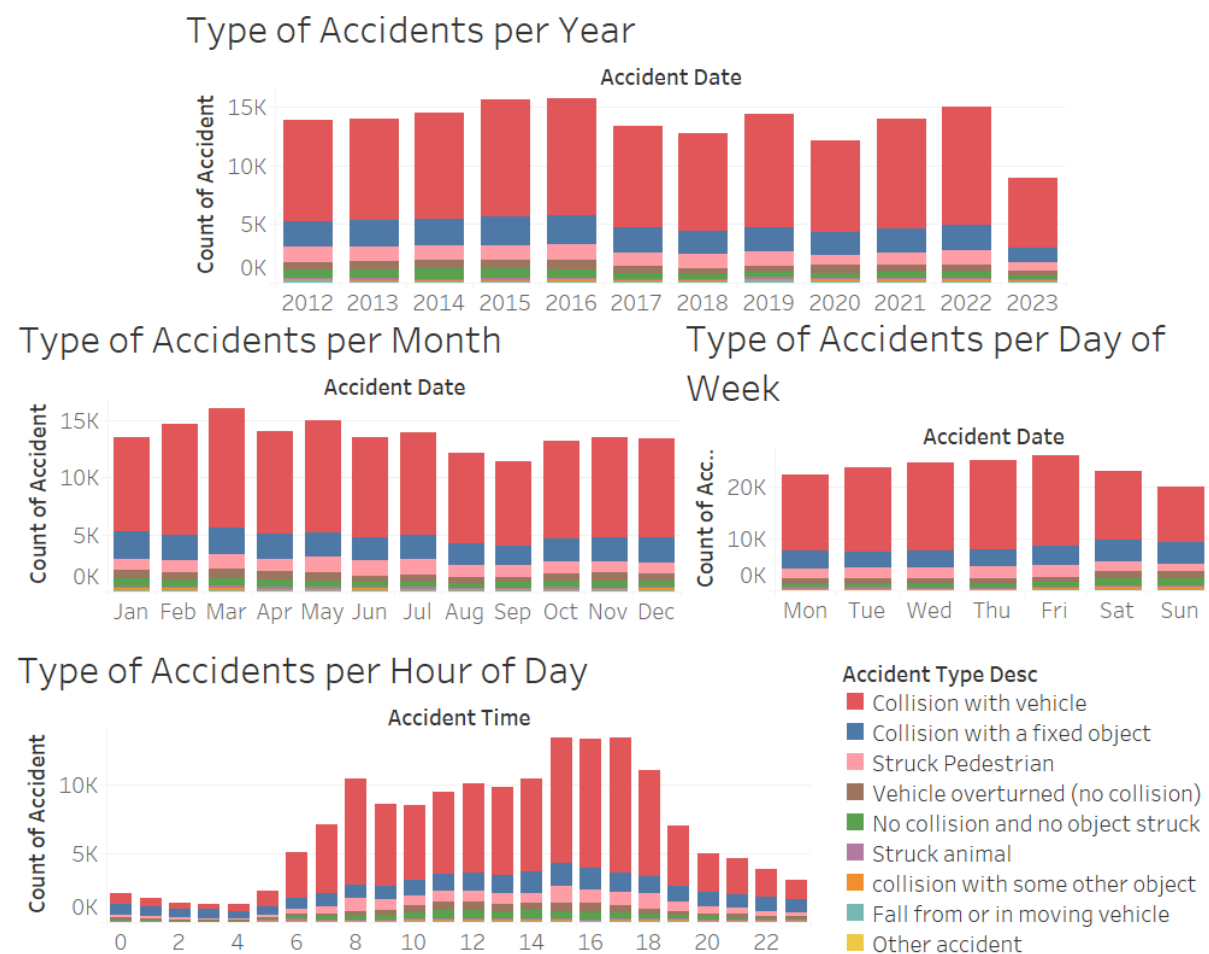


Figure 5 Type of Accidents Over Time (2022-2023)

By analysing data from 2012 to 2024, there were consistent patterns across different timeframes: year, month, day, and hour of the day. Notably, collisions with vehicles and fixed objects are the most frequent accident types throughout this period.

Regarding the yearly distribution, from 2012 to 2014, the number of accidents remained relatively stable. However, a noticeable increase occurred between 2015 and 2016, reaching its peak. Subsequently, there was a decline in the number of accidents, followed by another increase in 2019 and a subsequent rise again in 2022. Regarding the monthly distribution of accidents, March and May stand out as the months with the highest occurrences of accidents, while other months remained stable. These months may coincide with seasonal changes, increased travel activities, or other environmental factors contributing to elevated accident rates.

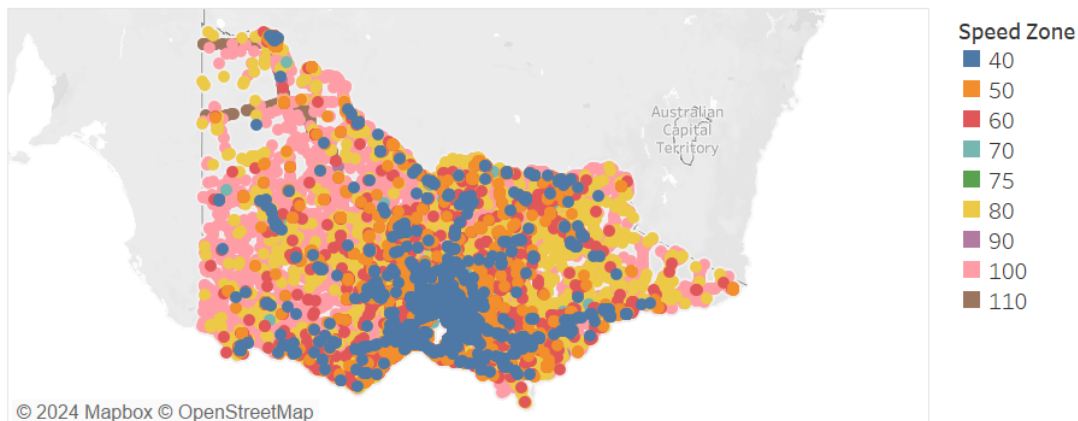
Regarding the day of the week, it reveals a tendency for accidents to be more prevalent during weekdays compared to weekends. This observation aligns with expectations, as weekdays typically experience higher

levels of vehicular traffic due to work commutes and other weekday activities. Furthermore, examining accidents by the hour of the day highlights two prominent peaks: one in the morning around 8:00 am and another in the late afternoon, between 3:00 pm and 6:00 pm. These peaks likely correspond to rush hour periods.

Question no 2

Compare and contrast where accidents occur geographically. Consider in particular the geometry, speed zone of the road and the urban/rural aspect of the location. How does the spatial mapping of the data and additional information about the road help support, challenge, or change your conclusions to the first question?

Speed Zone Analysis



Urban Analysis

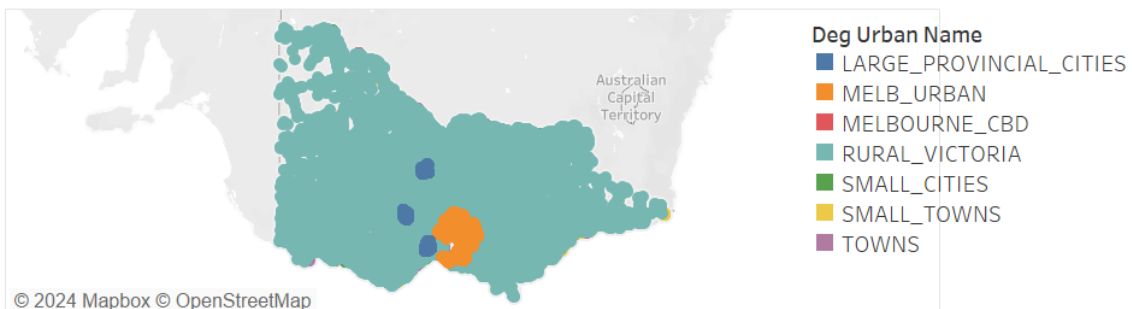
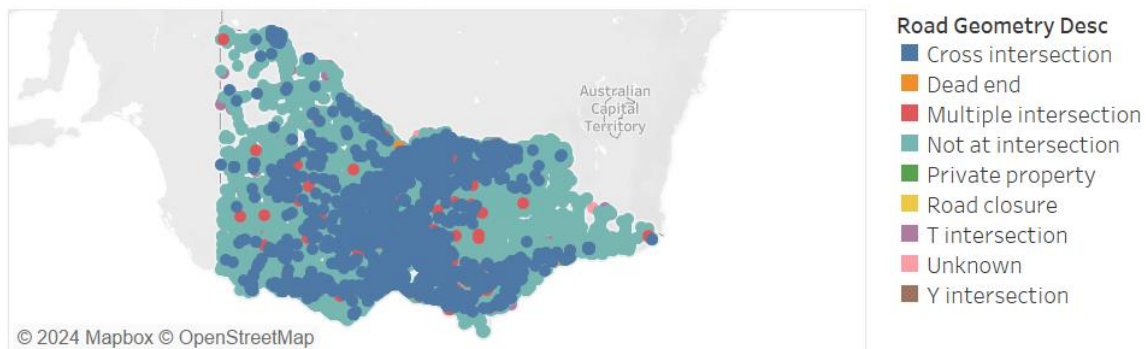


Figure 6 Spatial Analysis of Road Accidents (Speed Zone and Urban)

Based on the analysis of speed zones, it has been found that a considerable number of accidents are taking place in areas where the speed limit is around 100 and 110 km/h. This indicates that speed significantly contributes to road accidents (McPherson, 2024).

The map clearly illustrates that most vehicle crashes occur in rural Victoria. Consequently, we need to explore additional factors such as driver behaviour and road safety enforcement in rural regions to gain a comprehensive understanding of this phenomenon.

Road Geometry Analysis



Light Condition Analysis

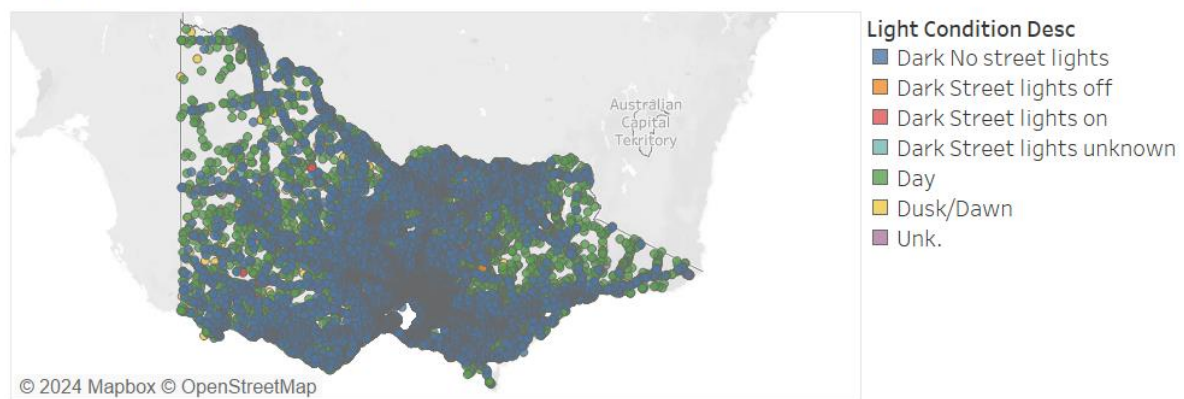


Figure 7 Spatial Analysis of Road Accidents (Road Geometry and Light Condition)

The Choropleth map also reveals that accidents predominantly occur outside intersections, particularly in areas with limited lighting, such as dark roads without streetlights. It proves that the most common case is because of minor mistakes such as being distracted or fatigued (QBE, 2017).

In summary, integrating spatial mapping with additional information about road infrastructure and environmental factors enhances understanding of where and why accidents occur geographically.

References

Horsch, A. (2021, February 15). *Detecting and Treating Outliers In Python — Part 3*. Medium.

<https://towardsdatascience.com/detecting-and-treating-outliers-in-python-part-3-dcb54abaf7b0>

McPherson, E. (2024, March 12). *The common factor first responders say they see in almost all fatal road crashes* [Review of *The common factor first responders say they see in almost all fatal road crashes*]. 9 News; 9 News. <https://www.9news.com.au/national/road-toll-australia-the-one-common-factor-first-responders-say-they-see-at-fatal-road-crashes/4b1c7ab2-1583-45ab-b517-97cfe43f8908>

QBE. (2017). *The most common causes of car accidents in Australia* / QBE AU. Australia.

<https://www.qbe.com/au/news/the-most-common-causes-of-car-accidents>