# Group 12 Exploratory Data Analysis
## DATA301

Wian Lusse (300489294), Nicholas Gibbs (300579601), Satrio Wiradika (300578313)

5 September 2023

## Contents

## Background and Data

Our project specifically focused on working with the Crash Analysis System (CAS) data sourced from Waka Kotahi, the New Zealand Transport Agency. We accessed the CAS data through the Open Data platform provided by Waka Kotahi offering a comprehensive dataset that is updated monthly containing detailed data about all traffic crashes reported by the NZ Police to Waka Kotahi. This data covers crashes on all New Zealand roadways or places where the public has legal access to a motor vehicle that has occurred across New Zealand since the 1st of January 2000. The dataset only includes crash variables that are non-personal data.

The CAS dataset is of great interest as it can be a resource to uncover insights into the domain of road safety. Specifically in our world where transportation plays a pivotal role in connecting societies, economies, and individuals, being able to comprehend the underlying dynamic of traffic crashes can be really informative. By using the CAS dataset, we are able to analyse the data from traffic crashes offering a different perspective on road safety, crash patterns, and the relationship between these traffic crashes and the geographic factors in New Zealand. At its core, the CAS data is much more than a repository of data, but a reflection of real-world events. Every entry within the dataset encapsulates a moment of impact, actions and consequences, and the complexity of vehicle interactions. By examining this collection of data, researchers will have the possibility to potentially unveil patterns within traffic crashes, and identify the factors that are contributing to their occurrences, to then create strategies to mitigate their impact. Specifically within the CAS dataset, its comprehensive nature is not being limited to a single dimension of analysis. Instead, it offers a vast range of variables that can be analysed. These variables span a wide range from the obvious elements such as weather conditions, vehicle types, and speed limit, to the more nuanced factors such as crash severity, and the region in which the traffic crash occurred. These wide ranges of variables can help researchers to inform and design targeted measures. For instance, if the data reveals a connection between a road with a higher incidence of accidents during certain weather conditions, additional efforts can be put in place to improve signage or alter the speed limits. Even for researchers who are unfamiliar with the dataset itself, they can answer inquiries such as 'What are the primary factors contributing to road crashes in New Zealand?' or 'How do crash rates differ across distinct geographical regions and evolve over time?.

The CAS data has numerical and categorical data types, structured with 70 columns and 821744 rows. Numerical attributes consist of various quantitative measures, including the unique identifier of an area unit, the number of lanes on the crash road, the unique identifier of a mesh block, and the speed limit in force at the crash site at the time of the crash may be a number, or LSZ for a limited speed zone. As to the categorical data, here are the features that we will use to answer our research question:

- 'crashSeverity' which represents The severity of a crash with possible values of 'F' (fatal), 'S' (serious), 'M' (minor), 'N' (non-injury) that are determined by the worst injury sustained in the crash at time of entry.
- 'weatherA' that indicates the weather at the crash time/place which includes 'Fine', 'Mist', 'Light Rain', 'Heavy Rain', 'Snow', and 'Unknown'.
- 'weatherB' which tells us additional information on the weather such as 'Frost', 'Strong Wind' or 'Unknown'.

Table 1: Count of Unique Values on Each Category

| crashSeverity | Count | weatherA | Count | weatherB | Count |
|---|---|---|---|---|---|
| Fatal Crash | 7589 | Fine | 635621 | Frost | 9254 |
| Minor Crash | 191336 | Hail or Sleet | 132 | Strong wind | 14389 |
| Non-Injury Crash | 575954 | Heavy rain | 33153 | Unknown | 798101 |
| Serious Crash | 46865 | Light rain | 124210 | | |
| | | Mist or Fog | 11306 | | |
| | | Snow | 1544 | | |
| | | Unknown | 15778 | | |

Table 2: Missing and region data summary

| | Count | Percentage | region | Count |
|---|---|---|---|---|
| crashSeverity | 0 | 0.0000000 | Auckland Region | 285346 |
| weatherA | 0 | 0.0000000 | Bay of Plenty Region | 47177 |
| weatherB | 0 | 0.0000000 | Canterbury Region | 82146 |
| region | 3188 | 0.0038796 | Gisborne Region | 9784 |
| | | | Hawke's Bay Region | 32388 |
| | | | Manawatū-Whanganui Region | 46329 |
| | | | Marlborough Region | 8266 |
| | | | Nelson Region | 8076 |
| | | | Northland Region | 33299 |
| | | | Otago Region | 44574 |
| | | | Southland Region | 20234 |
| | | | Taranaki Region | 18604 |
| | | | Tasman Region | 7541 |
| | | | Waikato Region | 87849 |
| | | | Wellington Region | 79725 |
| | | | West Coast Region | 7218 |

- 'region' which identifies the local government (LG) region. The boundaries match territorial local authority (TLA) boundaries.

We found a varying degree of missing information across the variables. While some variables exhibit high degree of completeness, others display more inconsistencies in the data. This does not necessarily imply a structured pattern, such as missing data clustered in specific time periods or geographic regions, but rather the lack of information that can be collected about the specific traffic crash. When specifically looking at the variables that we want to analyse, we can see that there is no missing data for crash severity, weatherA, and weatherB. However, for the region variable, there is a level of missing data with 3188 observations missing. This can be translated to a percentage of missing data of 0.0038796. In terms of errors within the specific variables that we wanted to analyse, there does not seem to be any errors. The observations of the different categories are consistent and there are no observations for categories in which they appear to be errors.

## Ethics, Privacy and Security

Ethical considerations hold an important role in our project as we delve into analysing the CAS dataset. While the dataset can be an extraordinary resource in helping analyse variables and generate insights, it presents us with ethical considerations that we must address. Given the sensitive nature of crash data, it is crucial that we exercise discretion in sharing insights that could inadvertently identify individuals involved in accidents. Finding a balance between informative reporting and preserving privacy is essential. Additionally, presenting the data in a way that avoids overemphasising accidents while focusing on generating useful insights is ethically crucial. Acknowledging the broader implications of our possible findings, particularly when suggesting insights, can ensure that our insights contribute positively to aspects such as road safety without causing much distress to those affected by road accidents. Another ethical responsibility lies in the manner in which we choose to convey the data to our audience. Committing to ethical practices ensures that we refrain from any form of sensationalism that can arise from the presentation of accident-related information. Instead, our approach centers on accentuating the importance for road safety improvements. This demands that the design of our communication is focused on the goal of mitigating accidents. We need to approach the task of suggesting insights with a heightened sense of responsibility. Our acknowledgement of the potential impact that our insights could have on road safety measures ensures a positive change without causing unwarranted distress to those who have been directly or indirectly affected by traffic crashes. Our project can be seen from an aspect of social responsibility, where every action we make in our analysis is aligned with the overarching goal of road safety and societal well-being.

Privacy concerns are also a consideration when dealing with the CAS dataset. Even though there is no personal information within the dataset, privacy concerns can still exist. Even if de-identified, there is still potential for the data to be linked back to the specific individual or disclose details into their personal life. The nature of the CAS dataset necessitates a comprehensive approach to ensure the mitigation of privacy risks. Anonymisation techniques such as data aggregation, suppression, and pseudonymisation, might have been employed to safeguard the identities of those mentioned in the dataset. However, the effectiveness of these techniques is not always absolute, and the complicated connections between the variables in the dataset could potentially allow for the re-identification of individuals, compromising their privacy. Furthermore, the potential for unintended consequences emerges when privacy concerns are not thoroughly addressed. The insights drawn from our analysis might inadvertently perpetuate biases. Ensuring that care is taken to ensure that the analytical process respects the privacy of the individuals involved in the dataset is a must. In light of these considerations, it is crucial to adopt a comprehensive privacy framework that encompasses the technical safeguards but also the ethical guidelines. Collaboration with ethics could help in identifying potential pitfalls, and the transparency in the methods that are used for data analysis even if indirectly are steps in ensuring privacy through the information encapsulated within the dataset.

To maintain the security of our project data and results, several steps can be taken, even though they have not been implemented for the purpose of this report. One of the foremost strategies for safeguarding sensitive information is the implementation of data encryption protocols. By encrypting data both during storage and transmission, we can prevent potential breaches and unauthorised access attempts. Encryption converts the data into an unreadable format. Without the appropriate decryption key, it is extremely challenging for malicious actors to interpret the information even if

they manage to gain access. Additionally, controlling and restricting access to the data repository is crucial. This can be achieved through the deployment of secure authentication and authorisation mechanisms. By implementing strong authentication methods such as multi-factor authentication and ensuring that only authorised personnel possess the necessary credentials, we can effectively prevent unauthorised entry to the data and results. Authorisation mechanisms can be employed to establish varying levels of access privileges based on roles within the project team, thereby minimizing the risk of inadvertent data exposure. In the event that an incident occurs such as data corruption, hardware failures, or cyberattacks, regular data backups emerge as another critical line of defence. Regularly backing up project data is a reliable means of preventing data loss and facilitating swift recovery. However, it is essential to store these backups in a secure manner. Adopting a best practice approach involves storing backups in off-site locations that are well-protected and disconnected from the primary network to safeguard against potential simultaneous compromise of both primary and backup data. Furthermore, ensuring that the secure collaboration among the project stakeholders is vital. Securing the data-sharing methods plays another crucial role in maintaining data confidentiality while enabling efficient teamwork. Secure file-sharing platforms equipped with end-to-end encryption and access controls provide a secure environment for sharing sensitive documents and information. Virtual private networks can also be employed to establish encrypted communication channels, particularly useful when collaborating with remote team members or external partners.

## Exploratory Data Analysis

### Exploring the Difference of Weather Conditions Across All the Regions in New Zealand

Adverse weather conditions play a crucial role in shaping the dynamics of road safety, significantly amplifying the potential for accidents to escalate into fatal or critical events. These environmental factors introduce a complex set of challenges that drivers must navigate through, demanding heightened vigilance, skill, and adaptability. The impact of severe weather encompassing heavy rainfall, snow, and fog, extends beyond just the inconvenience; it directly interacts with the systems that govern road safety, creating an environment that demands caution from vehicle users. According to a study conducted in 2019 (Fanny Malin 2019), the relative accident risks are increased for poor road weather conditions; however, they are highest for icy rain and slippery and very slippery road conditions. Especially those stemming from extreme weather, are more likely to result in grave consequences compared to accidents unaffected by adverse weather.
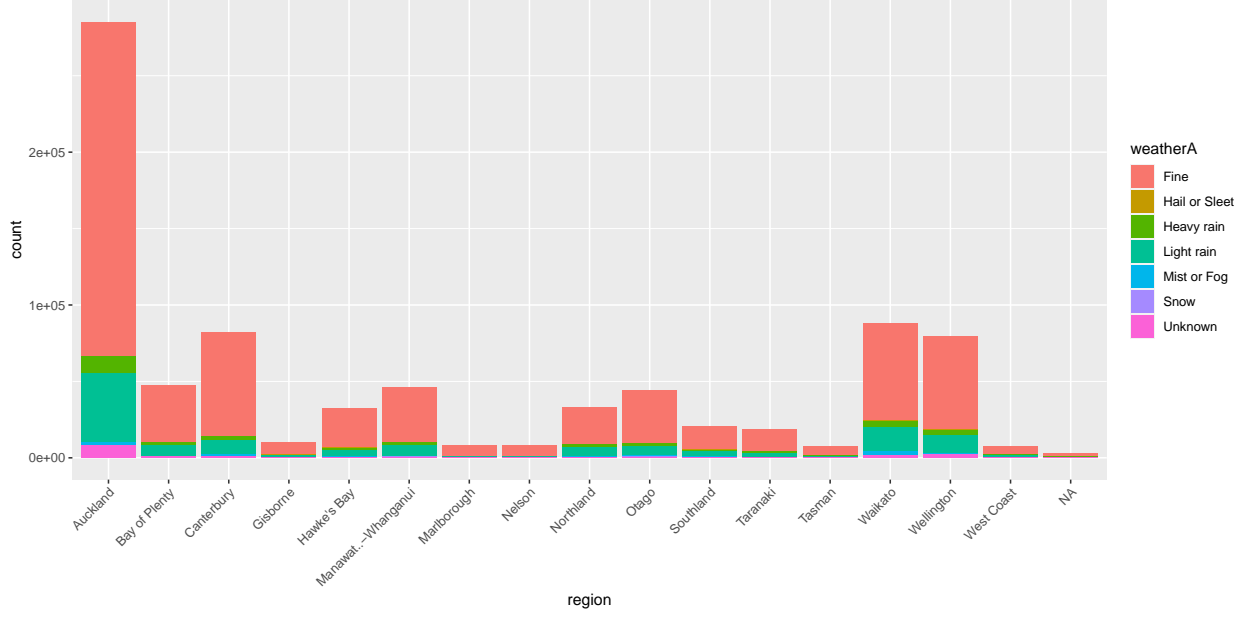
Figure 1: The Distribution of Weather Conditions in the Regions of New Zealand

The provided data presents weather-related information categorized by region. Several key patterns and conclusions can be drawn from this dataset. Firstly, the "Fine" weather condition is prevalent across most regions, with Auckland Region, Bay of Plenty Region, Canterbury Region, and many others reporting substantial counts under this category. Conversely, "Hail or Sleet" conditions are relatively rare across all regions, indicating infrequent occurrences of such extreme weather events. Secondly, "Heavy rain" and "Light rain" are prevalent conditions in many regions, suggesting a significant amount of precipitation is experienced throughout New Zealand. These regions include Auckland, Canterbury, Gisborne, Hawke's Bay, Manawatū-Whanganui, Northland, Otago, Taranaki, Tasman, Waikato, Wellington, and West Coast. Thirdly, "Mist or Fog" conditions are reported in various regions, but the counts are relatively lower compared to other weather conditions. Snowfall is infrequent, mainly occurring in limited quantities in certain regions like Auckland, Bay of Plenty, Canterbury, Gisborne, Hawke's Bay, Manawatū-Whanganui, Nelson, Otago, Southland, Taranaki, Tasman, Waikato, Wellington, and West Coast. Lastly, the "Unknown" category suggests that in some instances, weather conditions were not clearly identified or recorded. Overall, this data highlights the diverse weather conditions experienced across different regions in New Zealand, with variations in the prevalence of specific weather patterns.

**Is There Any Difference on the Level of Each Crash Severity Between the Regions?**
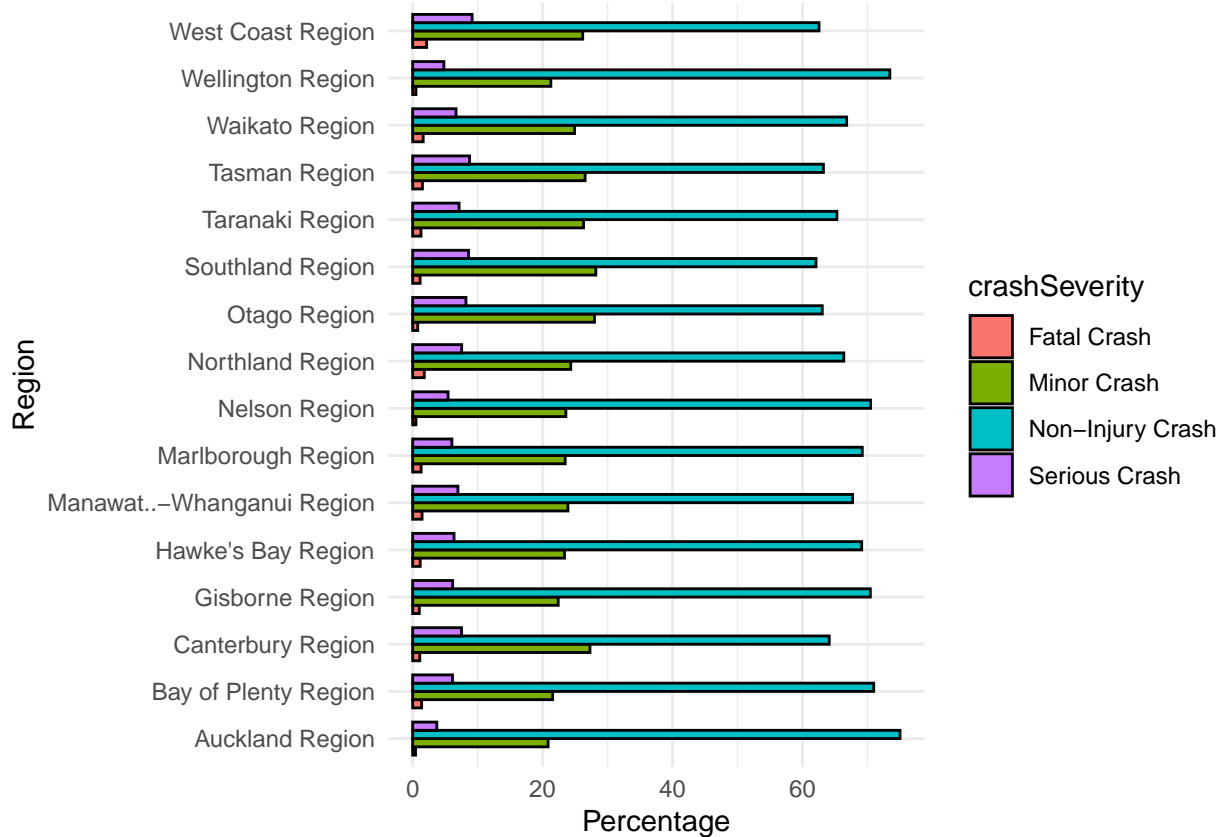
Figure 2: Percentage Grouped Bar Chart of Crash Severity by Region

The plot of the percentage of crash severity by region in New Zealand, allows us to see how the distribution varies between regions. Specifically, it presents a breakdown of road crash severities across the various regions in New Zealand. There is a notable variation in crash severity patterns across New Zealand regions, with some regions experiencing a higher prevalence of non-injury and minor crashes, while others have a relatively higher proportion of serious and fatal crashes. These differences may be influenced by factors such as local road conditions, traffic density, and enforcement efforts. Auckland stands out with a high percentage 75.01% of non-injury crashes, suggesting a higher frequency of less severe accidents. Serious and fatal crashes are comparatively lower. Bay of Plenty also has a substantial share of non-injury crashes 70.96% but a somewhat higher percentage of minor crashes 21.54% compared to Auckland, making it relatively safer. Canterbury has a significant number of non-injury crashes 64.11%, but it also experiences a higher proportion of serious crashes 7.51% compared to some other regions. Gisborne and Hawke's Bay both prioritize dealing with less severe incidents, with a focus on non-injury and minor crashes, resulting in lower percentages of serious and fatal crashes. Northland has a relatively higher percentage of serious crashes 7.53%, indicating a potential need for improved safety measures. Otago experiences a higher percentage of both minor and serious crashes compared to other regions, suggesting a need for heightened safety awareness. Tasman has a lower percentage of fatal and serious crashes, indicating a safer driving environment with a focus on minor and non-injury incidents. These regional variations may be influenced by factors such as local road conditions, traffic density, and law enforcement efforts. Further analysis may be needed to understand the underlying causes.

**Exploring the effect of wind on the severity of vehicle crashes**

Strong winds can play a role in increasing the risk of accidents and potentially be a cause of fatal or serious vehicle crashes. Strong winds can pose multiple hazards on the road. High winds can lead to sudden obstacles on the road and can also turn loose debris like rock into projectiles. High winds can also reduce vehicle control making it harder to stay in lanes safely or even be able to push large vehicles off the side of the road. According to this study in 2019 (Bhattachan et al. 2019), the fatality rate related to winds is almost twice as high as the rate in accident caused my weather conditions other than winds. This might suggest that wind-related accidents tend to be more lethal compared to accidents caused my other type of weather.
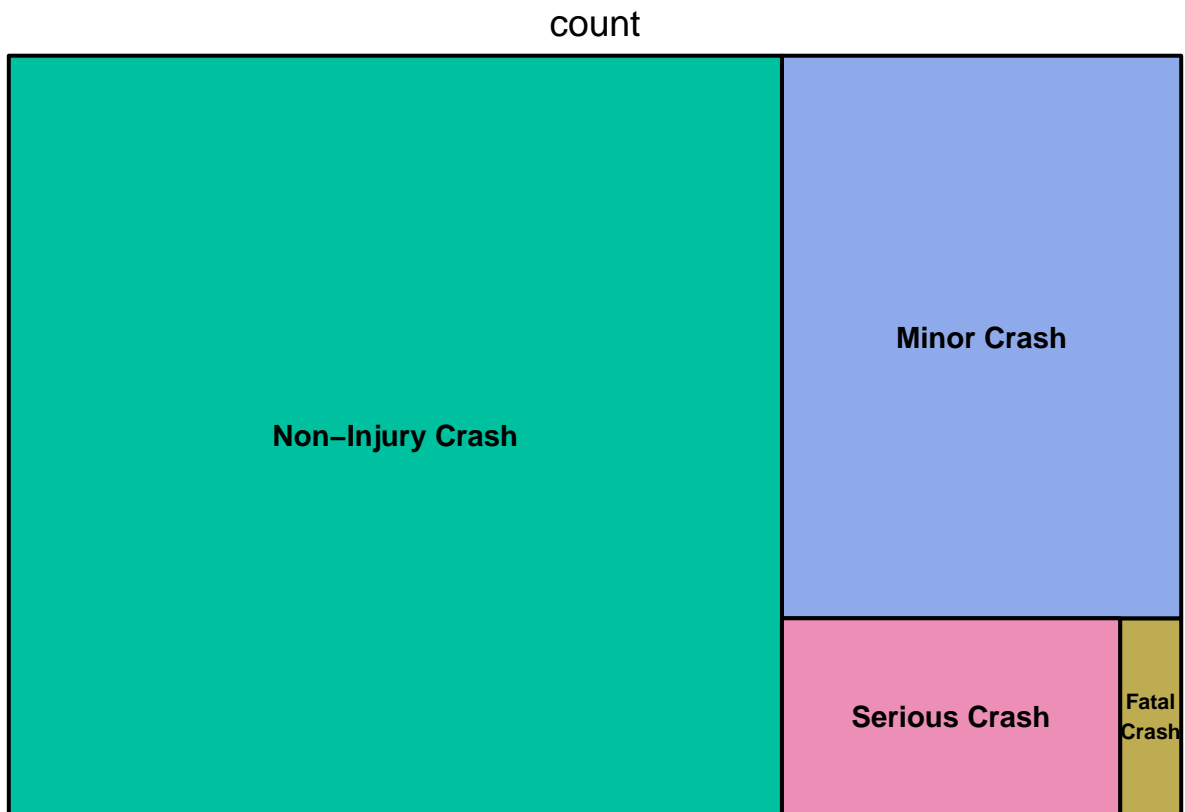


Figure 3: Crash Severity Across All Regions by Strong Winds

Table 3: Counts of Crash Severity Levels for Strong Wind Cases

| crashSeverity | count |
|---|---|
| Non-Injury Crash | 9489 |
| Minor Crash | 3632 |
| Serious Crash | 1074 |

| crashSeverity | count |
|---|---|
| Fatal Crash | 194 |

The data plotted in Figure 3 and displayed in Table 3 presents a breakdown of crash severity levels in cases affected by strong winds. The most prevalent outcome is non-injury crashes, with a count of 9489 cases, indicating that strong wind conditions more commonly result in accidents where no injuries are reported. Following this, minor crashes account for 3632 cases, suggesting a substantial number of incidents with minor injuries or property damage. Serious crashes are less frequent, with 1074 cases, indicating that strong winds can lead to more severe accidents, but these are relatively rare in comparison to non-injury and minor crashes. The least frequent outcome is fatal crashes, with a count of 194, signifying that fatalities due to strong wind-related accidents are relatively infrequent. In conclusion, this data underscores the importance of caution and preparedness during strong wind conditions, as they can lead to a range of crash severities, with non-injury and minor crashes being the most common outcomes, while serious and fatal crashes are less frequent but still significant concerns for road safety.
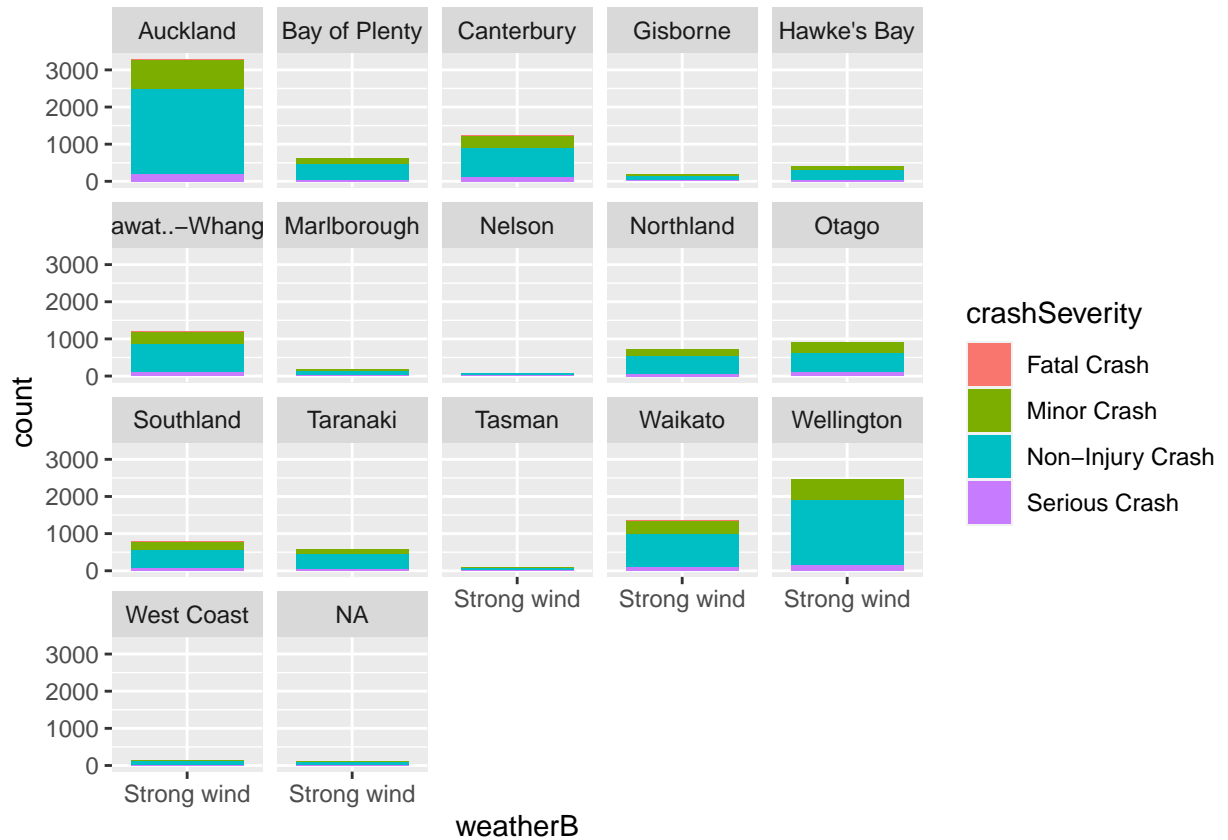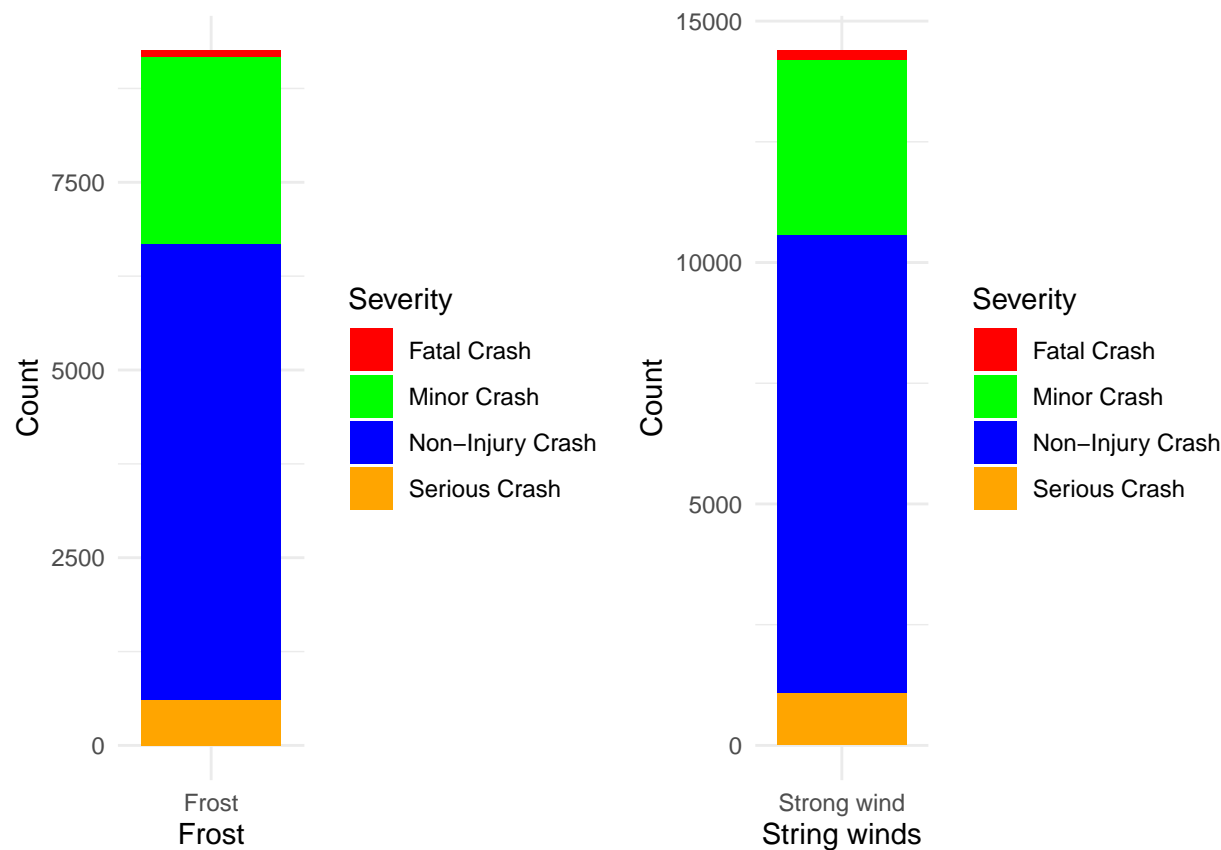


Figure 4: Crash Severity in Each Region for Strong Wind Condition

The data plotted in Figure 4 reveals the crash severity in various regions during strong wind conditions. In Auckland, strong winds resulted in 27 fatal crashes, while Bay of Plenty had 5,

Canterbury had 24, Gisborne had 2, and Hawke's Bay had 12 fatal crashes. Manawatū-Whanganui experienced 26 fatal crashes, Marlborough had 5, Northland had 8, Otago had 18, and Southland had 14 fatal crashes during strong winds. In contrast, minor crashes during strong wind conditions were more prevalent, with Auckland reporting 789, Bay of Plenty 152, Canterbury 340, Gisborne 42, Hawke's Bay 109, and Manawatū-Whanganui 317. The trend continued with other regions, including Nelson 22, Northland 190, Otago 275, Southland 227, Taranaki 127, Tasman 29, Waikato 368, Wellington 550, and West Coast 30. Furthermore, non-injury crashes were significantly higher in frequency, with Auckland leading at 2271, followed by Bay of Plenty at 418, Canterbury at 760, Gisborne at 130, and others. Serious crashes were less common, with Auckland having 198, Bay of Plenty 43, and Canterbury 116. Gisborne, Hawke's Bay, Manawatū-Whanganui, and other regions also reported varying numbers of serious crashes. These data suggest that strong wind conditions generally lead to a higher incidence of minor and non-injury crashes, while fatal and serious crashes are relatively rarer occurrences.

**Comparing light rain to strong winds**



Further analysis of crash severity data in different regions under frost and strong wind conditions provides additional insights. In frosty weather, there were 86 fatal crashes, indicating a significant risk to life, while strong wind conditions led to 194 fatal crashes, suggesting a higher fatality rate under strong winds. However, frost conditions also saw 604 serious crashes, which is noteworthy and potentially indicative of more severe accidents overall. In contrast, strong wind conditions had 1,074 serious crashes, indicating a higher propensity for serious accidents. Nonetheless, frosty weather conditions yielded a substantially larger number of non-injury crashes at 6,064 compared

to 9,489 non-injury crashes during strong wind conditions. This indicates that while strong winds may lead to more severe accidents, frosty weather results in a higher number of accidents without injuries.

## Individual Contributions

### Wian Lusse (300489294)

Wian Lusse made significant contributions to the project, particularly in the area of data analysis and visualization. Their role primarily involved handling references and creating wind plots. Wian's work on managing references was crucial for the project. They likely organized and cited sources, ensuring that the project was well-researched and credible. Proper referencing is vital in any research or data analysis project to give credit to original sources and avoid plagiarism. Also helping create wind plots suggests that Wian was involved in visualizing weather data related to wind patterns. Wind plots were essential in understanding how wind conditions affect crash severity, and in understanding any relationships or evidence in suggesting a difference between wind types and crash severity.

### Nicholas Gibbs (300579601)

Nicholas Gibbs played a pivotal role in setting up and maintaining the project's version control and contributed to weather data visualization. The creation of a GitHub repository indicates that they were responsible for managing the project's version control. This is a critical task in collaborative projects as it allows team members to work on code or documents simultaneously without conflicts. It also ensures that a history of changes is maintained, facilitating collaboration and troubleshooting. Like Wian, Nicholas was involved in the creation of weather plots and helped in the wind plots. Weather plots are essential for conveying complex weather data in a visually accessible format. Weather plots were essential in understanding how wind conditions affect crash severity, and in understanding any relationships or evidence in suggesting a difference between weather types and crash severity.

### Satrio Wiradikas

Satrio Wiradikas made valuable contributions, particularly in data management and weather data visualization. Satrio's role in importing the data suggests that they were responsible for gathering, cleaning, and preparing the raw data for analysis. Data import is a critical initial step in any data-driven project, as the quality and structure of the data greatly influence the subsequent analysis. Similar to Wian and Nicholas, Satrio worked on weather and wind plots. These plots contributed to the project's overall analysis and helped in visually representing weather patterns or trends.

## References

Bhattachan, Abinash, Gregory S Okin, Junzhe Zhang, Solomon Vimal, and Dennis P Lettenmaier. 2019. "Characterizing the Role of Wind and Dust in Traffic Accidents in California." *GeoHealth* 3 (10): 328–36.

Fanny Malin, Satu Innamaa, Ilkka Norros. 2019. "Accident Risk of Road and Weather Conditions on Different Road Types." *ScienceDirect* 122 (1): 181–88.