

# Análise da base US accidents

Ana Beatriz da Silva Maques

João Victor Pietchaki

Leonardo Eizo Sakai

Vitor Pignatari

2024-12-15

A principal atividade de um estatístico e cientista de dados é a análise de dados, quaisquer sejam as complexidades envolvidas. Tal consiste em um estudo pormenorizado de um conjunto de informações a fim de gerar conhecimento, instigar curiosidade, produzir ideias, detectar padrões a partir do que antes não tinha significado nenhum. Este presente artigo publica um conjunto de análises estatísticas sobre a base de dados US Accidents. Foi elaborado por estudantes do curso de Bacharelado em Estatística e Ciência de Dados da UFPR para a disciplina “Elementos de Programação para Estatística”, a fim de que aprimorem suas habilidades como analistas. As análises foram feitas utilizando a linguagem de programação R. Com ela, dados foram manipulados, análises exploratórias e estatísticas realizadas, gráficos elaborados, este presente relatório feito, e o dashboard desenvolvido. Dessa forma, comunicaremos informações de 7 milhões de linhas de forma intuitiva e clara, evidenciando tendências e detectando padrões entre a ocorrência de acidentes e as variáveis que os acompanham.

## Índices

|         |   |   |
|---------|---|---|
| 0.1     | Introdução . . . . .  | 2 |
| 0.2     | Materiais e Métodos . . . . .   | 2 |
| 0.2.1   | Fonte de Dados . . . . .  | 2 |
| 0.2.2   | A Base de Dados US_accidents . . . . .  | 3 |
| 0.2.2.1 | Grupos de Variáveis . . . . .   | 4 |
| 0.2.3   | Manipulação de dados . . . . .  | 4 |
| 0.2.3.1 | Pré-processando as variáveis do grupo ‘Tempo’ . . . . .                                     | 5 |
| 0.2.3.2 | Pré-Processando as variáveis do grupo posição geográfica . . . . .                          | 5 |
| 0.2.3.3 | Pré-processando as variáveis do grupo ‘Condições Climáticas’ . . . . .                      | 6 |
| 0.2.3.4 | Pré-processando as variáveis do grupo ‘Condições de Infraestrutura e Sinalização’ . . . . . | 7 |

|             |   |    |
|-------------|---|----|
| 0.2.3.5     | Pré-processando as variáveis do grupo “Duração do acidente”.                    | 8  |
| 0.2.4       | Técnicas Estatísticas Aplicadas . . . . .                                       | 8  |
| 0.2.4.1     | Técnicas de Análise para as Variáveis Temporais . . . . .                       | 8  |
| 0.2.4.2     | Técnicas de Análise para as Variáveis de Localização . . . . .                  | 11 |
| 0.2.4.3     | Técnicas de Análise para as Variáveis de Condições Climáticas . . . . .         | 12 |
| 0.2.4.4     | Técnicas de Análise para as Variáveis de Sinalização e Infraestrutura . . . . . | 14 |
| 0.2.4.5     | Técnicas de Análise para as Variáveis de Duração do acidente.                   | 14 |
| 0.3         | Resultados e Discussão . . . . .  | 15 |
| 0.3.1       | Resultados para as Variáveis Temporais . . . . .                                | 15 |
| 0.3.2       | Resultados para as Variáveis de Localização . . . . .                           | 22 |
| 0.3.3       | Resultados para as Variáveis de Condições Climáticas . . . . .                  | 23 |
| 0.3.4       | Resultados para as Variáveis de Sinalização e Infraestrutura . . . . .          | 23 |
| 0.3.5       | Resultados para as Variáveis de Duração do acidente . . . . .                   | 24 |
| Referências | . . . . .   | 25 |

## 0.1 Introdução

A análise de dados é uma atividade central para estatísticos e cientistas de dados, permitindo transformar informações aparentemente desorganizadas em conhecimento útil, detectar padrões e gerar insights. Este artigo apresenta um conjunto de análises estatísticas sobre a base de dados US Accidents, elaborado por estudantes do curso de Bacharelado em Estatística e Ciência de Dados da UFPR para a disciplina “Elementos de Programação para Estatística”. As análises foram conduzidas utilizando a linguagem de programação R, abrangendo etapas de manipulação de dados, análise exploratória, construção de gráficos, desenvolvimento de um dashboard interativo em Shiny e a produção deste relatório técnico em formato de artigo científico. Com base nos 7 milhões de registros da base, foram identificadas tendências e padrões na ocorrência de acidentes e suas variáveis associadas, comunicados de maneira clara e intuitiva. Este trabalho visa aprimorar as habilidades analíticas dos estudantes e contribuir para a compreensão dos dados relacionados a acidentes nos Estados Unidos.

## 0.2 Materiais e Métodos

### 0.2.1 Fonte de Dados

Para este projeto, utilizamos o linguagem de programação R para analisar a base de dados US Accidents, a partir de um aplicativo Shiny para a criação de um dashboard interativo, permitindo tanto uma visão abrangente dos dados, quanto informações específicas sobre o comportamento das variáveis inter-relacionada e autonomamente. A base de dados pode ser encontrada em <https://www.kaggle.com/datasets/sobhanmoosavi/us-accidents>, e foi coletada utilizando APIs de trânsito em tempo real, nos Estados Unidos entre os anos de 2016 e 2023

(Moosavi, Samavatian, Parthasarathy, & Ramnath, 2019), (Moosavi, Samavatian, Parthasarathy, Teodorescu, et al., 2019).

## 0.2.2 A Base de Dados US\_accidents

A base de dados consiste em 46 variáveis e aproximadamente 7,8 milhões de observações sobre acidentes de trânsito nos Estados Unidos entre fevereiro de 2016 e março de 2023, com informações sobre a severidade do acidente, posição geográfica, data e hora do ocorrido, descrição, endereço, condições climáticas, infraestrutura do local e condições de luminosidade no momento do ocorrido. O programa R foi utilizado com os pacotes tidyverse para a importação e limpeza dos dados e análise exploratória, ggplot para a criação de gráficos, leaflet para a criação de visualizações em mapas, e shiny para a criação do dashboard. Em conjunto com Shiny, foi utilizado o pacote bslibs, que usa versões mais atualizadas do Bootstrap, permitindo a criação de layouts mais modernos e flexíveis.

A fim de começar a entender e analisar os acidentes nos EUA, iniciamos lendo o banco de dados US\_Accidents\_March23.csv. Os dados foram coletados por Sobhan Moosavi entre 01/01/2016 e 04/01/2023. Para otimizar o tempo de carregamento, salvamos o arquivo original dos dados em um de extensão .Rdata.

```
1 # Carrega os dados originais
2 # load("../99_dados/acidentes_US.RData")
3 load("../99_dados/acidentes_US_sample.RData")
4 # load("../99_dados/acidentes_US_date.RData")
5 load("../99_dados/acidentes_US_date_sample.RData")
6 st_read("../99_dados/Estados_US/tl_2024_us_state.shp")
7 acidentes_US_sample %>% class()
```

A fim de começarmos a entender esse conjunto de dados, vamos verificar a estrutura dos elementos presentes no objeto `acidentes_US`.

```
1 acidentes_US_sample %>% glimpse()
```

O conjunto de dados dispõe de 7728394 observações e 46 variáveis. Cada coluna é uma variável que descreve um aspecto do acidente, ao passo que cada linha corresponde única e exclusivamente a uma observação de acidente.

```
1 acidentes_US_sample %>% head(10)
```

### 0.2.2.1 Grupos de Variáveis

Olhando em um primeiro momento, todas essas informações nos confundem. Para contornar esse problema, identificamos grupos de variáveis que podem ser analisadas em conjunto e excluimos outras que observamos não contribuir significativamente.

Escolhemos não trabalhar com as colunas `ID`, `Source`, `Timzone`, `Airport_code`, `Weather_Timestamp`, `Sunrise_Sunset`, `Nautical_Twilight`, `Astronomical_Twilight`, `zipcode`, `End_Lat`, `End_Lng`, `Street`, `County`, `Wind_Direction`, `Weather_Timestamp`, `Description`.

Para otimizar nossas análises, então, definimos 5 grandes grupos de variáveis: tempo, localização, condições climáticas, condições de infraestrutura e sinalização, descrição humana do acidente.

A coluna de `severity` não foi agrupada e será analisada interconjuntamente.

As variáveis que compõem cada grupo são:

- Tempo:
  - `Start_Time`, `End_Time`, `Civil_Twilight`;
- Localização:
  - `Start_Lat`, `Start_Lng`, `End_Lat`, `End_Lng`, `Distance(mi)`, `Street`, `City`, `County`, `State`, `Zipcode`, `Country`;
- Condições Climáticas:
  - `Temperature(F)`, `Wind_Chill(F)`, `Humidity(%)`, `Pressure(in)`, `Visibility(mi)`, `Wind_Speed(mph)`, `Precipitation(in)`, `Weather_Condition`;
- Condições de Infraestrutura e Sinalização:
  - `Bump`, `Crossing`, `Give_Way`, `Junction`, `No_Exit`, `Railway`, `Roundabout`, `Station`, `Stop`, `Traffic_Calming`, `Traffic_Signal`, `Turning_Loop`; ...
- Descrição Humana do Acidente:
  - `Description`.

### 0.2.3 Manipulação de dados

Definimos que cada uma das cinco abas seria de cada um dos grupos de variáveis. Para isso, tivemos de fazer um pré-processamento dos dados para:

- otimizar o tempo de processamento de scripts;
- filtrar colunas que não seriam utilizadas;

- tratar valores faltantes;
- transformar variáveis para formatos mais adequados;
- criar novas variáveis;
- sumarizar os dados.

### 0.2.3.1 Pré-processando as variáveis do grupo ‘Tempo’

O primeiro passo foi filtrar as colunas que fazem parte do grupo ‘Tempo’ e criar uma nova tabela com essas variáveis.

```

1 # Selecciona as variaveis de 'Tempo'
2 acidentes_US_date <- acidentes_US_sample[, .(ID, Severity, Start_Time, End_Time, Civil_Twili
3
4 acidentes_US_date <- acidentes_US_sample[, `:=`(Start_Time = Start_Time,
5                                     End_Time = End_Time,
6                                     Civil_Twilight = Civil_Twilight)]
7
8 acidentes_US_date <- acidentes_US_date[, .(ID, Severity, Start_Time, End_Time, Civil_Twilight)]

```

Com apenas as colunas de interesse, a tabela `acidentes_US_date` foi criada. A partir dela, criamos novas variáveis para enriquecer a análise dos dados.

Esse `data.table` `acidentes_US_date` foi salvo em um arquivo `.RData` para ser utilizado posteriormente. Ele será a tabela base para todas as análises e gráficos realizados para o grupo ‘Tempo’.

### 0.2.3.2 Pré-Processando as variáveis do grupo posição geográfica

Para a visualização com o mapa, foi consultada uma base externa no [site](#), sobre os *shapefiles* dos estados. Assim, tendo os limites dos estados, foi possível criar uma paleta de cores de acordo com os valores da quantidade de acidentes que cada estado teve.

Para a plotagem do mapa de calor na visualização geográfica, foi calculada a localização aproximada do acidente, arredondando a latitude e longitude para que a renderização do código fosse mais rápida.

```

1 # Alterando o nome, para não ter que alterar todas as variáveis após
2 acidentes_arredondado <- acidentes_US_sample
3
4 # Shapefiles
5 estados <- st_read("../99_dados/Estados_US/tl_2024_us_state.shp")
6
7 # Arredondando latitudes e longitudes

```

```

8 acidentes_arredondado$Start_Lat <- round(acidentes_arredondado$Start_Lat,2)
9 acidentes_arredondado$Start_Lng <- round(acidentes_arredondado$Start_Lng,2)
10
11 # Alterando a base shp para facilitar o entendimento
12 nomeestados <- as.data.frame(estados) %>% select(STUSPS,NAME)
13 names(nomeestados) <- c("STUSPS","NOMEESTADO")
14
15 # Juntando as informações já tratadas
16 acidentes <- left_join(acidentes_arredondado,nomeestados, by = c("State"="STUSPS"))
17
18 acidentes_p_estado <- acidentes %>%
19   group_by(State, NOMEESTADO) %>%
20   summarise(n = n(), .groups = "drop")
21
22
23 # Juntando com a geometria dos estados
24 acidentes_p_estado <- left_join(acidentes_p_estado, estados, by = c("State" = "STUSPS"))
25 acidentes_p_estado <- st_as_sf(acidentes_p_estado)
26
27 # Substituindo NA por 0 para estados sem acidentes
28 acidentes_p_estado$n[is.na(acidentes_p_estado$n)] <- 0
29
30 mapa_de_calor <- acidentes %>%
31   group_by(Start_Lat, Start_Lng) %>%
32   summarise(n = n(), .groups = "drop")
33
34 cores <- c("white", "red", "darkred")
35 palestado <- colorNumeric(palette = cores, domain = c(0, max(acidentes_p_estado$n, na.rm = T

```

### 0.2.3.3 Pré-processando as variáveis do grupo 'Condições Climáticas'

A variável 'Weather Condition' possuía muitos possíveis valores, o que a deixava difícil analisar. O primeiro passo foi agrupar os valores parecidos.

```

1 #Agrupando valores parecidos em 'Weather Condition'.
2
3 acidentes_US_sample <- acidentes_US_sample %>%
4   mutate(Weather_Cluster = case_when(
5     grepl("Rain|Drizzle|Shower|Showers|Light Rain|Heavy Rain|Rain Showers|Light Rain Showers
6     grepl("Snow|Blowing Snow|Light Snow|Heavy Snow|Snow Grains|Snow Showers|Snow and Thunder
7     grepl("Thunder|Tornado|Thunderstorm|Light Thunderstorm|Heavy Thunderstorms and Rain|Thun
8     grepl("Wind|Windy|Blowing Dust|Blowing Snow|Duststorm|Dust Whirls|Blowing Dust / Windy|W

```

```

9     grepl("Fog|Haze|Shallow Fog|Patches of Fog|Mist|Light Fog|Fog / Windy|Partial Fog", Weather_Condition, ignore.case = TRUE) ~ "Hazardous",
10     grepl("Clear|Fair|Mostly Cloudy|Partly Cloudy|Overcast|Cloudy", Weather_Condition, ignore.case = TRUE) ~ "Fair / Windy",
11     grepl("Smoke|Volcanic Ash|Haze", Weather_Condition, ignore.case = TRUE) ~ "Hazardous",
12     grepl("Sleet|Ice Pellets|Freezing Rain|Light Freezing Rain|Freezing Drizzle|Light Freezing Rain", Weather_Condition, ignore.case = TRUE) ~ "Hazardous",
13     grepl("Squalls|Hail|Small Hail|Thunder / Wintry Mix|Hail", Weather_Condition, ignore.case = TRUE) ~ "Hazardous",
14     grepl("Mist|Blowing Snow Nearby|Sand|Blowing Snow", Weather_Condition, ignore.case = TRUE) ~ "Hazardous",
15     grepl("Volcanic Ash|Blowing Sand|Sand", Weather_Condition, ignore.case = TRUE) ~ "Dust Storm",
16     grepl("Fair / Windy|Fair", Weather_Condition, ignore.case = TRUE) ~ "Fair / Windy",
17     TRUE ~ "Other"
18 ))

```

Depois, foi criada uma nova variável para identificar se choveu ou não.

```

1 acidentes_US_sample <- acidentes_US_sample %>%
2   mutate(Choveu = ifelse(Precipitation(in) > 0, "Sim", "Não"))

```

#### 0.2.3.4 Pré-processando as variáveis do grupo 'Condições de Infraestrutura e Sinalização'

```

1 # Seleciona as variáveis de interesse
2 poi <- acidentes_US_sample[, .(Severity, Start_Time, City, State, `Temperature(F)`,
3   Amenity, Bump, Crossing, Give_Way, Junction, No_Exit,
4   Railway, Roundabout, Station, Stop, Traffic_Calming,
5   Traffic_Signal, Turning_Loop)]
6
7 # Cria duas seleções de dados, contando a quantidade de acidentes com e sem cada ponto de interesse
8
9 poi_true <- poi %>%
10   select(Amenity:Turning_Loop) %>%
11   pivot_longer(Amenity:Turning_Loop, names_to = "PoI") %>%
12   filter(value == 1) %>%
13   group_by(PoI) %>%
14   summarise(t = n()) %>%
15   arrange(desc(t))
16
17 poi_false <- poi %>%
18   select(Amenity:Turning_Loop) %>%
19   pivot_longer(Amenity:Turning_Loop, names_to = "PoI") %>%
20   filter(value == FALSE) %>%
21   group_by(PoI) %>%
22   summarise(f = n()) %>%

```

```

23   arrange(f)
24
25   # Cria uma nova coluna com a soma dos pontos de interesse
26   poi_sum <- poi[, PoI_Sum := rowSums(.SD), .SDcols = c("Amenity", "Bump", "Crossing", "Give_Way",
27                                                         "Stop", "Traffic_Calming",
28                                                         "Traffic_Signal", "Turning_Loop")]

```

### 0.2.3.5 Pré-processando as variáveis do grupo “Duração do acidente”.

Primeiramente, foi preciso criar uma nova variável chamada “Duration\_minutes” e calcular a duração do acidente em minutos à partir das duas variáveis: Start Date e End Date.

```

1   acidentes_US_sample <- acidentes_US_sample %>%
2     mutate(Duration_minutes = as.numeric(difftime(End_Time, Start_Time, units = "mins")))

```

## 0.2.4 Técnicas Estatísticas Aplicadas

Como nossas análises foram baseadas em grupos de variáveis que compartilham contextos específicos, as técnicas estatísticas utilizadas variaram de uma para a outra de acordo com os respectivos objetivos.

### 0.2.4.1 Técnicas de Análise para as Variáveis Temporais

Para as variáveis temporais, foram considerados os diferentes níveis de granularidade das informações, como ano, mês, dia, hora, dia da semana e faixa de horário. A análise foi feita com base numa análise exploratória de dados, estatística descritiva, análise de dados categóricos, visualização de dados aplicada.

Um dos objetivos foi analisar a distribuição das frequências de acidentes entre os diferentes níveis de granularidade de tempo (ano, mês, mês-dia, dia), e verificar se existem padrões sazonais. Para isso, foram utilizados gráficos de barras e de pontos, bem como tabelas de frequência que foram geradas com o auxílio dos pacotes `tidyverse` e `ggplot2`.

```

1   acidentes_ano_barplot <- acidentes_US_sample_date[Civil_Twilight == "Day" | Civil_Twilight ==
2     ".(N,
3       Distinct_date = uniqueN(Date),
4       Media_diaria = round(.N/uniqueN(Date),0)),
5     by = Year][
6       order(Year)][
7       , Percent_variation := round((Media_diaria -
8   setorder(Year) %>% # order data set increasingly by year

```



```

9   as.data.frame() %>% # Transforming data.table into data.frame
10  ggplot() +
11  aes(x = Year, y = N,
12      text = paste("Ano:", Year,
13                  "<br>Qtd. Acidentes:", round(N/1000, 1), "k",
14                  "<br>Disa Obs. no Ano:", Distinct_date,
15                  "<br>Acidentes por Dia:", round(Media_diaria/1000, 1), "k")) +
16  geom_bar(stat = "identity", fill = "darkred") +
17  geom_text(aes(label = percent(Percent_variation)), size = 3) +
18  scale_y_continuous(
19    labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
20  ) +
21  labs(title = "Qtd. de Acidentes por Ano", x = "Ano", y = "Qtd. Acidentes") +
22  theme_minimal() +
23  theme(plot.title = element_text(hjust = 0.5)) +
24  plotly()
25  acidentes_ano_barplot %>% ggplotly(tooltip = "text")

```

Juntamente, para acrescentar mais informações à análise, foi feita uma sumarização dessas distribuições de frequências de acidentes por crepúsculo civil, a fim de verificar se existe uma relação entre a quantidade de acidentes e a luminosidade do dia.

```

1  acidentes_ano_barplot_CT <- acidentes_US_sample_date[Civil_Twilight == "Day" | Civil_Twilight == "Night"]
2    , N_year := sum(N), by = Year][
3    , Percent := round((N/N_year), 2)] %>%
4  setorder(Year, Civil_Twilight) %>% # order data set increasingly by year
5  as.data.frame() %>% # Transforming data.table into data.frame
6  ggplot() +
7  aes(x = Year, y = N,
8      text = paste("Ano:", Year,
9                  "<br>Qtd. Acidentes:", round(N/1000, 1), "k",
10                 "<br>% Total Acidentes:", round(100*Percent, 0), "%")) +
11  geom_bar(aes(fill = Civil_Twilight), stat = "identity") +
12  scale_y_continuous(
13    labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
14  ) +
15  scale_fill_manual(values = c('Day' = '#D9A404', 'Night' = '#0F3BBF')) +
16  labs(title = "Qtd. de Acidentes por Ano", x = "Ano", y = "Qtd. Acidentes") +
17  theme_minimal() +
18  theme(plot.title = element_text(hjust = 0.5))
19  acidentes_ano_barplot_CT %>% ggplotly(tooltip = "text") %>%
20  layout(

```

```

21 legend = list(
22   title = list(text = 'Tipo de Dia'), # Título da legenda
23   font = list(size = 14), # Tamanho da fonte
24   itemsizing = 'constant', # Define o tamanho fixo dos itens da legenda
25   traceorder = 'normal', # Ordem das legendas
26   marker = list(
27     size = 10
28   )
29 )
30 )

```

Isso foi feito para os seguintes níveis de granularidade de tempo: ano, mês e dia. Para este último, foi traçada uma curva de ajuste com o auxílio da função `geom_smooth` do pacote `ggplot2`. Obtivemos resultados interessantes e os discutiremos na próxima seção.

Alternativamente, foi feita uma análise entre as variáveis categóricas, dia da semana e faixa horária, a fim de explorar a relação entre elas. Para isso, foram utilizadas tabelas de contingência bem como matriz de calor, sempre com `ggplot2`. Essa técnica foi aplicada para verificar se existem horários e dias da semana com maior incidência de acidentes ou com variação na severidade dos acidentes.

```

1 # Cria uma tabela de contingencia de acidentes por dia da semana e faixa horaria
2 acidentes_US_sample_date[, .N, by = .(Weekday, Time_Range)] %>%
3   dcast(Time_Range ~ Weekday, value.var = "N", fill = 0) %>% setorder(-Time_Range)

```

```

1 # Cria uma matriz de calor para as variaveis dias da semana e faixa horaria
2 acidentes_data_matriz_AWKTR <- acidentes_US_sample_date[Civil_Twilight == "Day" | Civil_Twil.
3   setorder(Weekday, -Time_Range) %>% # order data set increasingly by weekday and time range
4   as.data.frame() %>% # Transforming data.table into data.frame
5   ggplot() +
6   aes(x = Weekday, y = Time_Range, fill = N) +
7   geom_tile() +
8   scale_fill_gradient(low = "white", high = "darkred") +
9   labs(title = "Agenda Diaria de Acidentes", x = "Dia da Semana", y = "Horario") +
10  theme_minimal() +
11  scale_x_discrete(position = "top") +
12  theme(
13    axis.title.x.top = element_text(),
14    axis.text.x.top = element_text(),
15    axis.line.x.top = element_line(),
16    axis.ticks.x.top = element_line()
17  )

```

```

18   )
19   acidentes_data_matriz_AWKTR

```

Ademais, a fim de investigar a severidade média dos acidentes ao longo do tempo, foi realizada uma análise de tendência da média de severidade dos acidentes por ano. Para isso, foi feito um gráfico de pontos associado a um ajuste de curva (`geom_smooth`).

De modo geral, para cada análise foi feita uma sumarização de dados, seguida de uma visualização dos resultados obtidos. Os códigos acima são apenas exemplos e, na seção de resultados e discussão, mais detalhes serão mostrados.

#### 0.2.4.2 Técnicas de Análise para as Variáveis de Localização

```

1   leaflet() %>%
2     addTiles() %>%
3     addPolygons(
4       data = acidentes_p_estado,
5       fillColor = ~palestado(n),
6       weight = 2,
7       color = "black",
8       fillOpacity = 0.4,
9       label = ~NAME,
10      popup = ~glue("<b>Estado: </b> {NAME}<br><b>Quantidade de acidentes: </b> {n}"),
11      highlightOptions = highlightOptions(
12        weight = 3,
13        color = "#666",
14        fillOpacity = 0.6,
15        bringToFront = TRUE
16      ),
17      group = "Nível estadual"
18    ) %>%
19    addLegend(
20      "bottomright",
21      pal = palestado,
22      values = acidentes_p_estado$n,
23      title = "Número de Acidentes",
24      opacity = 1,
25      group = "Nível estadual"
26    ) %>%
27    addHeatmap(
28      data = mapa_de_calor,
29      lng = ~`Start Lng`,

```

```

30     lat = ~`Start Lat`,
31     intensity = ~n,
32     blur = 20,
33     radius = 15,
34     group = "Mapa de calor"
35   ) %>%
36   addLayersControl(
37     baseGroups = c("Nível estadual", "Mapa de calor"),
38     position = "topright"
39   )

```

#### 0.2.4.3 Técnicas de Análise para as Variáveis de Condições Climáticas

Para as variáveis de condições climáticas, foram realizadas análises de frequência de acidente por condição climática.

```

1  condicoes_count <- acidentes_US_sample %>%
2    group_by(Weather_Cluster) %>%
3    tally(name = "Count")
4  ggplot(condicoes_count, aes(x = reorder(Weather_Cluster, -Count), y = Count, fill = Weather_Cluster)) +
5    geom_bar(stat = "identity") +
6    labs(title = "Número de Acidentes por Condição Climática", x = "Condição Climática", y = "Número de Acidentes") +
7    theme_minimal() +
8    theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

Este gráfico mostra que a maioria dos acidentes ocorrem quando o céu está limpo.

Depois, foi analisado a correlação entre as variáveis climáticas.

```

1  climatic_data <- acidentes_US_sample %>%
2    select(`Temperature(F)`, `Humidity(%)`, `Pressure(in)`, `Wind_Speed(mph)`, `Visibility(mi)`)
3    na.omit()
4
5  cor_matrix <- cor(climatic_data)
6  corrplot(cor_matrix, method = "color", type = "upper", tl.col = "black")

```

Foi analisado que, algumas delas possuem correlação positiva e outras, negativa.

No próximo passo, o objetivo era analisar o número de acidentes por Condição Climática. Para isso, foi preciso agrupar por condição climática, contar o número de acidentes e fazer um gráfico de barras.

```

1 # Número de acidentes por condição climática
2 condicoes_count <- acidentes_US_sample %>%
3   group_by(Weather_Cluster) %>%
4   tally(name = "Count")
5 ggplot(condicoes_count, aes(x = reorder(Weather_Cluster, -Count), y = Count, fill = Weather_Cluster)) +
6   geom_bar(stat = "identity") +
7   labs(title = "Número de Acidentes por Condição Climática", x = "Condição Climática", y = "Número de Acidentes") +
8   theme_minimal() +
9   theme(axis.text.x = element_text(angle = 45, hjust = 1))

```

Depois, foram analisadas as distribuições de precipitação e visibilidade.

```

1 ggplot(acidentes_US_sample, aes(x = `Precipitation(in)`)) +
2   geom_histogram(binwidth = 0.1, fill = "blue", color = "black") +
3   labs(title = "Distribuição de Precipitação", x = "Precipitação (in.)", y = "Contagem")
4
5
6 ggplot(acidentes_US_sample, aes(x = `Visibility(mi)`)) +
7   geom_histogram(binwidth = 1, fill = "green", color = "black") +
8   labs(title = "Distribuição de Visibilidade", x = "Visibilidade (milhas)", y = "Contagem")

```

Para analisar o número de acidentes por condição de chuva no dia foi feito um gráfico de barras.

```

1 accident_chuva <- acidentes_US_sample %>%
2   group_by(Choveu) %>%
3   summarise(Accident_Count = n())
4
5
6 ggplot(accident_chuva, aes(x = Choveu, y = Accident_Count, fill = Choveu)) +
7   geom_bar(stat = "identity") +
8   labs(
9     title = "Acidentes em Cidades com Chuva vs. Sem Chuva",
10    x = "Choveu?",
11    y = "Número de Acidentes"
12  ) +
13  theme_minimal()

```

Aqui podemos analisar que há um número muito maior de acidentes quando não chove do que quando chove.

Por último, foram realizadas análises da relação entre a precipitação(chuva) e a severidade do acidente.

```
1 ggplot(acidentes_US_sample, aes(x = `Precipitation(in)`, y = Severity)) +
2   geom_jitter(width = 0.2, alpha = 0.5) +
3   labs(
4     title = "Relação entre Precipitação e Severidade",
5     x = "Precipitação (mm ou in)",
6     y = "Severidade"
7   ) +
8   theme_minimal()
```

#### 0.2.4.4 Técnicas de Análise para as Variáveis de Sinalização e Infraestrutura

Para as variáveis de infraestrutura e sinalização, foram realizadas análises de frequência de acidentes com e sem cada ponto de interesse

```
1 barplot <- acidentes_US_sample %>% select(Amenity:Turning_Loop) %>%
2   pivot_longer(everything(), names_to = "poi", values_to = "value") %>%
3   ggplot(aes(x = poi, fill = value)) +
4   geom_bar() +
5   labs(title = "Pontos de Interesse",
6     x = NULL,
7     y = NULL,
8     fill = "Presença de poi",
9     legend.position = "bottom") +
10  theme_minimal() +
11  theme(axis.text.x = element_text(angle = 45, hjust = 1)) +
12  scale_y_continuous(labels = comma) +
13  scale_fill_manual(values = c("darkgreen", "darkred"),
14    labels = c("Não", "Sim"))
15
16 ggplotly(barplot)
```

Este gráfico torna evidente que a grande maioria dos acidentes não ocorre próximo a qualquer ponto de interesse, porém é notável que há um aumento de acidentes em locais onde fluxos diferentes se encontram, como cruzamentos e junções.

#### 0.2.4.5 Técnicas de Análise para as Variáveis de Duração do acidente.

Para analisar a duração dos acidentes, foi construído um boxplot para analisar a distribuição de acidentes por severidade e por duração, para verificar se tinha uma relação entre essas duas variáveis.

```
1 ggplot(acidentes_US_sample, aes(x = as.factor(Severity), y = Duration_minutes)) +
2   geom_boxplot(outlier.shape = NA) +
3   labs(
4     title = "Duração dos Acidentes por Severidade",
5     x = "Severidade",
6     y = "Duração (minutos)"
7   ) +
8   scale_y_continuous(limits = quantile(acidentes_US_sample$Duration_minutes, c(0.0, 0.9))) +
9   theme_minimal()
```

Como possuía muitos outliers e isso dificultava a visualização, optamos por retirá-los e calcular os limites a partir da densidade dos dados.

## 0.3 Resultados e Discussão

Com os dados processados e metodologias estatísticas aplicadas, vamos agora de fato analisar os dados e extrair informações relevantes. Far-lo-emos para cada grupo de variáveis, apresentando os resultados obtidos e discutindo-os.

### 0.3.1 Resultados para as Variáveis Temporais

Começando pelo básico, vamos observar a distribuição de acidentes ao longo dos anos.

```
1 acidentes_ano_barplot <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight == "Night"] %>%
2   summarise(N = n(),
3             Distinct_date = uniqueN(Date),
4             Media_diaria = round(.N/uniqueN(Date),0)),
5   by = Year][
6   order(Year)][
7   , Percent_variation := round((Media_diaria - Media_diaria_lag1)/Media_diaria_lag1*100,1)]
8 setorder(Year) %>% # order data set increasingly by year
9 as.data.frame() %>% # Transforming data.table into data.frame
10 ggplot() +
11 aes(x = Year, y = N,
12     text = paste("Ano:", Year,
13                 "<br>Qtd. Acidentes:", round(N/1000, 1), "k",
14                 "<br>Disa Obs. no Ano:", Distinct_date,
```

```

15         "<br>Acidentes por Dia:", round(Media_diaria/1000, 1), "k")) +
16     geom_bar(stat = "identity", fill = "darkred") +
17     geom_text(aes(label = percent(Percent_variation)), size = 3) +
18     scale_y_continuous(
19         labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
20     ) +
21     labs(title = "Qtd. de Acidentes por Ano", x = "Ano", y = "Qtd. Acidentes") +
22     theme_minimal() +
23     theme(plot.title = element_text(hjust = 0.5)) +
24     plotly()
25 acidentes_ano_barplot %>% ggplotly(tooltip = "text")

```

Não nos resta dúvida de que houve um aumento significativo na quantidade de acidentes ao longo dos anos, com aumentos percentuais de 56, 24, 7, 23, 33, 16 e -33 por cento entre cada um dos anos de 2016 até 2023. Observa-se, entretanto, que houve apenas 73 dias observados no último ano, o que pode ter influenciado na queda de 33 por cento na quantidade de acidentes, não nos permitindo tirar conclusões definitivas.

Agora, realizando a mesma investigação, porém observando apenas ao nível mês da granularidade do tempo:

```

1 acidentes_mes_barplot <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight == "Night"]
2     .(N,
3         Distinct_date = uniqueN(Date),
4         Media_diaria = round(.N/uniqueN(Date),0)),
5     by = Month][
6         order(Month)][
7         , Percent_variation := round((Media_diaria -
8     setorder(Month) %>% # order data set increasingly by month
9     as.data.frame() %>% # Transforming data.table into data.frame
10     ggplot() +
11     aes(x = Month, y = N,
12         text = paste("Mes:", Month,
13             "<br>Qtd. Acidentes:", round(N/1000, 1), "k",
14             "<br>Dias Obs. no Mes:", Distinct_date,
15             "<br>Acidentes por Dia:", round(Media_diaria/1000, 1), "k")) +
16     geom_text(aes(label = percent(Percent_variation)), size = 3) +
17     geom_bar(stat = "identity", fill = "darkred") +
18     scale_y_continuous(
19         labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
20     ) +
21     labs(title = "Qtd. de Acidentes por Mes", x = "Mes", y = "Qtd. Acidentes") +

```



```

22 theme_minimal() +
23 theme(plot.title = element_text(hjust = 0.5))
24 acidentes_mes_barplot %>% ggplotly(tooltip = "text")

```

Aqui, observamos que a quantidade de acidentes é maior nos meses de outubro, novembro e dezembro, com uma queda significativa a partir de janeiro até julho. Isso nos permite inferir um certo grau de sazonalidade entre os meses de um ano, com pico sempre no mês mais festivo: dezembro.

Por finalizar as análises de distribuição de frequências, analisamos os acidentes dia a dia:

```

1 data_timeseries <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight == "Night", .()
2   setorder(Date) %>% # order data set increasingly by date
3   as.data.frame() %>% # Transforming data.table into data.frame
4   ggplot() +
5   aes(x = Date, y = N) +
6   geom_bar(stat = "identity", alpha = 0.4, fill = "darkred", alpha = 0.6) +
7   geom_point(size = 0.5, alpha = 0.5, color = "darkred",
8             aes(text = paste("Data:", Date, "<br>Qtd. Acidentes:", N/1000, "k")))) +
9   geom_smooth(method = "gam", color = "darkred", size = 0.5) +
10  scale_y_continuous(
11    labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
12  ) +
13  labs(title = "Qtd. de Acidentes por Data", x = "Data", y = "Qtd. Acidentes") +
14  theme_minimal()
15 data_timeseries %>% ggplotly(tooltip = "text")

```

A partir deste gráfico de pontos e colunas (frequência de acidentes por data), podemos observar tanto a tendência crescente anual bem como a sazonalidade mensal. Além disso, e de forma mais interessante, é nítida a existência visual de uma banda superior e inferior: isto nos indica uma variabilidade muito grande na ocorrência de acidentes em datas muito próximas. Utilizando o zoom do gráfico, podemos investigar mais a fundo essas variabilidades: os dias com menos acidentes ocorrem na sua imensa maioria dois a dois, e os com mais, de cinco em cinco. Ou seja, a variabilidade pode ser explicada pelo fato de que a quantidade de acidentes é muito maior entre os dias da semana do que em finais de semana!

Além disso, neste gráfico fica claro uma queda abrupta dos acidentes durante a pandemia, seguida de um aumento expressivo em quantidade e variabilidade.

Okay, e quanto à faixa horária? Como os acidentes se distribuem? Vamos investigar com a agenda diária de acidentes:

```

1 acidentes_data_matriz_AWKTR <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight ==
2   setorder(Weekday, -Time_Range) %>% # order data set increasingly by weekday and time range
3   as.data.frame() %>% # Transforming data.table into data.frame
4   ggplot() +
5   aes(x = Weekday, y = Time_Range, fill = N) +
6   geom_tile() +
7   scale_fill_gradient(low = "white", high = "darkred") +
8   labs(title = "Agenda Diaria de Acidentes", x = "Dia da Semana", y = "Horario") +
9   theme_minimal() +
10  scale_x_discrete(position = "top") +
11  theme(
12    axis.title.x.top = element_text(),
13    axis.text.x.top = element_text(),
14    axis.line.x.top = element_line(),
15    axis.ticks.x.top = element_line()
16  )
17  )
18  acidentes_data_matriz_AWKTR

```

Aqui, podemos observar que a quantidade de acidentes é maior nos horários de pico, entre 7hrs e 9hrs e entre 16hrs e 18hrs, e que a quantidade de acidentes é maior durante a semana do que nos finais de semana. Isso nos permite inferir que a quantidade de acidentes é maior nos horários de pico e durante a semana.

Em uma de minhas intuitivas explorações, decidi investigar a relação entre a quantidade de acidentes e a luminosidade do dia (crepúsculo civil). Para isso, novamente realizei um gráfico de pontos e barras de acidentes por data, porém agora separando pelas categorias do crepúsculo civil:

```

1 acidentes_data_timeseries_CT <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight ==
2   setorder(Date, Civil_Twilight) %>% # order data set increasingly by date
3   as.data.frame() %>% # Transforming data.table into data.frame
4   ggplot() +
5   aes(x = Date, y = N,
6     text = paste("Data:", Date, "<br>Qtd. Acidentes:", N/1000, "k")) +
7   geom_bar(stat = "identity", alpha = 0.4) +
8   geom_point(aes(color = Civil_Twilight), alpha = 0.4, size = 0.4) +
9   scale_y_continuous(
10    labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
11  ) +
12  scale_color_manual(values = c('Day' = '#D9A404', 'Night' = '#0F3BBF')) +
13  labs(title = "Qtd. de Acidentes por Data",

```

```

14     x = "Data", y = "Qtd. Acidentes") +
15     theme_minimal()
16 acidentes_data_timeseries_CT %>% ggplotly(tooltip = "text") %>%
17     layout(
18       legend = list(
19         title = list(text = 'Crepusculo Civil'), # Título da legenda
20         font = list(size = 14), # Tamanho da fonte
21         itemsizing = 'constant', # Define o tamanho fixo dos itens da legenda
22         traceorder = 'normal', # Ordem das legendas
23         marker = list(
24           size = 10
25         )
26       )
27     )

```

A partir deste gráfico, conseguimos observar todos os resultados já vistos até aqui e mais um pouco: agora conseguimos observar que o crepúsculo civil apresenta uma sazonalidade mensal muito expressiva. Em geral, de outubro a dezembro, os acidentes noturnos atingem o seu pico e os acidentes diurnos os seus vales.

Entretanto, ao mesmo tempo, não podemos afirmar que essa sazonalidade é devida ao crepúsculo civil, devido ao fato de que as estações do ano influenciam significativamente a luminosidade do dia no inverno e verão. Para contornar essa dúvida, investigamos e colocamos o mesmo gráfico, entretanto, ao invés de categorizar pelo crepúsculo civil, categorizamos por faixa horária do dia: acidentes que aconteceram das 06:00 às 18:00 e das 18:00 às 06:00.

```

1 acidentes_data_timeseries_6_18 <- acidentes_US_date[Day_Period == "06:00:00 -> 18:00:00" | D
2   setorder(Date, Day_Period) %>% # order data set increasingly by date
3   as.data.frame() %>% # Transforming data.table into data.frame
4   ggplot() +
5   aes(x = Date, y = N,
6       text = paste("Data:", Date, "<br>Qtd. Acidentes:", N/1000, "k")) +
7   geom_bar(stat = "identity", alpha = 0.4) +
8   geom_point(aes(color = Day_Period), alpha = 0.4, size = 0.4) +
9   scale_y_continuous(
10     labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
11   ) +
12   scale_color_manual(values = c('06:00:00 -> 18:00:00' = '#D9A404', '18:00:00 -> 06:00:00' =
13     labs(title = "Qtd. de Acidentes por Data", x = "Data", y = "Qtd. Acidentes") +
14     theme_minimal()
15 acidentes_data_timeseries_6_18 %>% ggplotly(tooltip = "text") %>%
16     layout(

```

```

17 legend = list(
18   title = list(text = 'Período do Dia'), # Título da legenda
19   font = list(size = 14), # Tamanho da fonte
20   itemsizing = 'constant', # Define o tamanho fixo dos itens da legenda
21   traceorder = 'normal', # Ordem das legendas
22   marker = list(
23     size = 10
24   )
25 )
26 )

```

A partir deste gráfico, conseguimos observar que a expressividade da sazonalidade mensal não existe e, na verdade, é sutil em ambas as faixas horárias. Isso nos permite inferir que a sazonalidade observada no gráfico anterior é devida à luminosidade do dia e não às estações do ano. Entretanto, conseguimos perceber que diferença da quantidade de acidentes entre dias da semana e finais de semana permanecem em ambas as faixas horárias.

E como se comporta a severidade dos acidentes ao longo do tempo? Vamos investigar a tendência da média de severidade dos acidentes por data:

```

1 severity_data_points <- acidentes_US_date[(Civil_Twilight == "Day" | Civil_Twilight == "Night") &&
2   .(Severidade_Media = mean(Severity)),
3   by = .(Date)][order(Date)] %>%
4   ggplot() +
5   aes(x = Date, y = Severidade_Media) +
6   geom_point(size = 0.5, alpha = 0.1, color = "darkred",
7     aes(text = paste("Data:", Date,
8       "<br>Severidade Média:", round(Severidade_Media, 1)))) +
9   geom_smooth(method = "lm", size = 0.5, se = FALSE, color = "darkred") +
10  ylim(0, 4) +
11  labs(title = "Severidade Média por Data", x = "Data", y = "Severidade") +
12  theme_minimal()
13 severity_data_points %>% ggplotly(tooltip = "text")

```

Essa tendência linear decrescente pode ser explicada pela variação da quantidade de acidentes ao longo do tempo para cada uma das severidades. Observemo-las:

```

1 acidentes_ano_barplot <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight == "Night"]
2   .(N,
3     Distinct_date = uniqueN(Date),
4     Media_diaria = round(.N/uniqueN(Date),0)),
5   by = .(Year, Severity)][

```

```

6                                     order(Year,Severity)] %>%
7   setorder(Year) %>% # order data set increasingly by year
8   as.data.frame() %>% # Transforming data.table into data.frame
9   ggplot() +
10  aes(x = Year, y = N,
11      text = paste("Ano:", Year,
12                  "<br>Qtd. Acidentes:", round(N/1000, 1), "k",
13                  "<br>Disa Obs. no Ano:", Distinct_date,
14                  "<br>Acidentes por Dia:", round(Media_diaria/1000, 1), "k")) +
15  geom_bar(stat = "identity", fill = "darkred") +
16  scale_y_continuous(
17    labels = label_number(scale = 1/1000, suffix = "k") # Divide por 1000 e adiciona o sufixo
18  ) +
19  labs(title = "Qtd. de Acidentes por Ano", x = "Ano", y = "Qtd. Acidentes") +
20  facet_wrap(~Severity, scales = "free_y") +
21  theme_minimal() +
22  theme(plot.title = element_text(hjust = 0.5)) +
23  plotly()
24 acidentes_ano_barplot %>% ggplotly(tooltip = "text")

```

Aqui, podemos observar que a quantidade de acidentes é maior para a severidade 2, seguida pela 3, 4 e 1, respectivamente. Além disso, a quantidade de acidentes para a severidade 1 e 2 são as que mais crescem, ao passo que a quantidade de acidentes para a severidade 3 é a que mais decresce. Os acidentes de severidade 4, por outro lado, apresentam uma tendência de crescimento ao longo dos anos proporcionalmente menor. A explicação das causas da diminuição da severidade média não está no escopo desta análise, entretanto, podemos destacar que algumas hipóteses podem ser levantadas, como a melhoria da infraestrutura viária, a conscientização dos motoristas, a melhoria dos sistemas de segurança dos veículos, alteração de legislações de trânsito, entre outros.

Para fechar as análises de variáveis temporais, investigamos o padrão da média de severidade dos acidentes entre as faixas horárias e dias da semana, bem como ao longo dos anos:

```

1 acidentes_data_matriz_SWKTR <- acidentes_US_date[Civil_Twilight == "Day" | Civil_Twilight ==
2   setorder(Weekday, -Time_Range) %>% # order data set increasingly by weekday and time range
3   as.data.frame() %>% # Transforming data.table into data.frame
4   ggplot() +
5   aes(x = Weekday, y = Time_Range, fill = Severity) +
6   geom_tile() +
7   scale_fill_gradient(low = "white", high = "darkred") +
8   labs(title = "Agenda Diária da Severidade Média de Acidentes", x = "Dia da Semana", y = "H
9   theme_minimal() +

```

```

10 scale_x_discrete(position = "top") +
11 theme(
12   axis.title.x.top = element_text(),
13   axis.text.x.top = element_text(),
14   axis.line.x.top = element_line(),
15   axis.ticks.x.top = element_line()
16
17 )
18 acidentes_data_matriz_SWKTR

```

Destes, conseguimos observar que, mesmo que os acidentes ocorram menos durante os finais de semana, são eles os mais severos. Além disso, a severidade média dos acidentes tem diminuído ao longo dos anos!

### 0.3.2 Resultados para as Variáveis de Localização

A maioria dos acidentes ocorre em cidades maiores e mais populosas, bem como em estados com maior densidade demográfica. Isso pode ser atribuído ao maior volume de tráfego, maior número de veículos em circulação, e maior complexidade das vias urbanas nessas regiões. Além disso, áreas urbanas costumam apresentar condições mais propensas a acidentes, como interseções movimentadas, trânsito intenso e uma maior frequência de eventos climáticos adversos que afetam a visibilidade e a dirigibilidade. Fatores socioeconômicos também podem desempenhar um papel, já que cidades maiores tendem a concentrar atividades econômicas e culturais, atraindo um número significativo de pessoas e aumentando o risco de incidentes

```

1 leaflet() %>%
2   addTiles() %>%
3   addPolygons(
4     data = acidentes_p_estado,
5     fillColor = ~palestado(n),
6     weight = 2,
7     color = "black",
8     fillOpacity = 0.4,
9     label = ~NAME,
10    popup = ~glue("<b>Estado: </b> {NAME}<br><b>Quantidade de acidentes: </b> {n}"),
11    highlightOptions = highlightOptions(
12      weight = 3,
13      color = "#666",
14      fillOpacity = 0.6,
15      bringToFront = TRUE
16    ),

```

```

17     group = "Nível estadual"
18   ) %>%
19   addLegend(
20     "bottomright",
21     pal = palestado,
22     values = acidentes_p_estado$n,
23     title = "Número de Acidentes",
24     opacity = 1,
25     group = "Nível estadual"
26   ) %>%
27   addHeatmap(
28     data = mapa_de_calor,
29     lng = ~`Start Lng`,
30     lat = ~`Start Lat`,
31     intensity = ~n,
32     blur = 20,
33     radius = 15,
34     group = "Mapa de calor"
35   ) %>%
36   addLayersControl(
37     baseGroups = c("Nível estadual", "Mapa de calor"),
38     position = "topright"
39   )

```

### 0.3.3 Resultados para as Variáveis de Condições Climáticas

A análise das condições climáticas revelou que a maioria dos acidentes ocorreu em dias com céu limpo. A correlação entre variáveis climáticas, como temperatura, umidade e pressão, mostrou interações tanto positivas quanto negativas entre elas. A frequência de acidentes variou conforme as condições climáticas, com maior incidência em dias claros.

A distribuição de precipitação indicou que a maioria dos acidentes ocorreu sem chuva, e a visibilidade também foi um fator importante. Quando analisada a relação entre precipitação e severidade do acidente, não foi observada uma correlação clara, sugerindo que outros fatores influenciam a gravidade dos acidentes.

### 0.3.4 Resultados para as Variáveis de Sinalização e Infraestrutura

```

1   poi_Severity <- acidentes_US_sample[,c(3,30:42)]
2

```

```

3   poi_sum <- poi_Severity[, Total := rowSums(.SD),
4                                   .SDcols = c("Amenity", "Bump", "Crossing", "Give_Way", "Junction",
5                                               "No_Exit", "Railway", "Roundabout", "Station", "Stop",
6                                               "Traffic_Calming", "Traffic_Signal", "Turning_Loop")]
7
8   severityPoi_sum_table <- poi_sum[, .(Severity, Total)] %>% table
9   severityPoi_sum_df <- as.data.frame(severityPoi_sum_table)
10
11  severity_sumPoi <- severityPoi_sum_df %>% summarise(.by = Severity, sum_Poi = sum(Freq))
12  severityPoi_sum_df <- severityPoi_sum_df %>%
13    left_join(severity_sumPoi) %>% mutate(Rel_Freq = Freq/sum_Poi)
14
15
16  absolute_severityBars <- severityPoi_sum_df %>% as.data.frame %>%
17    ggplot(aes(x = Severity, y=Freq, fill = Total)) +
18    geom_bar(stat = "identity", position = "dodge") +
19    scale_y_continuous(labels = comma) +
20    labs(title = "Severidade de acidentes por ponto de interesse",
21         x = NULL,
22         y = NULL) +
23    theme_minimal()
24  absolute_severityBars

```

O gráfico de frequência absoluta da severidade por soma de pontos de interesse nos mostra que a maior quantidade de pontos ao redor de um acidente é 7, e também que a maioria dos acidentes ocorre em locais sem pontos de interesse, e ocorrem com frequência reduzida quanto mais pontos de interesse existem ao redor, porém pode apenas indicar que existem poucos lugares com muitos pontos de interesse. Também é que evidente que acidentes com grau 2 de severidade são extremamente mais comuns que os outros.

### 0.3.5 Resultados para as Variáveis de Duração do acidente

Os dados possuem muitos outliers e acidentes que duraram até 2 milhões de minutos, o que dificultou a nossa análise. Porém, de acordo com o gráfico mostrado nas técnicas estatísticas analisadas, podemos observar que os acidentes de severidade 2 e 4 têm uma variabilidade muito maior em questão do tempo. Podemos observar também, que a maior parte dos acidentes de severidade 1 e 3, costumam demorar entre 0 e 100 minutos.



## Referências

- Moosavi, S., Samavatian, M. H., Parthasarathy, S., & Ramnath, R. (2019). *A Countrywide Traffic Accident Dataset*. <https://arxiv.org/abs/1906.05409>
- Moosavi, S., Samavatian, M. H., Parthasarathy, S., Teodorescu, R., & Ramnath, R. (2019). Accident Risk Prediction based on Heterogeneous Sparse Data: New Dataset and Insights. *Proceedings of the 27th ACM SIGSPATIAL International Conference on Advances in Geographic Information Systems*, 33–42. <https://doi.org/10.1145/3347146.3359078>