

MIDTERM 615

Yang Xiao

2023-11-03

1. Storm data cleaning an EDA

1.1 cleaning and merge two dataset

```
library(readr)
library(dplyr)
```

```
##  
## 载入程辑包: 'dplyr'
```

```
## The following objects are masked from 'package:stats':  
##  
##     filter, lag
```

```
## The following objects are masked from 'package:base':  
##  
##     intersect, setdiff, setequal, union
```

```
library(ggplot2)
```

```
# Read the data for 2020 and 2021
storm_data_2020 <- read_csv("StormEvents_details-ftp_v1.0_d2020_c20230927.csv")
```

```
## Rows: 61279 Columns: 51
```

```
## ---- Column specification ----  
  
## Delimiter: ","
## chr (26): STATE, MONTH_NAME, EVENT_TYPE, CZ_TYPE, CZ_NAME, WFO, BEGIN_DATE_T...
## dbl (25): BEGIN_YEARMONTH, BEGIN_DAY, BEGIN_TIME, END_YEARMONTH, END_DAY, EN...
##  
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
storm_data_2021 <- read_csv("StormEvents_details-ftp_v1.0_d2021_c20231017.csv")
```

```
## Rows: 61389 Columns: 51
## └─ Column specification ──────────────────────────────────────────
## 
##   ┌─ Delimiter: ","
##   ┌─ chr (26): STATE, MONTH_NAME, EVENT_TYPE, CZ_TYPE, CZ_NAME, WFO, BEGIN_DATE_T...
##   ┌─ dbl (25): BEGIN_YEARMONTH, BEGIN_DAY, BEGIN_TIME, END_YEARMONTH, END_DAY, EN...
##   └─
##   ┌─ i Use `spec()` to retrieve the full column specification for this data.
##   ┌─ i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
# Filter for flood-related events
flood_events_2020 <- filter(storm_data_2020, EVENT_TYPE %in% c('Flash Flood', 'Flood', 'Coastal
Flood', 'Debris Flow', 'Lakeshore Flood'))
flood_events_2021 <- filter(storm_data_2021, EVENT_TYPE %in% c('Flash Flood', 'Flood', 'Coastal
Flood', 'Debris Flow', 'Lakeshore Flood'))

# Handling missing values for 2020 and 2021
columns_to_drop <- c('MAGNITUDE', 'MAGNITUDE_TYPE', 'CATEGORY', 'TOR_F_SCALE', 'TOR_LENGTH', 'T
OR_WIDTH', 'TOR_OTHER_WFO', 'TOR_OTHER_CZ_STATE', 'TOR_OTHER_CZ_FIPS', 'TOR_OTHER_CZ_NAME')
flood_events_2020_cleaned <- flood_events_2020[, !names(flood_events_2020) %in% columns_to_dro
p]
flood_events_2021_cleaned <- flood_events_2021[, !(names(flood_events_2021) %in% columns_to_dro
p)]

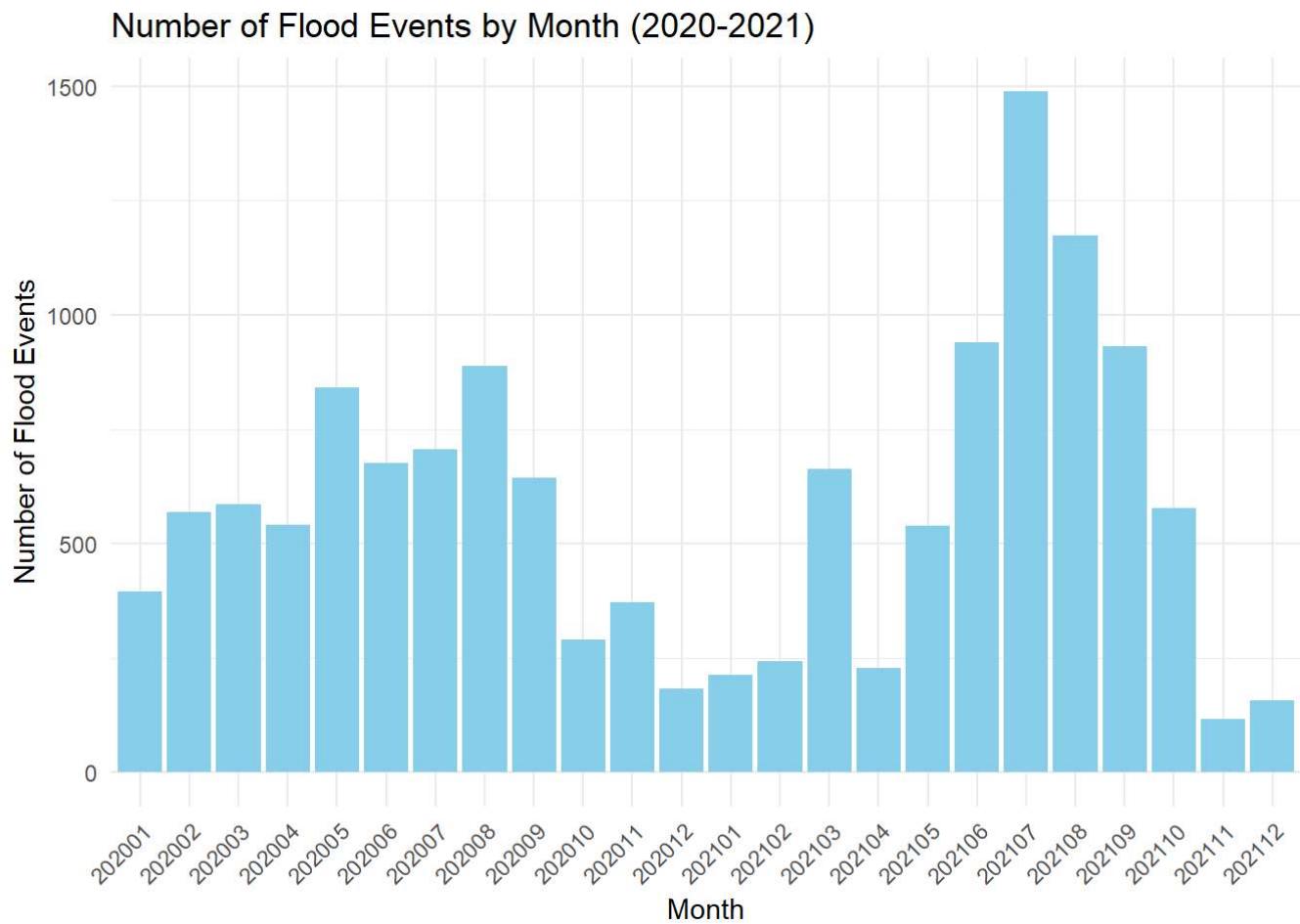
# Remove duplicates
flood_events_2020_cleaned <- distinct(flood_events_2020_cleaned)
flood_events_2021_cleaned <- distinct(flood_events_2021_cleaned)

# Merge the datasets from 2020 and 2021
flood_events_combined <- bind_rows(flood_events_2020_cleaned, flood_events_2021_cleaned)

# Extract year and month from 'BEGIN_YEARMONTH'
flood_events_combined$YEAR_MONTH <- substr(flood_events_combined$BEGIN_YEARMONTH, 1, 6)

# Count the number of events for each month and plot
monthly_events <- flood_events_combined %>%
  group_by(YEAR_MONTH) %>%
  summarise(Count = n()) %>%
  arrange(YEAR_MONTH)

ggplot(monthly_events, aes(x = YEAR_MONTH, y = Count)) +
  geom_bar(stat = "identity", fill = "skyblue") +
  theme_minimal() +
  labs(title = 'Number of Flood Events by Month (2020-2021)', x = 'Month', y = 'Number of Flood
Events') +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```



Most of the floods are concentrated in: 2020.01-2020.09 and 2021.05-2021.09.

1.2 Initial question

Is the Property damage and Duration of floods related to different continents?

- a. Property damage vs Duration
- b. Property damage vs State

1.3 EDA and Solution

(a) Property damage vs Duration

```
library(tidyverse)
```

```
## └─ Attaching core tidyverse packages ─────────────────────────────────── tidy
verse 2.0.0 ─
## ✓forcats    1.0.0      ✓stringr    1.5.0
## ✓lubridate  1.9.2      ✓tibble     3.2.1
## ✓purrr     1.0.2      ✓tidyverse  1.3.0
## └─ Conflicts ───────────────────────────────────
───────── tidyverse_conflicts() ─
## ✘dplyr::filter() masks stats::filter()
## ✘dplyr::lag()   masks stats::lag()
## ┌ Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to beco
me errors
```

```
library(readr)
library(dplyr)
library(ggplot2)
library(scales)
```

```
##
## 载入程辑包: 'scales'
##
## The following object is masked from 'package:purrr':
##   discard
##
## The following object is masked from 'package:readr':
##   col_factor
```

```
data<-flood_events_combined
# Function to convert yearmonthday and time to a datetime object
convert_to_datetime <- function(yearmonth, day, time) {
  year <- as.integer(yearmonth / 100)
  month <- yearmonth %% 100
  hour <- as.integer(time / 100)
  minute <- time %% 100
  as.POSIXct(paste(year, month, day, hour, minute), format="%Y %m %d %H %M")
}

# Apply the function to BEGIN and END columns to create datetime objects
data$BEGIN_DATETIME <- mapply(convert_to_datetime, data$BEGIN_YEARMONTH, data$BEGIN_DAY, data$BEGIN_TIME)
data$END_DATETIME <- mapply(convert_to_datetime, data$END_YEARMONTH, data$END_DAY, data$END_TIME)

# Calculate the duration in hours
data$DURATION_HOURS <- as.numeric(difftime(data$END_DATETIME, data$BEGIN_DATETIME, units="hours"))

# Selecting the specific columns
data <- data[, c("DURATION_HOURS", "DAMAGE_PROPERTY", "STATE")]

# Display the first few rows of the selected data
head(data)
```

```
## # A tibble: 6 × 3
##   DURATION_HOURS DAMAGE_PROPERTY STATE
##       <dbl> <chr>           <chr>
## 1         3    3.00K        WEST VIRGINIA
## 2         3    5.00K        VIRGINIA
## 3         1    0.00K        OHIO
## 4         1    0.00K        OHIO
## 5         1    0.00K        OHIO
## 6        10.8  2.00K        OHIO
```

```
# Function to convert damage property values to numeric
convert_damage <- function(damage) {
  if (grepl("K", damage)) {
    return(as.numeric(sub("K", "", damage)) * 1e3)
  } else if (grepl("M", damage)) {
    return(as.numeric(sub("M", "", damage)) * 1e6)
  } else if (grepl("B", damage)) {
    return(as.numeric(sub("B", "", damage)) * 1e9)
  } else {
    return(as.numeric(damage))
  }
}

# Apply the conversion function to the DAMAGE_PROPERTY column
data$DAMAGE_PROPERTY <- sapply(data$DAMAGE_PROPERTY, convert_damage)

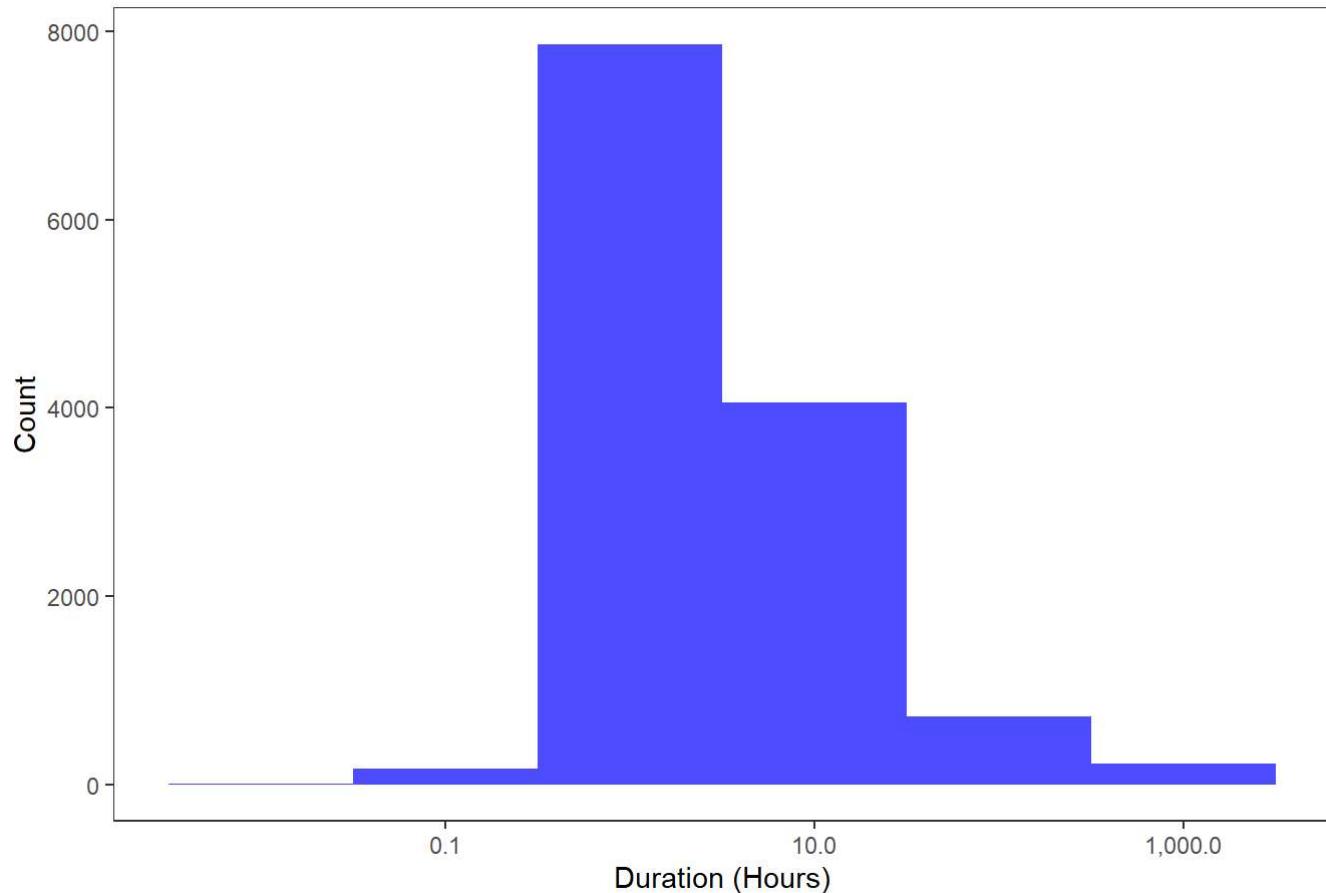
# Customize the theme for the plots
theme_set(theme_bw() + theme(panel.grid.major = element_blank(),
                             panel.grid.minor = element_blank(),
                             plot.title = element_text(hjust = 0.5)))
```

```
# Histogram and Boxplot of DURATION_HOURS
ggplot(data, aes(x=DURATION_HOURS)) +
  geom_histogram(binwidth = 1, fill = "blue", alpha = 0.7) +
  scale_x_log10(labels = scales::comma) +
  labs(title = "Histogram of DURATION_HOURS") +
  xlab("Duration (Hours)") +
  ylab("Count")
```

```
## Warning: Transformation introduced infinite values in continuous x-axis
```

```
## Warning: Removed 938 rows containing non-finite values (`stat_bin()`).
```

Histogram of DURATION_HOURS



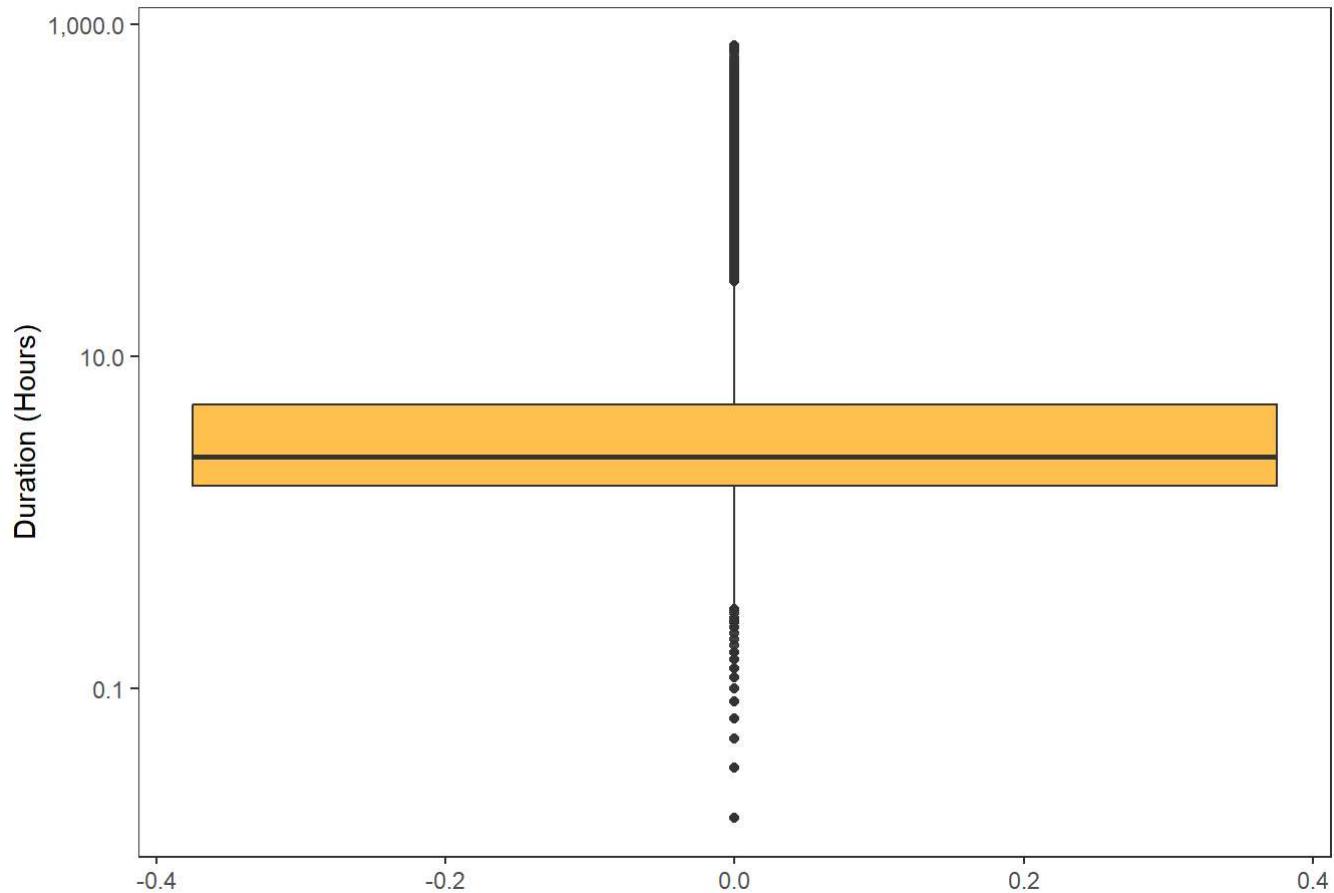
DURATION_HOURS shows a highly right-skewed distribution, with a large number of events having short durations and a few events with very long durations.

```
ggplot(data, aes(y=DURATION_HOURS)) +  
  geom_boxplot(fill = "orange", alpha = 0.7) +  
  scale_y_log10(labels = scales::comma) +  
  labs(title = "Boxplot of DURATION_HOURS") +  
  ylab("Duration (Hours)")
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Removed 938 rows containing non-finite values (`stat_boxplot()`).
```

Boxplot of DURATION_HOURS



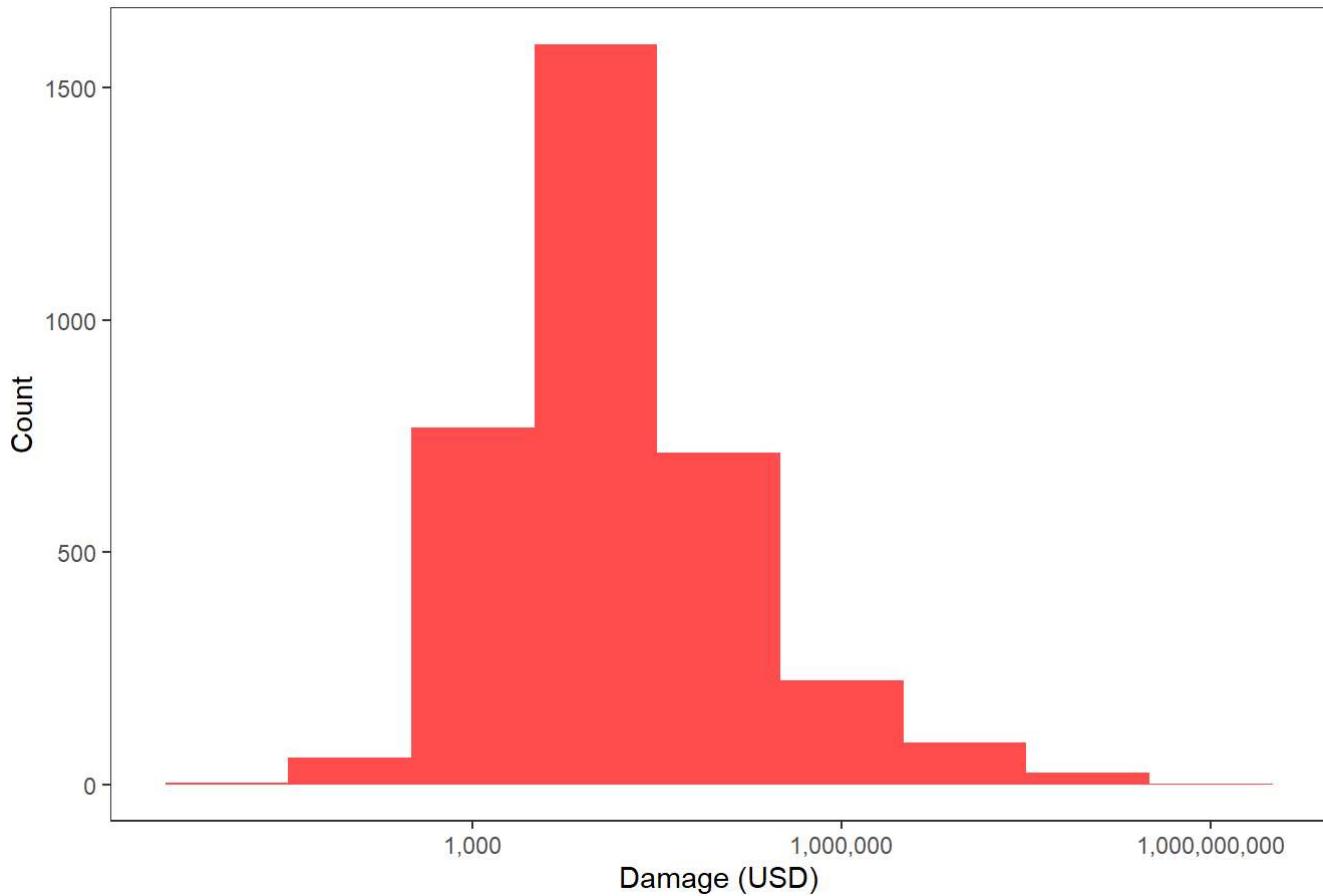
The boxplot for DURATION_HOURS reveals a significant number of outliers beyond the upper whisker, indicating some events with exceptionally long durations.

```
# Histogram and Boxplot of DAMAGE_PROPERTY
ggplot(data, aes(x=DAMAGE_PROPERTY)) +
  geom_histogram(binwidth = 1, fill = "red", alpha = 0.7) +
  scale_x_log10(labels = scales::comma) +
  labs(title = "Histogram of DAMAGE_PROPERTY") +
  xlab("Damage (USD)") +
  ylab("Count")
```

```
## Warning: Transformation introduced infinite values in continuous x-axis
```

```
## Warning: Removed 10490 rows containing non-finite values (`stat_bin()`).
```

Histogram of DAMAGE_PROPERTY



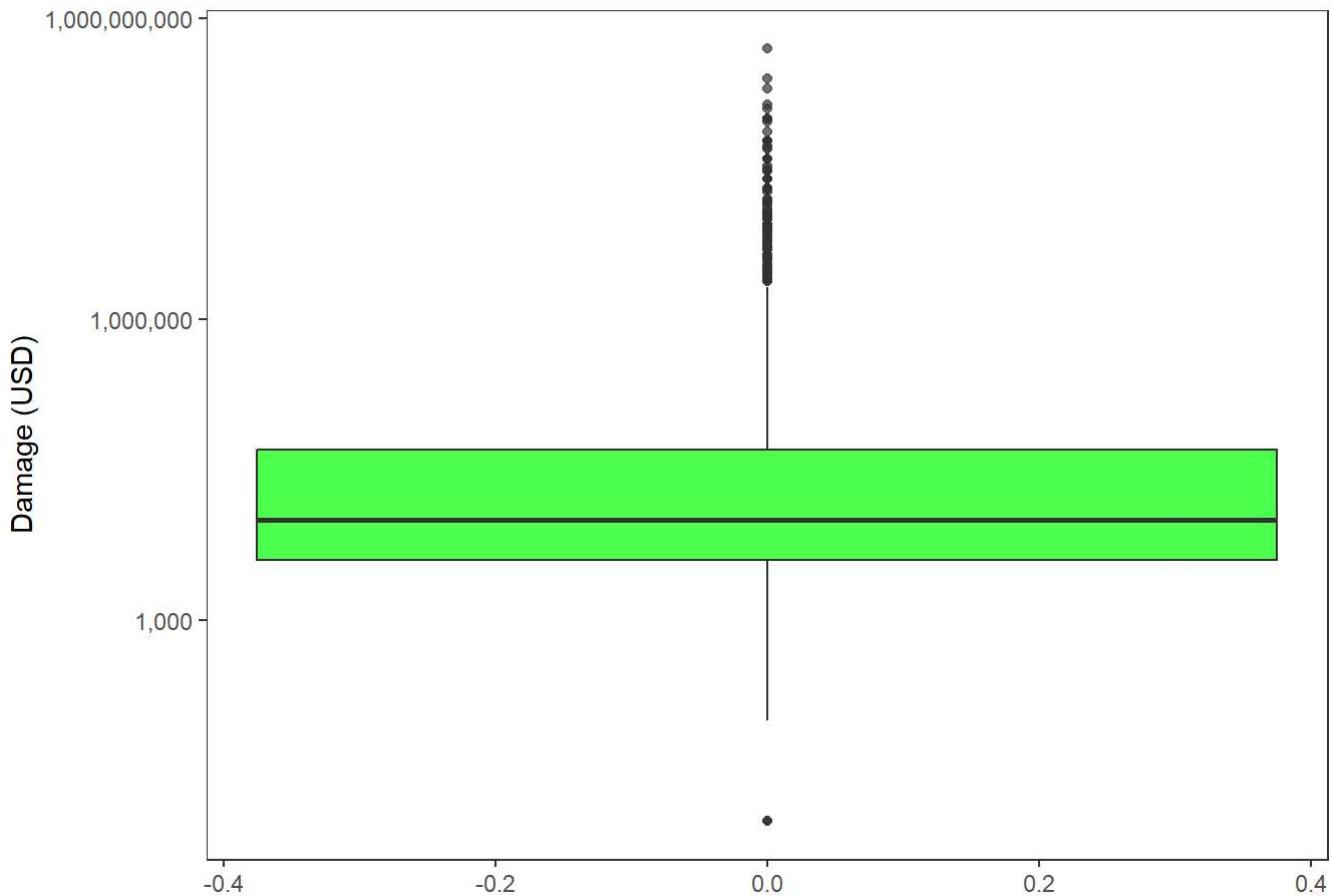
DAMAGE_PROPERTY also shows a highly right-skewed distribution. Most property damage amounts are at the lower end of the scale, with a few instances of very high damage.

```
ggplot(data, aes(y=DAMAGE_PROPERTY)) +  
  geom_boxplot(fill = "green", alpha = 0.7) +  
  scale_y_log10(labels = scales::comma) +  
  labs(title = "Boxplot of DAMAGE_PROPERTY") +  
  ylab("Damage (USD)")
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Removed 10490 rows containing non-finite values (`stat_boxplot()`).
```

Boxplot of DAMAGE_PROPERTY



The boxplot for DAMAGE_PROPERTY on a logarithmic scale indicates many outliers, with property damage varying significantly from one event to another.

(b) Property damage vs State

```
# Compute the correlation coefficient
correlation <- cor(data$DURATION_HOURS, data$DAMAGE_PROPERTY, use = "complete.obs")

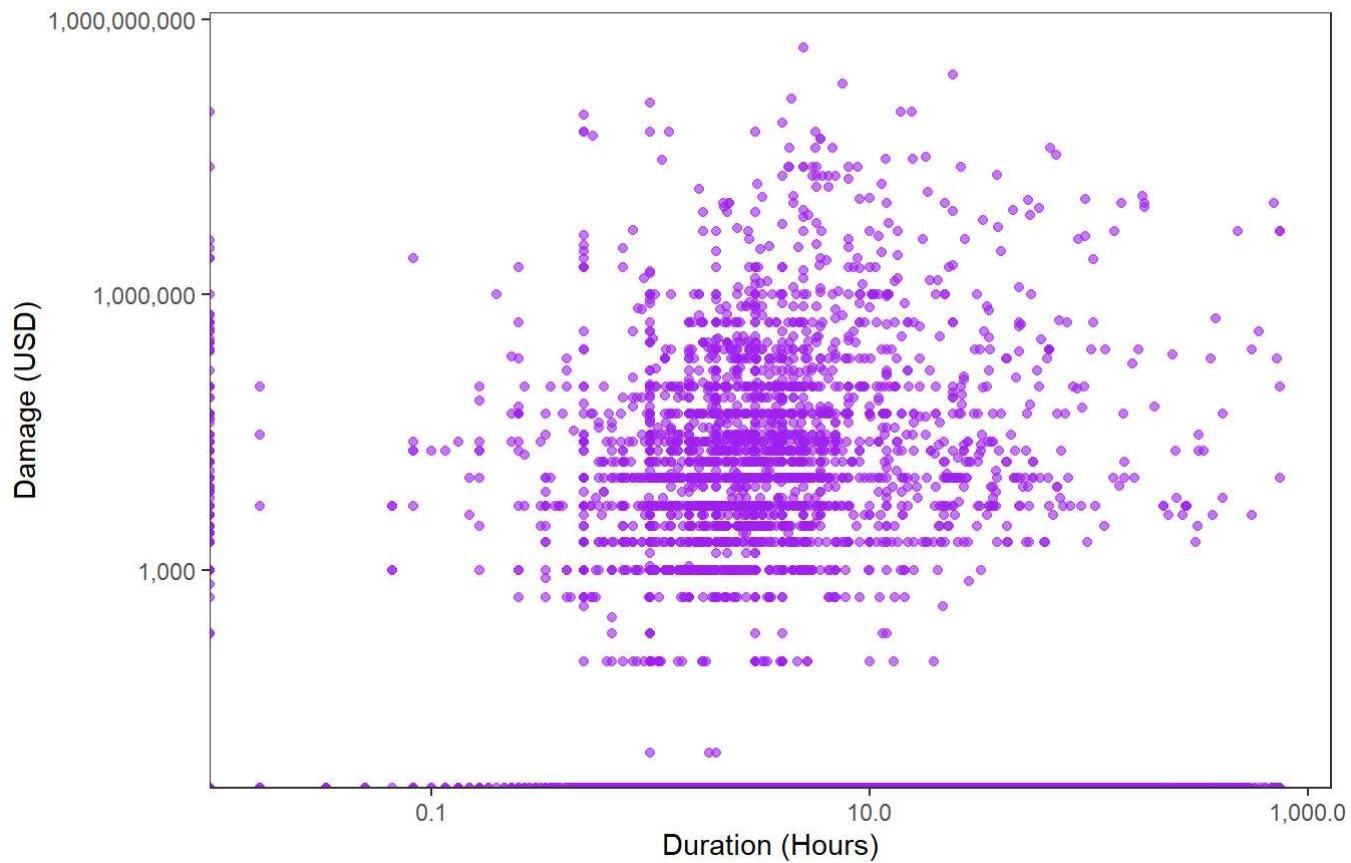
# Scatter plot of DURATION_HOURS vs DAMAGE_PROPERTY
ggplot(data, aes(x=DURATION_HOURS, y=DAMAGE_PROPERTY)) +
  geom_point(alpha = 0.6, color = "purple") +
  scale_x_log10(labels = scales::comma) +
  scale_y_log10(labels = scales::comma) +
  labs(title = paste("Scatter Plot of DURATION_HOURS vs DAMAGE_PROPERTY",
                     "\nPearson Correlation:", round(correlation, 2))) +
  xlab("Duration (Hours)") +
  ylab("Damage (USD)")
```

```
## Warning: Transformation introduced infinite values in continuous x-axis
```

```
## Warning: Transformation introduced infinite values in continuous y-axis
```

```
## Warning: Removed 130 rows containing missing values (`geom_point()`).
```

Scatter Plot of DURATION_HOURS vs DAMAGE_PROPERTY Pearson Correlation: 0



```
# Print the correlation coefficient
print(paste("Pearson Correlation Coefficient:", round(correlation, 4)))
```

```
## [1] "Pearson Correlation Coefficient: -9e-04"
```

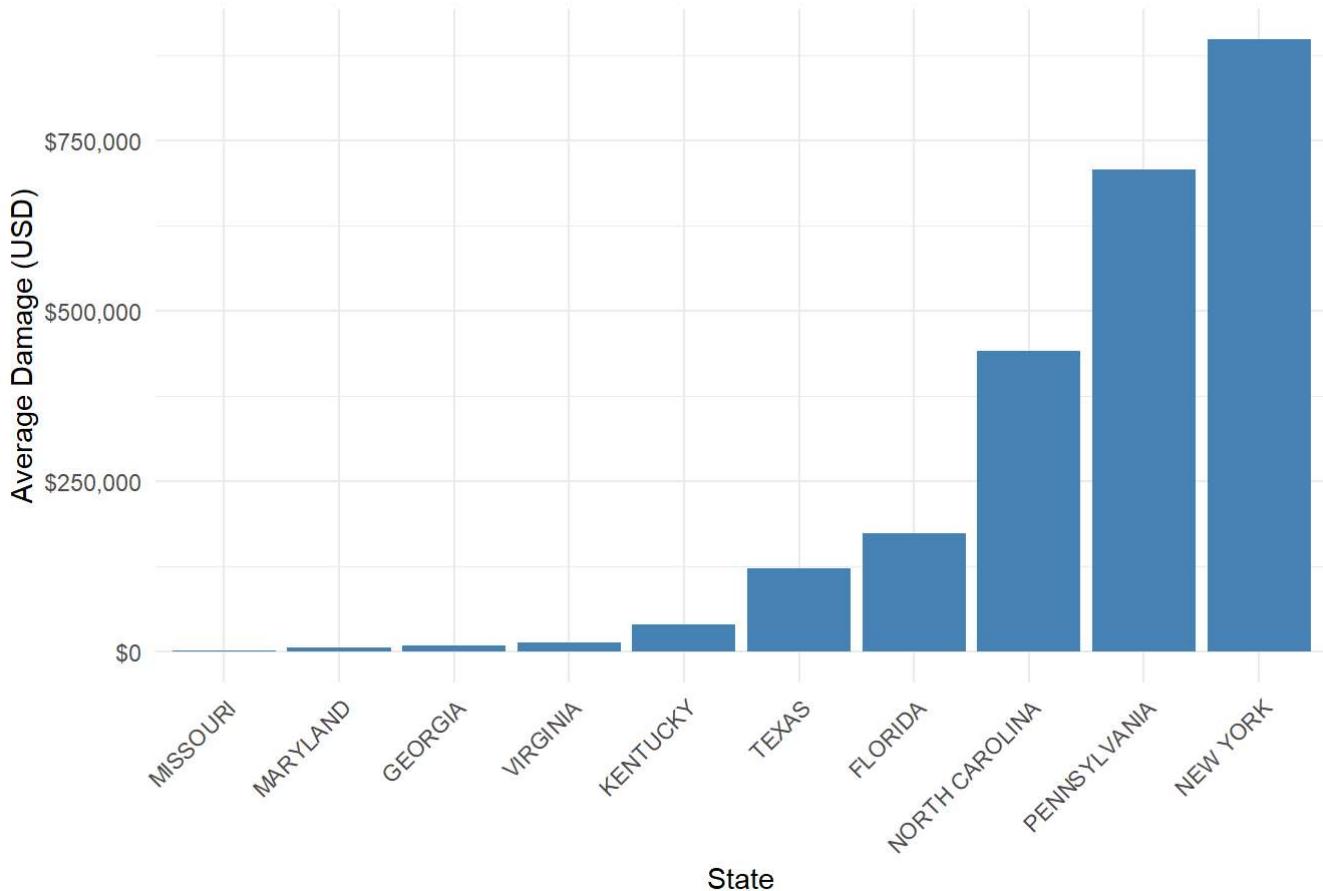
The scatter plot, with both axes on a logarithmic scale due to the wide range of values, also suggests that there is no clear pattern or relationship between the two variables. The points are widely dispersed, and no trend is visible.

```
# Group by STATE and summarize counts and average damage
state_summary <- data %>%
  group_by(STATE) %>%
  summarise(Count = n(),
            Average_Damage = mean(DAMAGE_PROPERTY, na.rm = TRUE)) %>%
  ungroup()

# Identify the top 10 states with the highest counts
top_states <- state_summary %>%
  arrange(desc(Count)) %>%
  slice_head(n = 10)

# Create a bar plot for the average DAMAGE_PROPERTY of the top states
ggplot(top_states, aes(x=reorder(STATE, Average_Damage), y=Average_Damage)) +
  geom_bar(stat="identity", fill="steelblue") +
  scale_y_continuous(labels = scales::dollar) +
  labs(title = "Average Property Damage by State (Top 10 by Count)",
       x = "State",
       y = "Average Damage (USD)") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1))
```

Average Property Damage by State (Top 10 by Count)



```
print(top_states)
```

```
## # A tibble: 10 × 3
##   STATE      Count Average_Damage
##   <chr>     <int>      <dbl>
## 1 VIRGINIA    1052     12950.
## 2 MISSOURI     797      1288.
## 3 TEXAS        751     122508.
## 4 KENTUCKY      714     40396.
## 5 NEW YORK     612     897879.
## 6 GEORGIA       529      9770.
## 7 NORTH CAROLINA 529    441778.
## 8 MARYLAND      527      6401.
## 9 PENNSYLVANIA   508    707848.
## 10 FLORIDA      491    174396.
```

The states with significant losses per flood on average are: TEXAS , FLORIDA ,NORTH CAROLINA ,PENNSYLVANIA,NEW YORK.The losses in these states are quite significant, with New York at the most.

2. Flooding data cleaning an EDA

2.1 Cleaning and merge two dataset

```
library(tidyverse)
library(readr)
library(lubridate)
library(ggplot2)

disaster_declarations_df <- read_csv('DisasterDeclarationsSummaries.csv')
```

```
## Rows: 64950 Columns: 25
## └─ Column specification ──────────────────────────────────────────────────
## 
## #> Delimiter: ","
## #> chr (10): femaDeclarationString, state, declarationType, incidentType, decl...
## #> dbl (9): disasterNumber, fyDeclared, ihProgramDeclared, iaProgramDeclared, ...
## #> dttm (6): declarationDate, incidentBeginDate, incidentEndDate, disasterClos...
## #>
## #> ┌ Use `spec()` to retrieve the full column specification for this data.
## #> ┌ Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
fema_web_summaries_df <- read_csv('FemaWebDisasterSummaries.csv')
```

```
## Rows: 3588 Columns: 14
## └─ Column specification ──────────────────────────────────────────
## 
##   ## Delimiter: ","
##   ## chr (2): hash, id
##   ## dbl (9): disasterNumber, totalNumberIaApproved, totalAmountIhpApproved, tot...
##   ## dttm (3): paLoadDate, iaLoadDate, lastRefresh
## 
##   ## ┌─ Use `spec()` to retrieve the full column specification for this data.
##   ## ┌─ Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```
flood_declarations_df <- disaster_declarations_df %>%
  filter(incidentType == 'Flood', fyDeclared %in% c(2020, 2021))

flood_disaster_numbers <- unique(flood_declarations_df$disasterNumber)

## Filter the FEMA web summaries data for the corresponding disaster numbers
flood_financials_df <- fema_web_summaries_df %>%
  filter(disasterNumber %in% flood_disaster_numbers)

## Merge the two datasets on the 'disasterNumber' column
combined_flood_data_df <- inner_join(flood_declarations_df, flood_financials_df, by = 'disasterNumber')

## Group by 'disasterNumber' and aggregate the financial data
combined_flood_data_aggregated <- combined_flood_data_df %>%
  group_by(disasterNumber) %>%
  summarise(
    totalNumberIaApproved = sum(totalNumberIaApproved, na.rm = TRUE),
    totalAmountIhpApproved = sum(totalAmountIhpApproved, na.rm = TRUE),
    totalAmountHaApproved = sum(totalAmountHaApproved, na.rm = TRUE),
    totalAmountOnaApproved = sum(totalAmountOnaApproved, na.rm = TRUE),
    totalObligatedAmountPa = sum(totalObligatedAmountPa, na.rm = TRUE),
    totalObligatedAmountCatAb = sum(totalObligatedAmountCatAb, na.rm = TRUE),
    totalObligatedAmountCatC2g = sum(totalObligatedAmountCatC2g, na.rm = TRUE),
    totalObligatedAmountHmgp = sum(totalObligatedAmountHmgp, na.rm = TRUE)
  ) %>%
  ungroup()

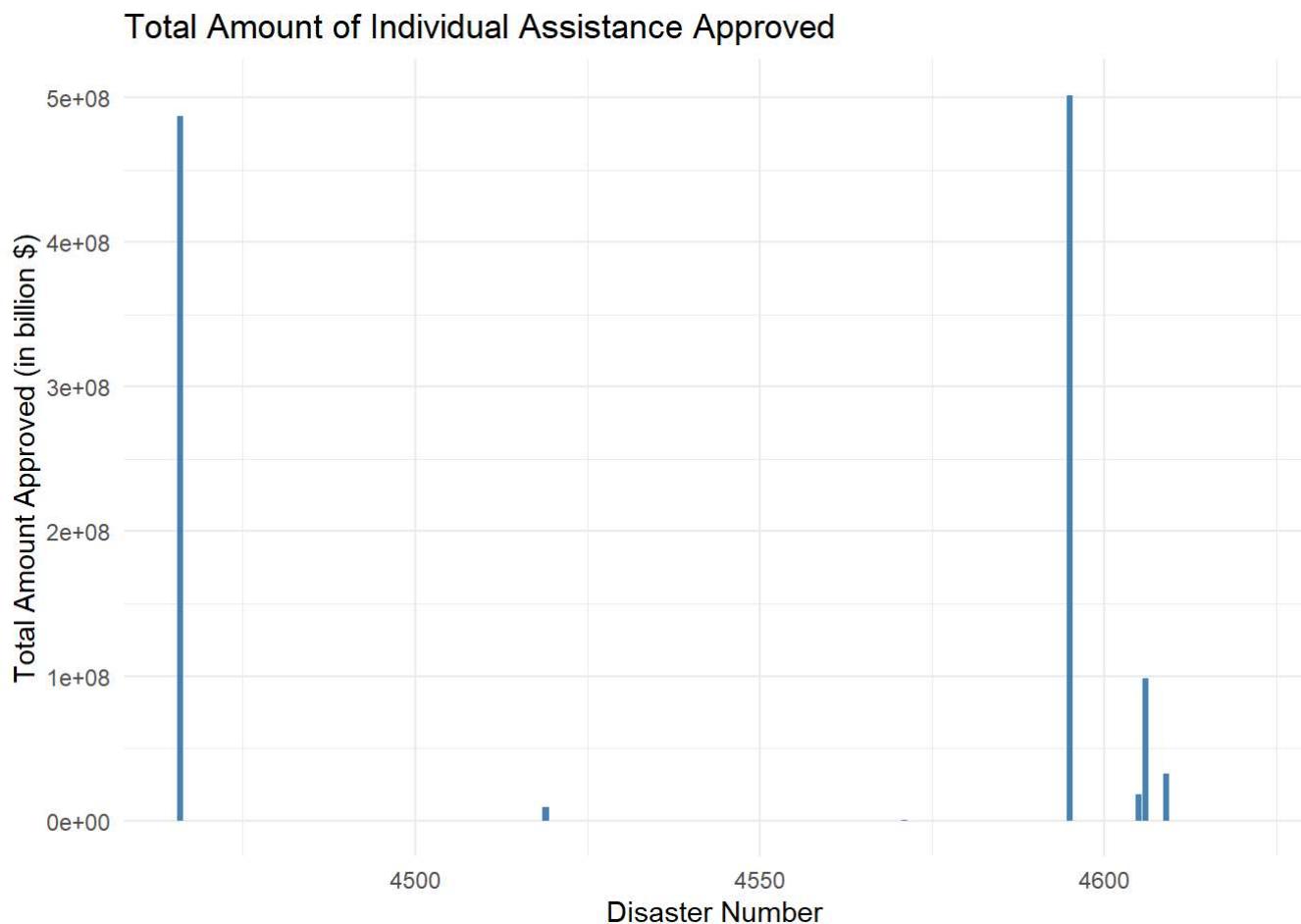
head(combined_flood_data_aggregated)
```

```
## # A tibble: 6 × 9
##   disasterNumber totalNumberIaApproved totalAmountIhpApproved
##       <dbl>           <dbl>           <dbl>
## 1        4466          78554          487195694.
## 2        4475             0              0
## 3        4477             0              0
## 4        4519            660          8969013.
## 5        4539             0              0
## 6        4553             0              0
## # ℹ 6 more variables: totalAmountHaApproved <dbl>,
## #   totalAmountOnaApproved <dbl>, totalObligatedAmountPa <dbl>,
## #   totalObligatedAmountCatAb <dbl>, totalObligatedAmountCatC2g <dbl>,
## #   totalObligatedAmountHmgp <dbl>
```

2.2 Plotting for Visualization

(a) Individual Assistance (IA)

```
ggplot(combined_flood_data_aggregated, aes(x = disasterNumber, y = totalAmountIhpApproved)) +
  geom_bar(stat = "identity", fill = "steelblue") +
  theme_minimal() +
  labs(title = "Total Amount of Individual Assistance Approved", x = "Disaster Number", y = "Total Amount Approved (in billion $)")
```

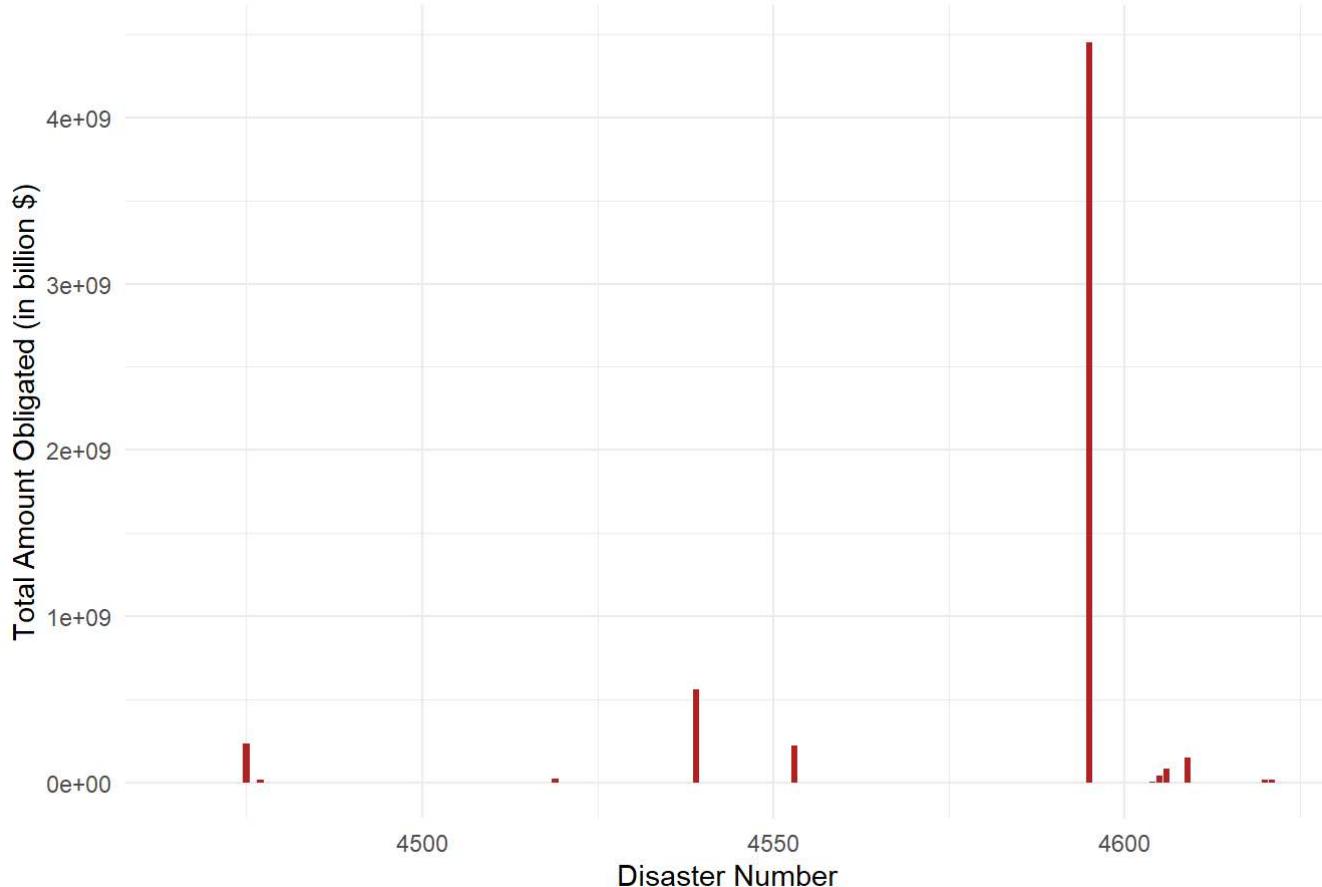


The first bar chart shows the total amount of Individual Assistance approved for each disaster.

(b) Public Assistance (PA)

```
ggplot(combined_flood_data_aggregated, aes(x = disasterNumber, y = totalObligatedAmountPa)) +
  geom_bar(stat = "identity", fill = "firebrick") +
  theme_minimal() +
  labs(title = "Total Obligated Amount for Public Assistance", x = "Disaster Number", y = "Total Amount Obligated (in billion $)")
```

Total Obligated Amount for Public Assistance

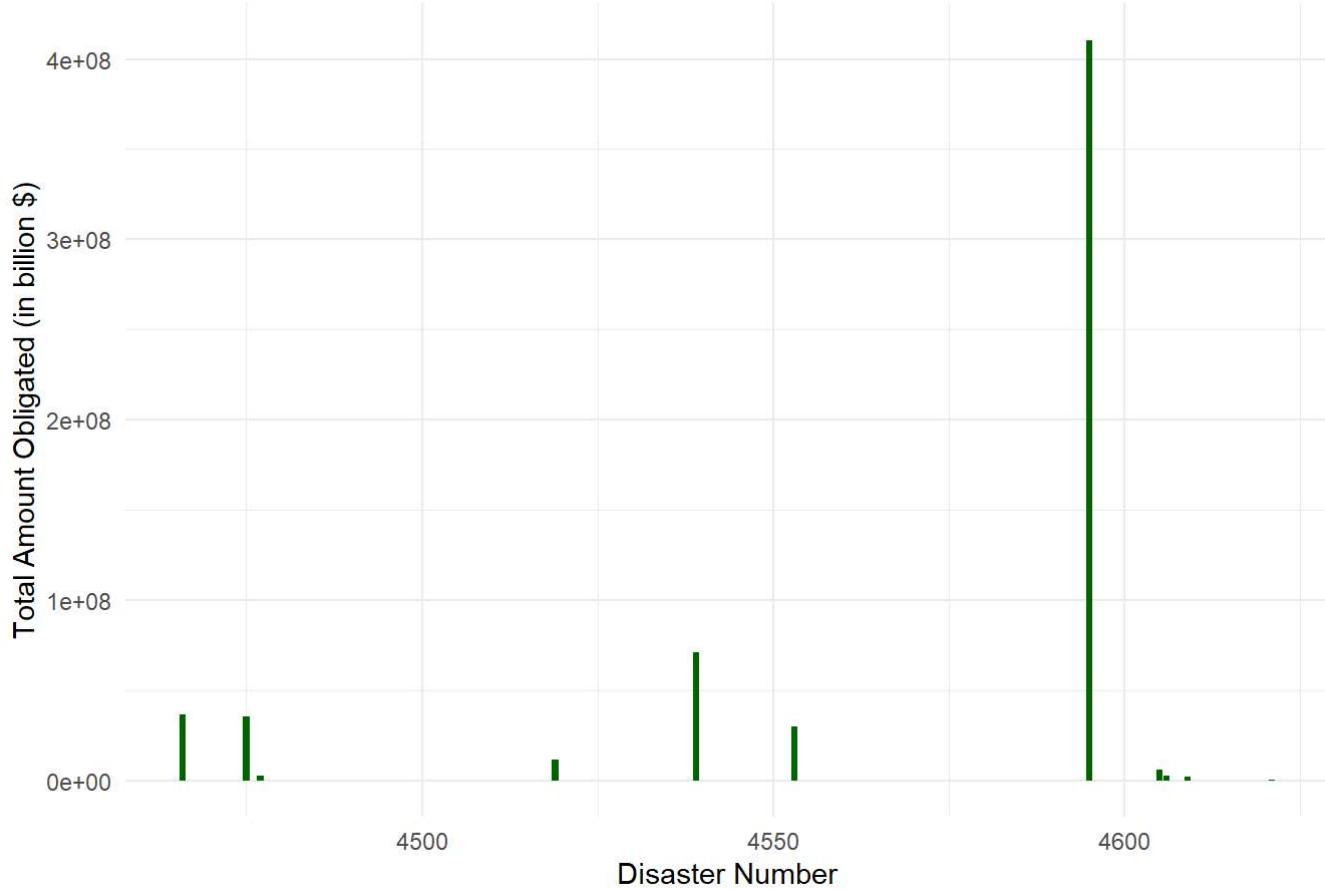


The second chart illustrates the total obligated amount for Public Assistance for each disaster, also in billions of dollars.

(c) Hazard Mitigation Grant Program (HMGP)

```
ggplot(combined_flood_data_aggregated, aes(x = disasterNumber, y = totalObligatedAmountHmgp)) +
  geom_bar(stat = "identity", fill = "darkgreen") +
  theme_minimal() +
  labs(title = "Total Obligated Amount for Hazard Mitigation Grant Program", x = "Disaster Number", y = "Total Amount Obligated (in billion $)")
```

Total Obligated Amount for Hazard Mitigation Grant Program



The third chart presents the total obligated amount for the Hazard Mitigation Grant Program for each disaster, in billions of dollars.

From these charts, we can observe significant differences in the financial assistance provided by different disaster events. Some disasters have a much greater financial impact and require more assistance than others.

3. Explore Top 10 States by Disaster Counts recent 10 years

3.1 Initial question

- How does the number of disasters change over time?
- And what are the change curves of the ten continents that have received the most disasters in the past 10 years?
- What does this indicate or reflect?

3.2 EDA and Solution

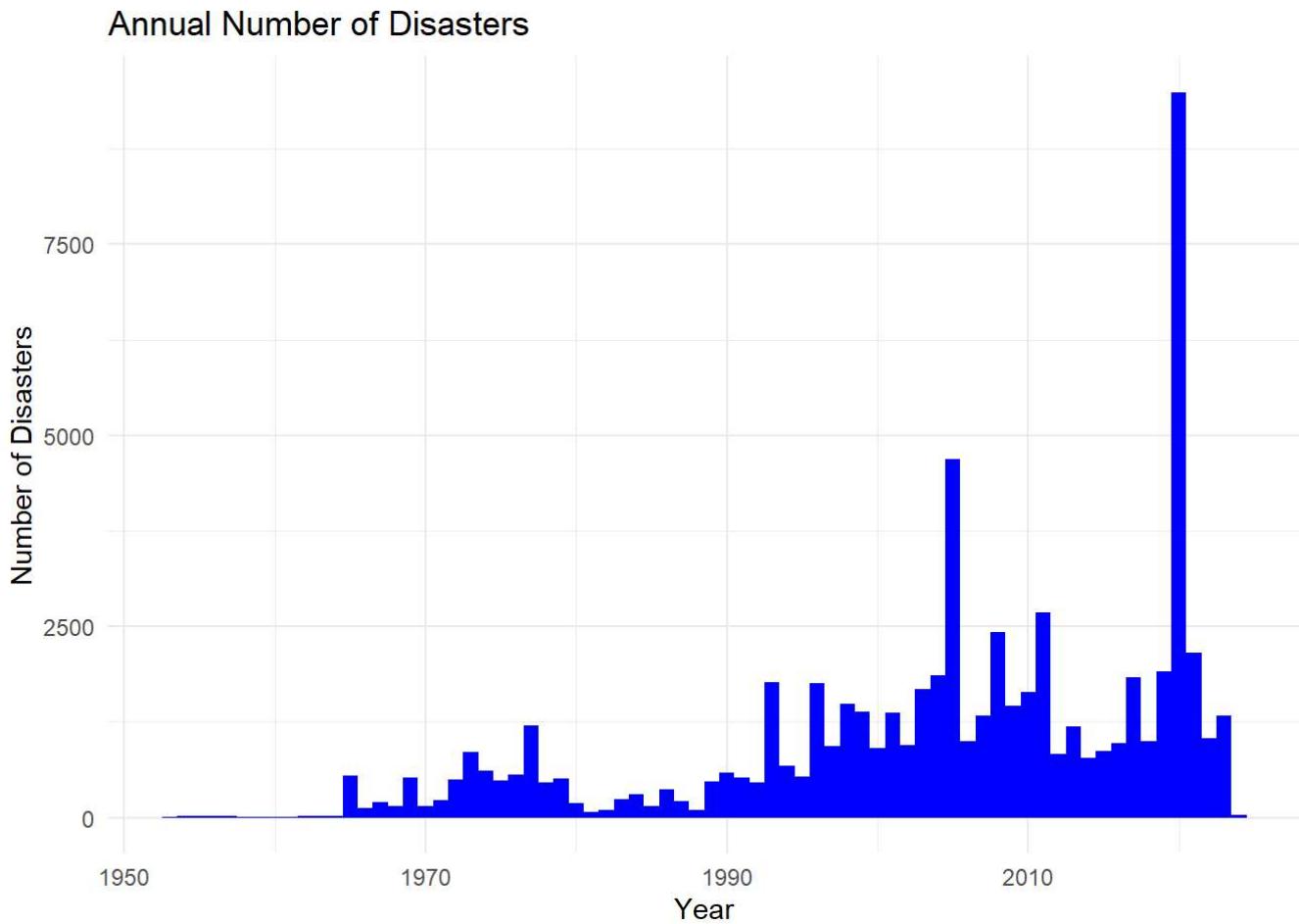
(a) Overall bar chart of disasters over time

```
# Descriptive statistics analysis
disaster_annual_summary <- disaster_declarations_df %>%
  group_by(fyDeclared) %>%
  summarise(
    totalDisasters = n(),
    disasterTypes = list(unique(incidentType))
  )

# Print the number and types of disasters for different years
print(disaster_annual_summary)
```

```
## # A tibble: 72 × 3
##   fyDeclared totalDisasters disasterTypes
##       <dbl>          <int>     <list>
## 1      1953            10 <chr [3]>
## 2      1954            14 <chr [5]>
## 3      1955            20 <chr [5]>
## 4      1956            18 <chr [5]>
## 5      1957            18 <chr [5]>
## 6      1958             5 <chr [2]>
## 7      1959             8 <chr [2]>
## 8      1960            13 <chr [7]>
## 9      1961            11 <chr [2]>
## 10     1962            16 <chr [2]>
## # ... 62 more rows
```

```
# plot for change of disasters per year
ggplot(disaster_declarations_df, aes(x = fyDeclared)) +
  geom_histogram(binwidth = 1, fill = "blue") +
  theme_minimal() +
  labs(title = "Annual Number of Disasters", x = "Year", y = "Number of Disasters")
```



According to the graph, it can be observed that the distribution of disasters is random and seems to have little to do with the year. However, compared to before and after 1950, the overall trend is still on the rise.

(b) Top 10 States

```
# Financial impact analysis - Distribution of financial assistance by year
financial_annual_summary <- combined_flood_data_aggregated %>%
  left_join(disaster_declarations_df %>% select(disasterNumber, fyDeclared), by = "disasterNumber") %>%
  group_by(fyDeclared) %>%
  summarise(
    totalIhpApproved = sum(totalAmountIhpApproved, na.rm = TRUE),
    totalPaObligated = sum(totalObligatedAmountPa, na.rm = TRUE),
    totalHmgpObligated = sum(totalObligatedAmountHmgp, na.rm = TRUE)
  )

print(financial_annual_summary)
```

```
## # A tibble: 2 × 4
##   fyDeclared  totalIhpApproved  totalPaObligated  totalHmgpObligated
##       <dbl>            <dbl>            <dbl>            <dbl>
## 1      2020        3446245911.     17056452040.    2601485703.
## 2      2021        26038765454.    224272604234.   20592840986.
```

```
disasters_2014_2024 <- disaster_declarations_df %>%
  filter(fyDeclared >= 2014, fyDeclared <= 2024)

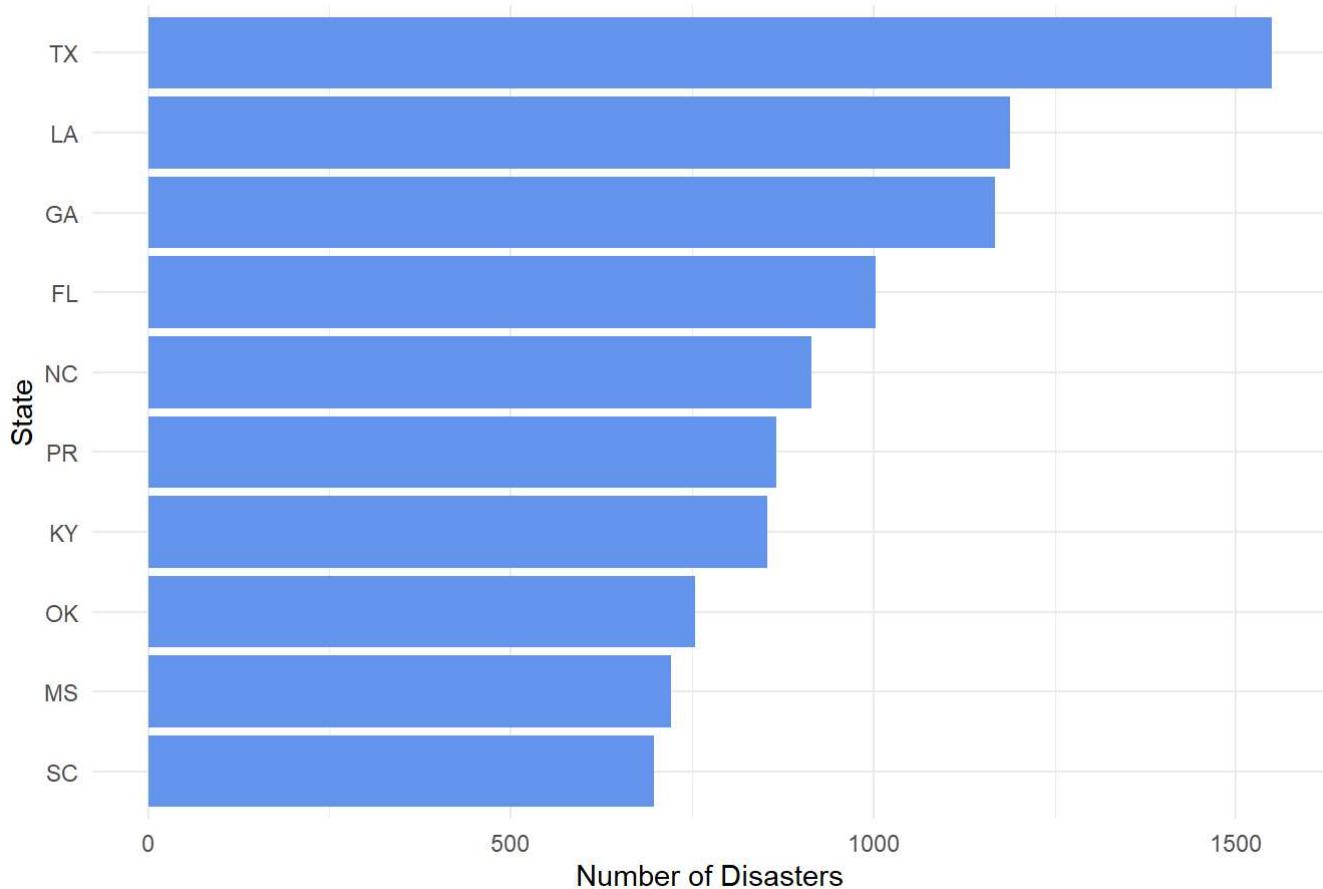
# Count the number of disasters by state
state_disaster_counts <- disasters_2014_2024 %>%
  count(state, sort = TRUE)

top10_states_disasters <- head(state_disaster_counts, 10)
top10_states_disasters
```

```
## # A tibble: 10 × 2
##   state     n
##   <chr> <int>
## 1 TX     1549
## 2 LA     1188
## 3 GA     1168
## 4 FL     1003
## 5 NC     914
## 6 PR     866
## 7 KY     854
## 8 OK     754
## 9 MS     721
## 10 SC    697
```

```
# Create a bar chart to display the number of disasters for the top 10 states
ggplot(top10_states_disasters, aes(x = reorder(state, n), y = n)) +
  geom_bar(stat = "identity", fill = "cornflowerblue") +
  coord_flip() + # For horizontal bars
  theme_minimal() +
  labs(title = "Top 10 States by Disaster Counts (2014-2024)", x = "State", y = "Number of Disasters")
```

Top 10 States by Disaster Counts (2014-2024)



We can obtain the ten continents with the highest number of disasters, which have been presented in the code and table. We can see that the continent with the most disasters in the past decade has been Texas(n=1549).

(c) Conclusion

This means that these ten continents have the highest number of disasters, and they should strengthen their disaster prevention measures and have increased their financial expenditure on disasters.

4. Population for whom poverty status is determined date cleaning and EDA

4.1 Data cleaning and merge the data from 2020 and 2021

```
# Load necessary libraries
library(readr)
library(dplyr)
library(ggplot2)

# Load the 2020 data
data_2020 <- read_csv("ACSST5Y2020.S1701-Data.csv")
```

```
## New names:
## Rows: 3222 Columns: 735
## --- Column specification
## -----
## ----- Delimiter: ","
## (734): GEO_ID, NAME, S1701_C01_001E, S1701_C01_001M, S1701_C01_001MA, S1... lgl
## (1): ...
## i Use `spec()` to retrieve the full column specification for this data. i
## Specify the column types or set `show_col_types = FALSE` to quiet this message.
## • ` ` -> `...`
```

```
# Select columns and rename for clarity
gender_poverty_2020 <- data_2020 %>%
  select(
    NAME,
    Total_Male_Population = S1701_C01_011E,
    Male_Population_Below_Poverty = S1701_C02_011E,
    Percent_Male_Population_Below_Poverty = S1701_C03_011E,
    Total_Female_Population = S1701_C01_012E,
    Female_Population_Below_Poverty = S1701_C02_012E,
    Percent_Female_Population_Below_Poverty = S1701_C03_012E
  ) %>%
  mutate(
    Total_Male_Population = as.numeric(Total_Male_Population),
    Male_Population_Below_Poverty = as.numeric(Male_Population_Below_Poverty),
    Total_Female_Population = as.numeric(Total_Female_Population),
    Female_Population_Below_Poverty = as.numeric(Female_Population_Below_Poverty)
  ) %>%
  mutate(
    Male_Poverty_Rate = Male_Population_Below_Poverty / Total_Male_Population,
    Female_Poverty_Rate = Female_Population_Below_Poverty / Total_Female_Population
  ) %>%
  na.omit()
```

```
## Warning: There were 4 warnings in `mutate()` .
## The first warning was:
## i In argument: `Total_Male_Population = as.numeric(Total_Male_Population)` .
## Caused by warning:
## ! 强制改变过程中产生了NA
## i Run `dplyr::last_dplyr_warnings()` to see the 3 remaining warnings.
```

```
data_2021 <- read_csv("ACSST5Y2021.S1701-Data.csv")
```

```
## New names:
## Rows: 3222 Columns: 747
## --- Column specification
## -----
## ----- Delimiter: ","
## (746): GEO_ID, NAME, S1701_C01_001E, S1701_C01_001EA, S1701_C01_001M, S1... lgl
## (1): ...
## i Use `spec()` to retrieve the full column specification for this data. i
## Specify the column types or set `show_col_types = FALSE` to quiet this message.
## • ` ` -> ` ... 747`
```

```
# Repeat the process for 2021 data
gender_poverty_2021 <- data_2021 %>%
  select(
    NAME,
    Total_Male_Population = S1701_C01_011E,
    Male_Population_Below_Poverty = S1701_C02_011E,
    Percent_Male_Population_Below_Poverty = S1701_C03_011E,
    Total_Female_Population = S1701_C01_012E,
    Female_Population_Below_Poverty = S1701_C02_012E,
    Percent_Female_Population_Below_Poverty = S1701_C03_012E
  ) %>%
  mutate(
    Total_Male_Population = as.numeric(Total_Male_Population),
    Male_Population_Below_Poverty = as.numeric(Male_Population_Below_Poverty),
    Total_Female_Population = as.numeric(Total_Female_Population),
    Female_Population_Below_Poverty = as.numeric(Female_Population_Below_Poverty)
  ) %>%
  mutate(
    Male_Poverty_Rate = Male_Population_Below_Poverty / Total_Male_Population,
    Female_Poverty_Rate = Female_Population_Below_Poverty / Total_Female_Population
  ) %>%
  na.omit()
```

```
## Warning: There were 4 warnings in `mutate()` .
## The first warning was:
## i In argument: `Total_Male_Population = as.numeric(Total_Male_Population)` .
## Caused by warning:
## ! 强制改变过程中产生了NA
## i Run `dplyr::last_dplyr_warnings()` to see the 3 remaining warnings.
```

```
# Merge the datasets by NAME
merged_poverty <- left_join(gender_poverty_2020, gender_poverty_2021, by = "NAME", suffix = c
  ("_2020", "_2021"))

head(merged_poverty)
```

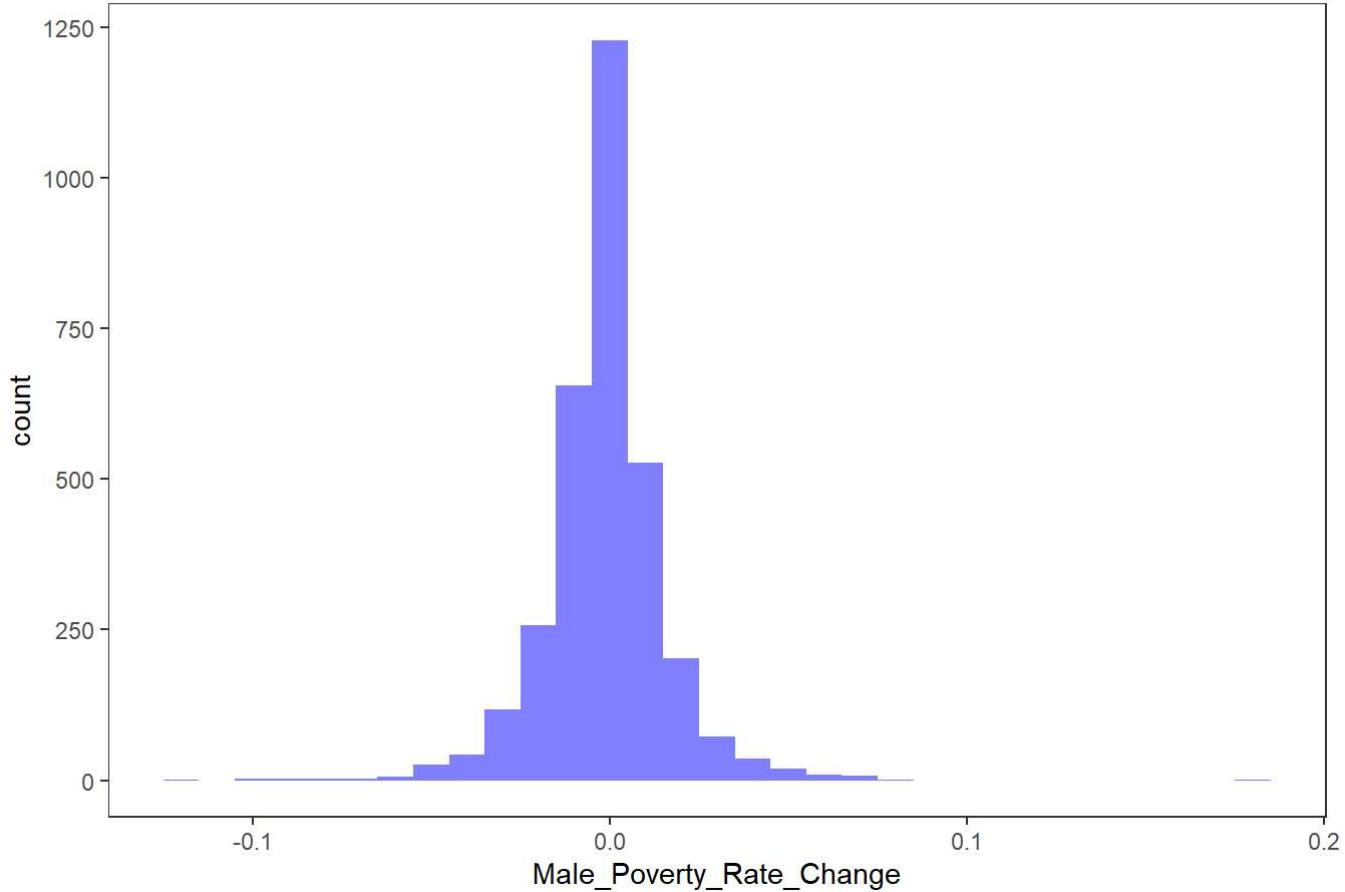
```
## # A tibble: 6 × 17
##   NAME      Total_Male_Population_2020¹ Male_Population_Below_Poverty_2020² Percent_Male_Population_Below_Poverty_2020³
##   <chr>          <dbl>                  <dbl> <chr>
## 1 Autauga ...       26781                 3417 12.8
## 2 Baldwin ...      103832                7803 7.5
## 3 Barbour ...      10346                 2509 24.3
## 4 Bibb Cou...       10507                 1865 17.8
## 5 Blount C...       28261                 3278 11.6
## 6 Bullock ...       5242                  1355 25.8
## # i abbreviated names: ¹Total_Male_Population_2020,
## #   ²Male_Population_Below_Poverty_2020,
## #   ³Percent_Male_Population_Below_Poverty_2020
## # i 13 more variables: Total_Female_Population_2020 <dbl>,
## #   Female_Population_Below_Poverty_2020 <dbl>,
## #   Percent_Female_Population_Below_Poverty_2020 <chr>,
## #   Male_Poverty_Rate_2020 <dbl>, Female_Poverty_Rate_2020 <dbl>, ...
```

4.2 Explore the sex and year for the poor people distribution

```
# Calculate the change in poverty rate
merged_poverty <- merged_poverty %>%
  mutate(
    Male_Poverty_Rate_Change = Male_Poverty_Rate_2021 - Male_Poverty_Rate_2020,
    Female_Poverty_Rate_Change = Female_Poverty_Rate_2021 - Female_Poverty_Rate_2020
  )

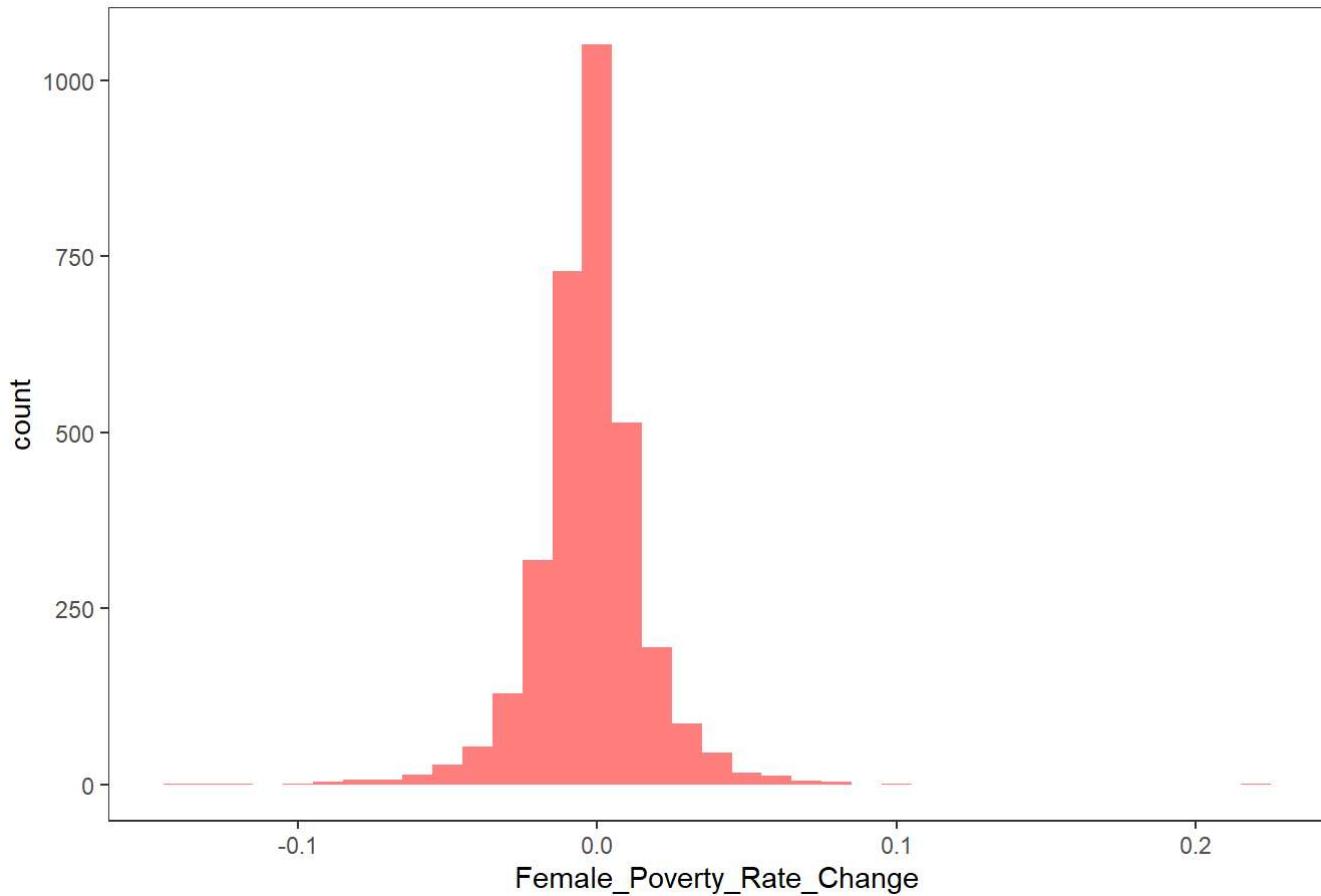
# Histogram of poverty rate change
ggplot(merged_poverty, aes(x = Male_Poverty_Rate_Change)) +
  geom_histogram(binwidth = 0.01, fill = "blue", alpha = 0.5) +
  labs(title = "Histogram of Male Poverty Rate Change 2020–2021")
```

Histogram of Male Poverty Rate Change 2020-2021



```
ggplot(merged_poverty, aes(x = Male_Poverty_Rate_Change)) +  
  geom_histogram(binwidth = 0.01, fill = "red", alpha = 0.5) +  
  labs(title = "Histogram of Male Poverty Rate Change 2020-2021")
```

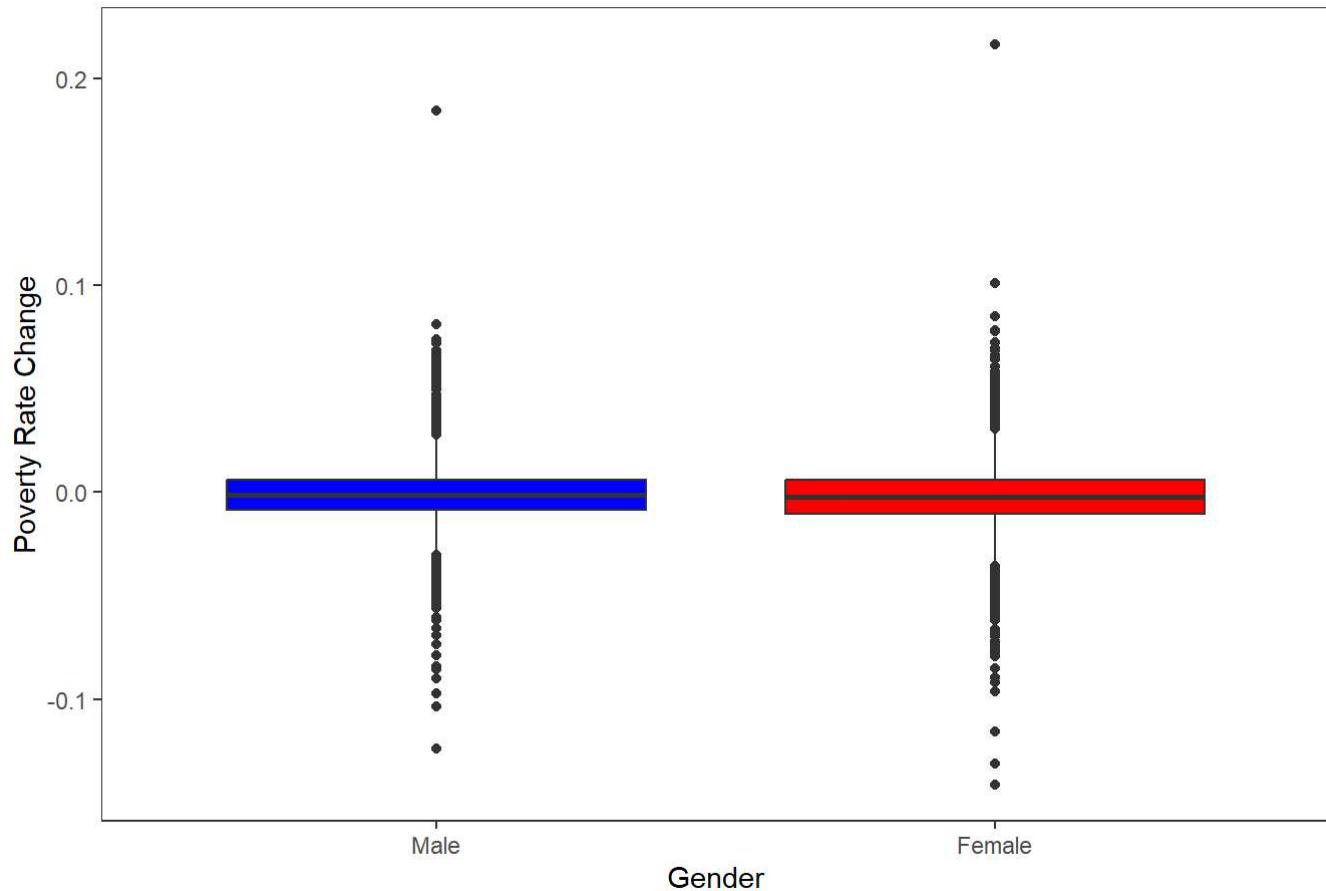
Histogram of Female Poverty Rate Change 2020-2021



The histogram shows the distribution of changes in male and female poverty rates from 2020 to 2021.

```
# Boxplot of poverty rate change
ggplot(merged_poverty) +
  geom_boxplot(aes(x = factor(0), y = Male_Poverty_Rate_Change), fill = "blue") +
  geom_boxplot(aes(x = factor(1), y = Female_Poverty_Rate_Change), fill = "red") +
  labs(title = "Boxplot of Poverty Rate Change by Gender 2020-2021") +
  xlab("Gender") +
  ylab("Poverty Rate Change") +
  scale_x_discrete(labels = c("Male", "Female"))
```

Boxplot of Poverty Rate Change by Gender 2020-2021



4.3 Conclusion

1. The histogram of the change in male poverty rate shows that most values are concentrated near zero, which means that for many regions, there is no significant change in male poverty rate.
2. The histogram of the change in female poverty rate also shows a similar pattern, but overall, it appears that the downward trend in female poverty rate is more pronounced.
3. The median change for both genders is close to zero, indicating that over half of the regions have experienced very small changes in poverty rates.
4. The quartile range for women is slightly wider than that for men, indicating that changes in female poverty rates are more dispersed across regions.

5. Combine Population for whom poverty and

flooding

5.1 Data merge and organization

```

library(tidyverse)
library(readr)

poverty_data_2020_df <- data_2020

# Filter out flood disasters
flood_disasters <- filter(disaster_declarations_df, incidentType == 'Flood')

columns_to_use <- c('NAME', 'S1701_C01_011E', 'S1701_C02_011E', 'S1701_C01_012E', 'S1701_C02_012E')
column_renames <- c('NAME', 'Total_Male_Population', 'Male_Population_Below_Poverty',
                     'Total_Female_Population', 'Female_Population_Below_Poverty')

gender_poverty_df <- poverty_data_2020_df %>%
  select(all_of(columns_to_use)) %>%
  rename_with(~ column_renames) %>%
  drop_na()

# Convert columns to numeric
for (col in column_renames[-1]) {
  gender_poverty_df[[col]] <- as.numeric(gender_poverty_df[[col]])
}

```

Warning: 强制改变过程中产生了NA

Warning: 强制改变过程中产生了NA

Warning: 强制改变过程中产生了NA

Warning: 强制改变过程中产生了NA

```

flood_disasters <- flood_disasters %>%
  mutate(County = str_replace(designatedArea, '\\\\(County\\\\)', ''),
        State = state)

gender_poverty_df <- gender_poverty_df %>%
  separate(NAME, into = c('County', 'State'), sep = ', ', extra = 'merge') %>%
  mutate(County = str_replace(County, ' County', ''))

```

Warning: Expected 2 pieces. Missing pieces filled with `NA` in 1 rows [1].

```

# A lookup table for state abbreviations and their full names
state_name_mapping <- data.frame(
  Abbreviation = c("AL", "AK", "AZ", "AR", "CA", "CO", "CT", "DE", "FL", "GA",
    "HI", "ID", "IL", "IN", "IA", "KS", "KY", "LA", "ME", "MD",
    "MA", "MI", "MN", "MS", "MO", "MT", "NE", "NV", "NH", "NJ",
    "NM", "NY", "NC", "ND", "OH", "OK", "OR", "PA", "RI", "SC",
    "SD", "TN", "TX", "UT", "VT", "VA", "WA", "WV", "WI", "WY"),
  FullName = c("Alabama", "Alaska", "Arizona", "Arkansas", "California", "Colorado",
    "Connecticut", "Delaware", "Florida", "Georgia",
    "Hawaii", "Idaho", "Illinois", "Indiana", "Iowa", "Kansas", "Kentucky",
    "Louisiana", "Maine", "Maryland", "Massachusetts", "Michigan", "Minnesota",
    "Mississippi", "Missouri", "Montana", "Nebraska", "Nevada", "New Hampshire",
    "New Jersey", "New Mexico", "New York", "North Carolina", "North Dakota",
    "Ohio", "Oklahoma", "Oregon", "Pennsylvania", "Rhode Island", "South Carolina",
    "South Dakota", "Tennessee", "Texas", "Utah", "Vermont", "Virginia",
    "Washington", "West Virginia", "Wisconsin", "Wyoming")
)

# Assuming flood_disasters is your dataframe and it has a column 'state' with state abbreviations
flood_disasters <- flood_disasters %>%
  left_join(state_name_mapping, by = c("state" = "Abbreviation")) %>%
  select(-state) %>%
  rename(state = FullName)

## Remove the na
gender_poverty_df <- gender_poverty_df[-1, ]
flood_disasters <- flood_disasters[-1, ]

# Merge the datasets

merged_data <- merge(flood_disasters, gender_poverty_df, by.x = c("state", "County"), by.y = c
  ("State", "County"), all.x = TRUE)
## select the cols we need

merged_data <- merged_data %>%
  select(
    County ,
    state,
    Total_Female_Population,
    Total_Male_Population,
    Male_Population_Below_Poverty,
    Female_Population_Below_Poverty
  )

# remove the na
merged_data <- na.omit(merged_data)

head(merged_data)

```

```
##          County      state Total_Female_Population Total_Male_Population
## 2001    Carroll     Iowa              9876                  9885
## 2070    Clinton     Iowa             23248                 22358
## 3790 Baltimore Maryland           424925                382754
## 3791 Baltimore Maryland           424925                382754
## 3792 Baltimore Maryland           424925                382754
## 5299 St. Louis Missouri          512126                462482
##          Male_Population_Below_Poverty Female_Population_Below_Poverty
## 2001                      725                   743
## 2070                     2719                  3512
## 3790                     33086                 41858
## 3791                     33086                 41858
## 3792                     33086                 41858
## 5299                     38585                 52052
```

5.2 Initial question

What conclusion can you draw by comparing the poverty rate in flood affected areas with the national average poverty rate and the average poverty rate in each state.

5.3 EDA and Solution

```
# Add the total columns to the gender_poverty_df
gender_poverty_df <- gender_poverty_df %>%
  mutate(Total_Population = Total_Male_Population + Total_Female_Population,
         Total_Population_Below_Poverty = Male_Population_Below_Poverty + Female_Population_Below_Poverty)

# Calculate the poverty rate for each county affected by floods
merged_data <- merged_data %>%
  mutate(Total_Population = Total_Male_Population + Total_Female_Population,
         Total_Population_Below_Poverty = Male_Population_Below_Poverty + Female_Population_Below_Poverty,
         Poverty_Rate = Total_Population_Below_Poverty / Total_Population)

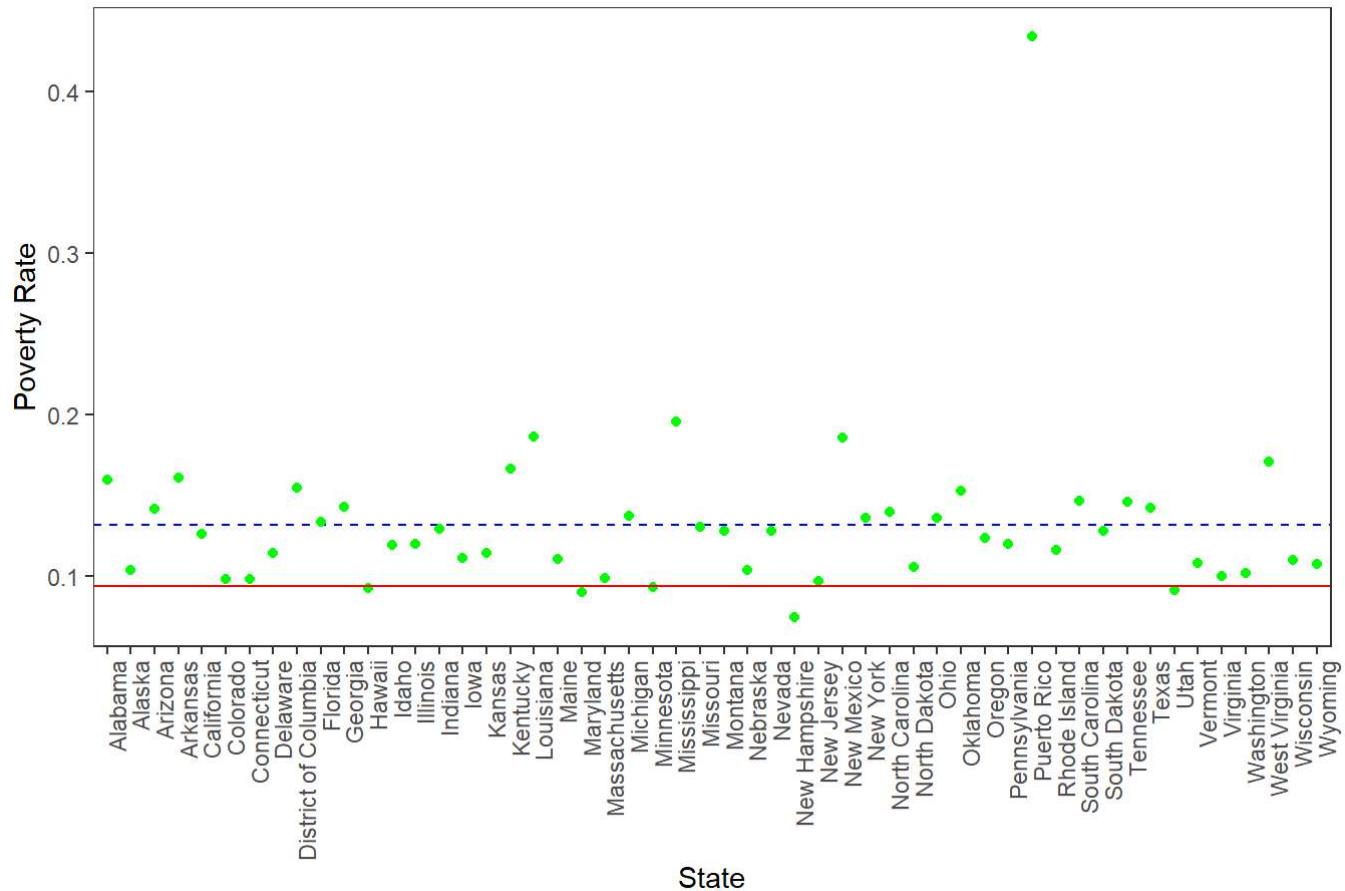
national_average_poverty_rate <- sum(gender_poverty_df$Total_Population_Below_Poverty) /
  sum(gender_poverty_df$Total_Population)

# Calculate the average poverty rate for each state
state_average_poverty_rates <- gender_poverty_df %>%
  group_by(State) %>%
  summarise(Total_Population_Below_Poverty = sum(Total_Population_Below_Poverty),
            Total_Population = sum(Total_Population)) %>%
  mutate(Poverty_Rate = Total_Population_Below_Poverty / Total_Population)

# Calculate the average poverty rate for counties affected by floods
average_poverty_rate_flood_affected <- mean(merged_data$Poverty_Rate)

# Plotting the average poverty rates for comparison
ggplot() +
  geom_hline(yintercept = national_average_poverty_rate, color = 'blue', linetype = 'dashed') +
  geom_point(data = state_average_poverty_rates, aes(x = State, y = Poverty_Rate), color = 'green') +
  geom_hline(yintercept = average_poverty_rate_flood_affected, color = 'red') +
  labs(x = 'State', y = 'Poverty Rate', title = 'Comparison of Poverty Rates: National, State, and Flood Affected Counties') +
  theme(axis.text.x = element_text(angle = 90, hjust = 1))
```

Comparison of Poverty Rates: National, State, and Flood Affected Counties



- In this chart, we can see the following content:
1. The blue dashed line represents the national average poverty rate.
 2. The green dots represent the average poverty rate for each state.
 3. The solid red line represents the average poverty rate of counties affected by floods.

Conclusion:

We can see that some states have poverty rates far above the national average, while the average poverty rate in flood affected areas is slightly higher than the national average.