

# DS5110 Homework 6

Ameya Santosh Gidh

2024-04-12

## Part A

### Problem 1

```
# Suppress startup messages from packages
suppressPackageStartupMessages(library(tidyverse))

# Load necessary libraries
library(readr)      # For reading data
library(dplyr)      # For data manipulation
library(ggplot2)    # For data visualization
library(tidyverse)  # Comprehensive data manipulation and visualization package
library(glmnet)     # For fitting Lasso and Elastic-Net regularized generalized linear models

## Loading required package: Matrix
##
## Attaching package: 'Matrix'
## The following objects are masked from 'package:tidyr':
##
##   expand, pack, unpack
## Loaded glmnet 4.1-8
library(stringr)    # For string manipulation
library(tokenizers) # For tokenization
library(lubridate)  # For working with dates
library(tidytext)   # For text mining with tidy data principles

# Check if tidyverse package is installed, if not, install it
if (!requireNamespace("tidyverse", quietly = TRUE)) {
  install.packages("tidyverse")
}

# Define file directory and path
trump_data <- read_csv(file.path("./", "twitter", "realDonaldTrump-20201106.csv"), col_types=cols(id=co

# Load stop words
data(stop_words)

# Define additional stop words
extra <- c("realdonaldtrump", "donaldtrump", "donald", "trump", "rt") # Define additional words to be c
```

```

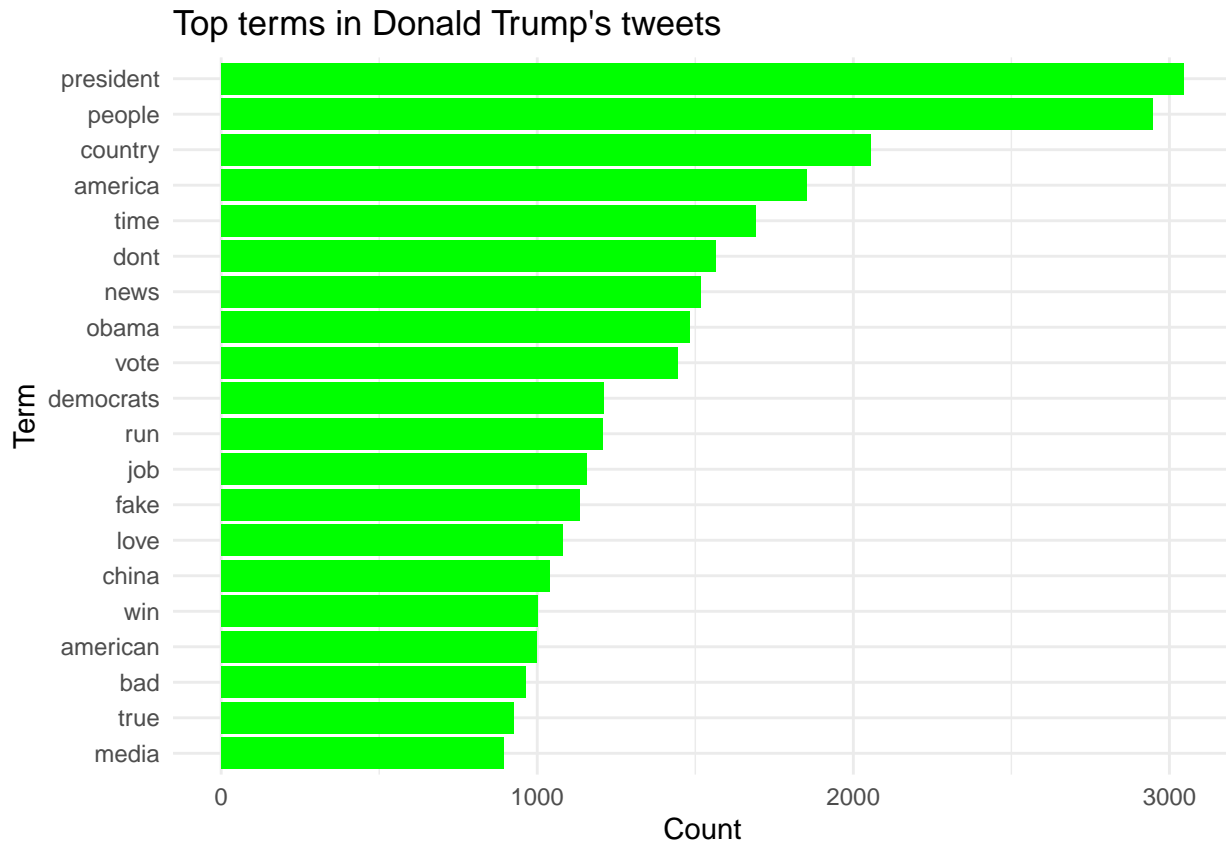
# Process tweets
tidy <- trump_data %>% # Start processing data using piping
  filter(!isRetweet, str_detect(text, "[:space:]")) %>% # Filter out retweets and tweets without spaces
  select(id, text, retweets, date) %>% # Select necessary columns
  unnest_tokens("word", text, token = "regex") %>% # Tokenize the text column into individual words
  mutate( # Replace words starting with @ with AT
    word=str_replace_all(word, "[:punct:]&&[~#]", ""), # Remove punctuation from words
  ) %>% # Replace AT with @
  anti_join(stop_words, by = "word", copy = TRUE) %>% # Remove stop words
  filter(
    !word %in% extra, # Filter out additional stop words
    !str_detect(word, "@"), # Filter out words containing @
    !str_detect(word, "amp"), # Filter out words containing "amp"
    !str_detect(word, "http"), # Filter out words containing "http"
    str_length(word) > 0L # Filter out empty words
  )

# Visualize top 20 most common terms
top_20_terms <- tidy %>% # Start processing tidy dataframe
  count(word) %>% # Count the frequency of each word
  top_n(20) # Select top 20 words

## Selecting by n

# Plot
ggplot(top_20_terms, aes(x = reorder(word, n), y = n)) + # Create ggplot with word frequency data
  geom_col(fill = "green") + # Add column plot with light blue color
  coord_flip() + # Flip the coordinates
  labs(x = "Term", y = "Count", title = "Top terms in Donald Trump's tweets") + # Add labels to axes and title
  theme_minimal() # Apply minimal theme

```



The bar chart illustrates the most frequently used terms in Donald Trump's tweets, where each term's frequency is represented by the length of the horizontal bars. The x-axis displays the count of each term, ranging from 0 to 3000, while the y-axis lists the terms themselves. The bars are shaded in varying tones of grey.

At the top of the list, the term "president" is notably the most frequent, followed by terms like "people," "country," "America," and "time." Further down, terms such as "news," "Obama," "vote," "Democrats," "run," and "job" also feature prominently. Additionally, terms like "fake," "love," "China," "win," "American," "bad," "true," and "media" are observed. Overall, the longest bar corresponds to the term "president," indicating its prevalence in Trump's tweets.

## Problem 2

```
# Load necessary libraries
library(lubridate) # For handling date data
library(ggplot2)   # For data visualization
library(viridis)   # For color scales

## Loading required package: viridisLite

# Top 20 terms for each year
top_20_terms <- tidy %>%
  mutate(year = as.integer(format(date, "%Y"))) %>% # Extract year from date
  filter(year >= 2015 & year <= 2020) %>% # Filter data for years 2015 to 2020
  count(word, year) %>% # Count occurrences of each word by year
  group_by(year) %>% # Group by year
  top_n(20) %>% # Select top 20 terms for each year
  ungroup()
```

```

## Selecting by n
# Identify terms that start with a hashtag
hashtag_terms <- top_20_terms %>%
  filter(stringr::str_detect(word, "^#")) # Filter terms starting with a hashtag

# Create a column for labels with hashtags for terms starting with a hashtag
top_20_terms$label <- top_20_terms$word

# Plotting top 10 terms
ggplot(top_20_terms, aes(x = reorder_within(word, n, year), y = n, fill = factor(year))) + # Start plot
  geom_col(color = "white", size = 0.25) + # Add column plot with white outline
  coord_flip() + # Flip the coordinates
  scale_x_reordered() + # Reorder x-axis labels
  scale_y_continuous(labels = NULL) + # Remove y-axis labels
  scale_fill_viridis(discrete = TRUE, option = "magma") + # Set color scale
  facet_wrap(~year, scales = "free", ncol = 2) + # Wrap facets by year, adjust scales, and arrange in t
  labs(x = "Term", y = "Count", title = "Top 20 Common Terms in Trump's Tweets by Year") + # Add labels
  theme_minimal() + # Apply minimal theme
  theme( # Customize the plot appearance
    text = element_text(size = 6), # Set text size
    axis.title = element_text(face = "bold", size = 20), # Set axis title appearance
    panel.spacing = unit(10, "pt"), # Set panel spacing
    plot.title = element_text(face = "bold", size = 14), # Set plot title appearance
    panel.background = element_rect(fill = "lightblue") # Set light blue background
  )

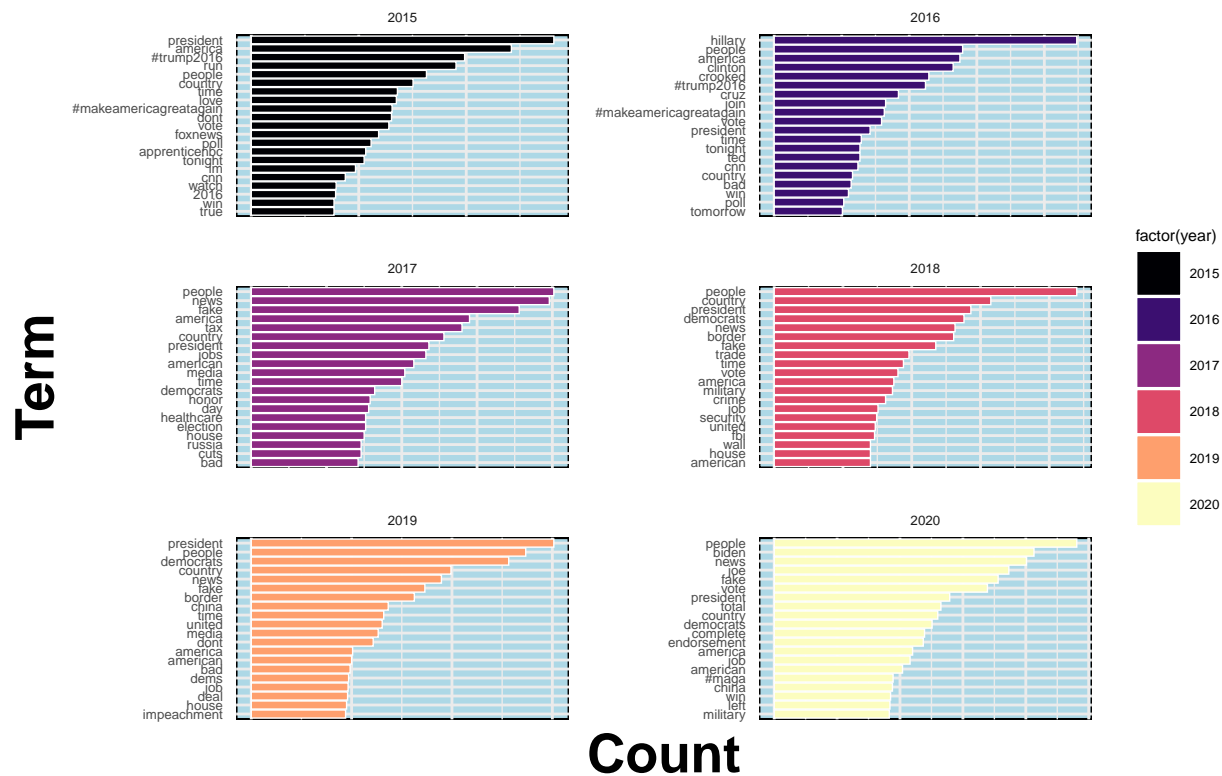
```

```

## Warning: Using `size` aesthetic for lines was deprecated in ggplot2 3.4.0.
## i Please use `linewidth` instead.
## This warning is displayed once every 8 hours.
## Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
## generated.

```

## Top 20 Common Terms in Trump's Tweets by Year



The bar charts titled “Top 20 Common Terms from Trump’s Tweets for Each Year” cover the years 2015 to 2020, showcasing the frequency of terms in Donald Trump’s tweets. Each bar represents the count of a specific term for a particular year.

In 2015 and 2016, Trump’s campaign hashtags such as #Trump2016 and #MakeAmericaGreatAgain are noticeable. The term “fake” recurs from 2017 to 2020, while “Biden” emerges in 2020 due to the election.

In 2015, the terms mainly revolve around Trump’s involvement in “The Apprentice” TV show and his presidential campaign, featuring terms like “apprentice” and “makeamericagreatagain.”

Transitioning to 2016, the focus shifts to his presidential campaign, with terms like “trump2016” and “makeamericagreatagain,” alongside references to his opponent Ted Cruz.

The terms in 2017 highlight political issues and events, including “fake news,” “tax reform,” and “obamacare.”

In 2018, political controversies and events are prominent, with terms like “witch hunt,” “mueller,” and “collusion.”

2019 revolves around political opposition and impeachment topics, such as “impeachment,” “mueller,” and “whistleblower.”

Lastly, 2020 is dominated by the COVID-19 pandemic and the presidential election, featuring prevalent words like “coronavirus,” “impeachment,” “biden,” and “covid.”

### Problem 3

```
# Load necessary libraries
library(tidytext) # For text mining with tidy data principles
library(dplyr)    # For data manipulation
library(ggplot2)  # For data visualization
```

```

# Assuming trump_tidy is already processed and in the correct format

# Create year variable from date
tidy <- tidy %>%
  mutate(year = lubridate::year(date)) # Add a new column "year" by extracting the year from the "date"

# Calculate TF-IDF for each term and year
trump_year_wise_tweets <- tidy %>%
  filter(year >= 2015 & year <= 2020) %>% # Filter data for the years 2015 to 2020
  count(word, year) %>% # Count the occurrences of each word by year
  bind_tf_idf(word, year, n) %>% # Calculate TF-IDF for each term and year
  arrange(desc(tf_idf)) # Arrange the data in descending order of TF-IDF score

# Visualize the top 20 characteristic terms for each year
top_terms_by_year <- trump_year_wise_tweets %>%
  group_by(year) %>% # Group the data by year
  top_n(20, wt = tf_idf) # Select the top 20 terms with the highest TF-IDF score for

# Plotting
ggplot(top_terms_by_year, aes(x = reorder_within(word, tf_idf, year), y = tf_idf, fill = factor(year)))
  geom_col(show.legend = FALSE) + # Add column plot without legend
  coord_flip() + # Flip the coordinates
  scale_x_reordered() + # Reorder x-axis labels
  scale_y_continuous(labels = NULL) + # Remove y-axis labels
  scale_fill_viridis(discrete = TRUE, option = "magma") + # Set color scale to light colors using Viri
  facet_wrap(~year, scales = "free", ncol = 2) + # Wrap facets by year, adjust scales, and ar
  labs(x = "Term", y = "tf-idf", # Add labels to axes
       title = "Top 20 Characteristic Terms from Trump's Tweets for Each Year") + # Add title
  theme(text = element_text(size = 6), # Set text size
        axis.title = element_text(face = "bold", size = 20), panel.spacing = unit(10, "pt"), # Customize ax
        plot.title = element_text(face = "bold", size = 12) ) # Customize plot title appearance

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :

```

```

## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <b8>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <b8>

```

[illegible]



[illegible]

[illegible]

```

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <b8>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <b8>

```



```

## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <ba>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <87>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbcsToSbcs': dot substituted for <b8>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <f0>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <9f>

## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#americafirst ' in 'mbcsToSbcs': dot substituted for
## <87>

```



[illegible]

```
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <87>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <b8>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <f0>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <9f>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <87>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <ba>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <f0>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <9f>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <87>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbcsToSbc': dot substituted for <b8>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <f0>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <9f>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <87>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <ba>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <f0>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <9f>  
  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbcsToSbc': dot substituted for  
## <87>
```



[illegible]

```
## Warning in grid.Call(as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <87>  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <ba>  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <f0>  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <9f>  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <87>  
## Warning in grid.Call(C_textBounds, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <b8>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <f0>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <9f>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <87>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <ba>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <f0>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <9f>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#americafirst ' in 'mbsToSbcs': dot substituted for  
## <87>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <f0>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <9f>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <87>  
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :  
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <ba>
```

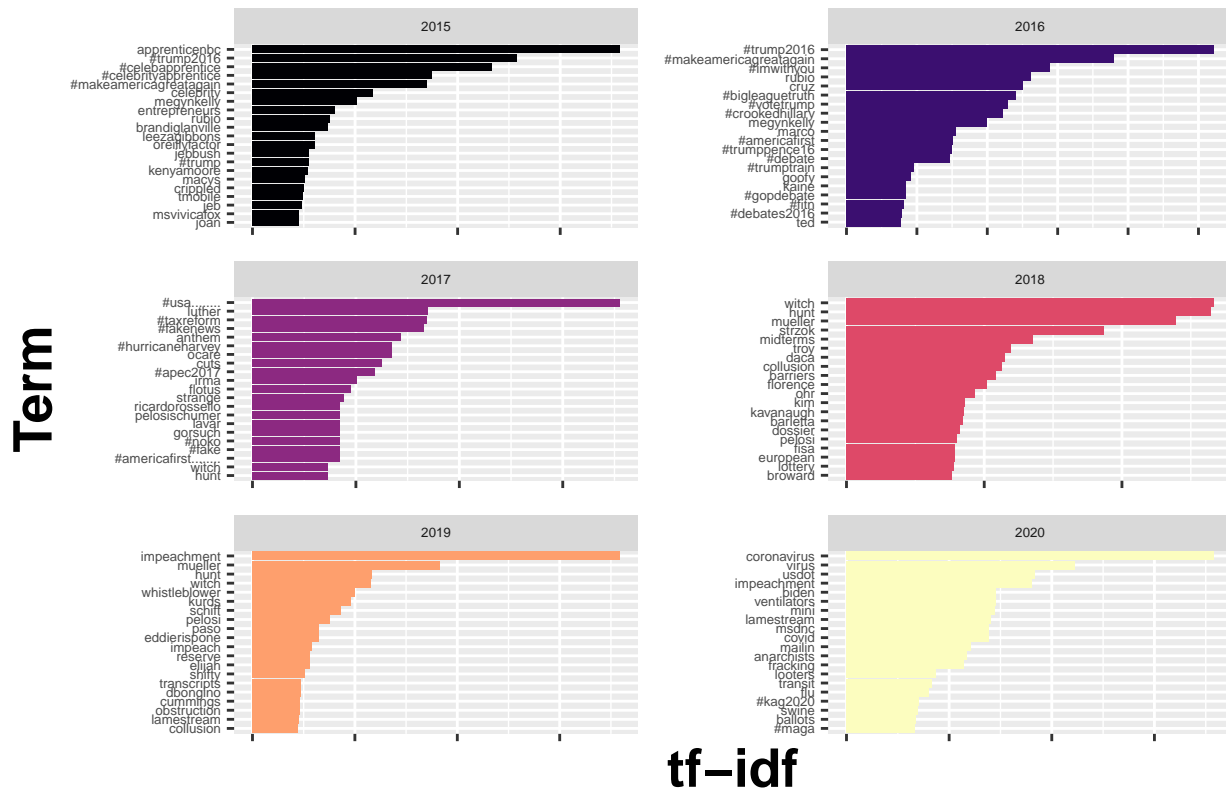
```
## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <f0>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <9f>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <87>

## Warning in grid.Call.graphics(C_text, as.graphicsAnnot(x$label), x$x, x$y, :
## conversion failure on '#usa ' in 'mbsToSbcs': dot substituted for <b8>
```

## Top 20 Characteristic Terms from Trump's Tweets for Each Year



The bar chart depicts the top 20 characteristic terms from Donald Trump's tweets for each year spanning from 2015 to 2020. Each year is represented by bars of different colors, where the length of each bar corresponds to the term frequency-inverse document frequency (tf-idf) score for that specific term and year.

In 2015, the chart emphasizes terms associated with Donald Trump's early campaign efforts and television appearances, such as "apprentice," "trump2016," and "makeamericagreatagain."

Transitioning to 2016, the terms primarily revolve around the presidential campaign, featuring terms like "trump2016," "makeamericagreatagain," and mentions of other political figures like "cruz" and "rubio."

In 2017, the focus shifts towards political and legislative issues, with significant terms including "fake news," "tax reform," and "obamacare."

By 2018, attention is drawn to political controversies, with terms like "witch hunt," "mueller," and "collusion" taking the forefront.

The terms from 2019 center around impeachment proceedings and political adversaries, featuring notable words such as "impeachment," "mueller," and "whistleblower."

Finally, in 2020, the terms are greatly influenced by the COVID-19 pandemic and political events, with “coronavirus,” “impeachment,” and “biden” being prominent terms.

## Part B

### Problem 4

```
# Filter the data to include only tweets from 2016-2020
filtered_data <- tidy %>%
  filter(year >= 2016 & year <= 2020)

# Create a document-term matrix
dtm <- filtered_data %>%
  count(id, word) %>%
  cast_sparse(id, word, n)

# Display dimensions of the document-term matrix
dim(dtm)

## [1] 18483 19440

# Extract the tweet IDs to get the target variable (# of retweets)
ids <- tibble(id=rownames(dtm))
ids <- left_join(ids, trump_data)

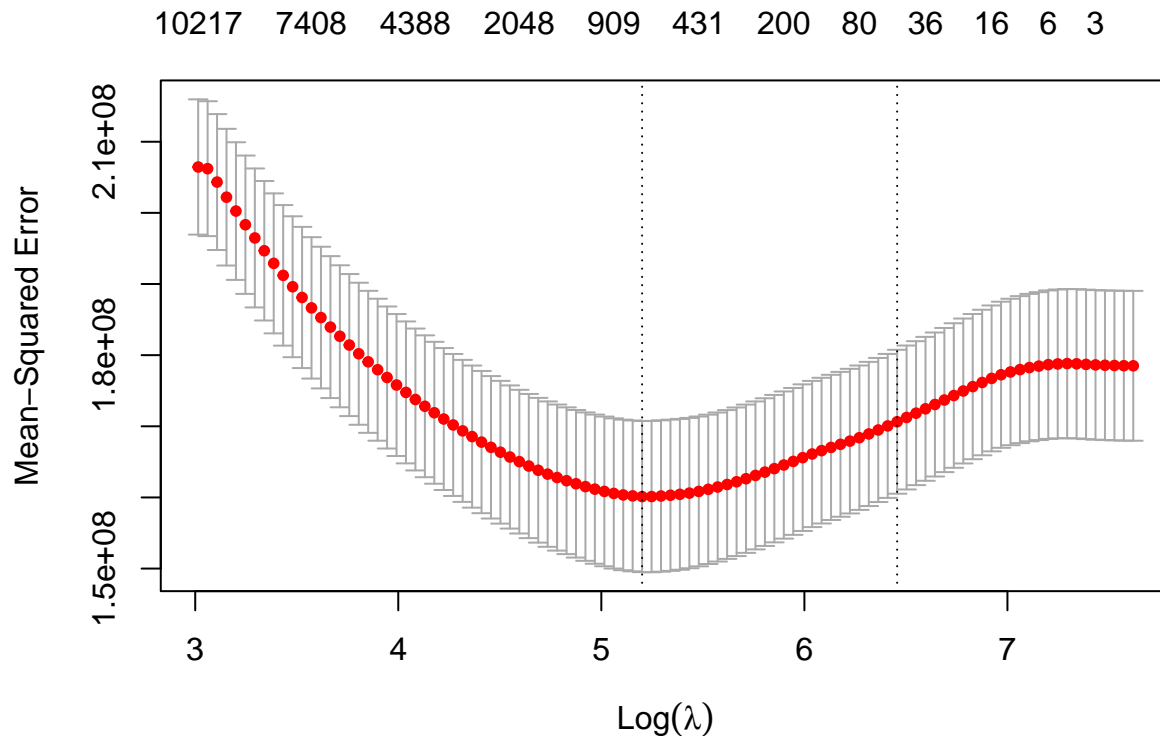
## Joining with `by = join_by(id)`

# Define the target variable
retweets <- ids$retweets

# Set seed for reproducibility
set.seed(2020)

# Fit Lasso and Elastic-Net regularized generalized linear models
fit_cv <- cv.glmnet(dtm, retweets)

# Visualize the cross-validation results
plot(fit_cv)
```



```
# Display cross-validation results
fit_cv
```

```
##
## Call:  cv.glmnet(x = dtm, y = retweets)
##
## Measure: Mean-Squared Error
##
##      Lambda Index   Measure      SE Nonzero
## min 181.4    53 160119617 10622183    714
## 1se 637.0    26 170638272 10113412    51
```

```
# Extract the selected value of lambda
selected_lambda <- fit_cv$lambda.min
cat("Selected value of lambda:", selected_lambda, "\n")
```

```
## Selected value of lambda: 181.4303
```

```
# Get the number of non-zero coefficients
coefficients <- coef(fit_cv, s = selected_lambda)
non_zero_coefficients <- sum(coefficients != 0)
cat("Number of non-zero coefficients:", non_zero_coefficients, "\n")
```

```
## Number of non-zero coefficients: 715
```

The model achieving the lowest Mean-Squared Error (MSE) opts for  $\lambda = 181.4$  and encompasses 719 terms. Moreover, within one standard error of the optimal model's MSE, the most parsimonious model selects  $\lambda = 637.0$ , employing merely 50 terms.

## Problem 5

```
# Extract coefficients from the best model
coefficients <- coef(fit_cv, s = "lambda.min")
```

```
# Convert coefficients to a tibble
coefficients <- tibble(word = rownames(coefficients), coef = as.numeric(coefficients))

# Select the top 15 coefficients
top_15_terms <- coefficients %>%
  top_n(15)
```

```
## Selecting by coef
```

```
# Plot the top 15 terms with their coefficients
```

```
library(ggplot2)
```

```
ggplot(top_15_terms, aes(x = reorder(word, coef), y = coef/10)) + # Create ggplot with top 15 terms and
```

```
  geom_col(fill = "orange", width = 0.5) + # Add column plot with orange color
```

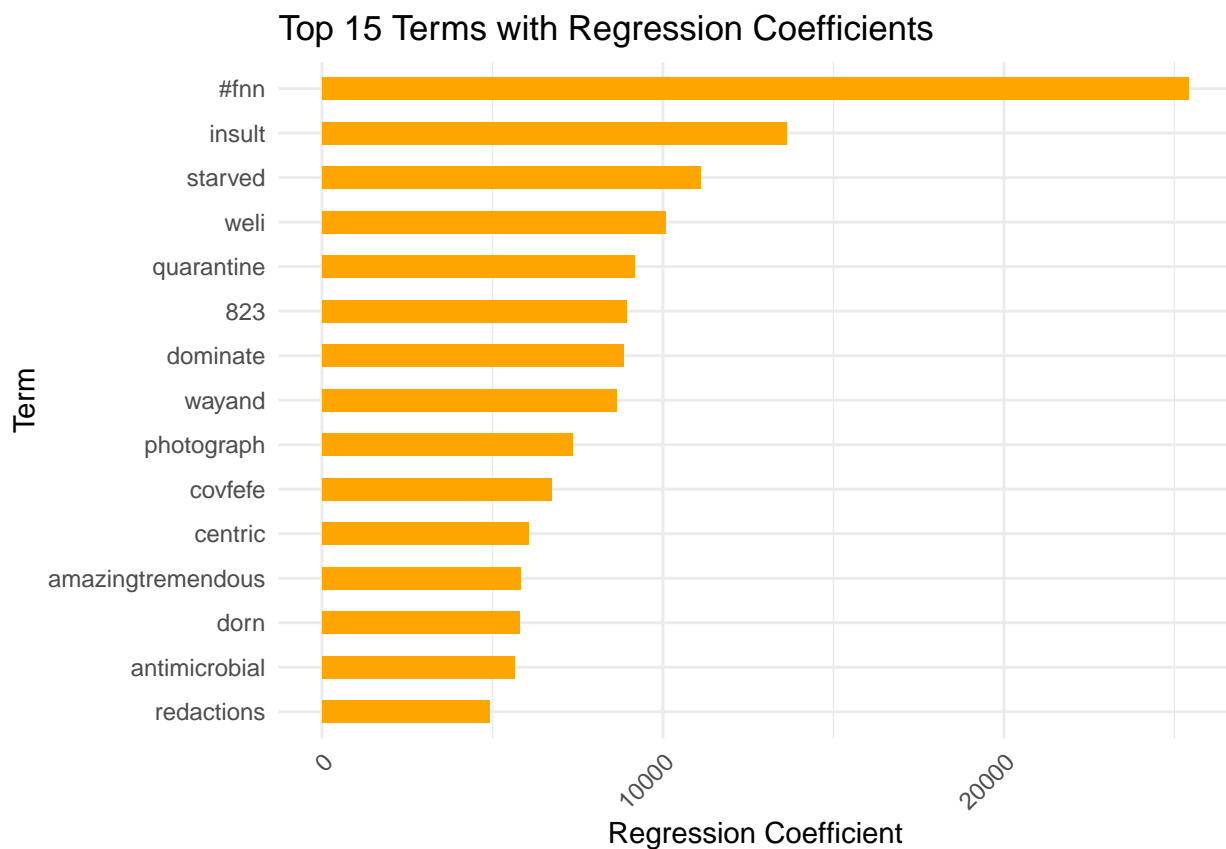
```
  coord_flip() + # Flip the coordinates
```

```
  labs(x = "Term", y = "Regression Coefficient", # Add labels to axes
```

```
        title = "Top 15 Terms with Regression Coefficients") + # Add title
```

```
  theme_minimal() + # Apply minimal theme
```

```
  theme(axis.text.x = element_text(angle = 45, hjust = 1)) # Rotate x-axis labels
```



The term “FNN,” which presumably refers to Fox News Network, emerges as the top term significantly correlated with a high number of retweets. Interestingly, the phrase “covfefe,” uniquely associated with Trump, also demonstrates a notable correlation with increased retweet counts. To delve deeper into the context surrounding these terms, we’ll analyze their usage within the original tweets.

```
trump_lower <- trump_data %>%
  filter(year(date) >= 2016) %>%
  mutate(text=str_to_lower(text))
```

```
trump_lower %>%
  filter(str_detect(text, "#fnn")) %>%
  select(date, text, retweets)
```

```
## # A tibble: 1 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2017-07-02 13:21:42 #fraudnewscnn #fnn https://t.co/wyunhjjujg 293109
```

```
trump_lower %>%
  filter(str_detect(text, "insult")) %>%
  select(date, text, retweets)
```

```
## # A tibble: 13 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2020-08-06 16:48:38 "rt @teamtrump: president @realdonaldtrump: "jo~ 12172
## 2 2020-08-06 21:36:37 "rt @christianwalk1r: joe biden's comfort with ~ 18885
## 3 2020-08-07 12:36:13 "rt @dcexaminer: \"joe biden this morning, he t~ 6688
## 4 2020-01-25 02:31:50 "rt @hawleymo: schiff & co end as they bega~ 11624
## 5 2020-01-28 04:16:55 "rt @sendansullivan: congressman nadler tried t~ 4872
## 6 2020-07-26 23:58:24 "rt @richlowry: the most insulting racist stere~ 3493
## 7 2016-09-10 12:47:18 "wow, hillary clinton was so insulting to my su~ 16931
## 8 2017-11-12 00:48:01 "\"why would kim jong-un insult me by calling m~ 217199
## 9 2018-11-09 21:10:46 "president macron of france has just suggested ~ 32187
## 10 2018-06-09 23:04:54 "pm justin Trudeau of Canada acted so meek and ~ 27539
## 11 2019-10-10 01:34:14 "\"i don't think it's a whistleblower at all. i ~ 17857
## 12 2019-06-25 14:42:30 "...iran's very ignorant and insulting stateme~ 20604
## 13 2019-02-21 16:09:57 ".@jussiesmollett - what about maga and the ten~ 45070
```

*# Filter the dataframe to include only tweets containing the word "starved"*

```
trump_lower %>%
  filter(str_detect(text, "starved")) %>%
  # Select relevant columns for analysis
  select(date, text, retweets)
```

```
## # A tibble: 1 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2018-01-03 00:49:19 north korean leader kim jong un just stated that~ 153496
```

*# Filter tweets containing the pattern "weli" using the lowercased text*

```
trump_lower %>%
  filter(str_detect(text, "weli")) %>%
  # Select relevant columns for display
  select(date, text, retweets)
```

```
## # A tibble: 1 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2020-10-03 03:31:34 going weli, i think! thank you to all. love!!! 139605
```

*# Filter the trump\_lower dataframe to include only tweets containing the word "quarantine"*

```
trump_lower %>%
  filter(str_detect(text, "quarantine")) %>%
  select(date, text, retweets)
```

```
## # A tibble: 8 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2020-10-02 02:44:21 hope hicks, who has been working so hard without~ 48126
## 2 2020-03-28 17:31:58 i am giving consideration to a quarantine of dev~ 30638
## 3 2020-03-29 00:19:31 ....federal government. a quarantine will not be~ 15108
## 4 2020-03-29 16:23:10 rt @ingrahamangle: smart that @realdonaldtrump l~ 6998
## 5 2020-03-06 15:47:58 spoke to governor @gavinnewsom early this mornin~ 11776
## 6 2020-03-22 09:42:35 rt @whnsc: psa: text messages and emails about n~ 2453
## 7 2020-03-23 12:32:55 rt @randpaul: senator rand paul has tested posit~ 18102
## 8 2020-10-02 04:54:06 tonight, @flotus and i tested positive for covid~ 408866
```

```
# Filter tweets containing the term "covfefe"
```

```
trump_lower %>%
  filter(str_detect(text, "covfefe")) %>%
  # Select specific columns for analysis
  select(date, text, retweets)
```

```
## # A tibble: 2 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2017-05-31 10:09:22 "\"who can figure out the true meaning of \""co~ 68046
## 2 2017-05-31 04:06:25 "despite the constant negative press covfefe" 127507
```

```
# Filter the 'trump_lower' dataset to include only rows where the text contains the pattern "way[:punct
```

```
trump_lower %>%
  filter(str_detect(text, "way[:punct:]+and")) %>%
  # Select specific columns for output
  select(date, text, retweets)
```

```
## # A tibble: 2 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2020-01-05 05:11:03 the united states just spent two trillion dollar~ 131915
## 2 2018-09-01 13:27:05 ....donald trump, and now we find out that there~ 18901
```

```
# Filter tweets containing the pattern "amazing[:punct:]+tremendous"
```

```
trump_lower %>%
  filter(str_detect(text, "amazing[:punct:]+tremendous")) %>%
  # Select columns for date, text, and retweets
  select(date, text, retweets)
```

```
## # A tibble: 1 x 3
##   date                text                retweets
##   <dtm>              <chr>              <dbl>
## 1 2020-10-03 17:19:38 doctors, nurses and all at the great walter reed~ 100144
```