

# Homework 6

Aidan Pizzo

10/6/21

```
#data_munged.R
library("stringr")

munge_data <- function()
{
  con <- file("../data/Presidential_Debate_Transcript.txt", "r")
  lines <- readLines(con)
  lines <- setdiff(lines, "")
  close(con)
  ind <- grep("\\\\(", lines)
  time <- vector("character", length(ind))
  speaker <- vector("character", length(ind))
  their_line <- vector("character", length(ind))
  get_time_speaker <- lines[ind]
  speaker_header <- substr(get_time_speaker, 1, 3)
  for (i in 1:length(ind))
  {
    if (speaker_header[i] == "Chr") #Chris Wallace
    {
      speaker[i] <- "Wallace"
    }
    else if (speaker_header[i] == "Pre") #President Trump
    {
      speaker[i] <- "Trump"
    }
    else if (speaker_header[i] == "Vic") #VP Biden
    {
      speaker[i] <- "Biden"
    }
    time[i] <- get_time(get_time_speaker[i])
    for (i in 1:length(time))
    {
      time[i] <- str_replace(time[i], ":", "min")
      time[i] <- str_replace(time[i], "\\\\)", "sec")
    }
    their_line <- lines[ind+1]
    out_df <- data.frame(time, speaker, line = their_line)
  }
  return(out_df)
}

get_time <- function( line )
```

```

{
  split <- strsplit(line, split = "")
  start_ind <- grep("\\(", split[[1]])+1
  end_ind <- grep("\\)", split[[1]])
  time <- split[[1]][start_ind:end_ind]
  string_time <- paste(time, collapse= "")
  return(string_time)
}

d1 <- munge_data()

## Warning in readLines(con): incomplete final line found on '../data/
## Presidential_Debate_Transcript.txt'

write.csv(d1, "debate.csv", row.names = F)

#analysis.R

library(wordcloud2)

get_word_counts <- function(d, speaker)
{
  speaker_lines_vec <- d[d$speaker == speaker,3]
  split <- strsplit(speaker_lines_vec, " ")
  split <- unlist(split)
  unique_words <- unique(split)
  num_times_said <- vector("numeric", length(unique_words))
  for (i in 1:length(unique_words))
  {
    num_times_said[i] <- length(grep(unique_words[i], split, fixed = T))
  }
  out_df <- data.frame(word = unique_words, count = num_times_said)
  return(out_df)
}

total_word_counts <- function(d, speakers)
{
  if (length(speakers) == 2)
  {
    d1 <- d[d$speaker == speakers[1] | d$speaker == speakers[2], ]
  }
  else if (length(speakers) == 3)
  {
    d1 <- d[d$speaker == speakers[1] | d$speaker == speakers[2]
           | d$speaker == speakers[3], ]
  }
  split <- strsplit(d1$line, " ")
  split <- unlist(split)
  unique_words <- unique(split)
  num_times_said <- vector("numeric", length(unique_words))
  for (i in 1:length(unique_words))
  {
    num_times_said[i] <- length(grep(unique_words[i], split, fixed = T))
  }
  out_df <- data.frame(word = unique_words, count <- num_times_said)
}

```

```

    return(out_df)
}

prepare_word_cloud <- function(d, speaker)
{
  if (length(speaker)>1)
  {
    df <- total_word_counts(d, speaker)
  }
  else
  {
    df <- get_word_counts(d, speaker)
  }
  common_words_df <- read.csv("../data/word_frequency.csv", header=T, stringsAsFactors = F, strip.white=T)
  common_words <- toupper(common_words_df[1:500,2]) #vector is returned with spaces in front
  common_words <- clean_words(common_words)
  uncommon_words <- vector("character")
  freq <- vector("numeric")
  for (i in 1:length(df$word))
  {
    if (is.element(toupper(df$word[i]), common_words) == FALSE) #add the word to uncommon if it isn't in common
    {
      uncommon_words <- c(uncommon_words, df$word[i])
      freq <- c(freq, df$count[df$word == df$word[i]])
    }
  }
  return(data.frame(uncommon_words, weight = freq))
}

clean_words <- function(words)
{
  for (i in 1:length(words))
  {
    split <- strsplit(words[i], "")
    split <- split[[1]][4:length(split[[1]])]
    the_word <- paste(split, collapse= "")
    words[i] <- the_word
  }
  return(words)
}

```

Wordclouds on next pages:

1. Trump
2. Biden
3. Trump and Biden

They're Well,  
been has Joe.  
States. you, fact, election. wouldn't it, S  
Court up. he, phenomenal well. all. she's  
outcome Mr. plenty won way, forward  
There's Supreme greatest Senats voted Now, States  
rid middle Affordable  
election  
that. doing  
Vice  
wants



[illegible]



