

Peer-graded Assignment: Milestone Report

02 Sept

1000

This milestone report is for the exploratory analysis of the t

The motivation for this project is to:

1. Demonstrate that i have downloaded the data and have successfully loaded it in

3. Report any interesting findings.

- ## Initialization, Data Reading and Summary

```
list.of.packages <- c("dplyr", "ggplot2", "stringi", "knitr", "tm", "pwr", "wordcloud", "R",
"ra", "survey", "RSQLite", "SnowballC", "Rweka", "plotly", "devtools", "slam", "tokenizers")
```

```
new.packages <- list.of.packages[!(list.of.packages %in% in
if(length(new.packages)) install.packages(new.packages, rep
```

```
library(dplyr)
library(ggplot2)
library(stringi)
library(tm)
library(wordcloud)
library(RColorBrewer)
library(SnowballC)
library(Meika)
library(knitr)
library(kableExtra)
library(pwr)
library(survey)
library(RSQLite)
library(plotly)
library(devtools)
library(slam)
library(tokenizers)
library(quantmod)
library(ngram)
library(NLPI)

# devtools::install_github("ropensci/plotly")
```

Twitter	159.36	2360148	30093410
---------	--------	---------	----------

Data Sampling

Clearly, the data sets are very large and it is only logical that we take a non-biased sample to conduct

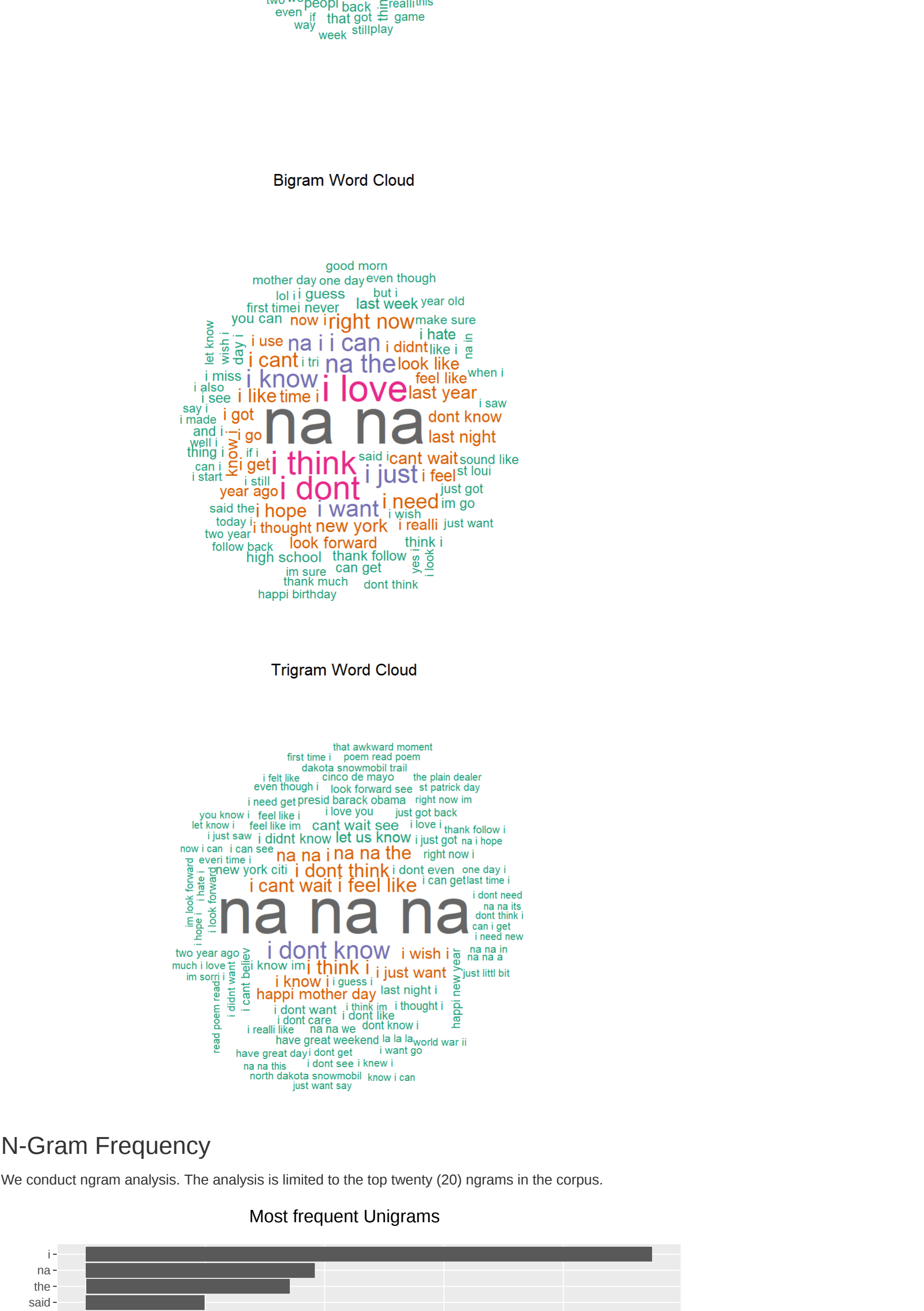
```
## [1] "The computed and combined sample size (lines) is 42695 with 1023145 words."
```

5. Stem the words (eg. "worked"
6. Convert to lower case

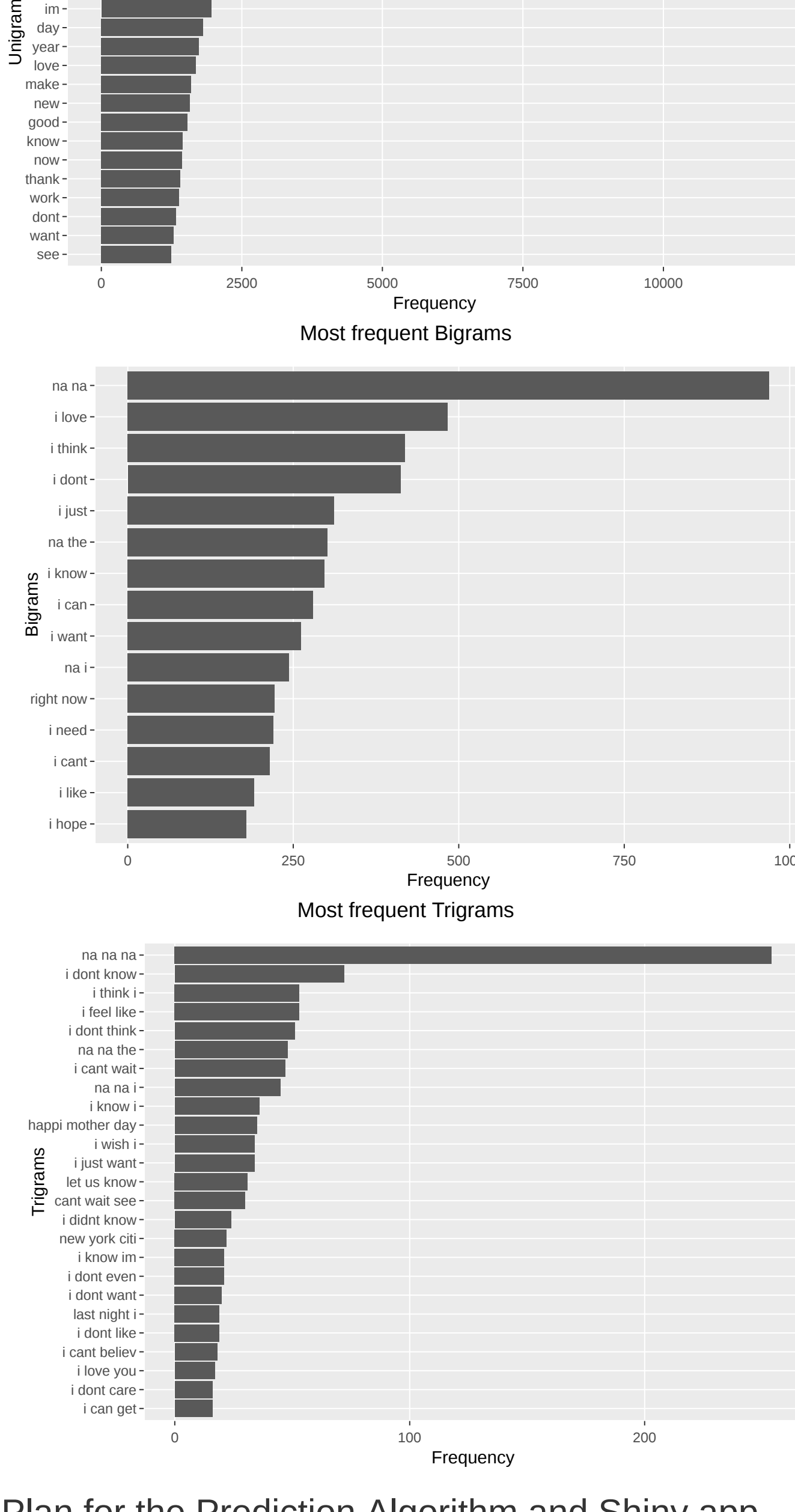
- identify all combinations of adjacent words (ngrams) in the documents. We then generate the word clouds.

but dont go one make you home

much he it's day last there
right on im na just e use let
its of the i get o go all cant
night need the get to your
and new said a year us show
feel come time can tr... what



	like	one
like	100%	95%
one	95%	100%

[illegible]

Resulting from exploratory analysis and the n-gram analysis, it is clear that sampling produced a reasonable approximation of the original text. Plotting the N-Grams for two and three word combinations mostly produced grammatically correct sentences.

1. Additional data cleaning
2. Better sample size selection technique for the training and testing data sets
3. Build the prediction model basing on the n-gram frequency on test data

4. Test the performance of the algorithm on test data using some of the following metrics:

5. Build a shiny app which is interactive and allows users to enter text for prediction of the next word(s)
 6. Experiment to optimize the model to handle unseen n-grams, increasing model performance.
 7. Get feedback ion the plan for creating a prediction algorithm and Shiny app
- ## ppendices
- ### ppendix - Setup
- ```
nitrr::opts_chunk$set(echo = FALSE)
```

```
ra", "survey", "RSQLi
new.packages <- list.
if(length(new.packages
```

```
library(dplyr)
```

```
library(ggplot2)
library(stringr)
library(tn)
library(wordcloud)
library(RColorBrewer)
library(SnowballC)
library(Rewek)
library(knitr)
library(kableExtra)
library(pwr)
library(survey)
library(RSQ Lite)
library(plotly)
library(devtools)
library(slam)
library(tokenizers)
library(quantdata)
library(ngram)
library(NLPIR)
devtools::install_github("ropensci/plotly")
```

```
on US summary <- matrix(0, nrow =
```

```

"no. of Lines", "no. of Words"))
for (i in 1:3) {
 con <- file(file.list[i], "rb")
 text[i,con] <- readLines(con, encoding = "UTF-8", skipNul = TRUE)
 close(con)
 en_US.summary[i,1] <- round(file.info(file.list[i])$size / 1024^2, 2)
 en_US.summary[i,2] <- length(text[i])
 en_US.summary[i,3] <- sum(str_count_words(text[i]))
}

```

```
kable_styling(bootstrap_options = "table")
```

```

set.seed(12345)

Get population sizes (lines)
blogs_pop <- length(text$blogs)
news_pop <- length(text$news)
twitter_pop <- length(text$twitter)

Create the sample
blogs_sample <- sample(text$blogs, 0.01*length(text$blogs))
news_sample <- sample(text$news, 0.01*length(text$news))
twitter_sample <- sample(text$twitter, 0.01*length(text$twitter))
sample_corpus_data <- c(blogs_sample, news_sample, twitter_sample)

sumLines <- length(sample_corpus_data)
sumWords <- sum(str_count(words(sample_corpus_data),
paste0("The computed and combined sample size (lines) is", sumLines, "with", sumWords, "words."))

Write the file to disc
writeLines(sample_corpus_data, "sample_corpus_data.txt")

```

```
corpusNew <- corpusNew %>%
 tm_map(removePunctuation) %>% #
```

```
tm_map(removeWords, stopwords("english")) %>% # remove stopwords
tm_map(stemDocument) %>% # Stem words
tm_map(tolower) %>% # convert to Lower case
tm_map(stripWhitespace) %>% # Strip white spaces
tm_map(PlainTextDocument) %>% Convert to plain text
```

---

## Appendix - Tokenise and Word Cloud

```
knitr::opts_chunk$set(echo = TRUE)
```

```
unigram <- NgramTokenizer(corpusNem, Weka_control(min = 1, max = 1))
bigram <- NgramTokenizer(corpusNem, Weka_control(min = 2, max = 2))
trigram <- NgramTokenizer(corpusNem, Weka_control(min = 3, max = 3))
```

```
unigram.df <- data.frame(table(unigram))
unigram.df <- unigram.df[order(unigram.df$Freq, de
```

```
Bigram
bigram.df <- data.frame(table(bigram))
bigram.df <- bigram.df[order(bigram.df$Freq, decreasing = TRUE),]

Trigram
trigram.df <- data.frame(table(trigram))
trigram.df <- trigram.df[order(trigram.df$Freq, decreasing = TRUE),]

Generate Word Clouds
col.palette <- brewer.pal(8, "Dark2") # Colour Palette to be used for all Word Clouds

Unigram Word Cloud
layout(matrix(c(1, 2), nrow=2), heights=c(1, 4))
par(mar=rep(8, 4))

plot.new()
text(x=0.5, y=0.5, "Unigram Word Cloud")
wordcloud(unigram.df$unigram, unigram.df$Freq, max.words=90, random.order = F, colors=col.palette, main="Title")

Bigram Word Cloud
layout(matrix(c(1, 2), nrow=2), heights=c(1, 4))
par(mar=rep(8, 4))

plot.new()
text(x=0.5, y=0.5, "Bigram Word Cloud")

wordcloud(bigram.df$bigram, bigram.df$Freq, max.words=90, random.order = F, colors=col.palette, main="Title")

Trigram Word Cloud
layout(matrix(c(1, 2), nrow=2), heights=c(1, 4))
par(mar=rep(8, 4))

plot.new()
text(x=0.5, y=0.5, "Trigram Word Cloud")
wordcloud(trigram.df$trigram, trigram.df$Freq, max.words=90, random.order = F, colors=col.palette, main="Title")
```

```
Tigram Frequency
ggTri <- ggplot(head(trigram_df, 25)) aes(reorder(
```