

Class Project 3 – Text Analysis with R

CSCI 6444 Introduction to big data
and analytics

Professor: Stephen Kaisler

Group 2

Shaival Vora (G37099269)

Nithin Raghava Aitha (G43945136)

Dataset Description

The data set given is a text file called Tarzan of the Apes by Edgar Rice Burroughs. For this project we are going to use chapters I – XV.

Objective

The goal of this project is to conduct text analytics on the extensive text document " Tarzan of the Apes " By employing different analytical techniques, our aim is to comprehend the text, discern its underlying sentiment, and draw insightful observations from it.

Packages

For this project we need some packages and we use "install.packages()" function to install them and "library()" function to load them.

```
# installing the required packages
install.packages('tm')
install.packages('dplyr')
install.packages('tidytext')
install.packages('ggplot2')
install.packages('wordcloud')

# loading the installed packages
library('tm')
library('dplyr')
library('tidytext')
library('ggplot2')
library('wordcloud')
```

These are the packages we initialize at the start but we do need a few other packages that are going to be installed and loaded in the further stage of the project.

Loading the Dataset and Creating a Corpus

We begin by setting up the directory using " setwd() " command.

```
> setwd("C:/Users/anith/OneDrive/Desktop/GW/sem 2/Big Data/project 3")
> getwd()
[1] "C:/Users/anith/OneDrive/Desktop/GW/sem 2/Big Data/project 3"
```

Then after setting up the directory, we load the dataset.

```
> Book <- readLines('TarzanOfTheApes.txt')
```

We then divide the chapters into individual documents to create a corpus. Each of the fifteen chapters (Chapter I-XV) is identified using the "which" method as demonstrated below. We segment the chapters based on their titles, such as "Chapter I," and save the content of each chapter in a separate variable. We can then review the output to confirm the accurate separation of chapters.

```

> index_chp1 <- which(Book == "Chapter I", arr.ind=TRUE)
> index_chp2 <- which(Book == "Chapter II", arr.ind=TRUE)
> index_chp3 <- which(Book == "Chapter III", arr.ind=TRUE)
> index_chp4 <- which(Book == "Chapter IV", arr.ind=TRUE)
> index_chp5 <- which(Book == "Chapter V", arr.ind=TRUE)
> index_chp6 <- which(Book == "Chapter VI", arr.ind=TRUE)
> index_chp7 <- which(Book == "Chapter VII", arr.ind=TRUE)
> index_chp8 <- which(Book == "Chapter VIII", arr.ind=TRUE)
> index_chp9 <- which(Book == "Chapter XI", arr.ind=TRUE)
> index_chp10 <- which(Book == "Chapter X", arr.ind=TRUE)
> index_chp11 <- which(Book == "Chapter XI", arr.ind=TRUE)
> index_chp12 <- which(Book == "Chapter XII", arr.ind=TRUE)
> index_chp13 <- which(Book == "Chapter XIII", arr.ind=TRUE)
> index_chp14 <- which(Book == "Chapter XIV", arr.ind=TRUE)
> index_chp15 <- which(Book == "Chapter XV", arr.ind=TRUE)
> index_chp16 <- which(Book == "Chapter XVI", arr.ind=TRUE)

> chp1 <- Book[(index_chp1+1):(index_chp2-1)]
> chp2 <- Book[(index_chp2+1):(index_chp3-1)]
> chp3 <- Book[(index_chp3+1):(index_chp4-1)]
> chp4 <- Book[(index_chp4+1):(index_chp5-1)]
> chp5 <- Book[(index_chp5+1):(index_chp6-1)]
> chp6 <- Book[(index_chp6+1):(index_chp7-1)]
> chp7 <- Book[(index_chp7+1):(index_chp8-1)]
> chp8 <- Book[(index_chp8+1):(index_chp9-1)]
> chp9 <- Book[(index_chp9+1):(index_chp10-1)]
> chp10 <- Book[(index_chp10+1):(index_chp11-1)]
> chp11 <- Book[(index_chp11+1):(index_chp12-1)]
> chp12 <- Book[(index_chp12+1):(index_chp13-1)]
> chp13 <- Book[(index_chp13+1):(index_chp14-1)]
> chp14 <- Book[(index_chp14+1):(index_chp15-1)]
> chp15 <- Book[(index_chp15+1):(index_chp16-1)]

```

We then proceed to create a directory for the chapters and write each chapter to a text file.

```

> dir.create("Chapters")
> # writing each chapter to a text file
> write.table(chp1, file = 'Chapters/chp1.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp2, file = 'Chapters/chp2.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp3, file = 'Chapters/chp3.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp4, file = 'Chapters/chp4.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp5, file = 'Chapters/chp5.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp6, file = 'Chapters/chp6.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp7, file = 'Chapters/chp7.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp8, file = 'Chapters/chp8.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp9, file = 'Chapters/chp9.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp10, file = 'Chapters/chp10.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp11, file = 'Chapters/chp11.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp12, file = 'Chapters/chp12.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp13, file = 'Chapters/chp13.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp14, file = 'Chapters/chp14.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)
> write.table(chp15, file = 'Chapters/chp15.txt', sep = '\t', row.names = FALSE, col.names = FALSE, quote = FALSE)

```



After separating the chapters, the next step is to create a VCorpus. We use the "VCorpus()" method from the "tm" package to accomplish this task and store the corpus in a variable named "CORP". We can then inspect the structure of the corpus using the str() method.

```
> CORP <- VCorpus(DirSource('./Chapters', ignore.case = TRUE, mode = 'text'))
> str(CORP)
Classes 'VCorpus', 'Corpus'  hidden list of 3
$ content:List of 15
..$ :List of 2
...$ content: chr [1:420] "" "Out to Sea" "" ...
...$ meta   :List of 7
....$ author      : chr(0)
....$ timestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
....$ description : chr(0)
....$ heading    : chr(0)
....$ id         : chr "chp1.txt"
....$ language   : chr "en"
....$ origin     : chr(0)
...- attr(*, "class")= chr "TextDocumentMeta"
...- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
..$ :List of 2
...$ content: chr [1:215] "" "The Fear-Phantom" "" ...
...$ meta   :List of 7
....$ author      : chr(0)
....$ timestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
....$ description : chr(0)
....$ heading    : chr(0)
....$ id         : chr "chp10.txt"
....$ language   : chr "en"
....$ origin     : chr(0)
...- attr(*, "class")= chr "TextDocumentMeta"
...- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
```

In this overview, we observe that our corpus comprises 15 documents, each representing a chapter, with their content reflecting the text they encompass. The content of each document ranges from approximately 190 to 900 characters. Additionally, we can explore other parameters such as "meta," which holds the metadata associated with each document.

Next, we generate initial Document Term Matrices (DTM). In these matrices, each row corresponds to a document, while each column corresponds to a term within the document. This setup enables us to analyze our data using vector and matrix algebra. Additionally, we create a Term Document Matrix (TDM), where rows represent terms and columns represent documents. These matrices are constructed using the DocumentTermMatrix() and TermDocumentMatrix() methods from the "tm" package.

DocumentTermMatrix:

```
> CORP_DTM <- DocumentTermMatrix(CORP)
> CORP_DTM
<<DocumentTermMatrix (documents: 15, terms: 7557)>>
Non-/sparse entries: 17106/96249
Sparsity : 85%
Maximal term length: 22
Weighting : term frequency (tf)
> inspect(CORP_DTM)
<<DocumentTermMatrix (documents: 15, terms: 7557)>>
Non-/sparse entries: 17106/96249
Sparsity : 85%
Maximal term length: 22
Weighting : term frequency (tf)
Sample :
  Terms
Docs and but for from had his that the was with
chp1.txt 98 28 40 29 34 54 55 240 56 32
chp11.txt 125 21 36 25 42 90 23 328 43 45
chp12.txt 89 21 28 12 44 66 37 199 41 21
chp13.txt 159 20 20 27 55 105 33 367 55 39
chp14.txt 104 24 27 23 35 39 39 325 62 26
chp2.txt 96 29 32 17 27 26 34 211 32 29
chp4.txt 74 17 20 20 33 49 29 200 45 29
chp5.txt 95 23 29 12 30 53 31 161 49 24
chp7.txt 126 26 29 28 37 93 35 373 56 37
chp8.txt 268 48 62 39 86 158 61 580 114 59
> str(CORP_DTM)
List of 6
$ i : int [1:17106] 1 1 1 1 1 1 1 1 1 ...
$ j : int [1:17106] 1 2 3 4 5 6 7 8 9 10 ...
$ v : num [1:17106] 1 1 1 1 1 1 1 1 1 ...
$ nrow : int 15
$ ncol : int 7557
$ dimnames:List of 2
..$ Docs : chr [1:15] "chp1.txt" "chp10.txt" "chp11.txt" "chp12.txt" ...
..$ Terms: chr [1:7557] "'ancient'" "'arf'" "'e's'" "'e;" ...
- attr(*, "class")= chr [1:2] "DocumentTermMatrix" "simple_triplet_matrix"
- attr(*, "weighting")= chr [1:2] "term frequency" "tf"
```

TermDocumentMatrix:

```
> CORP_TDM <- TermDocumentMatrix(CORP)
> CORP_TDM
<<TermDocumentMatrix (terms: 7557, documents: 15)>>
Non-/sparse entries: 17106/96249
Sparsity : 85%
Maximal term length: 22
Weighting : term frequency (tf)
> inspect(CORP_TDM)
<<TermDocumentMatrix (terms: 7557, documents: 15)>>
Non-/sparse entries: 17106/96249
Sparsity : 85%
Maximal term length: 22
Weighting : term frequency (tf)
Sample :
  Docs
Terms chp1.txt chp11.txt chp12.txt chp13.txt chp14.txt chp2.txt chp4.txt chp5.txt chp7.txt
and 98 125 89 159 104 96 74 95 126
but 28 21 21 20 24 29 17 23 26
for 40 36 28 20 27 32 20 29 29
from 29 25 12 27 23 17 20 12 28
had 34 42 44 55 35 27 33 30 37
his 54 90 66 105 39 26 49 53 93
that 55 23 37 33 39 34 29 31 35
the 240 328 199 367 325 211 200 161 373
was 56 43 41 55 62 32 45 49 56
with 32 45 21 39 26 29 24 37
> str(CORP_TDM)
List of 6
$ i : int [1:17106] 1 2 3 4 5 6 7 8 9 10 ...
$ j : int [1:17106] 1 1 1 1 1 1 1 1 1 ...
$ v : num [1:17106] 1 1 1 1 1 1 1 1 1 ...
$ nrow : int 15
$ ncol : int 7557
$ dimnames:List of 2
..$ Terms: chr [1:7557] "'ancient'" "'arf'" "'e's'" "'e;" ...
..$ Docs : chr [1:15] "chp1.txt" "chp10.txt" "chp11.txt" "chp12.txt" ...
- attr(*, "class")= chr [1:2] "TermDocumentMatrix" "simple_triplet_matrix"
- attr(*, "weighting")= chr [1:2] "term frequency" "tf"
```

We can see that the DTM and TDM are providing similar results in a different form.

Finding 10 longest words and sentences in each chapter

Before cleaning the corpus let us find 10 longest words and 10 longest sentences from all the chapters and individually from each chapter.

For this we need packages “dplyr”, “tidytext”, and “ggplot2” which we have installed in the beginning.

First we need to convert our corpus into tidy corpus using the tidy method.

```
> tidyCORP <- tidy(CORP)
> tidyCORP
# A tibble: 15 × 8
  author datetimestamp    description heading id    language origin text
  <gl>   <dttm>      <gl>    <gl>    <chr> <gl>    <chr>
1 NA     2024-05-04 20:54:09 NA     NA     chp1.txt en     NA     "\nout to Sea\n\nI...
2 NA     2024-05-04 20:54:09 NA     NA     chp10.txt en    NA     "\n\nThe Fear-Phantom...
3 NA     2024-05-04 20:54:09 NA     NA     chp11.txt en    NA     "\n'King of the Ap...
4 NA     2024-05-04 20:54:09 NA     NA     chp12.txt en    NA     "\nMan's Reason\n\n...
5 NA     2024-05-04 20:54:09 NA     NA     chp13.txt en    NA     "\nHis Own Kind\n\n...
6 NA     2024-05-04 20:54:09 NA     NA     chp14.txt en    NA     "\nAt the Mercy of t...
7 NA     2024-05-04 20:54:09 NA     NA     chp15.txt en    NA     "\n\nThe Forest God\n...
8 NA     2024-05-04 20:54:09 NA     NA     chp2.txt en     NA     "\n\nThe Savage Home\n...
9 NA     2024-05-04 20:54:09 NA     NA     chp3.txt en     NA     "\nLife and Death\n...
10 NA    2024-05-04 20:54:09 NA     NA     chp4.txt en     NA     "\n\nThe Apes\n\nIn...
11 NA    2024-05-04 20:54:09 NA     NA     chp5.txt en     NA     "\n\nThe White Ape\n...
12 NA    2024-05-04 20:54:09 NA     NA     chp6.txt en     NA     "\nJungle Battles\n...
13 NA    2024-05-04 20:54:09 NA     NA     chp7.txt en     NA     "\n\nThe Light of Know...
14 NA    2024-05-04 20:54:09 NA     NA     chp8.txt en     NA     "\n\nThe Tree-top Hunt...
15 NA    2024-05-04 20:54:09 NA     NA     chp9.txt en     NA     "\nChapter XI\n\n\..."
```

Then we find the 10 longest words by tokenizing the words.

```
> CORPWords<-tidyCORP %>% unnest_tokens(word, text) %>% select(id, word) %>% mutate(word_length=nchar(word)) %>% arrange(desc(word_length))
> CORPWords
# A tibble: 47,187 × 3
  id      word    word_length
  <chr>   <chr>     <int>
1 chp11.txt incomprehensible 16
2 chp3.txt responsibilities 16
3 chp10.txt destructiveness 15
4 chp11.txt experimentation 15
5 chp13.txt personification 15
6 chp2.txt notwithstanding 15
7 chp4.txt notwithstanding 15
8 chp5.txt resourcefulness 15
9 chp8.txt destructiveness 15
10 chp9.txt destructiveness 15
# i 47,177 more rows
# i Use `print(n = ...)` to see more rows
```

Here we can see the 10 longest words in all the chapters combined with the word “incomprehensible” being the longest with 16 words.

Similarly, by modifying the unnest_tokens method we find the 10 longest sentences.

```
> CORPSentences<-tidyCORP %>% unnest_tokens(sentence, text, token = "regex", pattern = "(?<!\\b\\p{L}r)\\.") %>% select(id, sentence) %>% mutate(sentence_length=nchar(sentence)) %>% arrange(desc(sentence_length))
> CORPSentences
# A tibble: 1,894 × 3
  id      sentence    sentence_length
  <chr>   <chr>          <int>
1 chp7.txt "\n\nfrom this primitive function has arisen, unquestionably, all t..." 633
2 chp3.txt "\n\nthe last entry in his diary was made the morning following her..." 610
3 chp8.txt "\n\nmany moons ago, when he had been much smaller, he had desired ..." 528
4 chp1.txt "\n\nfrom the records of the colonial office and from the dead man'..." 523
5 chp7.txt "\n\nsquatting upon his haunches on the table top in the cabin his ..." 507
6 chp3.txt "\n\nthe brilliant birds and the little monkeys had become accustom..." 492
7 chp12.txt "\nman's reason\n\n\nthere was one of the tribe of tarzan who quest..." 468
8 chp8.txt "\n\nand then lord greystoke wiped his greasy fingers upon his nake..." 450
9 chp5.txt "\n\nthe deep waters of the lake he had been taught by his wild mot..." 447
10 chp6.txt "\n\nhe put the book back in the cupboard and closed the door, for ..." 441
# i 1,884 more rows
# i Use `print(n = ...)` to see more rows
```

The regular expression above uses ‘.’ as a delimiter for identifying the longest sentences, assuming the period symbol is not preceded by any ‘r’ ending titles like, Dr., Mr. etc. This regular expression helps correctly split the text into sentences.

We can see that the longest sentence is from Chapter 7 containing a total length of 633 characters.

Now let us compute the 10 longest words and sentences for each chapter. We follow the same process converting each chapter into tidy corpus.

CORPUS Cleansing

Next, the data undergoes a cleaning process to eliminate stop words, numbers, and punctuation marks. Since our text is already in lowercase, we begin by removing quotes from the corpus using the `gsub()` function.

```
> CORP<-tm_map(CORP, content_transformer(gsub), pattern="", replacement="")
```

Following the removal of quotes, we define a function called "Remove_Num_Punct()" specifically designed to eliminate numbers and punctuation from the text. This function can be applied to the CORP corpus, resulting in the creation of our clean corpus named "CORP_clean."

```
> Remove_Num_Punct<-function(x) gsub("[[:alpha:]][[:space:]]*", "", x)
> CORP_clean <- tm_map(CORP, content_transformer(Remove_Num_Punct))
> CORP_clean
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 15
> str(CORP_clean)
Classes 'VCorpus', 'Corpus' hidden list of 3
$ content:List of 15
..$ :List of 2
... ..$ content: chr [1:420] "" "Out to Sea" "" ...
... ..$ meta :List of 7
... ... .$ author : chr(0)
... ... .$ timestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
... ... .$ description : chr(0)
... ... .$ heading : chr(0)
... ... .$ id : chr "chp1.txt"
... ... .$ language : chr "en"
... ... .$ origin : chr(0)
... ... ..- attr(*, "class")= chr "TextDocumentMeta"
... ... ..- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
..$ :List of 2
... ..$ content: chr [1:215] "" "The FearPhantom" "" ...
... ..$ meta :List of 7
... ... .$ author : chr(0)
... ... .$ timestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
... ... .$ description : chr(0)
... ... .$ heading : chr(0)
... ... .$ id : chr "chp10.txt"
... ... .$ language : chr "en"
... ... .$ origin : chr(0)
... ... ..- attr(*, "class")= chr "TextDocumentMeta"
... ... ..- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
..$ :List of 2
... ..$ content: chr [1:427] "" "King of the Apes" "" ...
... ..$ meta :List of 7
... ... .$ author : chr(0)
... ... .$ timestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
... ... .$ description : chr(0)
... ... .$ heading : chr(0)
... ... .$ id : chr "chn11.txt"
> inspect(CORP_clean)
<<VCorpus>>
Metadata: corpus specific: 0, document level (indexed): 0
Content: documents: 15
[[1]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 18862
[[2]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 9650
[[3]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 18652
[[4]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 14802
[[5]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 25186
[[6]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 17761
[[7]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 8427
[[8]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 17017
[[9]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 12204
[[10]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 13355
[[11]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 13961
[[12]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 12301
[[13]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 21204
[[14]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 39060
[[15]]
<<PlainTextDocument>>
Metadata: 7
Content: chars: 9669
```

We can see that we removed 515, 202, 947, 360, 702, 442, 208, 436, 267, 280, 310, 253, 449, 855, 202 characters from respective documents. Now, let us make sure that we convert our corpus into lowercase using tm_map() function.

```
> CORP_low<-tm_map(CORP_clean, tm::content_transformer(tolower))
> CORP_low
<<VCorpus>>
Metadata:  corpus specific: 0, document level (indexed): 0
Content:  documents: 15
> str(CORP_low)
Classes 'VCorpus', 'Corpus'  hidden list of 3
$ content:List of 15
..$ :List of 2
...$ content: chr [1:420] "" "out to sea" "" ...
...$ meta  :List of 7
...$ author   : chr(0)
...$ datetimestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
...$ description : chr(0)
...$ heading   : chr(0)
...$ id        : chr "chp1.txt"
...$ language   : chr "en"
...$ origin    : chr(0)
...- attr(*, "class")= chr "TextDocumentMeta"
...- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
..$ :List of 2
...$ content: chr [1:215] "" "the fearphantom" "" ...
...$ meta  :List of 7
...$ author   : chr(0)
...$ datetimestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
...$ description : chr(0)
...$ heading   : chr(0)
...$ id        : chr "chp10.txt"
...$ language   : chr "en"
...$ origin    : chr(0)
...- attr(*, "class")= chr "TextDocumentMeta"
...- attr(*, "class")= chr [1:2] "PlainTextDocument" "TextDocument"
..$ :List of 2
...$ content: chr [1:427] "" "king of the apes" "" ...
...$ meta  :List of 7
...$ author   : chr(0)
...$ datetimestamp: POSIXlt[1:1], format: "2024-05-04 20:54:09"
...$ description : chr(0)
...$ heading   : chr(0)
...$ id        : chr "chp11.txt"
...$ language   : chr "en"
```

The punctuations, numbers, quotes have been removed and everything is in lower case. Now, lets perform the DTM for the cleaned corpus.

```
> inspect(CORP_DTM)
<<DocumentTermMatrix (documents: 15, terms: 5425)>>
Non-/sparse entries: 14724/66651
Sparsity           : 82%
Maximal term length: 20
Weighting          : term frequency (tf)
Sample             :
Terms
Docs      and but for had him his that the was with
chp1.txt 104 33 40 34 17 54 56 240 57 32
chp11.txt 130 21 36 42 23 90 24 328 43 45
chp12.txt 91 21 29 45 21 67 39 199 42 21
chp13.txt 163 25 21 55 28 105 34 367 55 39
chp14.txt 107 26 28 35 25 39 40 325 62 26
chp2.txt 104 31 32 27 5 26 36 213 32 30
chp4.txt 75 17 20 33 11 49 29 200 45 29
chp5.txt 96 23 29 30 26 54 31 161 50 24
chp7.txt 128 27 30 37 22 93 36 373 56 37
chp8.txt 273 52 62 86 43 159 61 580 115 59
```

We can see that the sparsity has been reduced by 3% from 85 to 82. Although it is still high. So, we need to employ more methods to clean our corpus.

Now, let us remove stop words from the corpus.

```
> StopWords<-c(tm::stopwords("english"))
> StopWords
[1] "i"      "me"     "my"     "myself"  "we"     "our"    "ours"   "ourselves" "you"
[10] "your"   "yours"  "yourself" "yourselves" "he"     "him"    "his"    "himself"   "she"
[19] "her"    "hers"   "herself"  "it"      "its"    "itself" "itself"  "them"    "their"
[28] "theirs"  "themselves" "what"   "which"   "who"    "whom"   "this"   "that"    "these"
[37] "those"   "am"     "is"      "are"     "was"    "were"   "be"     "been"    "being"
[46] "have"   "has"    "had"     "having"  "do"     "does"   "did"    "doing"   "would"
[55] "should"  "could"  "ought"   "i'm"    "you're" "he's"   "she's"  "it's"    "we're"
[64] "they're" "i've"   "you've"  "we've"   "they've" "i'd"   "you'd" "he'd"   "she'd"
[73] "we'd"   "they'd" "i'll"    "you'll"  "he'll"  "she'll" "we'll"  "they'll" "isn't"
[82] "aren't"  "wasn't" "weren't" "hasn't" "haven't" "hadn't" "doesn't" "don't"   "didn't"
[91] "won't"   "wouldn't" "shan't"  "shouln't" "can't"   "cannot" "couldn't" "mustn't" "let's"
[100] "that's"  "who's"   "what's"   "here's"  "there's" "when's" "where's" "why's"   "how's"
[109] "a"       "an"     "the"     "and"     "but"    "if"     "on"     "because" "as"
[118] "until"   "while"  "of"      "at"      "by"     "for"    "with"   "about"   "against"
[127] "between" "into"   "through" "during"  "before"  "after"  "above"  "below"   "to"
[136] "from"   "up"     "down"   "in"     "out"    "on"     "off"    "over"   "under"
[145] "again"   "further" "then"    "once"   "here"   "there"  "when"   "where"   "why"
[154] "how"     "all"    "any"    "both"   "each"   "few"    "more"   "most"   "other"
[163] "some"   "such"   "no"     "nor"   "not"    "only"   "own"    "same"   "so"
[172] "than"   "too"    "very"
```

```
> CORP_StopWords<-tm::tm_map(CORP_low,tm::removeWords,StopWords)
> tm::inspect(CORP_StopWords[[1]])
<<PlainTextDocument>>
Metadata: 7
Content: chars: 13612
```

```
sea

story one business tell
may credit seductive influence old vintage upon
narrator beginning skeptical incredulity
days followed balance strange tale

convivial host discovered told much
prone doubtfulness foolish pride assumed task old
vintage commenced unearthed written evidence form
musty manuscript dry official records british colonial
office support many salient features remarkable
narrative

say story true witness happenings
portrays fact telling
taken fictitious names principal characters quite sufficiently
evidences sincerity belief may true

yellow mildewed pages diary man long dead
records colonial office dovetail perfectly narrative
convivial host give story painstakingly
pieced several various agencies

find credible will least one
acknowledging unique remarkable interesting

records colonial office dead mans diary
learn certain young english nobleman shall call john
clayton lord greystoke commissioned make peculiarly delicate
investigation conditions british west coast african colony
whose simple native inhabitants another european power known
recruiting soldiers native army used solely
```

In this step, we have utilized the tm package to remove various stop words from our corpus. By inspecting the first document, we can gain an understanding of the appearance of our corpus after the removal of stop words. As observed, the stop words have been successfully eliminated, enhancing the analyzability of our corpus.

Now let us create the Term Document Matrix again.

```
> CORP_StopWords_TDM<-tm::TermDocumentMatrix(CORP_StopWords)
> CORP_StopWords_TDM
<<TermDocumentMatrix (terms: 5327, documents: 15)>>
Non-/sparse entries: 13583/66322
Sparsity : 83%
Maximal term length: 20
Weighting : term frequency (tf)
> |
```

We can also find the frequent terms using `findFreqTerms` with a low freq of 5.

```
> Freq_Terms<-tm::findFreqTerms(CORP_StopWords_TDM, lowfreq = 5)
> Freq_Terms
[1] "able"      "accompanied" "accord"     "accorded"   "accustomed" "across"    "act"
[8] "added"     "advantage"   "adventure"  "africa"     "african"    "afternoon" "agility"
[15] "air"       "alarm"      "alice"      "alive"      "almost"     "alone"    "along"
[22] "aloud"     "already"    "also"       "always"    "among"      "amphitheater" "ancient"
[29] "anger"     "animal"     "animals"    "another"   "answer"     "answered"  "antagonist"
[36] "anthropoids" "ape"       "apeman"    "apes"      "apparent"   "apparently" "appeared"
[43] "approached" "approaching" "archer"    "arm"       "armed"     "arms"     "around"
[50] "arrow"      "arrows"     "article"   "asked"     "attack"    "attacking" "attempt"
[57] "attempted"  "attempting" "attention" "aught"     "authority" "away"     "awful"
[64] "babe"       "baby"       "back"      "bad"       "balance"   "band"     "bared"
[71] "battle"     "beach"      "bearing"   "beast"     "beasts"    "beat"     "beating"
[78] "beautiful"  "became"    "become"    "bed"       "began"     "behind"   "belly"
[85] "belongings" "belt"       "beneath"   "bent"     "beside"    "best"     "betokened"
[92] "better"     "beyond"    "bird"      "birds"    "black"     "blackness" "blacks"
[99] "blade"      "bloodthirsty" "bloody"    "blow"      "blows"     "boar"     "boat"
[106] "boats"      "bodies"    "body"      "bones"    "book"      "books"    "bore"
[113] "born"       "bosom"     "bow"       "box"      "boxes"    "boy"     "brain"
[120] "branch"     "branches"   "breast"    "bright"   "brilliant" "bring"    "british"
[127] "broke"      "broken"    "brothers"  "brought"  "brown"     "brute"    "brutes"
[134] "bugs"       "building"   "built"     "bull"      "bullet"    "cabin"    "call"
[141] "called"     "came"      "can"       "captain"  "carcass"   "care"     "carried"
[148] "carry"      "cast"      "caught"    "cauldron" "caused"    "cautiously" "ceased"
[155] "center"     "chain"     "challenge" "chance"   "change"    "charge"   "charged"
[162] "chatterering" "chest"    "child"    "children" "circle"    "clawing"  "clay"
[169] "clayton"    "claytons"   "clear"    "clearing" "close"    "closed"   "closely"
[176] "-----"     "-----"    "-----"    "-----"   "-----"    "-----"   "-----"
```

```
> length(Freq_Terms) > Freq_Terms[15]
[1] 1114
[1] "air"
[1] 3
```

We can see that there are 1114 words with low freq 5.

We can also check a specific term and the length of that term using `nchar()` function.

Now, we will unlist `tfList` and combine all term frequencies.

```
> tfList<-list()
> for(i in seq_along(CORP_Stopwords)){
+   tfList[[i]]<-termFreq(CORP_Stopwords[[i]])
+ }
> print(tfList[[1]])
  able      aboard      accentuated      accept      accord      accorded      achievement      acknowledging
  1          2          1                  1          1          1          1          1          1
  act      advised      affairs      afford      africa      african      aft      afternoon
  1          1          1          1          4          1          1          1
  agencies      aggregation      agin      ago      alice      almost      alone      along
  1          1          1          1          1          11          4          1          2
  also      amazed      ambition      ammunition      among      amount      anchor      ancient
  1          1          1          1          1          1          1          1
  angry      another      answered      antagonize      anxious      anything      apartment      appalled
  1          2          1          2          1          1          1          1          2
  apparent      appointed      appointment      appreciate      aft      arguments      arm      army
  2          1          1          2          1          1          1          2          3
  around      arrived      aruwimi      asayin      ask      asked      ass      assisting
  1          1          1          1          1          2          4          1          1
  associate      assumed      atlantic      atmosphere      attempt      attempted      attribute      aught
  1          1          2          1          1          1          2          1          2
  aunts      authority      average      averted      avoid      away      awful      back
  1          1          1          1          1          1          3          1          1
  backwards      bad      bags      balance      barkentine      battlefields      battleship      bear
  1          2          1          1          1          1          1          2          2
  beastly      became      beds      beginning      begin      behind      belaying      belief
  1          1          1          1          1          1          2          1
  believed      belongings      beneath      best      betrayed      better      big      billings
  1          1          1          1          2          1          1          2
  bit      black      blacks      blank      blasted      bloody      bloomin      blow
  2          7          1          1          1          1          1          1          2
  board      body      bore      borrowing      boxes      brasses      breath      bright
  2          1          1          1          1          1          1          1          1
  british      brothers      brought      brutal      brutality      brute      built      bullet
  9          1          3          2          2          2          1          1
  bullied      bullies      bully      business      busted      cabin      calculated      call
  1          1          1          2          1          4          1          1
  came      can      candor      captain      captains      captaining      care      career
  3          3          1          22          1          1          3          1
  careful      carelessly      carriage      carried      carry      caused      centered      certain
  1          1          1          2          1          2          1          1
```

Now, let us inspect the TDM again for the cleaned Corpus with no Stop Words.

```
> tm::inspect(CORP_Stopwords_TDM)
<<TermDocumentMatrix (terms: 5327, documents: 15)>>
Non-/sparse entries: 13583/66322
Sparsity : 83%
Maximal term length: 20
Weighting : term frequency (tf)
Sample :
    Docs
Terms   chp1.txt chp11.txt chp12.txt chp13.txt chp14.txt chp2.txt chp4.txt chp5.txt chp7.txt chp8.txt
apes      0       18      11       6      2       0       9      12      12      21
clayton   29       1       0      10      24      18       5       0       0       0
great      4      18      14       9      17       8      11      12      15      28
jungle     0       6       5      20      16       3       2       6      14      25
little     12      12      10      17       3      11      16      24      28      19
man        8       5       7      31      19      10       3       1       2      13
now        7       8       5       6      13       5       5       5      11      10
one        11      15      9      19      10       2      10       8      22      25
tarzan     0      39      43      33      14       0       0      20      22      67
upon      11      20      11      31      20      15      14      13      22      42
|
```

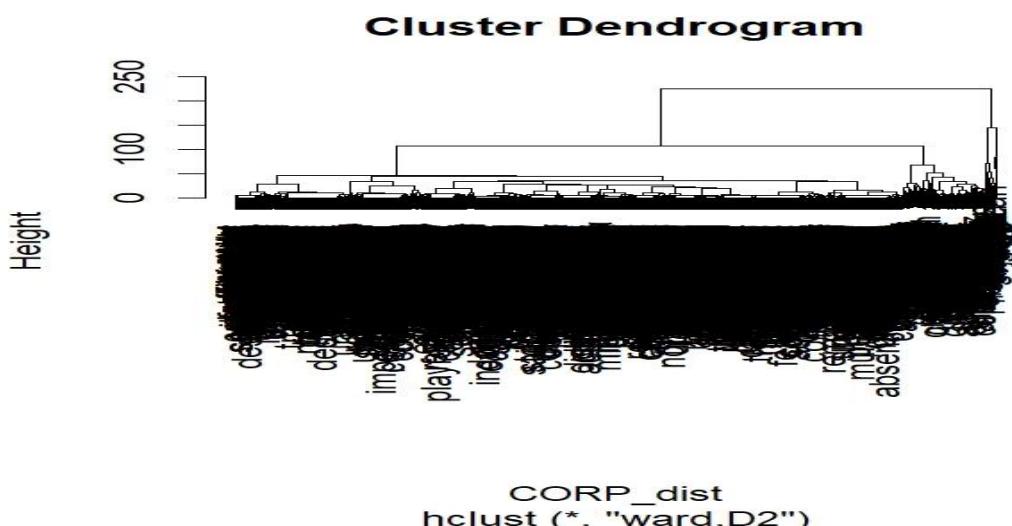
Dendrogram

A dendrogram in R is a graphical representation of hierarchical clustering. Hierarchical clustering is a method of cluster analysis which seeks to build a hierarchy of clusters.

To create a dendrogram in R, you typically use the hclust() function to perform hierarchical clustering on your data and then visualize the clustering using the plot() function with the hclust object as input.

```
> CORP_df <- as.data.frame(as.matrix(CORP_Stopwords_TDM))
> CORP_dist <- dist(CORP_df)
> CORP_dg <- hclust(CORP_dist, method = 'ward.D2')
> str(CORP_dg)
List of 7
 $ merge      : int [1:5326, 1:2] -1 -1481 -1981 -3107 -3157 ...
 $ height     : num [1:5326] 0 0 0 0 0 0 0 0 ...
 $ order      : int [1:5327] 1950 1960 4532 5279 4149 1078 1645 4871 428 557 ...
 $ labels     : chr [1:5327] "abandon" "abandoned" "abandoning" "abashed" ...
 $ method     : chr "ward.D2"
 $ call       : language hclust(d = CORP_dist, method = "ward.D2")
 $ dist.method: chr "euclidean"
 - attr(*, "class")= chr "hclust"
|
```

> plot(CORP_dg)

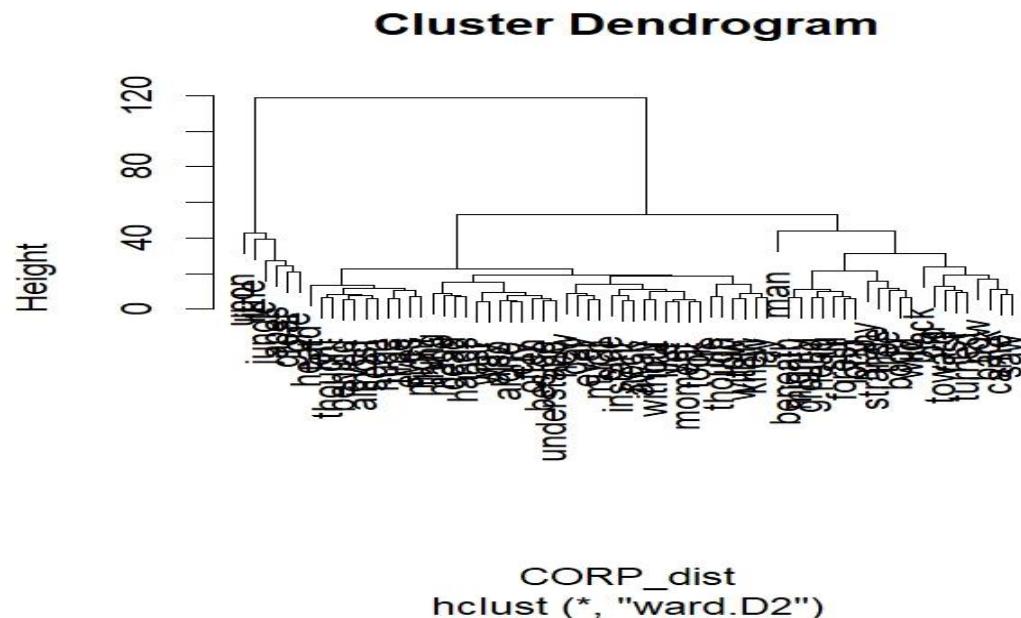


We noticed that our dendrogram has many labels, making analysis difficult. To solve this, we applied the `removeSparseTerms()` function from the `tm` package to eliminate sparse terms from the term-document matrix. After experimenting with various maximal allowed sparsity values, we found that a value of 0.2 produced a comprehensive dendrogram for analysis, resulting in a sparsity level of 7%.

We do not prefer 0% sparsity as that would lose a great percentage of the terms and we want as many relevant terms as possible for our analysis

```
> CORP_StopWords_TDM <- removeSparseTerms(CORP_StopWords_TDM, sparse = 0.2)
> #generating the new dendrogram
> CORP_df <- as.data.frame(as.matrix(CORP_StopWords_TDM))
> CORP_dist <- dist(CORP_df)
> CORP_dg <- hclust(CORP_dist, method = 'ward.D2')
> str(CORP_dg)
List of 7
 $ merge      : int [1:69, 1:2] -2 -30 -34 -43 -9 -3 -27 -5 -7 -17 ...
 $ height     : num [1:69] 4.12 4.12 4.12 4.12 4.8 ...
 $ order      : int [1:70] 65 35 28 4 21 50 25 32 58 7 ...
 $ labels     : chr [1:70] "almost" "alone" "also" "apes" ...
 $ method     : chr "ward.D2"
 $ call       : language hclust(d = CORP_dist, method = "ward.D2")
 $ dist.method: chr "euclidean"
 - attr(*, "class")= chr "hclust"
> tm::inspect((CORP_StopWords_TDM))
<<TermDocumentMatrix (terms: 70, documents: 15)>>
Non-/sparse entries: 978/72
Sparsity          : 7%
Maximal term length: 10
Weighting         : term frequency (tf)
Sample            :
    Docs
Terms   chp11.txt chp12.txt chp13.txt chp14.txt chp2.txt chp3.txt chp4.txt chp5.txt chp7.txt chp8.txt
apes     18        11        6        2        0        3        9        12        12        21
back     7         7        14        7        4        4        4        2         5        18
great    18        14        9        17        8        13        11        12        15        28
jungle    6         5        20        16        3        5        2        6        14        25
little   12        10        17        3        11        20        16        24        28        19
man      5         7        31        19        10        2        3        1         2        13
now      8         5        6        13        5        6        5        5        11        10
one      15        9        19        10        2        9        10        8        22        25
strange   7         6        4         4        0        5        6        0         7        17
upon     20        11        31        20        15        14        14        13        22        42
```

```
> plot(CORP_dg)
```



Word Cloud

A word cloud is a graphical representation of text data, where the size of each word corresponds to its frequency or importance. Words are typically displayed in different sizes and colors, with more prominent words indicating higher frequency or importance.

We need wordcloud package for this which we installed in the beginning.

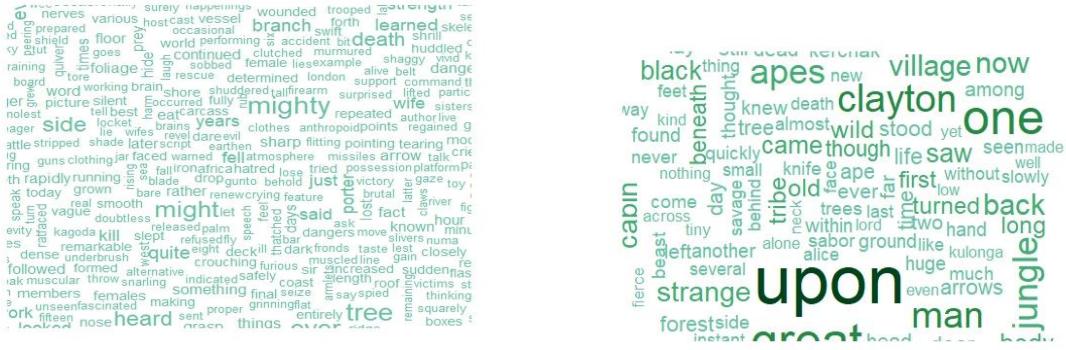
```
> tf_comb <- unlist(tfList)
> tf_summ <- tapply(tf_comb, names(tf_comb), sum)
> Words <- names(tf_summ)
> Words
[1] "abandon"      "abandoned"     "abandoning"    "abashed"      "abated"       "ability"
[7] "able"          "aboard"        "abode"         "abruptly"     "absence"      "absent"
[13] "absentmindedness" "absolute"     "absolutely"    "absorbed"     "abstruse"     "abundance"
[19] "abundant"      "abuses"        "abysmal"       "accentuate"   "accentuated"  "accept"
[25] "accession"      "accident"      "accidental"    "accomplishment" "accompañado"   "accompany"
[31] "accomplish"     "accomplished"  "accomplishment" "accord"        "accorded"      "account"
[37] "accounted"      "accuracy"      "accurately"    "accustom"     "accustomed"   "achievement"
[43] "acknowledging" "acquaintances" "acquirement"   "across"        "act"          "action"
[49] "active"         "add"           "added"         "adding"        "addition"     "additional"
[55] "addressed"     "adds"          "adequate"      "adjusted"     "adjusting"    "admiration"
[61] "admired"        "admit"         "admitted"      "adopted"      "adult"        "adults"
[67] "advance"        "advanced"      "advancing"     "advantage"    "adventure"    "adverbs"
[73] "adversary"      "advised"       "aerial"        "aerie"        "affairs"      "affected"
[79] "affection"      "afford"        "affrighted"   "affront"      "afraid"       "africa"
[85] "african"        "aft"            "afternoon"     "age"          "agencies"     "ages"
[91] "aggregation"   "agile"          "agility"        "agin"         "ago"          "agonized"
[97] "agony"          "agree"          "agreeable"     "agreedthat"   "ahead"        "aim"
[103] "aimed"          "aimlessly"     "aint"          "air"          "airraised"   "alarm"
[109] "alas"           "alert"          "alice"         "alight"       "alighting"   "alike"
[115] "alive"          "alla"           "allay"         "allegorical" "allied"       "allow"
[121] "allowed"        "almost"        "aloft"         "alone"        "along"        "aloud"
[127] "alphabet"       "alphabetical" "already"       "also"         "altardrum"   "alternately"
[133] "alternative"   "always"        "amazed"        "amazing"      "ambition"     "ambitions"
[139] "amen"           "american"     "amidst"        "amiss"        "amity"        "ammunition"
[145] "among"          "amongst"       "amount"        "amphitheater" "amuck"        "amused"
[151] "amusement"     "ancestor"     "ancestors"     "ancestry"     "anchor"       "anchored"
[157] "ancient"        "andoh"         "anger"         "angered"     "angry"        "animal"
[163] "animals"        "animosity"     "ankle"         "anklet"       "anklets"      "announced"
[169] "annoy"          "another"       "anothera"      "answer"       "answered"     "answering"
[175] "answers"        "antagonist"   "antagonize"   "antelope"     "anthropoid"   "anthropoids"
[181] "anticipation"  "anxiety"       "anxious"       "anyone"       "anything"     "anyway"
[187] "apart"          "apartment"    "apathy"        "ape"          "apechild"    "apehood"
[193] "aneman"         "anemans"       "aperture"     "apes"         "apesno"       "aplenty"
```

We gather unique words by using the term frequencies from all 15 documents. By unlisting the tfList variable, we combine all the term frequencies. Then, we sum up the term frequencies for each unique term. Finally, we store the set of words into the words variable.

Next, we utilize the wordcloud() function to create the word cloud. We provide the summed up term frequencies and the unique words as input. However, we encountered warnings indicating that certain words are too lengthy to be accommodated within the word cloud and are consequently omitted. So, we produce another word cloud while setting min.freq = 30, ensuring that only words appearing at least 30 times in the corpus are displayed.

```
> pal <- brewer.pal(9, 'BuGn')
> str(pal)
chr [1:9] "#F7FCFD" "#E5F5F9" "#CCECE6" "#99D8C9" "#66C2A4" "#41AE76" "#238B45" "#006D2C" "#00441B"
> CORP_WC <- wordcloud(Words, tf_summ, colors = pal[-(1:4)])
There were 50 or more warnings (use warnings() to see the first 50)

> CORP_WC <- wordcloud(Words, tf_summ, min.freq = 30, colors = pal[-(1:4)])
```



The second corpus allows us to understand the wordcloud better.

Applying a different pallet



Applying some functions on the Corpus

Following the rubric, we go through several packages like “quanteda”, “syuzhet”, etc. to perform further analysis.

```
install.packages('quanteda')
library('quanteda')
```

We use the CORP_clean copy of the corpus to perform further analysis.
Printing out the first ten lines of the first document:

```
> CORP_text <- CORP_clean[[1]]
> CORP_text$content[1:10]
[1] ""
[2] "out to Sea"
[3] ""
[4] ""
[5] "I had this story from one who had no business to tell it to me or to"
[6] "any other I may credit the seductive influence of an old vintage upon"
[7] "the narrator for the beginning of it and my own skeptical incredulity"
[8] "during the days that followed for the balance of the strange tale"
[9] ""
[10] "When my convivial host discovered that he had told me so much and that"
```

Let us apply a few methods from the quanteda package on our corpus.

First, we apply tokenization. We do this by using the lapply() function to apply tokenization to each of the 15 documents of the corpus.

```
> CORP_tokens <- lapply(CORP_clean,function(x) quanteda::tokens(x$content))
> CORP_tokens
$chp1.txt
Tokens consisting of 420 documents.
text1 :
character(0)

text2 :
[1] "out" "to" "Sea"

text3 :
character(0)

text4 :
character(0)

text5 :
[1] "I"       "had"     "this"    "story"   "from"    "one"     "who"     "had"     "no"      "business"
[11] "to"      "tell"    ...
[ ... and 5 more ]

text6 :
[1] "any"     "other"   "I"       "may"     "credit"  "the"     "seductive" "influence" "of"
[10] "an"      "old"     "vintage" ...
[ ... and 1 more ]

[ reached max_ndoc ... 414 more documents ]

$chp10.txt
Tokens consisting of 215 documents.
text1 :
character(0)

text2 :
[1] "The"     "FearPhantom"

text3 :
character(0)
```

When we examine the structure of CORP_tokens using str(), we see a huge output for all the 15 documents.

Kwic() helps locate keywords in context from the text.

```
> kwic(CORP_tokens[[1]], pattern = "great")
Keyword-in-context with 4 matches.
[text144, 13] fierce black mustachios and a | great |
[text153, 4] blank at the | great | mountain of muscle towering before
[text211, 3] as the | great | lines of a British battleship
[text239, 1] | great | battleship The old fellow was

> kwic(CORP_tokens[[2]], pattern = "great")
Keyword-in-context with 4 matches.
[text45, 2] him | great | respect for the little sharp
[text48, 7] length he came to a | great | tree heavy laden with thick
[text199, 7] glances behind them from their | great | rolling eyes
[text202, 1] | great | tree There was much in

> kwic(CORP_tokens[[3]], pattern = "great")
Keyword-in-context with 18 matches.
[text81, 2] the | great | tree and as before he
[text136, 10] their poor victim to a | great | post near the center of
[text212, 5] of arrows beneath the | great | tree at the end of
[text216, 6] Silently he climbed to a | great | height until he found a
[text245, 5] they had offended some | great | god by placing their village
[text247, 7] was daily placed below the | great | tree from whence the arrows
[text261, 2] The | great | yellow eyes were fixed upon
[text271, 13] one side and as the | great | cat
[text279, 9] Apes went down beneath the | great | body of his enemy but
[text284, 8] he wriggled from beneath the | great | weight and as he
[text295, 1] | great | anthropoids
[text306, 5] Deftly he removed the | great | pelt for he had practiced
[text341, 7] With a frightful roar the | great | beast sprang among the
[text351, 4] Come down Tarzan | great | killer cried Kerchak Come down
[text366, 7] backdrawn snarling lips exposed his | great | fighting fangs and his
[text371, 6] feet of height and his | great | rolling sinews seemed pitifully inadequate
[text403, 8] apeman close to him his | great | jaws sought Tarzans
[text413, 5] a shuddering tremor the | great | body stiffened for an instant
```

We now use the dfm() function to create sparse document-feature matrices for all the chapters. Here, we also remove the stopwords using dfm_remove().

```
> CORP_DFM_Chp1 <- quanteda::dfm(CORP_tokens[[1]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp2 <- quanteda::dfm(CORP_tokens[[2]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp3 <- quanteda::dfm(CORP_tokens[[3]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp4 <- quanteda::dfm(CORP_tokens[[4]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp5 <- quanteda::dfm(CORP_tokens[[5]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp6 <- quanteda::dfm(CORP_tokens[[6]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp7 <- quanteda::dfm(CORP_tokens[[7]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp8 <- quanteda::dfm(CORP_tokens[[8]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp9 <- quanteda::dfm(CORP_tokens[[9]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp10 <- quanteda::dfm(CORP_tokens[[10]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp11 <- quanteda::dfm(CORP_tokens[[11]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp12 <- quanteda::dfm(CORP_tokens[[12]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp13 <- quanteda::dfm(CORP_tokens[[13]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp14 <- quanteda::dfm(CORP_tokens[[14]]) %>% dfm_remove(StopWords)
> CORP_DFM_Chp15 <- quanteda::dfm(CORP_tokens[[15]]) %>% dfm_remove(StopWords)
```

```

> str(CORP_DFM_Chp1)
Formal class 'dfm' [package "quanteda"] with 8 slots
..@ docvars :data.frame': 420 obs. of 3 variables:
... ..$ docname_ : chr [1:420] "text1" "text2" "text3" "text4" ...
... ..$ docid_ : Factor w/ 420 levels "text1","text2",...: 1 2 3 4 5 6 7 8 9 10 ...
... ..$ segid_ : int [1:420] 1 1 1 1 1 1 1 1 1 ...
..@ meta   :List of 3
... ..$ system:List of 5
... ...$ package-version:Classes 'package_version', 'numeric_version' hidden list of 1
... ...$ r-version  :Classes 'R_system_version', 'package_version', 'numeric_version' hidden list of 1
... ...$ system   : Named chr [1:3] "Windows" "x86-64" "anith"
... ...$ attr(*, "names")= chr [1:3] "sysname" "machine" "user"
... ..$ directory : chr "C:/Users/anith/OneDrive/Desktop/GW/sem 2/Big Data/project 3"
... ..$ created   : Date[1:1], format: "2024-05-04"
..@ object:List of 10
... ..$ unit    : chr "documents"
... ..$ what    : chr "word"
... ..$ tokenizer: chr "tokenize_word4"
... ..$ ngram   : int 1
... ..$ skip    : int 0
... ..$ concatenator: chr ""
... ..$ weight_tf :List of 3
... ...$ scheme: chr "count"
... ...$ base   : NULL
... ...$ k      : NULL
... ..$ weight_df :List of 5
... ...$ scheme  : chr "unary"
... ...$ base    : NULL
... ...$ c      : NULL
... ..$ smoothing: NULL
... ..$ threshold: NULL
... ..$ smooth   : num 0
... ..$ summary  :List of 2
... ...$ hash: chr(0)
... ...$ data: NULL
..@ user   : list()
..@ i      : int [1:1713] 1 115 305 4 16 23 50 4 26 53 ...
..@ p      : int [1:1044] 0 37 18 20 21 31 32 33 34 ...
..@ Dim    : int [1:2] 420 1043
..@ dimnames:List of 2
... ..$ : chr "dfm"
... ..$ : chr "corpus"

> str(CORP_DFM_Chp6)
Formal class 'dfm' [package "quanteda"] with 8 slots
..@ docvars :data.frame': 397 obs. of 3 variables:
... ..$ docname_ : chr [1:397] "text1" "text2" "text3" "text4" ...
... ..$ docid_ : Factor w/ 397 levels "text1","text2",...: 1 2 3 4 5 6 7 8 9 10 ...
... ..$ segid_ : int [1:397] 1 1 1 1 1 1 1 1 1 ...
..@ meta   :List of 3
... ..$ system:List of 5
... ...$ package-version:Classes 'package_version', 'numeric_version' hidden list of 1
... ...$ r-version  :Classes 'R_system_version', 'package_version', 'numeric_version' hidden list of 1
... ...$ system   : Named chr [1:3] "Windows" "x86-64" "anith"
... ...$ attr(*, "names")= chr [1:3] "sysname" "machine" "user"
... ..$ directory : chr "C:/users/anith/OneDrive/Desktop/GW/sem 2/Big Daata/project 3"
... ..$ created   : Date[1:1], format: "2024-05-04"
..@ object:List of 10
... ..$ unit    : chr "documents"
... ..$ what    : chr "word"
... ..$ tokenizer: chr "tokenize_word4"
... ..$ ngram   : int 1
... ..$ skip    : int 0
... ..$ concatenator: chr ""
... ..$ weight_tf :List of 3
... ...$ scheme: chr "count"
... ...$ base   : NULL
... ...$ k      : NULL
... ..$ weight_df :List of 5
... ...$ scheme  : chr "unary"
... ...$ base    : NULL
... ...$ c      : NULL
... ..$ smoothing: NULL
... ..$ threshold: NULL
... ..$ smooth   : num 0
... ..$ summary  :List of 2
... ...$ hash: chr(0)
... ...$ data: NULL

```

Next, we get the frequency of terms in the DFM for the chapters.

```
> CORP_Freq_Chp1 <- quanteda::docfreq(CORP_DFM_Chp1)
> CORP_Freq_Chp2 <- quanteda::docfreq(CORP_DFM_Chp2)
> CORP_Freq_Chp3 <- quanteda::docfreq(CORP_DFM_Chp3)
> CORP_Freq_Chp4 <- quanteda::docfreq(CORP_DFM_Chp4)
> CORP_Freq_Chp5 <- quanteda::docfreq(CORP_DFM_Chp5)
> CORP_Freq_Chp6 <- quanteda::docfreq(CORP_DFM_Chp6)
> CORP_Freq_Chp7 <- quanteda::docfreq(CORP_DFM_Chp7)
> CORP_Freq_Chp8 <- quanteda::docfreq(CORP_DFM_Chp8)
> CORP_Freq_Chp9 <- quanteda::docfreq(CORP_DFM_Chp9)
> CORP_Freq_Chp10 <- quanteda::docfreq(CORP_DFM_Chp10)
> CORP_Freq_Chp11 <- quanteda::docfreq(CORP_DFM_Chp11)
> CORP_Freq_Chp12 <- quanteda::docfreq(CORP_DFM_Chp12)
> CORP_Freq_Chp13 <- quanteda::docfreq(CORP_DFM_Chp13)
> CORP_Freq_Chp14 <- quanteda::docfreq(CORP_DFM_Chp14)
> CORP_Freq_Chp15 <- quanteda::docfreq(CORP_DFM_Chp15)
```

```
> str(CORP_Freq_Chp3)
  ...
  Named int [1:1000] 1 18 1 1 2 7 2 3 1 1 ...
  - attr(*, "names")= chr [1:1000] "king" "apes" "yet" "dark" ...
> CORP_Freq_Chp3
      king          apes          yet          dark          reached        tribe        though
      1             18            1             1              2             7             2
      stopped       exhumed       devour       remains         wild           boar          cached
      3             1             1             1              4             1             1
      preceding     day          take          kulongas        bow           arrows         tree
      1             5             1             1              4             10            7
      top           hidden        wellladen      tarzan         dropped        branches        midst
      1             3             1             1              37            2             1
      kerchak       swelling      chest          narrated        glories        adventure      exhibited
      12            2             1             1              1             1             1
      spoils         conquest      grunted        turned          away           jealous        strange
      1             1             1             1              7             1             6
      member         band          little         evil            brain          sought         excuse
      1             1             1             1              12            2             1
      wreak          hatred        upon          next           practicing      first          gleam
      1             3             1             1              20            5             2
      dawn           lost          nearly        every          bolt           shot          finally
      2             1             3             1              3             1             2             1
      learned        guide         shafts         fair           accuracy      ere           month
      5             1             1             1              1              1             5             1
      passed         mean          proficiency    cost           entire         supply        continued
      2             1             1             1              1              1             2             1
      find           hunting       good          vicinity        beach          varied         archery
      3             6             1             1              2              2             1             1
      practice       investigation fathers        choice         store          books          period
      1             2             1             1              1              1             2             1
      young          english       lord          found          back           one          cupboards
      7             2             5             1              5              7             15            1
      cabin          small         metal         box            key           lock          moments
      2             4             2             1              4              2             1             1
      experimentation rewarded     successful    opening        receptacle      faded         photograph
      1             2             2             1              2              1             1             3
      smooth         faced         man           golden        locket         studded        diamonds
      2             2             5             1              1              2             1             1
      reward         gold         chain         buttons        book          examined        minutes
```

Now, let us assign weights to these words.

```
> CORP_Weight_Chp1 <- quanteda::dfm_weight(CORP_DFM_Chp1)
> CORP_Weight_Chp2 <- quanteda::dfm_weight(CORP_DFM_Chp2)
> CORP_Weight_Chp3 <- quanteda::dfm_weight(CORP_DFM_Chp3)
> CORP_Weight_Chp4 <- quanteda::dfm_weight(CORP_DFM_Chp4)
> CORP_Weight_Chp5 <- quanteda::dfm_weight(CORP_DFM_Chp5)
> CORP_Weight_Chp6 <- quanteda::dfm_weight(CORP_DFM_Chp6)
> CORP_Weight_Chp7 <- quanteda::dfm_weight(CORP_DFM_Chp7)
> CORP_Weight_Chp8 <- quanteda::dfm_weight(CORP_DFM_Chp8)
> CORP_Weight_Chp9 <- quanteda::dfm_weight(CORP_DFM_Chp9)
> CORP_Weight_Chp10 <- quanteda::dfm_weight(CORP_DFM_Chp10)
> CORP_Weight_Chp11 <- quanteda::dfm_weight(CORP_DFM_Chp11)
> CORP_Weight_Chp12 <- quanteda::dfm_weight(CORP_DFM_Chp12)
> CORP_Weight_Chp13 <- quanteda::dfm_weight(CORP_DFM_Chp13)
> CORP_Weight_Chp14 <- quanteda::dfm_weight(CORP_DFM_Chp14)
> CORP_Weight_Chp15 <- quanteda::dfm_weight(CORP_DFM_Chp15)
> str(CORP_Weight_Chp3)
Formal class 'dfm' [package "quanteda"] with 8 slots
..@ docvars :data.frame': 427 obs. of 3 variables:
.. ..$ docname_ : chr [1:427] "text1" "text2" "text3" "text4" ...
.. ..$ docid_ : Factor w/ 427 levels "text1","text2",...: 1 2 3 4 5 6 7 8 9 10 ...
.. ..$ segid_ : int [1:427] 1 1 1 1 1 1 1 1 1 ...
..@ meta :List of 3
.. ..$ system:List of 5
.. .. ..$ package-version:Classes 'package_version', 'numeric_version' hidden list of 1
.. .. ..$ : int [1:3] 4 0 2
.. .. ..$ r-version :Classes 'R_system_version', 'package_version', 'numeric_version' hidden list of 1
.. .. ..$ : int [1:3] 4 3 2
.. .. ..$ system : Named chr [1:3] "Windows" "x86-64" "anith"
.. .. .. ..- attr(*, "names")= chr [1:3] "sysname" "machine" "user"
.. .. ..$ directory : chr "C:/Users/anith/OneDrive/Desktop/GW/sem 2/Big Data/project 3"
.. .. ..$ created : Date[1:1], format: "2024-05-04"
.. .. ..$ object:List of 10
.. .. ..$ unit : chr "documents"
.. .. ..$ what : chr "word"
.. .. ..$ tokenizer : chr "tokenize_word4"
.. .. ..$ ngram : int 1
.. .. ..$ skip : int 0
.. .. ..$ concatenator: chr "_"
.. .. ..$ weight_tf :List of 3
.. .. .. ..$ scheme: chr "count"
.. .. .. ..$ base : NULL
.. .. .. ..$ k : NULL
.. .. .. ..$ weight_df :List of 5
.. .. .. .. ..$ scheme : chr "unary"
.. .. .. .. ..$ base : NULL
.. .. .. .. ..$ c : NULL
.. .. .. .. ..$ smoothing: NULL
.. .. .. .. ..$ threshold: NULL
.. .. .. .. ..$ smooth : num 0
.. .. .. .. ..$ summary :List of 2
.. .. .. .. ..$ hash: chr(0)
.. .. .. .. ..$ data: NULL

> CORP_Weight_Chp3
Document-feature matrix of: 427 documents, 1,000 features (99.59% sparse) and 0 docvars.
  features
docs   king apes yet dark reached tribe though stopped exhumed devour
text1    0    0    0    0      0    0    0    0    0    0
text2    1    1    0    0      0    0    0    0    0    0
text3    0    0    0    0      0    0    0    0    0    0
text4    0    0    0    0      0    0    0    0    0    0
text5    0    0    1    1      1    1    1    0    0    0
text6    0    0    0    0      0    0    0    1    1    1
[ reached max_ndoc ... 421 more documents, reached max_nfeat ... 990 more features ]
```

Now, let us compute the tf-idf score. This helps us weigh a dfm by term frequency-inverse document frequency (tf-idf), with full control over options.

```
> CORP_TFIDF_Chp1 <- quanteda::dfm_tfidf(CORP_DFM_Chp1,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp2 <- quanteda::dfm_tfidf(CORP_DFM_Chp2,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp3 <- quanteda::dfm_tfidf(CORP_DFM_Chp3,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp4 <- quanteda::dfm_tfidf(CORP_DFM_Chp4,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp5 <- quanteda::dfm_tfidf(CORP_DFM_Chp5,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp6 <- quanteda::dfm_tfidf(CORP_DFM_Chp6,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp7 <- quanteda::dfm_tfidf(CORP_DFM_Chp7,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp8 <- quanteda::dfm_tfidf(CORP_DFM_Chp8,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp9 <- quanteda::dfm_tfidf(CORP_DFM_Chp9,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp10 <- quanteda::dfm_tfidf(CORP_DFM_Chp10,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp11 <- quanteda::dfm_tfidf(CORP_DFM_Chp11,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp12 <- quanteda::dfm_tfidf(CORP_DFM_Chp12,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp13 <- quanteda::dfm_tfidf(CORP_DFM_Chp13,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp14 <- quanteda::dfm_tfidf(CORP_DFM_Chp14,scheme_tf = "count", scheme_df="inverse")
> CORP_TFIDF_Chp15 <- quanteda::dfm_tfidf(CORP_DFM_Chp15,scheme_tf = "count", scheme_df="inverse")

> str(CORP_TFIDF_Chp3)
Formal class 'dfm' [package "quanteda"] with 8 slots
..@ docvars :data.frame': 427 obs. of 3 variables:
.. ..$ docname_ : chr [1:427] "text1" "text2" "text3" "text4" ...
.. ..$ docid_ : Factor w/ 427 levels "text1","text2",...: 1 2 3 4 5 6 7 8 9 10 ...
.. ..$ segid_ : int [1:427] 1 1 1 1 1 1 1 1 1 ...
..@ meta :List of 3
.. ..$ system:List of 5
.. .. ..$ package-version:Classes 'package_version', 'numeric_version' hidden list of 1
.. .. ..$. : int [1:3] 4 0 2
.. .. ..$ r-version :Classes 'R_system_version', 'package_version', 'numeric_version' hidden list of 1
.. .. ..$. : int [1:3] 4 3 2
.. .. ..$ system : Named chr [1:3] "Windows" "x86-64" "anith"
.. .. ..-. attr(*, "names")= chr [1:3] "sysname" "machine" "user"
.. .. ..$ directory : chr "C:/Users/anith/OneDrive/Desktop/GW/sem 2/Big Data/project 3"
.. .. ..$ created : Date[1:1], format: "2024-05-04"
.. ..$ object:List of 10
.. .. ..$ unit : chr "documents"
.. .. ..$ what : chr "word"
.. .. ..$ tokenizer : chr "tokenize_word4"
.. .. ..$ ngram : int 1
.. .. ..$ skip : int 0
.. .. ..$ concatenator: chr "_"
.. .. ..$ weight_tf :List of 3
.. .. .. ..$ scheme: chr "count"
.. .. .. ..$ base : NULL
.. .. .. ..$ k : NULL
.. .. ..$ weight_df :List of 2
.. .. .. ..$ scheme: chr "inverse"
.. .. .. ..$ base : num 10
.. .. ..$ smooth : num 0
.. .. ..$ summary :List of 2
.. .. .. ..$ hash: chr(0)
.. .. .. ..$ data: NULL
.. .. ..$ user : list()
..@ i : int [1:1748] 1 1 26 87 111 131 143 161 238 250 ...
..@ p : int [1:1001] 0 1 19 20 21 23 30 32 35 36 ...
..@ Dim : int [1:2] 427 1000
..@ Dimnames:List of 2
.. ..$ docs : chr [1:427] "text1" "text2" "text3" "text4" ...

> CORP_TFIDF_Chp3
Document-feature matrix of: 427 documents, 1,000 features (99.59% sparse) and 0 docvars.
  features
docs      king      apes      yet      dark reached tribe though stopped exhumed devour
text1 0       0       0       0       0       0       0       0       0       0       0
text2 2.630428 1.375155 0       0       0       0       0       0       0       0       0
text3 0       0       0       0       0       0       0       0       0       0       0
text4 0       0       0       0       0       0       0       0       0       0       0
text5 0       2.630428 2.630428 2.329398 1.78533 2.329398 2.153307 0       0       0
text6 0       0       0       0       0       0       0       0       0       2.630428 2.630428
[ reached max_ndoc ... 421 more documents, reached max_nfeat ... 990 more features ]
```

Now we install the “syuzhet” package and make it a library.

```
> install.packages('syuzhet')
also installing the dependencies 'textshape', 'zoo', 'dtt'

trying URL 'https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/textshape_1.7.5.tgz'
Content type 'application/x-gzip' length 521774 bytes (509 KB)
=====
downloaded 509 KB

trying URL 'https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/zoo_1.8-12.tgz'
Content type 'application/x-gzip' length 1026435 bytes (1002 KB)
=====
downloaded 1002 KB

trying URL 'https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/dtt_0.1-2.tgz'
Content type 'application/x-gzip' length 19738 bytes (19 KB)
=====
downloaded 19 KB

trying URL 'https://cran.rstudio.com/bin/macosx/big-sur-arm64/contrib/4.3/syuzhet_1.0.7.tgz'
Content type 'application/x-gzip' length 3118163 bytes (3.0 MB)
=====
downloaded 3.0 MB

The downloaded binary packages are in
  /var/folders/yn/h35xrlns7ld2b_5gl156brbw0000gn/T//RtmpAMGjzI/downloaded_packages
> library('syuzhet')
>
```

We can extract text as a data frame. Let us apply a few functions from the syuzhet package on our document.

```
> CORP_df <- as.data.frame(CORP_text$content)
> CORP_df
               CORP_text$content
1                         Out to Sea
2
3
4
5     I had this story from one who had no business to tell it to me or to
6     any other I may credit the seductive influence of an old vintage upon
7     the narrator for the beginning of it and my own skeptical incredulity
8     during the days that followed for the balance of the strange tale
9
10    When my convivial host discovered that he had told me so much and that
11    I was prone to doubtfulness his foolish pride assumed the task the old
12    vintage had commenced and so he unearthed written evidence in the form
13    of musty manuscript and dry official records of the British Colonial
14    Office to support many of the salient features of his remarkable
15                           narrative
16
17    I do not say the story is true for I did not witness the happenings
18    which it portrays but the fact that in the telling of it to you I have
19    taken fictitious names for the principal characters quite sufficiently
20    evidences the sincerity of my own belief that it MAY be true
21
22    The yellow mildewed pages of the diary of a man long dead and the
23    records of the Colonial Office dovetail perfectly with the narrative of
24    my convivial host and so I give you the story as I painstakingly
25          pieced it out from these several various agencies
26
27    If you do not find it credible you will at least be as one with me in
28          acknowledging that it is unique remarkable and interesting
29
30    From the records of the Colonial Office and from the dead mans diary
31    we learn that a certain young English nobleman whom we shall call John
32    Clayton Lord Greystoke was commissioned to make a peculiarly delicate
33    investigation of conditions in a British West Coast African Colony from
34    whose simple native inhabitants another European power was known to be
```

Let's get the sentences in TarzanOfTheApes.txt, we first read the file as one large string as shown with the help of `get_text_as_string()` function.

```
420
> #read the text document as a string
> CORP_asString <- get_text_as_string("TarzanOfTheApes.txt")
> CORP_asString
Tarzan of the Apes By Edgar Rice Burroughs           CONTENTS      I Out to Sea      II The Savage H
e Battles    VII The Light of Knowledge    VIII The Tree-top Hunter    IX Man and Man    X The Fear-P
XIV At the Mercy of the Jungle    XV The Forest God    XVI "Most Remarkable"    XVII Burials    XVIII Th
luge of Torture    XXII The Search Party    XXIII Brother Men    XXIV Lost Treasure    XXV The Outpost of t
lusion    Chapter I Out to Sea    I had this story from one who had no business to tell it to me, or to any ot
the beginning of it, and my own skeptical incredulity during the days that followed for the balance of the stra
I was prone to doubtfulness, his foolish pride assumed the task the old vintage had commenced, and so he unear
he British Colonial Office to support many of the salient features of his remarkable narrative. I do not say t
act that in the telling of it to you I have taken fictitious names for the principal characters quite sufficien
ewed pages of the diary of a man long dead, and the records of the Colonial Office dovetail perfectly with the
ed it out from these several various agencies. If you do not find it credible you will at least be as one with
rds of the Colonial Office and from the dead man's diary we learn that a certain young English nobleman, whom w
llicate investigation of conditions in a British West Coast African Colony from whose simple native inhabitants
ch it used solely for the forcible collection of rubber and ivory from the savage tribes along the Congo and th
en were enticed away through the medium of fair and glowing promises, but that few if any ever returned to thei
cks were held in virtual slavery, since after their terms of enlistment expired their ignorance was imposed upo
rve. And so the Colonial Office appointed John Clayton to a new post in British West Africa, but his confident
ack British subjects by the officers of a friendly European power. Why he was sent, is, however, of little mom
ach his destination. Clayton was the type of Englishman that one likes best to associate with the noblest monu
rile man--mentally, morally, and physically. In stature he was above the average height; his eyes were gray, h
ced by his years of army training. Political ambition had caused him to seek transference from the army to the
```

Many lines have been omitted.

Now we get the sentences using the `get_sentences()` function from the `syuzhet` package. It parses the string into individual sentences.

```
> #getting sentences from the string
> CORP_sentences <- get_sentences(CORP_asString)
> CORP_sentences[1:50]
[1] "Tarzan of the Apes By Edgar Rice Burroughs           CONTENTS      I Out to Sea      II The Sav
Jungle Battles    VII The Light of Knowledge    VIII The Tree-top Hunter    IX Man and Man    X The F
nd    XIV At the Mercy of the Jungle    XV The Forest God    XVI "Most Remarkable"    XVII Burials
The Village of Torture    XXII The Search Party    XXIII Brother Men    XXIV Lost Treasure    XXV The Outpo
I Conclusion    Chapter I Out to Sea    I had this story from one who had no business to tell it to me, or to
[2] "I may credit the seductive influence of an old vintage upon the narrator for the beginning of it, and my
nge tale."
[3] "When my convivial host discovered that he had told me so much, and that I was prone to doubtfulness, his
ritten evidence in the form of musty manuscript, and dry official records of the British Colonial Office to sup
[4] "I do not say the story is true, for I did not witness the happenings which it portrays, but the fact that
ters quite sufficiently evidences the sincerity of my own belief that it MAY be true."
[5] "The yellow, mildewed pages of the diary of a man long dead, and the records of the Colonial Office doveta
s I painstakingly pieced it out from these several various agencies."
[6] "If you do not find it credible you will at least be as one with me in acknowledging that it is unique, re
[7] "From the records of the Colonial Office and from the dead man's diary we learn that a certain young Engli
ake a peculiarly delicate investigation of conditions in a British West Coast African Colony from whose simple
ts native army, which it used solely for the forcible collection of rubber and ivory from the savage tribes alo
[8] "The natives of the British Colony complained that many of their young men were enticed away through the m
lies."
[9] "The Englishmen in Africa went even further, saying that these poor blacks were held in virtual slavery, s
eir white officers, and they were told that they had yet several years to serve."
[10] "And so the Colonial Office appointed John Clayton to a new post in British West Africa, but his confident
ack British subjects by the officers of a friendly European power."
[11] "Why he was sent, is, however, of little moment to this story, for he never made an investigation, nor, in
[12] "Clayton was the type of Englishman that one likes best to associate with the noblest monuments of histori
ly, morally, and physically."
[13] "In stature he was above the average height; his eyes were gray, his features regular and strong; his carr
```

We now proceed to extract the sentiments from these sentences. To do this, we utilize the `get_sentiment()` method that takes a vector of strings and one of four methods for sentiment extraction, "syuzhet", "bing", "afinn", and "nrc". We experiment with all of these methods to observe the results

```

L49] "though the fellow's tone was surly, his words were evidently well meant."
[50] "Ere he had scarce finished his little speech he had turned and was limping off toward the forecastle with
> str(CORP_sentences)
chr [1:4312] "Tarzan of the Apes By Edgar Rice Burroughs"
> |

```

CONTENTS

I Out to Sea

II

We first call `get_sentiment()` with the default value of “syuzhet”

```

> #analyzing the sentiment using syuzhet method
> CORP_syuzhet <- get_sentiment(CORP_sentences, "syuzhet")
> CORP_syuzhet
[1]  0.90 -0.25  1.10  2.40  0.10  2.60  1.25  1.15 -3.15  0.55  0.00  3.50  2.00  2.95  1.10  2.75
[17] -0.35 -0.75  1.00  0.00  0.50  0.00  0.00  0.55  0.80  0.00 -2.25 -0.25 -0.25  0.40 -0.50  0.20
[33] -0.15  0.00 -0.50 -1.00 -0.75 -2.15 -0.50 -0.25 -1.05 -1.50 -1.85 -1.00 -1.20  1.00  1.00 -0.50
[49]  0.20 -1.05 -0.30  0.10 -0.85 -2.25 -0.50 -0.40 -1.15 -1.75  0.25  1.10 -1.00  1.65  0.00 -0.25
[65]  0.50  1.05 -0.10  1.15  0.00  0.00  0.00 -0.80  0.80 -1.00  0.00 -1.00  0.40 -1.00 -0.50  0.00
[81]  0.05  0.00  0.40  0.75 -0.75 -1.00 -0.90  0.15  0.00  0.05  0.50 -1.00 -3.90  1.25  0.00 -0.25
[97] -0.75  0.75 -0.50 -0.35 -0.25  0.00 -0.55  0.40 -0.60  0.80  0.80 -1.75 -0.50  0.00 -2.65  1.00
[113]  0.00 -1.65  1.15 -0.10 -1.45  0.25  1.60 -0.10  0.25  0.85  0.50 -0.50 -0.90 -0.90 -0.75 -0.50
[129]  0.80 -0.55  0.00  0.00  1.15 -2.65  0.80  0.00  0.60  0.50 -0.50  0.00  1.30 -0.50  0.80  0.80
[145]  0.00 -1.40 -0.75 -3.00  0.50 -0.25 -0.55 -1.00 -0.50  0.40 -0.40 -0.50 -0.20 -0.85 -0.40 -2.80
[161] -3.20 -0.50 -2.00 -1.00 -0.60 -2.45  0.00  0.00  0.00 -0.75 -0.75  0.50  2.45  0.00  0.00  0.90
[177] -0.50 -2.10 -0.75  0.55 -1.05  0.00  0.00 -0.75  0.55 -2.15  0.65 -0.75 -0.25  1.40  0.00  0.35
[193] -1.50 -2.00  0.75  0.50  0.80 -0.10  0.75  2.55 -0.50 -2.25 -4.10  0.50  2.95  0.10  0.40  0.80
[209]  1.00  0.55  1.35  0.50  0.80 -1.00 -0.85 -0.75 -0.50  1.70 -1.05 -1.65 -1.40  0.00  1.55  0.80
[225]  0.25  0.25 -0.35  0.65  0.25 -1.35 -1.75 -0.75 -1.30 -0.90 -0.80  2.90 -2.00  0.00  0.00  0.25
[241]  1.25 -2.00  0.00 -0.50  3.20 -0.25  0.75  0.00  2.20  1.50 -0.70  1.30  2.00  0.65 -0.50  1.85
[257]  0.55  0.00  0.40  0.40  0.40  1.55  4.15 -3.05  1.30  0.50  0.00  0.00  1.10 -0.25  0.00 -1.60
[273] -1.25 -0.75  1.50  0.80  0.10 -0.40 -2.25  0.00  3.25 -2.50  1.85  2.10  0.05 -0.05  0.00  0.00
[289]  2.00  1.35  0.00  2.05 -0.85  0.75  1.25 -1.50  0.00  1.85  3.25 -1.15 -2.50 -0.40  1.25  0.35
[305] -2.45 -1.05  0.25 -1.25  0.05 -3.35 -1.00  0.40 -1.50  0.50 -2.75 -0.85  0.65  0.00  0.00 -0.45
[321] -2.75 -0.70 -3.00 -0.65  0.00  0.25 -1.50 -0.45 -1.50  1.35  0.00 -1.25 -1.00 -0.60  0.00 -1.75
[337]  3.15 -1.80  0.25  1.50  0.10  0.00  0.50 -0.45 -1.00  0.50 -2.65  2.20 -3.30  1.60  0.85 -0.50
[353]  0.00  0.00 -0.25 -2.85 -8.15 -1.50 -2.05 -1.50 -1.40 -0.75 -1.20 -2.00  0.35 -2.25 -0.40 -3.75
[369] -1.40 -1.95  0.10 -0.90  0.50  0.85 -2.10  0.50  0.40 -0.05  1.80  0.25  0.50 -2.25  0.40 -0.90
[385] -0.25  0.00 -0.50 -0.90  1.35 -0.40  0.00 -0.50 -1.85 -2.85  0.00 -0.85 -0.55 -0.60 -1.50 -1.60
[401] -0.25 -2.00  1.25 -0.05  3.75  0.15 -1.75  0.25 -0.70 -2.30 -2.15 -0.30  0.30  0.35 -2.65 -0.40
[417] -0.05  1.15 -1.40  0.85 -0.50 -1.00 -1.35 -0.40 -3.00 -0.90 -0.80 -0.90  0.85  0.30  0.00 -0.15

```

If we try the “nrc” method, we get the following results.

```

[977] 0.40 0.85 -1.25 0.60 -0.25 -1.60 0.10 -2.75 -0.95 0.15 0.00 -0.25 -0.40 0.50 0.50 0.80
[993] 0.50 0.25 -0.75 -1.50 0.50 -0.65 -1.25 0.00
[ reached getOption("max.print") -- omitted 3312 entries ]
> #analyzing the sentiment using nrc method
> CORP_nrc <- get_sentiment(CORP_sentences, "nrc")
> CORP_nrc
[1] 7 2 2 2 1 4 3 3 -3 -2 1 3 2 3 2 3 0 -1 0 -1 -1 0 -1 1 2 0 -1 1 -1 1 -1 1 -1 1
[33] 0 1 0 -1 -1 -2 -2 -3 1 0 1 0 -2 0 1 -1 -2 0 1 2 1 -4 0 -1 -2 -2 2 4 0 0 -2 1
[65] 1 5 1 2 0 0 0 -1 1 -1 0 -1 1 -1 0 5 1 -1 0 0 1 1 0 1 0 -4 2 -1 1
[97] -1 1 0 0 1 1 -1 1 0 0 0 -2 -1 1 -2 0 1 -1 2 2 1 2 0 0 0 0 1 0 -1 -1 -1
[129] 0 -1 -1 0 0 0 0 0 -1 -1 1 -1 0 3 0 0 1 -3 2 -1 -2 1 -1 3 -1 0 0 -1 0 -3
[161] -3 -1 0 -2 1 -3 1 0 0 -2 -1 1 4 0 0 2 1 0 -1 0 0 0 1 -1 -1 -3 0 0 -1 1 0 2
[193] -3 -6 2 -1 0 1 2 2 -1 -3 -8 1 3 1 -2 0 3 0 0 0 0 -1 -3 -2 -1 3 -1 -3 -2 -1 0 0
[225] 0 1 2 0 -1 -4 -3 -1 -1 0 0 4 -2 0 0 -1 1 -1 -1 0 3 0 1 0 2 1 -4 0 1 0 -1 3
[257] -2 0 -1 0 0 1 7 -2 0 0 -1 0 2 0 -1 1 -3 -1 1 0 2 -1 -2 -1 1 -2 2 2 -1 0 -1 0
[289] 1 0 -1 3 -1 0 1 -4 1 -1 1 0 -2 -1 1 2 -3 -2 0 -2 0 -5 0 1 -2 0 -5 0 0 0 0 1
[321] -4 -2 -2 0 1 0 -3 0 -3 2 0 -1 -2 0 0 0 5 1 1 3 1 0 1 2 1 1 -2 1 -4 2 2 1
[353] 0 0 1 -3 -7 -4 2 -1 -3 -1 1 -1 1 -2 1 -4 0 -2 1 0 2 2 -3 0 2 -1 5 0 0 -2 1 0
[385] -1 2 0 1 0 1 0 -1 0 -5 0 1 -3 0 0 0 -2 -3 3 1 7 -1 -1 2 0 -1 -4 1 0 -1 -2 -1
[417] -1 -1 -2 0 0 -1 -1 0 -3 -1 -2 -3 1 -2 1 2 -1 1 0 0 0 0 -2 -2 4 -1 3 1 -2 1 0 1
[449] 1 0 0 -1 0 1 3 2 3 -1 1 3 4 1 0 0 -2 -2 -3 1 1 1 0 -1 0 1 3 -1 2 1 0
[481] 0 1 -1 -1 -1 0 2 0 2 -1 1 -3 0 -5 0 -1 -2 0 0 1 1 4 0 -2 0 -1 -1 -4 2 0 1 -3
[513] -1 -1 -1 0 -2 1 3 0 1 0 0 1 -2 -1 1 -1 2 1 1 -1 3 1 0 -3 -1 2 -1 -1 -3 -1 2 2
[545] 1 2 0 1 0 -1 2 2 3 -2 0 -2 0 -1 1 3 0 -1 -2 -2 1 1 0 0 0 2 0 -1 1 1 0 1
[577] -1 0 1 -2 -1 0 0 0 -2 -1 -1 -2 -2 0 2 1 -2 -1 -1 -3 0 0 -3 -3 -4 0 2 1 -2 0 1 -4
[609] -1 -4 -5 -3 -1 -1 0 0 0 -1 0 -2 -1 -2 -1 -1 -5 1 1 0 1 1 -3 2 0 0 1 1 1 2 2 1
[641] 2 1 -1 0 0 0 1 -1 2 2 0 2 -1 0 0 2 2 0 1 3 1 1 4 -2 1 0 0 1 6 1 -2 0
[673] 1 -2 1 -2 0 0 1 -6 1 0 1 2 1 2 -3 1 0 -1 0 1 2 -1 1 -1 1 -3 -3 -3 -1 0 -1 -5
[705] 0 -1 -1 0 -2 -2 1 -5 -3 0 0 -1 -1 -4 0 7 -2 1 -1 0 4 -2 -6 -2 -1 0 -1 0 -2 -1 -1
[737] -1 1 2 0 -1 1 -2 1 2 -1 0 -2 -2 -3 -1 0 0 1 1 -1 -1 2 1 3 0 -1 0 0 -2 -1 -7
[769] 2 2 0 0 0 1 0 1 2 1 2 0 0 2 1 0 1 1 2 1 1 4 2 2 0 1 2 0

```

For “afinn” we get the following results.

```

L reached getOption("max.print") -- omitted 3312 entries ]
> #analyzing the sentiment using afinn method
> CORP_afinn <- get_sentiment(CORP_sentences, "afinn")
> CORP_afinn
[1] -2 -3 5 4 2 4 -2 1 -8 0 1 9 6 4 1 10 0 -2 1 0 1 0 0 0
[25] 0 -1 -6 2 0 0 0 -1 0 0 -1 -3 -3 -1 0 4 -1 -3 1 -1 -1 3 1 1
[49] 0 0 -1 2 4 -5 -2 0 -2 0 1 0 0 1 -2 3 3 -1 -1 3 0 0 0 0 0
[73] 1 0 0 0 0 0 -2 0 -4 0 0 0 2 -2 -3 -3 0 3 -2 2 0 -5 3 -1 0
[97] -2 2 -2 -3 3 0 -2 0 0 0 0 -2 0 0 -4 2 0 -2 3 -2 0 5 2 -3
[121] -1 3 0 -1 -5 -2 -2 -2 2 -3 0 0 5 -3 0 0 0 3 1 -1 1 -1 1 0
[145] 0 0 -1 -10 5 0 0 -3 0 -4 0 -2 -4 0 0 -5 -6 0 -3 -4 -2 -3 0 0
[169] 0 -3 -3 2 4 0 0 2 -1 -6 -1 0 -3 0 0 -2 0 -2 3 -2 0 0 0 0 0
[193] 0 -5 3 1 0 4 4 2 -1 -2 -11 0 4 0 0 2 1 0 3 0 0 -2 0 -1
[217] -1 2 -1 -5 -9 0 3 0 2 2 -6 2 1 -3 -3 -2 0 -4 -2 6 -2 0 0 0
[241] 2 -3 0 -3 7 -2 3 0 2 4 -2 2 6 1 -1 5 -1 0 4 -1 0 1 15 -7
[265] 4 4 0 0 3 0 0 0 0 -2 5 0 2 0 -1 3 13 -6 0 6 2 -1 0 0
[289] 5 -1 0 7 2 3 3 -4 3 7 3 0 -7 0 5 3 -5 3 3 -1 -3 -6 -2 0
[313] -5 3 -6 1 0 0 -4 1 -7 -2 -6 -3 2 0 -3 -2 -2 3 0 -2 -5 -2 0 -5
[337] 12 -5 2 3 0 0 1 0 0 5 -1 -1 -3 0 2 -1 0 0 -1 -2 -19 -1 -2 -2
[361] 1 1 -2 -2 -7 -4 3 -5 -2 -3 1 2 3 -1 -4 3 0 -3 9 4 5 -4 0 -5
[385] -3 0 3 -2 2 -3 1 1 -5 -6 -1 0 -5 0 -3 -1 1 3 0 -5 9 1 -1 -1
[409] 1 -3 -3 3 -2 0 1 0 1 7 0 3 -2 0 0 0 -4 -5 0 -6 -2 0 0 -5
[433] -6 -2 0 0 -5 1 0 -4 4 -4 4 -2 -1 2 3 1 3 -2 1 -1 0 1 -1 7
[457] 7 0 2 1 3 2 2 2 3 1 0 -6 1 0 -3 -2 -3 0 0 4 -2 3 2 0
[481] 3 -1 0 -3 1 1 3 4 8 0 0 -5 0 -10 -1 -4 -5 -7 -2 -2 0 2 -2 -4
[505] -3 -5 -1 0 0 4 2 -3 -1 -6 1 0 -4 0 3 1 2 0 3 -5 2 2 12 -4
[529] -1 -3 2 0 1 -2 0 -2 2 0 0 -3 -10 2 1 2 0 3 0 -1 3 3 7 2
[553] 2 -6 0 -3 -3 0 2 0 0 -3 -7 -7 0 -1 2 0 0 4 0 0 2 0 0 0

```

For “bing” we get the following results.

```

> #analyzing the sentiment using bing method
> CORP_bing <- get_sentiment(CORP_sentences,"bing")
> CORP_bing
[1] -3 -2  1  1  0  3 -2  3 -3  1  0  5  3  2  0  3 -2 -1  0  1  1  0  0 -1  0  0 -3 -1 -1  0 -1  0
[33]  0  0 -1 -1 -3 -1  0 -2 -2  0 -2  0  0 -2  1 -1  0 -1  0 -3 -1 -1 -2 -2  0  1 -1  0  0  0
[65]  1  0  0  1  0  0  0  0  0  0  0  0  0 -1  0  1  0  0  1  0 -1 -1  1 -1  0  0 -1 -4  1  0 -1
[97] -1  1 -1 -1  1  0 -1  0  0  1  1 -1 -1  0 -2  0  0 -1  1  1 -3 -1  1 -1  1 -1  1 -1 -2  0 -1  0
[129]  0 -1  0  0  1 -3  1  0  1  1 -1  0  1  0  0  0  0 -1 -2 -4  2  0 -1 -2 -1 -1  0 -1 -1  0  0 -2
[161] -1 -1 -3 -1  0 -4  0  0  0 -1 -1  0  0  0  0  1 -1 -2 -2  1 -1  0  0 -1  1 -2  0 -1  0  1  0 -1
[193] -3 -1  1  2  1  1  1  3 -1 -4 -6  1  2  0  2  0  1  2  3  1  1 -1 -1  0  0  1 -1 -3 -1  0  2  0
[225]  0  1 -2 -1  1  1 -3 -1 -1  1  0  3  0  0  0  1  1 -2  1  0  1 -2  1  0  2  2 -2  1  2  1  0  2
[257]  1  0  2  1  1 -1  4 -1  1  1  0  0  1 -1  0 -2 -2 -1  2  1  0  0 -2  1  5 -1  1  2  0  0  0
[289]  0  0  0  2 -1  1  2 -4  1  2 -1 -2 -1 -2  2  0 -4 -2  0 -1  0 -5  0 -1 -2  1 -3  0 -1  0  0 -1
[321] -1 -1 -3 -2  1  1 -3 -1 -2  2 -1  0 -2  0 -1 -4  4 -3  0  3  0  0  0  0  0  2 -1 -2 -2  0  1 -1
[353]  0  0  0 -4 -9 -1 -1 -2 -1 -1 -2 -2 -3 -3  0 -5 -1 -2  0 -1  0  1 -3  1  0 -2  4  1  0 -2  0 -2
[385]  0  0 -1 -1  1 -1  0 -5 -4 -2  0 -2 -2  0 -2 -2  0 -1 -1 -3  4  1 -2  0 -1 -4 -4  0 -1 -1 -1  0
[417] -1  1 -2  0 -1  0 -2  0 -3 -2  0  0  0  1  0 -1 -4 -1 -2 -1 -2  0  0 -3  3 -1 -1  0 -1 -1  1  1
[449]  1 -1  1 -1  0  0  0  3  2  0  1 -1  5  2  0 -1  1 -2 -1 -3  1  0 -2 -1  0  0 -1  3  0  1  1  0
[481]  1 -1 -1 -1  1  0  1  2  2  0  1 -2  1 -7  1 -1 -3 -1  1  1  0  2 -1 -2 -2 -3 -1 -3  2 -1  1 -1
[513]  0 -2  0  1 -1  0  2  1 -2  0  1 -1  0 -1  2 -2  0 -2  1 -1 -1  1 -1 -1  0  0 -1  0 -3  1  0  2
[545]  0  2  0  1  0  0  2  1  0 -1  0  1 -3  0  1 -2  0 -1 -3 -3  0 -2  0  1  0 -1  0  0  0  0  0
[577] -5  0  0 -1 -1  1  1 -1 -2  0  0 -1 -2  1  2  1 -1 -1 -2 -3  1  1 -2 -3 -2  0 -1  2 -6  1  3 -1
[609] -1 -4 -3 -1 -2 -1 -2 -1  0 -3  0 -2 -1 -4  0  0 -2  3  2  2  0  0 -3  4 -2  0  0 -1  1  1  3
[641] -7 -1 -1  0 -1 -1  0 -1 -3 -1  0  3  0 -1 -2  0  0  0 -1  3  0  2  0 -2  1  0  0  2  3  0 -1 -3
[673]  1 -3  3 -2 -1  0  0 -7  1  2  0  0  0 -2 -2  0  0 -2  0 -2  1  0  0 -1  0 -1 -1 -3 -1 -1 -3 -4
[705] -1 -1  0 -1 -2 -3 -2  1 -4 -3 -1 -2  1  0 -1  0  1 -2  2 -1 -1  1  0 -6 -1  0  0  1  1 -3 -1 -2
[737] -1  0  2 -2  0  0  0  0  1  0  0 -2 -5 -1  0  0  0  1  1 -2 -3  0 -2  0  0 -2  0  1 -1  0 -6
[769] -2 -1  1  1 -1 -1  1  1  1 -1 -1  1 -3  0 -1 -1  1  1  0  0 -1 -4  0 -1 -3 -4 -2 -4  0  0  1  0
[801] -2 -1  0  0 -2 -1  0 -1  1 -1  0  2  0 -1  1  0  2  0  2 -1  0  1  1  0 -1 -1  0 -2  2 -1  1 -1

```

Let us look at the sentiment dictionaries for these four methods

```

> #sentiment dictionary for syuzhet method
> CORP_Dictionary_syuzhet <- get_sentiment_dictionary()
> CORP_Dictionary_syuzhet
      word value
1    abandon -0.75
2  abandoned -0.50
3   abandoner -0.25
4 abandonment -0.25
5    abandons -1.00
6    abducted -1.00
7    abduction -0.50
8   abductions -1.00
9     aberrant -0.60
10   aberration -0.80
11     abhor -0.50
12  abhorred -1.00
13  abhorrent -0.50
14    abhors -1.00
15   abilities  0.60
16     ability  0.50
17    abject -1.00
18    ablaze -0.25
19   abnormal -0.50

```

```

L reached 'max' / getOption("max.print") -- omitted 10248 row
> #sentiment dictionary for nrc method
> CORP_Dictionary_nrc <- get_sentiment_dictionary("nrc")
> CORP_Dictionary_nrc
  lang      word sentiment value
1 english     abba  positive   1
2 english    ability  positive   1
3 english  abovementioned  positive   1
4 english    absolute  positive   1
5 english  absolution  positive   1
6 english   absorbed  positive   1
7 english  abundance  positive   1
8 english   abundant  positive   1
9 english   academic  positive   1
10 english  academy  positive   1
11 english acceptable  positive   1
12 english acceptance  positive   1
13 english accessible  positive   1
14 english accolade  positive   1
15 english accommodation  positive   1
16 english accompaniment  positive   1
17 english accomplish  positive   1
18 english accomplished  positive   1
19 english accomplishment  positive   1
20 english accord  positive   1
21 english accountability  positive   1
22 english accountable  positive   1
23 english accredited  positive   1
24 english   accueil  positive   1
25 english   accurate  positive   1

```

```

> #sentiment dictionary for bing method
> CORP_Dictionary_bing <- get_sentiment_dictionary("bing")
> CORP_Dictionary_bing
      word value
1          a+    1
2        abound    1
3       abounds    1
4     abundance    1
5     abundant    1
6   accessable    1
7     accessible    1
8      acclaim    1
9     acclaimed    1
10    acclamation    1
11      accolade    1
12     accolades    1
13   accommodative    1
14   accomodative    1
15     accomplish    1
16     accomplished    1
17   accomplishment    1
18  accomplishments    1
19      accurate    1
20   accurately    1
21     achievable    1
22     achievement    1
23    achievements    1
24     achievable    1
25      acumen    1
26    adentable    1

```

```

L reached 'max' / getOption("max.print") -- omitted 6289 rows
> #sentiment dictionary for afinn method
> CORP_Dictionary_afinn <- get_sentiment_dictionary("afinn")
> CORP_Dictionary_afinn
      word value
1     abandon -2
2 abandoned -2
3   abandons -2
4    abducted -2
5   abduction -2
6  abductions -2
7      abhor -3
8   abhorred -3
9   abhorrent -3
10     abhors -3
11 abilities  2
12   ability  2
13    aboard  1
14   aborted -1
15    aborts -1
16   absentee -1
17  absentees -1
18    absolve  2
19   absolved  2
20   absolves  2
21  absolving  2
22   absorbed  1
23     abuse -3
24   abused -3
25   abuses -3

```

We can notice that each of these methods have different criteria and rules that assign sentiment to words, and their dictionaries have different lengths.

Let us sum the values of the sentiment vectors in order to get a measure of the overall emotional valence in the text:

```

L reached 'max' / getOption("max.print") -- omitted 2882 rows
> #sum of sentiment score for syuzhet method
> CORP_Sum_syuzhet <- sum(CORP_syuzhet)
> CORP_Sum_syuzhet
[1] -469.75
> #sum of sentiment score for nrc method
> CORP_Sum_nrc <- sum(CORP_nrc)
> CORP_Sum_nrc
[1] 0
> #sum of sentiment score for bing method
> CORP_Sum_bing <- sum(CORP_bing)
> CORP_Sum_bing
[1] -1183
> #sum of sentiment score for afinn method
> CORP_Sum_afinn <- sum(CORP_afinn)
> CORP_Sum_afinn
[1] -909
>

```

We can observe that except “nrc”, the other methods say that the content of the text is negative. All the other three hand, gives a very strong negative number for the sentiment. To do sentiment analysis, it is important to understand the basis of the sentiment dictionary we use.

```

L� J -509
> #mean of sentiment score for syuzhet method
> CORP_Mean_syuzhet <- mean(CORP_syuzhet)
> CORP_Mean_syuzhet
[1] -0.1089402
> #mean of sentiment score for nrc method
> CORP_Mean_nrc <- mean(CORP_nrc)
> CORP_Mean_nrc
[1] 0
> #mean of sentiment score for bing method
> CORP_Mean_bing <- mean(CORP_bing)
> CORP_Mean_bing
[1] -0.2743506
> #mean of sentiment score for afinn method
> CORP_Mean_afinn <- mean(CORP_afinn)
> CORP_Mean_afinn
[1] -0.2108071
>

```

We can then examine the summary() results of the sentiments generated from these dictionaries to understand the distribution of sentiment.

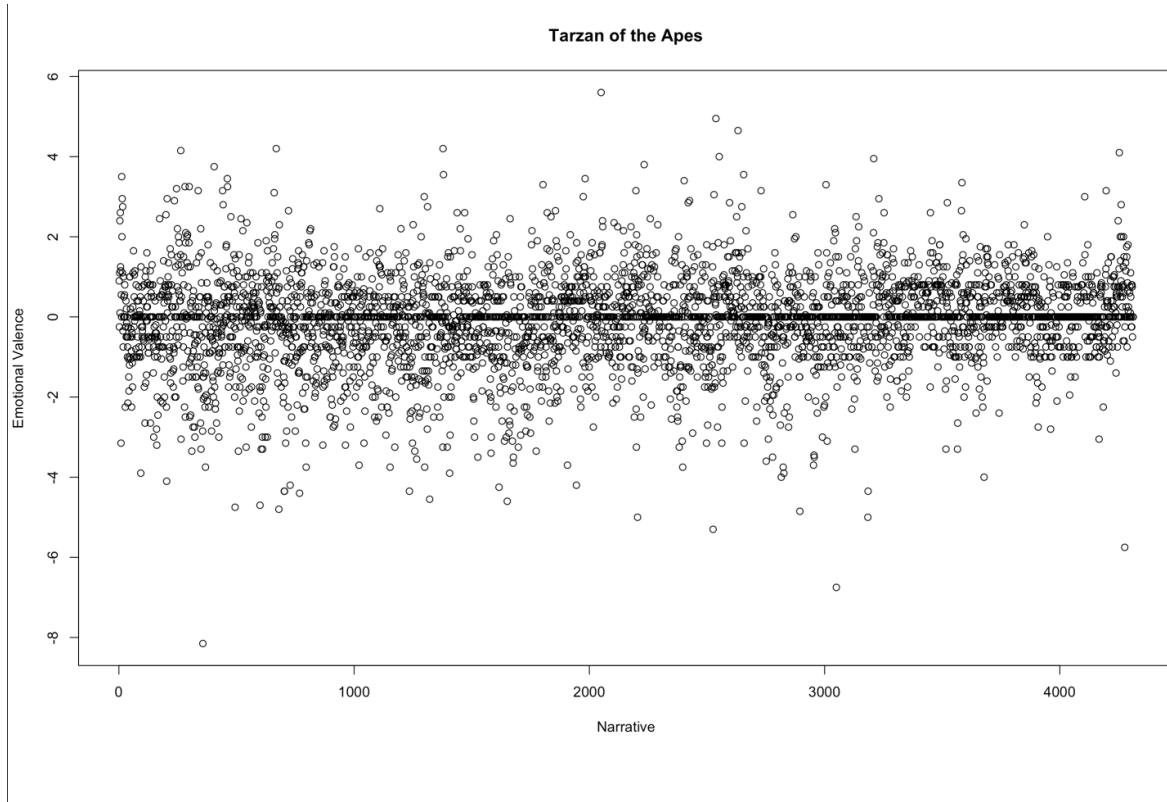
```

> #summaries of the sentiments generated
> summary(CORP_syuzhet)
   Min. 1st Qu. Median      Mean 3rd Qu.      Max.
-8.1500 -0.6000  0.0000 -0.1089  0.5000  5.6000
> summary(CORP_nrc)
   Min. 1st Qu. Median      Mean 3rd Qu.      Max.
     -8       -1        0         0        1        7
> summary(CORP_bing)
   Min. 1st Qu. Median      Mean 3rd Qu.      Max.
-9.0000 -1.0000  0.0000 -0.2744  0.0000  6.0000
> summary(CORP_afinn)
   Min. 1st Qu. Median      Mean 3rd Qu.      Max.
-20.0000 -1.0000  0.0000 -0.2108  1.0000 16.0000
> |

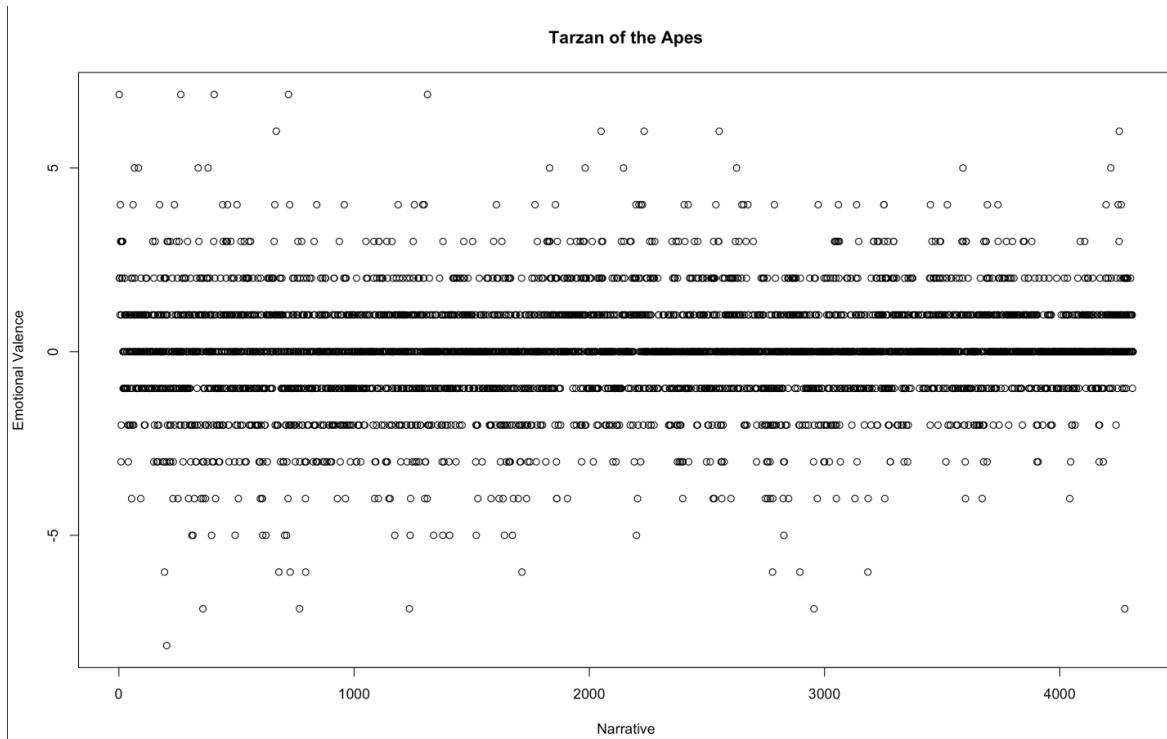
```

However, to further understand how the sentiment of text changes from the beginning of the text till the end, let us plot the values.

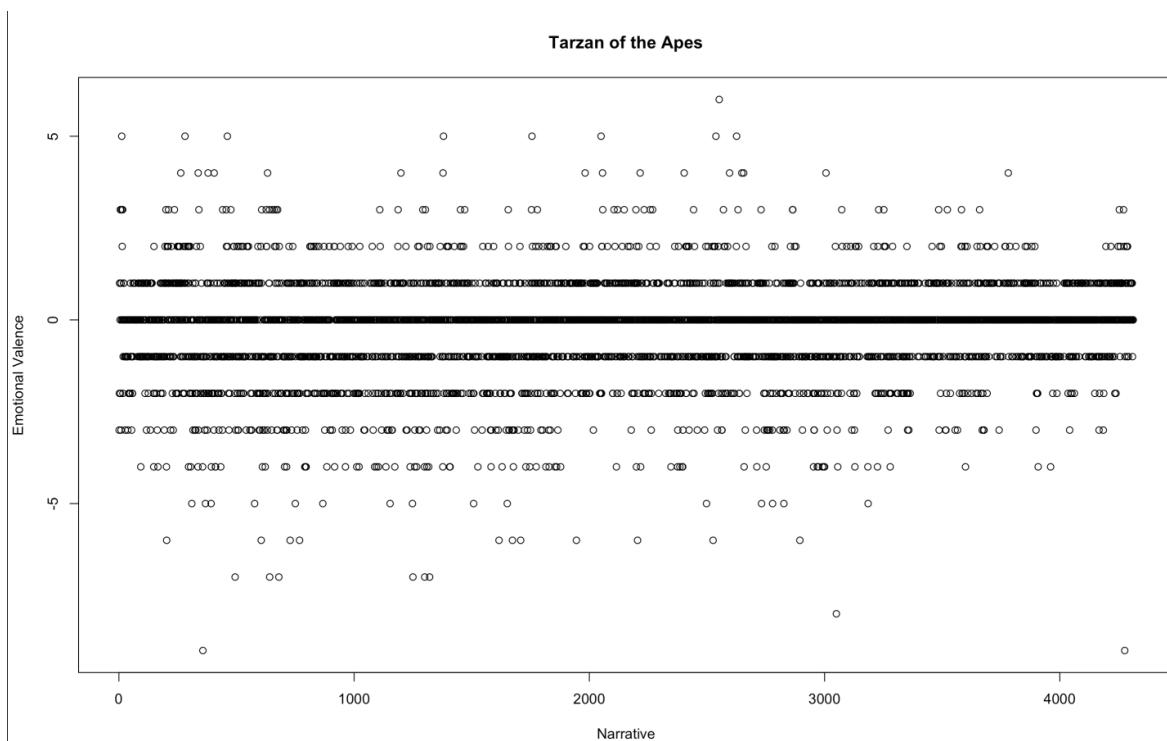
Syuzhet: plot(CORP_syuzhet, main="Tarzan of the Apes", xlab="Narrative", ylab="Emotional Valence")



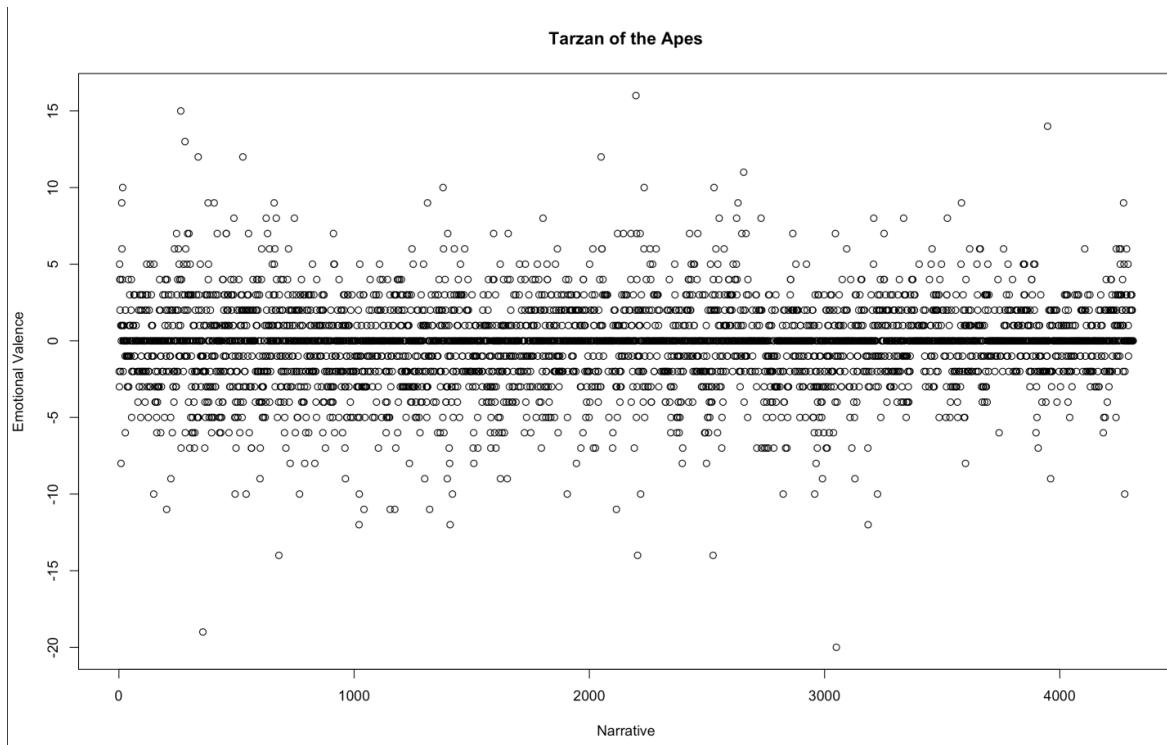
```
NRC: plot(CORP_nrc, main="Tarzan of the Apes", xlab="Narrative", ylab="Emotional Valence")
```



```
Bing plot(CORP_bing, main="Tarzan of the Apes", xlab="Narrative", ylab="Emotional Valence")
```



```
Afinn: plot(CORP_afinn, main="Tarzan of the Apes", xlab="Narrative", ylab="Emotional Valence")
```

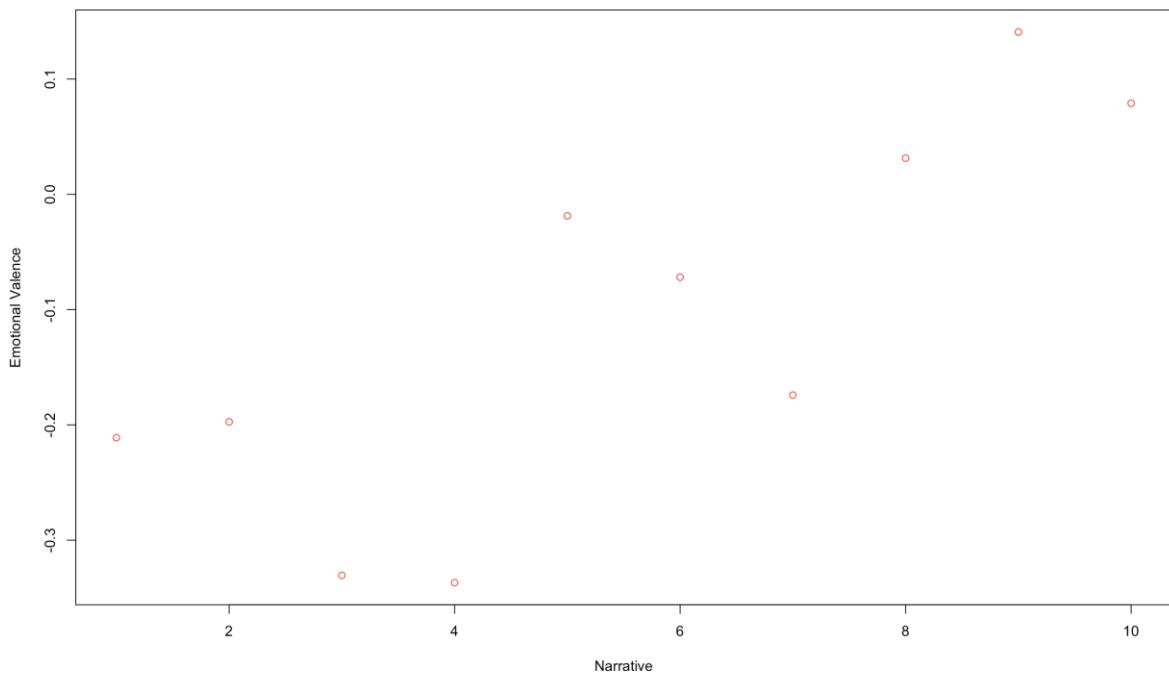


Let us now use the `get_percentage_values()` function from the `syuzhet` package to compare the shape of one trajectory to another. We have plotted values for 10, 20 and 50 bins respectively.

10 bins

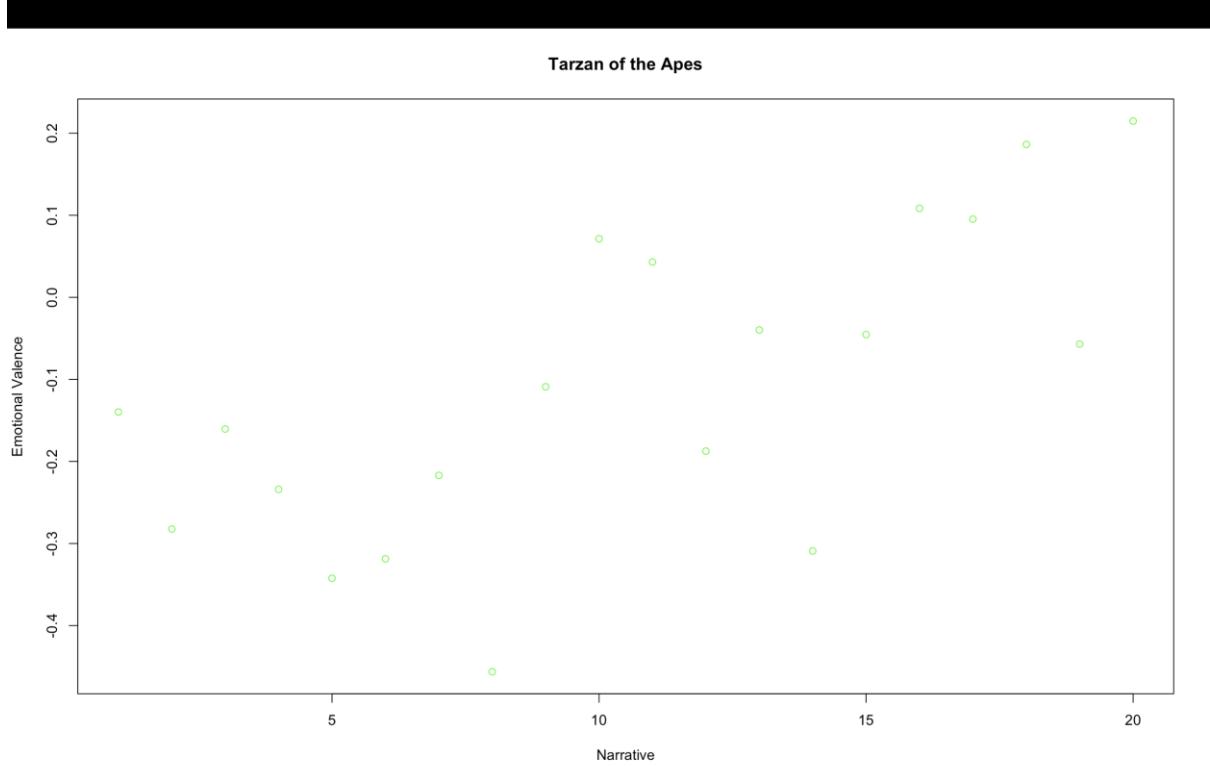
```
> #plotting the distribution of emotional valence, grouped by bins
> CORP_SentimentPctValue <- get_percentage_values(CORP_syuzhet , bins=10)
> structure(CORP_SentimentPctValue)
  1          2          3          4          5          6          7          8
-0.21111111 -0.19733179 -0.33051044 -0.33689095 -0.01867749 -0.07192575 -0.17412993  0.03132251
  9          10
  0.14071926  0.07893519
> plot(CORP_SentimentPctValue, main="Tarzan of the Apes" , xlab="Narrative" , ylab="Emotional Valence" , col = "red")
> |
```

Tarzan of the Apes



20 bins

```
> CORP_SentimentPctValue <- get_percentage_values(CORP_syuzhet , bins=20)
> structure(CORP_SentimentPctValue)
  1      2      3      4      5      6      7      8
-0.13981481 -0.28240741 -0.16046512 -0.23402778 -0.34232558 -0.31875000 -0.21697674 -0.45625000
  9      10     11     12     13     14     15     16
-0.10906977  0.07129630  0.04305556 -0.18744186 -0.03981481 -0.30906977 -0.04537037  0.10837209
 17     18     19     20
 0.09537037  0.18627907 -0.05694444  0.21481481
> plot(CORP_SentimentPctValue, main="Tarzan of the Apes" , xlab="Narrative" , ylab="Emotional Valence" , col="green")
>
```



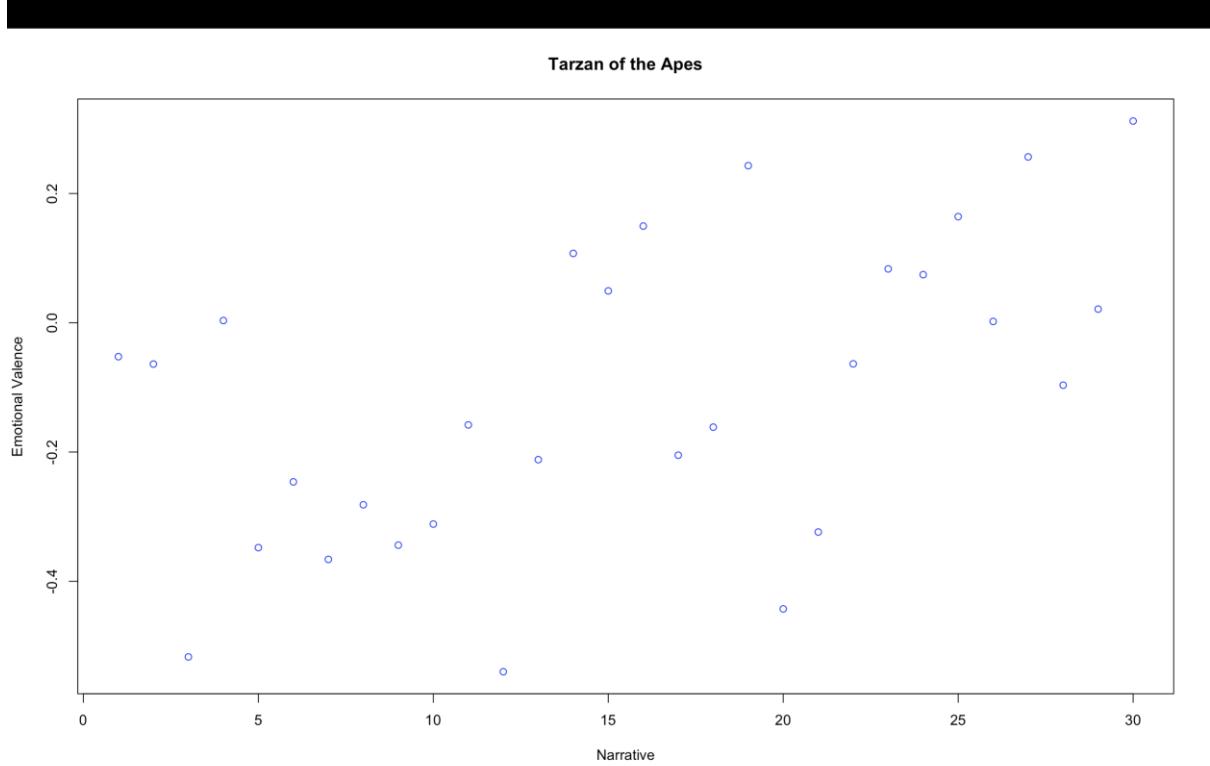
30 bins

```

> CORP_SentimentPctValue <- get_percentage_values(CORP_syuzhet , bins=30)
> structure(CORP_SentimentPctValue)
  1      2      3      4      5      6      7      8
-0.052430556 -0.063888889 -0.517013889  0.003496503 -0.347916667 -0.246180556 -0.366083916 -0.281597222
  9      10     11     12     13     14     15     16
-0.344097222 -0.311458333 -0.158041958 -0.539930556 -0.211805556  0.107342657  0.049305556  0.149652778
 17     18     19     20     21     22     23     24
-0.204895105 -0.161458333  0.243055556 -0.442708333 -0.323776224 -0.063541667  0.083333333  0.074475524
 25     26     27     28     29     30
 0.164236111  0.002083333  0.256643357 -0.096527778  0.021180556  0.312152778
> plot(CORP_SentimentPctValue, main="Tarzan of the Apes" , xlab="Narrative" , ylab="Emotional Valence" , col="blue")
>

```

Fmotions



Through these plots, we understand that the sentiment shifts a lot throughout the text, from being very high in the beginning, to low in the middle, and high again towards the end. This can imply volatile and more emotional content in the book.

Discussion

This project gave us practical, hands-on experience using Text Analytics. Following the rubric, we had a path to follow for this project, yet the challenges along the way, such as dealing with multiple documents, or accommodating large text, reducing sparsity etc. are all concepts we learned while trying to solve a problem. For the same reason, we now understand why Data Science is an empirical science. We learned about data preprocessing and how important it is to clean the data before performing the analysis, because stopwords, punctuation, quote marks, and other elements do not contribute to the analysis and must be removed as they may skew the analysis we are attempting to perform. We learned various packages for Text Analytics and the interesting capabilities they support, including creating wordclouds, retrieving information about the text, assigning sentiments etc. Through all these bits and pieces of analysis, in the end we are able to gain a somewhat rough image of the text we are looking at.

The Text Analytics project was a fantastic learning experience, and as the third project, it built on the information we obtained from the prior ones. It helps to solidify analytical methodologies, cognitive processes, and troubleshooting. It also demonstrated the ease with which R can be used to execute high-quality, comprehensive data analytics.

To summarize, after completing this project, we are confident in working with R, performing text analytics on large datasets, cleaning data, extracting useful information, generating plots, wordclouds, understanding sentiments, and applying different functions and measures/metrics

to documents that allow us to draw different inferences about the problem domain, i.e. the text, as well as simplifying text for better analysis.