

Text Mining of Dan Brown

Shih-Ching, Huang

May 5, 2018

This R Markdown document is made interactive using Shiny.

Make Packages Available

Here I load all the packages necessary for this project. The first thing I need to do is load some packages that I'm going to be using. I use pacman simply to manage packages, and tm is a text mining and that will give us most of our functionality. SnowballC adds some additional text analysis, and dplyr is for manipulating data and for arranging the code using pipes, where the output of one command feeds directly into the input of another one.

```
library(shiny)
library(wordcloud)
library(devtools)
library(tidyverse)
library(stringr)
library(tidytext)
library(dplyr)
library(reshape2)
library(igraph)
library(ggraph)
library(memoise)
if (packageVersion("devtools") < 1.6) {
  install.packages("devtools")
}
pacman::p_load(pacman, tm, SnowballC, dplyr)
```

Import Three Books

The books this project is going to do textmining and sentiment analysis on are three novels by Dan Brown - "Angels & Demons", "The Da Vinci Code", and "The Lost Symbol". I'll start by importing book data, which is the full content of the three books. I have everything in the same directory, so there's no need to give a specific file path. I've already removed the metadata at the beginning and the end of the documents, so all that's left is the novels themselves.

```
# "Angels & Demons" by Dan Brown, published 2000
bookAAD <- readLines('ANGELS AND DEMONS.txt')

# "The Da Vinci Code" by Dan Brown, published 2003
bookDVC <- readLines('The Da Vinci Code.txt')

# "The Lost Symbol" by Dan Brown, published 2009
bookTLS <- readLines('The Lost Symbol.txt')
```

I'll begin by giving the data of every single book respectively, and then I'll compare their features in a set of 2 books and 3 later. First, I'm going to create a Corpus, which is a body of text for each book. I'll begin by creating what I call a preliminary corpus, because I'm going to do some later clean-up on it. These

commands come from tm, for text mining. I'm going to remove the punctuation, any numbers, change everything to lowercase, and remove stopwords. Stopwords are words such as "the", "I", "but", which are usually meaningless when doing text mining.

I'm also going to stem the documents, and what that does is it takes a word like, "Stop" and it takes the variations of it, "Stops, stopped, stopping," and it cuts off those end parts and leaves us with just the beginning, "Stop."

```
# CORPUS FOR ANGELS & DEMONS

# Preliminary corpus
corpusAAD <- Corpus(VectorSource(bookAAD)) %>%
  tm_map(removePunctuation) %>%
  tm_map(removeNumbers) %>%
  tm_map(content_transformer(tolower)) %>%
  tm_map(removeWords, stopwords("english")) %>%
  tm_map(stripWhitespace) %>%
  tm_map(stemDocument)

# Create term-document matrices & remove sparse terms
tdmAAD <- DocumentTermMatrix(corpusAAD) %>%
  removeSparseTerms(1 - (5/length(corpusAAD)))
```

Word Frequencies

Now I'm going to get absolute frequencies for each word, and then relative frequencies

```
# Calculate and sort by word frequencies
word.freqAAD <- sort(colSums(as.matrix(tdmAAD)),
  decreasing = T)

# Create frequency table
tableAAD <- data.frame(word = names(word.freqAAD),
  absolute.frequency = word.freqAAD,
  relative.frequency =
    word.freqAAD/length(word.freqAAD))

# Remove the words from the row names
rownames(tableAAD) <- NULL

# Show the 10 most common words
head(tableAAD, 10)
```

##	word	absolute.frequency	relative.frequency
## 1	angel	211	0.8373016
## 2	god	126	0.5000000
## 3	heaven	71	0.2817460
## 4	earth	62	0.2460317
## 5	bodi	55	0.2182540
## 6	will	51	0.2023810
## 7	one	42	0.1666667
## 8	bibl	42	0.1666667
## 9	know	40	0.1587302
## 10	time	39	0.1547619

As we can see in the table, Dan Brown uses “angel” 211 times and the relative frequency of “angel” is about 0.84.

I’m now going to create a csv file that has the most common words together with their absolute and relative frequencies. The file name will be AAD_1000 in which AAD stands for Angels And Demons. The file will be saved to the same directory where I have my other documents.

```
# Export the 1000 most common words in CSV files
write.csv(tableAAD[1:1000, ], "AAD_1000.csv")
```

I’ll repeat the same steps described above on the other two books in the following codes.

```
# CORPUS FOR THE DA VINCI CODE
```

```
corpusDVC <- Corpus(VectorSource(bookDVC)) %>%
  tm_map(removePunctuation) %>%
  tm_map(removeNumbers) %>%
  tm_map(content_transformer(tolower)) %>%
  tm_map(removeWords, stopwords("english")) %>%
  tm_map(stripWhitespace) %>%
  tm_map(stemDocument)
tdmDVC <- DocumentTermMatrix(corpusDVC) %>%
  removeSparseTerms(1 - (5/length(corpusDVC)))
word.freqDVC <- sort(colSums(as.matrix(tdmDVC)),
  decreasing = T)
tableDVC <- data.frame(word = names(word.freqDVC),
  absolute.frequency = word.freqDVC,
  relative.frequency =
    word.freqDVC/length(word.freqDVC))
rownames(tableDVC) <- NULL
head(tableDVC, 10)
```

##	word	absolute.frequency	relative.frequency
## 1	langdon	1579	0.6082435
## 2	sophi	1127	0.4341294
## 3	teab	601	0.2315100
## 4	said	536	0.2064715
## 5	now	430	0.1656394
## 6	look	418	0.1610169
## 7	fach	398	0.1533128
## 8	one	325	0.1251926
## 9	back	295	0.1136364
## 10	grail	290	0.1117103

```
write.csv(tableDVC[1:1000, ], "DVC_1000.csv")
```

```
# CORPUS FOR THE LOST SYMBOL
```

```
corpusTLS <- Corpus(VectorSource(bookTLS)) %>%
  tm_map(removePunctuation) %>%
  tm_map(removeNumbers) %>%
  tm_map(content_transformer(tolower)) %>%
  tm_map(removeWords, stopwords("english")) %>%
  tm_map(stripWhitespace) %>%
  tm_map(stemDocument)
tdmTLS <- DocumentTermMatrix(corpusTLS) %>%
  removeSparseTerms(1 - (5/length(corpusTLS)))
```

```
word.freqTLS <- sort(colSums(as.matrix(tdmTLS)),
                     decreasing = T)
tableTLS <- data.frame(word = names(word.freqTLS),
                       absolute.frequency = word.freqTLS,
                       relative.frequency =
                         word.freqTLS/length(word.freqTLS))
rownames(tableTLS) <- NULL
head(tableTLS, 10)
```

```
##      word absolute.frequency relative.frequency
## 1  langdon             1365          0.5000000
## 2 katherin              762          0.2791209
## 3   said               696          0.2549451
## 4   now               555          0.2032967
## 5  peter               553          0.2025641
## 6   man               441          0.1615385
## 7   look               437          0.1600733
## 8 pyramid              415          0.1520147
## 9 solomon              394          0.1443223
## 10  one               389          0.1424908
```

```
write.csv(tableTLS[1:1000, ], "TLS_1000.csv")
```

Here's the part where I'll compare their features in a set of 2 books to find out the most distinctive words. I'm going to create one called dProp, which is for a difference in proportions. Now, in this case, I'm simply taking the difference, a subtraction.

"Angels & Demons" vs "The Da Vinci Code"

```
# MOST DISTINCTIVE WORDS #####

# Set number of digits for output
options(digits = 2)

# Compare relative frequencies (via subtraction)
# ("Angels & Demons" vs "The Da Vinci Code")
AADvsDVC <- tableAAD %>%
  merge(tableDVC, by = "word") %>%
  mutate(dProp =
    relative.frequency.x -
    relative.frequency.y,
    dAbs = abs(dProp)) %>%
  arrange(desc(dAbs)) %>%
  rename(AAD.freq = absolute.frequency.x,
         AAD.prop = relative.frequency.x,
         DVC.freq = absolute.frequency.y,
         DVC.prop = relative.frequency.y)

# Show the 10 most distinctive terms
head(AADvsDVC, 10)
```

```
##      word AAD.freq AAD.prop DVC.freq DVC.prop dProp dAbs
## 1  angel      211   0.837      6   0.0023  0.83 0.83
## 2   god      126   0.500     105   0.0404  0.46 0.46
## 3 heaven      71   0.282      22   0.0085  0.27 0.27
## 4  earth      62   0.246      36   0.0139  0.23 0.23
```

```
## 5     bodi      55    0.218      76    0.0293  0.19 0.19
## 6     said      11    0.044     536    0.2065 -0.16 0.16
## 7     bibl      42    0.167      26    0.0100  0.16 0.16
## 8     creat     33    0.131      25    0.0096  0.12 0.12
## 9     will      51    0.202     228    0.0878  0.11 0.11
## 10    lord      28    0.111      18    0.0069  0.10 0.10
```

```
# Save full table to CSV
write.csv(AADvsDVC, "AAD vs DVC.csv")
```

As we can see in the table above, “angel” appears 211 times in Angels & Demons, while it only appears 6 times in The Da Vinci Code, which makes sense given the story. That’s why it has a positive dProp, or difference in proportions. The full table is going to be saved as a csv file in the same directory.

I’ll continue the same steps for the other two sets.

“Angels & Demons” vs “The Lost Symbol”

```
# ("Angels & Demons" vs "The Lost Symbol")
AADvsTLS <- tableAAD %>%
  merge(tableTLS, by = "word") %>%
  mutate(dProp =
    relative.frequency.x -
    relative.frequency.y,
    dAbs = abs(dProp)) %>%
  arrange(desc(dAbs)) %>%
  rename(AAD.freq = absolute.frequency.x,
    AAD.prop = relative.frequency.x,
    TLS.freq = absolute.frequency.y,
    TLS.prop = relative.frequency.y)
head(AADvsTLS, 10)
```

```
##      word AAD.freq AAD.prop TLS.freq TLS.prop dProp dAbs
## 1   angel     211    0.837      13    0.0048  0.83 0.83
## 2    god     126    0.500     156    0.0571  0.44 0.44
## 3 heaven      71    0.282      48    0.0176  0.26 0.26
## 4  earth      62    0.246      63    0.0231  0.22 0.22
## 5   said      11    0.044     696    0.2549 -0.21 0.21
## 6  peter       5    0.020     553    0.2026 -0.18 0.18
## 7   bodi      55    0.218     142    0.0520  0.17 0.17
## 8   bibl      42    0.167      46    0.0168  0.15 0.15
## 9   will      51    0.202     227    0.0832  0.12 0.12
## 10  jesus      31    0.123      24    0.0088  0.11 0.11
```

```
write.csv(AADvsTLS, "AAD vs TLS.csv")
```

“The Da Vinci Code” vs “The Lost Symbol”

```
# ("The Da Vinci Code" vs "The Lost Symbol")
DVCvsTLS <- tableDVC %>%
  merge(tableTLS, by = "word") %>%
  mutate(dProp =
    relative.frequency.x -
    relative.frequency.y,
    dAbs = abs(dProp)) %>%
  arrange(desc(dAbs)) %>%
  rename(DVC.freq = absolute.frequency.x,
    DVC.prop = relative.frequency.x,
```

```
TLS.freq = absolute.frequency.y,
TLS.freq = relative.frequency.y)
head(DVCvsTLS, 10)
```

```
##      word DVC.freq DVC.prop TLS.freq TLS.freq dProp dAbs
## 1   peter      13   0.0050     553   0.2026 -0.198 0.198
## 2  solomon      17   0.0065     394   0.1443 -0.138 0.138
## 3 pyramid      40   0.0154     415   0.1520 -0.137 0.137
## 4   mason      16   0.0062     320   0.1172 -0.111 0.111
## 5 langdon    1579   0.6082    1365   0.5000  0.108 0.108
## 6  church     236   0.0909      11   0.0040  0.087 0.087
## 7 ancient      68   0.0262     260   0.0952 -0.069 0.069
## 8     man     254   0.0978     441   0.1615 -0.064 0.064
## 9 brother      14   0.0054     166   0.0608 -0.055 0.055
## 10 teacher    146   0.0562       9   0.0033  0.053 0.053
```

```
write.csv(DVCvsTLS, "DVC vs TLS.csv")
```

Here's the part where I'll compare their features in a set of 3 books

```
# Three BOOKS DATA
titles <- c("Angels & Demons", "The Da Vinci Code", "The Lost Symbol")

books <- list(bookAAD, bookDVC, bookTLS)

##Each book is an array in which each value in the array is a chapter
series <- tibble()
for(i in seq_along(titles)) {

  temp <- tibble(chapter = seq_along(books[[i]]),
                 text = books[[i]]) %>%
    unnest_tokens(word, text) %>%
    ##Here we tokenize each chapter into words
    mutate(book = titles[i]) %>%
    select(book, everything())

  series <- rbind(series, temp)
}

# set factor to keep books in order of publication
series$book <- factor(series$book, levels = rev(titles))
series
```

```
## # A tibble: 316,249 x 3
##   book      chapter word
##   <fct>      <int> <chr>
## 1 Angels & Demons      2 angels
## 2 Angels & Demons      3 and
## 3 Angels & Demons      3 demons
## 4 Angels & Demons      5 their
## 5 Angels & Demons      5 nature
## 6 Angels & Demons      5 origin
## 7 Angels & Demons      5 mtntsfr
## 8 Angels & Demons      6 ond
## 9 Angels & Demons      6 classification
## 10 Angels & Demons     10 four
## # ... with 316,239 more rows
```

We can get counts for each word using the count function.

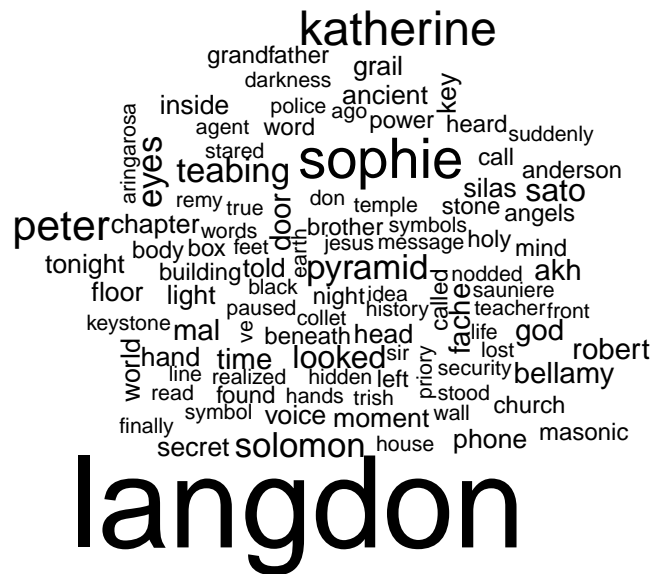
```
series %>% count(word, sort = TRUE)
```

```
## # A tibble: 16,946 x 2
##   word      n
##   <chr> <int>
## 1 the    22156
## 2 of      7647
## 3 to      7110
## 4 a       7011
## 5 and     6517
## 6 in      4475
## 7 he      4167
## 8 was     4044
## 9 his     3493
## 10 had    3064
## # ... with 16,936 more rows
```

Many of the words in the top 10 most common words are stopwords. I'm going to remove the stopwords here and make a word cloud.

```
# Remove stopwords and make a word cloud
series$book <- factor(series$book, levels = rev(titles))
series %>%
  anti_join(stop_words) %>%
  count(word) %>%
  with(wordcloud(word, n, max.words = 100))
```

```
## Joining, by = "word"
```



Sentiment derived from the NRC

```
(hp_nrc <- series %>%  
  inner_join(get_sentiments("nrc")) %>%  
  group_by(book, chapter, sentiment))
```

```
## Joining, by = "word"

## # A tibble: 70,115 x 4
## # Groups:   book, chapter, sentiment [48,688]
##   book          chapter word      sentiment
##   <fct>          <int> <chr>      <chr>
## 1 Angels & Demons      10 radio      positive
## 2 Angels & Demons      32 throne     positive
## 3 Angels & Demons      32 throne     trust
## 4 Angels & Demons      33 elders     positive
## 5 Angels & Demons      33 elders     trust
## 6 Angels & Demons      38 subject    negative
## 7 Angels & Demons      41 spirit      positive
## 8 Angels & Demons      45 heavenly  anticipation
## 9 Angels & Demons      45 heavenly  joy
## 10 Angels & Demons     45 heavenly  positive
## # ... with 70,105 more rows
```


Sentiment analysis using the AFINN dictionary

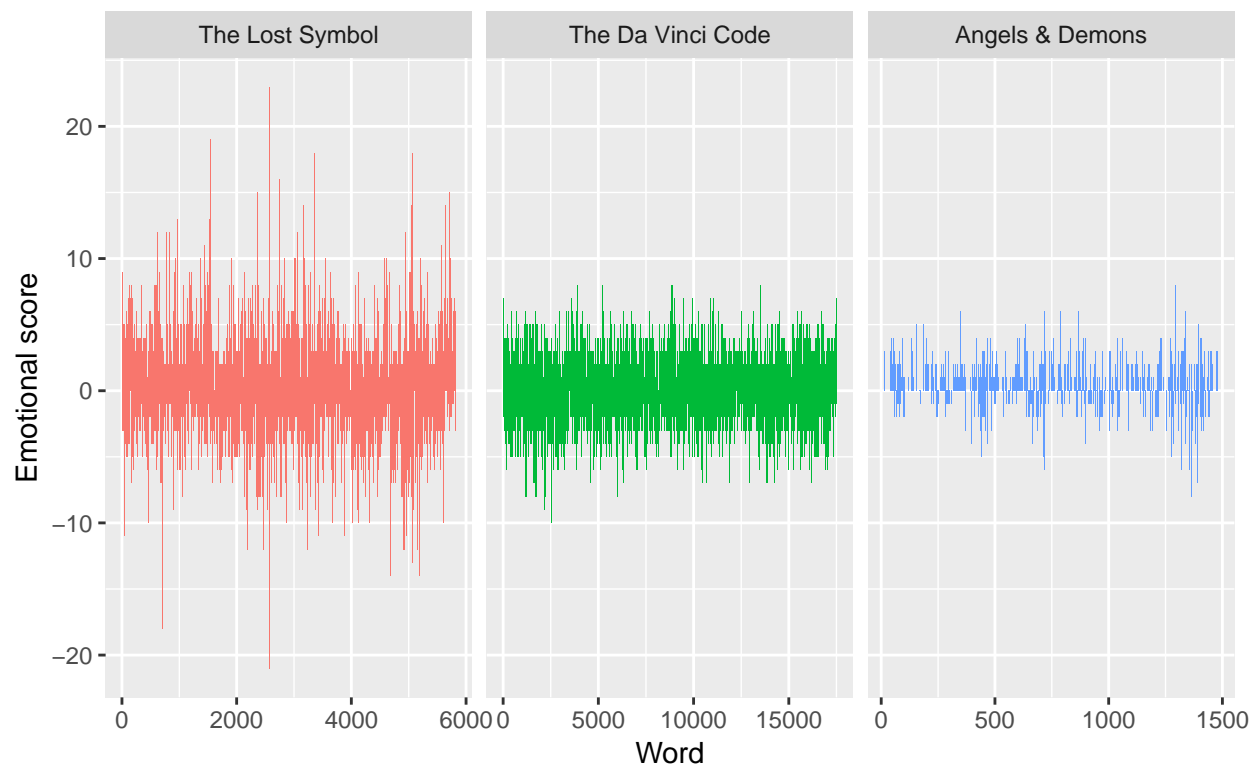
Here I want to visualize the positive/negative sentiment for each book over time using the AFINN dictionary.

```
series %>%  
  inner_join(get_sentiments("afinn")) %>%  
  group_by(book, chapter) %>%  
  summarize(score = sum(score)) %>%  
  ggplot(aes(chapter, score, fill = book)) +  
  geom_col() +  
  facet_wrap(~ book, scales = "free_x") +  
  labs(title = "Emotional Arc of the Three Books",  
       subtitle = "AFINN sentiment dictionary",  
       x = "Word",  
       y = "Emotional score") +  
  theme(legend.position = "none")
```

Joining, by = "word"

Emotional Arc of the Three Books

AFINN sentiment dictionary



Here I'm going to show the cumulative score

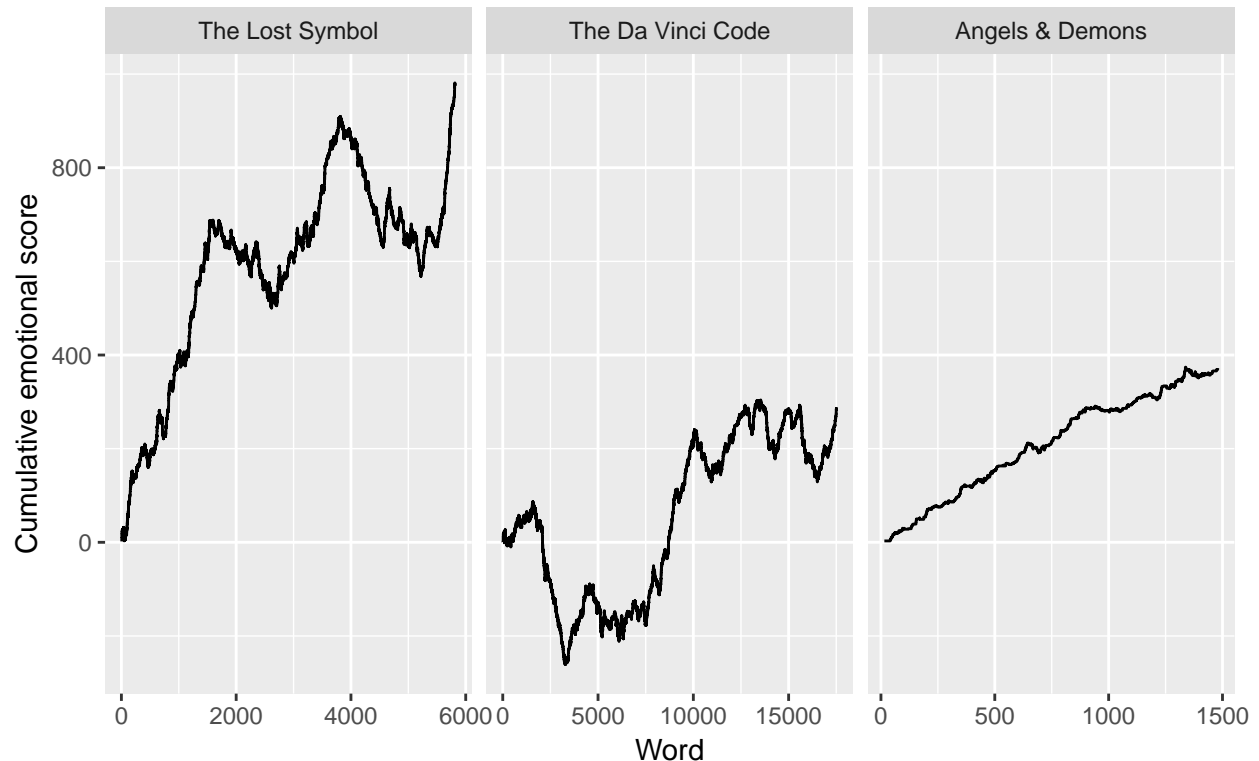
```
# Cumulative score  
series %>%  
  inner_join(get_sentiments("afinn")) %>%  
  group_by(book) %>%  
  mutate(cumscore = cumsum(score)) %>%  
  ggplot(aes(chapter, cumscore, fill = book)) +  
  geom_step() +
```

```
facet_wrap(~ book, scales = "free_x") +
labs(title = "Emotional Arc of the Three Books",
      subtitle = "AFINN sentiment dictionary",
      x = "Word",
      y = "Cumulative emotional score")
```

```
## Joining, by = "word"
```

Emotional Arc of the Three Books

AFINN sentiment dictionary

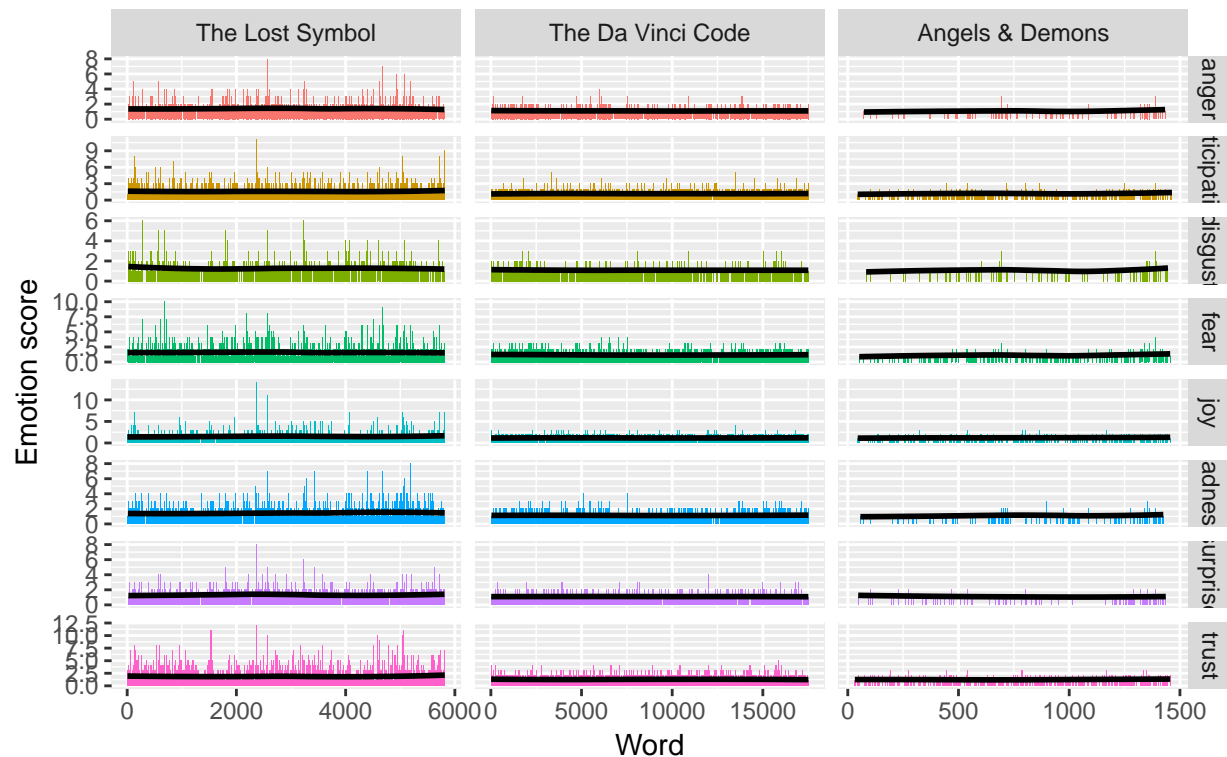


After this, I want to visualize the sentimental content of each chapter in each book using the NRC dictionary.

```
hp_nrc %>%
count(sentiment, book, chapter) %>%
filter(!(sentiment %in% c("positive", "negative"))) %>%
# create area plot
ggplot(aes(x = chapter, y = n)) +
geom_col(aes(fill = sentiment)) +
# add black smoothing line without standard error
geom_smooth(aes(fill = sentiment), method = "loess", se = F, col = 'black') +
theme(legend.position = 'none') +
labs(x = "Word", y = "Emotion score", # add labels
      title = "Emotions during the books",
      subtitle = "Using tidytext and the nrc sentiment dictionary") +
# separate plots per sentiment and book and free up x-axes
facet_grid(sentiment ~ book, scale = "free")
```

Emotions during the books

Using tidytext and the nrc sentiment dictionary



```
### This chunk takes longer to execute
```

Another sentiment analysis

```
series %>%
  right_join(get_sentiments("nrc")) %>%
  filter(!is.na(sentiment)) %>%
  count(sentiment, sort = TRUE)
```

```
## Joining, by = "word"
```

```
## # A tibble: 10 x 2
##   sentiment      n
##   <chr>      <int>
## 1 positive  16767
## 2 negative  12053
## 3 trust     11092
## 4 anticipation 7358
## 5 fear       7126
## 6 joy        5776
## 7 sadness    5462
## 8 anger      4599
## 9 surprise   3764
## 10 disgust   3256
```

The 'bing' lexicon only classifies words as positive or negative.

```
series %>%
  right_join(get_sentiments("bing")) %>%
  filter(!is.na(sentiment)) %>%
  count(sentiment, sort = TRUE)
```

```
## Joining, by = "word"
```

```
## # A tibble: 2 x 2
##   sentiment      n
##   <chr>      <int>
## 1 negative  12203
## 2 positive   9988
```

Next I'm going to Use the the 'bing' lexicon for sentiment analysis and make a comparison cloud.

```
series %>%
  inner_join(get_sentiments("bing")) %>%
  count(word, sentiment, sort = TRUE) %>%
  acast(word ~ sentiment, value.var = "n", fill = 0) %>%
  comparison.cloud(colors = c("#F8766D", "#00BFC4"),
    max.words = 50)
```

```
## Joining, by = "word"
```

negative



positive

The comparison above still contains stopwords, and now I'm going to remove them and make a new cloud. I also use colors o separate words that are positive or negative. We can see here that character names don't appear in the following word cloud, because 'bing' doesn't classify names as positive or negative.

```

series %>%
  anti_join(stop_words) %>%
  inner_join(get_sentiments("bing")) %>%
  count(word, sentiment, sort = TRUE) %>%
  acast(word ~ sentiment, value.var = "n", fill = 0) %>%
  comparison.cloud(colors = c("#F8766D", "#00BFC4"),
                   max.words = 50)

## Joining, by = "word"
## Joining, by = "word"

## Warning in comparison.cloud(., colors = c("#F8766D", "#00BFC4"), max.words
## = 50): perfectly could not be fit on page. It will not be plotted.

## Warning in comparison.cloud(., colors = c("#F8766D", "#00BFC4"), max.words
## = 50): correct could not be fit on page. It will not be plotted.

## Warning in comparison.cloud(., colors = c("#F8766D", "#00BFC4"), max.words
## = 50): eager could not be fit on page. It will not be plotted.

```



Calculating Sentiment Score

```

series %>%
  group_by(book) %>%
  mutate(word_count = 1:n(),
         index = word_count %/% 500 + 1) %>%
  inner_join(get_sentiments("bing")) %>%
  count(book, index = index, sentiment) %>%

```

```

ungroup() %>%
spread(sentiment, n, fill = 0) %>%
mutate(sentiment = positive - negative,
       book = factor(book, levels = titles)) %>%
ggplot(aes(index, sentiment, fill = book)) +
geom_bar(alpha = 0.5, stat = "identity", show.legend = FALSE) +
facet_wrap(~ book, ncol = 2, scales = "free_x")

```

```
## Joining, by = "word"
```



Inputs and Outputs

You can embed Shiny inputs and outputs in your document. Outputs are automatically updated whenever inputs change. This demonstrates how a standard R plot can be made interactive by wrapping it in the Shiny `renderPlot` function. The `selectInput` and `sliderInput` functions create the input widgets used to drive the plot. However, I can't markdown the PDF file with shiny, so I create another RMD file and Rcode file which you can run them directly to see the results. Also, you can see ui. file and server file in the folder. But I don't know why it can't work. (Version Problem maybe)

Conclusion

I use text mining to analyze these three novels and find out sentiment variation, word frequency to see the difference among them.