

Further Free Response Exploration

Jason Cairns

May 3, 2019

1 Background

Adding to the previous exploration, there were subtleties felt worth exploring further; namely, the base bigram count, and the `textrank_keywords` function, as well as their relation.

2 Bigram Count

Let's look at the most common bigrams in the Cancer Society data

```
1 library(tidyverse)
2 library(readxl)
3 library(tidytext)
4
5 cancersoc <- read_excel("../data/raw/Cancer Soc.xlsx")
6
7 cancersoc %>%
8   mutate(id = row_number()) %>%
9   unnest_tokens(bigram, sq6a, token = "ngrams", n = 2) %>%
10  count(bigram, sort = TRUE) %>%
11  head
```

bigram	n
nil	589
to help	284
to give	69
a good	67
help others	65
good cause	61

The previous attempt at assessing bigrams used the following code:

```
1 cancer_soc_bigrams <-
2   cancersoc %>%
3   select(sq6a) %>%
4   mutate(id = row_number()) %>%
5   unnest_tokens(bigram, sq6a, token = "ngrams", n = 2) %>%
6   na.omit()
```

```

7
8 bigrams_separated <- cancer_soc_bigrams %>%
9   separate(bigram, c("word1", "word2"), sep = " ")
10
11 bigrams_filtered <- bigrams_separated %>%
12   filter(!word1 %in% stop_words$word) %>%
13   filter(!word2 %in% stop_words$word)
14
15 bigrams_united <- bigrams_filtered %>%
16   unite(bigram, word1, word2, sep = " ")
17
18 # new bigram counts:
19 bigram_counts <- bigrams_united %>%
20   count(bigram, sort = TRUE)
21
22 head(bigram_counts)

```

bigram	n
salvation army	19
raise money	9
donated money	5
red cross	5
spare time	5
cancer society	4

These results are completely different. On closer inspection, it seems that stopwords removal got rid of a lot of bigrams that would otherwise be useful. This raises the question of whether stopwords removal is so essential when ngrams are needed?

3 Further exploration of `textrank_keywords`

The results of `textrank_keywords` appeared extremely useful, but the lingering question was how they were generated. The documentation of `textrank_keywords` was helpful in this respect, showing that

In order to find relevant keywords, the `textrank` algorithm constructs a word network. This network is constructed by looking which words follow one another. A link is set up between two words if they follow one another, the link gets a higher weight if these 2 words occur more frequently next to each other in the text. On top of the resulting network the 'Pagerank' algorithm is applied to get the importance of each word. The top 1/3 of all these words are kept and are considered relevant. After this, a keywords table is constructed by combining the relevant words together if they appear following one another in the text.