

TEXT MINING WITH R

Research Data Management Support

ABOUT TEXT MINING

- Text mining refers to the process of extracting (*mining*) information and insights from text;
- Text mining can be extremely useful when looking for any sort of pattern, trend, or relationships in large volumes of text data (articles, documents, emails, social media posts, etc);
- The main challenge of text mining is obtaining meaningful information from unstructured and ambiguous material.

R PACKAGES FOR TEXT MINING

```
1 library(tidyverse)
2 library(tidytext)
3 library(wordcloud)
```

- **tidyverse**: this is an “opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures”;
- **tidytext**: an R package for text mining based on the tidy data principles;
- **wordcloud**: a package to generate word cloud plots.

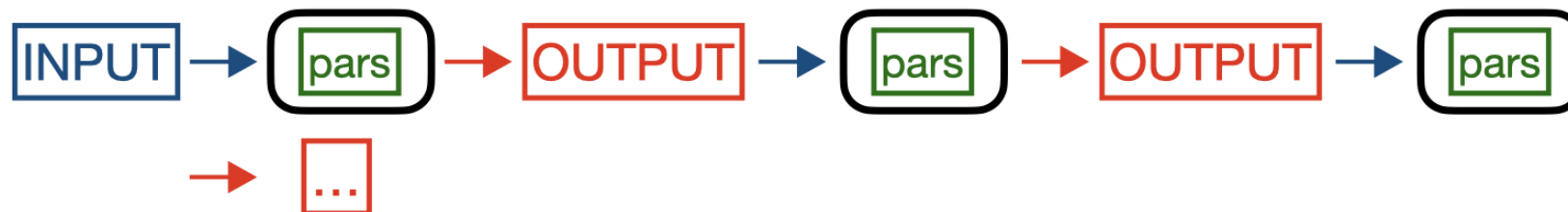
TIDYVERSE PIPELINE



A GENERAL FUNCTION



NESTED FUNCTIONS



PLAIN R SYNTAX

```
1 output1 <- func1(data, pars1)
2 output2 <- func2(output1, pars2)
3 output3 <- func3(output2, pars3)
4 output4 <- func4(output3, pars4)
```

or

```
1 output4 <-
2   func4(func3(func2(func1(data, pars1), pars2), pars3), pars4)
```

TIDYVERSE SYNTAX

```
1 output4 <- data %>%
2   func1(pars1) %>%
3   func2(pars2) %>%
4   func3(pars3) %>%
5   func4(pars4)
```

READING DATA

```
1 data_file_name <- '../../data/ianalyzer_query.csv'
2
3 data_df <- read_delim(data_file_name,
4   delim = ";",
5   escape_double = FALSE,
6   col_types = cols(`date-pub` = col_date(format = "%B %d, %Y"),
7     issue = col_integer()), trim_ws = TRUE)
8
9 print(nrow(data_df))
```

```
[1] 1532
```

```
1 print(colnames(data_df))
```

```
[1] "author"    "category" "content"  "date-pub" "edition"  "issue"    "query"
[8] "title"     "volume"
```



if you are in trouble with
`read_delim()`, get help from R
studio

TOKENIZATION

- Tokenization is process of dividing a string of text into meaningful units called **tokens**;
- A token can be a word, a phrase, a paragraph, or a single character depending on the nature of our analysis;
- In R tokenization is performed using the tidytext function `unnest_tokens()`.

TOKENIZATION

```
1 tidy_content <- data_df %>% unnest_tokens(word, content, token="words")
2
3 tidy_content
```

```
# A tibble: 1,549,578 × 9
```

	author	category	`date-pub`	edition	issue	query	title	volume
word	<chr>	<chr>	<date>	<lgl>	<int>	<chr>	<chr>	<lgl>
<chr>								
1	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
from								
2	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
3	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
spec...								
4	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
corr...								
5	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
6	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA
this								
7	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time..	Euro..	NA

CLEANING UP DATA

Checking if the column *issue* has any Na

```
1 are_there_na <- any(is.na(tidy_content$issue))  
2 are_there_na
```

```
[1] TRUE
```

let's clean up

```
1 tidy_content <- tidy_content[!is.na(tidy_content$issue), ]
```

and let's check again

```
1 are_there_na <- any(is.na(tidy_content$issue))  
2 are_there_na
```

```
[1] FALSE
```

REMOVING STOP WORDS

Unstructured data can contain a lot of irrelevant information. The most common words in a text are words that have very little meaning, such as “the”, “and”, “a”, etc. These words are referred to as **stop words** and removing stop words from text (in a way or another) is a fundamental step of text mining.

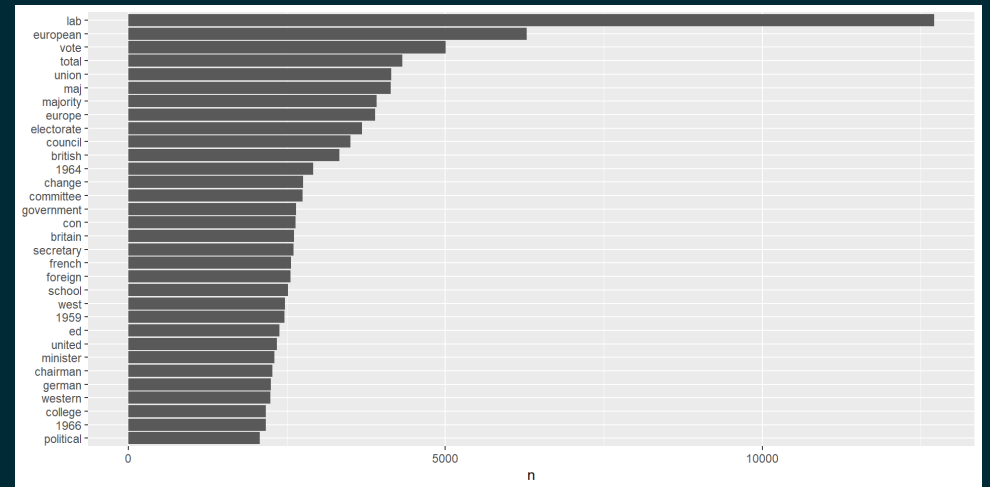
```
1 data(stop_words)
2
3 tidy_clean_content <- tidy_content %>% anti_join(stop_words)
4
5 tidy_clean_content
```

```
# A tibble: 801,754 × 9
```

	author	category	`date-pub`	edition	issue	query	title	volume
word	<chr>	<chr>	<date>	<lgl>	<int>	<chr>	<chr>	<lgl>
	<chr>							
1	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
spec...								
2	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
corr...								
3	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
junc...								
4	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
focus								
5	['FROM A SPECIAL ...	['News']	1962-11-05	NA	55540	time...	Euro...	NA

COUNTING WORDS

```
1 word_count <- tidy_clean_content %>%
2   count(word) %>%
3   filter(n > 2000) %>%
4   mutate(word = reorder(word, n))
5
6 word_count_plot <-
7   word_count %>%
8   ggplot(aes(n, word)) +
9   geom_col() +
10  labs(y = NULL)
11
12 word_count_plot
```



WORD CLOUD VISUALIZATION

```
1 word_cloud_plot <-  
2   word_count %>%  
3   with(wordcloud(word, n))
```

```
1 word_cloud_plot
```

NULL

lab government european
secretary council school
german chairman
con western committee 1959
political west electorate
vote total britained
1964 british 1966 minister
college union europe foreign
french maj united
majority change

SENTIMENT ANALYSIS

- **sentiment analysis** has the goal of systematically identify, extract, quantify, and study affective states and subjective information from text;
- Sentiment analysis is based on the assumption that we can view a text as a combination of individual words (the text sentiment will be the sum of the sentiment of its individual words);
- To perform sentiment analysis, we need a reference database of words called **lexicon** assigning a sentiment to each word.

LEXICON AND JOY WORDS

```
1 nrc_lexicon_df <- read.table("../../lexicons/NRC_lexicon.txt", header = FALSE)
2
3 joy_words <- nrc_lexicon_df %>%
4   filter(emotion == "joy", score == 1)
5
6 joy_words
```

	word	emotion	score
1	absolution	joy	1
2	abundance	joy	1
3	abundant	joy	1
4	accolade	joy	1
5	accompaniment	joy	1
6	accomplish	joy	1
7	accomplished	joy	1
8	achieve	joy	1
9	achievement	joy	1
10	acrobat	joy	1
11	admirable	joy	1
12	admiration	joy	1
13	adorable	joy	1
14	adoration	joy	1
15	adorn	joy	1

COMPUTING JOY WORDS FRACTION

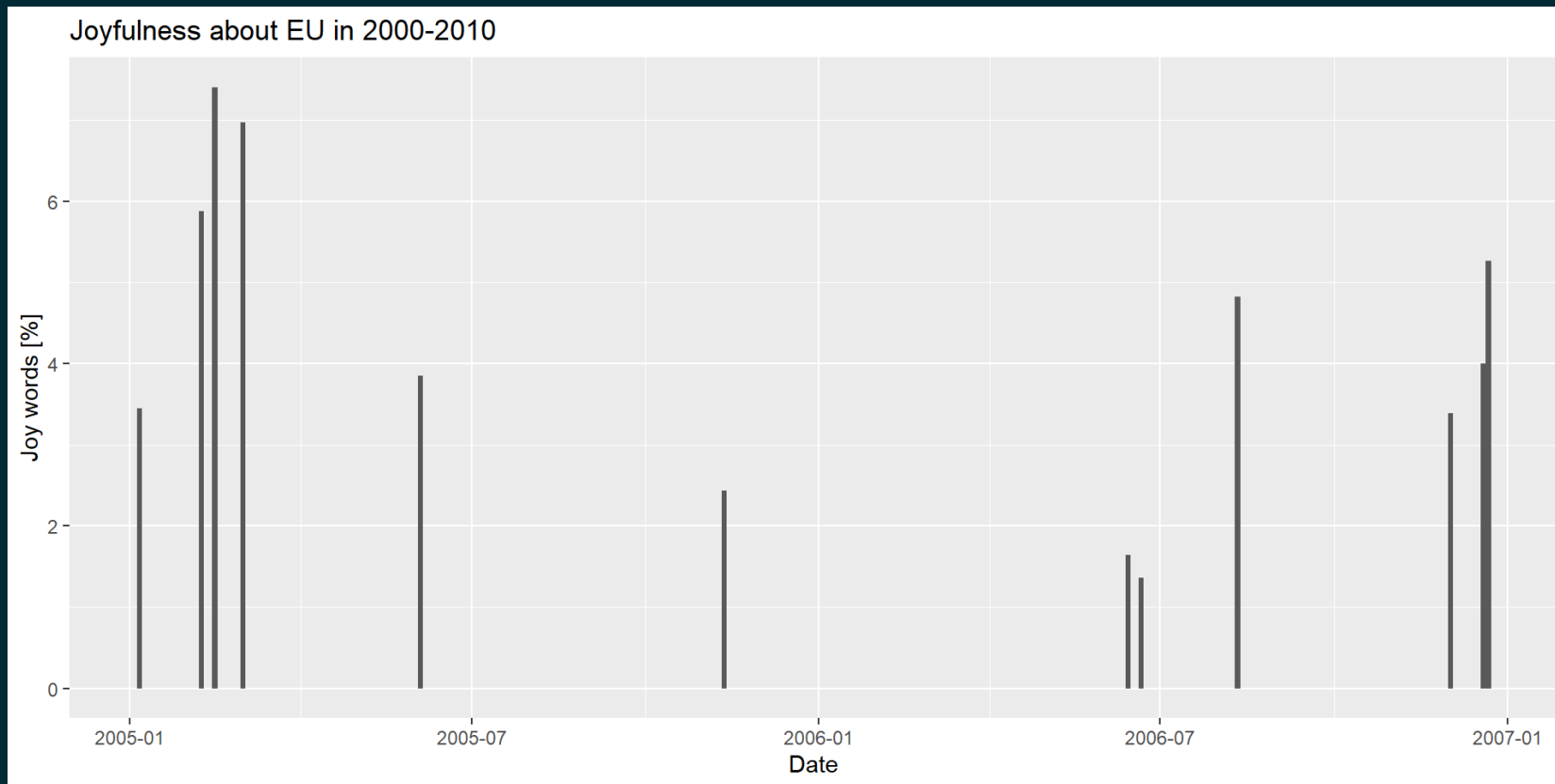
$$Frac_{joy}(issue) = \frac{\text{Number of joy words per issue}}{\text{Number of words per issue}} * 100$$

```
1 issue_df <- tidy_clean_content %>%
2   filter(`date-pub` >= '2000-01-01' & `date-pub` < '2010-01-01') %>%
3   group_by(issue) %>%
4   reframe(words_per_issue = n(), date = `date-pub`) %>%
5   unique()
6
7 issue_joy_df <- tidy_clean_content %>%
8   filter(`date-pub` >= '2000-01-01' & `date-pub` < '2010-01-01') %>%
9   inner_join(joy_words) %>%
10  group_by(issue) %>%
11  reframe(joy_words_per_issue = n())
12
13 issue_tot_df <- merge(issue_df, issue_joy_df, by='issue')
```

```

1 percent_of_joy_plot <-
2   issue_tot_df %>%
3   mutate(per_cent_joy=joy_words_per_issue/words_per_issue*100) %>%
4   ggplot(aes(x = date, y = per_cent_joy) )+
5   geom_col() +
6   labs(x = "Date", y = "Joy words [%]",
7         title = "Joyfulness about EU in 2000-2010")
8
9 percent_of_joy_plot

```



COMPUTING “TOTAL JOY” FRACTION

$$Frac_{joy} = \frac{\text{Number of joy words}}{\text{Number of words}} * 100[\%]$$

```
1 distinct_words <- tidy_clean_content %>%  
2   distinct(word)  
3  
4 total_dis_words <- distinct_words %>%  
5   nrow()  
6 total_dis_joy_words <- distinct_words %>%  
7   inner_join(joy_words, by='word') %>%  
8   nrow()  
9  
10 total_joy <- (total_dis_joy_words/total_dis_words)*100  
11 print(paste(total_joy, ' [%]'))
```

```
[1] "0.534795712569965 [%]"
```

ANALYZING WORD AND DOCUMENT FREQUENCY: TF-IDF

$$\text{tf-idf} = \text{term frequency} * \text{idf}$$

$$\text{idf}(\text{term}) = \log \left(\frac{n_{\text{documents}}}{n_{\text{documents containing term}}} \right)$$

COMPUTING TERM FREQUENCY

- Let's compute and store in two DataFrames the frequency of occurrence of each word and the total number of words per issue.

```
1 issue_words <- data_df %>%
2   unnest_tokens(word, content) %>%
3   count(issue, word)
4
5 issue_words <- na.omit(issue_words)
6
7 total_words <- issue_words %>%
8   group_by(issue) %>%
9   summarize(total = sum(n))
10
11 issue_total_words <- left_join(issue_words, total_words) %>%
12   arrange(desc(issue))
```

COMPUTING TERM FREQUENCY

```
1 unique_issues <- issue_total_words %>%
2   filter(total>10000) %>%
3   distinct(issue)
4
5 first_6_unique_issues <- unique_issues %>% slice(1:6)
6
7 issue_total_words6 <- issue_total_words %>%
8   semi_join(first_6_unique_issues, by="issue") %>%
9   mutate(issue=as.character(issue))
10
11 freq_per_issue_plot <-
12   issue_total_words6 %>%
13   ggplot(aes(n/total, fill = issue)) +
14   geom_histogram(show.legend = FALSE) +
15   xlim(NA, 0.0005) +
16   facet_wrap(~issue, ncol = 2, scales = "free_y")
17
18 freq_per_issue_plot
```

COMPUTING TERM FREQUENCY



COMPUTING AND DISPLAYING TF-IDF

```
1 issue_tf_idf <- issue_words %>%  
2   bind_tf_idf(word, issue, n)  
3  
4 issue_tf_idf %>%  
5   arrange(desc(tf))
```

```
# A tibble: 574,385 × 6
```

	issue	word	n	tf	idf	tf_idf
	<int>	<chr>	<int>	<dbl>	<dbl>	<dbl>
1	68302	the	8	0.170	0	0
2	52204	the	31	0.150	0	0
3	53256	the	28	0.146	0	0
4	53191	the	58	0.146	0	0
5	53078	the	27	0.141	0	0
6	57761	the	18	0.132	0	0
7	53284	the	58	0.130	0	0
8	53077	the	26	0.13	0	0
9	61094	the	49	0.130	0	0
10	53175	the	76	0.130	0	0

```
# i 574,375 more rows
```

COMPUTING AND DISPLAYING TF-IDF

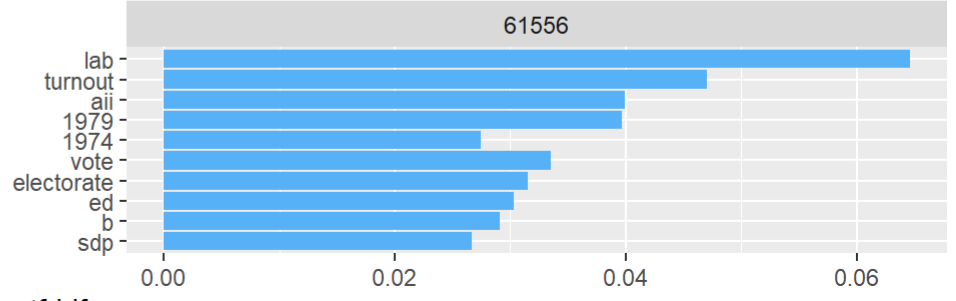
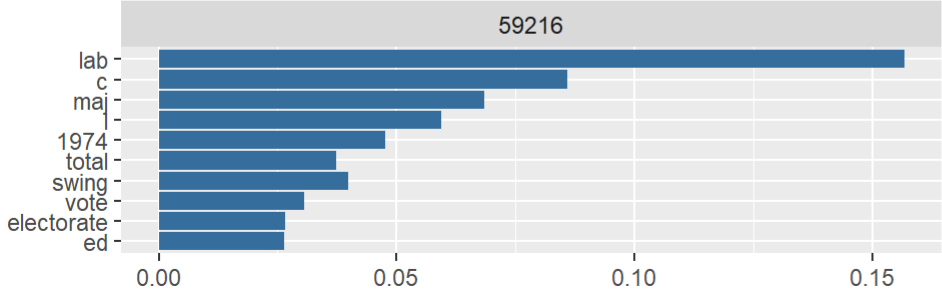
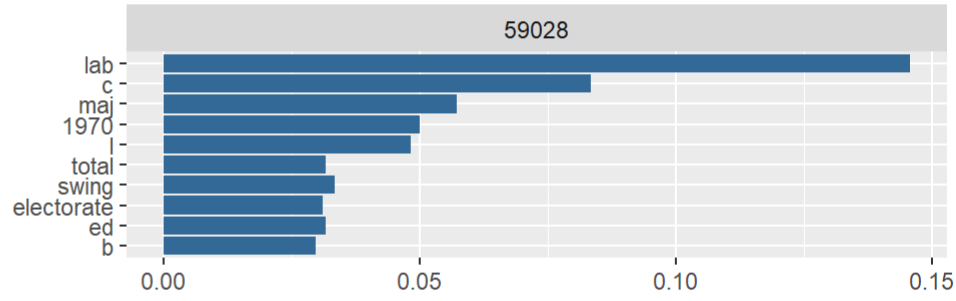
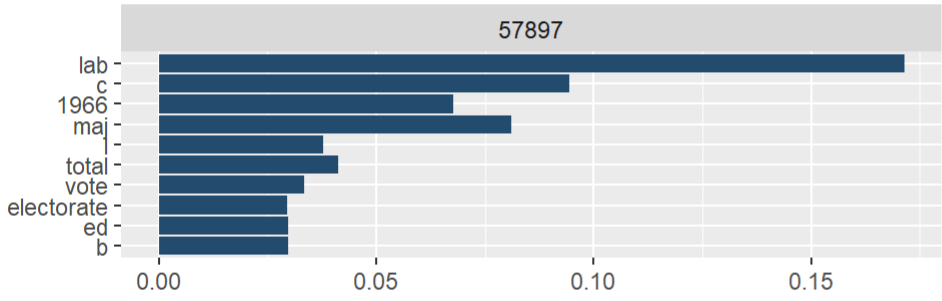
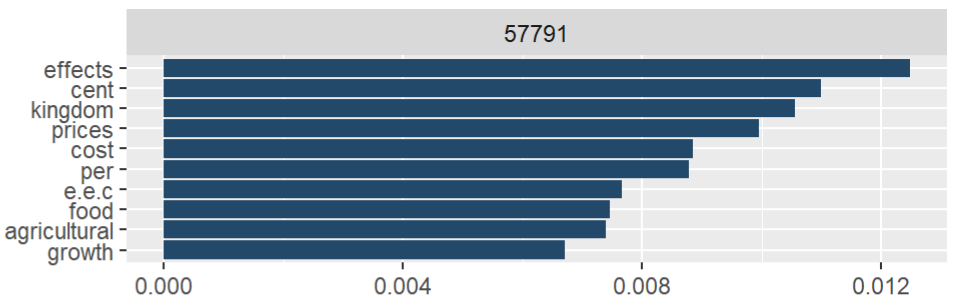
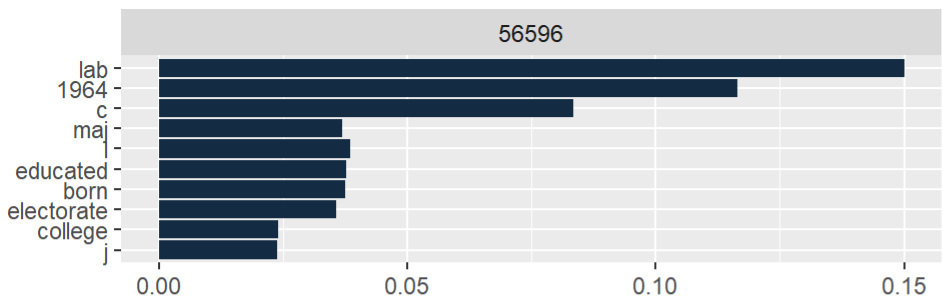
```
1 issue_tf_idf %>%  
2   arrange(desc(tf_idf))
```

```
# A tibble: 574,385 × 6  
  issue word      n    tf    idf tf_idf  
  <int> <chr>  <int> <dbl> <dbl> <dbl>  
1  68732 cod      21 0.0463  7.17  0.332  
2  68277 bosnian    2 0.0435  7.17  0.312  
3  68405 code      2 0.0426  5.38  0.229  
4  68873 croatia    4 0.0317  7.17  0.228  
5  68873 rehn      4 0.0317  7.17  0.228  
6  55541 merlot     3 0.0283  7.17  0.203  
7  68578 flag      2 0.0455  4.40  0.200  
8  68890 ceausescu  2 0.0278  7.17  0.199  
9  68277 agents     2 0.0435  4.53  0.197  
10 68302 wording    2 0.0426  4.61  0.196  
# i 574,375 more rows
```

COMPUTING AND DISPLAYING TF-IDF

```
1 issue_tf_idf %>%
2   semi_join(first_6_unique_issues, by="issue") %>%
3   group_by(issue) %>%
4   slice_max(tf_idf, n = 10) %>%
5   ungroup() %>%
6   ggplot(aes(tf_idf, fct_reorder(word, tf_idf), fill = issue)) +
7   geom_col(show.legend = FALSE) +
8   facet_wrap(~issue, scales="free", ncol = 2) +
9   labs(x = "tf-idf", y = NULL)
```


COMPUTING AND DISPLAYING TF-IDF



tf-idf

RELATIONSHIPS BETWEEN WORDS

- We can tokenize text so to obtain groups of n words or ngrams;
- An ngram is just a contiguous sequence of n items;
- In R ngrams are made using the tidytext function `unnest_tokens()`.

RELATIONSHIPS BETWEEN WORDS

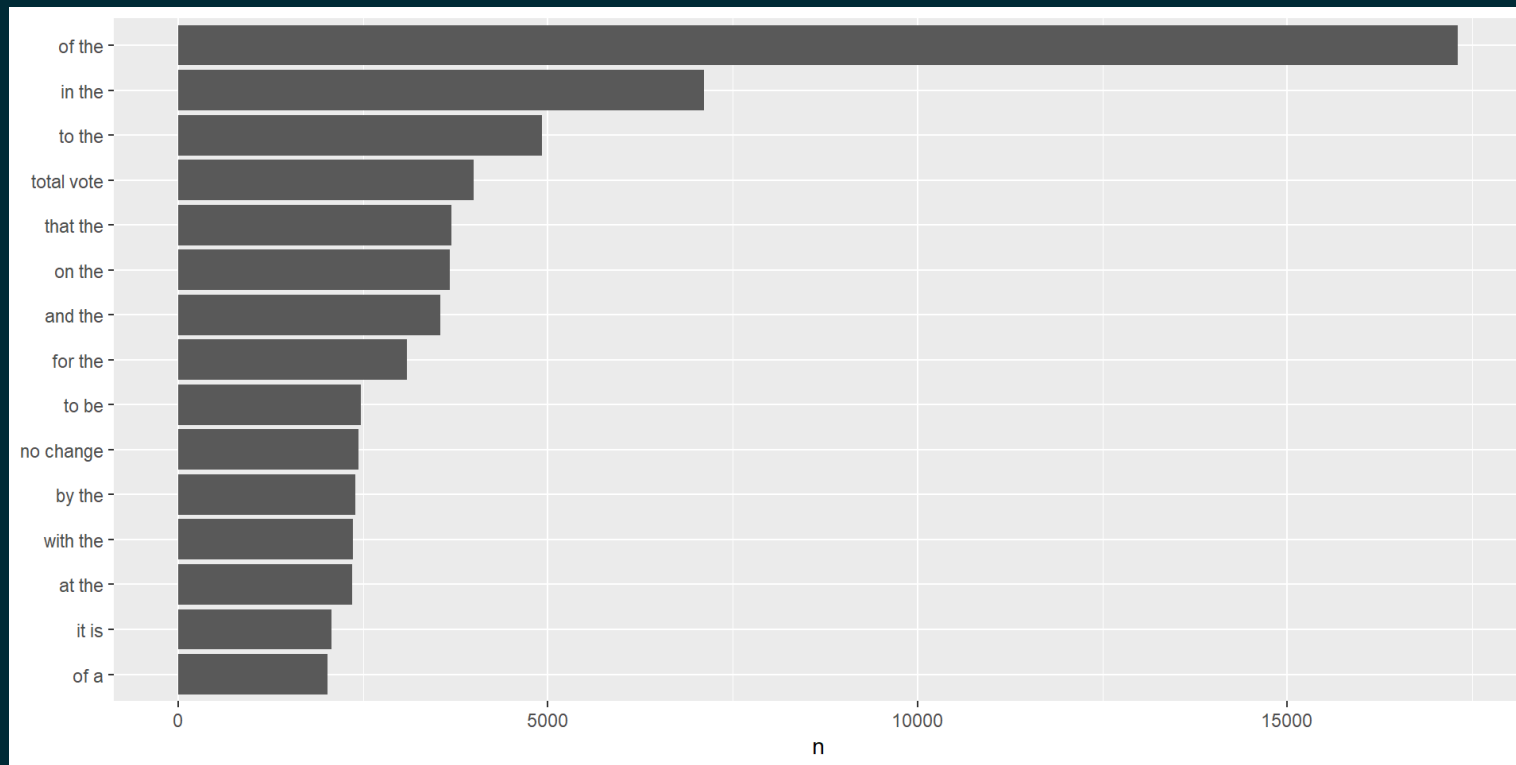
```
1 tidy_content_rel <- data_df %>%  
2   unnest_tokens(bigram, content, token="ngrams", n=2)  
3  
4 tidy_content_rel
```

```
# A tibble: 1,548,046 × 9
```

	author	category	`date-pub`	edition	issue	query	title	volume
bigram								
	<chr>	<chr>	<date>	<lgl>	<int>	<chr>	<chr>	<lgl>
<chr>								
1	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
a								from
2	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
a								a
spe...								
3	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
speci...								
4	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
corre...								
5	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA
th...								at
6	['FROM A SPECIAL...	['News']	1962-11-05	NA	55540	time...	Euro...	NA

RELATIONSHIPS BETWEEN WORDS

```
1 tidy_content_rel %>%  
2   count(bigram, sort = TRUE) %>%  
3   filter(n > 2000) %>%  
4   mutate(bigram = reorder(bigram, n)) %>%  
5   ggplot(aes(n, bigram)) +  
6   geom_col() +  
7   labs(y = NULL)
```



CLEANING UP BIGRAMS

```
1 bigrams_separated <- tidy_content_rel %>%
2   separate(bigram, c("word1", "word2"), sep = " ")
3
4 bigrams_filtered <- bigrams_separated %>%
5   filter(!word1 %in% stop_words$word) %>%
6   filter(!word2 %in% stop_words$word)
7
8 tidy_content_rel_clean <- bigrams_filtered %>%
9   unite(bigram, word1, word2, sep = " ")
```

PLOTTING BIAGRAMS

```
1 tidy_content_rel_clean %>%  
2   count(bigram, sort = TRUE) %>%  
3   filter(n > 500) %>%  
4   mutate(bigram = reorder(bigram, n)) %>%  
5   ggplot(aes(n, bigram)) +  
6   geom_col() +  
7   labs(y = NULL)
```

