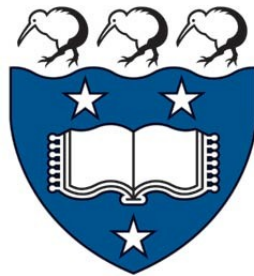


Text Analytics

Jason Peter Cairns

Supervised by Chris Wild



Bachelor of Science (Honours)
Department of Statistics
The University of Auckland
New Zealand

Todo list

should this be an abstract?	9
---------------------------------------	---

Acknowledgements

Contents

Listings	5
Tables	6
Figures	7
1 Introduction	9
1.1 Intention	9
1.2 Background: Text Analytics (incl. examples)	9
1.2.1 common functions: sentiment, summarisation, scoring	9
1.2.2 Existing Systems	10
1.3 Background: inZight	10
1.3.1 What inZight is - capabilities, popularity, etc.	10
1.3.2 how our program fits in - shiny, inzicht lite etc.	10
1.4 Literature Review (existing packages in R)	11
1.4.1 Copy over from notes, flesh out a bit	11
1.4.2 Praise tidytext book, complain about the package . . .	11
1.5 Scope of work	11
2 Text Analytics Prolusion	12
2.1 overview	12
2.1.1 Explain broadness of term	12
2.1.2 compile glossary from terms here	12
2.1.3 Areas of text analytics in a data science framework . .	12
2.1.4 what we have done	12
2.1.5 what we haven't done	12
2.2 terms	12
2.2.1 terms and their centrality	14
2.2.2 generalisation: n-grams, sentences etc.	14
2.3 Historical Background	14
2.3.1 computer science vs statistics - reflection in data science	14
2.4 Processing	14
2.4.1 why process	14

2.4.2	stopwords, lemmatisation etc.	14
2.4.3	modelling vs db joins - more info in notes	14
2.5	scores & statistics	14
2.5.1	why compute scores & statistics	14
2.5.2	scoring - tf-idf, word count	14
2.5.3	Suggestions for further research - more on the statistics of words	14
2.5.4	recount the book of John text analysis	14
2.6	Sentiment	14
2.6.1	why sentiment	14
2.6.2	Process of sentiment	14
2.6.3	sentiment modelling vs db joins	14
2.6.4	our implementation and why	14
2.6.5	reviews	14
2.6.6	issues	14
2.7	Summarisation	14
2.7.1	why compute summarisations	14
2.7.2	lexrank, textrank - include notes on lexrank	14
2.7.3	other methods	14
2.7.4	reddit bot example	14
2.8	what we didn't do (yet)	14
2.8.1	topic modelling	14
2.8.2	Term correlation	14
2.8.3	modelling based on linguistic features	14
2.9	Visualisation	14
2.9.1	talk about score vs structure	14
2.9.2	complain about tag clouds	14
2.9.3	talk about ggpage	14
2.9.4	discuss our experimentations with some alternative visualisations	14
3	Program Structure & Development	15
3.0.1	why R	15
3.0.2	Why Shiny	15
3.0.3	why tidyverse	15
3.0.4	Git	15
3.0.5	possible future: datatables, futures, etc.	15
3.0.6	why functional	15
3.0.7	Why lossless data	15
3.1	Program Architecture	15
3.1.1	Why structure it like it has been	15
3.1.2	make graph of architecture	15
3.1.3	Describe package and package creation	15

3.1.4	following three sections copy and paste from the notes - buffing up as necessary	15
3.1.5	include screenshots	15
3.2	Preparation	15
3.2.1	Importing	16
3.2.2	Object Preparation	18
3.2.3	Filtering	19
3.2.4	Lemmatisation	19
3.2.5	Stemming	19
3.2.6	Stopwords	20
3.2.7	Formatting	20
3.2.8	Sectioning	22
3.2.9	Grouping	23
3.3	Insight	24
3.3.1	Term Insight	24
3.3.2	Aggregate Insight	28
3.3.3	Wrapper	29
3.4	Visualisation	31
3.4.1	Rank	31
3.4.2	Score	31
3.4.3	Distribution	32
3.4.4	Structure	33
3.4.5	Wrapper	34
3.5	Application	35
4	Conclusion	41
4.1	Summary	41
4.1.1	summarise successes	41
4.1.2	summarise failures	41
4.1.3	general thoughts on the topic	41
4.2	Recommendations	41
4.2.1	educational potential of text analytics	41
4.2.2	what else remains	41
5	Appendix	42
	Glossary	69
	Index	69
	Bibliography	70

List of source codes

1	Import txt	16
2	Import csv	17
3	Import excel	17
4	Import files	18
5	Prepare text	18
6	Manage stopwords	20
7	Format data	22
8	Detect and add sections	23
9	Determine Term Frequencies	24
10	Get bigrams functionally	24
11	Concatenate walks	25
12	Concatenate walks	25
13	Bigram example	25
14	Get n-grams	26
15	Get n-grams frequencies	26
16	Determine Key Words with Textrank	27
17	Determine Term Sentiments	27
18	Determine the Moving Average Term Sentiment	28
19	Determine the term count over some aggregate	28
20	Determine the Key Sections	29
21	Determine the Aggregate Sentiments	29
22	Insight functions wrapper	30
23	Bind aggregate terms	31
24	Create Bar Plot	32
25	Create Word Cloud	32
26	Create Histogram	33
27	Create Kernel Density Estimate Plot	33
28	Create Time Series Plot	34
29	Create Page View Plot	34
30	Create Visualisation	35
31	UI and Server of Shiny Application	40

List of Tables

3.1	Primary data structure format	18
3.2	Formatting Logic for Stopwords and Lemmatisation	21

List of Figures

Chapter 1

Introduction

1.1 Intention

Text Analytics serves to glean insights from a body of text. Within the broad category of text analytics, we seek to answer questions about what the text is communicating, what is felt about it, and how this information is structured. In this dissertation, we demonstrate the creation of a user-friendly program to perform text analytics functions using modern R with the Shiny web application framework. In a literate style, we illustrate top-down the structure of such a program, as well as the data structures and computational processes that have established their value for such a program.

should
this be
an ab-
stract?

1.2 Background: Text Analytics (incl. examples)

1.2.1 common functions: sentiment, summarisation, scoring

Text Analytics is comprised of a variety of processes and techniques to extract information from text. The text almost always requires some initial processing. Some of the following functions have proven utility, and are expanded upon in chapter 2;

- **Sentiment:** In order to answer what emotions are conveyed in a text, sentiment analysis is commonly performed. The technique yields some measure of what is represented in an emotional sense by the text, with a range different methods and their associated outputs allowing for different forms of the analysis. Sentiment analysis won't pick up the subtle nuances that a human reader would, but generally gives reasonable output over the extent of a text.
- **Associated Words:** The meaning of a text is dependent on the struc-

ture between and within words. Looking at how words are associated, through correlation, common sequences, visualisation of sections, etc., allow for a clear high-level assessment of the associations between words. The higher level not only saves individual efforts, but can demonstrate any emergent properties inherent to a text, in a way that a direct reading won't necessarily reveal.

- **Summarisation:** Automation of an executive summary, or a list of key words, typically falls under the purview of summarisation. The primary aim is to rank and select the most “representative” words or sentences from a text. A few major techniques dominate, being somewhat complex in nature. The results seem to be surprisingly representative of a text.
- **Feature Counts:** The simplest quantitative measure is very often the most informative; from simple word counts, to selective counts of sentences within groups, counting features can reveal how much written weighting is given to various elements, aiding insight into content, structure, and sentiment simultaneously.

1.2.2 Existing Systems

There are several existing systems in the field of Text Analytics. The field was initially nurtured as a sub-field of Computer Science, being computationally-dependent in nature. More recently, there has been increasing statistical interest. The existing systems reflect this; most older text analytics programs were Artificial Intelligence focussed, being experimental in nature, typically composed in lisp. More recently, major statistical programs have been incorporating text analytic features, with a few smaller text-analytic specific programs appearing. SAS, SPSS, and R are all examples of major statistical processing systems, with recent additions of text analytics capabilities. An overview of R packages aiding in text analytics will be given in section 1.4.

1.3 Background: inZight

1.3.1 What iNZight is - capabilities, popularity, etc.

1.3.2 how our program fits in - shiny, inzight lite etc.

Our program will form part of the suite of modules extending iNZight. It provides a simple GUI interface to rapidly perform common text analyses. The primary audience are those learning the fundamentals and potential of text analysis and statistics, which could include students of the traditional text analytics fields of Statistics and Computer Science, but can and should include students of Linguistics, Communications, Law, History, and any

other text-based field. Beyond the educational aspects of the program, it is fully functional for practical use for general text analysis.

1.4 Literature Review (existing packages in R)

1.4.1 Copy over from notes, flesh out a bit

1.4.2 Praise tidytext book, complain about the package

1.5 Scope of work

While the total scope possible for text analytics is enormous, our time in creating this program is not. Thus, it is essential that we limit the scope. There are two primary areas with which we created the limitations: Text type, and analysis type.

By limiting the forms of text we work with, we can spend less effort on consideration of every single possible import and transformation case, and more time on the actual design of analysis. The simplest means with which to create the limitation exists in allowing only import of particular text files — in this case, we allow for flat `.txt` files, as well as tabular `.csv` and `.xlsx` files. What we do not provide (though by design leaving open the future possibility of inclusion) is access in-program to common text sources through their API, such as Twitter or Project Gutenberg.

Through focussing on dictionary-based, rather than model-based analyses, we have avoided much of the associated complexities. An example of this follows. It is common to categorise words based on their grammatical category, then use models that take this into account. By avoiding that (again, keeping the design flexible enough to incorporate this in the future), we have been able to get far more functionality implemented in a shorter amount of time, with the analyses still performing soundly. Additionally, we focus on the general audience, as it is typically more advanced, linguistically-trained users who would make intelligent use of such analyses.

Chapter 2

Text Analytics Prolusion

2.1 overview

Most importantly, words must be extracted, serving as the basic unit of analysis, from which more complex items may be derived.

2.1.1 Explain broadness of term

2.1.2 compile glossary from terms here

2.1.3 Areas of text analytics in a data science framework

2.1.4 what we have done

2.1.5 what we haven't done

2.2 terms

term

2.2.1 terms and their centrality

2.2.2 generalisation: n-grams, sentences etc.

2.3 Historical Background

2.3.1 computer science vs statistics - reflection in data science

2.4 Processing

2.4.1 why process

2.4.2 stopwords, lemmatisation etc.

2.4.3 modelling vs db joins - more info in notes

2.5 scores & statistics

2.5.1 why compute scores & statistics

2.5.2 scoring - tf-idf, word count

2.5.3 Suggestions for further research - more on the statistics of words

2.5.4 recount the book of John text analysis

2.6 Sentiment

2.6.1 why sentiment

2.6.2 Process of sentiment

2.6.3 sentiment modelling vs db joins

2.6.4 our implementation and why

2.6.5 reviews

2.6.6 issues

2.7 Summarisation

2.7.1 why compute summarisations

2.7.2 lexrank, textrank - include notes on lexrank

2.7.3 other methods

2.7.4 reddit bot example

2.8 what we didn't do (yet)

2.8.1 topic modelling

2.8.2 Term correlation

2.8.3 modelling based on linguistic features

2.9 Visualisation

Chapter 3

Program Structure & Development

3.0.1 why R

3.0.2 Why Shiny

3.0.3 why tidyverse

3.0.4 Git

3.0.5 possible future: datatables, futures, etc.

3.0.6 why functional

3.0.7 Why lossless data

3.1 Program Architecture

3.1.1 Why structure it like it has been

3.1.2 make graph of architecture

3.1.3 Describe package and package creation

3.1.4 following three sections copy and paste from the notes
- buffing up as necessary

3.1.5 include screenshots

3.2 Preparation

The first step in all text analysis is to import the text data and wrangle it into a data structure suitable for statistical analysis.

3.2.1 Importing

Text must first be brought in from an outside source to be useful for the program. The import functions are such that all text from different files exist in dataframes of equivalent structure. The primary differences are that each row of an imported `.txt` file corresponds to a single line, whereas each row of an imported tabular file corresponds to the row of the tabular file. Importantly for tabular files, the column of the text intended for analysis must be given the header of “text” prior to import. This condition will be relaxed later.-

Import `.txt`

The following is the simple function used in the import of `.txt` files:

```
1 #' Import text file
2 #'
3 #' @param filepath a string indicating the relative or absolute
4 #'   filepath of the file to import
5 #'
6 #' @return a [tibble][tibble::tibble-package] of each row
7 #'   corresponding to a line of the text file, with the column named
8 #'   "text"
9 import_txt <- function(filepath){
10   paste(readLines(filepath), collapse = "\n") %>%
11     tibble::tibble(text=.)
12 }
```

Listing 1: Import txt

Import `.csv`

CSV is a plaintext tabular format, with columns typically delimited by commas, and rows by new lines. A particular point of difference in the importation of tabular data and regular plaintext is that the text of interest for the analysis should be (as per tidy principles) in one column, with the rest being additional information that can be used for grouping or filtering. Thus, additional user input is required, in the specification of which column is the text column of interest. The following function is effectively just a wrapper around

```
1 #' Import csv file
2 #'
3 #' @param filepath a string indicating the relative or absolute
4 #'   filepath of the file to import
5 #'
6 #' @return a [tibble][tibble::tibble-package] of each row
7 #'   corresponding to a line of the text file, with the column named
8 #'   "text"
9 import_csv <- function(filepath){
10   readr::read_csv(filepath) ## %>%
11     ## table_textcol()
12 }
```

Listing 2: Import csv

Import Excel

Unfortunately, much data exists in the Microsoft Excel format, but this must be catered for. As tabular data, it is treated equivalently to csv, with a wrapper around `readr::read_excel()`

```
1 #' Import excel file
2 #'
3 #' @param filepath a string indicating the relative or absolute
4 #'   filepath of the file to import
5 #'
6 #' @return a [tibble][tibble::tibble-package] of each row
7 #'   corresponding to a line of the text file, with the column
8 #'   named "text"
9 import_excel <- function(filepath){
10   readxl::read_excel(filepath) ## %>%
11   ## table_textcol()
12 }
```

Listing 3: Import excel

Import Wrapper for Arbitrary Number of Files

To have just one function required to import files, we define two functions; one that imports any file, and one making use of it to import multiple files. The import of multiple files is no trivial task; the program must shape them in such a way that they retain identification, and fit into the same datastructure together.

The base wrapper function takes in the filename, and other relevant information, handling the importation process. It also stamps in the name of the document as a column.

The base file import is generalised to multiple files with a multiple import function: this will be our sole import function

```
1 #' Base case for file import
2 #'
3 #' @param filepath string filepath of file for import
4 #'
5 #' @return imported file with document id
6 import_base_file <- function(filepath){
7   filetype <- get_filetype(filepath)
8   filename <- basename(filepath)
9   if (filetype == "csv"){
10     imported <- import_csv(filepath)
11   } else if (filetype == "xlsx" | filetype == "xls") {
12     imported <- import_excel(filepath)
13   } else {
14     imported <- import_txt(filepath)
15   }
16   imported %>%
17     dplyr::mutate(doc_id = filename)
```

```

18 }
19
20 #' Import any number of files
21 #'
22 #' @param filepaths char vector of filepaths
23 #'
24 #' @return a [tibble][tibble::tibble-package] imported files with
25 #' document id
26 #'
27 #' @export
28 import_files <- function(filepaths){
29   filepaths %>%
30   purrr::map(import_base_file) %>%
31   dplyr::bind_rows()
32 }

```

Listing 4: Import files

3.2.2 Object Preparation

From the imported files, we work at transforming their representations into a lossless and efficient data structure that any analysis can make use of. Our solution to the essential constraint of losslessness is to separate and ID by each word in a dataframe. To do this, we take the line ID, the sentence ID, then the word ID, producing a dataframe that takes the following form:

line_id	sentence_id	word_id	word
1	1	1	the
1	1	2	quick
2	1	3	brown

Table 3.1: Primary data structure format

The reason for the ID columns is the preservation of the structure of the text; If required, the original text can be reconstructed in entirety, sans minor punctuation differences. The following function automatically formats any data of the format returned by the initial import functions.

```

1 text_prep <- function(data){
2   data %>%
3     tidytext::unnest_tokens(output = sentence, input = text,
4                             token = "sentences", to_lower = FALSE) %>%
5     dplyr::mutate(sentence_id = dplyr::row_number()) %>%
6     dplyr::group_by(sentence_id, add=TRUE) %>%
7     dplyr::group_modify(~ {
8       .x %>%
9         tidytext::unnest_tokens(output = word, input = sentence,
10                                token = "words", to_lower=FALSE) %>%
11         dplyr::mutate(word_id = dplyr::row_number())
12     }) %>%
13     ungroup_by("sentence_id")
14 }

```

Listing 5: Prepare text

3.2.3 Filtering

Filtering of text is implemented directly with the `dplyr::filter()` function, directly in the server of the shiny app. Filtering can take place multiple times throughout an analysis. The program is flexible enough such that after some initial analytics have been done in the insight layer, preparation can be returned to and the text can be filtered on based on features seen in the analytics.

3.2.4 Lemmatisation

Lemmatisation is effectively the process of getting words into dictionary form. It is a very complex, stochastic procedure, as natural languages don't follow consistent and clear rules all the time. Hence, models have to be used. Despite the burden, it is generally worthwhile to lemmatise words for analytics, as there are many cases of words not being considered significant, purely due to taking so many different forms relative to others. Additionally, stopwords work better when considering just the lemmatised form, rather than attempting to exhaustively cover every possible form of a word. `textstem` is an R package allowing for easy lemmatisation, with its function `lemmatize_words()` transforming a vector of words into their lemmatised forms (thus being compatible with `mutate()` straight out of the box). We have the lemmatisation in this program managed completely by this single function in the server end of the shiny app. The package `Udpipe` was another option, but it requires downloading model files, and performs far more in depth linguistic determinations such as parts-of-speech tagging, that we don't need at this point. It is worth noting that, like stopwords, there are different dictionaries available for the lemmatisation process, but we will use the default, as testing has shown it to be the simplest to set up and just as reliable as the rest.

3.2.5 Stemming

Stemming is far simpler than lemmatisation, being the removal of word endings. This doesn't require as complex a model, as it is deterministic. It is not quite as effective, as the base word ending is not concatenated back on at the tail, so we are left with word stumps and morphemes. However, it may sometimes be useful when the lemmatisation model isn't working effectively, and `textstem` provides the capability with `stem_words()`. We have not implemented this yet, as it is not as essential to an analysis when lemmatisation is already available.

3.2.6 Stopwords

We make use of dictionary-form stopwords, allowing for the input of both developed lexicons as well as user input. Two functions compose stopwords in the program: `get_sw()`, which gathers user input, queries the selected lexicon and combines the two, and `determine_stopwords()` which adds a boolean `TRUE | FALSE` column to the input dataframe. The following code listing defines `get_sw()` and `determine_stopwords()`:

```
1  #' Gets stopwords from a default list and user-provided list
2  #'
3  #' @param lexicon a string name of a stopwords list, one of "smart",
4  #'       "snowball", or "onix"
5  #'
6  #' @param addl user defined character vector of additional stopwords,
7  #'       each element being a stopwords
8  #'
9  #' @return a [tibble][tibble::tibble-package] with one column named "word"
10 get_sw <- function(lexicon = "snowball", addl = NA){
11   addl_char <- as.character(addl)
12   tidytext::get_stopwords(source = lexicon) %>%
13     dplyr::select(word) %>%
14     dplyr::bind_rows(., tibble::tibble(word = addl_char)) %>%
15     stats::na.omit() %>%
16     purrr::as_vector() %>%
17     tolower() %>%
18     as.character()
19 }
20
21 #' determine stopwords status
22 #'
23 #' @param .data vector of words
24 #'
25 #' @param ... arguments of get_sw
26 #'
27 #' @return a [tibble][tibble::tibble-package] equivalent to the input
28 #'   dataframe, with an additional stopwords column
29 #'
30 #' @export
31 determine_stopwords <- function(.data, ...){
32   sw_list <- get_sw(...)
33   .data %in% sw_list
34 }
```

Listing 6: Manage stopwords

3.2.7 Formatting

The final component in preparation is to format the prepared object with the correct attributes to have formatting automated. We define a wrapper that takes all combinations of stopwords and lemmatisation options and intelligently connects them for the “insight column” in a dataframe, which the insight is performed upon. For the purpose of standard interoperability, e.g., with `ggpage`, we name this column “text”.

At the heart of this function is an `ifexp()` that encodes the following lo-

gic involving the interaction of stopwords and lemmatisation, to enable the correct output text based on stopword and lemmatisation options;

	Stopwords True	Stopwords False
Lemmatise True	Lemmatise, determine stopwords on lemmatisation, perform insight on lemmas sans stopwords	Lemmatise, perform insight on lemmas
Lemmatise False	Determine stopwords on original words (no lemmatisation), perform insight on words sans stopwords	Perform insight on original words

Table 3.2: Formatting Logic for Stopwords and Lemmatisation

Based on the combination, stopword filtering and lemmatisation take place inside the function, defined as the following:

```

1  #' takes imported one-line-per-row data and prepares it for later analysis
2  #'
3  #' @param .data tibble with one line of text per row
4  #'
5  #' @param lemmatize boolean, whether to lemmatize or not
6  #'
7  #' @param stopwords boolean, whether to remove stopwords or not
8  #'
9  #' @param sw_lexicon string, lexicon with which to remove stopwords
10 #'
11 #' @param addl_stopwords char vector of user-supplied stopwords
12 #'
13 #' @return a [tibble][tibble::tibble-package] with one token per line,
14 #'   stopwords removed leaving NA values, column for analysis named
15 #'   "text"
16 #'
17 #' @export
18 format_data <- function(.data, lemmatize=TRUE, stopwords=TRUE,
19   sw_lexicon="snowball", addl_stopwords=NA){
20   formatted <- .data %>%
21     text_prep()
22
23   text <- ifexp(lemmatize,
24     ifexp(stopwords,
25       dplyr::mutate(formatted,
26         lemma = tolower(textstem::lemmatize_words(word)),
27         stopword = determine_stopwords(lemma,
28           sw_lexicon,
29           addl_stopwords),
30         text = dplyr::if_else(stopword,
31           as.character(NA),
32           lemma)),
33       dplyr::mutate(formatted,
34         lemma = tolower(textstem::lemmatize_words(word)),
35         text = lemma)),
36     ifexp(stopwords,
37       dplyr::mutate(formatted,
38         stopword = determine_stopwords(word,

```

```

39                                     sw_lexicon,
40                                     addl_stopwords),
41     text = dplyr::if_else(stopword,
42                           as.character(NA),
43                           word)),
44     dplyr::mutate(formatted, text = word)))
45   return(text)
46 }

```

Listing 7: Format data

3.2.8 Sectioning

Plaintext, as might exist as a Gutenberg Download, differs from more complex representations in many ways, including a lack of sectioning — for example, chapters require a specific search in order to jump to them. Here, I compose a closure that searches and sections text based on a Regular Expression intended to capture a particular section. Several functions are created from that. At a later date, advanced users could be given the option to compose their own regular expressions for sectioning.

```

1  #' creates a search closure to section text
2  #'
3  #' @param search a string regexp for the term to seperate on, e.g. "Chapter"
4  #'
5  #' @return closure over search expression
6  get_search <- function(search){
7    function(.data){
8      .data %>%
9        stringr::str_detect(search) %>%
10       purrr::accumulate(sum, na.rm=TRUE)
11    }
12  }
13
14  #' sections text based on chapters
15  #'
16  #' @param .data vector to section
17  #'
18  #' @return vector of same length as .data with chapter numbers
19  #'
20  #' @export
21  get_chapters <- get_search("^\\s*[Cc][Hh][Aa]?[Pp][Tt]([Ee][Rr])?")
22
23  #' sections text based on parts
24  #'
25  #' @param .data vector to section
26  #'
27  #' @return vector of same length as .data with part numbers
28  #'
29  #' @export
30  get_parts <- get_search("^\\s*[Pp]([Aa][Rr])?[Tt]")
31
32  #' sections text based on sections
33  #'
34  #' @param .data vector to section
35  #'
36  #' @return vector of same length as .data with section numbers
37  #'

```

```

38 #' @export
39 get_sections <- get_search("^\\s*([Ss][Ss])|([Ss][Ee][Cc][Tt][Ii][Oo][Nn])")
40
41 #' sections text based on cantos
42 #'
43 #' @param .data vector to section
44 #'
45 #' @return vector of same length as .data with canto numbers
46 #'
47 #' @export
48 get_cantos <- get_search("(?i)canto (XC|XL|L?X{0,3})(IX|IV|V?I{0,3})$")
49
50 #' sections text based on book
51 #'
52 #' @param .data vector to section
53 #'
54 #' @return vector of same length as .data with book numbers
55 #'
56 #' @export
57 get_books <- get_search("(?i)book$")
58
59 #' Adds section column to dataframe
60 #'
61 #' @param .data dataframe formatted as per output of prep process
62 #'
63 #' @param section_by character name of what to section over
64 #'
65 #' @return input dataframe with additional section column
66 #'
67 #' @export
68 section <- function(.data, section_by){
69   sec_table <- list("chapter" = get_chapters,
70                     "part" = get_parts,
71                     "section" = get_sections,
72                     "canto" = get_cantos,
73                     "book" = get_books)
74   .data %>%
75     dplyr::mutate(! section_by := sec_table[[section_by]](word))
76 }

```

Listing 8: Detect and add sections

3.2.9 Grouping

Grouping is an essential, killer feature of our app. The implementation is to run a `dplyr::group_by()` command in the shiny server on the prepared object, over user-specified groups, and all further insights and visualisations are performed groupwise. This allows for immediate and clear comparisons between groups.

Like filtering, after some initial analytics have been done in the insight layer, preparation can be returned to and the text can be grouped on based on the analytics.

3.3 Insight

3.3.1 Term Insight

Term Frequency

```
1 #' Determine term frequency
2 #'
3 #' @param .data character vector of terms
4 #'
5 #' @return numeric vector of term frequencies
6 #'
7 #' @export
8 term_freq <- function(.data){
9   .data %>%
10     tibble::enframe() %>%
11     dplyr::add_count(value) %>%
12     dplyr::mutate(n = dplyr::if_else(is.na(value),
13                                   as.integer(NA),
14                                   n)) %>%
15     dplyr::pull(n)
16 }
```

Listing 9: Determine Term Frequencies

n-Grams

```
1 #' Determine bigrams
2 #'
3 #' @param .data character vector of words
4 #'
5 #' @return character vector of bigrams
6 #'
7 get_bigram <- function(.data){
8   concat_walk(.data, c(.data[-1], NA))
9 }
```

Listing 10: Get bigrams functionally

```
1 #' concat list 1 and 2 at index, skipping NA values
2 #'
3 #' @param i numeric index to assess index at
4 #'
5 #' @param list1 list or vector for first token
6 #'
7 #' @param list2 list or vector for second token
8 #'
9 #' @return paste of list1 and list2 at index i, skipping NA's
10 concat_walk_i <- function(i, list1, list2){
11   ifelse(length(list2) < i | is.na(list1[i]),
12         as.character(NA),
13   ifelse(!(is.na(list1[i]) | is.na(list2[i])),
14         paste(list1[i], list2[i]),
15         concat_walk_i(i, list1, list2[-1])))
16 }
17
18 #' concat list 1 and 2, moving past NA values
19 #'
```

```

20 #' @param list1 list or vector for first bigram token
21 #'
22 #' @param list2 list or vector for second bigram token
23 #'
24 #' @return paste of list1 and list2, skipping NA's
25 concat_walk <- function(list1, list2){
26   stopifnot(length(list1) == length(list2))
27   sapply(seq_along(list1), concat_walk_i, list1, list2)
28 }

```

Listing 11: Concatenate walks

```

1 #' concat list 1 and 2 at index, skipping NA values
2 #'
3 #' @param i numeric index to assess index at
4 #'
5 #' @param list1 list or vector for first token
6 #'
7 #' @param list2 list or vector for second token
8 #'
9 #' @return paste of list1 and list2 at index i, skipping NA's
10 concat_walk_i <- function(i, list1, list2){
11   ifelse(length(list2) < i | is.na(list1[i]),
12         as.character(NA),
13   ifelse(!(is.na(list1[i]) | is.na(list2[i])),
14         paste(list1[i], list2[i]),
15         concat_walk_i(i, list1, list2[-1])))
16 }
17
18 #' concat list 1 and 2, moving past NA values
19 #'
20 #' @param list1 list or vector for first bigram token
21 #'
22 #' @param list2 list or vector for second bigram token
23 #'
24 #' @return paste of list1 and list2, skipping NA's
25 concat_walk <- function(list1, list2){
26   stopifnot(length(list1) == length(list2))
27   sapply(seq_along(list1), concat_walk_i, list1, list2)
28 }

```

Listing 12: Concatenate walks

```

1 x <- c(1, 2, NA, 4, 5, NA, 7, NA)
2 get_bigram(x)

```

Listing 13: Bigram example

```

1 #' Returns the n-grams, skipping NA values
2 #'
3 #' @param .data vector to get n-grams from
4 #'
5 #' @param n number of n-grams to attain
6 #'
7 #' @return n-gram vector without NA values
8 #'
9 #' @export
10 get_ngram <- function(.data, n){
11   main_n <- n
12   ngrams <- rep(NA_character_, length(.data))
13   for (i in seq_along(.data)){

```

```

14         if (i - 1 > length(.data) - n) break
15         if (is.na(.data[i])){
16             next
17         } else {
18             ngram <- .data[i]
19         }
20         j <- i + 1
21         while(n > 1){
22             if (j > length(.data)){
23                 ngram <- NA_character_
24                 break
25             }
26             if (is.na(.data[j])){
27                 j <- j + 1
28             } else {
29                 ngram <- paste(ngram, .data[j])
30                 n <- n - 1
31                 j <- j + 1
32             }
33         }
34         ngrams[i] <- ngram
35         n <- main_n
36     }
37     return(ngrams)
38 }

```

Listing 14: Get n-grams

```

1  #' NOT FOR PRODUCTION - STILL IN TESING. Returns the count of n-grams, skipping NA values
2  #'
3  #' @param .data vector to get n-grams from
4  #'
5  #' @param n number of n-grams to attain
6  #'
7  #' @return count of each associated n-gram
8  #'
9  #' @export
10 ngram_freq <- function(.data, n){
11     term_freq(get_ngram(.data, n))
12 }

```

Listing 15: Get n-grams frequencies

Key Words

```

1  #' Determine textrank score for vector of words
2  #'
3  #' @param .data character vector of words
4  #'
5  #' @param summ_method method to use for summarisation: textrank or
6  #'     lexrakn. Doesn't do anything yet
7  #'
8  #' @return vector of scores for each word
9  #'
10 #' @export
11 keywords_tr <- function(.data, summ_method){
12     relevent <- !is.na(.data)
13     tr <- textrank::textrank_keywords(.data, relevent, p=+Inf)
14     score <- tr$pagerank$vector %>% tibble::enframe()
15     data <- .data %>% tibble::enframe("number", "name")

```

```

16   dplyr::full_join(data, score, by="name") %>%
17   dplyr::pull(value)
18 }

```

Listing 16: Determine Key Words with Textrank

Term Sentiment

```

1  #' Determine sentiment of terms
2  #'
3  #' @param .data vector of terms
4  #'
5  #' @param lexicon sentiment lexicon to use, based on the corpus
6  #'   provided by tidytext
7  #'
8  #' @return vector with sentiment score of each word in the vector
9  #'
10 #' @export
11 term_sentiment <- function(.data, lexicon="afinn"){
12   data <- tibble::enframe(.data, "number", "word")
13   tidytext::get_sentiments(lexicon) %>%
14     dplyr::select(word, value) %>%
15     dplyr::right_join(data, by="word") %>%
16     dplyr::pull(value)
17 }

```

Listing 17: Determine Term Sentiments

Moving Average Term Sentiment

```

1  #' Determine the lagged sentiment of terms
2  #'
3  #' @param .data vector of terms
4  #'
5  #' @param lexicon sentiment lexicon to use, based on the corpus
6  #'   provided by tidytext
7  #'
8  #' @param lag how many (inclusive) terms to compute statistic over
9  #'
10 #' @param statistic base statistic used to summarise the data, capable
11 #'   of taking an na.rm argument
12 #'
13 #' @return vector with lagged sentiment score of each term in the input vector
14 #'
15 #' @export
16 ma_term_sentiment <- function(.data, lexicon="afinn", lag = 10, statistic = mean){
17   sents <- term_sentiment(.data, lexicon)
18   ## lagged_sents <- rep(NA, length(sents))
19   ## for (i in seq(lag, length(sents))){
20   ##   lagged_sents[i] <- statistic(sents[(seq(i - lag + 1, i))], na.rm = TRUE)
21   ## }
22   ## lagged_sents
23   c(rep(NA, lag - 1),
24     sapply(seq(lag, length(sents)),
25            function(i){x <- statistic(sents[(seq(i - lag + 1, i))],
26                                     na.rm = TRUE)
27                          ifelse(is.nan(x),
28                                NA,

```

```

29                                     x)
30             )))
31 }

```

Listing 18: Determine the Moving Average Term Sentiment

3.3.2 Aggregate Insight

Term Count

```

1  #' Determine the number of terms at each aggregate level
2  #'
3  #' @param .data character vector of terms
4  #'
5  #' @param aggregate_on vector to split .data on for insight
6  #'
7  #' @return vector of number of terms for each aggregate level, same
8  #'         length as .data
9  #'
10 #' @export
11 term_count <- function(.data, aggregate_on){
12   split(.data, aggregate_on) %>%
13     purrr::map(function(x){rep(length(x), length(x))}) %>%
14     dplyr::combine()
15 }

```

Listing 19: Determine the term count over some aggregate

Key Sections

```

1  #' get score for key sentences as per Lexrank
2  #'
3  #' @param .data character vector of words
4  #'
5  #' @param summ_method method to use for summarisation: textrank or
6  #'         lexrank. Doesn't do anything yet
7  #'
8  #' @param aggregate_on vector to aggregate .data over; ideally, sentence_id
9  #'
10 #' @return lexrank scores of aggregates
11 #'
12 #' @export
13 key_aggregates <- function(.data, aggregate_on, summ_method){
14   ## prepare .data for lexrank
15   base <- tibble::tibble(word = !! .data, aggregate = aggregate_on)
16   aggregated <- base %>%
17     dplyr::group_by(aggregate) %>%
18     stats::na.omit() %>%
19     dplyr::summarise(sentence = paste(word, collapse = " ")) %>%
20     dplyr::mutate(sentence = paste0(sentence, "."))
21   ## lexrank
22   lr <- aggregated %>%
23     dplyr::pull(sentence) %>%
24     lexRankr::lexRank(., n=length(.), removePunc = FALSE, returnTies = FALSE,
25       removeNum = FALSE, toLower = FALSE, stemWords = FALSE,
26       rmStopWords = FALSE, Verbose = TRUE)
27   ## match lexrank output to .data
28   lr %>%

```

```

29     dplyr::distinct(sentence, .keep_all = TRUE) %>%
30     dplyr::full_join(aggregated, by="sentence") %>%
31     dplyr::full_join(base, by="aggregate") %>%
32     dplyr::arrange(aggregate) %>%
33     dplyr::pull(value)
34 }

```

Listing 20: Determine the Key Sections

Aggregate Sentiment

```

1  #' Get statistics for sentiment over some group, such as sentence.
2  #'
3  #' @param .data character vector of words
4  #'
5  #' @param aggregate_on vector to aggregate .data over; ideally,
6  #'   sentence_id, but could be chapter, document, etc.
7  #'
8  #' @param lexicon as per term sentiment
9  #'
10 #' @param statistic function that accepts na.rm argument; e.g. mean,
11 #'   median, sd.
12 #'
13 #' @return sentiment of same length as input vector aggregated over the aggregate_on vector
14 #'
15 #' @export
16 aggregate_sentiment <- function(.data, aggregate_on, lexicon = "afinn", statistic = mean){
17   tibble::enframe(.data, "nil1", "word") %>%
18     dplyr::bind_cols(tibble::enframe(aggregate_on, "nil2", "aggregate")) %>%
19     dplyr::select(word, aggregate) %>%
20     dplyr::mutate(sentiment = term_sentiment(word, lexicon)) %>%
21     dplyr::group_by(aggregate) %>%
22     dplyr::mutate(aggregate_sentiment =
23       (function(.x){
24         rep(statistic(.x, na.rm = TRUE), length(.x))
25       })(sentiment)) %>%
26     dplyr::pull(aggregate_sentiment)
27 }

```

Listing 21: Determine the Aggregate Sentiments

Word Correlation

Term Frequency — Inverse Document Frequency

Topic Modelling

3.3.3 Wrapper

```

1  #' perform group-aware term operations on the data
2  #'
3  #' @param .data dataframe of terms as per output of format_data
4  #'
5  #' @param operations character vector of term operations to perform
6  #'
7  #' @param ... additional arguments to the operation - only sensible for singular operations
8  #'

```

```

9  #' @return .data with operation columns added
10 #'
11 #' @export
12 get_term_insight <- function(.data, operations, ...){
13   opstable <- list("Term Frequency" = term_freq,
14                   "n-gram Frequency" = ngram_freq,
15                   "n-grams" = get_ngram,
16                   "Key Words" = keywords_tr,
17                   "Term Sentiment" = term_sentiment,
18                   "Moving Average Term Sentiment" = ma_term_sentiment)
19   ops <- opstable[operations]
20   lapply(seq(length(ops)),
21         function(x){
22           name <- dplyr::sym(names(ops[x]))
23           operation <- ops[x][[1]]
24           df <- dplyr::mutate(.data,
25                             !!name := operation(text, ...))
26           df[names(ops[x])]
27         }) %>%
28     dplyr::bind_cols(.data, .)
29 }
30
31 #' perform group-aware aggregate operations on the data
32 #'
33 #' @param .data dataframe of terms as per output of format_data
34 #'
35 #' @param operations character vector of operations to perform
36 #'
37 #' @param aggregate_on character name of the column to perform aggregate operations on
38 #'
39 #' @param ... additional arguments to the operation - only sensible for singular operations
40 #'
41 #' @return .data with operation columns added
42 #'
43 #' @export
44 get_aggregate_insight <- function(.data, operations, aggregate_on, ...){
45   opstable <- list("Aggregated Term Count" = term_count,
46                   "Key Sections" = key_aggregates,
47                   "Aggregated Sentiment" = aggregate_sentiment,
48                   "Bound Aggregates" = bind_aggregation)
49   ops <- opstable[operations]
50   lapply(seq(length(ops)),
51         function(x){
52           name <- dplyr::sym(names(ops[x]))
53           operation <- ops[x][[1]]
54           agg_on <- dplyr::sym(aggregate_on)
55           df <- if (names(ops[x]) == "Bound Aggregates"){
56             dplyr::mutate(.data,
57                           !!name := operation(word, !! agg_on))
58           } else {
59             dplyr::mutate(.data,
60                           !!name := operation(text, !! agg_on, ...))
61           }
62           df[names(ops[x])]
63         }) %>%
64     dplyr::bind_cols(.data, .)
65 }

```

Listing 22: Insight functions wrapper

```

1  #' bind aggregate terms together
2  #'

```

```

3  #' @param data vector of terms
4  #'
5  #' @param aggregate_on vector of aggregations
6  #'
7  #' @return data with every aggregation bound, as in a sentence
8  #'
9  #'
10 bind_aggregation <- function(data, aggregate_on){
11   tibble::tibble(data, agg = aggregate_on) %>%
12     dplyr::group_by(agg) %>%
13     dplyr::mutate(bound = paste(data, collapse = " ")) %>%
14     dplyr::pull(bound)
15 }

```

Listing 23: Bind aggregate terms

3.4 Visualisation

3.4.1 Rank

3.4.2 Score

Bar Plot

```

1  #' output a ggplot column graph of the top texts from some insight function
2  #'
3  #' @param .data a dataframe containing "text" and insight columns as
4  #'           per the output of the get_(term|aggregate)_insight wrapper
5  #'           function
6  #'
7  #' @param y symbol name of the column insight was
8  #'           outputted to
9  #'
10 #' @param x symbol name of column for insight labels
11 #'
12 #' @param n number of bars to display
13 #'
14 #' @param desc bool: show bars in descending order
15 #'
16 #' @export
17 score_barplot <- function(.data, y, n = 15,
18                           x = text, desc = FALSE){
19   wrap <- 50
20   text <- dplyr::enquo(x)
21   insight_col <- dplyr::enquo(y)
22   .data %>%
23     dplyr::distinct(!! text, .keep_all=TRUE) %>%
24     dplyr::arrange(dplyr::desc(!! insight_col)) %>%
25     dplyr::group_modify(~{.x %>% head(n)}) %>%
26     dplyr::ungroup() %>%
27     dplyr::mutate(text = forcats::fct_reorder(shorten(!! text, wrap),
28                                               !! insight_col,
29                                               .desc = desc)) %>%
30     ggplot2::ggplot(ggplot2::aes(x = text)) +
31     ggplot2::geom_col(ggplot2::aes(y = !! insight_col)) +
32     ggplot2::coord_flip()
33 }
34

```



```

35 #' Shorten some text up to n characters
36 #'
37 #' @param .data character vector
38 #'
39 #' @param n wrap length of text
40 #'
41 #' @return shortened form of .data
42 #'
43 shorten <- function(.data, n){
44   ifelse(nchar(.data) > n,
45         paste(substr(.data, 1, n), "...", sep = ""),
46         .data)
47 }

```

Listing 24: Create Bar Plot

Word Cloud

```

1 #' output a ggplot wordcloud graph of the top texts from some insight function
2 #'
3 #' @param .data a dataframe containing "text" and insight columns as
4 #'   per the output of the get_(term|aggregate)_insight wrapper
5 #'   function
6 #'
7 #' @param y symbol name of the column insight was
8 #'   outputted to
9 #'
10 #' @param x symbol name of column for insight labels
11 #'
12 #' @param n number of words to display
13 #'
14 #' @param shape character: shape of the wordcloud
15 #'
16 #' @export
17 score_wordcloud <- function(.data, y, n = 15,
18                             x = text, shape = "circle"){
19   text <- dplyr::enquo(x)
20   insight_col <- dplyr::enquo(y)
21   .data %>%
22     dplyr::distinct(! text, .keep_all=TRUE) %>%
23     dplyr::arrange(dplyr::desc(! insight_col)) %>%
24     dplyr::group_modify(~{.x %>% head(n)}) %>%
25     dplyr::ungroup() %>%
26     ggplot2::ggplot(ggplot2::aes(label = stringr::str_wrap(! text, 30), size = ! insight_col)) +
27     ggwordcloud::geom_text_wordcloud(shape = shape, rm_outside = TRUE) +
28     ggplot2::scale_size_area(max_size = 24)
29 }

```

Listing 25: Create Word Cloud

3.4.3 Distribution

Histogram

```

1 #' output a histogram of the distribution of some function of words
2 #'
3 #' @param .data the standard dataframe, modified so the last column
4 #'   is the output of some insight function (eg. output from

```

```

5 #'      term_freq)
6 #'
7 #' @param col_name symbol name of the column insight was
8 #'      performed on
9 dist_hist <- function(.data, col_name){
10   q_col_name <- dplyr::enquo(col_name)
11   .data %>%
12     ggplot2::ggplot(ggplot2::aes(x = !! q_col_name)) +
13     ggplot2::geom_histogram()
14 }

```

Listing 26: Create Histogram

Density

```

1 #' output a histogram of the distribution of some function of words
2 #'
3 #' @param .data the standard dataframe, modified so the last column
4 #'      is the output of some insight function (eg. output from
5 #'      term_freq)
6 #'
7 #' @param col_name symbol name of the column insight was
8 #'      performed on
9 dist_density <- function(.data, col_name){
10   q_col_name <- dplyr::enquo(col_name)
11   .data %>%
12     ggplot2::ggplot(ggplot2::aes(x = !! q_col_name)) +
13     ggplot2::geom_density()
14 }

```

Listing 27: Create Kernel Density Estimate Plot

Box Plot

3.4.4 Structure

Time Series

```

1 #' output a ggplot time series plot of some insight function
2 #'
3 #' @param .data a dataframe containing "text" and insight columns as
4 #'      per the output of the get_(term/aggregate)_insight wrapper
5 #'      function
6 #'
7 #' @param y symbol name of the column insight was
8 #'      outputted to
9 #'
10 #' @export
11 struct_time_series <- function(.data, y){
12   q_y <- dplyr::enquo(y)
13   .data %>%
14     dplyr::filter(!is.na(!! q_y)) %>%
15     dplyr::mutate(term = seq_along(!! q_y)) %>%
16     ggplot2::ggplot(ggplot2::aes(term, !! q_y)) +
17     ggplot2::geom_line(na.rm = TRUE)
18 }

```

Listing 28: Create Time Series Plot

Page View

```
1 #' Colours a ggpage based on an insight function
2 #'
3 #' @param .data a dataframe containing "word" and insight columns as
4 #'   per the output of the get_(term|aggregate)_insight wrapper
5 #'   function
6 #'
7 #' @param col_name symbol name of the insight column intended to
8 #'   colour plot
9 #'
10 #' @param num_terms the number of terms to visualise
11 #'
12 #' @param term_index which term to start the visualisation from
13 #'
14 #' @param palette determine coloration of palette (not yet implemented)
15 #'
16 #' @return ggplot object as per ggpage
17 #'
18 #' @export
19 struct_pageview <- function(.data, col_name, num_terms, term_index, palette){
20   end <- min(nrow(.data), term_index + num_terms)
21   q_col_name <- dplyr::enquo(col_name)
22   .data[seq(term_index, end),] %>%
23     dplyr::pull(word) %>%
24     ggpage::ggpage_build() %>%
25     dplyr::bind_cols(.data[seq(term_index, end),]) %>%
26     ggpage::ggpage_plot(ggplot2::aes(fill = !! q_col_name)) ## +
27 }
```

Listing 29: Create Page View Plot

3.4.5 Wrapper

```
1 #' create a group-aware visualisation
2 #'
3 #' @param .data the standard dataframe, modified so the last column
4 #'   is the output of some insight function (eg. output from
5 #'   term_freq)
6 #'
7 #' @param vis character name of visualisation function
8 #'
9 #' @param col character name of the column to get insight from
10 #'
11 #' @param facet_by character name of the column to facet by
12 #'
13 #' @param scale_fixed force scales to be fixed in a facet
14 #'
15 #' @param ... additional arguments to the visualisation
16 #'
17 #' @export
18 get_vis <- function(.data, vis, col, facet_by="", scale_fixed = TRUE, ...){
19   vistable <- list("Page View" = struct_pageview,
20     "Time Series" = struct_time_series,
```

```

21         "Bar" = score_barplot,
22         "Density" = dist_density,
23         "Histogram" = dist_hist,
24         "Word Cloud" = score_wordcloud)
25     y <- dplyr::sym(col)
26     chart <- vistable[[vis]](.data, !! y, ...)
27     if (shiny::isTruthy(facet_by)){
28         facet_name <- dplyr::sym(facet_by)
29         q_facet_name <- dplyr::enquo(facet_name)
30         return(chart + ggplot2::facet_wrap(ggplot2::vars(!! q_facet_name),
31                                           scales = ifelse(vis == "struct_pageview" | scale_fixed,
32                                                         "fixed",
33                                                         "free")))
34     } else {
35         return(chart)
36     }
37 }
38 }

```

Listing 30: Create Visualisation

3.5 Application

```

1  library(shiny)
2  library(inzightta)
3  library(rlang)
4
5  ui <- navbarPage("iNZight Text Analytics",
6                  tabPanel("Processing",
7                          sidebarLayout(
8                              sidebarPanel(
9                                  tags$p("Import"),
10                                 fileInput("file1", "Choose File(s)",
11                                           multiple = TRUE,
12                                           accept = c("text/csv",
13                                                       "text/comma-separated-values,text/plain",
14                                                       ".csv", ".xlsx", ".xls")),
15                                 tags$hr(),
16                                 tags$p("Process"),
17                                 checkboxInput("lemmatise", "Lemmatise"),
18                                 uiOutput("sw_lexicon"),
19                                 checkboxInput("stopwords", "Stopwords"),
20                                 actionButton("prep_button", "Prepare Text"),
21                                 selectInput("section_by", "Section By",
22                                           list("", "chapter", "part", "section", "canto", "book")),
23                                 uiOutput("vars_to_filter"),
24                                 textInput("filter_pred", "value to match", ""),
25                                 mainPanel(
26                                     tableOutput("table")))),
27                          tabPanel("Visualisation",
28                              sidebarLayout(
29                                  sidebarPanel(selectInput("what_vis",
30                                                            "Select what you want to Visualise",
31                                                            list("Term Frequency",
32                                                                "n-gram Frequency",
33                                                                "Key Words",
34                                                                "Term Sentiment",
35                                                                "Moving Average Term Sentiment",
36                                                                "Aggregated Term Count",

```

```

37                                     "Key Sections",
38                                     "Aggregated Sentiment")),
39     uiOutput("group_by"),
40     uiOutput("insight_options"),
41     uiOutput("vis_options"),
42     uiOutput("vis_facet_by"),
43     actionButton("plot_button", "PLOT!!!"),
44     mainPanel(
45       plotOutput("plot")))
46   )
47
48 server <- function(input, output) {
49   imported <- reactive({
50     inzhightta::import_files(input$file1$datapath)})
51   prepped <- eventReactive(input$prep_button, {
52     imported() %>%
53       format_data(input$lemmatise, input$stopwords, input$sw_lexicon, NA)})
54   sectioned <- reactive({
55     data <- prepped()
56     if (isTruthy(input$section_by)){
57       data <- data %>%
58         section(input$section_by)}
59     data})
60   filtered <- reactive({
61     data <- sectioned()
62     if (isTruthy(input$filter_var) &
63         isTruthy(input$filter_pred)){
64       data <- data %>%
65         dplyr::filter(!! dplyr::sym(input$filter_var) == input$filter_pred)}
66     data})
67   grouped <- reactive({
68     data <- filtered()
69     if (isTruthy(input$group_var)){
70       data <- data %>%
71         dplyr::group_by(!! dplyr::sym(input$group_var))}
72     data})
73   output$table <- renderTable({
74     filtered() %>% head(300)})
75   output$sw_lexicon <- renderUI(selectInput("sw_lexicon", "Select the Stopword Lexicon",
76     stopwords::stopwords_getsources()))
77   output$vars_to_filter <- renderUI(selectInput("filter_var",
78     "select which column to apply filtering to",
79     c("", names(sectioned())) %||% c("")))
80   output$group_by <- renderUI(selectInput("group_var",
81     "select which columns to group on",
82     c("", names(filtered())) %||% c("")))
83   output$insight_options <- renderUI({
84     switch(input$what_vis,
85       "Term Frequency" = selectInput("vis_type",
86         "Select how to Visualise it",
87         list("Bar",
88             "Word Cloud",
89             "Page View",
90             "Time Series",
91             "Density",
92             "Histogram")),
93       "n-gram Frequency" = tagList(selectInput("vis_type",
94         "Select how to Visualise it",
95         list("Bar",
96             "Word Cloud",
97             "Page View",
98             "Time Series",

```

```

99         "Density",
100         "Histogram")),
101     sliderInput("n_gram",
102         "n-gram count",
103         2, 8, 2)),
104     "Key Words" = tagList(selectInput("vis_type",
105         "Select how to Visualise it",
106         list("Bar",
107             "Word Cloud",
108             "Page View",
109             "Time Series",
110             "Density",
111             "Histogram")),
112         selectInput("summ_method",
113             "Method of summary generation",
114             list("TextRank", "LexRank"))),
115     "Term Sentiment" = tagList(selectInput("vis_type",
116         "Select how to Visualise it",
117         list("Page View",
118             "Word Cloud",
119             "Time Series",
120             "Bar",
121             "Density",
122             "Histogram")),
123         selectInput("sent_lex",
124             "Lexicon for Sentiment Dictionary",
125             list("afinn", "bing",
126                 "loughran", "nrc"))),
127     "Moving Average Term Sentiment" = tagList(selectInput("vis_type",
128         "Select how to Visualise it",
129         list("Time Series",
130             "Word Cloud",
131             "Page View",
132             "Bar",
133             "Density",
134             "Histogram")),
135         sliderInput("term_sent_lag",
136             "Lag Length for Calculation of Moving Average",
137             3,500,50),
138         selectInput("sent_lex",
139             "Lexicon for Sentiment Dictionary",
140             list("afinn", "bing",
141                 "loughran", "nrc"))),
142     "Aggregated Term Count" = tagList(selectInput("vis_type",
143         "Select how to Visualise it",
144         list("Bar",
145             "Word Cloud",
146             "Page View",
147             "Time Series",
148             "Density",
149             "Histogram")),
150         selectInput("agg_var",
151             "Select which variable to aggregate on",
152             c("", names(grouped())) %||% c("")),
153     "Key Sections" = tagList(selectInput("vis_type",
154         "Select how to Visualise it",
155         list("Bar",
156             "Word Cloud",
157             "Page View",
158             "Time Series",
159             "Density",
160             "Histogram")),

```

```

161         selectInput("summ_method",
162                     "Method of summary generation",
163                     list("TextRank", "LexRank")),
164         selectInput("agg_var",
165                     "Select which variable to aggregate on",
166                     c("", names(grouped())) %||% c("")),
167         "Aggregated Sentiment" = tagList(selectInput("vis_type",
168                                                     "Select how to Visualise it",
169                                                     list("Page View",
170                                                         "Word Cloud",
171                                                         "Time Series",
172                                                         "Bar",
173                                                         "Density",
174                                                         "Histogram")),
175                                         selectInput("sent_lex",
176                                                     "Lexicon for Sentiment Dictionary",
177                                                     list("afinn", "bing",
178                                                         "loughran", "nrc")),
179                                         selectInput("agg_var",
180                                                     "Select which variable to aggregate on",
181                                                     c("", names(grouped())) %||% c("")))))}
182
183     insightful <- reactive({
184       switch(input$what_vis,
185             "Term Frequency" = get_term_insight(grouped(),
186                                                  input$what_vis),
187             "n-gram Frequency" = get_term_insight(grouped(),
188                                                  c("n-grams", "n-gram Frequency"),
189                                                  input$n_gram),
190             "Key Words" = get_term_insight(grouped(),
191                                             input$what_vis,
192                                             input$summ_method),
193             "Term Sentiment" = get_term_insight(grouped(),
194                                                  input$what_vis,
195                                                  input$sent_lex),
196             "Moving Average Term Sentiment" = get_term_insight(grouped(),
197                                                                 input$what_vis,
198                                                                 input$sent_lex,
199                                                                 input$term_sent_lag),
200             "Aggregated Term Count" = get_aggregate_insight(grouped(),
201                                                             c("Bound Aggregates", input$what_vis),
202                                                             input$agg_var),
203             "Key Sections" = get_aggregate_insight(grouped(),
204                                                     c("Bound Aggregates", input$what_vis),
205                                                     input$agg_var,
206                                                     input$summ_method),
207             "Aggregated Sentiment" = get_aggregate_insight(grouped(),
208                                                             c("Bound Aggregates", input$what_vis),
209                                                             input$agg_var,
210                                                             input$sent_lex)))}
211
212     output$vis_options <- renderUI({
213       switch(input$vis_type,
214             "Word Cloud" = tagList(sliderInput("num_terms",
215                                                  "Select the number of terms to visualise",
216                                                  3, 50, 15),
217                                     selectInput("wordcloud_shape",
218                                                  "Select the shape of the wordcloud",
219                                                  list("circle",
220                                                      "cardioid",
221                                                      "diamond",
222                                                      "square",
223                                                      "triangle-forward",
224                                                      "triangle-upright",

```

```

223                                     "pentagon",
224                                     "star"))),
225 "Page View" = tagList(sliderInput("num_terms",
226                                 "Select the number of terms to visualise",
227                                 3, 400, 100),
228 sliderInput("term_index",
229             "Select the point to begin visualisation from",
230             1, nrow(insighted()), 1),
231 selectInput("palette",
232             "Select the colour palette type",
233             list("Sequential", "Diverging"))),
234 "Bar" = tagList(sliderInput("num_terms", "Select the number of terms to visualise",
235                             2,50,15)))})
236 output$vis_facet_by <- renderUI(tagList(selectInput("vis_facet",
237                                                     "select which variable to facet on",
238                                                     c("", names(grouped())) %||% c("")),
239 checkboxInput("scale_fixed", "Scale Fixed", value=TRUE)))
240
241 visualisation <- reactive({
242   switch(input$vis_type,
243     "Word Cloud" = switch(input$what_vis,
244                           "n-gram Frequency" = get_vis(
245                             insightful(),
246                             input$vis_type,
247                             input$what_vis,
248                             input$vis_facet,
249                             input$scale_fixed,
250                             input$num_terms,
251                             x = `n-grams`,
252                             shape = input$wordcloud_shape),
253                           "Aggregated Term Count" = ,
254                           "Key Sections" = ,
255                           "Aggregated Sentiment" = get_vis(
256                             insightful(),
257                             input$vis_type,
258                             input$what_vis,
259                             input$vis_facet,
260                             input$scale_fixed,
261                             input$num_terms,
262                             x = `Bound Aggregates`,
263                             shape = input$wordcloud_shape),
264                             get_vis(insighted(),
265                                     input$vis_type,
266                                     input$what_vis,
267                                     input$vis_facet,
268                                     input$scale_fixed,
269                                     input$num_terms,
270                                     shape = input$wordcloud_shape)),
271     "Page View" = get_vis(insighted(), input$vis_type,
272                           input$what_vis,
273                           input$vis_facet,
274                           input$scale_fixed,
275                           input$num_terms,
276                           input$term_index,
277                           palette = input$palette),
278     "Time Series" = get_vis(insighted(), input$vis_type,
279                             input$what_vis,
280                             input$vis_facet,
281                             input$scale_fixed),
282     "Bar" = switch(input$what_vis,
283                   "n-gram Frequency" = get_vis(insighted(),
284                                                 input$vis_type,

```



```

285                                     input$vis_facet,
286                                     input$scale_fixed,
287                                     input$num_terms,
288                                     x = `n-grams`),
289     "Aggregated Term Count" =,
290     "Key Sections" =,
291     "Aggregated Sentiment" = get_vis(insighted(),
292                                     input$vis_type,
293                                     input$what_vis,
294                                     input$vis_facet,
295                                     input$scale_fixed,
296                                     input$num_terms,
297                                     x = `Bound Aggregates`),
298     get_vis(insighted(), input$vis_type,
299             input$what_vis,
300             input$vis_facet,
301             input$scale_fixed,
302             input$num_terms)),
303     "Density" = get_vis(insighted(), input$vis_type,
304                         input$what_vis, input$vis_facet,
305                         input$scale_fixed),
306     "Histogram" = get_vis(insighted(), input$vis_type,
307                           input$what_vis,
308                           input$vis_facet,
309                           input$scale_fixed)))
310     output$plot <- renderPlot({
311       visualisation()})
312   }
313
314   # Create Shiny app ----
315   shinyApp(ui, server)

```

Listing 31: UI and Server of Shiny Application

Chapter 4

Conclusion

4.1 Summary

4.1.1 summarise successes

4.1.2 summarise failures

4.1.3 general thoughts on the topic

4.2 Recommendations

4.2.1 educational potential of text analytics

4.2.2 what else remains

Chapter 5

Appendix

The following pages are a copy of the documentation for the R package created as a part of this dissertation. They were automatically generated through the Roxygen2 system.

Package ‘inzightta’

September 24, 2019

Title iNZight Text Analytics

Version 0.0.0.9000

Description Provides text analytics functions for the importation, analysis, and visualisation of text. This package is designed specifically for output in the shiny program, with the analytical functions all working well with dplyr tools.

License GPL-3

Encoding UTF-8

LazyData true

Imports readr,
tibble,
stringr,
dplyr,
readxl,
purrr,
tidytext,
textstem,
magrittr,
stats,
textrank,
lexRankr,
ggpage,
ggplot2,
forcats,
shiny,
ggwordcloud

RoxygenNote 6.1.1

NeedsCompilation no

Author Jason Cairns [aut, cre]

Maintainer Jason Cairns <jcai849@aucklanduni.ac.nz>

R topics documented:

aggregate_sentiment	3
bind_aggregation	3
concat_walk	4
concat_walk_i	4
determine_stopwords	5
dist_density	5
dist_hist	6
format_data	6
get_aggregate_insight	7
get_bigram	7
get_books	8
get_cantos	8
get_chapters	9
get_filetype	9
get_ngram	10
get_parts	10
get_search	11
get_sections	11
get_sw	12
get_term_insight	12
get_valid_input	13
get_vis	13
ifexp	14
import_base_file	14
import_csv	15
import_excel	15
import_files	16
import_txt	16
keywords_tr	17
key_aggregates	17
ma_term_sentiment	18
ngram_freq	18
score_barplot	19
score_wordcloud	19
section	20
shorten	20
struct_pageview	21
struct_time_series	21
table_textcol	22
term_cooccurrence	22
term_corr	23
term_count	23
term_freq	24
term_sentiment	24
ungroup_by	25

aggregate_sentiment	<i>Get statistics for sentiment over some group, such as sentence.</i>
---------------------	--

Description

Get statistics for sentiment over some group, such as sentence.

Usage

```
aggregate_sentiment(.data, aggregate_on, lexicon = "afinn",  
  statistic = mean)
```

Arguments

.data	character vector of words
aggregate_on	vector to aggregate .data over; ideally, sentence_id, but could be chapter, document, etc.
lexicon	as per term sentiment
statistic	function that accepts na.rm argument; e.g. mean, median, sd.

Value

sentiment of same length as input vector aggregated over the aggregate_on vector

bind_aggregation	<i>bind aggregate terms together</i>
------------------	--------------------------------------

Description

bind aggregate terms together

Usage

```
bind_aggregation(data, aggregate_on)
```

Arguments

data	vector of terms
aggregate_on	vector of aggregations

Value

data with every aggregation bound, as in a sentence

concat_walk	<i>concat list 1 and 2, moving past NA values</i>
-------------	---

Description

concat list 1 and 2, moving past NA values

Usage

```
concat_walk(list1, list2)
```

Arguments

list1	list or vector for first bigram token
list2	list or vector for second bigram token

Value

paste of list1 and list2, skipping NA's

concat_walk_i	<i>concat list 1 and 2 at index, skipping NA values</i>
---------------	---

Description

concat list 1 and 2 at index, skipping NA values

Usage

```
concat_walk_i(i, list1, list2)
```

Arguments

i	numeric index to assess index at
list1	list or vector for first token
list2	list or vector for second token

Value

paste of list1 and list2 at index i, skipping NA's

determine_stopwords	<i>determine stopword status</i>
---------------------	----------------------------------

Description

determine stopword status

Usage

```
determine_stopwords(.data, ...)
```

Arguments

.data	vector of words
...	arguments of get_sw

Value

a [tibble][tibble::tibble-package] equivalent to the input dataframe, with an additional stopword column

dist_density	<i>output a histogram of the distribution of some function of words</i>
--------------	---

Description

output a histogram of the distribution of some function of words

Usage

```
dist_density(.data, col_name)
```

Arguments

.data	the standard dataframe, modified so the last column is the output of some insight function (eg. output from term_freq)
col_name	symbol name of the column insight was performed on

dist_hist	<i>output a histogram of the distribution of some function of words</i>
-----------	---

Description

output a histogram of the distribution of some function of words

Usage

```
dist_hist(.data, col_name)
```

Arguments

.data	the standard dataframe, modified so the last column is the output of some insight function (eg. output from term_freq)
col_name	symbol name of the column insight was performed on

format_data	<i>takes imported one-line-per-row data and prepares it for later analysis</i>
-------------	--

Description

takes imported one-line-per-row data and prepares it for later analysis

Usage

```
format_data(.data, lemmatize = TRUE, stopwords = TRUE,  
  sw_lexicon = "snowball", addl_stopwords = NA)
```

Arguments

.data	tibble with one line of text per row
lemmatize	boolean, whether to lemmatize or not
stopwords	boolean, whether to remove stopwords or not
sw_lexicon	string, lexicon with which to remove stopwords
addl_stopwords	char vector of user-supplied stopwords

Value

a [tibble][tibble::tibble-package] with one token per line, stopwords removed leaving NA values, column for analysis named "text"

get_aggregate_insight	<i>perform group-aware aggregate operations on the data</i>
-----------------------	---

Description

perform group-aware aggregate operations on the data

Usage

```
get_aggregate_insight(.data, operations, aggregate_on, ...)
```

Arguments

.data	dataframe of terms as per output of format_data
operations	character vector of operations to perform
aggregate_on	character name of the column to perform aggregate operations on
...	additional arguments to the operation - only sensible for singular operations

Value

.data with operation columns added

get_bigram	<i>Determine bigrams</i>
------------	--------------------------

Description

Determine bigrams

Usage

```
get_bigram(.data)
```

Arguments

.data	character vector of words
-------	---------------------------

Value

character vector of bigrams

get_books	<i>sections text based on book</i>
-----------	------------------------------------

Description

sections text based on book

Usage

get_books(.data)

Arguments

.data vector to section

Value

vector of same length as .data with book numbers

get_cantos	<i>sections text based on cantos</i>
------------	--------------------------------------

Description

sections text based on cantos

Usage

get_cantos(.data)

Arguments

.data vector to section

Value

vector of same length as .data with canto numbers

get_chapters	<i>sections text based on chapters</i>
--------------	--

Description

sections text based on chapters

Usage

```
get_chapters(.data)
```

Arguments

.data	vector to section
-------	-------------------

Value

vector of same length as .data with chapter numbers

get_filetype	<i>Get filetype</i>
--------------	---------------------

Description

Get filetype

Usage

```
get_filetype(filepath)
```

Arguments

filepath	string filepath of document
----------	-----------------------------

Value

filetype (string) - NA if no extension

get_ngram	Returns the n-grams, skipping NA values
-----------	---

Description

Returns the n-grams, skipping NA values

Usage

```
get_ngram(.data, n)
```

Arguments

.data	vector to get n-grams from
n	number of n-grams to attain

Value

n-gram vector without NA values

get_parts	sections text based on parts
-----------	------------------------------

Description

sections text based on parts

Usage

```
get_parts(.data)
```

Arguments

.data	vector to section
-------	-------------------

Value

vector of same length as .data with part numbers

get_search	<i>creates a search closure to section text</i>
------------	---

Description

creates a search closure to section text

Usage

```
get_search(search)
```

Arguments

search a string regexp for the term to seperate on, e.g. "Chapter"

Value

closure over search expression

get_sections	<i>sections text based on sections</i>
--------------	--

Description

sections text based on sections

Usage

```
get_sections(.data)
```

Arguments

.data vector to section

Value

vector of same length as .data with section numbers

get_sw	<i>Gets stopwords from a default list and user-provided list</i>
--------	--

Description

Gets stopwords from a default list and user-provided list

Usage

```
get_sw(lexicon = "snowball", addl = NA)
```

Arguments

lexicon	a string name of a stopwords list, one of "smart", "snowball", or "onix"
addl	user defined character vector of additional stopwords, each element being a stop-word

Value

a [tibble][tibble::tibble-package] with one column named "word"

get_term_insight	<i>perform group-aware term operations on the data</i>
------------------	--

Description

perform group-aware term operations on the data

Usage

```
get_term_insight(.data, operations, ...)
```

Arguments

.data	dataframe of terms as per output of format_data
operations	character vector of term operations to perform
...	additional arguments to the operation - only sensible for singular operations

Value

.data with operation columns added

get_valid_input	<i>helper function to get valid input (recursively)</i>
-----------------	---

Description

helper function to get valid input (recursively)

Usage

```
get_valid_input(options, init = TRUE)
```

Arguments

options	vector of options that valid input should be drawn from
init	whether this is the initial attempt, used only as recursive information

Value

readline output that exists in the vector of options

get_vis	<i>create a group-aware visualisation</i>
---------	---

Description

create a group-aware visualisation

Usage

```
get_vis(.data, vis, col, facet_by = "", scale_fixed = TRUE, ...)
```

Arguments

.data	the standard dataframe, modified so the last column is the output of some insight function (eg. output from term_freq)
vis	character name of visualisation function
col	character name of the column to get insight from
facet_by	character name of the column to facet by
scale_fixed	force scales to be fixed in a facet
...	additional arguments to the visualisation

ifexp	<i>scheme-like if expression, without restriction of returning same-size table of .test, as ifelse() does</i>
-------	---

Description

scheme-like if expression, without restriction of returning same-size table of .test, as ifelse() does

Usage

ifexp(.test, true, false)

Arguments

- | | |
|-------|---|
| .test | predicate to test |
| true | expression to return if .test evals to TRUE |
| false | expression to return if .test evals to TRUE |

Value

either true or false

import_base_file	<i>Base case for file import</i>
------------------	----------------------------------

Description

Base case for file import

Usage

import_base_file(filepath)

Arguments

- | | |
|----------|------------------------------------|
| filepath | string filepath of file for import |
|----------|------------------------------------|

Value

imported file with document id

import_csv	<i>Import csv file</i>
------------	------------------------

Description

Import csv file

Usage

```
import_csv(filepath)
```

Arguments

filepath a string indicating the relative or absolute filepath of the file to import

Value

a [tibble][tibble::tibble-package] of each row corresponding to a line of the text file, with the column named "text"

import_excel	<i>Import excel file</i>
--------------	--------------------------

Description

Import excel file

Usage

```
import_excel(filepath)
```

Arguments

filepath a string indicating the relative or absolute filepath of the file to import

Value

a [tibble][tibble::tibble-package] of each row corresponding to a line of the text file, with the column named "text"

import_files	<i>Import any number of files</i>
--------------	-----------------------------------

Description

Import any number of files

Usage

```
import_files(filepaths)
```

Arguments

filepaths char vector of filepaths

Value

a [tibble][tibble::tibble-package] imported files with document id

import_txt	<i>Import text file</i>
------------	-------------------------

Description

Import text file

Usage

```
import_txt(filepath)
```

Arguments

filepath a string indicating the relative or absolute filepath of the file to import

Value

a [tibble][tibble::tibble-package] of each row corresponding to a line of the text file, with the column named "text"

keywords_tr	<i>Determine textrank score for vector of words</i>
-------------	---

Description

Determine textrank score for vector of words

Usage

```
keywords_tr(.data, summ_method)
```

Arguments

.data	character vector of words
summ_method	method to use for summarisation: textrank or lexrank. Doesn't do anything yet

Value

vector of scores for each word

key_aggregates	<i>get score for key sentences as per Lexrank</i>
----------------	---

Description

get score for key sentences as per Lexrank

Usage

```
key_aggregates(.data, aggregate_on, summ_method)
```

Arguments

.data	character vector of words
aggregate_on	vector to aggregate .data over; ideally, sentence_id
summ_method	method to use for summarisation: textrank or lexrank. Doesn't do anything yet

Value

lexrank scores of aggregates

ma_term_sentiment	<i>Determine the lagged sentiment of terms</i>
-------------------	--

Description

Determine the lagged sentiment of terms

Usage

```
ma_term_sentiment(.data, lexicon = "afinn", lag = 10,
  statistic = mean)
```

Arguments

.data	vector of terms
lexicon	sentiment lexicon to use, based on the corpus provided by tidytext
lag	how many (inclusive) terms to compute statistic over
statistic	base statistic used to summarise the data, capable of taking an na.rm argument

Value

vector with lagged sentiment score of each term in the input vector

ngram_freq	<i>NOT FOR PRODUCTION - STILL IN TESING. Returns the count of n-grams, skipping NA values</i>
------------	---

Description

NOT FOR PRODUCTION - STILL IN TESING. Returns the count of n-grams, skipping NA values

Usage

```
ngram_freq(.data, n)
```

Arguments

.data	vector to get n-grams from
n	number of n-grams to attain

Value

count of each associated n-gram

score_barplot	<i>output a ggplot column graph of the top texts from some insight function</i>
---------------	---

Description

output a ggplot column graph of the top texts from some insight function

Usage

```
score_barplot(.data, y, n = 15, x = text, desc = FALSE)
```

Arguments

.data	a dataframe containing "text" and insight columns as per the output of the get_(termlaggregate)_insight wrapper function
y	symbol name of the column insight was outputted to
n	number of bars to display
x	symbol name of column for insight labels
desc	bool: show bars in descending order

score_wordcloud	<i>output a ggplot wordcloud graph of the top texts from some insight function</i>
-----------------	--

Description

output a ggplot wordcloud graph of the top texts from some insight function

Usage

```
score_wordcloud(.data, y, n = 15, x = text, shape = "circle")
```

Arguments

.data	a dataframe containing "text" and insight columns as per the output of the get_(termlaggregate)_insight wrapper function
y	symbol name of the column insight was outputted to
n	number of words to display
x	symbol name of column for insight labels
shape	character: shape of the wordcloud

section	<i>Adds section column to dataframe</i>
---------	---

Description

Adds section column to dataframe

Usage

```
section(.data, section_by)
```

Arguments

.data	dataframe formatted as per output of prep process
section_by	character name of what to section over

Value

input dataframe with additional section column

shorten	<i>Shorten some text up to n characters</i>
---------	---

Description

Shorten some text up to n characters

Usage

```
shorten(.data, n)
```

Arguments

.data	character vector
n	wrap length of text

Value

shortened form of .data

struct_pageview	<i>Colours a ggpage based on an insight function</i>
-----------------	--

Description

Colours a ggpage based on an insight function

Usage

```
struct_pageview(.data, col_name, num_terms, term_index, palette)
```

Arguments

.data	a dataframe containing "word" and insight columns as per the output of the get_(termlaggregate)_insight wrapper function
col_name	symbol name of the insight column intended to colour plot
num_terms	the number of terms to visualise
term_index	which term to start the visualisation from
palette	determine coloration of palette (not yet implemented)

Value

ggplot object as per ggpage

struct_time_series	<i>output a ggplot time series plot of some insight function</i>
--------------------	--

Description

output a ggplot time series plot of some insight function

Usage

```
struct_time_series(.data, y)
```

Arguments

.data	a dataframe containing "text" and insight columns as per the output of the get_(termlaggregate)_insight wrapper function
y	symbol name of the column insight was outputted to

table_textcol	<i>Interactively determine and automatically mark the text column of a table</i>
---------------	--

Description

Interactively determine and automatically mark the text column of a table

Usage

```
table_textcol(data)
```

Arguments

data dataframe with column requiring marking

Value

same dataframe with text column renamed to "text"

term_cooccurrence	<i>Determine term cooccurrences - extremely slow</i>
-------------------	--

Description

Determine term cooccurrences - extremely slow

Usage

```
term_cooccurrence(.data, term, aggregate_on)
```

Arguments

.data character vector of terms
term character to find correlations with
aggregate_on vector to aggregate .data over; ideally, sentence_id, but could be chapter, document, etc.

Value

numeric vector of term correlations as per phi_coef

term_corr	<i>Determine term correlations - extremely slow</i>
-----------	---

Description

Determine term correlations - extremely slow

Usage

```
term_corr(.data, term, aggregate_on)
```

Arguments

.data	character vector of terms
term	character to find correlations with
aggregate_on	vector to aggregate .data over; ideally, sentence_id, but could be chapter, document, etc.

Value

numeric vector of term correlations as per phi_coef

term_count	<i>Determine the number of terms at each aggregate level</i>
------------	--

Description

Determine the number of terms at each aggregate level

Usage

```
term_count(.data, aggregate_on)
```

Arguments

.data	character vector of terms
aggregate_on	vector to split .data on for insight

Value

vector of number of terms for each aggregate level, same length as .data

term_freq	<i>Determine term frequency</i>
-----------	---------------------------------

Description

Determine term frequency

Usage

```
term_freq(.data)
```

Arguments

.data	character vector of terms
-------	---------------------------

Value

numeric vector of term frequencies

term_sentiment	<i>Determine sentiment of terms</i>
----------------	-------------------------------------

Description

Determine sentiment of terms

Usage

```
term_sentiment(.data, lexicon = "afinn")
```

Arguments

.data	vector of terms
lexicon	sentiment lexicon to use, based on the corpus provided by tidytext

Value

vector with sentiment score of each word in the vector

ungroup_by	<i>helper function to ungroup for dplyr. functions equivalently to group_by() but with standard (string) evaluation</i>
------------	---

Description

helper function to ungroup for dplyr. functions equivalently to group_by() but with standard (string) evaluation

Usage

```
ungroup_by(x, ...)
```

Arguments

x	tibble to perform function on
...	string of groups to ungroup on

Value

x with ... no longer grouped upon

Index

aggregate_sentiment, [3](#)

bind_aggregation, [3](#)

concat_walk, [4](#)
concat_walk_i, [4](#)

determine_stopwords, [5](#)
dist_density, [5](#)
dist_hist, [6](#)

format_data, [6](#)

get_aggregate_insight, [7](#)
get_bigram, [7](#)
get_books, [8](#)
get_cantos, [8](#)
get_chapters, [9](#)
get_filetype, [9](#)
get_ngram, [10](#)
get_parts, [10](#)
get_search, [11](#)
get_sections, [11](#)
get_sw, [12](#)
get_term_insight, [12](#)
get_valid_input, [13](#)
get_vis, [13](#)

ifexp, [14](#)
import_base_file, [14](#)
import_csv, [15](#)
import_excel, [15](#)
import_files, [16](#)
import_txt, [16](#)

key_aggregates, [17](#)
keywords_tr, [17](#)

ma_term_sentiment, [18](#)

ngram_freq, [18](#)

score_barplot, [19](#)
score_wordcloud, [19](#)
section, [20](#)
shorten, [20](#)
struct_pageview, [21](#)
struct_time_series, [21](#)

table_textcol, [22](#)
term_cooccurrence, [22](#)
term_corr, [23](#)
term_count, [23](#)
term_freq, [24](#)
term_sentiment, [24](#)

ungroup_by, [25](#)

Glossary

term “a word or expression that has a precise meaning in some uses or is peculiar to a science, art, profession, or subject”[1] — here text analysts have capitalised on the generalisation of “term” to include subcomponents or aggregations of words. 9

Bibliography

- [1] Merriam-Webster Dictionary, ed. *Term — Definition of Term*. 17th Aug. 2019. URL: <https://www.merriam-webster.com/dictionary/term> (cit. on p. 69).