

Guide 6.2 Introduction to text mining in R

October 2, 2019

1 Guía 6.2: Introducción a la minería de datos textuales en R

Computación 2, IES. Profesor: Eduardo Jorquera, eduardo.jorquera@postgrado.uv.cl

2 Ordenando los datos de Jane Austen

Las novelas de la escritora están en el paquete `janeaustenr` de R. El paquete mantiene el formato de una línea por fila, donde cada línea hace referencia a cada línea literal de un libro impreso. Usemos la función `mutate` para anotar los números de línea para poder guardar las líneas en el formato original y `chapter` (usando `regex`) para hallar dónde están los capítulos.

```
In [2]: library(janeaustenr)
library(dplyr)
library(stringr)

original_books <- austen_books() %>%
  group_by(book) %>%
  mutate(linenumber = row_number(),
         chapter = cumsum(str_detect(text, regex("^chapter [\\divxlc]",
                                                ignore_case = TRUE)))) %>%
  ungroup()

original_books
```

	text <chr>	book <fct>
	SENSE AND SENSIBILITY	Sense & Sensibility
	by Jane Austen	Sense & Sensibility
	(1811)	Sense & Sensibility
	CHAPTER 1	Sense & Sensibility
	The family of Dashwood had long been settled in Sussex. Their estate was large, and their residence was at Norland Park, in the centre of their property, where, for many generations, they had lived in so respectable a manner as to engage the general good opinion of their surrounding acquaintance. The late owner of this estate was a single man, who lived to a very advanced age, and who for many years of his life, had a constant companion and housekeeper in his sister. But her death, which happened ten years before his own, produced a great alteration in his home; for to supply her loss, he invited and received into his house the family of his nephew Mr. Henry Dashwood, the legal inheritor of the Norland estate, and the person to whom he intended to bequeath it. In the society of his nephew and niece, and their children, the old Gentleman's days were comfortably spent. His attachment to them all increased. The constant attention of Mr. and Mrs. Henry Dashwood to his wishes, which proceeded not merely from interest, but from goodness of heart, gave him every degree of solid comfort which his age could receive; and the cheerfulness of the children added a relish to his existence.	Sense & Sensibility
A tibble: 73422 CE 4		
	Her recent good offices by Anne had been enough in themselves, and their marriage, instead of depriving her of one friend, secured her two. She was their earliest visitor in their settled life; and Captain Wentworth, by putting her in the way of recovering her husband's property in the West Indies, by writing for her, acting for her, and seeing her through all the petty difficulties of the case with the activity and exertion of a fearless man and a determined friend, fully requited the services which she had rendered, or ever meant to render, to his wife.	Persuasion
	Mrs Smith's enjoyments were not spoiled by this improvement of income, with some improvement of health, and the acquisition of such friends to be often with, for her cheerfulness and mental alacrity did not fail her; and while these prime supplies of good remained, she might have bid defiance even to greater accessions of worldly prosperity. She might have been absolutely rich and perfectly healthy, and yet be happy. Her spring of felicity was in the glow of her spirits, as her friend Anne's was in the warmth of her heart. Anne was tenderness itself, and she had the full worth of it in Captain Wentworth's	Persuasion

Para trabajar con este dataset ordenado, necesitamos reestructurarlo en el formato *un token por fila*, lo cual se puede hacer con la función `unnest_tokens`:

```
In [4]: library(tidytext)
        tidy_books <- original_books %>%
          unnest_tokens(word, text)

        head(tidy_books)
```

	book <fct>	linenumber <int>	chapter <int>	word <chr>
A tibble: 6 × 4	Sense & Sensibility	1	0	sense
	Sense & Sensibility	1	0	and
	Sense & Sensibility	1	0	sensibility
	Sense & Sensibility	3	0	by
	Sense & Sensibility	3	0	jane
	Sense & Sensibility	3	0	austen

Esta función usa el paquete `tokenizers` para separar cada línea de texto en el dataframe original a tokens. La tokenización por defecto es por palabra, pero se pueden incluir opciones que permita hacerlo por sentencia, párrafos, líneas, o separación respecto a un patrón de regex (expresiones regulares).

Ahora que los datos están en el formato *un término por fila*, podemos usar herramientas de manipulación con `dplyr`. Usualmente, en el análisis de texto queremos quitar las palabras vacías; éstas son, palabras que no aportan información. Podemos quitar estas palabras, con un `anti_join`:

```
In [6]: data(stop_words)

        tidy_books <- tidy_books %>%
          anti_join(stop_words)
```

Joining, by = "word"

Las palabras vacías en el paquete `tidytext` contienen palabras vacías de tres léxicos. Podemos usarlos todo junto, como lo tenemos aquí, o filtrar (`filter()`) para sólo usar un conjunto de palabras vacías si es más apropiado para ciertos análisis.

También podemos usar `count`, de `dplyr` para encontrar las palabras más comunes en todos los libros.

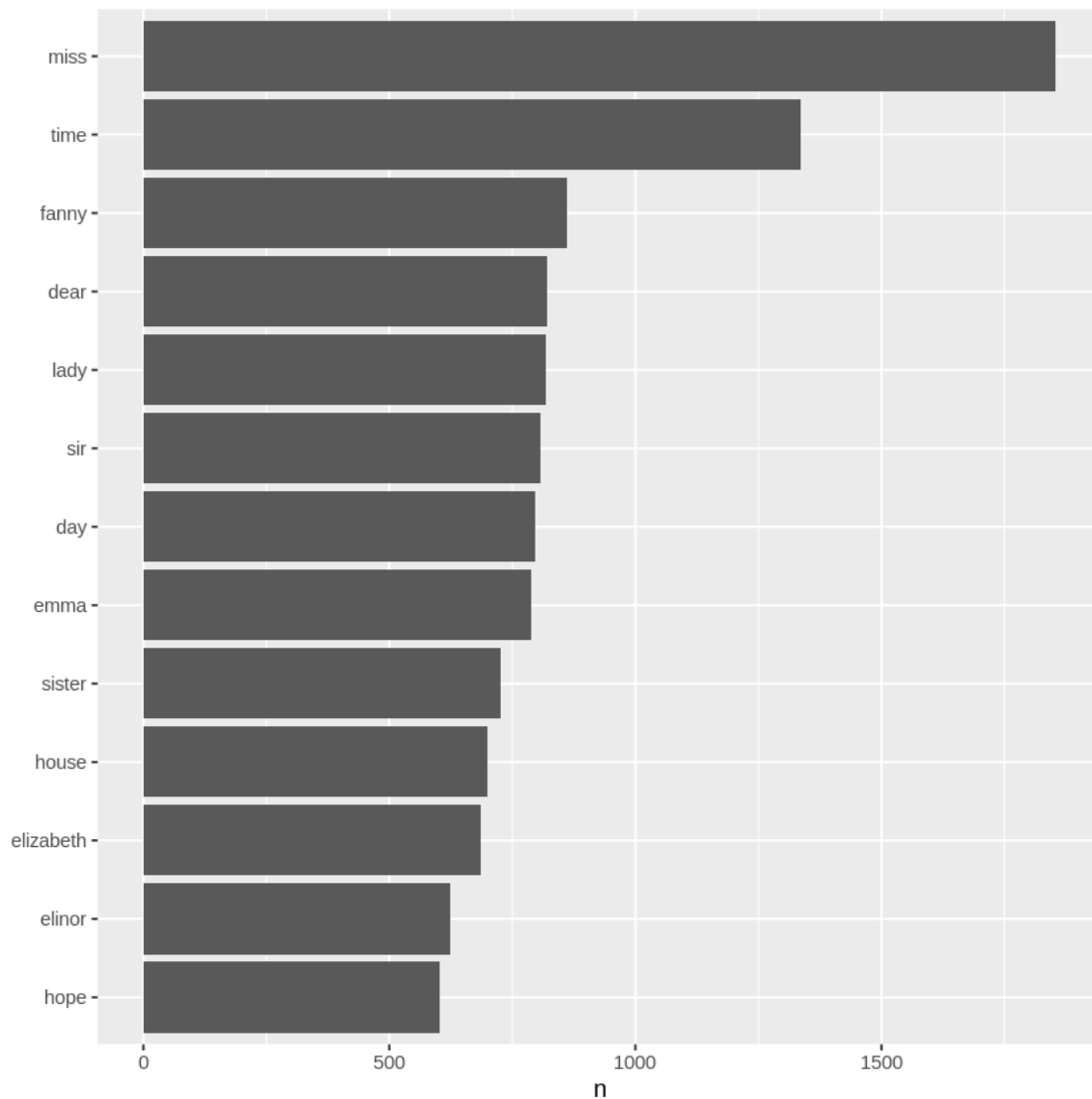
```
In [9]: head(tidy_books %>% count(word, sort = TRUE) )
```

	word <chr>	n <int>
A tibble: 6 × 2	miss	1855
	time	1337
	fanny	862
	dear	822
	lady	817
	sir	806

Ya que usamos herramientas de datos ordenados, nuestro conteo es almacenado como datos ordenados. Esto nos permite ir directamente a `ggplot2`, por ejemplo:

```
In [10]: library(ggplot2)

tidy_books %>%
  count(word, sort = TRUE) %>%
  filter(n > 600) %>%
  mutate(word = reorder(word, n)) %>%
  ggplot(aes(word, n)) +
  geom_col() +
  xlab(NULL) +
  coord_flip()
```



Note que la función `austen_books()` entregó inmediatamente el texto que queremos analizar, pero para otros casos hay que limpiar el texto lo suficiente. Nota que hay más textos disponibles para análisis en paquetes como `gutenbergr`.

Para obtener la frecuencia de las palabras, podemos usar los principios de los datos ordenados. Usemos algunas colecciones de ciencia ficción y novelas fantásticas de H.G. Wells: The Time Machine, The War of the Worlds, The Invisible Man, and The Island of Doctor Moreau. Podemos acceder a estos trabajos de la siguiente manera:

```
In [11]: library(gutenbergr)
```

```
hgwells <- gutenbergr_download(c(35, 36, 5230, 159))
```

Determining mirror for Project Gutenberg from <http://www.gutenberg.org/robot/harvest>
Using mirror <http://aleph.gutenberg.org>

```
In [12]: tidy_hgwells <- hgwells %>%  
  unnest_tokens(word, text) %>%  
  anti_join(stop_words)
```

Joining, by = "word"

3 Actividad:

¿Cuáles son las palabras más comunes en las novelas guardadas de H.G. Wells?

```
In [16]: bronte <- gutenbergr_download(c(1260, 768, 969, 9182, 767))  
tidy_bronte <- bronte %>%  
  unnest_tokens(word, text) %>%  
  anti_join(stop_words)
```

Joining, by = "word"

```
In [ ]:
```

Ahora calculemos la frecuencia para cada palabra de los trabajos de Jane Austen, hermanas Brontë y H.G. Wells, juntando dataframes. Podemos usar spread y gather de tidyr para cambiar la forma nuestro dataframe, y así sólo necesitamos hacer un gráfico para comparar estas tres novelas.

```
In [17]: library(tidyr)
```

```
frequency <- bind_rows(mutate(tidy_bronte, author = "Brontë Sisters"),  
  mutate(tidy_hgwells, author = "H.G. Wells"),  
  mutate(tidy_books, author = "Jane Austen")) %>%  
  mutate(word = str_extract(word, "[a-z']+")) %>%  
  count(author, word) %>%  
  group_by(author) %>%  
  mutate(proportion = n / sum(n)) %>%  
  select(-n) %>%  
  spread(author, proportion) %>%  
  gather(author, proportion, `Brontë Sisters`, `H.G. Wells`)
```

Usamos `str_extract()` ya que los textos tienen codificación UTF-8, y hay algunas palabras con guión bajo que indican cierto énfasis (como cursivas). El tokenizador cuenta separadamente las palabras, pero no queremos contar cualquier "*palabra*" de manera separada a "cualquier" otra palabra, como vimos anteriormente, por lo que usaremos `str_extract()`

```
In [18]: library(scales)
```

```
# expect a warning about rows with missing values being removed
ggplot(frequency, aes(x = proportion, y = `Jane Austen`, color = abs(`Jane Austen` - 1)
  geom_abline(color = "gray40", lty = 2) +
  geom_jitter(alpha = 0.1, size = 2.5, width = 0.3, height = 0.3) +
  geom_text(aes(label = word), check_overlap = TRUE, vjust = 1.5) +
  scale_x_log10(labels = percent_format()) +
  scale_y_log10(labels = percent_format()) +
  scale_color_gradient(limits = c(0, 0.001), low = "darkslategray4", high = "gray75")
  facet_wrap(~author, ncol = 2) +
  theme(legend.position="none") +
  labs(y = "Jane Austen", x = NULL)
```

Warning message:

Removed 41357 rows containing missing values (geom_point).Warning message:

Removed 41359 rows containing missing values (geom_text).



In []: