

Programas

```
load("data/programas.RData")

library(tidytext)

## Warning: package 'tidytext' was built under R version 3.6.1

library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(tidyr)
library(stringr)
library(scales)
library(readxl)
library(ggplot2)

stop_words_es <- read.csv("data/stop-words-spanish.csv", stringsAsFactors = FALSE, fileEncoding = "utf-8")
sentiments <- read.csv("data/sentimientos.csv", stringsAsFactors = FALSE, fileEncoding = "utf-8")

tidy_prog = function(txt){
  text_df <- tibble(line = 1:length(txt), text = txt)
  text_df = text_df %>%
    unnest_tokens(word, text) %>%
    mutate(word = gsub("[[:punct:]]", "", word)) %>%
    mutate(word = gsub("[[:digit:]]", "", word)) %>%
    anti_join(stop_words_es) %>%
    filter(!word %in% c("aa", "", "as"))
}

for (i in partidos){
  assign(paste0("tidy_", i), tidy_prog(get(i)))
}

## Joining, by = "word"

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

prog = function(prog, name){
  aux = data.frame("text" = prog, "partido" = name)
}
```

```

for (i in partidos){
  assign(paste0("prog_",i), prog(get(i),i))
}
programas <- bind_rows(prog_ca, prog_fa,prog_pc,prog_pg, prog_pi, prog_pn, prog_up)

## Warning in bind_rows_(x, .id): Unequal factor levels: coercing to character
## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): Unequal factor levels: coercing to character
## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

## Warning in bind_rows_(x, .id): binding character and factor vector,
## coercing into character vector

tidy_books<- bind_rows(mutate(tidy_ca, partido = "ca"),
  mutate(tidy_pc, partido = "pc"),
  mutate(tidy_pg, partido = "pg"),
  mutate(tidy_pi, partido = "pi"),
  mutate(tidy_pn, partido = "pn"),

```

```

mutate(tidy_up, partido = "up"),
mutate(tidy_fa, partido = "fa"))

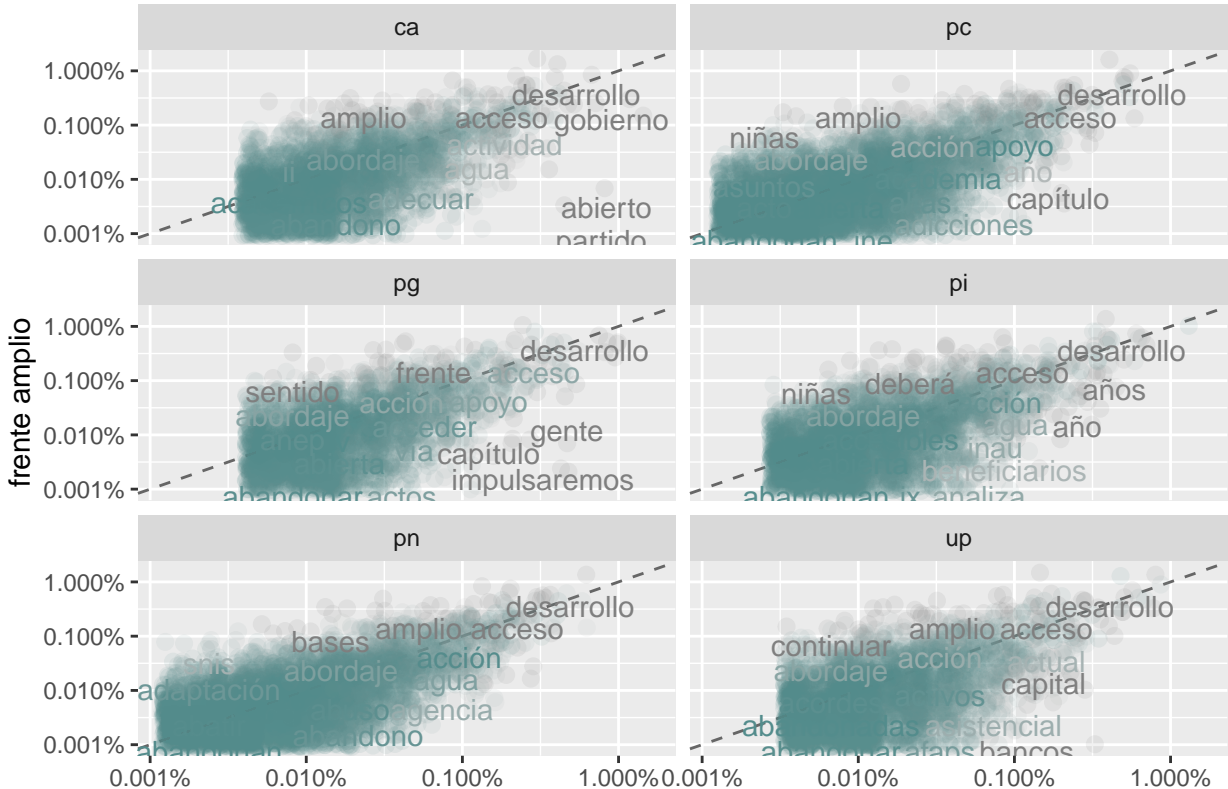
frequency <- bind_rows(mutate(tidy_ca, partido = "ca"),
                        mutate(tidy_pc, partido = "pc"),
                        mutate(tidy_pg, partido = "pg"),
                        mutate(tidy_pi, partido = "pi"),
                        mutate(tidy_pn, partido = "pn"),
                        mutate(tidy_up, partido = "up"),
                        mutate(tidy_fa, partido = "fa")) %>%
# mutate(word = str_extract(word, "[a-z']+")) %>%
count(partido, word) %>%
group_by(partido) %>%
mutate(proportion = n / sum(n)) %>%
select(-n) %>%
spread(partido, proportion) %>%
# melt(id.vars="word", variable.name="partido",value.name = "frec" )
gather(partido, proportion, c("ca","pc","pg","pi","pn","up"))

#comparacion contra el FA de los partidos de proporcion de palabras
ggplot(frequency, aes(x = proportion, y = `fa`, color = abs(`fa` - proportion))) +
  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(~partido, ncol = 2) +
  theme(legend.position="none") +
  labs(title="Proporción de frecuencia de palabras con el FA", y = "frente amplio", x = NULL)

## Warning: Removed 93244 rows containing missing values (geom_point).
## Warning: Removed 93244 rows containing missing values (geom_text).

```

Proporción de frecuencia de palabras con el FA

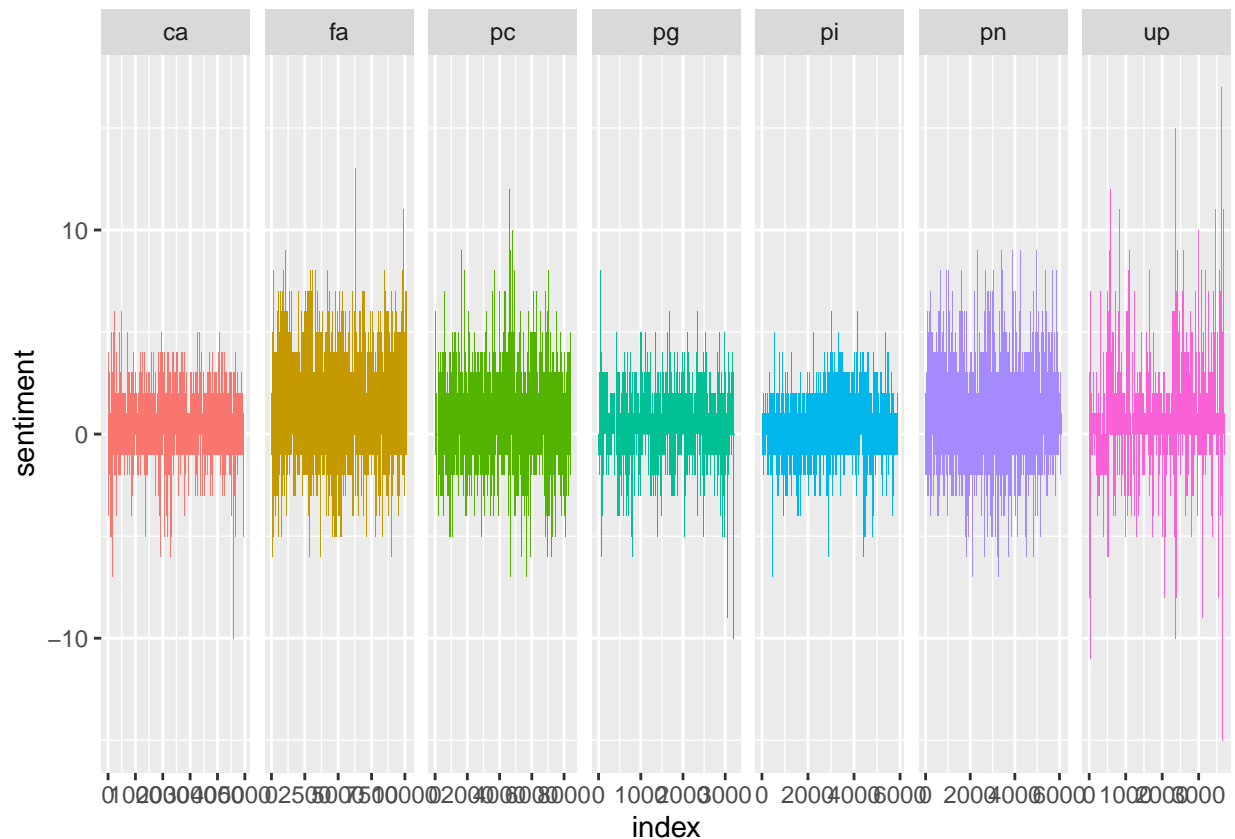


```
sent=sentiments %>%
# filter(Joy==1) %>%
  filter(Spanish..es!="NO TRANSLATION") %>%
  rename(word=Spanish..es., positive=Positive, negative=Negative) %>%
  select(word, positive,negative) %>%
  gather(sentiment,value,positive:negative) %>%
  filter(value==1)
```

```
sent_fa=tidy_books %>%
  inner_join(sent) %>%
  count(partido, index = line %/% 1, sentiment) %>%
  spread(sentiment, n, fill = 0) %>%
  mutate(sentiment = positive - negative)
```

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

```
ggplot(sent_fa, aes(index, sentiment, fill = partido)) +  
  geom_col(show.legend = FALSE)+  
  facet_grid(~partido, scales="free")
```



```
book_words <- tidy_books %>%
  select(partido, word) %>%
  count(partido, word, sort = TRUE)

total_words <- book_words %>%
  group_by(partido) %>%
  summarize(total = sum(n))

book_words <- left_join(book_words, total_words)
```

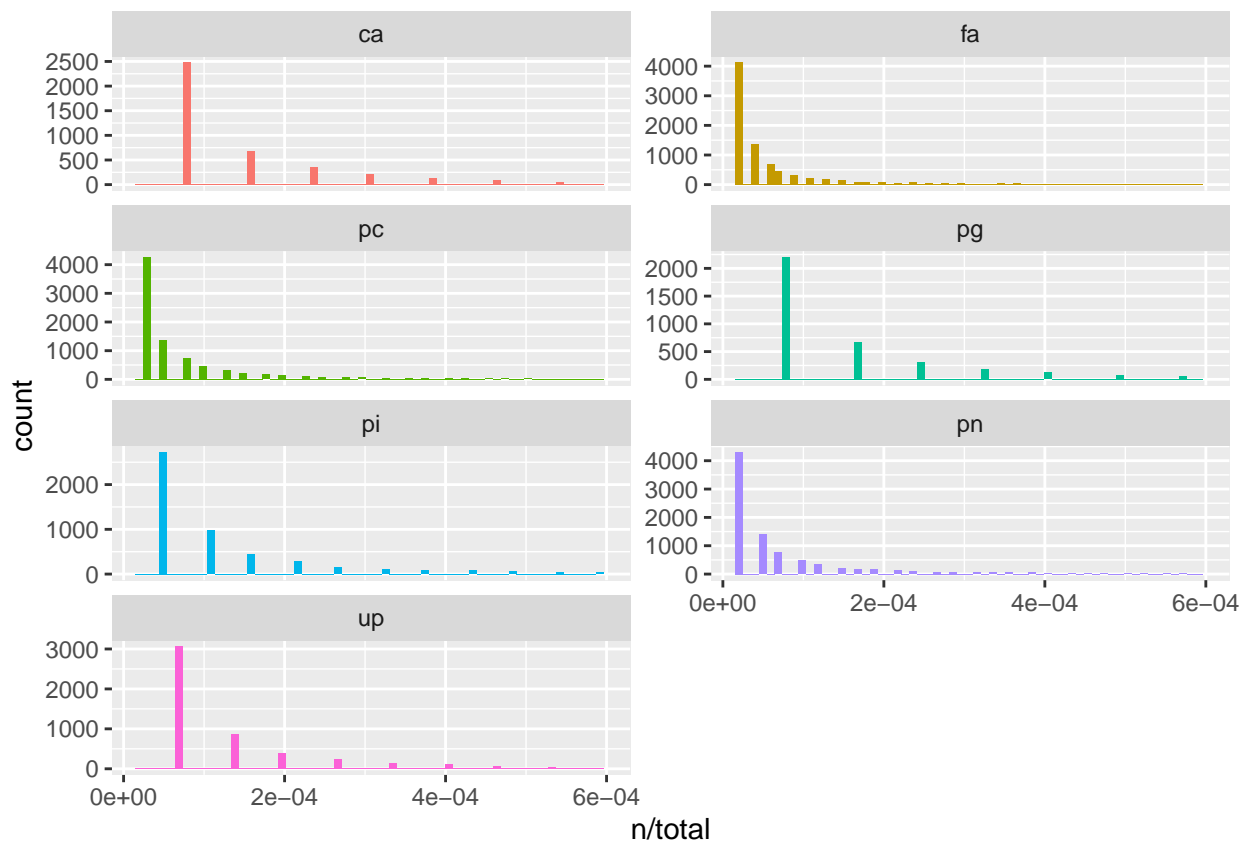
```
## Joining, by = "partido"
```

```
freq_by_rank <- book_words %>%
  group_by(partido) %>%
  mutate(rank = row_number(),
         `term frequency` = n/total)

ggplot(book_words, aes(n/total, fill = partido)) +
  geom_histogram(show.legend = FALSE, bins=60) +
  xlim(NA, 0.0006) +
  facet_wrap(~partido, ncol = 2, scales = "free_y")
```

```
## Warning: Removed 2153 rows containing non-finite values (stat_bin).
```

```
## Warning: Removed 7 rows containing missing values (geom_bar).
```



```
programas_bigrams <- programas %>%
  unnest_tokens(bigram, text, token = "ngrams", n = 2, to_lower = TRUE)
```

```
bigrams_separated <- programas_bigrams %>%
  separate(bigram, c("word1", "word2"), sep = " ")
```

```
bigrams_filtered <- bigrams_separated %>%
  filter(!word1 %in% stop_words_es$word) %>%
  filter(!word2 %in% stop_words_es$word) %>%
  filter(!word1 %in% c("aa", "", "as", "pág")) %>%
  filter(!word2 %in% c("aa", "", "as", "pág")) %>%
  filter(!(partido=="fa" & word1=="frente")) %>%
  filter(!(partido=="fa" & word1=="amplio")) %>%
  filter(!(partido=="fa" & word2=="frente")) %>%
  filter(!(partido=="fa" & word2=="amplio")) %>%
  filter(!(partido=="ca" & word1=="cabildo")) %>%
  filter(!(partido=="ca" & word1=="abierto")) %>%
  filter(!(partido=="ca" & word2=="cabildo")) %>%
  filter(!(partido=="ca" & word2=="abierto")) %>%
  filter(!(partido=="pc" & word1=="partido")) %>%
  filter(!(partido=="pc" & word1=="colorado")) %>%
  filter(!(partido=="pc" & word2=="partido")) %>%
  filter(!(partido=="pc" & word2=="colorado")) %>%
  filter(!(partido=="pi" & word1=="partido")) %>%
  filter(!(partido=="pi" & word1=="independiente")) %>%
```

```

filter(!(partido=="pi" & word2=="partido")) %>%
filter(!(partido=="pi" & word2=="independiente")) %>%
filter(!(partido=="pn" & word1=="partido")) %>%
filter(!(partido=="pn" & word1=="nacional")) %>%
filter(!(partido=="pn" & word2=="partido")) %>%
filter(!(partido=="pn" & word2=="nacional")) %>%
filter(!(partido=="up" & word1=="unidad")) %>%
filter(!(partido=="up" & word1=="popular")) %>%
filter(!(partido=="up" & word2=="unidad")) %>%
filter(!(partido=="up" & word2=="popular"))

# new bigram counts:
bigram_counts <- bigrams_filtered %>%
# group_by(partido) %>%
  count(word1, word2, sort = TRUE)

nombres_partidos = c("cabildo abierto","frente amplio","partido colorado", "partido independiente",
                     "partido nacional","unidad popular")

bigram_tf_idf <- bigrams_filtered %>%
  unite(bigram, word1, word2, sep = " ") %>%
  filter(!bigram %in% nombres_partidos) %>%
  count(partido, bigram) %>%
  bind_tf_idf(bigram, partido, n) %>%
  arrange(desc(tf_idf))

bigram_tf_idf %>%
  group_by(partido) %>%
  top_n(10, tf_idf) %>%
  ungroup() %>%
  ggplot(aes(reorder(bigram,tf_idf),tf_idf, fill = partido)) +
  geom_bar(stat="identity") +theme(legend.position = "none")+
  facet_wrap(~partido, ncol = 2, scales = "free")+
  coord_flip()

```



```

bigram_graph <- bigrams_filtered %>%
  filter(partido==i) %>%
  select(-partido) %>%
  count(word1, word2, sort = TRUE) %>%
  filter(n>5) %>%
  graph_from_data_frame()

a <- grid::arrow(type = "closed", length = unit(.15, "inches"))

fig=ggraph(bigram_graph, layout = "fr") +
  geom_edge_link(aes(edge_alpha = n), show.legend = FALSE,
    arrow = a, end_cap = circle(.07, 'inches')) +
  geom_node_point(color = "lightblue", size = 5) +
  geom_node_text(aes(label = name), vjust = 1, hjust = 1) +
  theme_void()
print(fig)
}

```

