

# hw3

Juan Vila

11/25/2019

## Problem Set 3

### Part 1 and 2

```
## Warning: package 'tibble' was built under R version 3.5.2
## Warning: package 'tidyr' was built under R version 3.5.2
## Warning: package 'dplyr' was built under R version 3.5.2
## Warning: package 'tidytext' was built under R version 3.5.2
# LOADING TEXTS
#1.
## FOR MAC:
texts <- file.path("~", "Documents", "GitHub", "problem-set-3", "Party Platforms Data")

# Now we can create our raw corpus, which we will preprocess in a moment
docs <- VCorpus(DirSource(texts))

docs <- tm_map(docs, removeNumbers)
docs <- tm_map(docs, removePunctuation)
docs <- tm_map(docs, content_transformer(tolower))
docs <- tm_map(docs, removeWords, stopwords("english"))

# Maybe we need a bit more cleaning of unique characters
for (j in seq(docs)) {
  docs[[j]] <- gsub("/", " ", docs[[j]])
  docs[[j]] <- gsub("'", " ", docs[[j]])
  docs[[j]] <- gsub("-", " ", docs[[j]])
  docs[[j]] <- gsub("\\\\", " ", docs[[j]])
  docs[[j]] <- gsub("@", " ", docs[[j]])
  docs[[j]] <- gsub("\u2028", " ", docs[[j]]) # an ascii character that does not translate
}

for (j in seq(docs)) {
  docs[[j]] <- gsub("federal government", "federal_government", docs[[j]])
  docs[[j]] <- gsub("united states", "united_states", docs[[j]])
}
(docs <- tm_map(docs, PlainTextDocument))

## <<VCorpus>>
## Metadata: corpus specific: 0, document level (indexed): 0
## Content: documents: 3
```

```
docs <- tm_map(docs, stripWhitespace)
# manually removing words too for your specific purposes (e.g., our "representative" or "honourable" ex

docs <- tm_map(docs, removeWords, c("will","must","also"))
docs <- tm_map(docs, removeWords, c("democrats","republicans"))
```

We clean the following items: 1. numbers 2. Punctuation 3. We transform the data into lower letters 4. Stop words.

Also, we eliminate words that are repited but they don't have content in the analysis because are words that refer to the actions that they will do as will and must. Also, we eliminate the word also due to that the role that have in the phrases is to connect them. Additionally, we see that democrats and republicans also have a high frequency especially in the case of democrats, this could imply that is their plataform is much more self-centered or trying to create a brand. Nevertheless, we decide to eliminate both, because as if even this word are high in the ranking this maybe bias the analysis of the content of the words. Finally, we decide also after analyzing some results to consider some words as one, for example Federal Government and United States, because this are one concept.

## Part 3

After doing this we analyze the basic frequencies found in both corpus. We analyze the sixth most repeated word in the Corpus. Now we proceed to analyze them by each party.

### Sixth more used words by Democrats

##	health	support	believe	people	americans	american
##	130	123	117	111	94	86

In the case of the Democrats we found something very interesting they use three words that in the context of the text means the same, we are talking of people, american and americans which rank from fourth to six. Also, we can see that the most used words are health, support and believe. The word health is used in the context of healthcare policy, support in any policy focused in help or bussiness, students or person in any kind of things as paying loans, helping farmers, first home buyers, etc.

### Sixth more used words by Republicans

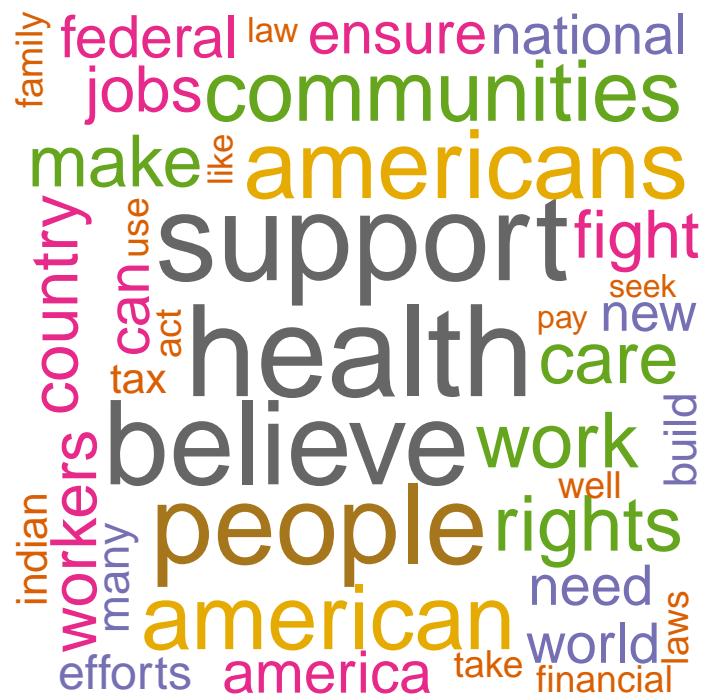
##	american	government	federal	support	people	national
##	121	110	106	100	98	83

In the case of the republicans they use also people(5th), national(6th) and american(1th), this is interesting taking into consideration that not even in the most 15 words used the word amaricans(19th) is not appear. Specially taking into consideration the fact that the Republican Ideological agenda is always focused in a individualized view of how to construct a society and they nationalistic views. In other hand, we see that they made reference to "government"(2nd), and "federal"(3th), appears beetween the places 2nd until 3th, the context where they used this word reading them in the corpus is criticizing big goverment. For example, "Big government undermines federalism through more than 1,100 grants-in-aid programs that consume more than one-sixth of the federal budget for matters that should be the exclusive responsibility of the states.". Finally, in the case of the case of support that have in the context depends in how the republican party supports different policies.

Now if we build the wordcloud for the Republican and Democratic corpus, we found the follwing result. Besides the six first we see the are words as ensure, take, jobs, work, rights, need, student, climate that star to show what are the elements that build the Democratic Plataform that made sense to the things that we

expect that a democrats would speak for and what are the agents and policy which their efforts if they are government will be focus. The same is true in the case of the republicans. Because we found words as economic, adminsitration, economy, private,tax, religuous, military,etc. This made sanse for what is expected how the policy that the Republican Party want to implement, policy related to free market economics, where the military have a predominant role in society and where the religious right, it's free to educate young people as they please.

### Democrats agenda Frequency Wordcloud



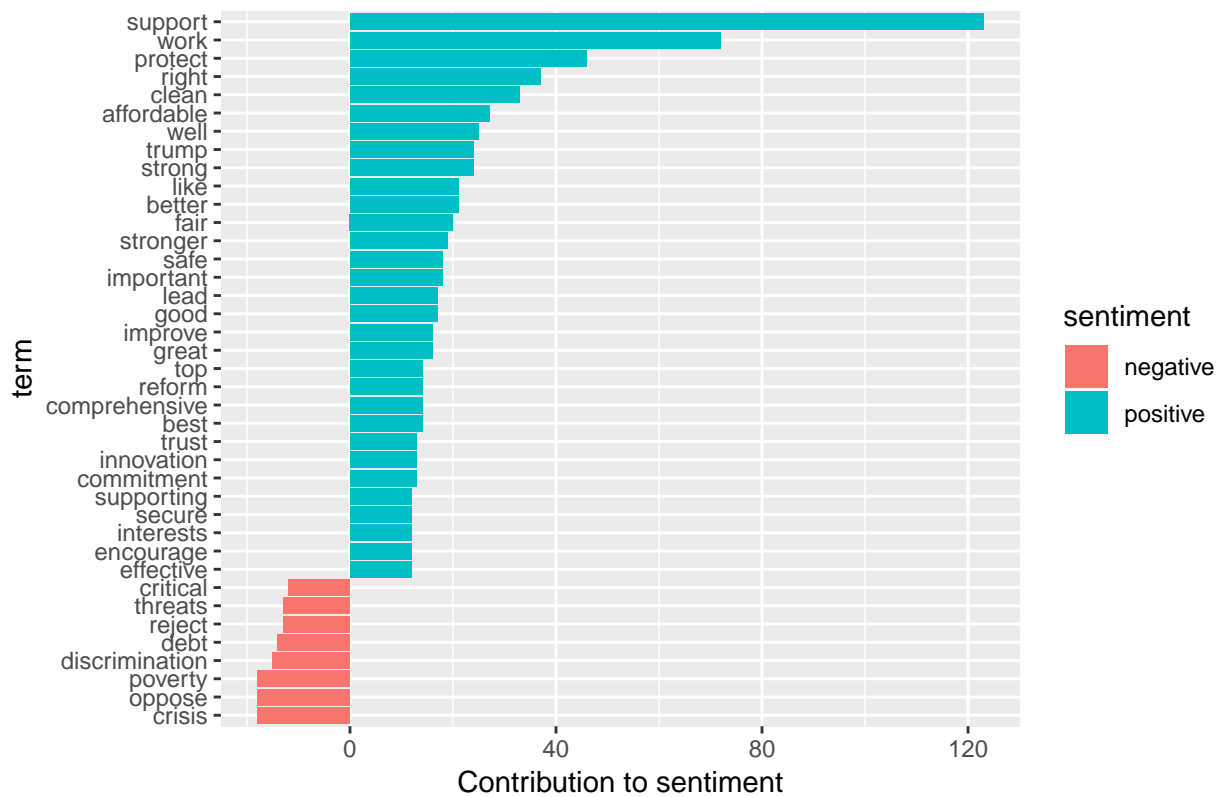
[illegible]

Bing

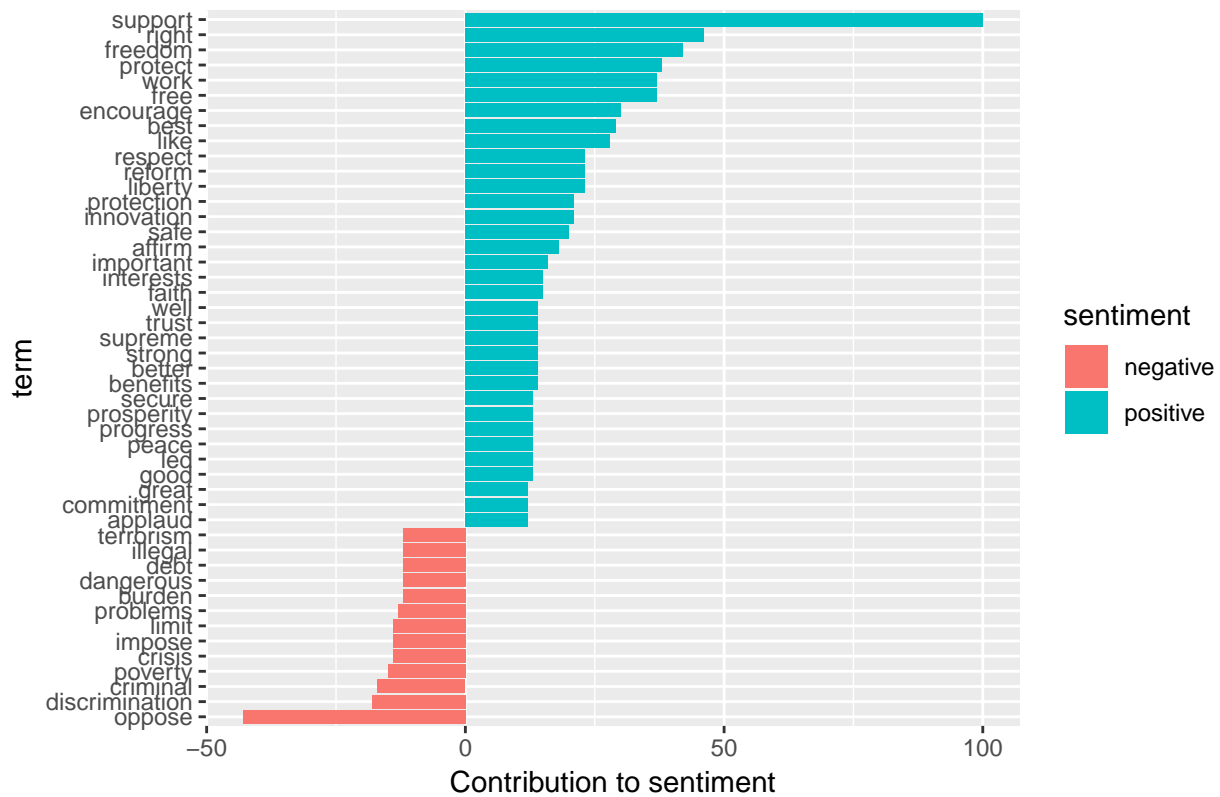
```
## Current value sentiment Rep. Party: -182
```

4

Dem. Platform Contribution to Sentiment (More than 12 counts)

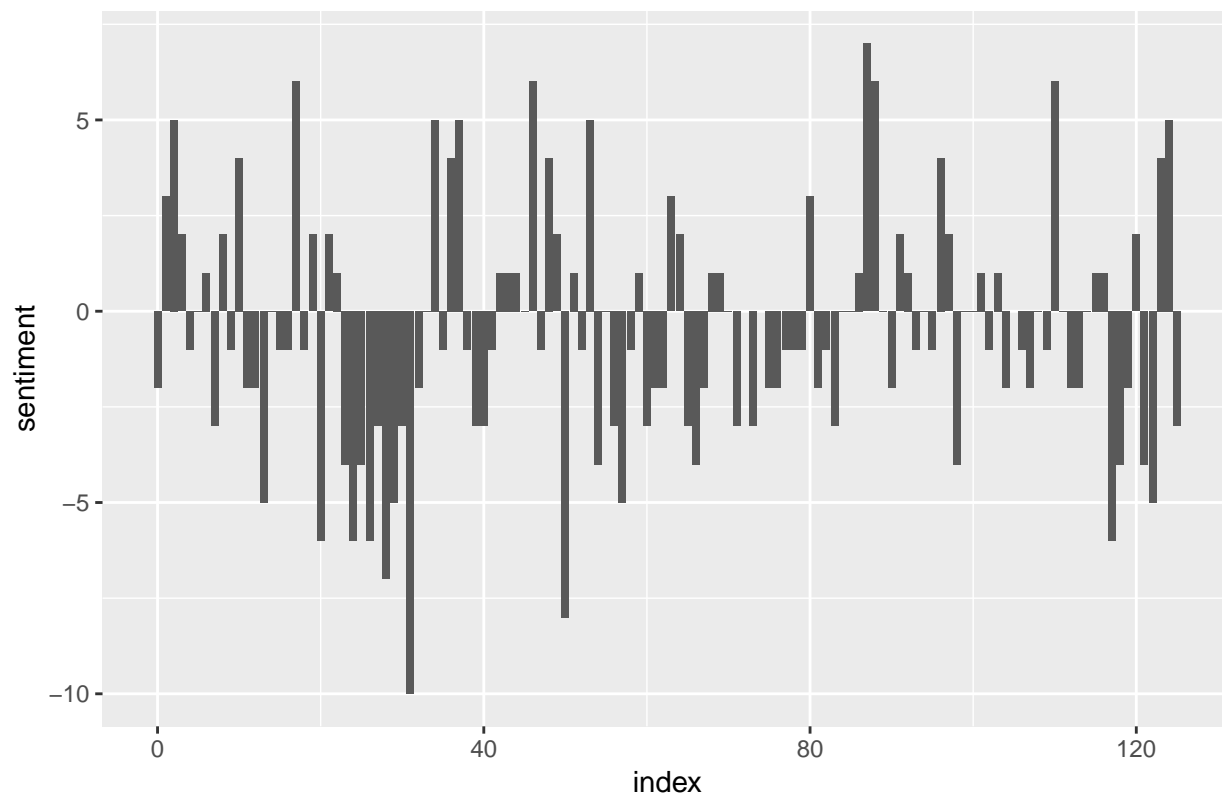


Rep. Platform Contribution to Sentiment (More than 12 counts)

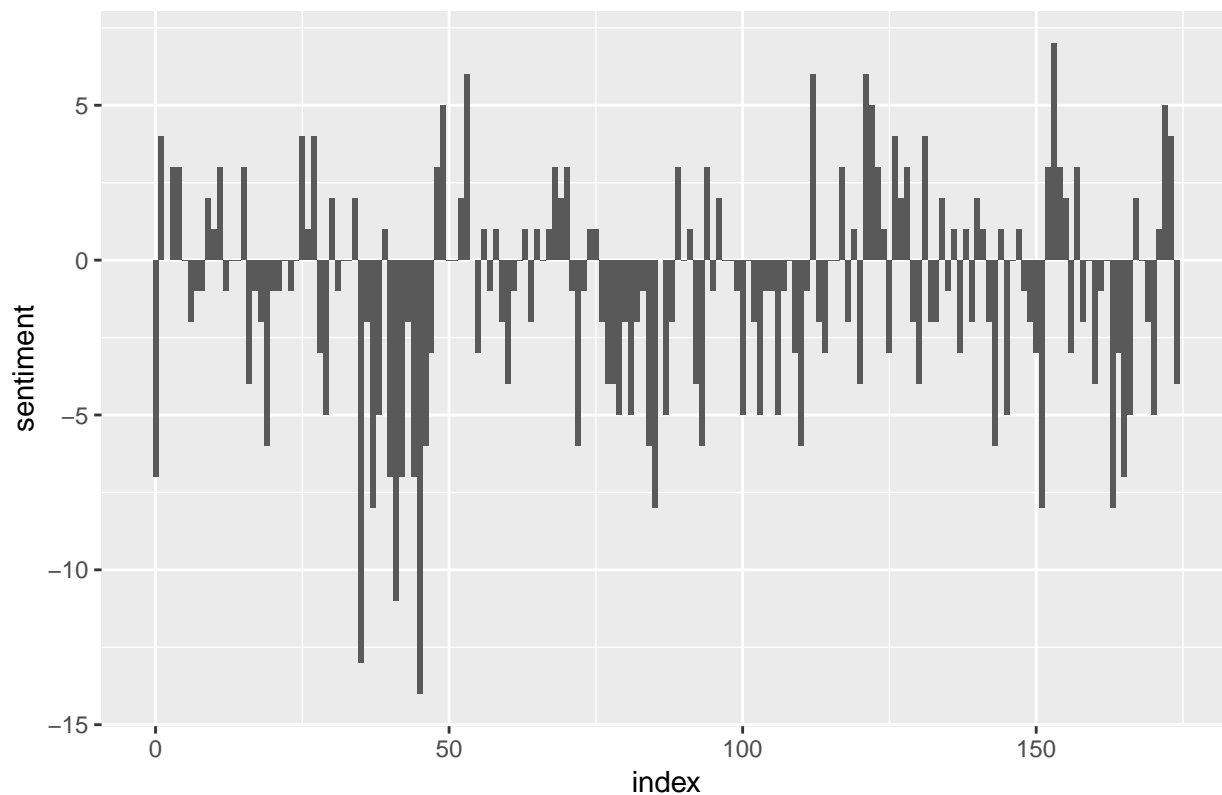


This overall negativity is noted when we bunch the merged word with the dictionary and see how the score evolves through the document, we see that are parts where the maximum negativity comes in the indexes 30 for Democrats and 50 for Republicans, but in this parts Republicans are more way negative than Democrats.

### Using Bing dictionary, change in sensitivity in Dem. Platform



## Using Bing dictionary, change in sensitivity in Rep. Platform



## AFINN

One of the major limitations of the Bing dictionary is that not take into consideration that even if a word is negative or positive they can have a different weight depending of the word. For taking that into consideration we now use the AFINN dictionary. At first, we found that the value sentiment for this dictionary are positive in aggregate, but we still see that the Democratic Party look is more positive in general.

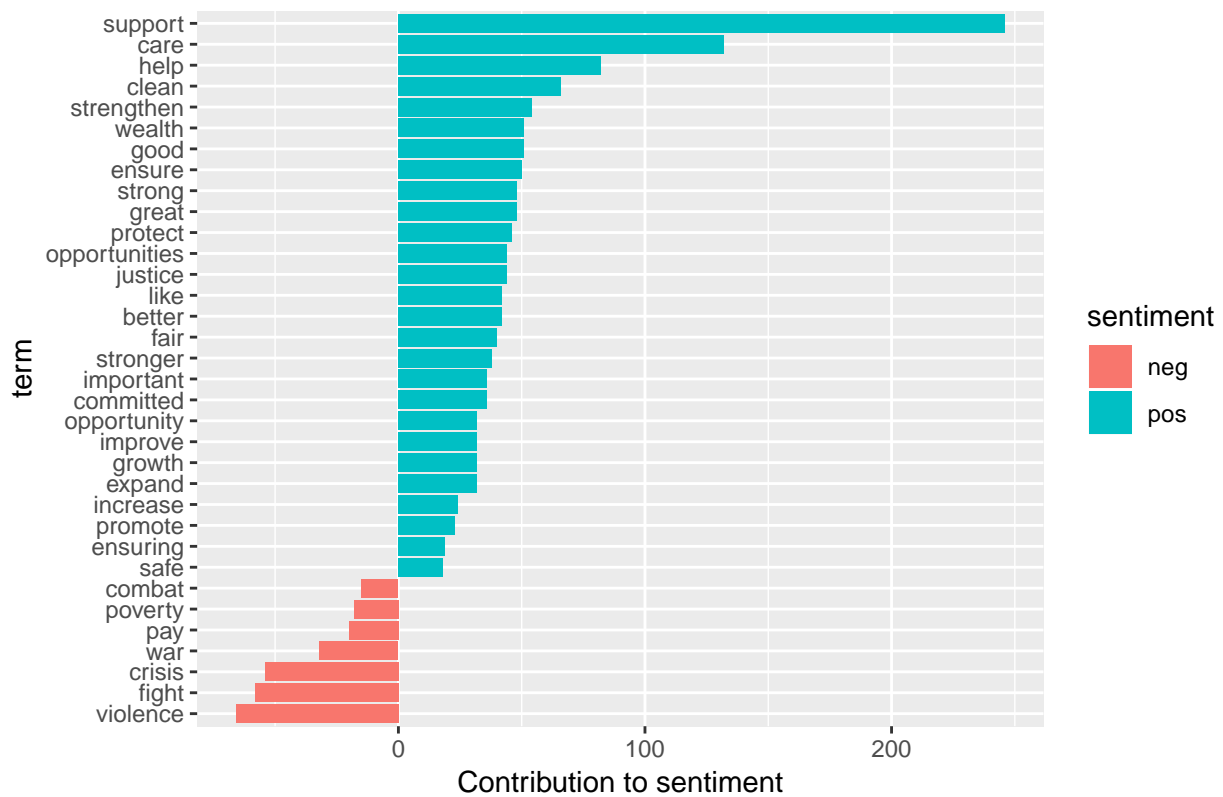
## Current value sentiment Dem. Party: 1270

## Current value sentiment Rep. Party: 891

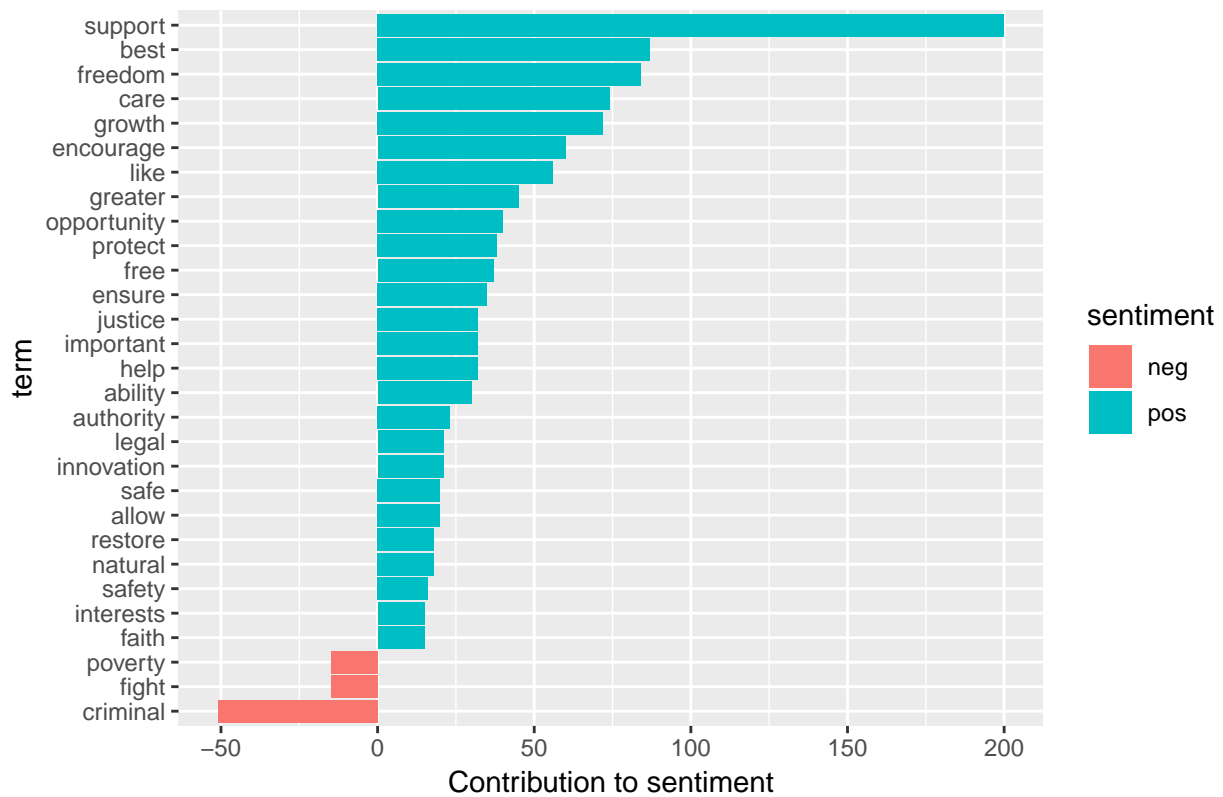
If we see the added sentiment by type of word with more of 15 counts we see a switch for the results of this same graphs in the previous part. We found that in general Democrats use negative words that have a much higher score of negativity than Republicans, for example Violence that is the more aggregated negative word that is used by the Democrats, is not used by Republicans. In other hand, Republicans tend to use positive words what have a higher count and weight as support, freedom and best.

In the case of Democrats we see a great amount of usage of words as strengthen, protect, opportunities and ensure, that is consistent with what we think about democrats, in the sense that as they are more pro big government in a certain way they are looking for more protection of states of nature and look through government intervention.

Dem. Plataform Contribution to Sentiment (More than 15 counts)



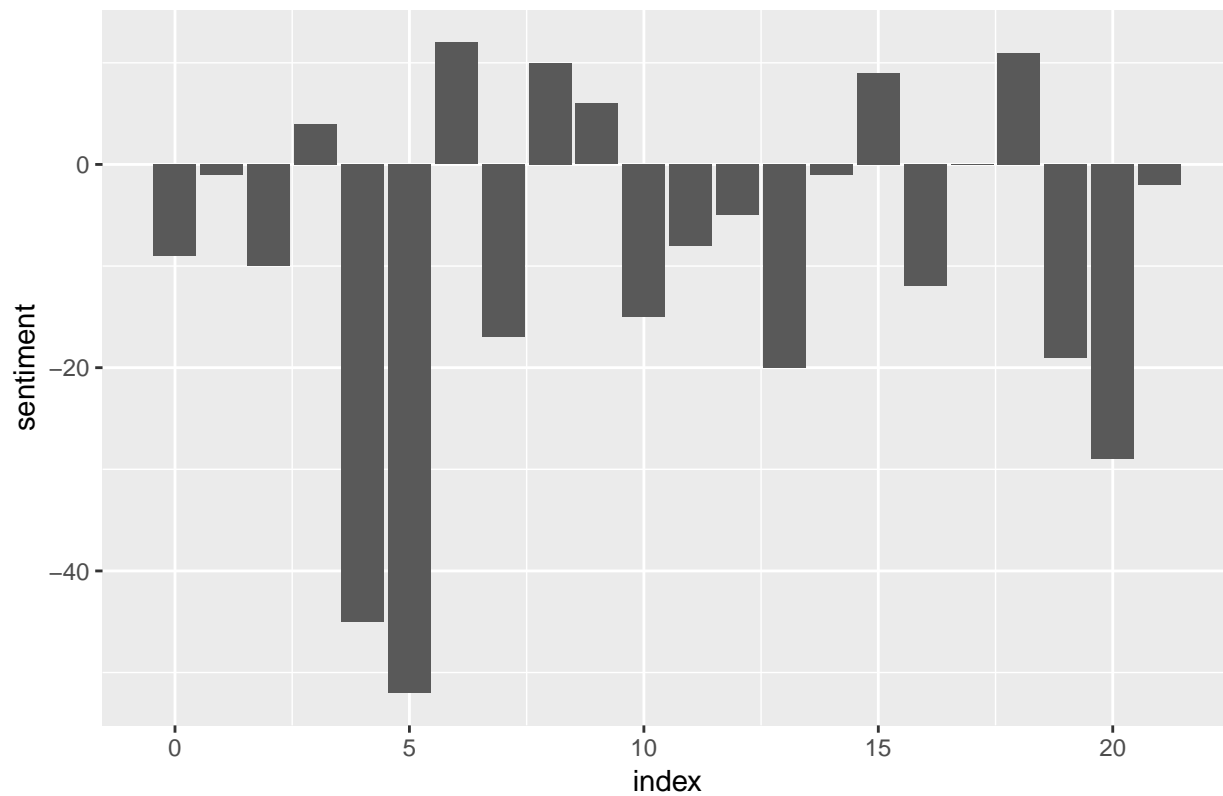
Rep. Plataform Contribution to Sentiment (More than 15 counts)

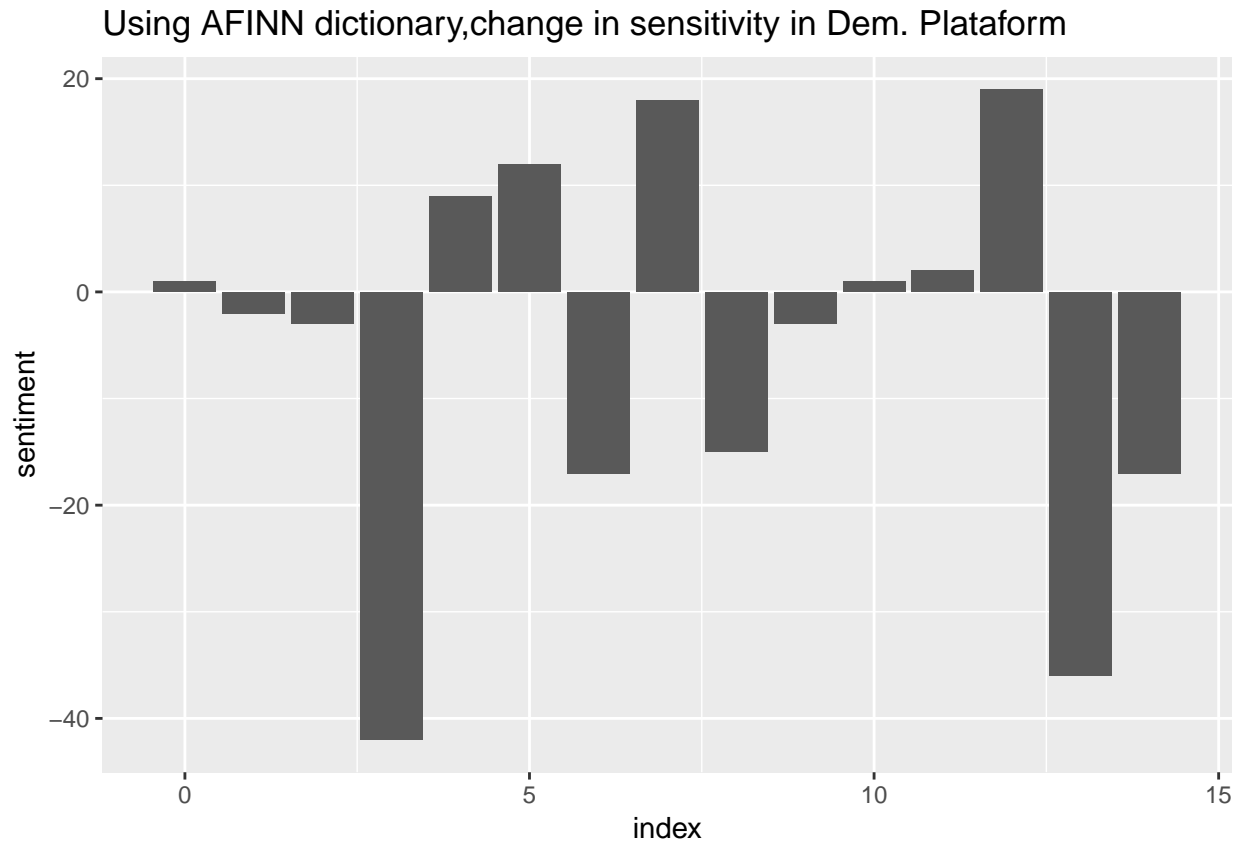




Finally, when we analyze how the the sentiment change through the document. We found that are parts of the plataforms where the overall sentiment is negative. Nevertheless, we keep finding as in the previous part that Republicans can be much more negative than Democrats in certain parts of the documents.

### Using AFINN dictionary,change in sensitivity in Rep. Plataform



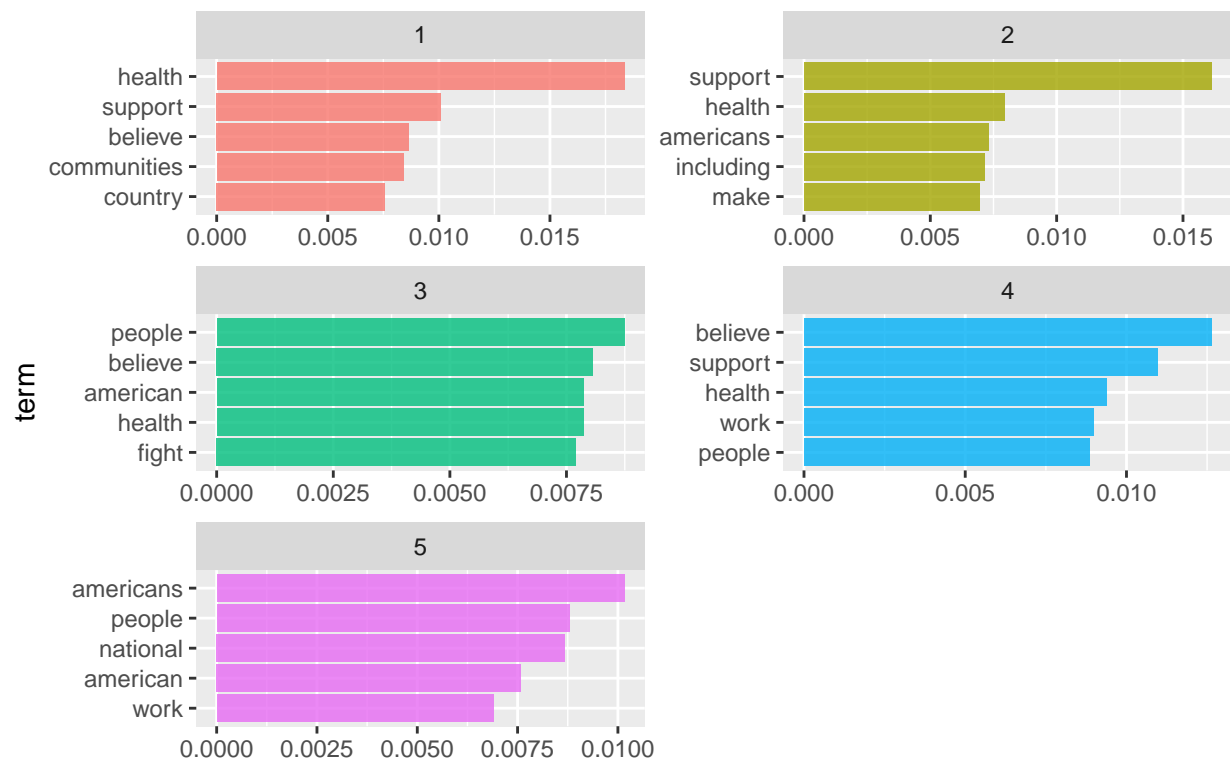


From this part we can see that in part the evidence presented is kind of contradictory when we look the overall use of words positive and negative. This contradiction is found in the result that republicans tend to use less negative words with a higher weight than Democrats. Nevertheless, this contradiction is only in this point and the other elements presented show that Democrats are more optimistic about the future.

## Part 6 and 7

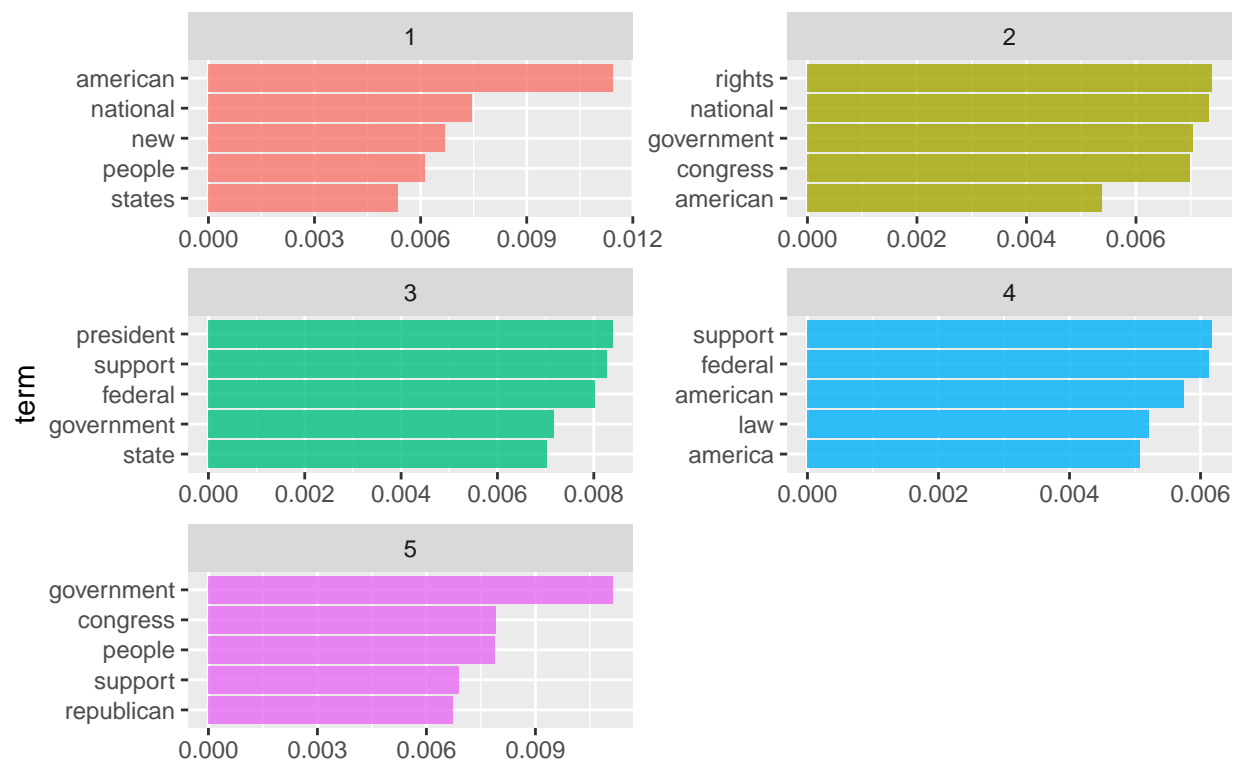
We can see that in general both platforms have different views. In the case of democrats we found that they have focus in concepts as Health, Support, Work and People. This made us think that the platform of the Democrat party is more close to protect people and communities. In the case of the Republican Party, we can see that is more about the relationship between government and people. Which now this made sense in the case of them. Because, a great part of the Republican agenda is “protect” people from the Government.

## Topics Dem. Plataform (K=5)



beta

## Topics Rep. Plataform (K=5)



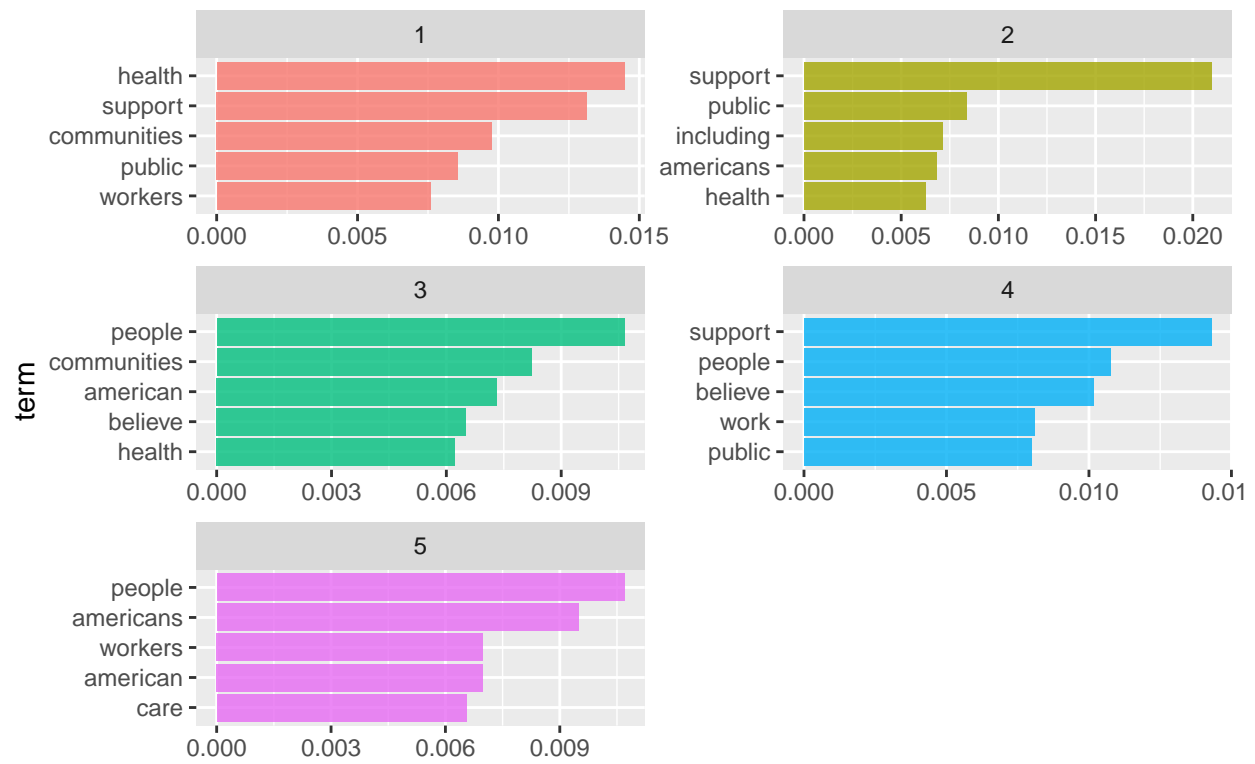
beta

## Part 8

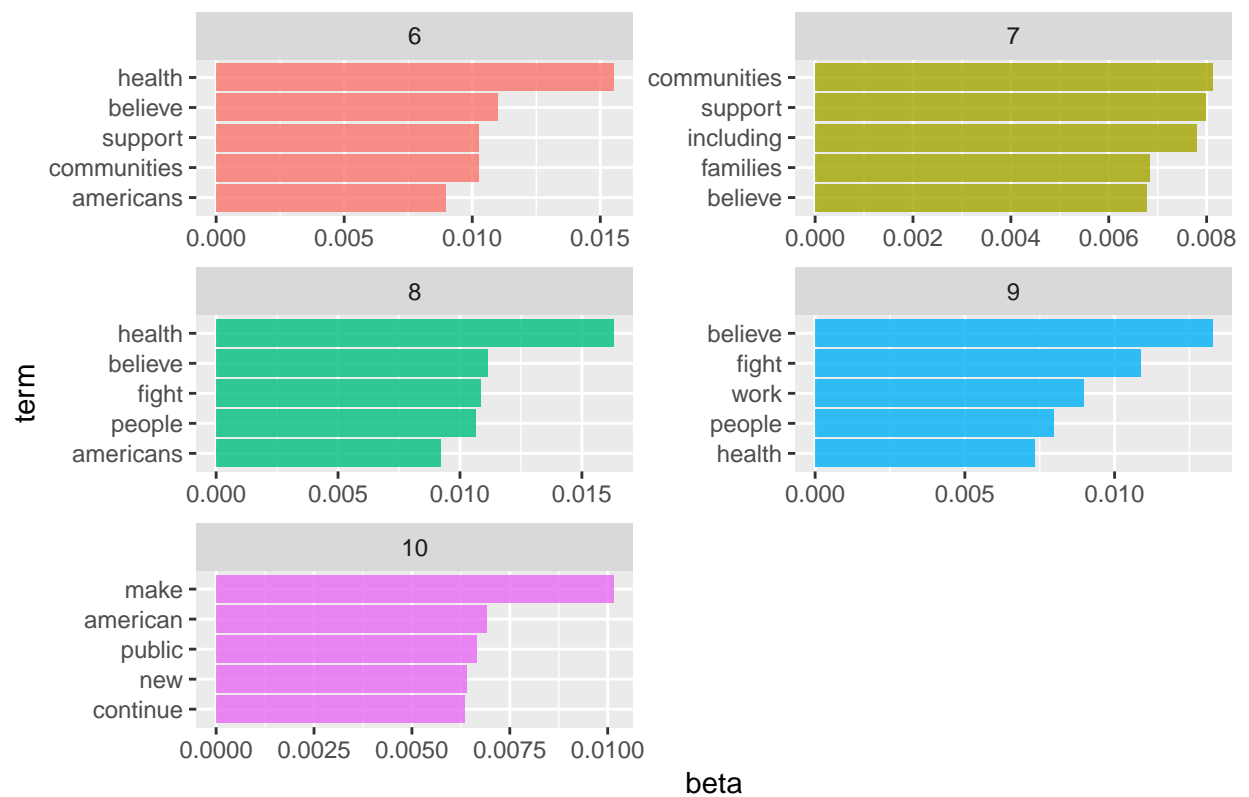
In general, if we look  $K=10$  for both partys we notice that the terms that appears are mostly lineal combinations of the ones that that thet appear when we use  $K=5$ , this would imply that is redundant to use  $K=10$ , because we can have a more parsimonius explication of what we had found. For  $K=25$ , is exactly the same situation.

K=10

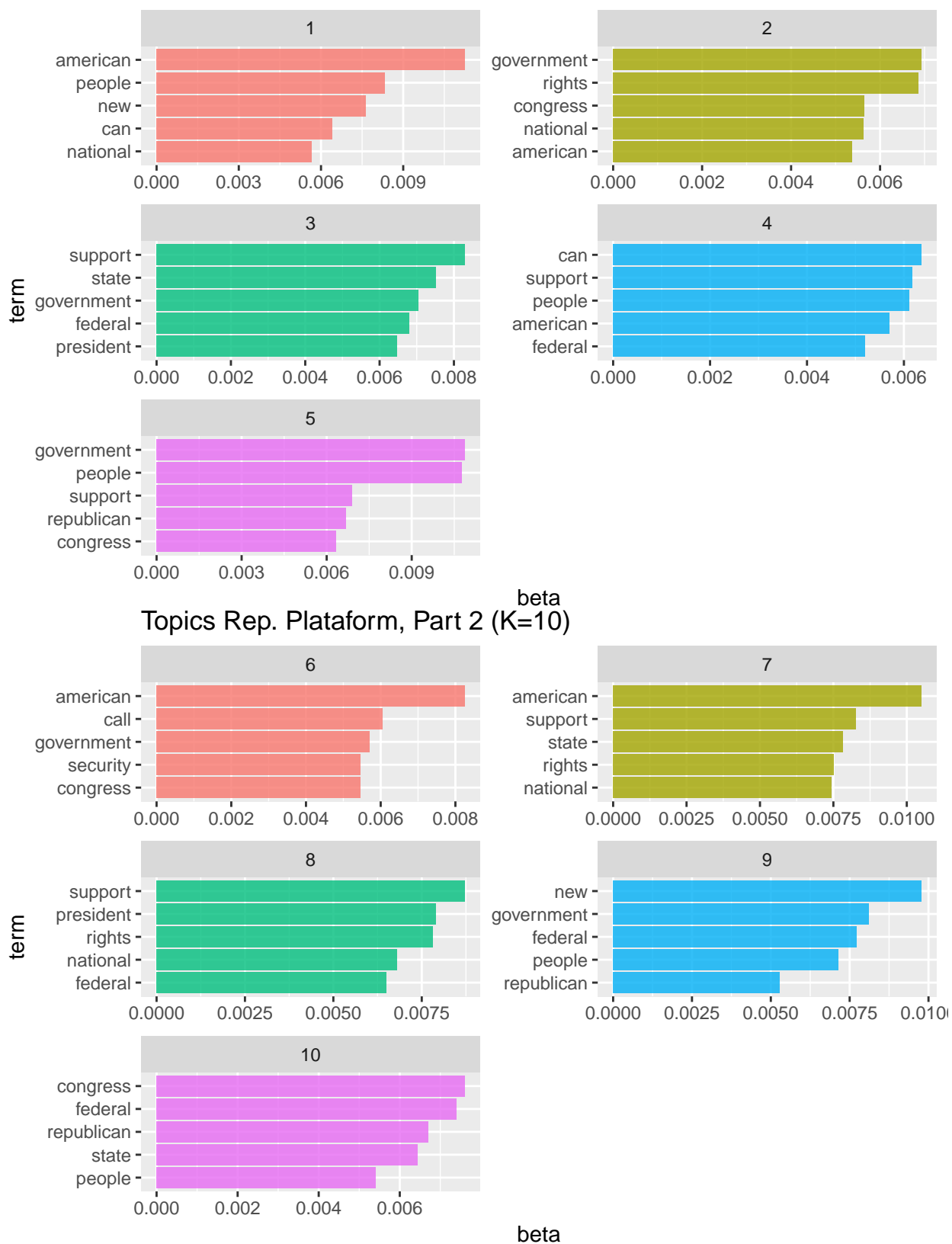
### Topics Dem. Plataform. Part 1 (K=10)



### Topics Dem. Plataform. Part 2 (K=10)

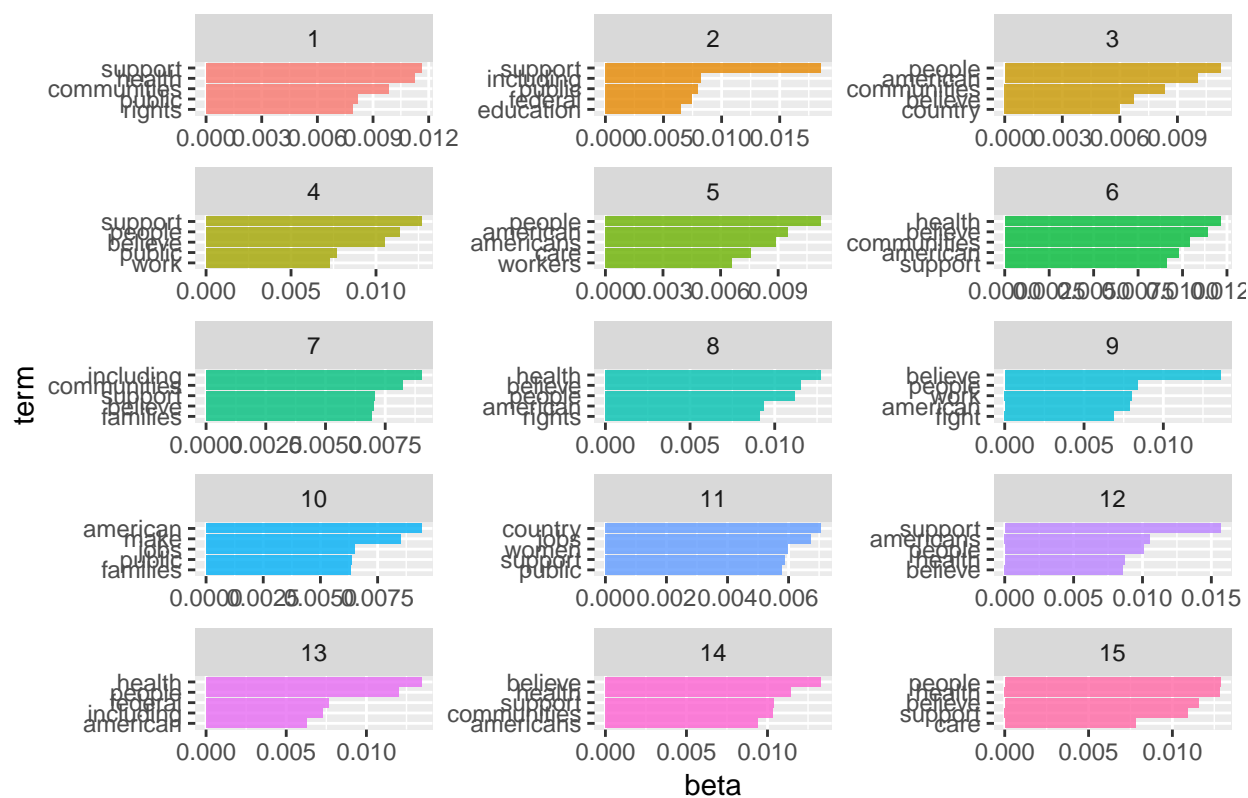


Topics Rep. Plataform, Part 2 (K=10)

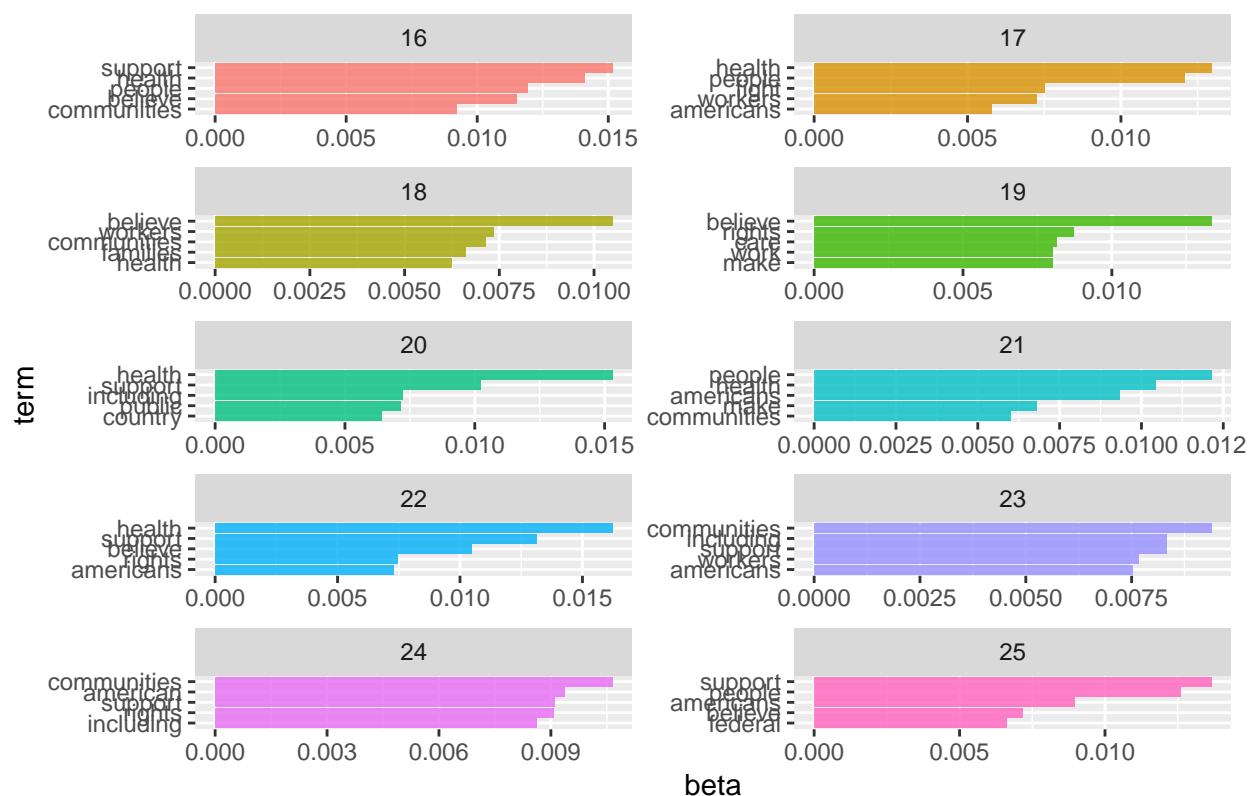


K=25

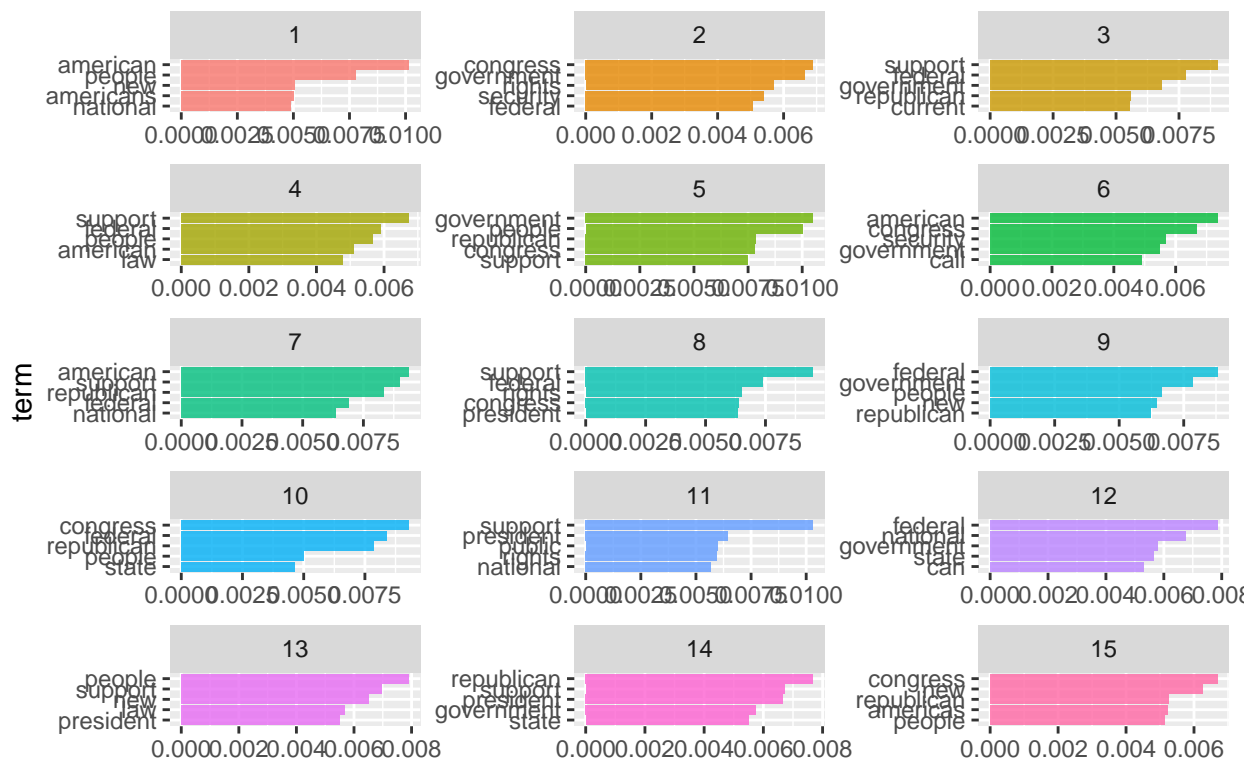
### Topics Dem. Plataform (K=25), part 1



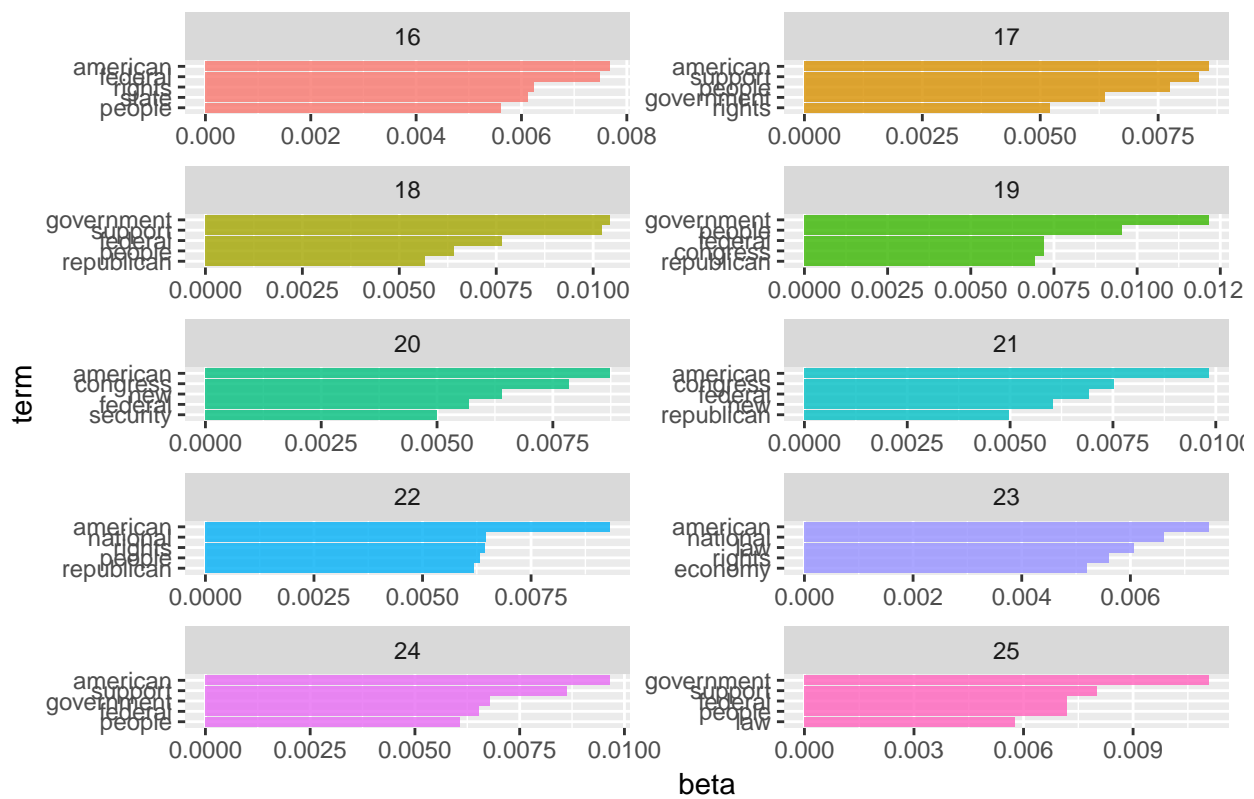
### Topics Dem. Plataform (K=25), part 2



## Topics Rep. Platform (K=25), Part 1



## Topics Rep. Platform. Part 2 (K=25)

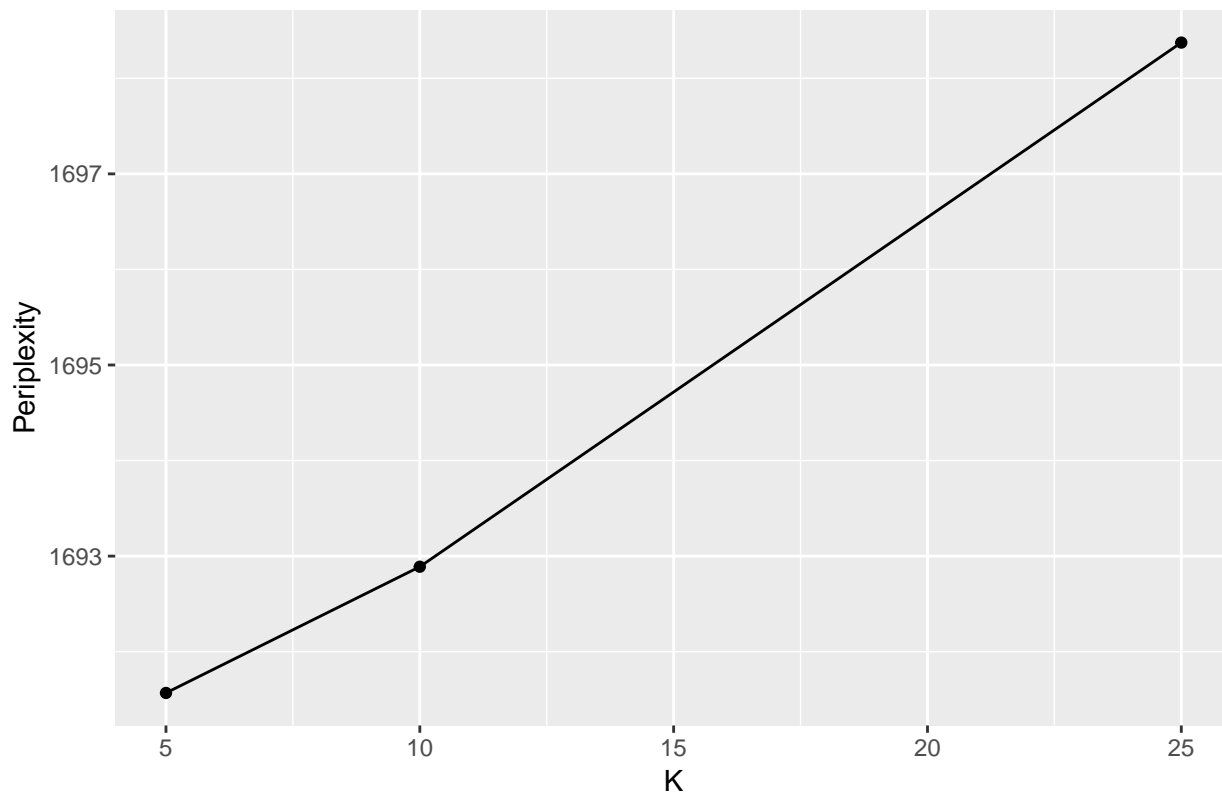


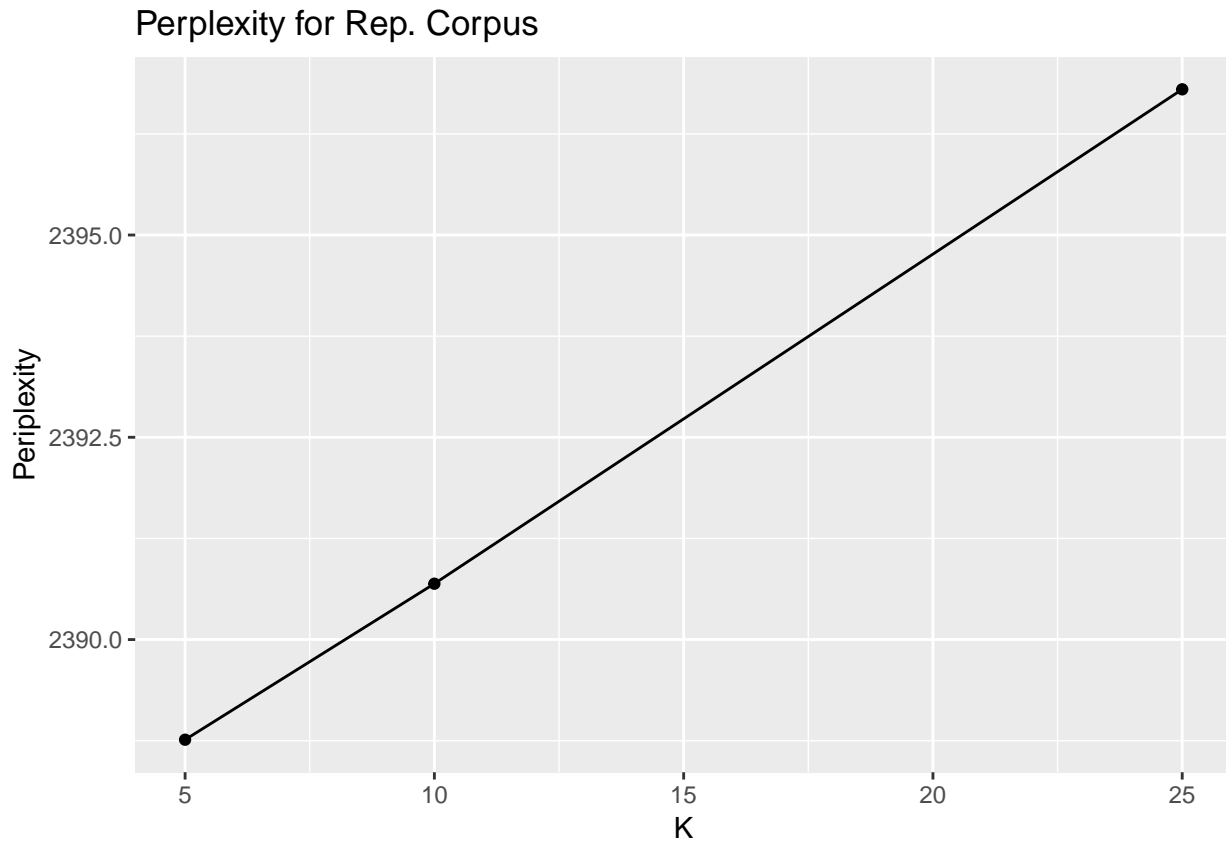


## Part 9

We can see that for any of both platforms, the optimal K is 5. This is related to the issue that as we see in the previous parts that the concepts in each focal point that we had found are redundant as we increase k.

### Perplexity for Dem. Corpus





## Part 10

As we see in the previous part, we can see that  $k=10$  it's not the optimal number of  $K$  because the perplexity index minimize at  $k=5$ . When we made the graph of bars of the different focal points that we found is that at  $K=10$ , they tend to intersect to the concepts what we find in  $K=5$ , then from an efficiency standpoint using 10  $k$ 's doesn't make sense.

## Part 11

In this case I would support the democrats, because the overall tone of that the government must protect the peoples and communities make sense with my own valoric framework. This is notice in the types of words that are selected in the different analysis that we made through this documents. For example, care, health, community between others are in connection with my own personal belief. In other hand, the general tone of the Republican Corpus of protecting people from government is not correct in the current context. This is notorious when we see the sentiment analysis, at least from my perspective one explanation of the more negativity of the democrats in the usage of high negative words is current context of how the public affairs are today, the diagnostic is grim and as is grim requires the usage of this type of words.

In other comments, I notice lately the word believe in the Democrat Corpus was used in a kind similar as Must in Republican, so I should have removed it from the analysis.