

KLonData

Text analysis of republican presidential debates

Here we analyze the republican presidential debate transcripts, focusing on three candidates in particular: Trump, Cruz and Rubio.

First let's answer the following question:

How often does Trump get mentioned by name by the other two candidates during a debate, versus how often do the other two candidates mention each other's name ?

The transcripts were downloaded from: <http://www.presidency.ucsb.edu/debates.php> (<http://www.presidency.ucsb.edu/debates.php>), from August 6th, 2015 to March 21st, 2016.

Here's the code:

```
1 library(tm)
2 library(SnowballC)
3 library(wordcloud)
4 library(RColorBrewer)
5
6 ref_matrix = function (date){
7
8   #read character data
9   text = scan(paste0('debate_', date, '.txt'), what='x', quote=NULL)
10 }
```

```

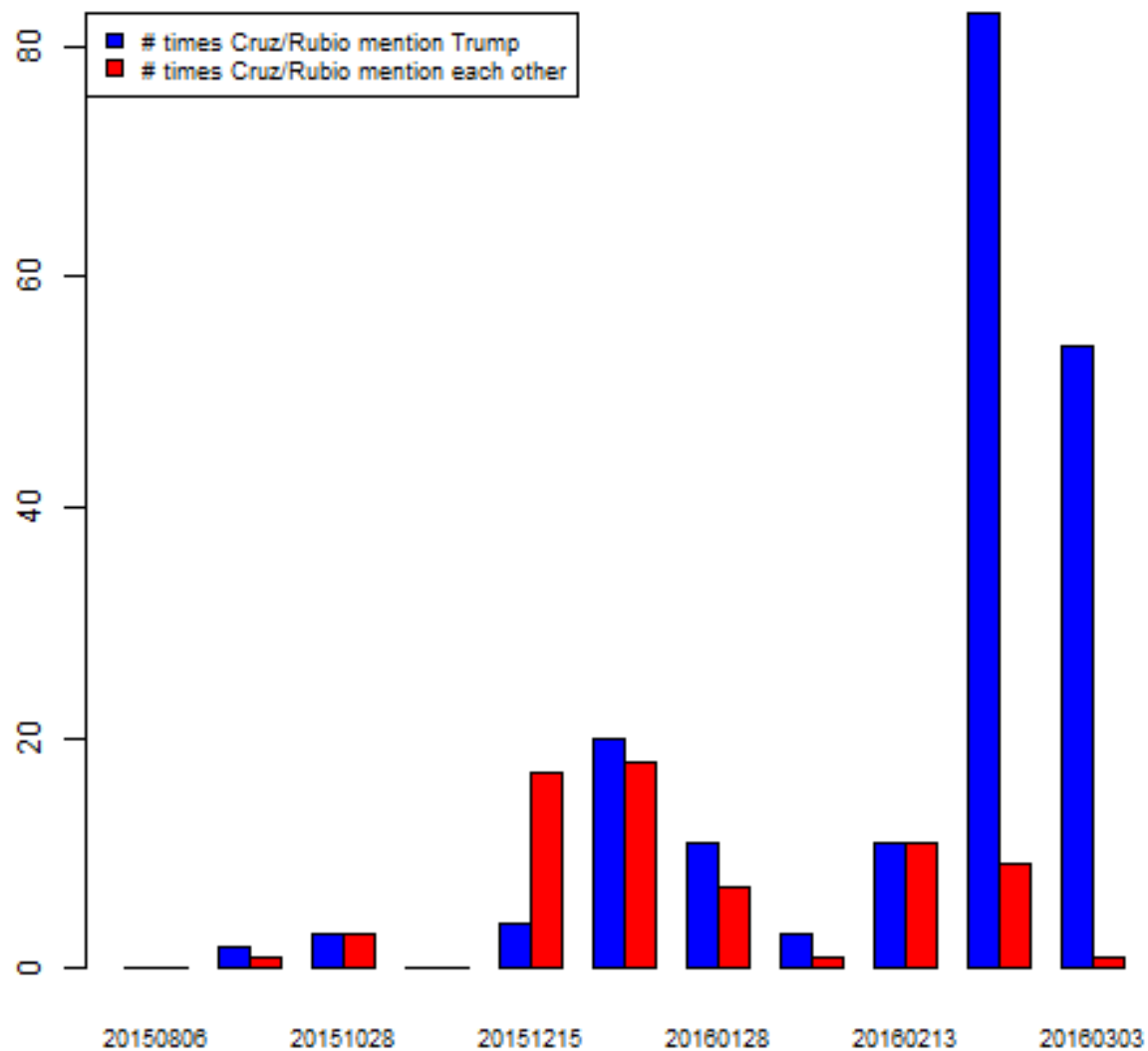
11 speakers =
12   c( CRUZ   = '',
13      TRUMP  = '',
14      RUBIO  = ''
15   )
16
17 #assign text to the right speaker
18 for(word in text){
19   #if word ends with :
20   if(substr(word, nchar(word), nchar(word))==':'){
21     #if word corresponds to one of the speakers of interest
22     if(word %in% paste0(names(speakers), ':')){
23       #set current speaker
24       currentSpeaker = substr(word, 1, nchar(word)-1)
25     }
26     else{
27       #if the current speaker is not one of the speakers of interest, set it to NA
28       currentSpeaker = NA
29     }
30   }
31   else if(!is.na(currentSpeaker)){
32     #if the current speaker is of interest, save what he is saying
33     speakers[currentSpeaker] = paste(speakers[currentSpeaker], word)
34   }
35 }
36
37 #preprocess text
38 prez = Corpus(VectorSource(speakers))
39 prez = tm_map(prez, tolower)
40 prez = tm_map(prez, removeWords, stopwords('english'))
41 #remove additional unwanted words
42 prez = Corpus(VectorSource(speakers))
43 prez = tm_map(prez, tolower)
44 prez = tm_map(prez, removeWords, stopwords('english'))
45 #remove additional unwanted words
46 prez = tm_map(prez, removeWords, c('will', 'going', 'applause', 'get', 'say', 'want', 'let', 'c',
47                                     'got', 'one', 'two', 'also', 'ever', 'even', 'need', 'every',
48
49 prez = tm_map(prez, removePunctuation, preserve_intra_word_dashes = FALSE)
50 prez = tm_map(prez, stemDocument)
51 prez = tm_map(prez, stripWhitespace)

```

[illegible]

```
91 ref_list = lapply(dates, ref_matrix)
92
93 names(ref_list) = dates
94
95 trump = sapply(ref_list, function(df) sum(df[rownames(df) != 'TRUMP', 'TRUMP']))
96 cruz_rubio = sapply(ref_list, function(df) sum(df[rownames(df) == 'CRUZ', 'RUBIO'], df[rownames(d
97
98 m = t(as.matrix(data.frame(trump, cruz_rubio)))
99 barplot(m, main='Number of times Cruz/Rubio mention Trump vs each other', beside=TRUE, col=c('blu
100     legend=c('# times Cruz/Rubio mention Trump', '# times Cruz/Rubio mention each other'),
101     args.legend=list(x='topleft', cex=0.75), cex.names=0.75)
```

Number of times Cruz/Rubio mention Trump vs each other



(<https://klondata.wordpress.com/2016/03/06/7/unnamed-chunk-1-1/>)

We can see that at the beginning of the race, the candidates really didn't refer to each other much at all.

Things change around the debate of December 15th, 2015, where Cruz and Rubio refer to each other significantly more. Then in the next debate and every other debate afterwards, Cruz and Rubio collectively refer to Trump more often than they refer to each other.

In particular, in the last two debates, of February 25th, 2016 and March 3rd, 2016, they mention each other 10 times in total, where as they mention Trump 137 times !

Let's now turn our attention to the words themselves.

We're gonna change the code a little bit in order to collect all the transcripts in a single character string:

```

1 | dates = c(20150806, 20150916, 20151028, 20151110, 20151215,
2 |          20160114, 20160128, 20160206, 20160213, 20160225,
3 |          20160303)
4 | read_transcript = function(date){scan(paste0('debate_', date, '.txt'), what='x', quote=NULL)}
5 |
6 | #read and collate all transcripts
7 | text = unlist(sapply(dates, read_transcript))

```

After reusing the same bit of code as in the beginning to preprocess the text, we can answer a few interesting questions:

For each candidate, what are the top 50 most frequent words ?

```

1 | get indexes of 50 most frequent words
2 | indexes = apply(dtm, 1, function(v) head(order(v, decreasing=TRUE), 50))
3 |
4 | # find the 50 most frequent words
5 | freq_words = apply(indexes, 2, function(v) colnames(dtm)[v])
6 | freq_words

```

```

1 | ##          CRUZ          TRUMP          RUBIO
2 | ##      [1,] 'donald'      'people'      'people'
3 | ##      [2,] 'president'   'country' 'president'

```

4	##	[3,]	'people'	'just'	'country'
5	##	[4,]	'now'	'think'	'now'
6	##	[5,]	'tax'	'said'	'america'
7	##	[6,]	'obama'	'now'	'states'
8	##	[7,]	'country'	'tell'	'united'
9	##	[8,]	'said'	'right'	'just'
10	##	[9,]	'question'	'great'	'world'
11	##	[10,]	'right'	'like'	'first'
12	##	[11,]	'back'	'way'	'issue'
13	##	[12,]	'just'	'look'	'think'
14	##	[13,]	'america'	'back'	'like'
15	##	[14,]	'washington'	'take'	'american'
16	##	[15,]	'years'	'lot'	'way'
17	##	[16,]	'court'	'much'	'important'
18	##	[17,]	'american'	'come'	'make'
19	##	[18,]	'clinton'	'make'	'years'
20	##	[19,]	'hillary'	'thing'	'immigration'
21	##	[20,]	'tell'	'never'	'money'
22	##	[21,]	'fight'	'years'	'tax'
23	##	[22,]	'law'	'china'	'isis'
24	##	[23,]	'amnesty'	'world'	'see'
25	##	[24,]	'isis'	'first'	'said'
26	##	[25,]	'day'	'talking'	'someone'
27	##	[26,]	'first'	'good'	'time'
28	##	[27,]	'look'	'money'	'believe'
29	##	[28,]	'think'	'win'	'government'
30	##	[29,]	'like'	'everybody'	'never'
31	##	[30,]	'state'	'trade'	'back'
32	##	[31,]	'crosstalk'	'care'	'barack'
33	##	[32,]	'new'	'something'	'military'
34	##	[33,]	'barack'	'deal'	'things'
35	##	[34,]	'flat'	'wall'	'right'
36	##	[35,]	'keep'	'really'	'americans'
37	##	[36,]	'percent'	'problem'	'economy'
38	##	[37,]	'plan'	'believe'	'made'
39	##	[38,]	'women'	'billion'	'obama'
40	##	[39,]	'bill'	'saying'	'today'
41	##	[40,]	'radical'	'crosstalk'	'fact'
42	##	[41,]	'stage'	'big'	'able'
43	##	[42,]	'government'	'president'	'hillary'

```

44  ##      [43,] 'business'      'time'      'question'
45  ##      [44,] 'john'         'wrong'     'donald'
46  ##      [45,] 'marco'        'done'      'clinton'
47  ##      [46,] 'everyone'     'excuse'    'plan'
48  ##      [47,] 'islamic'      'jobs'      'point'
49  ##      [48,] 'millions'     'laughter'  'support'
50  ##      [49,] 'immigration'  'far'       'better'
51  ##      [50,] 'men'          'jeb'       'place'

```

It would actually be interesting to see for each candidate, the words in his top-50 that are unique to him, i.e. that are not in the other candidates' top-50

```

1  find the number of times each word appears in the matrix
2  word_count = apply(freq_words, c(1,2), function(x) sum(x == freq_words))
3
4  #keep those words that appear only once
5  unique_words = word_count == 1
6
7  l = lapply(colnames(freq_words), function(name) freq_words[unique_words[,name], name])
8  names(l) = colnames(freq_words)
9  l

```

```

1  ## $CRUZ
2  ##      [1] 'washington' 'court'      'fight'      'law'        'amnesty'
3  ##      [6] 'day'        'state'      'new'        'flat'       'keep'
4  ##     [11] 'percent'    'women'      'bill'       'radical'    'stage'
5  ##     [16] 'business'   'john'       'marco'      'everyone'   'islamic'
6  ##     [21] 'millions'   'men'
7  ##
8  ## $TRUMP
9  ##      [1] 'great'      'take'      'lot'        'much'       'come'
10 ##      [6] 'thing'      'china'     'talking'    'good'       'win'
11 ##     [11] 'everybody'  'trade'     'care'       'something'  'deal'
12 ##     [16] 'wall'       'really'    'problem'    'billion'    'saying'
13 ##     [21] 'big'        'wrong'     'done'       'excuse'     'jobs'
14 ##     [26] 'laughter'   'far'       'jeb'
15 ##
16 ## $RUBIO

```



```
17 ## [1] 'states'      'united'      'issue'       'important'   'see'
18 ## [6] 'someone'    'military'    'things'      'americans'  'economy'
19 ## [11] 'made'       'today'      'fact'        'able'        'point'
20 ## [16] 'support'    'better'     'place'
```

Let's go ahead and make word clouds out of that:

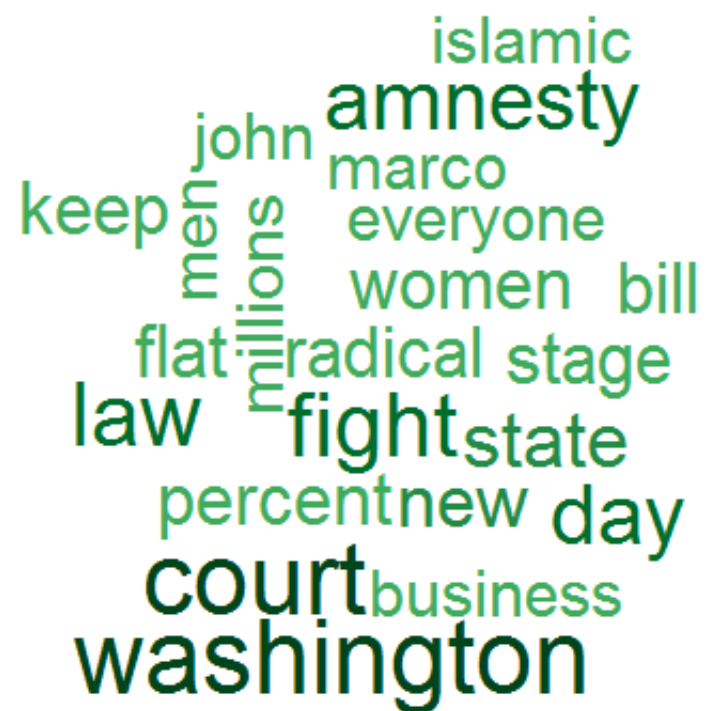
Trump:

```
1 freq = dtm['TRUMP', 1$TRUMP]
2 wordcloud(1$TRUMP, as.vector(freq), colors=brewer.pal(9, 'Blues'))
```



Cruz:

```
1 | freq = dtm['CRUZ', 1$CRUZ]  
2 | wordcloud(1$CRUZ, as.vector(freq), colors=brewer.pal(9, 'Greens'))
```



Rubio:

```
1 | freq = dtm['RUBIO', 1$RUBIO]  
2 | wordcloud(1$RUBIO, as.vector(freq), colors=brewer.pal(9, 'Oranges'))
```



For each candidate, what is the average word length ?

```
1 | #number of words by speaker  
2 | nb_words = apply(dtm, 1, sum)
```

```

3 | #word lengths
4 | word_lengths = sapply(colnames(dtm), nchar)
5 | #transform word count into total character count matrix
6 | character_counts = t(apply(dtm, 1, function(v) v * word_lengths))
7 | #total character count by speaker
8 | total_character_counts = apply(character_counts, 1, function(v) sum(v))
9 | #divide total character count by number of words
10 | round(total_character_counts / nb_words, digits=1)

1 | ## CRUZ TRUMP RUBIO
2 | ## 6.3 5.9 6.3

```

Finally, for each candidate, how diversified is their vocabulary ?

To quantify this, we are gonna count the number of unique words per 1000 words:

```

1 | apply(dtm, 1, function(v) round(sum(v != 0)/sum(v)*1000, digits=1))

1 | ## CRUZ TRUMP RUBIO
2 | ## 248.8 176.3 218.0

```

We see that Trump uses fewer and shorter words than his opponents.

Posted in [R](#) and tagged [R](#) on [March 6, 2016](#) by [Karim L.](#) [Leave a comment](#)

[BLOG AT WORDPRESS.COM. THE SUITS THEME.](#)