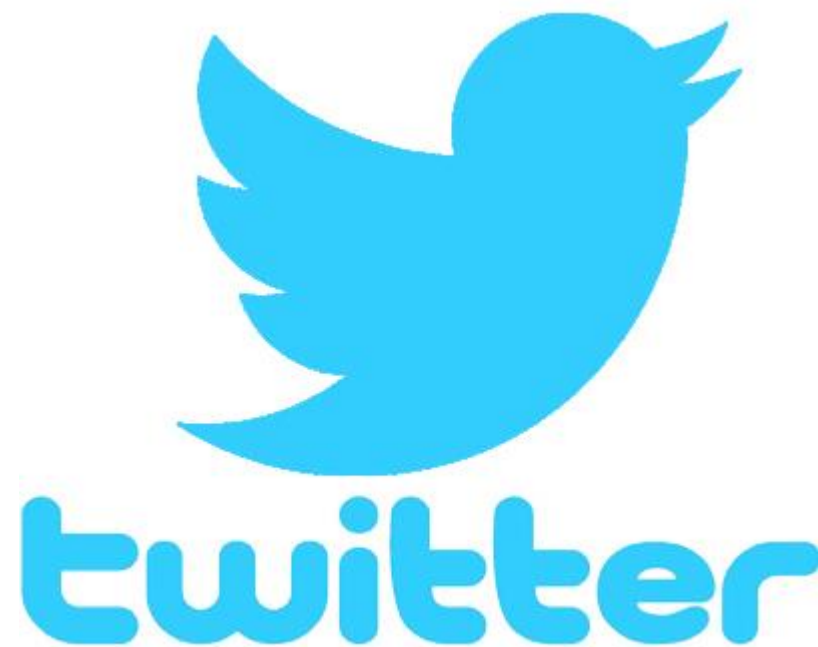


PROVIDER
NEW WORK
TEXTUAL
ENTITY
MANAGEMENT
RETRIEVAL
CUSTOMER
TEXT
LARGE
APPLICATION
INFORMATION
DATASETS
ANALYTICS
PROCESS
SOFTWARE
STATISTICAL
LITERATURE
ABSTRACTED
PREDICTIVE
NATURAL
AFFECTIVITY
SYNTACTIC
TOOL MINING
EXTRACTION
BUSINESS
SOCIAL
ANNOTATION
INTELLIGENCE
DOCUMENT
LEARNING
LINGUISTIC
SUITE
MEDIAL
LANGUAGE
PARSING
SPECIFIC
DATA
NEW
SEMANTIC
CATEGORIZATION
SUMMARIZATION
SENTIMENT
VISUALIZATION
CLUSTER
ATTITUDINAL
SCIENCE
COMPUTATIONAL
AFFINITIES
CLUSTER

TEXT MINING

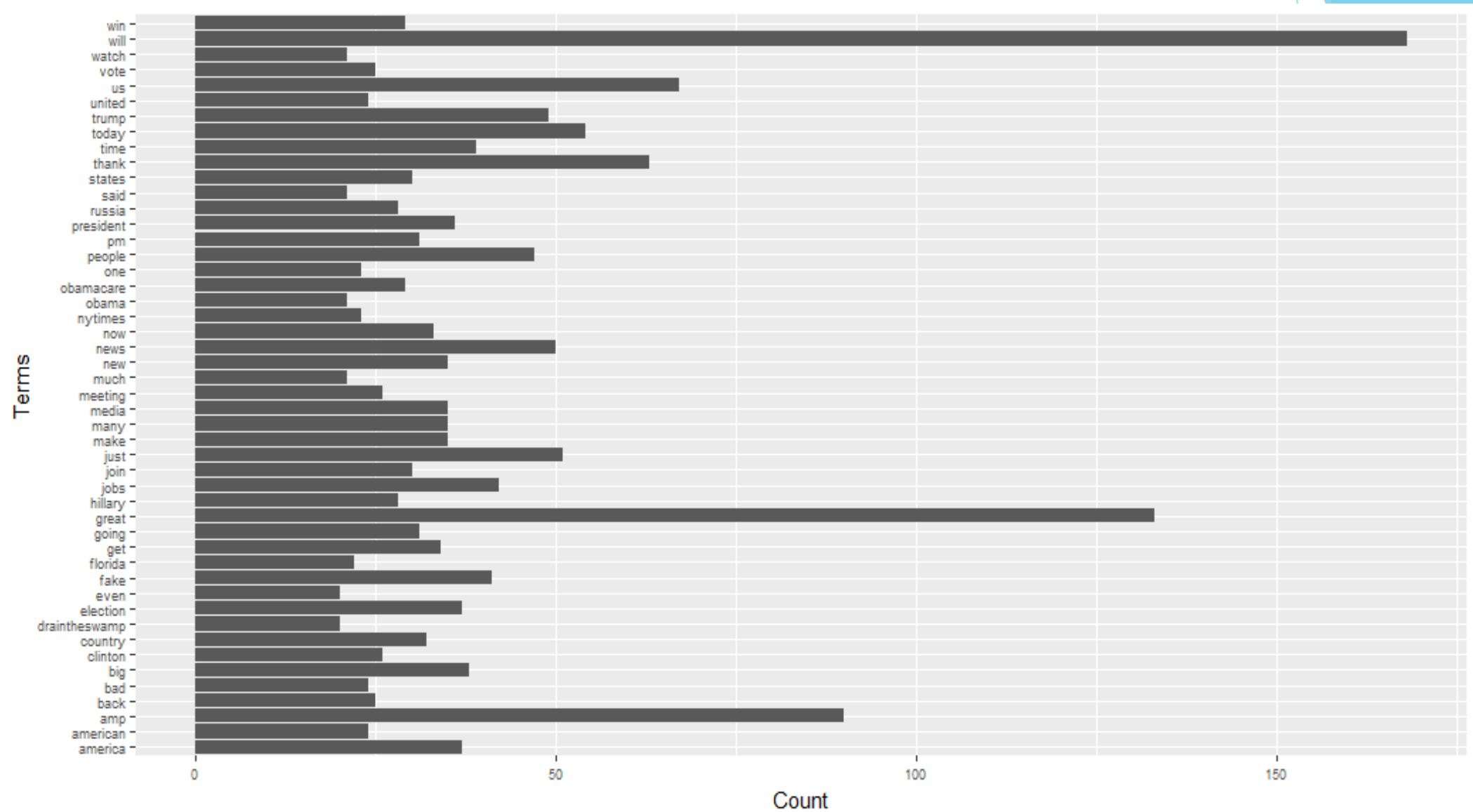


Text Mining in R (*Twitter Analysis*)

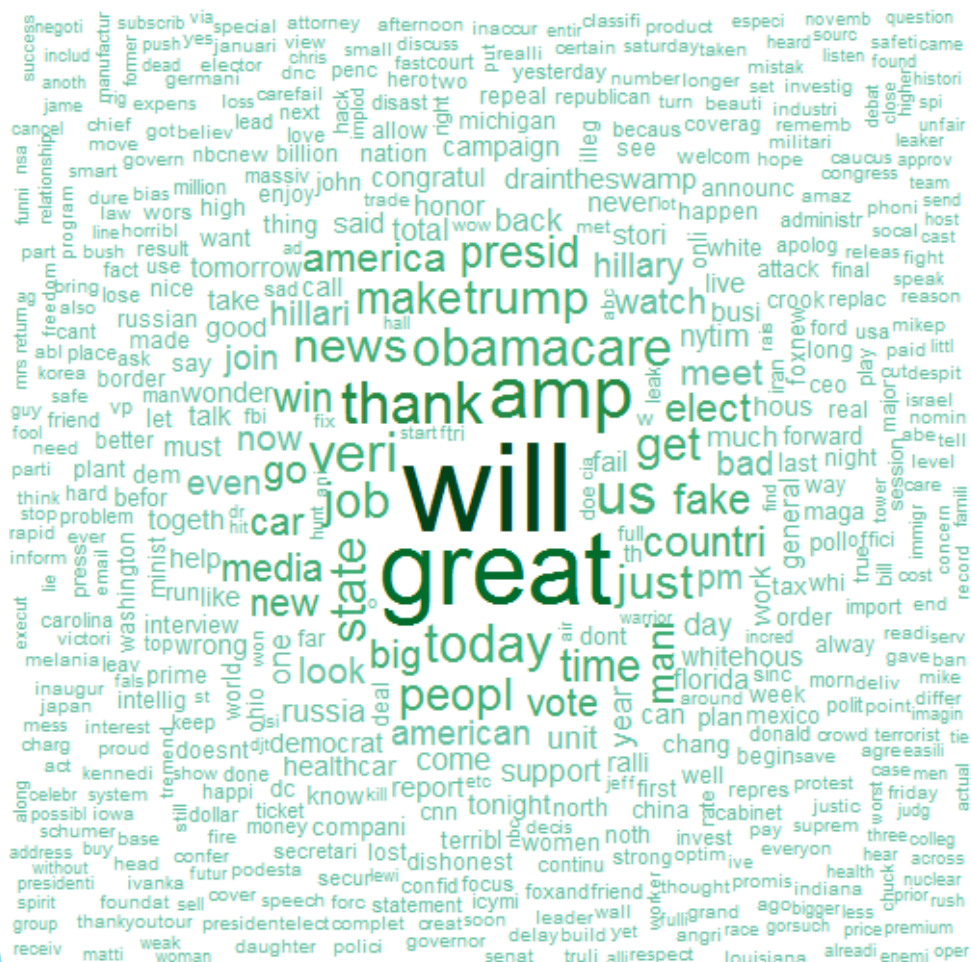
- ▶ Text mining on Donald Trump's Twitter account
- ▶ Freq words: vote, thank, maga, Obamacare, Hillary, nytimes, jobs, great, florida, drain the swamp, fake, election, Russia
- ▶ Wordcount :
 - ▶ Hillary Clinton : 60
 - ▶ Obamacare : 33
 - ▶ Jobs : 34
 - ▶ Russia : 35
 - ▶ Drain the swamp : 20



Frequency of the words



Word Cloud in R



Background in R

```
1 library(twitter)
2 library(ROAuth)
3
4 api_key <- "XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX"
5 api_secret <- "XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX"
6 access_token <- "XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX"
7 access_token_secret <- "XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX"
8
9 setup_twitter_oauth(api_key, api_secret, access_token, access_token_secret)
10
11 tweets <- userTimeline("realDonaldTrump", n = 100)
12 n.tweet <- length(tweets)
13
14 # convert tweets to a data frame
15 tweets.df <- twListToDF(tweets)
16
17 # tweet #190
18 tweets.df[10, c("id", "created", "screenName", "replyToSN", "favoriteCount", "retweetCount", "longitude", "latitude")]
19
20 # print tweet #190 and make text fit for slide width
21 writeLines(strwrap(tweets.df$text[10], 60))
22
23 library(tm)
24 # build a corpus, and specify the source to be character vectors
25 myCorpus <- Corpus(VectorSource(tweets.df$text))
26
27
28 # remove URLs
29 removeURL <- function(x) gsub("http[^\s:]*", "", x)
30 myCorpus <- tm_map(myCorpus, content_transformer(removeURL))
31
32 # remove anything other than English letters or space
33 removeNumPunct <- function(x) gsub("[^\p{L}\p{S}]", "", x)
34 myCorpus <- tm_map(myCorpus, content_transformer(removeNumPunct))
35
36 # convert to lower case
37 myCorpus <- tm_map(myCorpus, content_transformer(tolower))
38
39 # remove stopwords
40 myStopwords <- c(setdiff(stopwords("english"), c("r", "big")), "use", "see", "used", "via", "amp")
41 myCorpus <- tm_map(myCorpus, removeWords, myStopwords)
42 # remove extra whitespace
43 myCorpus <- tm_map(myCorpus, stripWhitespace)
44
45 # keep a copy for stem completion later
46 myCorpusCopy <- myCorpus
47
48
49 myCorpus <- tm_map(myCorpus, stemDocument) # stem words
50 writeLines(strwrap(myCorpus[[190]]$content, 60))
51
52 stemCompletion2 <- function(x, dictionary) {
53   x <- unlist(strsplit(as.character(x), " "))
54   x <- x[x != ""]
55   x <- stemCompletion(x, dictionary=dictionary)
56   x <- paste(x, sep=" ", collapse=" ")
57 }
58
59 92:37 (Top Level) ↕
```