

Text analysis of Trump's tweets confirms he writes only the (angrier) Android half

I don't normally post about politics (I'm not particularly savvy about polling, which is where data science has had the largest impact on politics). But this weekend I saw a hypothesis about Donald Trump's twitter account that simply begged to be investigated with data:



When Trump wishes the Olympic team good luck, he's tweeting from his iPhone. When he's insulting a rival, he's usually tweeting from an Android. Is this an artifact showing which tweets are Trump's own and which are by some handler?

Others have explored Trump's timeline and noticed this tends to hold up- and Trump himself does indeed tweet from a Samsung Galaxy. But how could we examine it quantitatively? I've been writing about text mining and sentiment analysis recently, particularly during my development of the tidytext R package with Julia Silge, and this is a great opportunity to apply it again.

My analysis, shown below, concludes that **the Android and iPhone tweets are clearly from different people**, posting during different times of day and using hashtags, links, and retweets in distinct ways. What's more, we can see that **the Android tweets are angrier and more negative**, while the iPhone tweets tend to be benign announcements and pictures. Overall I'd agree with [@tvaziri](#)'s analysis: this lets us tell the difference between the campaign's tweets (iPhone) and Trump's own (Android).

The dataset

First we'll retrieve the content of Donald Trump's timeline using the `userTimeline` function in the `twitter` package:¹

```
library(dplyr)
library(purrr)
library(twitteR)
```

```
# You'd need to set global options with an authenticated app
setup_twitter_oauth(getOption("twitter_consumer_key"),
                    getOption("twitter_consumer_secret"),
                    getOption("twitter_access_token"),
                    getOption("twitter_access_token_secret"))

# We can request only 3200 tweets at a time; it will return fewer
# depending on the API
trump_tweets <- userTimeline("realDonaldTrump", n = 3200)
trump_tweets_df <- tbl_df(map_df(trump_tweets, as.data.frame))
```

```
# if you want to follow along without setting up Twitter authentication,
# just use my dataset:
load(url("http://varianceexplained.org/files/trump_tweets_df.rda"))
```

We clean this data a bit, extracting the source application. (We're looking only at the iPhone and Android tweets- a much smaller number are from the web client or iPad).

```
library(tidyr)

tweets <- trump_tweets_df %>%
  select(id, statusSource, text, created) %>%
  extract(statusSource, "source", "Twitter for (.*)<" ) %>%
  filter(source %in% c("iPhone", "Android"))
```

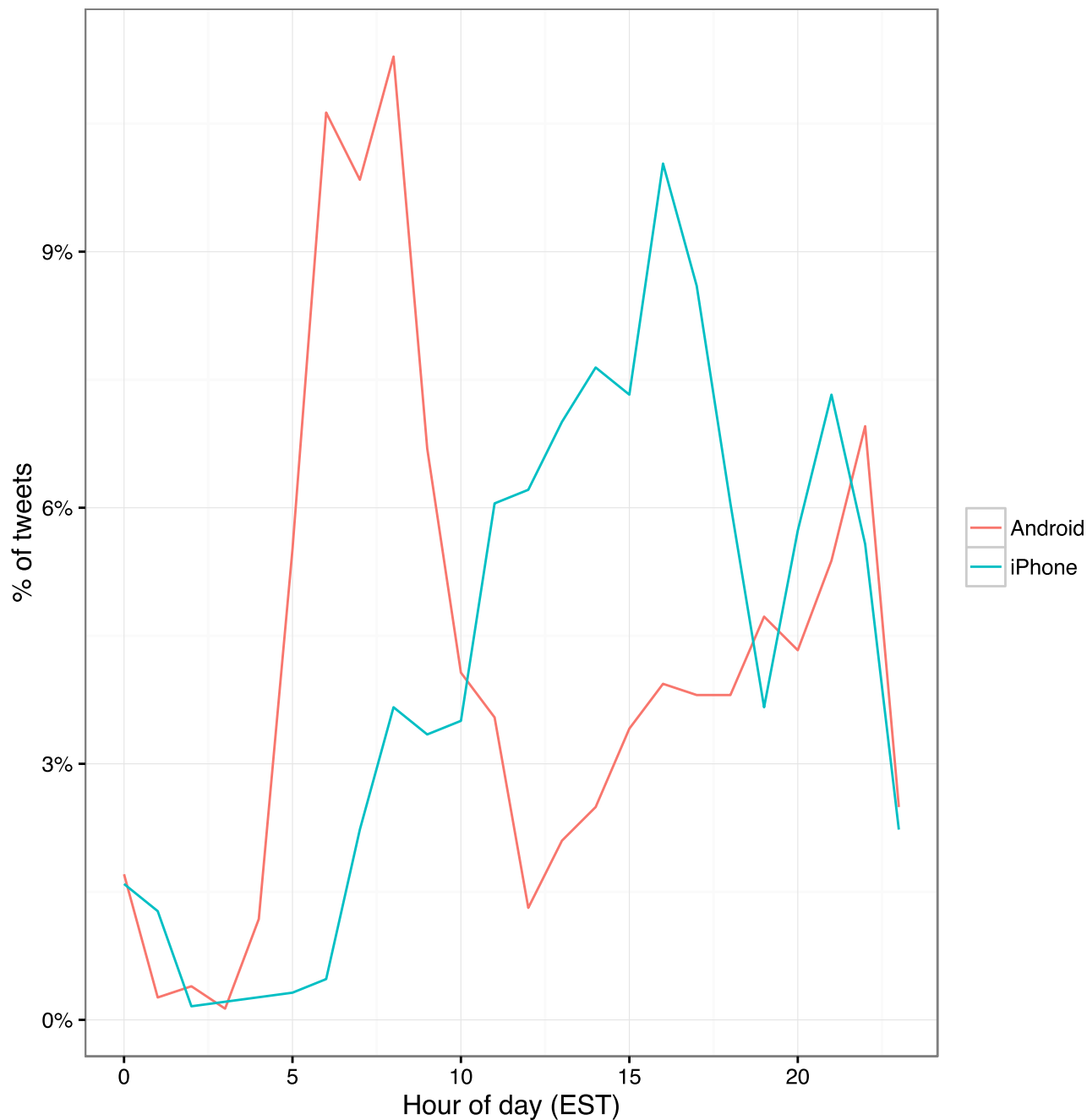
Overall, this includes 628 tweets from iPhone, and 762 tweets from Android.

One consideration is what time of day the tweets occur, which we'd expect to be a "signature" of their user. Here we can certainly spot a difference:

```
library(lubridate)
library(scales)

tweets %>%
  count(source, hour = hour(with_tz(created, "EST"))) %>%
  mutate(percent = n / sum(n)) %>%
  ggplot(aes(hour, percent, color = source)) +
  geom_line() +
```

```
scale_y_continuous(labels = percent_format()) +
labs(x = "Hour of day (EST)",
     y = "% of tweets",
     color = "")
```



Trump on the Android does a lot more tweeting in the morning, while the campaign posts from the iPhone more in the afternoon and early evening.

Another place we can spot a difference is in Trump's anachronistic behavior of "manually retweeting" people by copy-pasting their tweets, then surrounding them with quotation marks:



Donald J. Trump realDonaldTrump

Follow

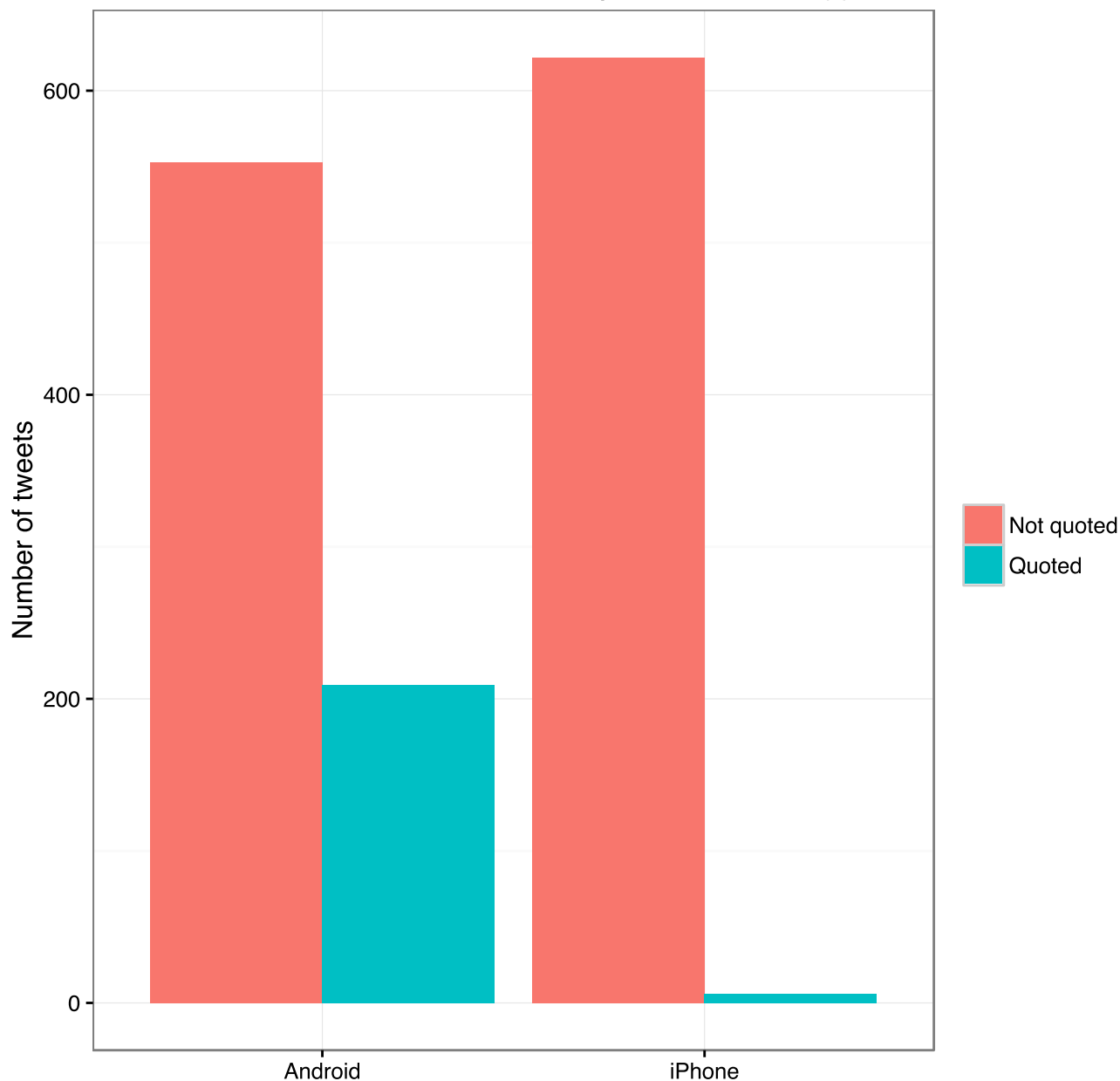
"@trumplican2016: @realDonaldTrump @DavidWohl stay the course mr trump your message is resonating with the PEOPLE"

9:00 PM - 27 Jul 2016

↩️ ↻️ 6,770 ❤️ 24,853

Almost all of these quoted tweets are posted from the Android:

Whether tweets start with a quotation mark (")



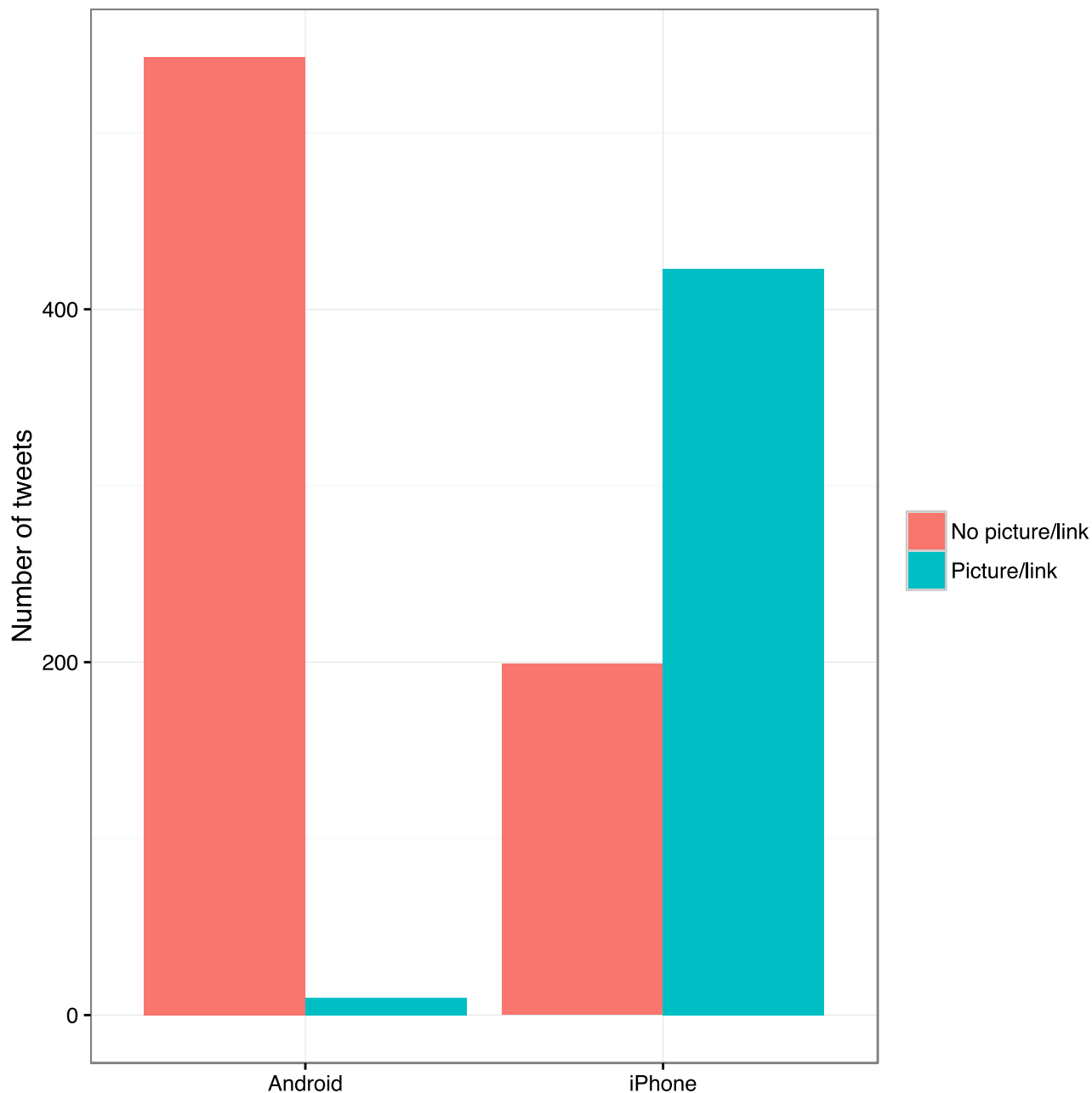
In the remaining by-word analyses in this text, I'll filter these quoted tweets out (since they contain text from followers that may not be representative of Trump's own tweets).

Somewhere else we can see a difference involves sharing links or pictures in tweets.

```
tweet_picture_counts <- tweets %>%
  filter(!str_detect(text, '^\"')) %>%
  count(source,
```

```
picture = ifelse(str_detect(text, "t.co"),
                 "Picture/link", "No picture/link")

ggplot(tweet_picture_counts, aes(source, n, fill = picture)) +
  geom_bar(stat = "identity", position = "dodge") +
  labs(x = "", y = "Number of tweets", fill = "")
```



It turns out tweets from the iPhone were **38 times as likely to contain either a picture or a link**. This also makes sense with our narrative: the iPhone (presumably run by the campaign) tends to write “announcement” tweets about events, like this:



Donald J. Trump realDonaldTrump

Follow

Thank you Windham, New Hampshire! [#TrumpPence16](#)[#MAGA](#)

7:19 PM - 6 Aug 2016 · Windham, NH, United States

6,243 20,851

While Android (Trump himself) tends to write picture-less tweets like:



Donald J. Trump realDonaldTrump

Follow

The media is going crazy. They totally distort so many things on purpose. Crimea, nuclear, "the baby" and so much more. Very dishonest!

2:31 PM - 7 Aug 2016

13,987 41,776

Comparison of words

Now that we're sure there's a difference between these two accounts, what can we say about the difference in the *content*? We'll use the [tidytext](#) package that [Julia Silge](#) and I developed.

We start by dividing into individual words using the `unnest_tokens` function (see [this vignette](#) for more), and removing some common "stopwords"²:

```
library(tidytext)

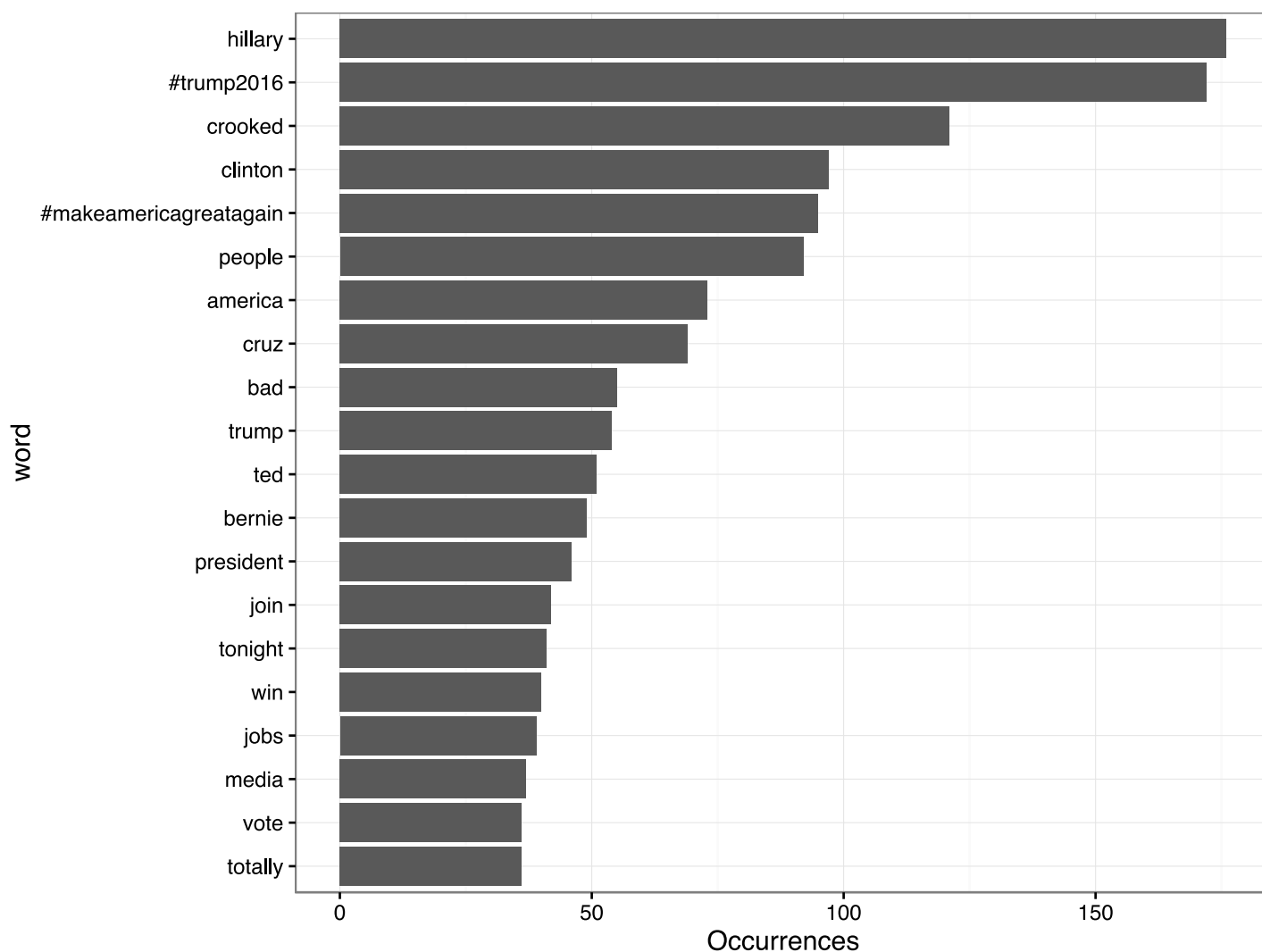
reg <- "([A-Za-z\\d#@']|'?![A-Za-z\\d#@'])"

tweet_words <- tweets %>%
  filter(!str_detect(text, '^\"')) %>%
  mutate(text = str_replace_all(text, "https://t.co/[A-Za-z\\d]+&"; "")) %>%
  unnest_tokens(word, text, token = "regex", pattern = reg) %>%
  filter(!word %in% stop_words$word,
         str_detect(word, "[a-z]"))

tweet_words
```

```
## # A tibble: 8,753 x 4
##           id source      created      word
##           <chr> <chr>      <time>    <chr>
## 1 676494179216805888 iPhone 2015-12-14 20:09:15 record
## 2 676494179216805888 iPhone 2015-12-14 20:09:15 health
## 3 676494179216805888 iPhone 2015-12-14 20:09:15 #makeamericagreatagain
## 4 676494179216805888 iPhone 2015-12-14 20:09:15 #trump2016
## 5 676509769562251264 iPhone 2015-12-14 21:11:12 accolade
## 6 676509769562251264 iPhone 2015-12-14 21:11:12 @trumpgolf
## 7 676509769562251264 iPhone 2015-12-14 21:11:12 highly
## 8 676509769562251264 iPhone 2015-12-14 21:11:12 respected
## 9 676509769562251264 iPhone 2015-12-14 21:11:12 golf
## 10 676509769562251264 iPhone 2015-12-14 21:11:12 odyssey
## # ... with 8,743 more rows
```

What were the most common words in Trump's tweets overall?



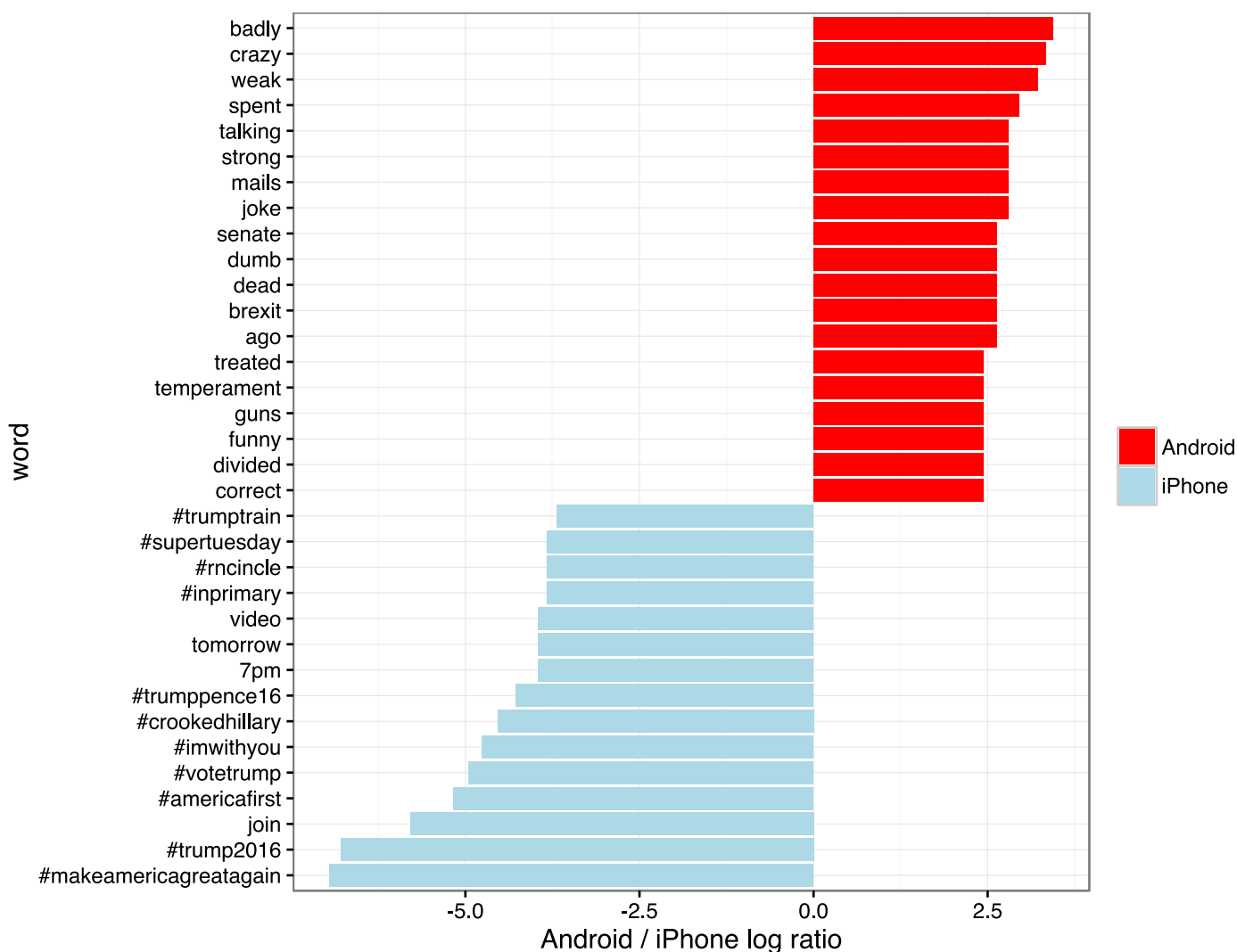
These should look familiar for anyone who has seen the feed. Now let's consider which words are most common from the Android relative to the iPhone, and vice versa. We'll use the simple measure of log odds ratio, calculated for each word as:³

$$\log_2\left(\frac{\frac{\# \text{ in Android} + 1}{\text{Total Android} + 1}}{\frac{\# \text{ in iPhone} + 1}{\text{Total iPhone} + 1}}\right)$$

$\log_2(\# \text{ in Android} + 1 \text{Total Android} + 1 \# \text{ in iPhone} + 1 \text{Total iPhone} + 1)$

```
android_iphone_ratios <- tweet_words %>%
  count(word, source) %>%
  filter(sum(n) >= 5) %>%
  spread(source, n, fill = 0) %>%
  ungroup() %>%
  mutate_each(funs((. + 1) / sum(. + 1)), -word) %>%
  mutate(logratio = log2(Android / iPhone)) %>%
  arrange(desc(logratio))
```

Which are the words most likely to be from Android and most likely from iPhone?



A few observations:

- **Most hashtags come from the iPhone.** Indeed, almost no tweets from Trump's Android contained hashtags, with some rare exceptions like [this one](#). (This is true only because we

filtered out the quoted “retweets”, as Trump does sometimes quote tweets [like this](#) that contain hashtags).

- **Words like “join” and “tomorrow”, and times like “7pm”, also came only from the iPhone.** The iPhone is clearly responsible for event announcements like [this one](#) (“Join me in Houston, Texas tomorrow night at 7pm!”)
- **A lot of “emotionally charged” words, like “badly”, “crazy”, “weak”, and “dumb”, were overwhelmingly more common on Android.** This supports the original hypothesis that this is the “angrier” or more hyperbolic account.

Sentiment analysis: Trump’s tweets are much more negative than his campaign’s

Since we’ve observed a difference in sentiment between the Android and iPhone tweets, let’s try quantifying it. We’ll work with the [NRC Word-Emotion Association](#) lexicon, available from the tidytext package, which associates words with 10 sentiments: **positive, negative, anger, anticipation, disgust, fear, joy, sadness, surprise, and trust.**

```
nrc <- sentiments %>%
  filter(lexicon == "nrc") %>%
  dplyr::select(word, sentiment)

nrc
```

```
## # A tibble: 13,901 x 2
##       word sentiment
##       <chr>      <chr>
## 1    abacus      trust
## 2  abandon      fear
## 3  abandon  negative
## 4  abandon  sadness
## 5 abandoned    anger
## 6 abandoned    fear
## 7 abandoned  negative
## 8 abandoned  sadness
## 9 abandonment  anger
## 10 abandonment  fear
## # ... with 13,891 more rows
```

To measure the sentiment of the Android and iPhone tweets, we can count the number of words in each category:

```

sources <- tweet_words %>%
  group_by(source) %>%
  mutate(total_words = n()) %>%
  ungroup() %>%
  distinct(id, source, total_words)

by_source_sentiment <- tweet_words %>%
  inner_join(nrc, by = "word") %>%
  count(sentiment, id) %>%
  ungroup() %>%
  complete(sentiment, id, fill = list(n = 0)) %>%
  inner_join(sources) %>%
  group_by(source, sentiment, total_words) %>%
  summarize(words = sum(n)) %>%
  ungroup()

head(by_source_sentiment)

```

```

## # A tibble: 6 x 4
##   source      sentiment total_words words
##   <chr>      <chr>      <int> <dbl>
## 1 Android      anger          4901    321
## 2 Android anticipation    4901    256
## 3 Android    disgust      4901    207
## 4 Android      fear      4901    268
## 5 Android      joy       4901    199
## 6 Android   negative      4901    560

```

(For example, we see that 321 of the 4901 words in the Android tweets were associated with “anger”). We then want to measure how much more likely the Android account is to use an emotionally-charged term relative to the iPhone account. Since this is count data, we can use a Poisson test to measure the difference:

```

library(broom)

sentiment_differences <- by_source_sentiment %>%
  group_by(sentiment) %>%
  do(tidy(poisson.test(.$words, .$total_words)))

sentiment_differences

```

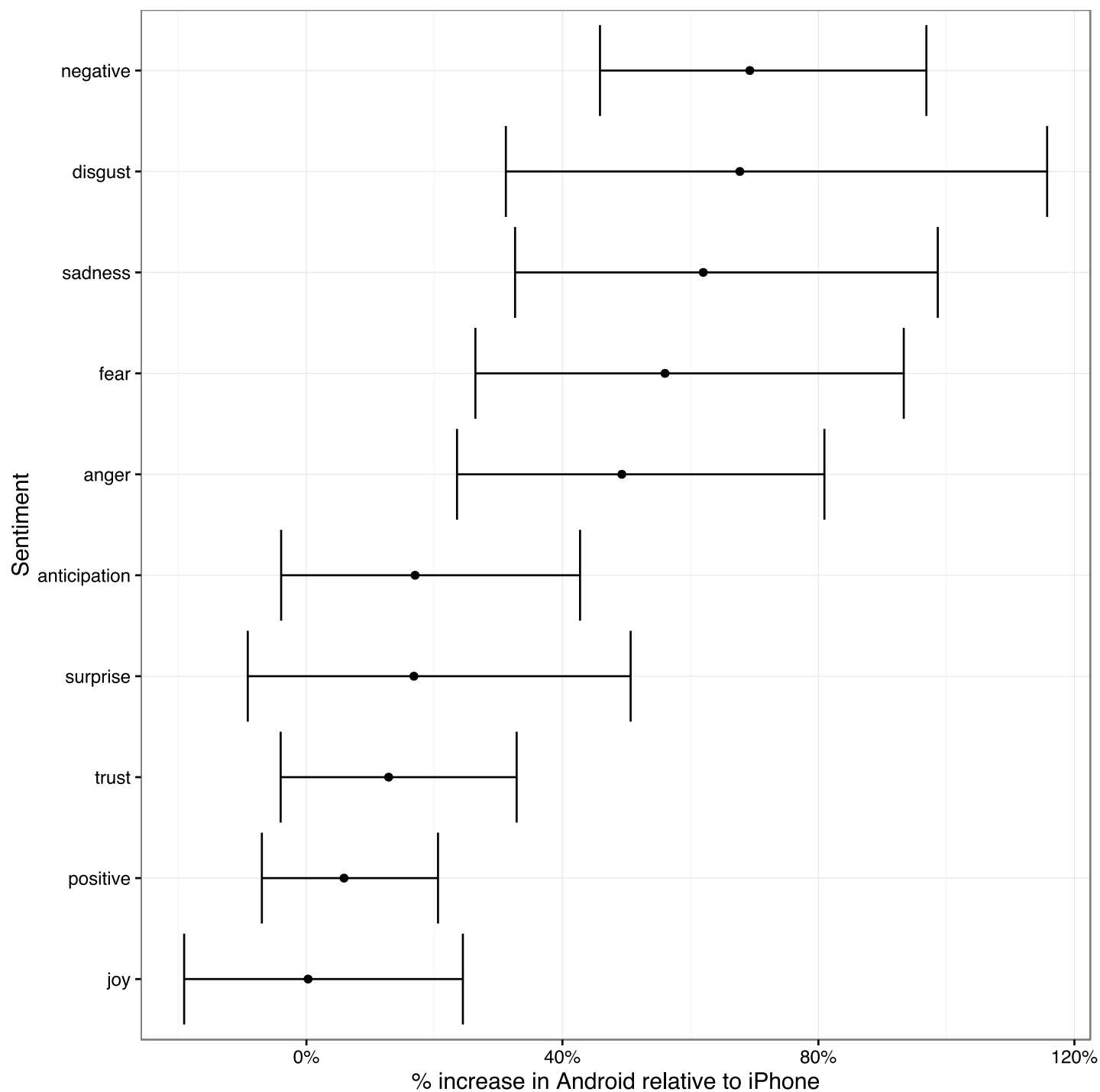
```

## Source: local data frame [10 x 9]
## Groups: sentiment [10]
##
##   sentiment estimate statistic      p.value parameter  conf.low
##   <chr>      <dbl>      <dbl>      <dbl>      <dbl>      <dbl>
## 1      anger  1.492863      321 2.193242e-05  274.3619  1.2353162

```

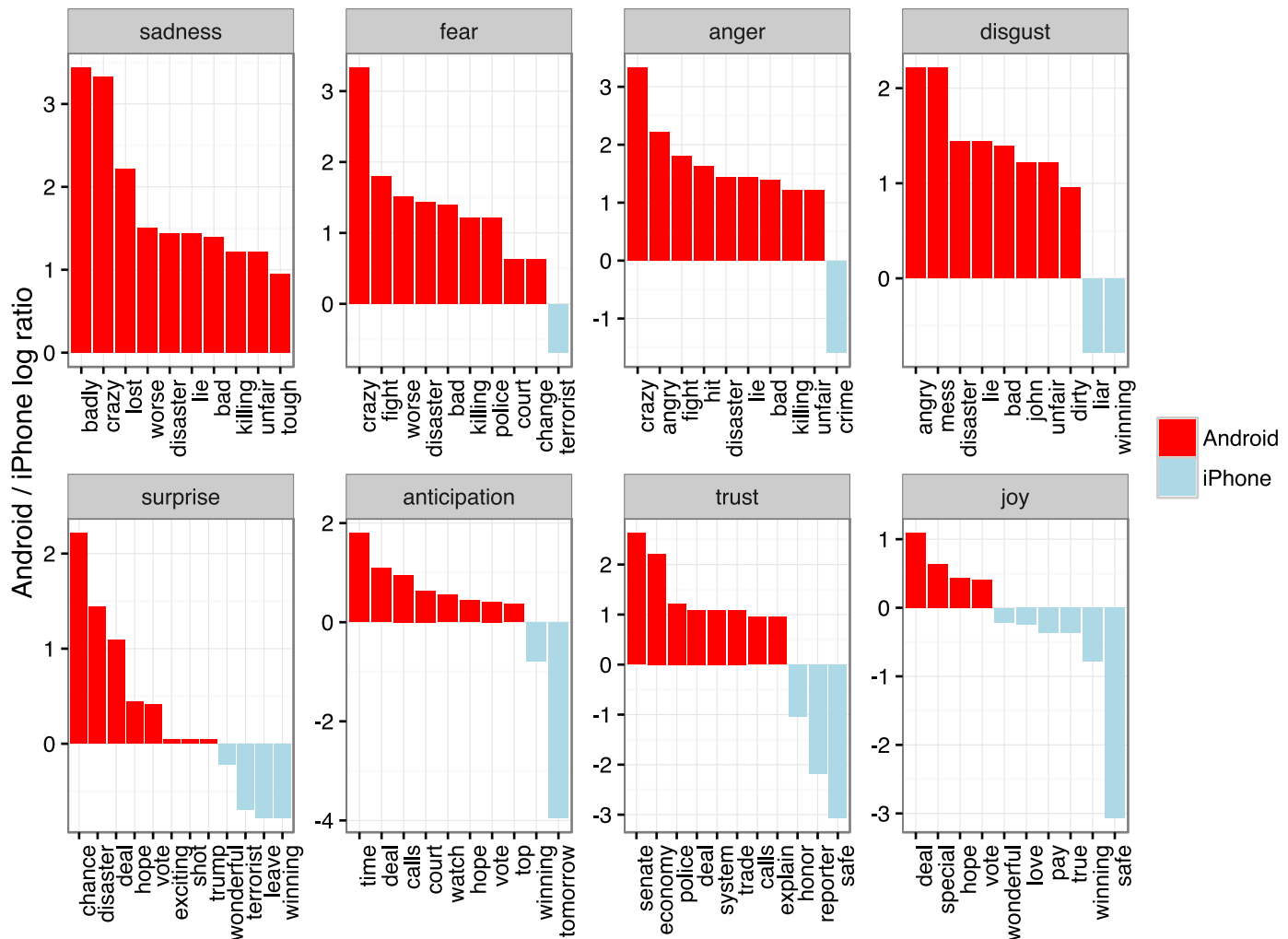
```
## 2 anticipation 1.169804      256 1.191668e-01 239.6467 0.9604950
## 3 disgust 1.677259      207 1.777434e-05 170.2164 1.3116238
## 4 fear 1.560280      268 1.886129e-05 225.6487 1.2640494
## 5 joy 1.002605      199 1.000000e+00 198.7724 0.8089357
## 6 negative 1.692841      560 7.094486e-13 459.1363 1.4586926
## 7 positive 1.058760      555 3.820571e-01 541.4449 0.9303732
## 8 sadness 1.620044      303 1.150493e-06 251.9650 1.3260252
## 9 surprise 1.167925      159 2.174483e-01 148.9393 0.9083517
## 10 trust 1.128482      369 1.471929e-01 350.5114 0.9597478
## # ... with 3 more variables: conf.high <dbl>, method <fctr>,
## # alternative <fctr>
```

And we can visualize it with a 95% confidence interval:



Thus, Trump's Android account uses about 40-80% more words related to **disgust**, **sadness**, **fear**, **anger**, and other “negative” sentiments than the iPhone account does. (The positive emotions weren't different to a statistically significant extent).

We're especially interested in which words drove this different in sentiment. Let's consider the words with the largest changes within each category:



This confirms that lots of words annotated as negative sentiments (with a few exceptions like “crime” and “terrorist”) are more common in Trump's Android tweets than the campaign's iPhone tweets.

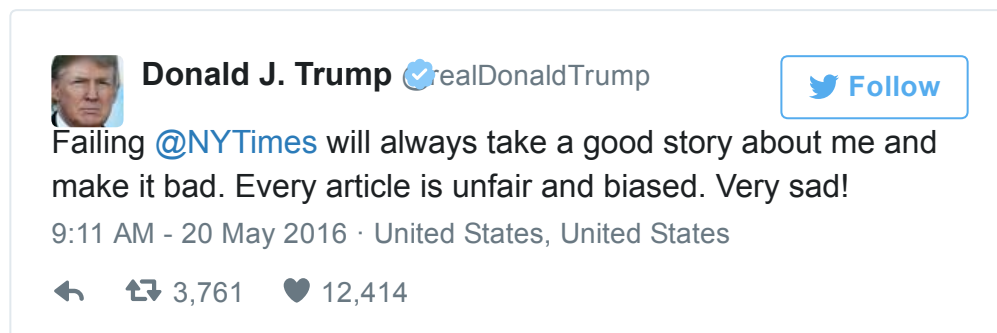
Conclusion: the ghost in the political machine

I was fascinated by the recent [New Yorker article](#) about Tony Schwartz, Trump's ghostwriter for The Art of the Deal. Of particular interest was how Schwartz imitated Trump's voice and philosophy:

In his journal, Schwartz describes the process of trying to make Trump's voice palatable in the book. It was kind of “a trick,” he writes, to mimic Trump's blunt, staccato, no-apologies delivery while making him seem almost boyishly appealing.... Looking back at the text now, Schwartz says, “I created a character far more winning than Trump actually is.”

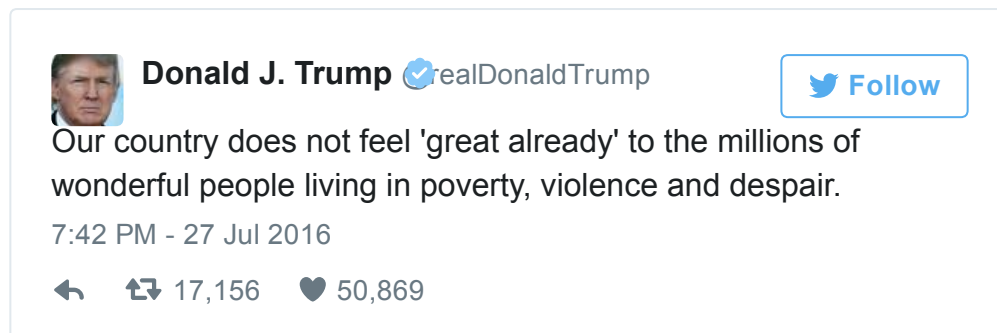
Like any journalism, data journalism is ultimately about human interest, and there's one human I'm interested in: who is writing these iPhone tweets?

The majority of the tweets from the iPhone are fairly benign declarations. But consider cases like these, both posted from an iPhone:



These tweets certainly sound like the Trump we all know. Maybe our above analysis isn't complete: maybe Trump has sometimes, however rarely, tweeted from an iPhone (perhaps dictating, or just using it when his own battery ran out). But what if our hypothesis is right, and these weren't authored by the candidate- just someone trying their best to sound like him?

Or what about tweets like this (also iPhone), which defend Trump's slogan- but doesn't really sound like something he'd write?



A lot has been written about Trump's mental state. But I'd really rather get inside the head of this anonymous staffer, whose job is to imitate Trump's unique cadence ("Very sad!"), or to put a positive spin on it, to millions of followers. Are they a true believer, or just a cog in a political machine, mixing whatever mainstream appeal they can into the @realDonaldTrump concoction? Like Tony Schwartz, will they one day regret their involvement?

1. To keep the post concise I don't show all of the code, especially code that generates figures. But you can find the full code [here](#). ↵

2. We had to use a custom regular expression for Twitter, since typical tokenizers would split the # off of hashtags and @ off of usernames. We also removed links and ampersands (`&`) from the text. ↵
 3. The “plus ones,” called Laplace smoothing are to avoid dividing by zero and to put more trust in common words. ↵
-



David Robinson

Data Scientist at Stack Overflow, works in R and Python.

[Email](#) [Twitter](#) [Github](#) [Stack Overflow](#)

Subscribe

[Subscribe to this blog](#)

Recommended Blogs

- [R Bloggers](#)
- [RStudio Blog](#)
- [R4Stats](#)
- [Simply Statistics](#)

Text analysis of Trump's tweets confirms he writes only the (angrier) Android half was published on August 09, 2016.