

# Brexit through the Lens of Data Science

Kenneth Benoit  
(joint work with Akitaka Matsuo)  
Department of Methodology  
LSE



## Text mining and data/social science

- Treats text as “data” that informs us about the authors of the text - the text itself is only incidental
- Key is some form of comparison, to gain insights about differences
- Involves the application of statistical models to judge differences using probability statements

## Analyzing Brexit through Twitter

- EU-funded project
- We collected some 26 million Tweets from January - July 2016
- capture based on #hashtags, @usernames, and search terms

## Search terms

### Hashtags:

#betterdealforbritain  
#betteroffout  
#brexit  
#euref  
#eureferendum  
#eusummit  
#getoutnow  
#leaveeu  
#no2eu  
#notoEU  
#strongerin  
#ukineu  
#voteleave  
#wewantout  
#yes2eu  
#yestoEU  
  
brexit

### Username:

@vote\_leave  
@brexitwatch  
@eureferendum  
@ukandeu  
@notoEU  
@leavehq  
@ukineu  
@leaveeuofficial  
@ukleave\_eu  
@strongerin  
@yesforeurope  
@grassroots\_out  
@stronger\_in

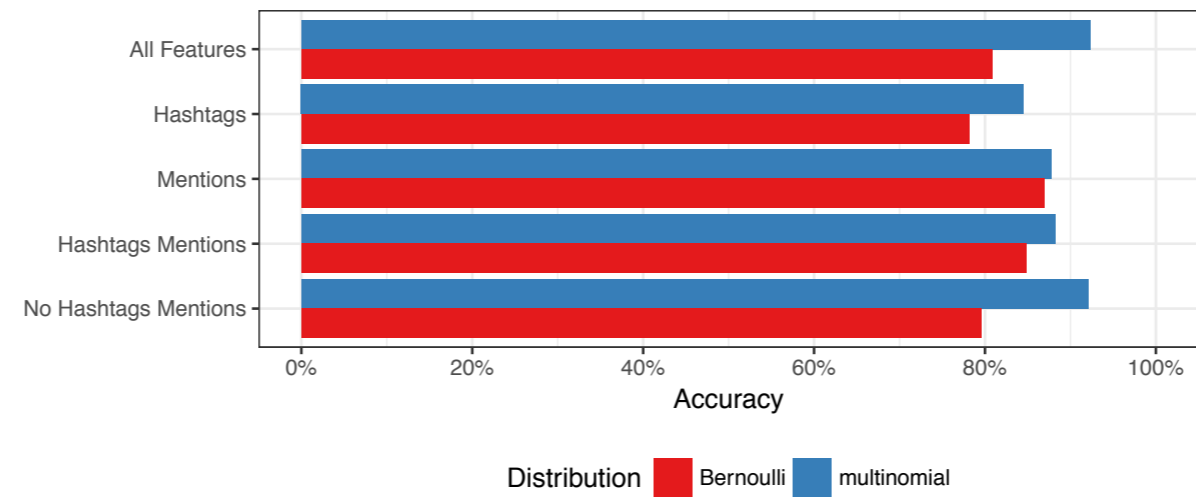
## Users in data

- 3.6M unique users
- number of tweets:
  - average: 7.2
  - median: 1 (more than 50% had only one tweet)
  - max: 81.1K

## First application: Predicting Leave v. Remain users

- Method: Naive Bayes classifier
- Data source: combined tweet corpus at user level
- Creating training data
  - Select “power-users” (more than 100 tweets in the corpus, 15K users)
  - Check the use of pre-determined set of “leave” and “remain” hashtags
  - Calculate the difference in the use of leave and remain hashtags. Construct training data from top and bottom 10% of power users

## Predictive accuracy based on feature types



- Use bag of words method (unigrams of #hashtags and @user\_mentions, remove features used less than 20 users)

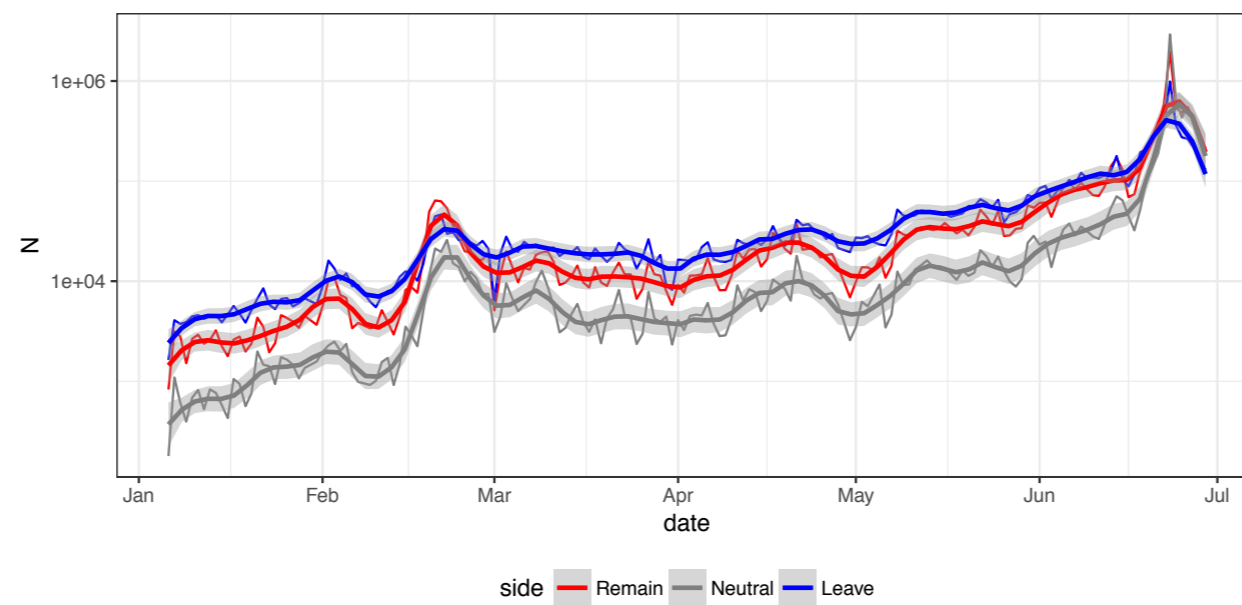
our version

## Predicting Leave v. Remain

	N	%
Remain	9,780,223	36.93%
Neutral	7,786,297	29.40%
Leave	8,914,207	33.66%

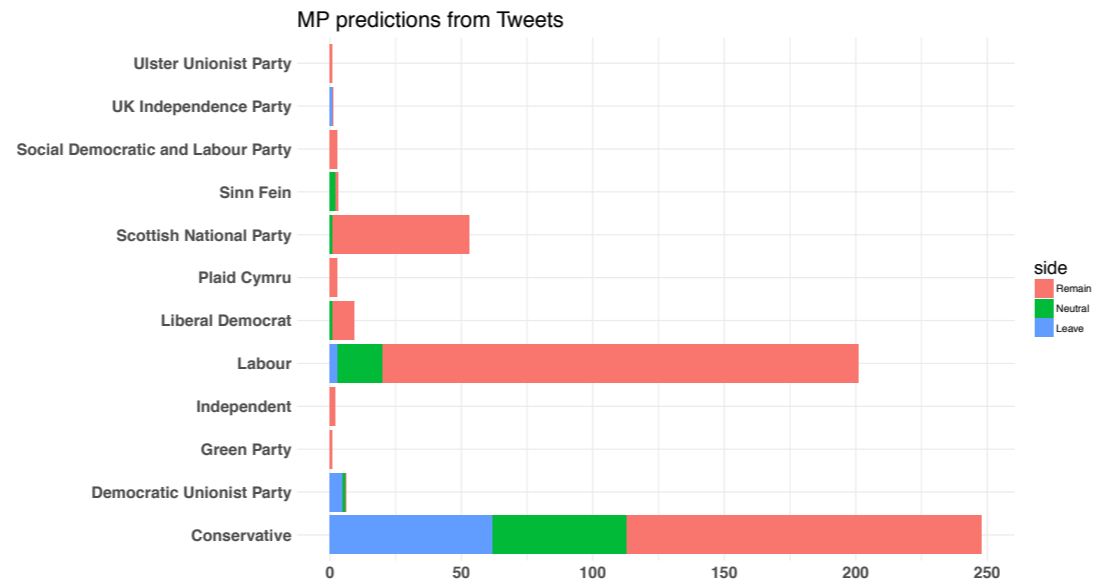


## Patterns of posts, Leave v. Remain

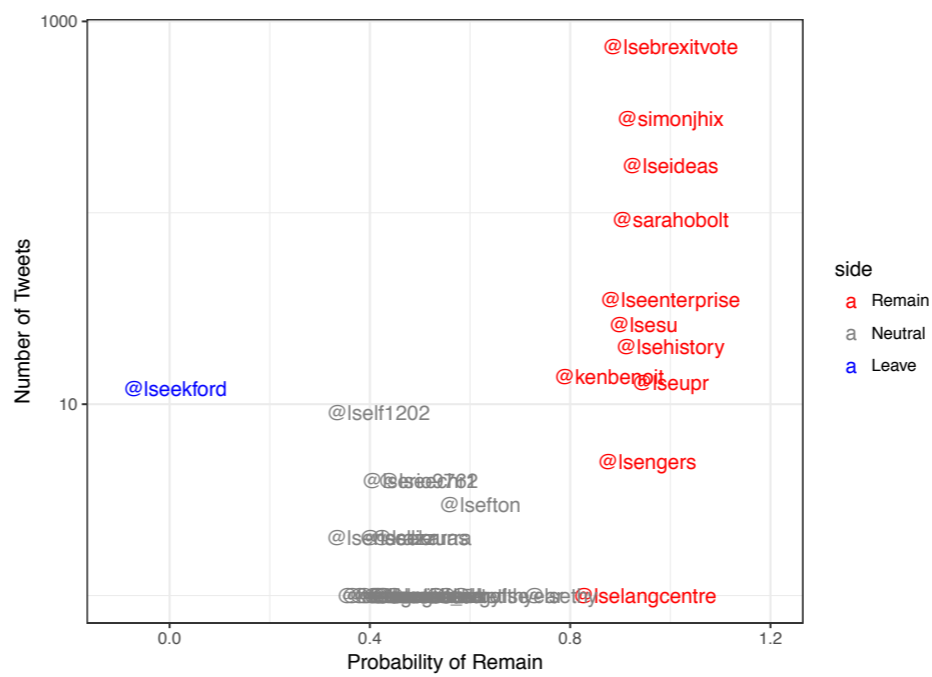


our version

## MP side predictions by party



A few well-known  
Twitter users



Profile header for **linda seekford** (@lseekford) with 4,078 tweets, 232 following, 123 followers, and 3,376 likes. Includes a "Follow" button and a blue header bar.

Profile information for **linda seekford** (@lseekford), joined October 2012. Includes a "Tweet to linda seekford" button and a gallery of 40 photos and videos.

Tweets tab showing a tweet from **linda seekford** (@lseekford) dated Jan 20: "Sad liberals spewing hate". Includes a retweet and a quote tweet from **AJ+** (@ajplus) with a link to a Periscope broadcast.

Tweet from **linda seekford** (@lseekford) dated Jan 20: "Thank God we survived the last 8 yrs".

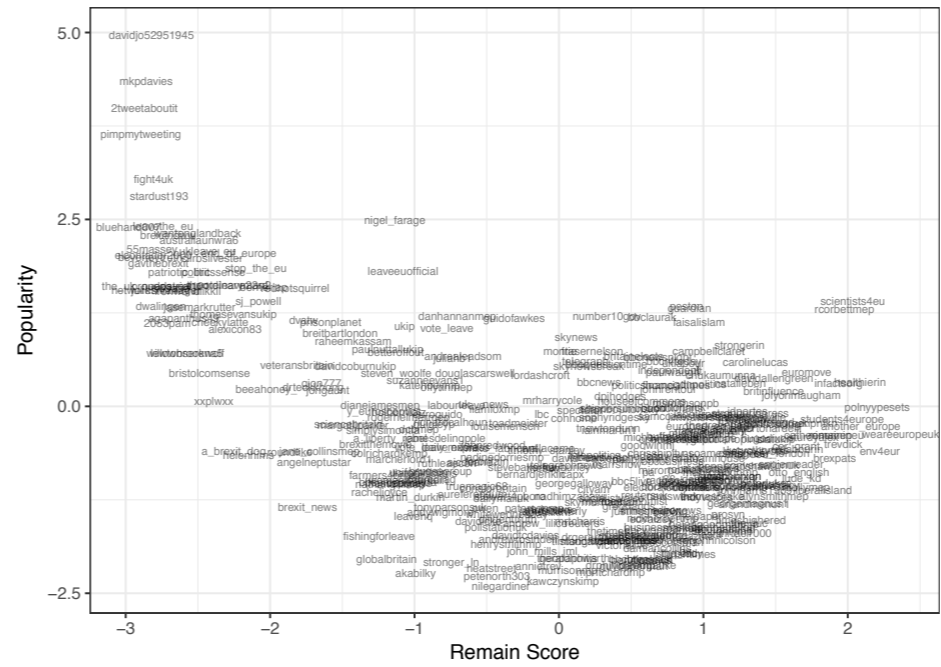
Tweet from **linda seekford** (@lseekford) dated Jan 17: "John McCain is, a snake. Go away you old goat".

Tweet from **linda seekford** (@lseekford) dated Jan 20: "Mass impeachment for the traitor democracts".

# Leave and Remain Hashtags

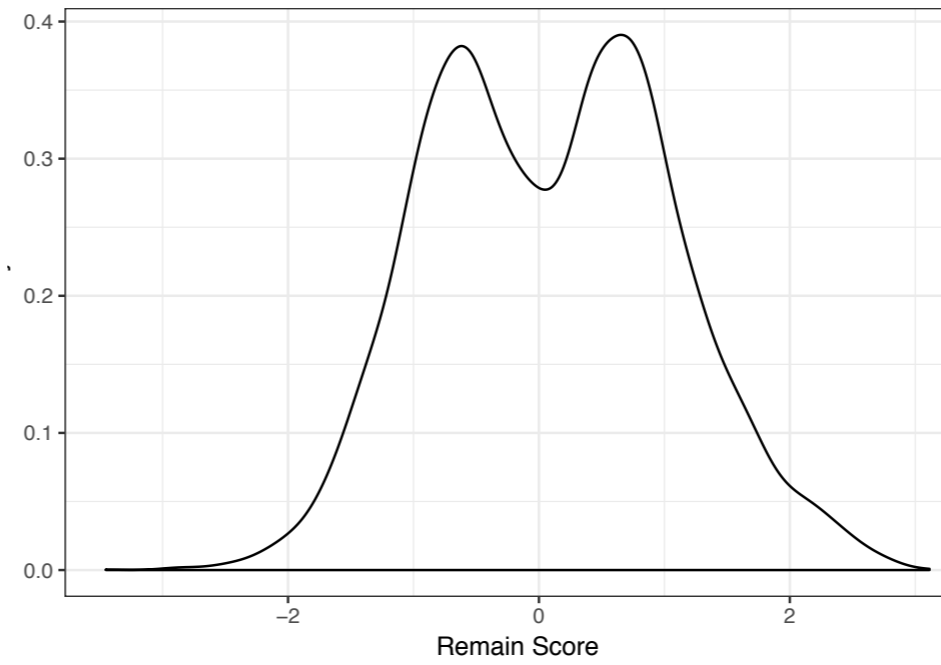


Followership network analysis, followers' positions (Barberà 2015)



Plot of top 300 hashtags. Cross side edges are highlighted. There are some connections, but the number of edges is smaller than the next figure.

Followership  
network  
analysis,  
followers  
distribution  
(Barberà  
2015)



Plot of top 300 hashtags. Cross side edges are highlighted. There are some connections, but the number of edges is smaller than the next figure.

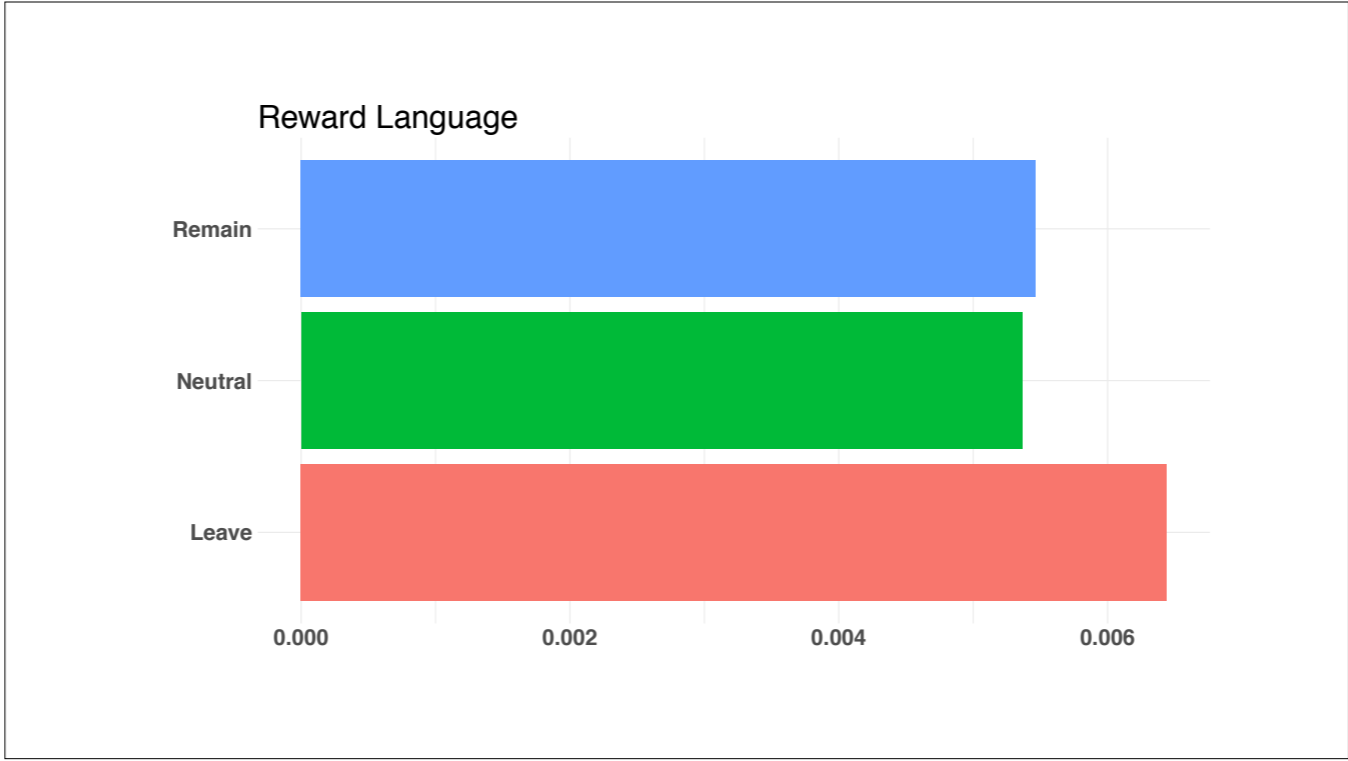
## Sentiment Analysis

- Looks up terms from the Linguistic Inquiry and Word Count, a psychological dictionary
- Contains categories about:
  - positive and negative emotion
  - politics
  - power
  - quantitative language
  - tentative language
  - sadness
  - future v. past orientation

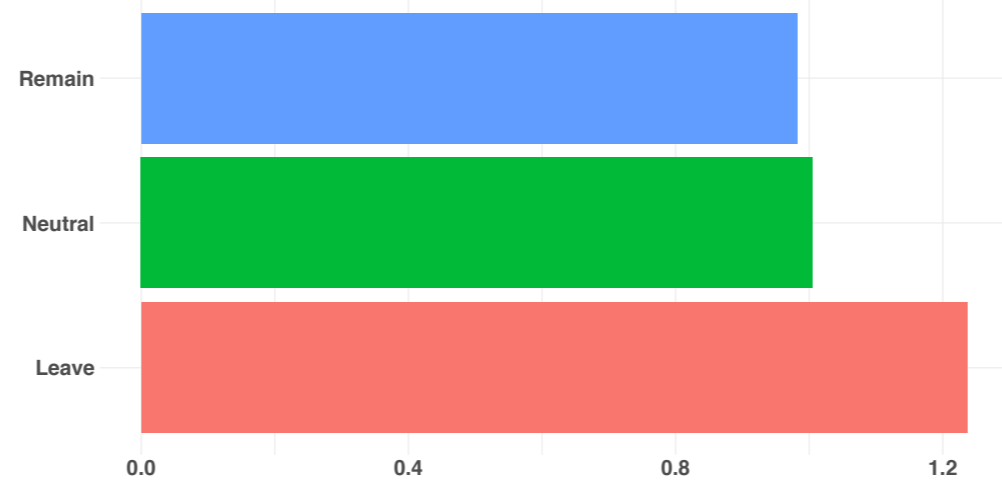


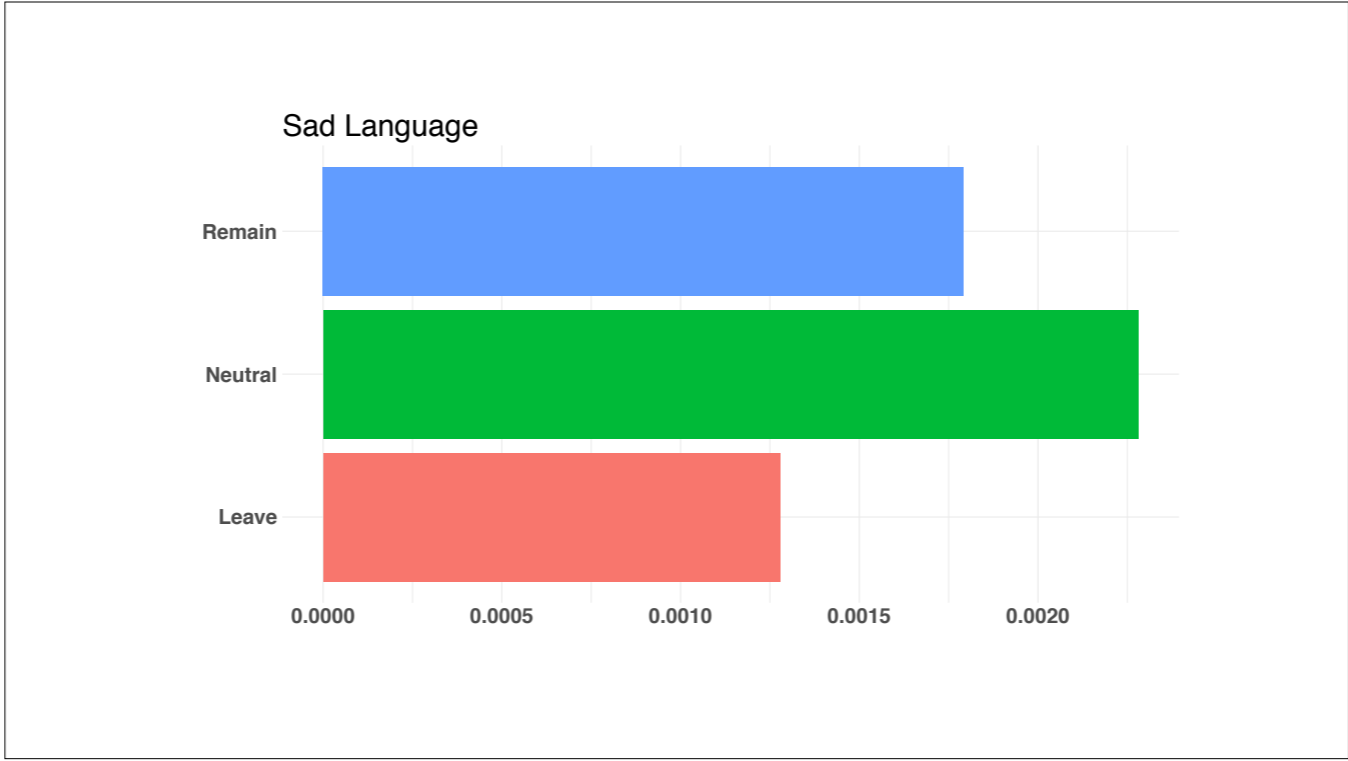
## example: "reward" language

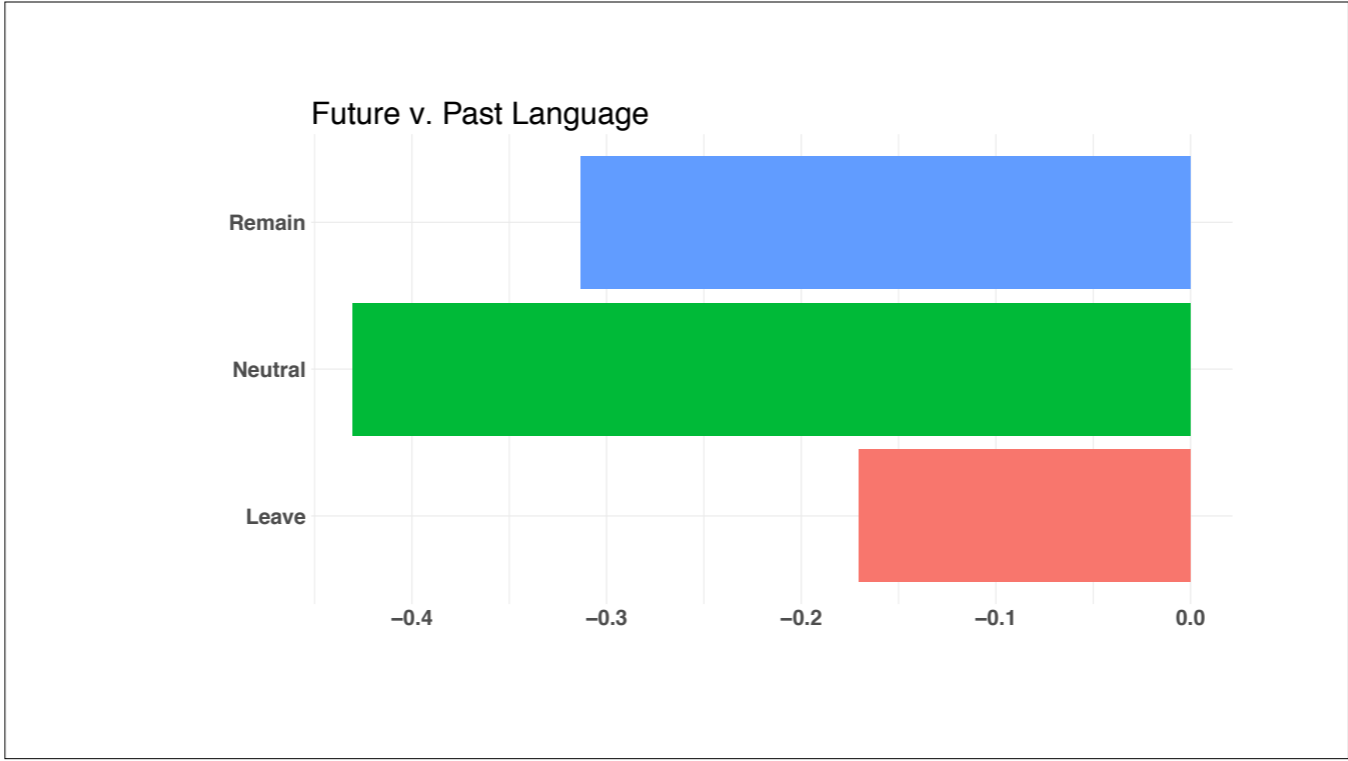
```
> data_dictionary_liwc[["reward"]]
[1] "access*" "accrue*" "accumul*" "achievable" "achieve*" "achiev*"
[7] "acquir*" "add" "added" "adding" "adds" "advanc*"
[13] "advantag*" "adventur*" "amass*" "approach" "approach*" "approach*"
[19] "approaching" "award*" "benefit" "benefits" "best" "bet"
[25] "bets" "better" "betting" "bold" "bonus*" "confidence"
[31] "confident" "confidently" "crave" "craving" "dare" "dared"
[37] "dares" "daring" "desir*" "eager" "eagerly" "eagerness"
[43] "earn" "earned" "earning" "earnings" "earns" "enthus*"
[49] "excite" "excited" "excitedly" "excitement" "exciting" "fearless*"
[55] "fulfill*" "gain*" "get" "gets" "getting" "goal*"
[61] "good" "got" "gotten" "great" "greed*" "invigor*"
[67] "jackpot*" "luck" "lucky" "obtain" "obtainable" "obtained"
[73] "obtaining" "obtains" "opportun*" "optimal*" "optimism" "optimistic"
[79] "perfect" "perfected" "perfecting" "perfection" "perfectly" "plus"
[85] "positive" "positively" "positives" "positivi*" "prize*" "profit*"
[91] "promot*" "reward*" "score*" "scoring" "seize*" "snag*"
[97] "steal*" "stole" "succeed*" "success" "successes" "successful"
[103] "successfully" "surpass*" "take" "taken" "takes" "taking"
[109] "took" "triumph*" "victor*" "wager" "wagered" "wagering"
[115] "wagers" "willing" "win" "winn*" "wins" "won"
```

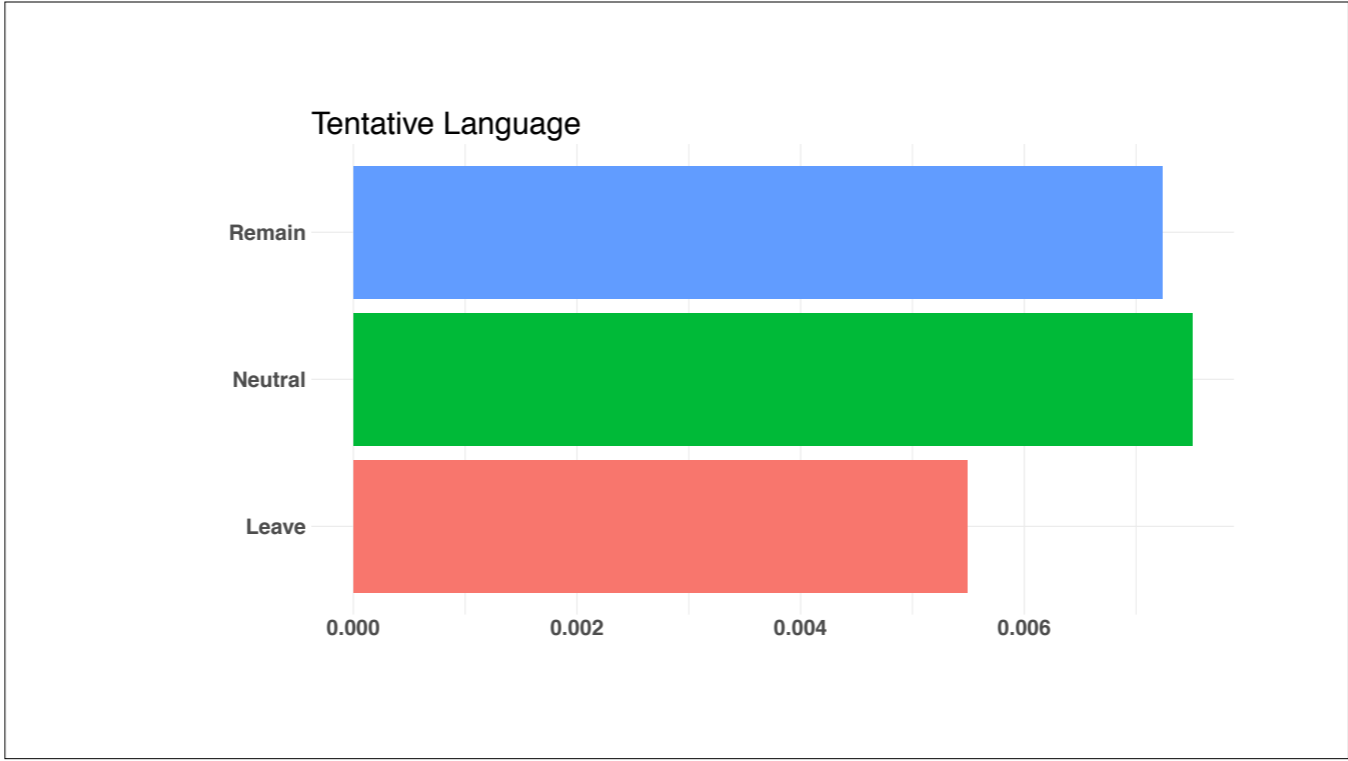


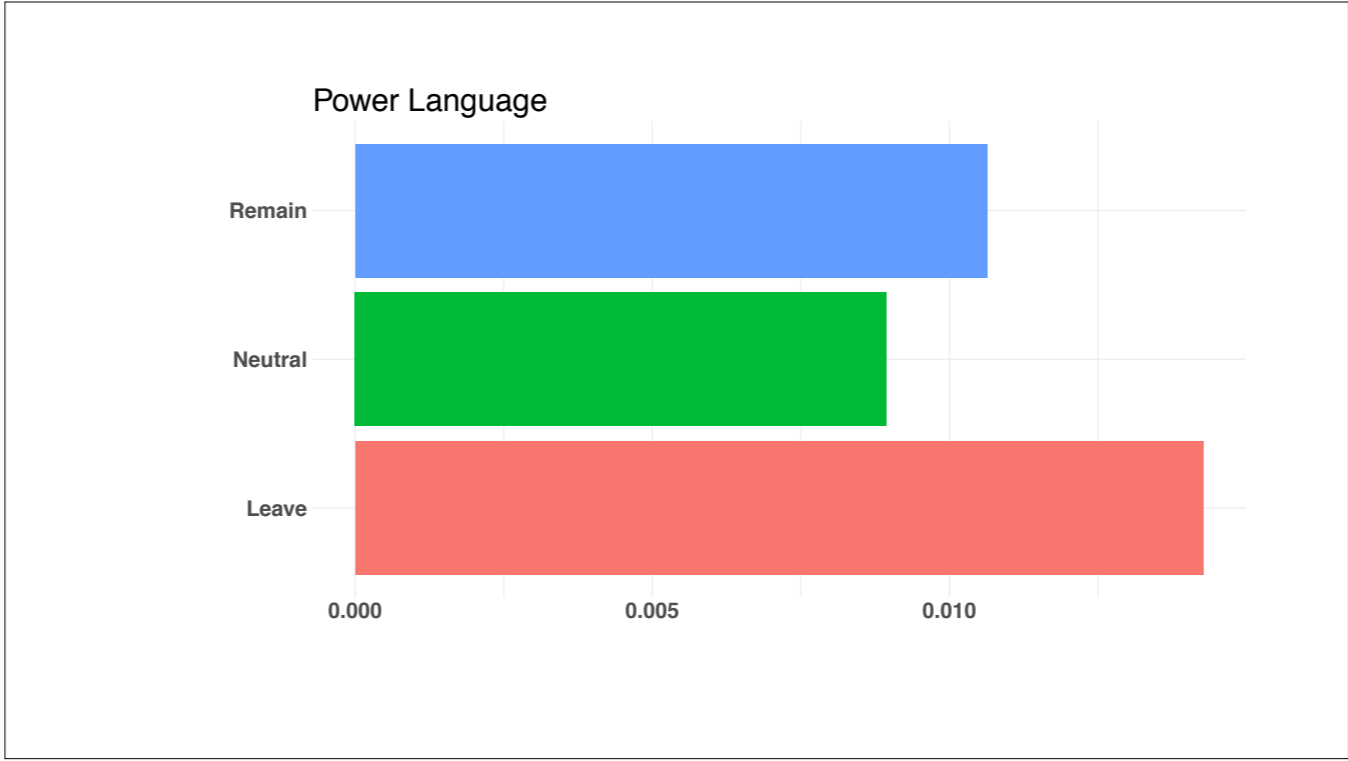
Positive v. Negative Emotion Language

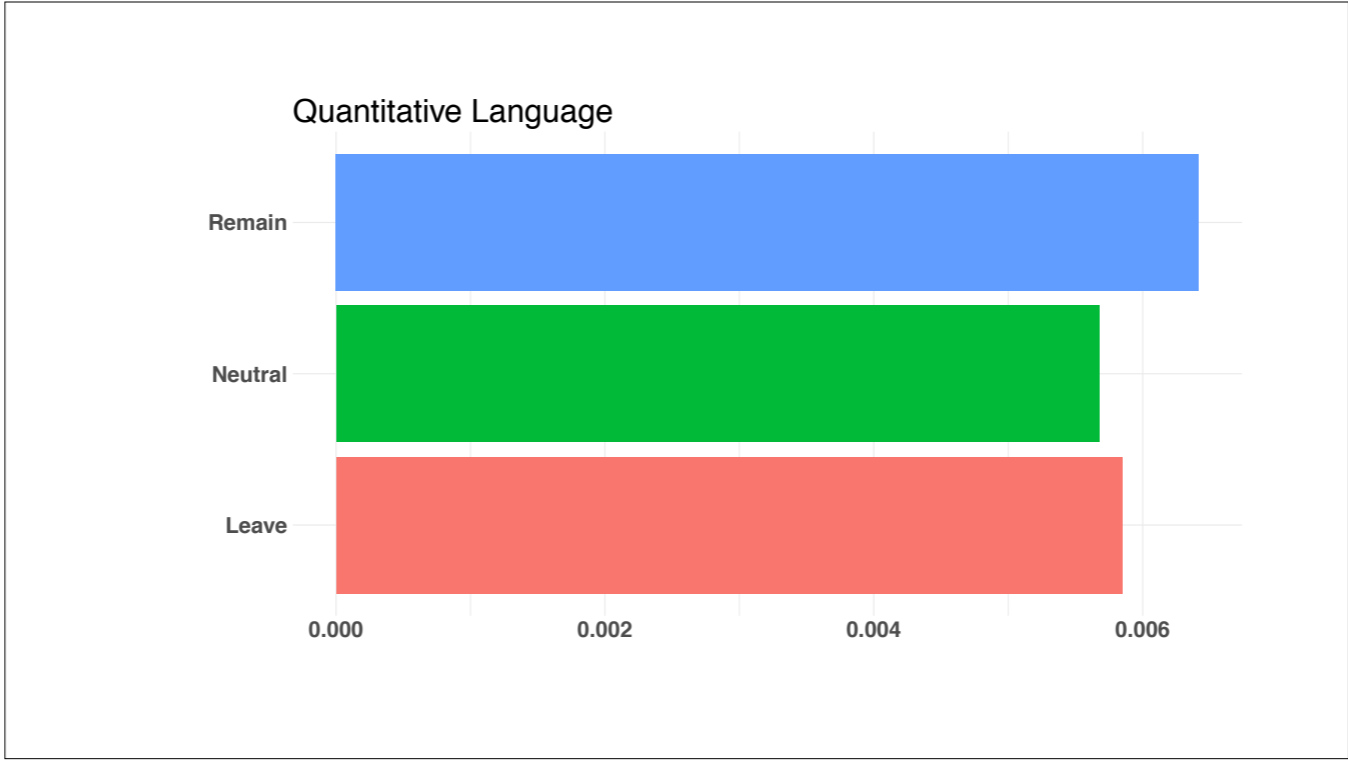








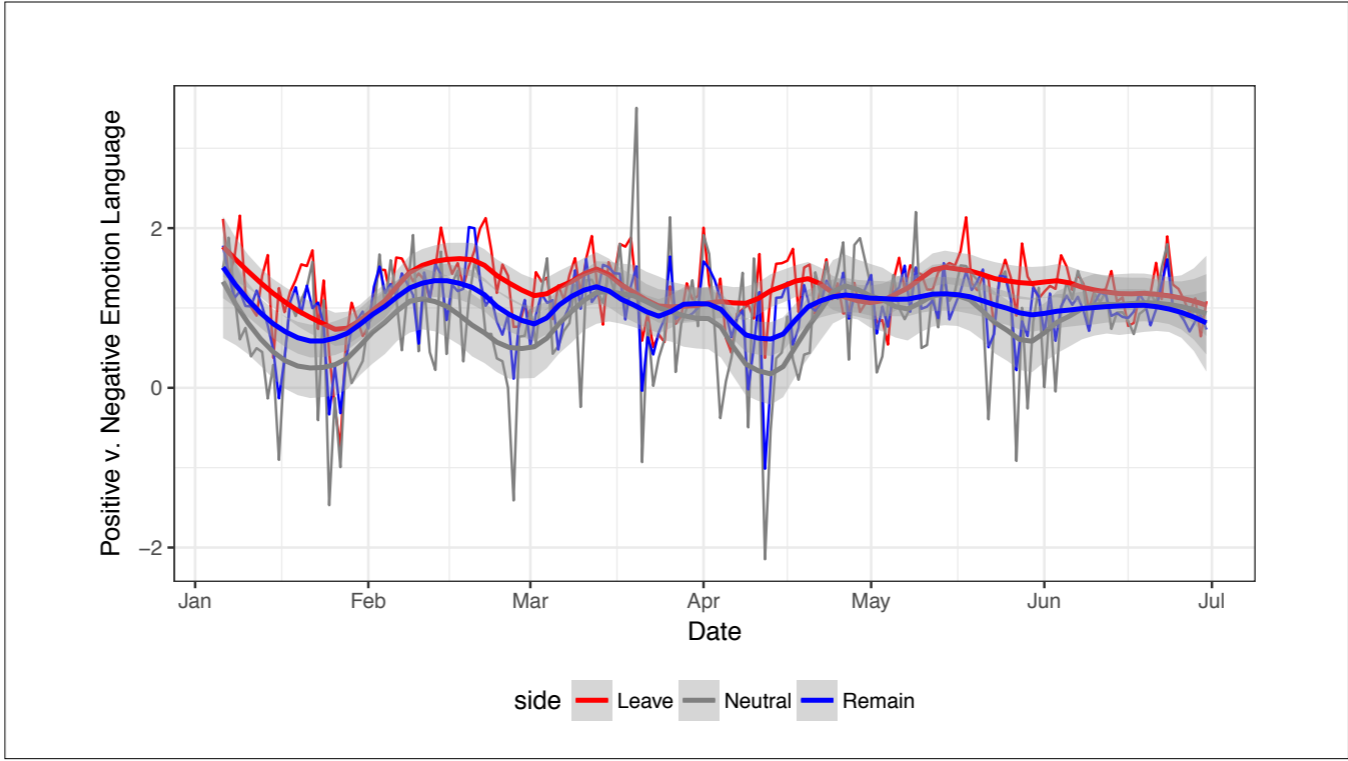


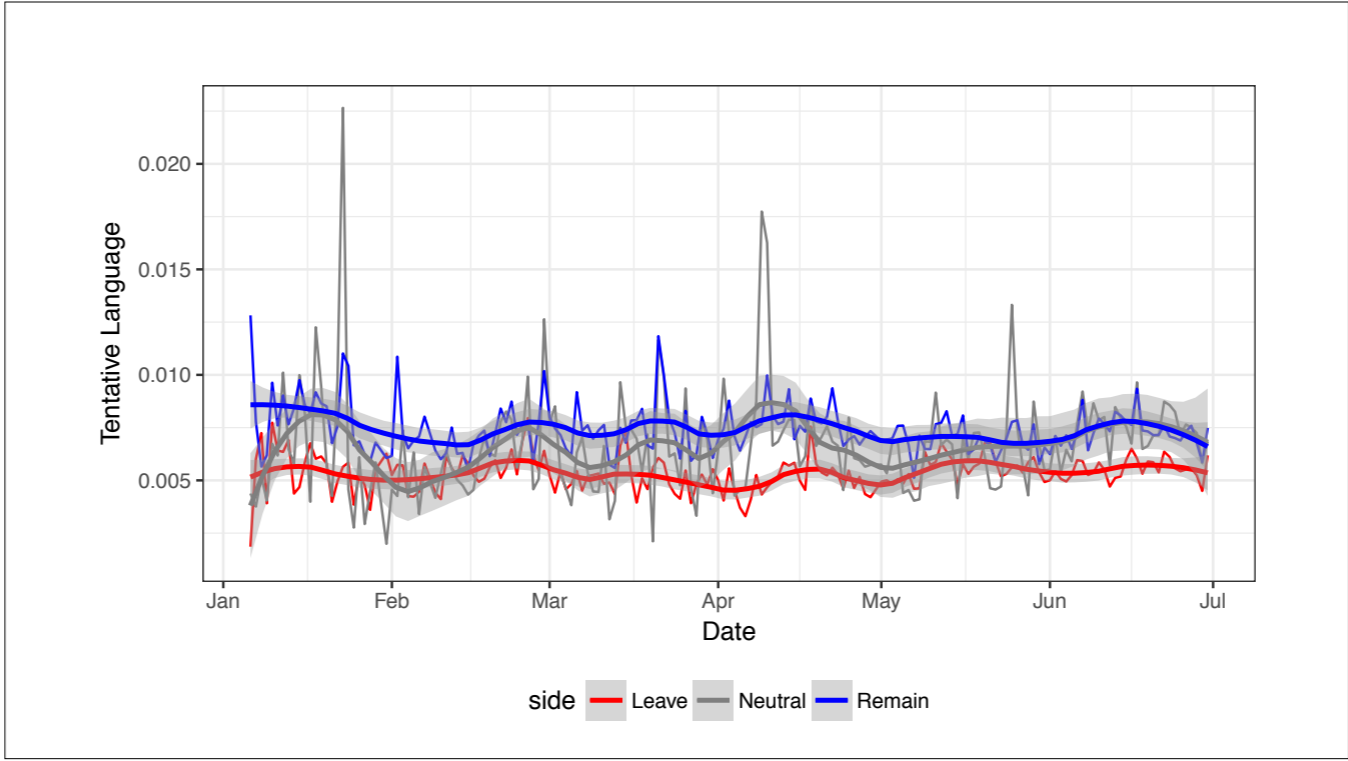




## Sentiment Analysis conclusions: Leave was more

- reward-oriented
- positive
- assertive of power
- less quantitative
- less tentative
- less sadness
- future-, versus past-, oriented



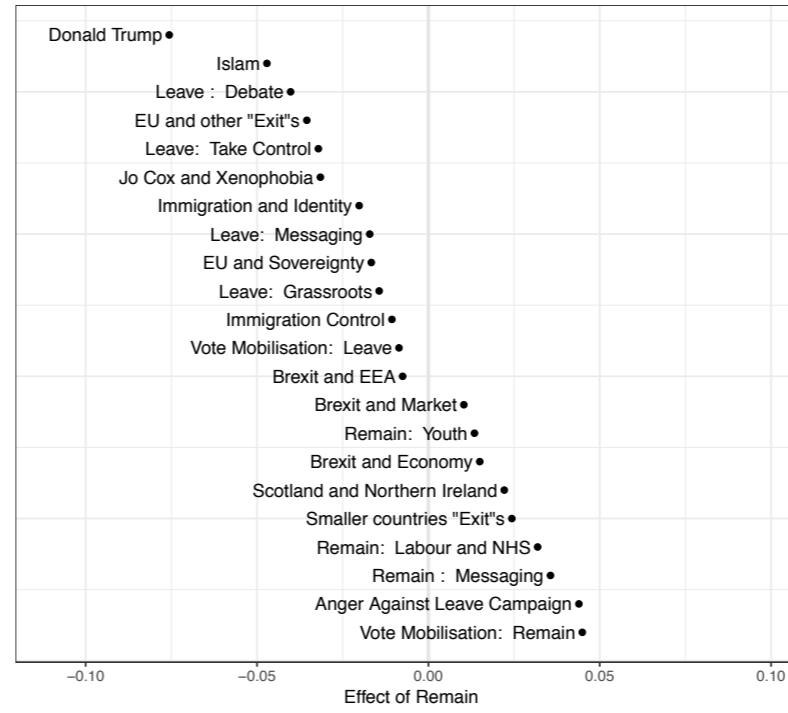


## Topic models

- Using unsupervised machine learning
  - detect topics in twitter conversation
  - topic distributions across sides
- Data
  - tweets from Leave and Remain accounts identified from NB classification
  - combine tweets from each account
  - 700,000 accounts are included
- Method
  - Structural Topic Model (STM) by Roberts, Stewart, and Tingley (2014)
  - Estimate models with 10-40 topics (incremented by 5)
  - covariate in topic prevalence: predicted sides of accounts

## STM Results

- Selected topics from 40 topics model





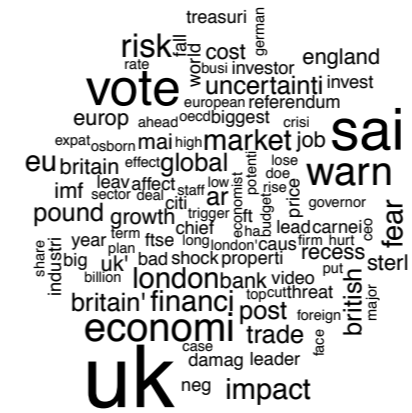


## STM Results (examples of remain economic topics)

Brexit and Market



Brexit and Economy





## Topics by side

Remain	Leave
Brexit ideologues	Brexit Movie
Economic consequence of Brexit	Brexit, UKIP, Leave EU
Encourage participation	David Cameron and Brussels
Exchange rate	David Camerons lies
Financial risk	Debate and discussions
Globalization and migration	Economic/Financial Sovereignty
Ireland	Free trade and migrants
Obama in London	General leave argument
Stock market risk	Jobs and social security
StrongerIn	Leave campaign
Tabloid (Trump, Queen, Jo Cox)	News discussion
Talking about articles on Brexit	Reasons to Leave EU
The city, and big business	Take back control
Voter registration	Undemocratic EU
	Vote Leave

# Networks of topics



## Summary: Text analysis was used to

- predict the side of the user
- determining the most common hashtags by side
- using followership networks to estimate “ideology”
- measuring sentiment using dictionaries
- analyzing the topics that were discussed
- mapping connections between topics via networks

# How? Using software written in R (and C++ and Python)



CRAN 1.3.0 downloads 8055/month downloads 136K build passing build passing codecov 88% DOI 10.5281/zenodo.1219063  
JOSS Under Review

Aug 12, 2012 – Jun 14, 2018

Contributions: Commits

Contributions to master, excluding merge commits

