

Quantitative text analysis: Social Media Data

Blake Miller

MY 459: Quantitative Text Analysis

March 25, 2019

Course website: lse-my459.github.io

1. Overview and Fundamentals
2. Descriptive Statistical Methods for Text Analysis
3. Automated Dictionary Methods
4. Machine Learning for Texts
5. Supervised Scaling Models for Texts
6. *Reading Week*
7. Unsupervised Models for Scaling Texts
8. Similarity and Clustering Methods
9. Topic models
10. Word embeddings
11. Working with Social Media

Evaluation

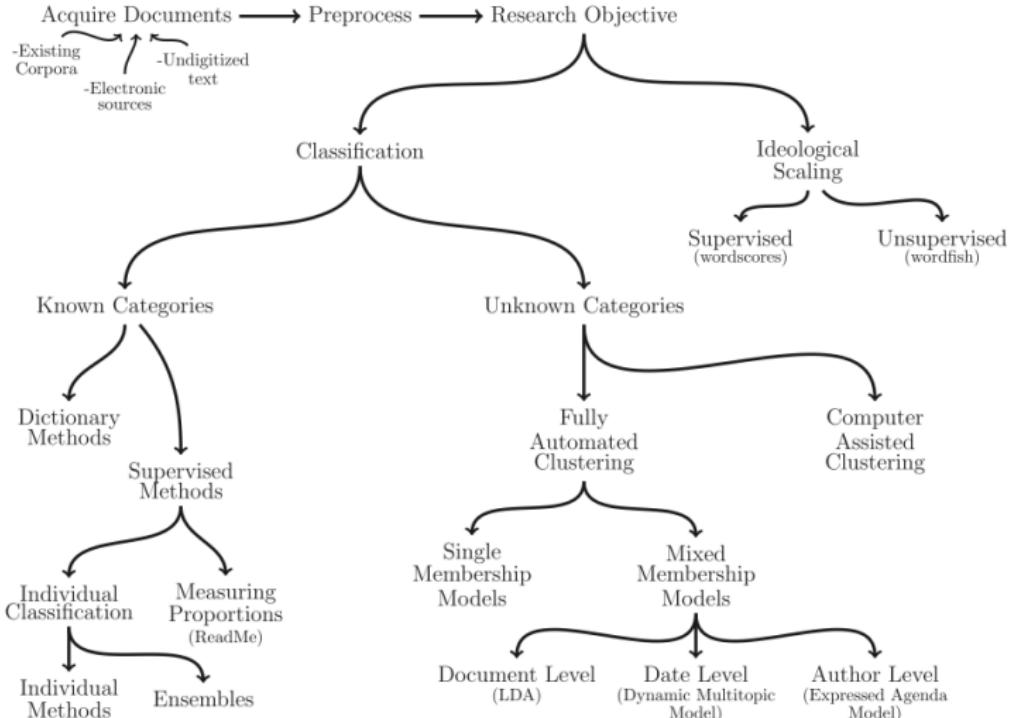
Formative coursework:

- ▶ Five problem sets (60% of course grade)
- ▶ Project:
 - ▶ Original analysis of texts using methods covered in class
 - ▶ It can replicate or extend a published work
 - ▶ 3,000 words (5,000 for MY559), due at the beginning of ST (May 3rd, 5pm)
 - ▶ 40% of course grade
 - ▶ Additional instructions available on Moodle
 - ▶ Submission via Moodle

Assessment criteria

- ▶ **70–100:** Very Good to Excellent (Distinction).
 - ▶ Perceptive, focused use of a good depth of material with a critical edge. Original ideas or structure of argument.
- ▶ **60–69:** Good (Merit)
 - ▶ Perceptive understanding of the issues plus a coherent well-read and stylish treatment though lacking originality
- ▶ **50–59:** Satisfactory (Pass)
 - ▶ A “correct” answer based largely on lecture material. Little detail or originality but presented in adequate framework. Small factual errors allowed.
- ▶ **30–49:** Unsatisfactory (Fail)
- ▶ **0–29:** Unsatisfactory (Bad fail)
 - ▶ Based entirely on lecture material but unstructured and with increasing error component. Concepts are disordered or flawed. Poor presentation. Errors of concept and scope or poor in knowledge, structure and expression.

Overview of text as data methods



Outline

- ▶ Social media data
 - ▶ Motivation
 - ▶ Opportunities and challenges
 - ▶ Twitter data
 - ▶ Facebook data
- ▶ Extracting data from PDF files

Why social media data?

- ▶ Volume: 500M registered users, 400M tweets per day (March 2013), Facebook has 1.15billion users, on average post 36 times a month — coverage and representation
- ▶ Real time — new data is available publicly immediately on current events
- ▶ Metadata — geographic location, user device, profile, timestamp and other metadata is accessible.



George Takei

March 28 at 10:10pm · 4

Who's with me.



Like · Comment · Share

408,735 people like this.

66,990 shares



Bon Alimago

@karma_thief



I need a hug. I have never been so traumatized by a television show.
#gameofthrones

Reply · Retweet · Favorite · More

RETWEETS

356

FAVORITES

110



10:06 PM - 2 Jun 2013



how do i convert to

how do i convert to judaism

how do i convert to islam

how do i convert to catholicism

how do i convert to pdf

VIA 9GAG.COM

Press Enter to search.



Justin Bieber

@justinbieber



I make music. I love music.

Reply · Retweet · Favorite · More

RETWEETS

54,213

FAVORITES

59,205



10:09 PM - 7 Apr 2014



dustin curtis

@dcurtis



Follow

"At any moment, Justin Bieber uses 3% of our infrastructure. Racks of servers are dedicated to him. - A guy who works at Twitter

RETWEETS

1,528

FAVORITES

267



8:56 PM - 6 Sep 2010



...



Dmitry Medvedev @MedvedevRussiaE

Follow

The harmonious development of Crimea and Sevastopol as part of our state is one of the main objectives of the Russian Government

[Reply](#) [Retweet](#) [Favorite](#) [More](#)

RETWEETS 144 FAVORITES 57



10:39 AM - 21 Mar 2014



The New York Times
April 2

"Much of the foreign media coverage has distorted the reality of my country and the facts surrounding the events," writes Nicolás Maduro, the president of Venezuela, in Opinion: <http://nyti.ms/1gP5o2I>

[Like](#) · [Comment](#) · [Share](#)

57

262 people like this.

[Top Comments](#) ▾



Elizabeth Warren shared a link.
January 16

I'm not giving up on our fight to extend unemployment benefits. Watch my interview with Now With Alex Wagner about why we need to keep fighting.



Warren: This is the moment to back on economy
www.msnbc.com

President Obama faces one huge problem with his effort to improve the economy: an opposition party

[Like](#) · [Comment](#) · [Share](#)

15,483 720 1,041



Jackie Walorski
@RepWalorski

Follow

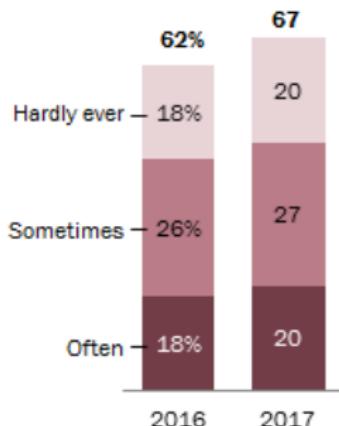
Today, a representative from my office will be meeting with constituents in Goshen. For more details, visit walorski.house.gov/services/upcom...

[Reply](#) [Retweet](#) [Favorite](#) [More](#)

11:22 AM - 8 Apr 2014

In 2017, two-thirds of U.S. adults get news from social media

% of U.S. adults who get news from social media sites ...



Source: Survey conducted Aug. 8-21, 2017.
"News Use Across Social Media Platforms 2017"

PEW RESEARCH CENTER

- ▶ 67% of Americans get news on social media (Pew Research)
- ▶ 58% of EU citizens active on social media & find it useful to get news on national political matters (Eurobarometer, Fall 2017)
- ▶ Social media: top source of news for U.S. young adults (Pew)

Outline

- ▶ Social media data
 - ▶ Motivation
 - ▶ Opportunities and challenges
 - ▶ Twitter data
 - ▶ Facebook data
- ▶ Extracting data from PDF files

Social media data

What are the main advantages of using social media data to study human behavior?

1. **Unobtrusive** data collection at scale, e.g. in study of networks, censorship
2. **Homogeneity** in data format across actors, countries, and over time, e.g. in study of political rhetoric
3. Temporal and spatial data **granularity**, e.g. in study of geographic segregation
4. Increasing **representativeness** of social media users, e.g. in study of political elites

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ Interpersonal networks
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ Social capital and interpersonal communication
 - ▶ Political attitudes and behavior

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ Interpersonal networks
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ Social capital and interpersonal communication
 - ▶ Political attitudes and behavior

Behavior, opinions, and latent traits

- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion



Regular Article

Every tweet counts? How sentiment analysis of social media can improve our knowledge of citizens' political preferences with an application to Italy and France

new media & society

2014, Vol. 16(2) 340–358

© The Author(s) 2013

Reprints and permissions:

sagepub.co.uk/journalsPermissions.nav

DOI: 10.1177/146144813480466

nms.sagepub.com



Andrea Ceron, Luigi Curini, Stefano M Iacus

Università degli Studi di Milano, Italy

Giuseppe Porro

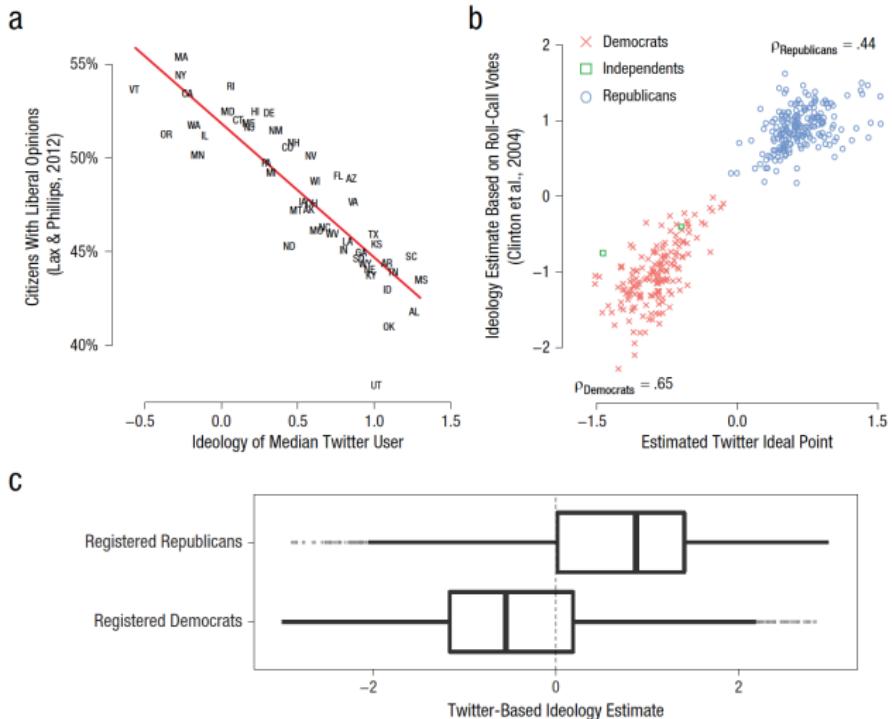
Università degli Studi dell'Insubria, Italy

AJPS

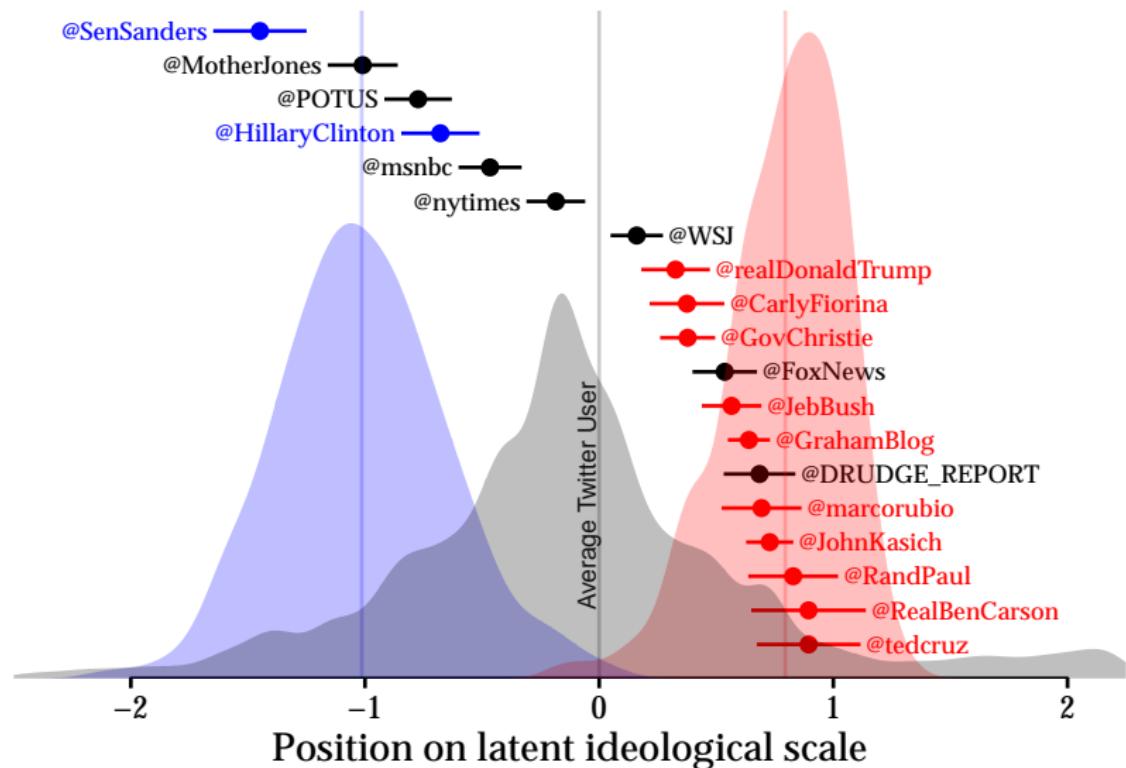
AMERICAN JOURNAL
of POLITICAL SCIENCE

Behavior, opinions, and latent traits

→ Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, . . .



Estimating political ideology using Twitter networks



Barberá “Who is the most conservative Republican candidate for president?” *The Monkey Cage / The Washington Post*, June 16 2015

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ **Interpersonal networks**
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ Social capital and interpersonal communication
 - ▶ Political attitudes and behavior

Interpersonal networks

- ▶ Political behavior is social, strongly influenced by peers

Today is Election Day [What's this?](#) • [close](#)

 Find your polling place on the U.S. Politics Page and click the "I Voted" button to tell your friends you voted.

I Voted

 **f** Jaime Settle, Jason Jones, and 18 other friends have voted.

Bond et al, 2012, “A 61-million-person experiment in social influence and political mobilization”, *Nature*

- ▶ Costly to measure network structure
- ▶ High overlap across online and offline social networks

OPEN  ACCESS Freely available online

PLOS ONE

Inferring Tie Strength from Online Directed Behavior

Jason J. Jones^{1,2*}, Jaime E. Settle², Robert M. Bond², Christopher J. Fariss², Cameron Marlow³, James H. Fowler^{1,2}

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data

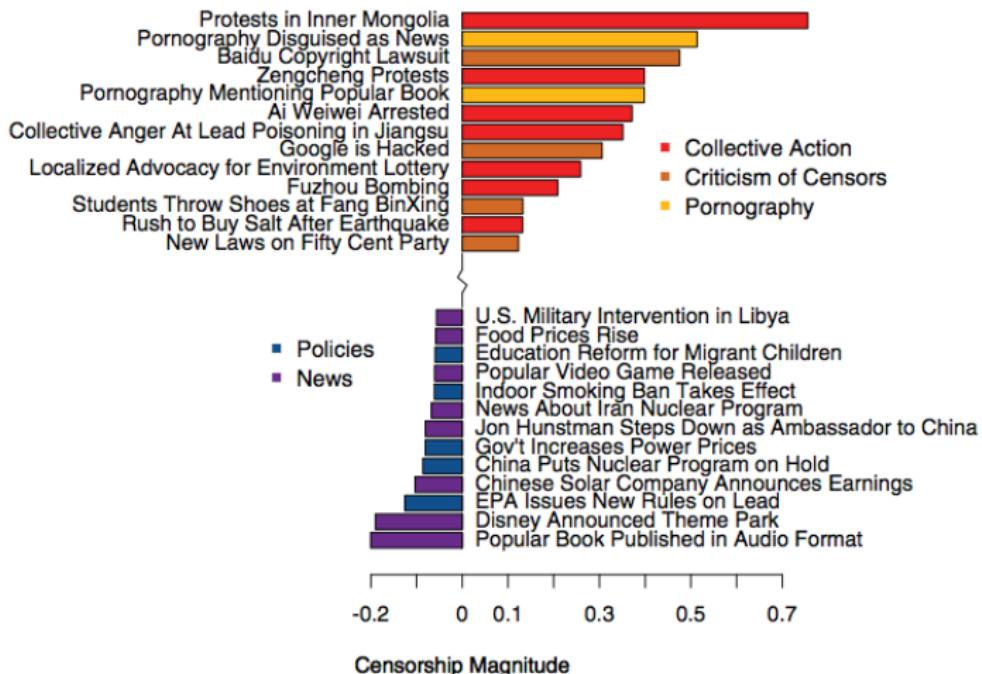
- ▶ Behavior, opinions, and latent traits
- ▶ Interpersonal networks
- ▶ Elite behavior
- ▶ Affordable field experiments

2. How social media affects social behavior

- ▶ Collective action and social movements
- ▶ Political campaigns
- ▶ Social capital and interpersonal communication
- ▶ Political attitudes and behavior

Elite behavior

- Authoritarian governments' response to threat of collective action



King et al, 2013, "How Censorship in China Allows Government Criticism but Silences Collective Expression", *APSR*

- Estimation of conflict intensity in real time

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data

- ▶ Behavior, opinions, and latent traits
- ▶ Interpersonal networks
- ▶ Elite behavior
- ▶ Affordable field experiments

2. How social media affects social behavior

- ▶ Collective action and social movements
- ▶ Political campaigns
- ▶ Social capital and interpersonal communication
- ▶ Political attitudes and behavior

Affordable field experiments



[Political Behavior](#)

September 2017, Volume 39, Issue 3, pp 629–649 | [Cite as](#)

Tweetment Effects on the Tweeted: Experimentally Reducing Racist Harassment

Authors

Kevin Munger

Original Paper

First Online: 11 November 2016

Authors and affiliations

2.7k

Shares

12k

Downloads

3

Citations



13 Sep 2015

@██████████ don't be a n̄ger



Rasheed

@Rasheed██████████

@██████████ Hey man, just remember that there are real people who are hurt when you harass them with that kind of language

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data

- ▶ Behavior, opinions, and latent traits
- ▶ Interpersonal networks
- ▶ Elite behavior
- ▶ Affordable field experiments

2. How social media affects social behavior

- ▶ **Collective action and social movements**
- ▶ Political campaigns
- ▶ Social capital and interpersonal communication
- ▶ Political attitudes and behavior



its
ental
on EU
DEMAND BETTER

BREXIT WRECK IT?
ITALY
SUPPORTS
BRITAIN
TO
REMAIN

TUNBRIDGE WELLS?
BREXIT WRECK IT?
ITALY
SUPPORTS
BRITAIN
TO
REMAIN

FEEL DUTIED

TO VOTE

FOR A MAJOR

CHANCE

FOR ALL

THE PEOPLE

PUT IT TO THE PEOPLE
DEMOCRATIC VOTE

I'M LEAVING
TO DEMAND A
PEOPLES VOTE

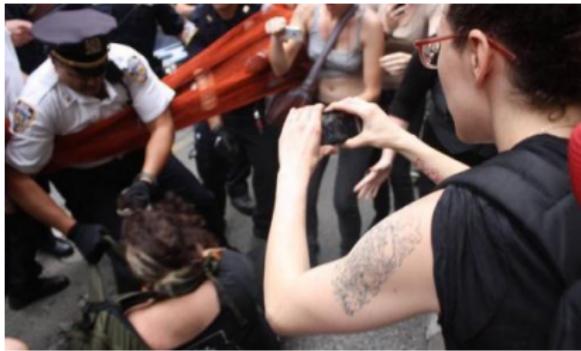
STOP BREXIT NOW
GIVE PEOPLE
THE FIN
I WANT A SAY

**I WANT A SAY
ON BREXIT**
LABOURSY.EU

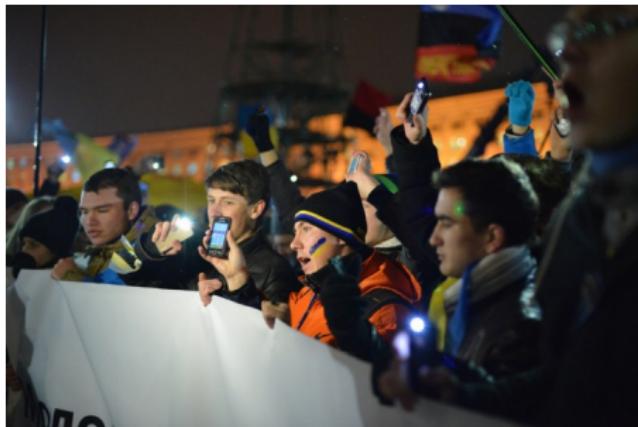
EXHIBITION
SIZZLING AND SHINING



#OccupyGezi



#OccupyWallStreet



#Euromaidan



#Indignados



slacktivism?

Why the revolution will not be tweeted

When the sit-in movement spread from Greensboro throughout the South, it did not spread indiscriminately. It spread to those cities which had preexisting “movement centers” – a core of dedicated and trained activists ready to turn the “fever” into action.

The kind of activism associated with social media isn’t like this at all. [...] Social networks are effective at increasing participation – by lessening the level of motivation that participation requires.

Gladwell, Small Change (New Yorker)

You can’t simply join a revolution any time you want, contribute a comma to a random revolutionary decree, rephrase the guillotine manual, and then slack off for months. Revolutions prize centralization and require fully committed leaders, strict discipline, absolute dedication, and strong relationships.

When every node on the network can send a message to all other nodes, confusion is the new default equilibrium.

Morozov, The Net Delusion: The Dark Side of Internet Freedom

The critical periphery



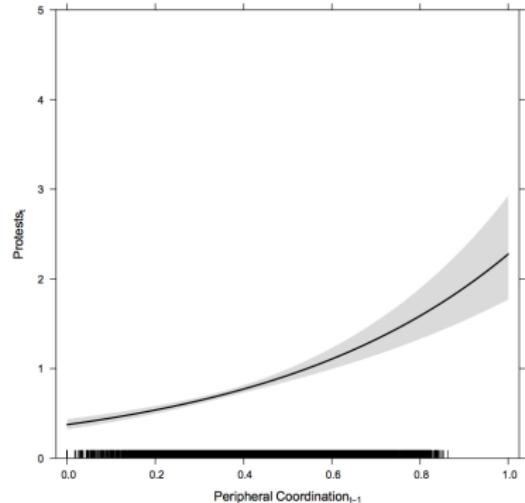
RESEARCH ARTICLE

The Critical Periphery in the Growth of Social Protests

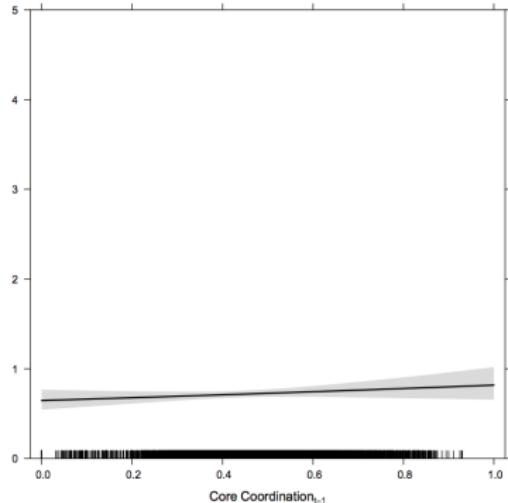
Pablo Barberá^{1*}, Ning Wang², Richard Bonneau^{3,4}, John T. Jost^{1,5,6}, Jonathan Nagler⁶, Joshua Tucker⁶, Sandra González-Bailón^{7*}

- ▶ Structure of online protest networks:
 1. Core: committed minority of resourceful protesters
 2. Periphery: majority of less motivated individuals
- ▶ Our argument: key role of peripheral participants
 1. Increase reach of protest messages (positional effect)
 2. Large contribution to overall activity (size effect)

Peripheral mobilization during the Arab Spring



(a) Increase in protest as peripheral coordination increases



(b) Coordination does not come through core individuals

Steinert-Threlkeld (APSR 2017) "Spontaneous Collective Action"

FROM LIBERATION TO TURMOIL: SOCIAL MEDIA AND DEMOCRACY

*Joshua A. Tucker, Yannis Theocharis, Margaret E. Roberts,
and Pablo Barberá*

"How can one technology – social media – simultaneously give rise to hopes for liberation in authoritarian regimes, be used for repression by these same regimes, and be harnessed by antisystem actors in democracy? We present a simple framework for reconciling these contradictory developments based on two propositions: 1) that social media give voice to those previously excluded from political discussion by traditional media, and 2) that although social media democratize access to information, the platforms themselves are neither inherently democratic nor nondemocratic, but represent a tool political actors can use for a variety of goals, including, paradoxically, illiberal goals."

Journal of Democracy, 2017

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ Interpersonal networks
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ Social capital and interpersonal communication
 - ▶ Political attitudes and behavior



Barack Obama

@BarackObama



Follow

Four more years.



RETWEETS

756,411

FAVORITES

288,867



11:16 PM - 6 Nov 2012

Sections ≡

The Washington Post

Search



Sign In

Post Politics

**By the end of the 2012 campaign,
every Mitt Romney tweet had to be
approved by 22 people**

Political persuasion

Social media as a new campaign tool:

"Let me tell you about Twitter. I think that maybe I wouldn't be here if it wasn't for Twitter. [...] Twitter is a wonderful thing for me, because I get the word out... I might not be here talking to you right now as president if I didn't have an honest way of getting the word out."

Donald Trump, March 16, 2017 (Fox News)

- ▶ Diminished **gatekeeping** role of journalists
 - ▶ Part of a trend towards citizen journalism (Goode, 2009)
- ▶ Information is contextualized within **social layer**
 - ▶ Messing and Westwood (2012): social cues can be as important as partisan cues to explain news consumption through social media
- ▶ **Real-time broadcasting** in reaction to events
 - ▶ e.g. *dual screening* (Vaccari et al, 2015)
- ▶ **Micro-targeting**
 - ▶ Affects how campaigns perceive voters (Hersh, 2015), but unclear if effective in mobilizing or persuading voters

Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ Interpersonal networks
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ **Social capital and interpersonal communication**
 - ▶ Political attitudes and behavior

Social capital

- ▶ Social connections are essential in democratic societies, but online interactions do not facilitate creation and strengthening of social capital (Putnam, 2001)
- ▶ Online networking sites facilitate and transform how social ties are established

Tweeting Alone? An Analysis of Bridging and Bonding Social Capital in Online Networks

American Politics Research

1–31

© The Author(s) 2014

Reprints and permissions:

sagepub.com/journalsPermissions.nav

DOI: 10.1177/1532673X14557942

apr.sagepub.com



**Javier Sajuria¹, Jennifer vanHeerde-Hudson¹,
David Hudson¹, Niheer Dasandi¹, and Yannis
Theocharis²**

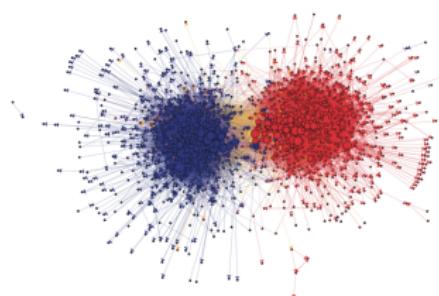
Social media research

Two different approaches in the growing field of social media research:

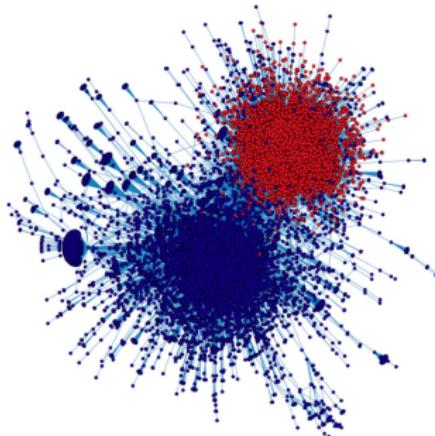
1. Social media as a new source of data
 - ▶ Behavior, opinions, and latent traits
 - ▶ Interpersonal networks
 - ▶ Elite behavior
 - ▶ Affordable field experiments
2. How social media affects social behavior
 - ▶ Collective action and social movements
 - ▶ Political campaigns
 - ▶ Social capital and interpersonal communication
 - ▶ Political attitudes and behavior

Social media as echo chambers?

- ▶ communities of like-minded individuals (homophily, influence)



Adamic and Glance (2005)

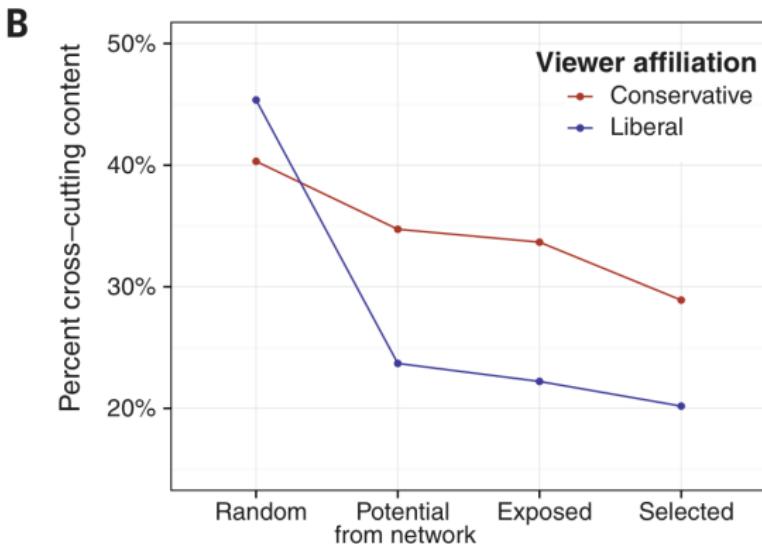


Conover et al (2012)

- ▶ ...generates selective exposure to congenial information
- ▶ ...reinforced by ranking algorithms – “filter bubble” (Parisier)
- ▶ ...increases political polarization (Sunstein, Prior)

Social media as echo chambers?

Fig. 3. Cross-cutting content at each stage in the diffusion process. (A) Illustration of how algorithmic ranking and individual choice affect the proportion of ideologically cross-cutting content that individuals encounter. Gray circles illustrate the content present at each stage in the media exposure process. Red circles indicate conservatives, and blue circles indicate liberals. (B) Average ideological diversity of content (i) shared by random others (random), (ii) shared by friends (potential from network), (iii) actually appeared in users' News Feeds (exposed), and (iv) users clicked on (selected).



Bakshy, Messing, & Adamic (2015) "Exposure to ideologically diverse news and opinion on Facebook". *Science*.

Fake news? Misinformation?



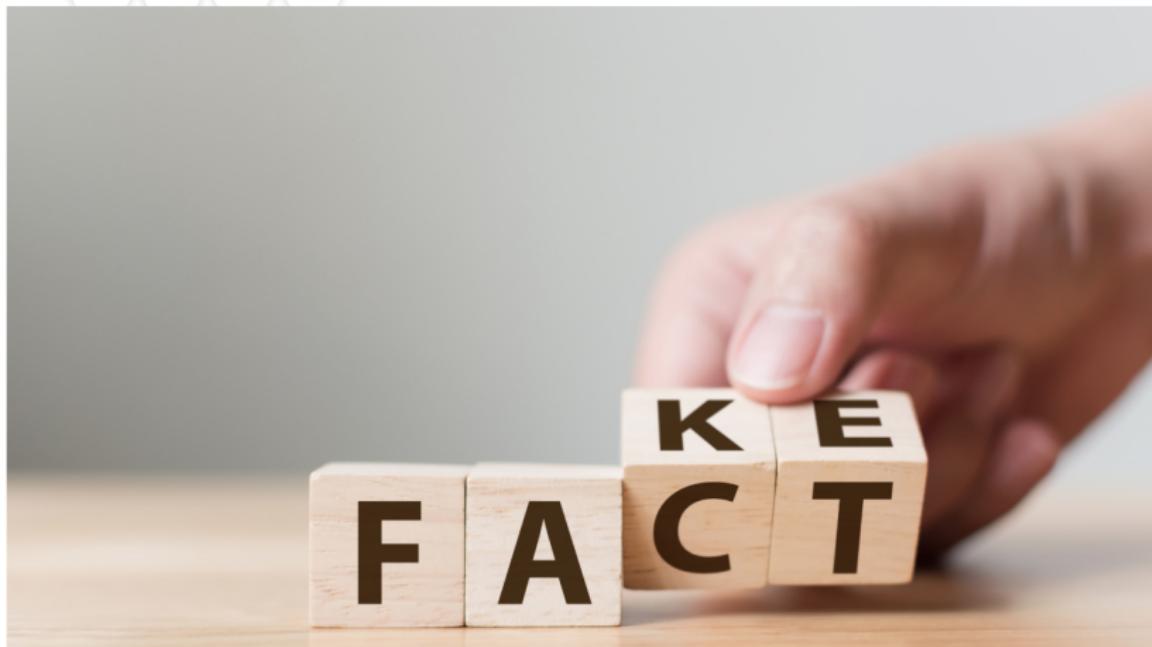
- ▶ Who consumes it? Guess et al (2018, 2019); Grinberg et al (2019)
 - ▶ Web tracking data: 25% Americans visited fake news websites during the 2016 campaigns
 - ▶ High concentration: 1% of users exposed to 80% of fake news
 - ▶ Older, conservative people more likely to be exposed
 - ▶ Facebook key vector of exposure
 - ▶ Fact-check does not reach consumers of misinformation
- ▶ Does it matter? Allcott and Gentzkow (2017):
 - ▶ Survey experiment with real and placebo fake news stories
 - ▶ Most people do not remember seeing fake news stories
 - ▶ Unlikely to affect citizens' behavior

Social Media, Political Polarization, and Political Disinformation: A Review of the Scientific Literature

By Joshua A. Tucker, Andrew Guess, Pablo Barberá, Cristian Vaccari, Alexandra Siegel, Sergey Sanovich, Denis Stukal, and Brendan Nyhan



SHARE



Social media research

Two different approaches in the growing field of social media research:

1. Social media as a new source of data

- ▶ Behavior, opinions, and latent traits
- ▶ Interpersonal networks
- ▶ Elite behavior
- ▶ Affordable field experiments

2. How social media affects social behavior

- ▶ Collective action and social movements
- ▶ Political campaigns
- ▶ Social capital and interpersonal communication
- ▶ Political attitudes and behavior

What are the most important challenges when working with social media data?

Big data, big bias?

SOCIAL SCIENCES

Social media for large studies of behavior

Large-scale studies of human behavior in social media need to be held to higher methodological standards

By Derek Ruths^{1*} and Jürgen Pfeffer²

On 3 November 1948, the day after Harry Truman won the United States presidential elections, the *Chicago Tribune* published one of the most famous erroneous headlines in newspaper history: "Dewey Defeats Truman" (1, 2). The headline was informed by telephone surveys, which had inadver-

different social media platforms (8). For instance, Instagram is "especially appealing to adults aged 18 to 29, African-American, Latinos, women, urban residents" (9) whereas Pinterest is dominated by females, aged 25 to 34, with an average annual household income of \$100,000 (10). These sampling biases are rarely corrected for (if even acknowledged).

Proprietary algorithms for public data. Platform-specific sampling problems, for example, the highest-volume source of pub-

The rise of "embedded researchers who have special relationships with providers that give them access to platform-specific data, algorithms, and resources" is creating a diverse media research community. Such researchers, for example, can see a platform's workings and make accommodations that may not be able to reveal their own or the data used to generate their findings.

Ruths and Pfeffer, 2015, "Social media for large studies of behavior", *Science*

Big data, big bias?

Sources of bias (Ruths and Pfeffer, 2015; Lazer et al, 2017)

- ▶ Population bias
 - ▶ Sociodemographic characteristics are correlated with presence on social media
- ▶ Self-selection within samples
 - ▶ Partisans more likely to post about politics (Barberá & Rivero, 2014)
- ▶ Proprietary algorithms for public data
 - ▶ Twitter API does not always return 100% of publicly available tweets (Morstatter et al, 2014)
- ▶ Human behavior and online platform design
 - ▶ e.g. *Google Flu* (Lazer et al, 2014)

Big data, big bias?

Reducing biases and flaws in social media data

DATA COLLECTION

- 1. Quantifies platform-specific biases (platform design, user base, platform-specific behavior, platform storage policies)
- 2. Quantifies biases of available data (access constraints, platform-side filtering)
- 3. Quantifies proxy population biases/mismatches

METHODS

- 4. Applies filters/corrects for nonhuman accounts in data
- 5. Accounts for platform and proxy population biases
 - a. Corrects for platform-specific and proxy population biases
OR
 - b. Tests robustness of findings
- 6. Accounts for platform-specific algorithms
 - a. Shows results for more than one platform
OR
 - b. Shows results for time-separated data sets from the same platform
- 7. For new methods: compares results to existing methods on the same data
- 8. For new social phenomena or methods or classifiers: reports performance on two or more distinct data sets (one of which was not used during classifier development or design)

Issues in evaluating data from social media. Large-scale social media studies of human behavior should i address issues listed and discussed herein (further discussion in supplementary materials).

Ruths and Pfeffer, 2015, "Social media for large studies of behavior",
Science

Spam and bots



"Follow your coordinators. We need to start tweeting, all at the same time, using the hashtag #ItsTimeForMexico... and don't forget to retweet tweets from the candidate's account..."

Unidentified PRI campaign manager
minutes before the May 8, 2012 Mexican Presidential debate

Spam and bots



Ferrara et al, 2016, *Communications of the ACM*

The privacy paradox

Online data present a paradox in the protection of privacy: Data are at once too revealing in terms of privacy protection, yet also not revealing enough in terms of providing the demographic background information needed by social scientists.

Golder & Macy, Digital footprints, 2014

Ethical concerns

1. Shifting notion of *informed consent*

PNAS

Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer^{a,1}, Jamie E. Guillory^{b,2}, and Jeffrey T. Hancock^{b,c}

^aCore Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and Departments of ^bCommunication and ^cInformation Science, Cornell University, Ithaca, NY 14853

Edited by Susan T. Fiske, Princeton University, Princeton, NJ, and approved March 25, 2014 (received for review October 23, 2013)

Emotional states can be transferred to others via emotional contagion, leading people to experience the same emotions without their awareness. Emotional contagion is well established in laboratory experiments, with people transferring positive and negative emotions to others. Data from a large real-world social network, collected over a 20-y period suggests that longer-lasting moods (e.g., depression, happiness) can be transferred through networks [Fowler JH, Christakis NA (2008) *BMJ* 337:a2338], although the results are controversial. In an experiment with people who use Facebook, we test whether emotional contagion occurs

demonstrated that (i) emotional contagion occurs via text-based computer-mediated communication (7); (ii) contagion of psychological and physiological qualities has been suggested based on correlational data for social networks generally (7, 8); and (iii) people's emotional expressions on Facebook predict friends' emotional expressions, even days later (7) (although some shared experiences may in fact last several days). To date, however, there is no experimental evidence that emotions or moods are contagious in the absence of direct interaction between experiencer and target.

On Facebook, people frequently express emotions, which are

2. Most personal data can be de-anonymized

[Ethics and Information Technology](#)

December 2010, Volume 12, [Issue 4](#), pp 313–325

“But the data is already public”: on the ethics of research in Facebook

Outline

- ▶ Social media data
 - ▶ Motivation
 - ▶ Opportunities and challenges
 - ▶ Twitter data
 - ▶ Facebook data
- ▶ Extracting data from PDF files

Twitter APIs

Two different methods to collect Twitter data:

1. REST API:

- ▶ Queries for specific information about users and tweets
- ▶ Search recent tweets
- ▶ Examples: user profile, list of followers and friends, tweets generated by a given user ("timeline"), users lists, etc.
- ▶ R library: tweetscores (also twitteR, rtweet)

2. Streaming API:

- ▶ Connect to the "stream" of tweets as they are being published
- ▶ Three streaming APIs:
 - 2.1 Filter stream: tweets filtered by keywords
 - 2.2 Geo stream: tweets filtered by location
 - 2.3 Sample stream: 1% random sample of tweets
- ▶ R library: streamR

Important limitation: tweets can only be downloaded in real time
(exception: user timelines, $\sim 3,200$ most recent tweets are available)

Anatomy of a tweet

 **Barack Obama** 
@BarackObama

Four more years.

◀ ▶ ★ ...



RETWEETS FAVORITES
756,411 **288,867**



11:16 PM - 6 Nov 2012

Anatomy of a tweet

Tweets are stored in JSON format:

```
{ "created_at": "Wed Nov 07 04:16:18 +0000 2012",
  "id": 266031293945503744,
  "text": "Four more years. http://t.co/bAJE6Vom",
  "source": "web",
  "user": {
    "id": 813286,
    "name": "Barack Obama",
    "screen_name": "BarackObama",
    "location": "Washington, DC",
    "description": "This account is run by Organizing for Action staff.  
Tweets from the President are signed -bo.",
    "url": "http://t.co/8aJ56Jcemr",
    "protected": false,
    "followers_count": 54873124,
    "friends_count": 654580,
    "listed_count": 202495,
    "created_at": "Mon Mar 05 22:08:25 +0000 2007",
    "time_zone": "Eastern Time (US & Canada)",
    "statuses_count": 10687,
    "lang": "en" },
  "coordinates": null,
  "retweet_count": 756411,
  "favorite_count": 288867,
  "lang": "en"
}
```

Sampling bias?

[Morstatter](#) et al, 2013, *ICWSM*, “Is the Sample Good Enough? Comparing Data from Twitter’s Streaming API with Twitter’s Firehose”:

- ▶ 1% random sample from Streaming API is not truly random
- ▶ Less popular hashtags, users, topics... less likely to be sampled
- ▶ But for keyword-based samples, bias is not as important

[González-Bailón](#) et al, 2014, *Social Networks*, “Assessing the bias in samples of large online networks”:

- ▶ Small samples collected by filtering with a subset of relevant hashtags can be biased
- ▶ Central, most active users are more likely to be sampled
- ▶ Data collected via search (REST) API more biased than those collected with Streaming API

Outline

- ▶ Social media data
 - ▶ Motivation
 - ▶ Opportunities and challenges
 - ▶ Twitter data
 - ▶ Facebook data
- ▶ Extracting data from PDF files

Collecting Facebook data

Facebook used to allow access to public pages' data through the [Graph API](#):

1. Posts on public pages and groups
2. Likes, reactions, comments, replies...

Currently not available.

Aggregate-level statistics available through the FB Marketing API.
See the code by Connor Gilroy (UW)

Access to other (anonymized) data used in published studies
requires permission from Facebook or from users.

Social Science One as a new model for academic partnerships
with Facebook.

Outline

- ▶ Social media data
 - ▶ Motivation
 - ▶ Opportunities and challenges
 - ▶ Twitter data
 - ▶ Facebook data
- ▶ Extracting data from PDF files