

# MY560 Workshop: Collecting and Analyzing Social Media Data

**Pablo Barberá**

London School of Economics

[www.pablobarbera.com](http://www.pablobarbera.com)

Workshop website:

[pablobarbera.com/social-media-workshop](http://pablobarbera.com/social-media-workshop)







George Takei

March 28 at 10:10pm ·

Who's with me.



Like · Comment · Share

408,735 people like this.

66,990 shares



George Takei

March 28 at 10:10pm ·

Who's with me.



Like · Comment · Share

408,735 people like this.

66,990 shares



Bon Alimago  
@karma\_thief

Follow

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



George Takei

March 28 at 10:10pm ·

Who's with me.



Like · Comment · Share

408,735 people like this.

66,990 shares



Bon Alimago  
@karma\_thief

Follow

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

Reply Retweet Favorite More

RETWEETS 356 FAVORITES 110

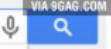


10:06 PM - 2 Jun 2013



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**



Press Enter to search.

VIA 9GAG.COM



George Takei

March 28 at 10:10pm ·

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**

More

RETWEETS 356    FAVORITES 110



10:06 PM - 2 Jun 2013



Google

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.

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



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



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](http://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



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](http://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...](http://walorski.house.gov/services/upcom...)

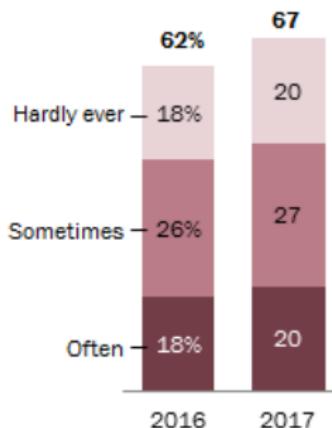
Reply Retweet Favorite More

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 ...*



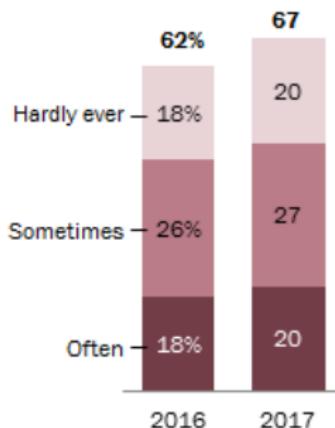
- ▶ 62% of Americans get news on social media (Pew)

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

---

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

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



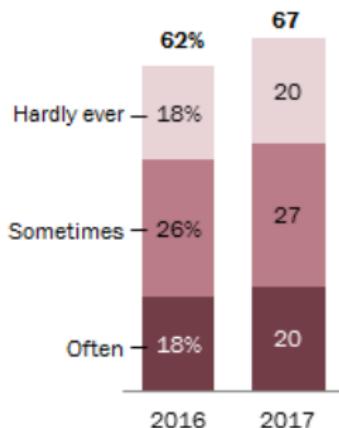
- ▶ 62% of Americans get news on social media (Pew)
- ▶ 27% of online EU citizens use social media to get news on national political matters (Eurobarometer, Fall 2012)

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

---

## 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

---

- ▶ 62% of Americans get news on social media (Pew)
- ▶ 27% of online EU citizens use social media to get news on national political matters (Eurobarometer, Fall 2012)
- ▶ Social media: top source of news for U.S. young adults (Pew)



## Shift in communication patterns



**Shift in communication patterns**



**Digital footprints of human behavior**

Hello!

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability
  - ▶ Social network analysis, Bayesian statistics, text as data methods

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability
  - ▶ Social network analysis, Bayesian statistics, text as data methods
  - ▶ Author of R packages to analyze data from social media

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability
  - ▶ Social network analysis, Bayesian statistics, text as data methods
  - ▶ Author of R packages to analyze data from social media
- ▶ Contact:

## About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability
  - ▶ Social network analysis, Bayesian statistics, text as data methods
  - ▶ Author of R packages to analyze data from social media
- ▶ Contact:
  - ▶ P.Barbera@lse.ac.uk

# About me

- ▶ Assistant Professor of Computational Social Science in the Methodology Department at LSE
- ▶ Previously Assistant Prof. at Univ. of Southern California
- ▶ PhD in Politics, New York University (2015)
- ▶ Data Science Fellow at NYU, 2015–2016
- ▶ My research:
  - ▶ Social media and politics, comparative electoral behavior, corruption and accountability
  - ▶ Social network analysis, Bayesian statistics, text as data methods
  - ▶ Author of R packages to analyze data from social media
- ▶ Contact:
  - ▶ P.Barbera@lse.ac.uk
  - ▶ www.pablobarbera.com

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

Session 2, 14–16:00

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

Session 2, 14–16:00

- ▶ Guided coding session: Collecting Twitter data from the REST API

# Today's workshop

Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

Session 2, 14–16:00

- ▶ Guided coding session: Collecting Twitter data from the REST API
- ▶ Coding challenge 2: Twitter's REST API

# Today's workshop

## Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

## Session 2, 14–16:00

- ▶ Guided coding session: Collecting Twitter data from the REST API
- ▶ Coding challenge 2: Twitter's REST API
- ▶ Guided coding session: Collecting Facebook data from the Graph API

# Today's workshop

## Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

## Session 2, 14–16:00

- ▶ Guided coding session: Collecting Twitter data from the REST API
- ▶ Coding challenge 2: Twitter's REST API
- ▶ Guided coding session: Collecting Facebook data from the Graph API
- ▶ Application: Dictionary methods applied to social media

# Today's workshop

## Session 1, 10–12:00

- ▶ Social media research: opportunities and challenges
- ▶ Guided coding session: collecting Twitter data from the Streaming API
- ▶ Challenge 1: interacting with Twitter's Streaming API

## Session 2, 14–16:00

- ▶ Guided coding session: Collecting Twitter data from the REST API
- ▶ Coding challenge 2: Twitter's REST API
- ▶ Guided coding session: Collecting Facebook data from the Graph API
- ▶ Application: Dictionary methods applied to social media
- ▶ Coding challenge 3: Facebook's Graph API

# Social media research

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

1. Social media as a new source of information
  - ▶ Behavior, opinions, and latent traits
  - ▶ Interpersonal networks
  - ▶ Elite behavior
  - ▶ Affordable field experiments

# Social media research

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

1. Social media as a new source of information
  - ▶ 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 information
  - ▶ 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, ...

## Behavior, opinions, and latent traits

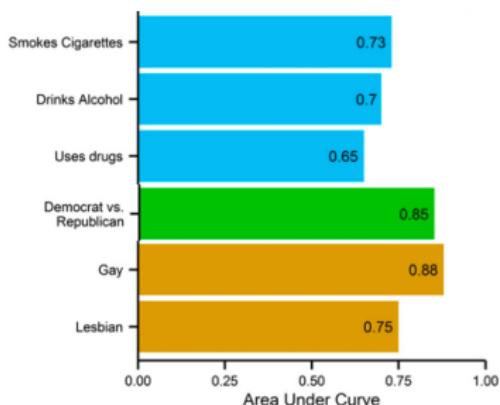
- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion  
Beauchamp (AJPS 2016): “Predicting and Interpolating State-level Polls using Twitter Textual Data”

## Behavior, opinions, and latent traits

- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion
- Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, ...

# Behavior, opinions, and latent traits

- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion
- Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, ...



Kosinski et al, 2013, "Private traits and attributes are predictable from digital records of human behavior", *PNAS* (also personality, *PNAS* 2015)

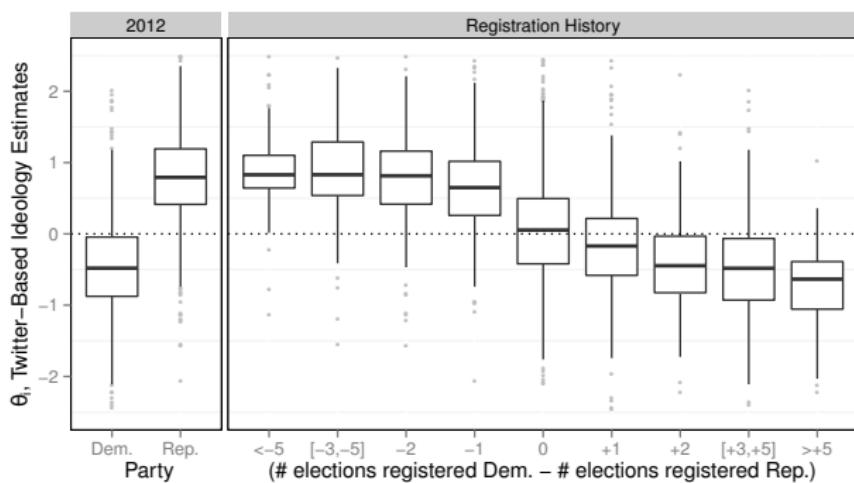
Fig. 2. Prediction accuracy of classification for dichotomous/dichotomized attributes expressed by the AUC.

## Behavior, opinions, and latent traits

- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion
- Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, ...

# Behavior, opinions, and latent traits

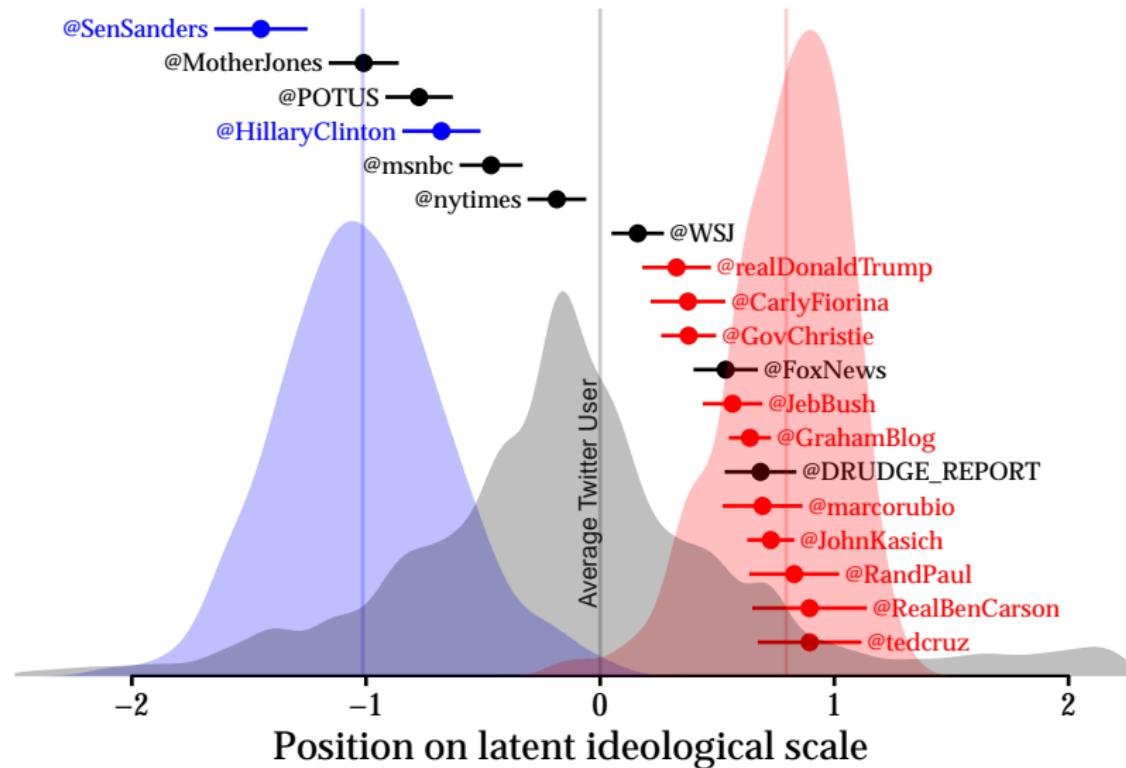
- ▶ Digital footprints: check-ins, conversations, geolocated pictures, likes, shares, retweets, ...
- Non-intrusive measurement of behavior and public opinion
- Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, ...



Data: 2,360 Twitter accounts, matched with Ohio voter file.

Barberá, 2015, "Birds of the Same Feather Tweet Together. Bayesian Ideal Point Estimation Using Twitter Data", *Political Analysis*

# 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 information
  - ▶ 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.

0 1 1 5 5 3 7 6  
People on Facebook 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*

## Interpersonal networks

- ▶ Political behavior is social, strongly influenced by peers
- ▶ Costly to measure network structure

# Interpersonal networks

- ▶ Political behavior is social, strongly influenced by peers
- ▶ Costly to measure network structure
- ▶ High overlap across online and offline social networks

OPEN  ACCESS Freely available online



## Inferring Tie Strength from Online Directed Behavior

Jason J. Jones<sup>1,2\*</sup>, Jaime E. Settle<sup>2</sup>, Robert M. Bond<sup>2</sup>, Christopher J. Fariss<sup>2</sup>, Cameron Marlow<sup>3</sup>, James H. Fowler<sup>1,2</sup>

**1** Medical Genetics Division, University of California, San Diego, La Jolla, California, United States of America, **2** Political Science Department, University of California, San Diego, La Jolla, California, United States of America, **3** Data Science, Facebook, Inc., Menlo Park, California, United States of America

### Abstract

Some social connections are stronger than others. People have not only friends, but also best friends. Social scientists have long recognized this characteristic of social connections and researchers frequently use the term *tie strength* to refer to this concept. We used online interaction data (specifically, Facebook interactions) to successfully identify real-world strong ties. Ground truth was established by asking users themselves to name their closest friends in real life. We found the frequency of online interaction was diagnostic of strong ties, and interaction frequency was much more useful diagnostically than were attributes of the user or the user's friends. More private communications (messages) were not necessarily more informative than public communications (comments, wall posts, and other interactions).

Jones et al, 2013, “Inferring Tie Strength from Online Directed Behavior”, *PLOS One*

# Interpersonal networks

- ▶ Political behavior is social, strongly influenced by peers
- ▶ Costly to measure network structure
- ▶ High overlap across online and offline social networks
- ▶ Online and offline ties are similar in nature

The screenshot shows the homepage of the American Political Science Review (APSR) website. At the top, the journal's name "American Political Science Review" is displayed next to a thumbnail image of the journal cover. To the right is the logo for "apsa" (American Political Science Association). Below the header, there are three navigation tabs: "Article" (which is underlined in blue), "Supplementary materials", and "Metrics". A large, light gray search bar is positioned below these tabs. Underneath the search bar, the text "Volume 111, Issue 3 August 2017, pp. 502-521" is visible. The main content area features the title of the article: "Testing Social Science Network Theories with Online Network Data: An Evaluation of External Validity" by JAMES BISBEE (a1) and JENNIFER M. LARSON (a1). Below the title, the DOI link "https://doi.org/10.1017/S0003055417000120" and the publication date "Published online: 13 June 2017" are provided.

American Political Science Review

apsa  
AMERICAN POLITICAL SCIENCE ASSOCIATION

Article    Supplementary materials    Metrics

Volume 111, Issue 3 August 2017, pp. 502-521

**Testing Social Science Network Theories with Online Network Data: An Evaluation of External Validity**

JAMES BISBEE (a1) and JENNIFER M. LARSON (a1)

<https://doi.org/10.1017/S0003055417000120> Published online: 13 June 2017

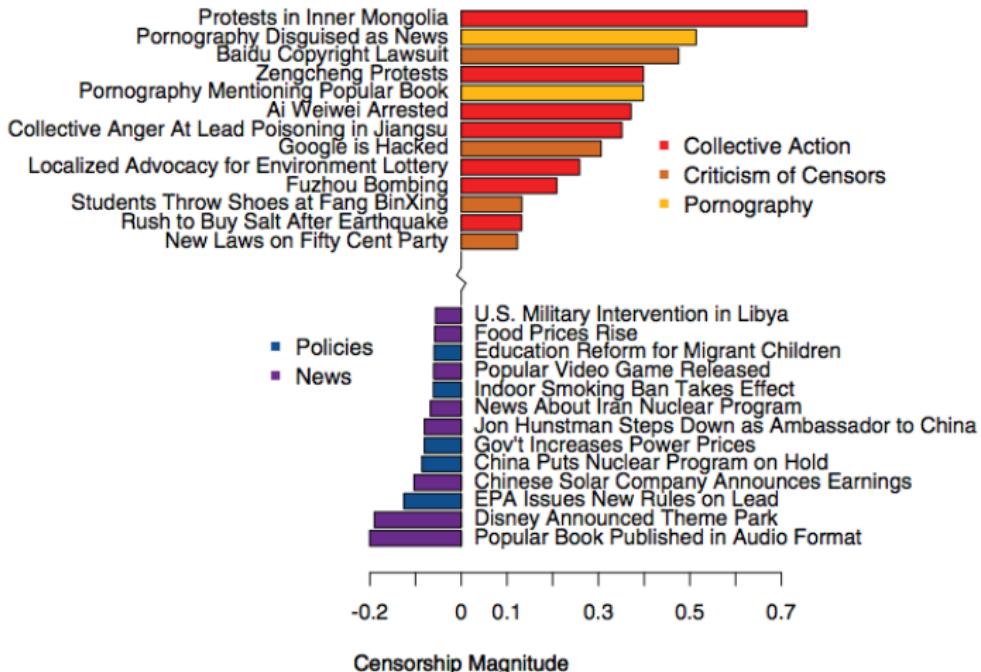
# Social media research

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

1. Social media as a new source of information
  - ▶ 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*

## Elite behavior

- ▶ Authoritarian governments' response to threat of collective action
- ▶ Estimation of conflict intensity in real time

---

Journal of Conflict Resolution  
55(6) 938-969

© The Author(s) 2011

Reprints and permission:

sagepub.com/journalsPermissions.nav

DOI: 10.1177/0022002711408014

<http://jcr.sagepub.com>



# Using Social Media to Measure Conflict Dynamics: An Application to the 2008–2009 Gaza Conflict

Thomas Zeitzoff<sup>1</sup>

# Elite behavior

- ▶ Authoritarian governments' response to threat of collective action
- ▶ Estimation of conflict intensity in real time
- ▶ How elected officials communicate with constituents

FEBRUARY 23, 2017



## For members of 114th Congress, partisan criticism ruled on Facebook



Facebook posts from members of the 114th Congress attracted more attention when they contained disagreement with the opposing party than when they expressed bipartisanship, according to a Pew Research Center analysis of over 100,000 posts.

# Social media research

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

1. Social media as a new source of information
  - ▶ 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 [Authors and affiliations](#)

Kevin Munger [✉](#)

Original Paper  
First Online: 11 November 2016

 2.7k Shares

 12k Downloads

 3 Citations

# Social media research

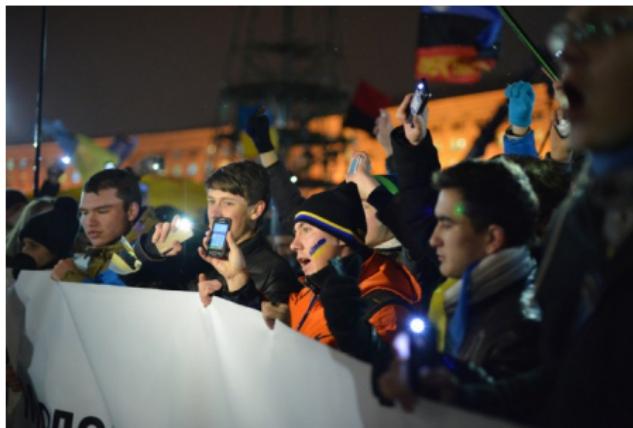
Two different approaches in the growing field of social media research:

1. Social media as a new source of information
  - ▶ 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





#OccupyGezi



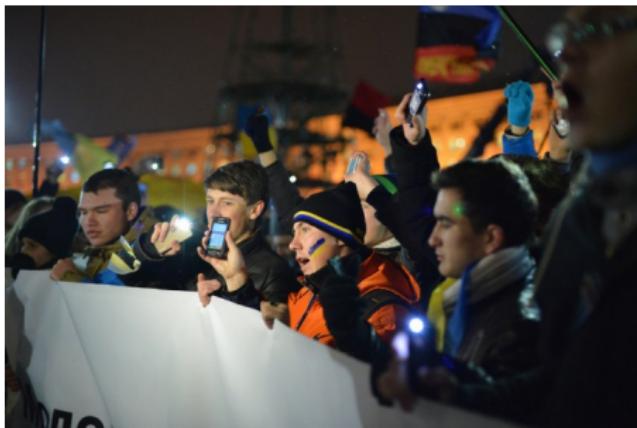
#Euromaidan



#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)**

# 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<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ Structure of online protest networks:

# The critical periphery



RESEARCH ARTICLE

## The Critical Periphery in the Growth of Social Protests

Pablo Barberá<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ Structure of online protest networks:
  1. Core: committed minority of resourceful protesters

# The critical periphery



RESEARCH ARTICLE

## The Critical Periphery in the Growth of Social Protests

Pablo Barberá<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ Structure of online protest networks:
  1. **Core**: committed minority of resourceful protesters
  2. **Periphery**: majority of less motivated individuals

# The critical periphery



RESEARCH ARTICLE

## The Critical Periphery in the Growth of Social Protests

Pablo Barberá<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ 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

# The critical periphery



RESEARCH ARTICLE

## The Critical Periphery in the Growth of Social Protests

Pablo Barberá<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ 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)

# The critical periphery



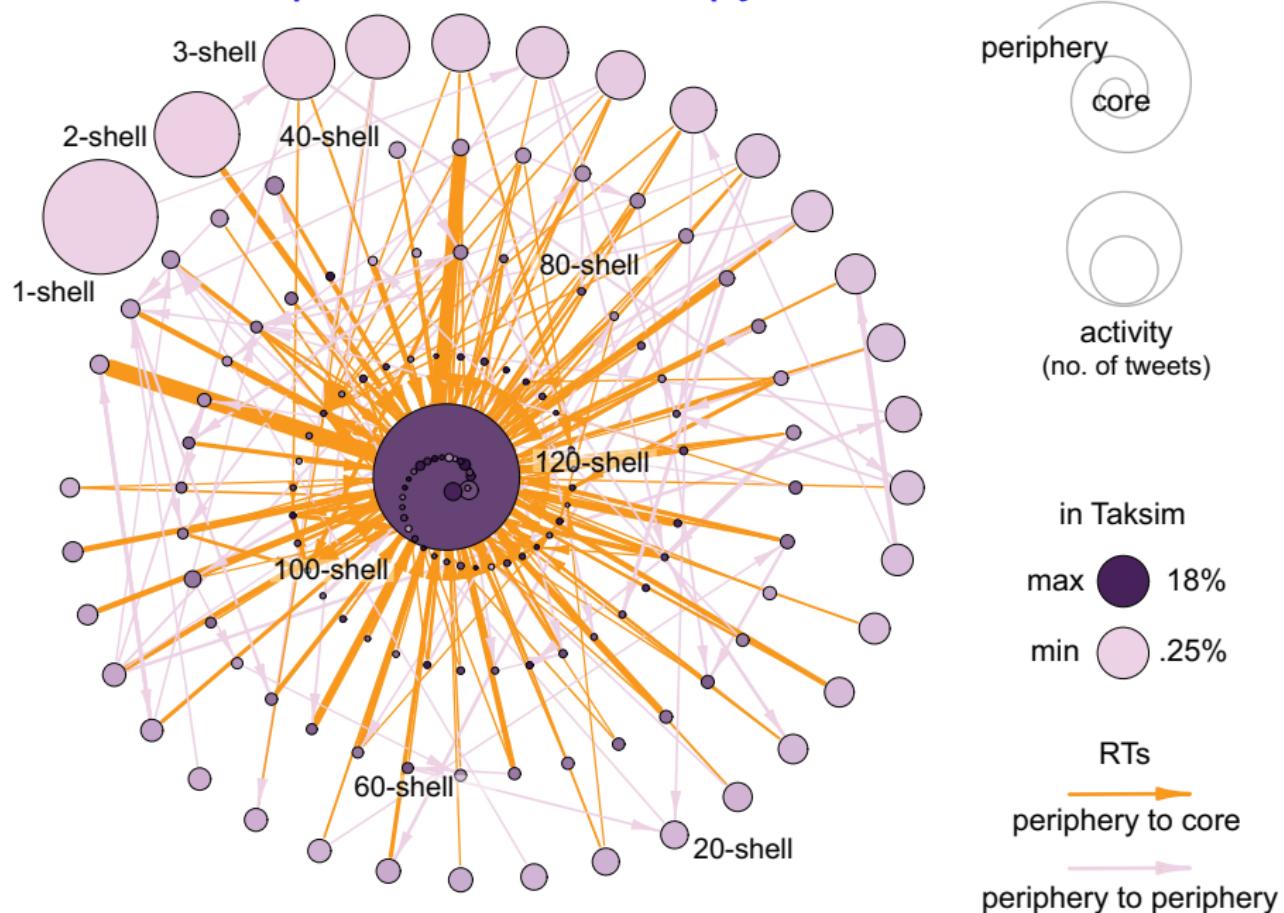
RESEARCH ARTICLE

## The Critical Periphery in the Growth of Social Protests

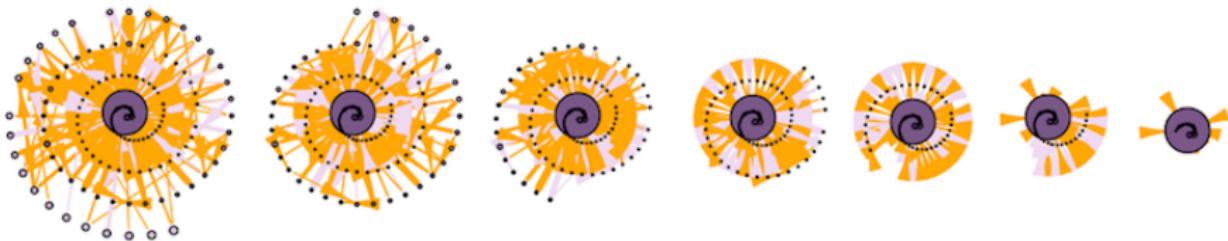
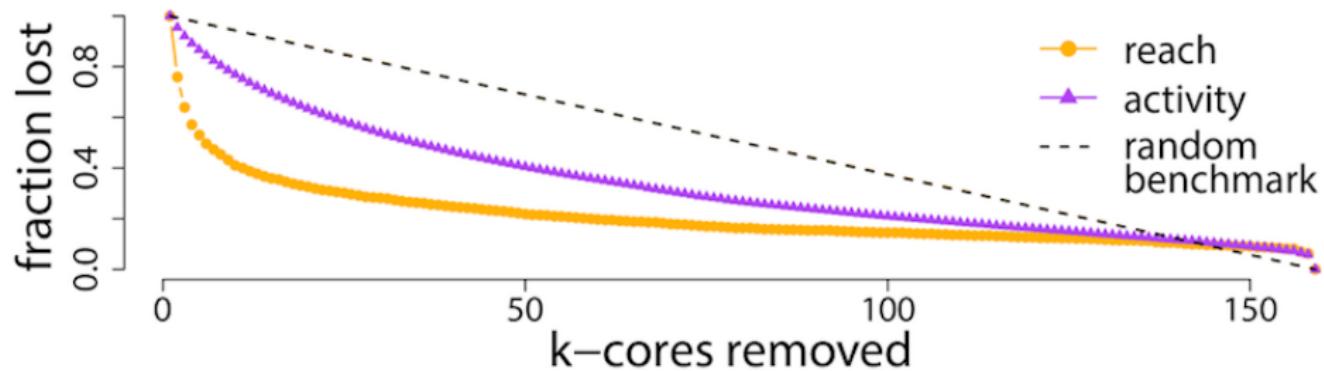
Pablo Barberá<sup>1\*</sup>, Ning Wang<sup>2</sup>, Richard Bonneau<sup>3,4</sup>, John T. Jost<sup>1,5,6</sup>, Jonathan Nagler<sup>6</sup>, Joshua Tucker<sup>6</sup>, Sandra González-Bailón<sup>7\*</sup>

- ▶ 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)

# k-core decomposition of #OccupyGezi network



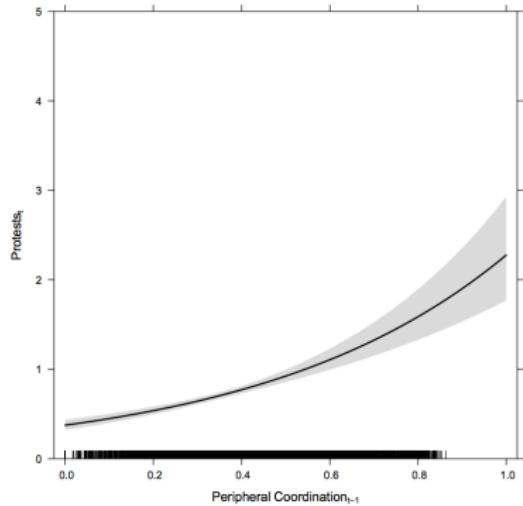
## Relative importance of core and periphery



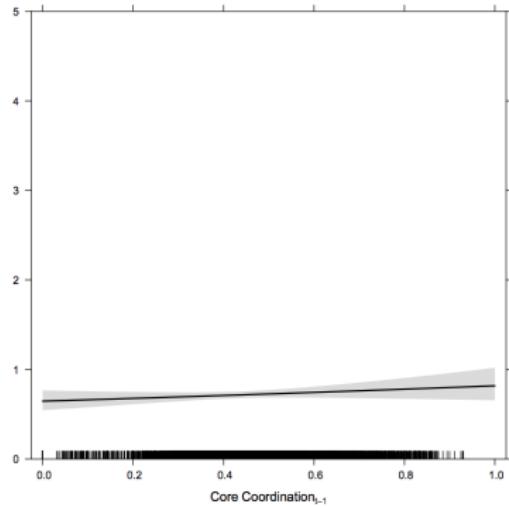
reach: aggregate size of participants' audience

activity: total number of protest messages published (not only RTs)

# 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”

# Social media research

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

1. Social media as a new source of information
  - ▶ 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)**

# 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

# 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)

# 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**

# 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

# 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

# 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)

# 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**

# 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 information
  - ▶ 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)

## 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

# 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](http://sagepub.com/journalsPermissions.nav)

DOI: 10.1177/1532673X14557942

[apr.sagepub.com](http://apr.sagepub.com)



**Javier Sajuria<sup>1</sup>, Jennifer vanHeerde-Hudson<sup>1</sup>,  
David Hudson<sup>1</sup>, Niheer Dasandi<sup>1</sup>, and Yannis  
Theocharis<sup>2</sup>**

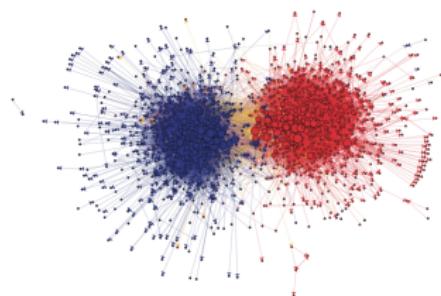
# Social media research

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

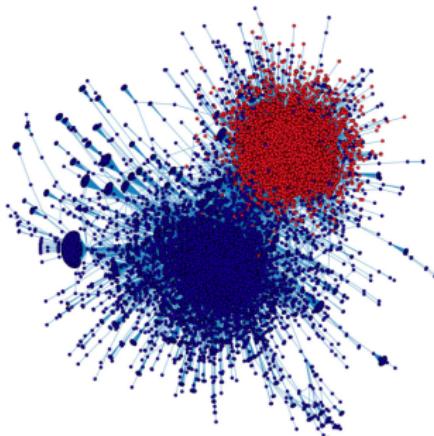
1. Social media as a new source of information
  - ▶ Behavior, opinions, and latent traits
  - ▶ Interpersonal networks
  - ▶ Elite behavior
  - ▶ Affordable field experiments
2. How social media affects social behavior
  - ▶ Collective action and social movements
  - ▶ Political campaigns
  - ▶ Social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

# Social media as echo chambers?

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



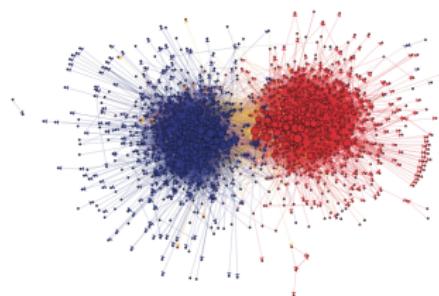
Adamic and Glance (2005)



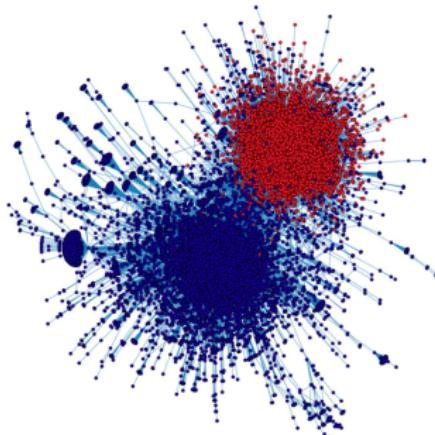
Conover et al (2012)

# Social media as echo chambers?

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



Adamic and Glance (2005)

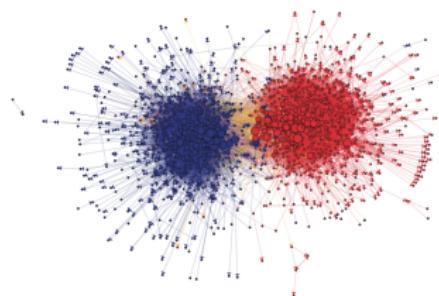


Conover et al (2012)

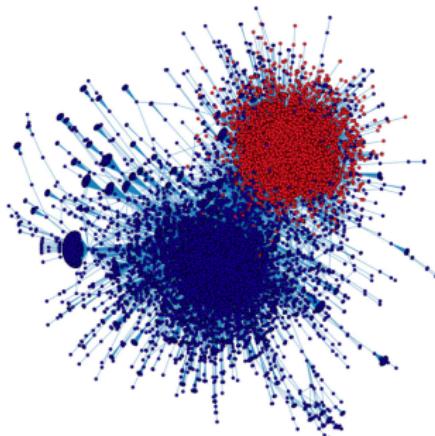
- ▶ ...generates selective exposure to congenial information
- ▶ ...reinforced by ranking algorithms – “filter bubble” (Parisier)

# Social media as echo chambers?

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



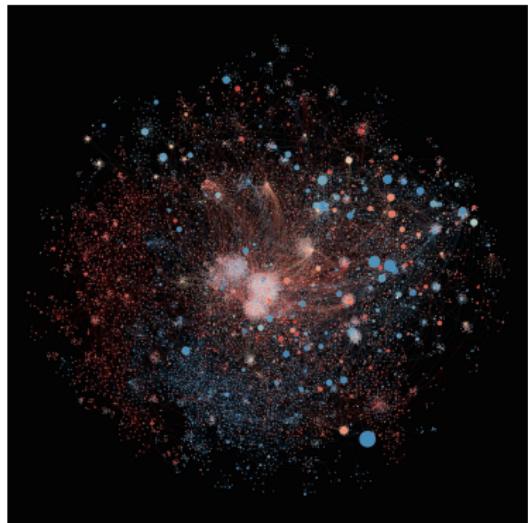
Adamic and Glance (2005)



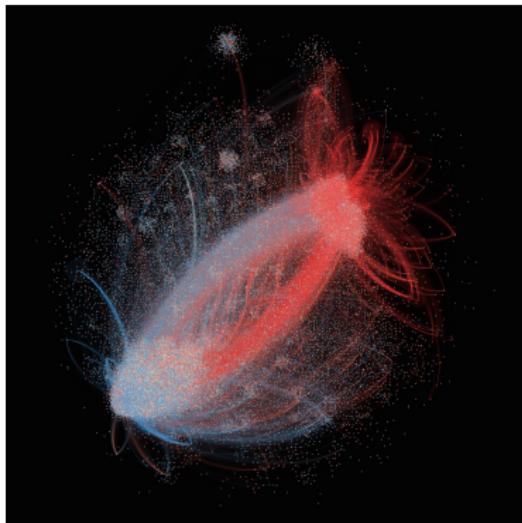
Conover et al (2012)

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

# Social media as echo chambers?



2013 SuperBowl



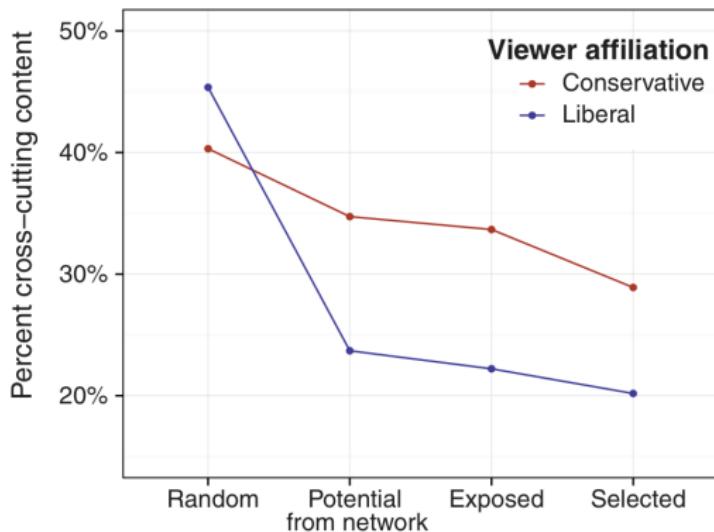
2012 Election

Barberá et al (2015) "Tweeting From Left to Right: Is Online Political Communication More Than an Echo Chamber?" *Psychological Science*

# 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).

B



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

# Social media research

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

1. Social media as a new source of information
  - ▶ 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?

# Social media data and social science: challenges

1. Big data, big bias?

# Social media data and social science: challenges

1. Big data, big bias?
2. The end of theory?

# Social media data and social science: challenges

1. Big data, big bias?
2. The end of theory?
3. Spam and bots

# Social media data and social science: challenges

1. Big data, big bias?
2. The end of theory?
3. Spam and bots
4. The privacy paradox

# Social media data and social science: challenges

1. Big data, big bias?
2. The end of theory?
3. Spam and bots
4. The privacy paradox
5. Generalizing from online to offline behavior

# Social media data and social science: challenges

1. Big data, big bias?
2. The end of theory?
3. Spam and bots
4. The privacy paradox
5. Generalizing from online to offline behavior
6. Ethical concerns

# 1. 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<sup>1\*</sup> and Jürgen Pfeffer<sup>2</sup>

**O**n 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 commercial 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

# 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

# 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

# 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)

# 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

# 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)

# 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

# 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)

# 1. 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*

## 2. The end of theory?

*Petabytes allow us to say: “Correlation is enough.” We can stop looking for models. We can analyze the data without hypotheses about what it might show. We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot.*

**Chris Anderson**, [Wired](#), June 2008

## 2. The end of theory?

*Petabytes allow us to say: "Correlation is enough." We can stop looking for models. We can analyze the data without hypotheses about what it might show. We can throw the numbers into the biggest computing clusters the world has ever seen and let statistical algorithms find patterns where science cannot.*

**Chris Anderson**, *Wired*, June 2008

*Correlations are a way of catching a scientist's attention, but the models and mechanisms that explain them are how we make the predictions that not only advance science, but generate practical applications.*

**John Timmer**, *Ars Technica*, June 2008

(Big) social media data as a complement - not a substitute - for theoretical work and careful causal inference.

### 3. 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*

### 3. Spam and bots



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

## 4. 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*

## 5. Generalizing from online to offline behavior

What makes online behavior different:

- ▶ Platform affordances may distort behavior

## 5. Generalizing from online to offline behavior

What makes online behavior different:

- ▶ Platform affordances may distort behavior
- ▶ Tools extend innate capacities (e.g. Dunbar's number)

## 5. Generalizing from online to offline behavior

What makes online behavior different:

- ▶ Platform affordances may distort behavior
- ▶ Tools extend innate capacities (e.g. Dunbar's number)
- ▶ Anonymity encourages vitriol

# 6. Ethical concerns

## 1. Shifting notion of *informed consent*

PNAS

### Experimental evidence of massive-scale emotional contagion through social networks

Adam D. I. Kramer<sup>a,1</sup>, Jamie E. Guillory<sup>b,2</sup>, and Jeffrey T. Hancock<sup>b,c</sup>

<sup>a</sup>Core Data Science Team, Facebook, Inc., Menlo Park, CA 94025; and Departments of <sup>b</sup>Communication and <sup>c</sup>Information 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

## 6. Ethical concerns

1. Shifting notion of *informed consent*
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

Authors

[Authors and affiliations](#)

Michael Zimmer 

Article

First Online: 04 June 2010

DOI: [10.1007/s10676-010-9227-5](https://doi.org/10.1007/s10676-010-9227-5)

Cite this article as:

Zimmer, M. Ethics Inf Technol (2010) 12:  
313. doi:10.1007/s10676-010-9227-5

144

27

8.3k

Citations

Shares

Downloads

## 6. Ethical concerns

1. Shifting notion of *informed consent*
2. Most personal data can be de-anonymized
3. Inequalities in data access

## 6. Ethical concerns

1. Shifting notion of *informed consent*
2. Most personal data can be de-anonymized
3. Inequalities in data access

“Ethical concerns must be weighed against the value of social research with appropriate steps taken to protect individual privacy” (Shah et al, 2015)

# Twitter data

# Twitter APIs

Two different methods to collect Twitter data:

1. REST API:

# Twitter APIs

Two different methods to collect Twitter data:

1. REST API:

- ▶ Queries for specific information about users and tweets

# Twitter APIs

Two different methods to collect Twitter data:

## 1. REST API:

- ▶ Queries for specific information about users and tweets
- ▶ Search recent tweets

# 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.

# 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)

# 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:

# 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

# 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:

# 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

# 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

# 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

# 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

# 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, ~ 3,200 most recent tweets are available)

# Anatomy of a tweet



Barack Obama

@BarackObama



Follow

Four more years.



RETWEETS

**756,411**

FAVORITES

**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"
}
```

## Streaming API

- ▶ Recommended method to collect tweets

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:

## Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample

## Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.
- ▶ My workflow:

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.
- ▶ My workflow:
  - ▶ Amazon EC2, cloud computing

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.
- ▶ My workflow:
  - ▶ Amazon EC2, cloud computing
  - ▶ Cron jobs to restart R scripts every hour.

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.
- ▶ My workflow:
  - ▶ Amazon EC2, cloud computing
  - ▶ Cron jobs to restart R scripts every hour.
  - ▶ Save tweets in .json files, one per day.

# Streaming API

- ▶ Recommended method to collect tweets
- ▶ Potential issues:
  - ▶ Filter streams have same rate limit as spritzer: when volume reaches 1% of all tweets, it will return random sample
  - ▶ Stream connections tend to die spontaneously. Restart regularly.
- ▶ My workflow:
  - ▶ Amazon EC2, cloud computing
  - ▶ Cron jobs to restart R scripts every hour.
  - ▶ Save tweets in .json files, one per day.
  - ▶ Will show some examples later

## 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

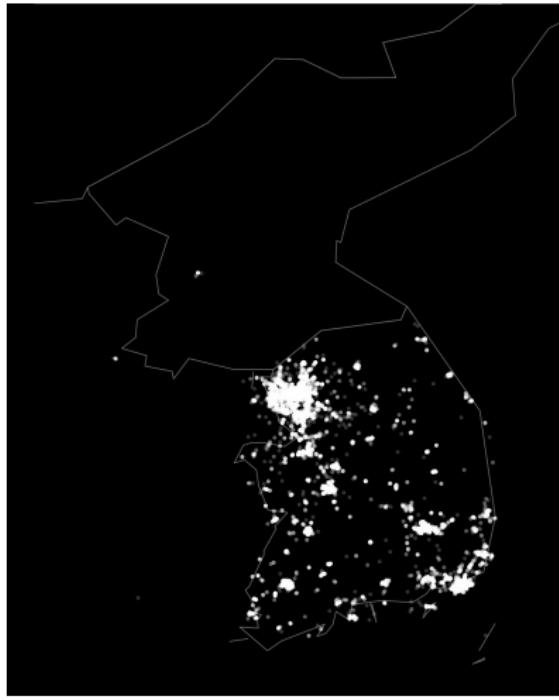
## Sampling bias?

[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



Tweets from Korea: 40k tweets collected in 2014 (left)  
Korean peninsula at night, 2003 (right). Source: NASA.

# Who is tweeting from North Korea?



**North Korea English**  
@uriminzok\_engl  
An English translation of @uriminzok - the official North Korea Twitter feed  
[uriminzokkiri.com](http://uriminzokkiri.com)

671 TWEETS    940 FOLLOWING    129 FOLLOWERS

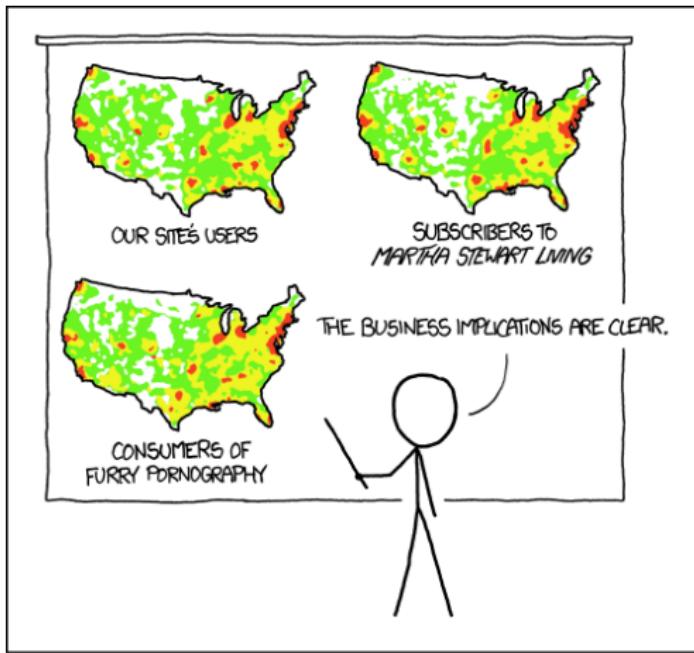
[Follow](#)

### Tweets

 **North Korea English** @uriminzok\_engl 13h  
Beloved Comrade Kim Jong-eun to stay in the national light industry competition attended by Code speeches do was [goo.gl/eJWsJ](http://goo.gl/eJWsJ)  
[Expand](#)

Twitter user: [@uriminzok\\_engl](#)

But remember...



PET PEEVE #208:  
GEOGRAPHIC PROFILE MAPS WHICH ARE  
BASICALLY JUST POPULATION MAPS

# Facebook data

## Collecting Facebook data

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

1. Posts on public pages and groups

## Collecting Facebook data

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

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

## Collecting Facebook data

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

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

Some public user data (gender, location) was available through previous versions of the API (not anymore)

## Collecting Facebook data

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

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

Some public user data (gender, location) was available through previous versions of the API (not anymore)

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

## Collecting Facebook data

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

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

Some public user data (gender, location) was available through previous versions of the API (not anymore)

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

## Collecting Facebook data

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

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

Some public user data (gender, location) was available through previous versions of the API (not anymore)

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

R library: [Rfacebook](#)

## Login details: RStudio Server

RStudio Server URL:

`rstudio.pablobarbera.com`

user = `userXX` and password = `passwordXX`

where XX is your assigned number