

# Social Media

Kenneth Benoit & Pablo Barberá

MY 459: Quantitative Text Analysis

March 19th, 2018

Course website: [lse-my459.github.io](https://lse-my459.github.io)

# Outline

- ▶ Social media data
  - ▶ Motivation
  - ▶ Opportunities and challenges
  - ▶ APIs
  - ▶ Twitter data
  - ▶ Facebook data

## 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

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

VIA 9GAG.COM

Press Enter to search.



Justin Bieber

@justinbieber

Follow

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](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...)

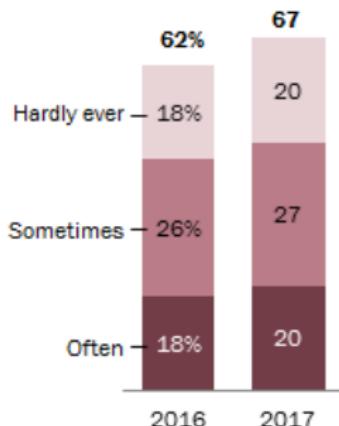
[Reply](#) [Retweet](#) [Favorite](#) [More](#)

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)
- ▶ 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)

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

# Outline

- ▶ Social media data
  - ▶ Motivation
  - ▶ Opportunities and challenges
  - ▶ APIs
  - ▶ Twitter data
  - ▶ Facebook data

# 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

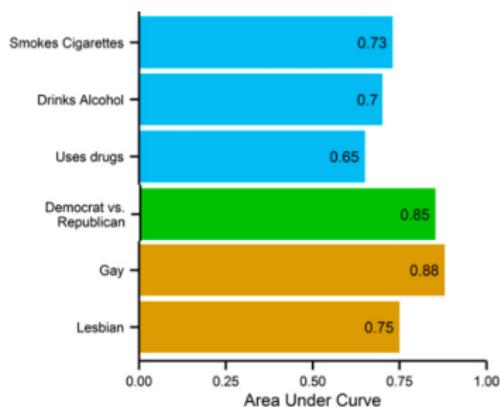
# Social media research

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, ...
- Non-intrusive measurement of behavior and public opinion  
Beauchamp (AJPS 2016): "Predicting and Interpolating State-level Polls using Twitter Textual Data"
- Inference of latent traits: political knowledge, ideology, personal traits, socially undesirable behavior, ...

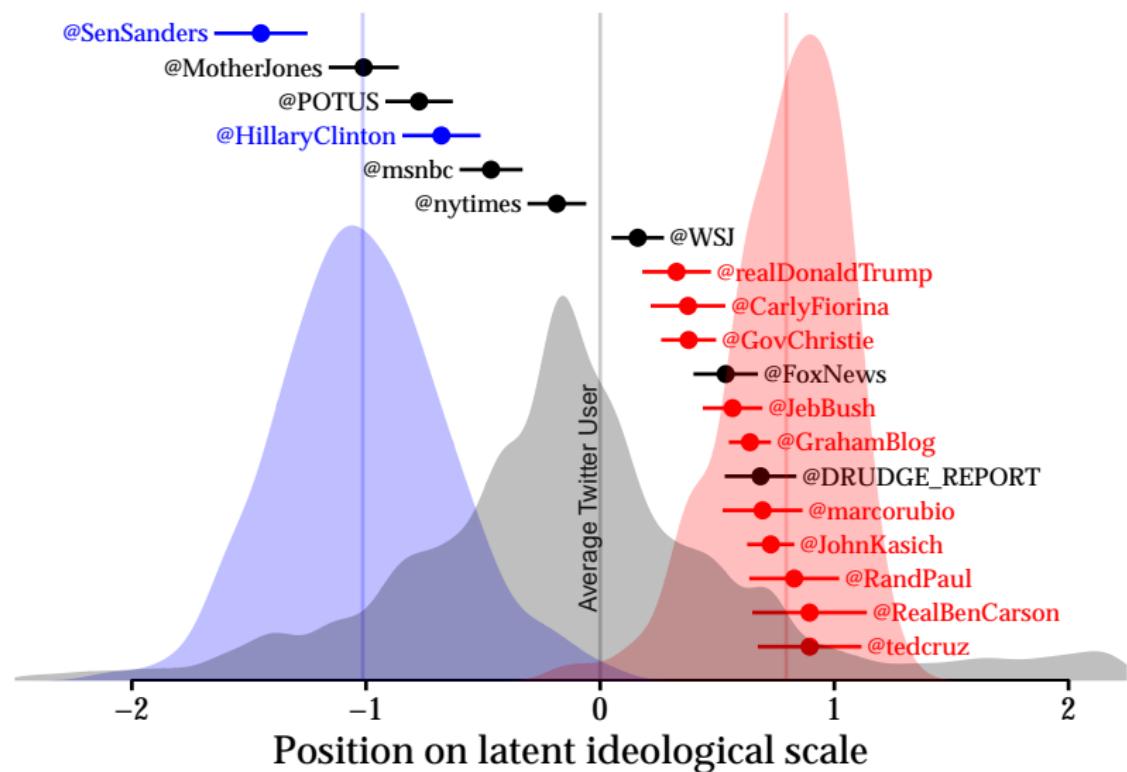


Kosinki et al, 2013, "Private traits and attributes are predictable from digital records of human behavior", *PNAS* (also personality, *PNAS* 2015)

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

2012	Registration History
99	+

# 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.

**I Voted**

 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

# Social media research

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

1. Social media as a new source of information

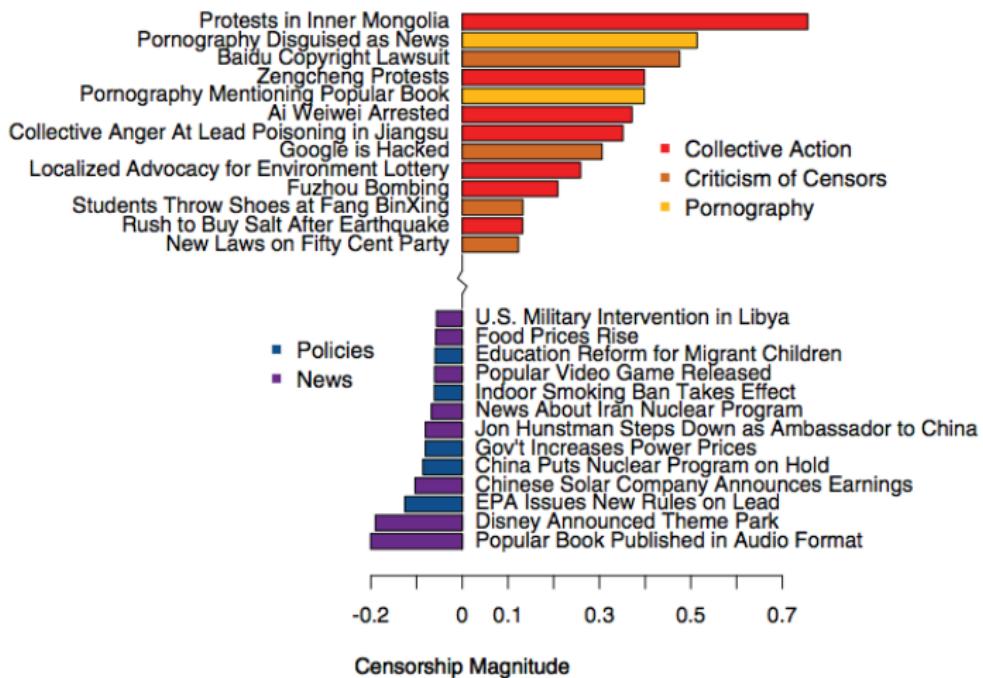
- ▶ Behavior, opinions, and latent traits
- ▶ Interpersonal networks
- ▶ Elite behavior
- ▶ Affordable field experiments

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*

# 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

The screenshot shows a journal article page from the Springer website. The article is titled "Tweetment Effects on the Tweeted: Experimentally Reducing Racist Harassment" by Kevin Munger. It was published in September 2017, Volume 39, Issue 3, pp 629–649. The page includes sections for authors, author affiliations, and metrics like shares, downloads, and citations.

**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

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

# 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 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
7. Social media as text

# 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 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)

# 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**

*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
- ▶ Tools extend innate capacities (e.g. Dunbar's number)
- ▶ Anonymity encourages vitriol

## 6. Ethical concerns

### 1. Shifting notion of *informed consent*



## 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

### 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

## Challenges of social media as source of text

- ▶ Large amounts of data
  - ▶ Storage problems
  - ▶ Analysis problems
- ▶ Language is informal and often non-textual (emojis, links, images) - and slang, txtspk, emoticons :-(
- ▶ Text in multiple languages, with URLs and hashtags
- ▶ Short length of text leads to sparse DFM

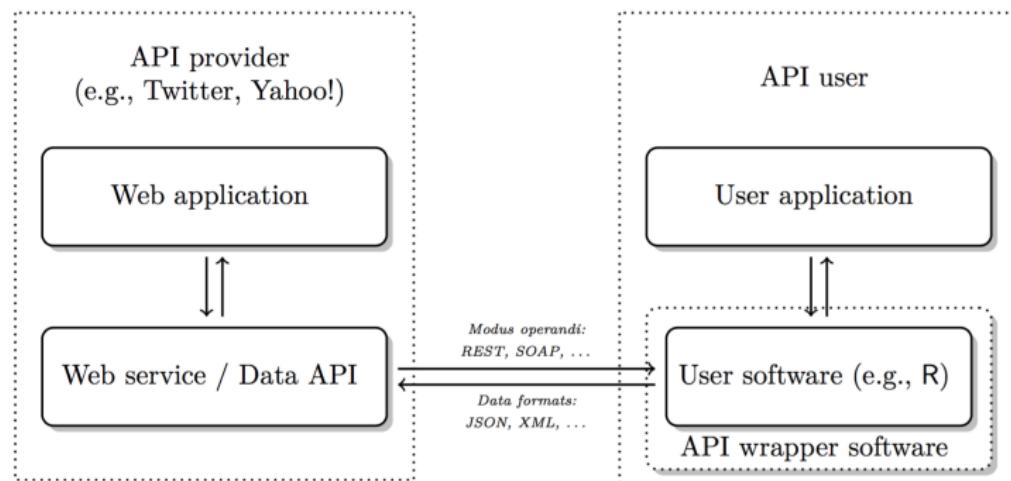
# Outline

- ▶ Social media data
  - ▶ Motivation
  - ▶ Opportunities and challenges
  - ▶ APIs
  - ▶ Twitter data
  - ▶ Facebook data

# APIs

API = Application Programming Interface; a set of structured http requests that return data in a lightweight format.

HTTP = Hypertext Transfer Protocol; how browsers and e-mail clients communicate with servers.



**Source:** Munzert et al, 2014, Figure 9.8

# APIs

Types of APIs:

1. RESTful APIs: queries for static information at current moment (e.g. user profiles, posts, etc.)
2. Streaming APIs: changes in users' data in real time (e.g. new tweets, weather alerts...)

APIs generally have extensive [documentation](#):

- ▶ Written for developers, so must be understandable for humans
- ▶ What to look for: [endpoints](#) and [parameters](#).

Most APIs are [rate-limited](#):

- ▶ Restrictions on number of API calls by user/IP address and period of time.
- ▶ Commercial APIs may impose a monthly fee

# Connecting with an API

Constructing a REST API call:

- ▶ Baseline URL **endpoint**:

`https://maps.googleapis.com/maps/api/geocode/json`

- ▶ Parameters: `?address=budapest`
- ▶ Authentication token (optional): `&key=XXXXX`

From R, use `httr` package to make GET request:

```
library(httr)
r <- GET(
  "https://maps.googleapis.com/maps/api/geocode/json",
  query=list(address="budapest"))
```

If request was successful, returned code will be 200, where 4xx indicates client errors and 5xx indicates server errors.

If you need to attach data, use POST request.

# JSON

Response is often in JSON format (Javascript Object Notation).

- ▶ Type: `content(r, "text")`
- ▶ Data stored in key-value pairs. Why? Lightweight, more flexible than traditional table format.
- ▶ Curly brackets embrace objects; square brackets enclose arrays (vectors)
- ▶ Use `fromJSON` function from `jsonlite` package to read JSON data into R
- ▶ But many packages have their own specific functions to read data in JSON format; `content(r, "parsed")`

## Authentication

- ▶ Many APIs require an access key or token
- ▶ An alternative, open standard is called OAuth
- ▶ Connections without sharing username or password, only temporary tokens that can be refreshed
- ▶ `httr` package in R implements most cases (examples)

## R packages

Before starting a new project, worth checking if there's already an R package for that API. Where to look?

- ▶ CRAN Web Technologies Task View (but only packages released in CRAN)
- ▶ GitHub (including unreleased packages and most recent versions of packages)
- ▶ rOpenSci Consortium

Also see this great list of APIs in case you need inspiration.

# Why APIs?

## Advantages:

- ▶ ‘Pure’ data collection: avoid malformed HTML, no legal issues, clear data structures, more trust in data collection...
- ▶ Standardized data access procedures: transparency, replicability
- ▶ Robustness: benefits from ‘wisdom of the crowds’

## Disadvantages

- ▶ They’re not too common (yet!)
- ▶ Dependency on API providers
- ▶ Lack of natural connection to R

# Outline

- ▶ Social media data
  - ▶ Motivation
  - ▶ Opportunities and challenges
  - ▶ APIs
  - ▶ Twitter data
  - ▶ Facebook data

# 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"
}
```

# 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

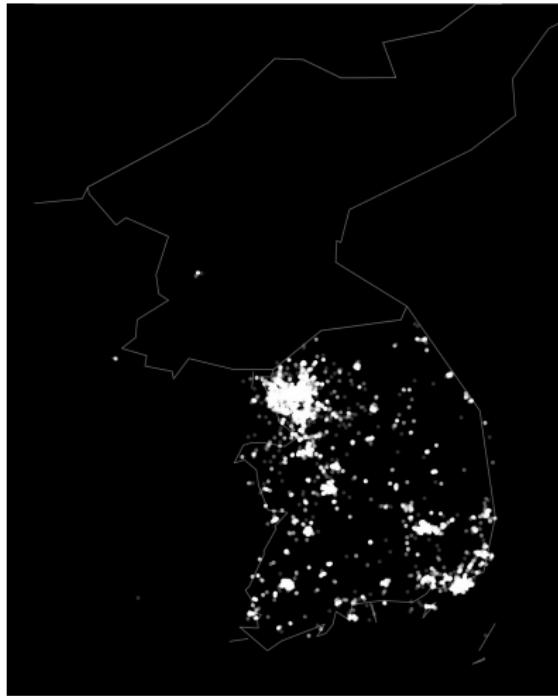
## Sampling bias?

[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

[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?



A screenshot of a Twitter profile card. At the top is the North Korean flag. Below it, the handle **North Korea English** is displayed in large white letters, followed by the URL [uriminzokkiri.com](http://uriminzokkiri.com). A bio below the handle reads: "An English translation of @uriminzok - the official North Korea Twitter feed". The card shows 671 tweets, 940 accounts followed, and 129 followers. There are "Follow" and "Profile" buttons. The main section is titled "Tweets" and shows one tweet from the user.

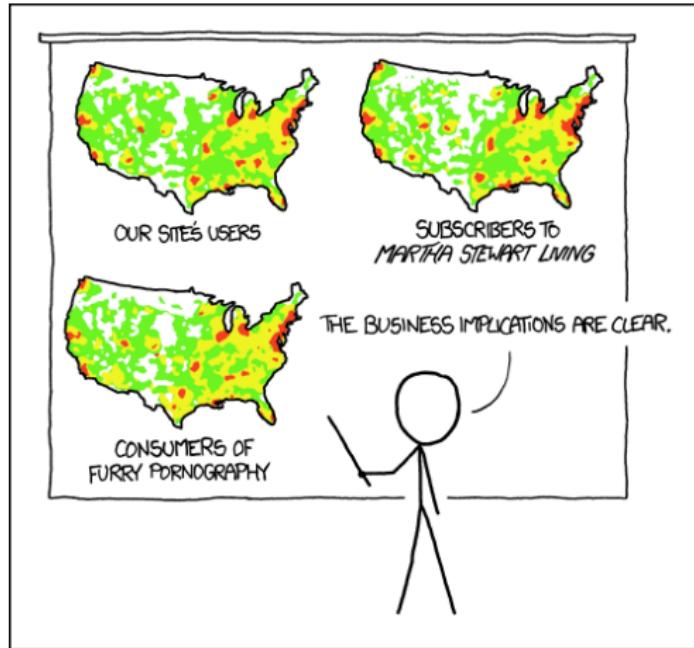
**Tweets**

 **North Korea English** @uriminzok\_engl 13h  
Beloved Comrade Kim Jong-un 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

# Outline

- ▶ Social media data
  - ▶ Motivation
  - ▶ Opportunities and challenges
  - ▶ APIs
  - ▶ Twitter data
  - ▶ Facebook data

## 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](#)