

# Why you can, can't & shouldn't do with social media data

Rachael Tatman  
JSM 2018



Growth in number of papers with "Social Media" as a topic  
(Source: [Semantic Scholar](#))

# What even is social media?

- An online platform...
- Containing user generated content...
- That allows users to interact with each other...
- And to create profiles\*

#irc

@rctatman



\*Obar, Jonathan A.; Wildman, Steve (2015). "Social media definition and the governance challenge: An introduction to the special issue". Telecommunications policy. 39(9): 745–750.

## Working with social media data, you...

### Can

- Increase the size and speed of research projects
- Ask new types of questions

### Can't

- Know for sure who you're looking at
- Violate developer's agreements

### Shouldn't

- Violate reasonable expectations of privacy
- Enable potentially unethical projects

## Can

- Increase the size and speed of research projects
- Ask new types of questions

## Can't

Working with social media data, you...

- Know for sure who you're looking at
- Violate developer's agreements

## Shouldn't

- Violate reasonable expectations of privacy
- Enable potentially unethical projects

## Dictionary of American Regional English (DARE)

Data collection

48 years (1965 - 2013)

Size of team

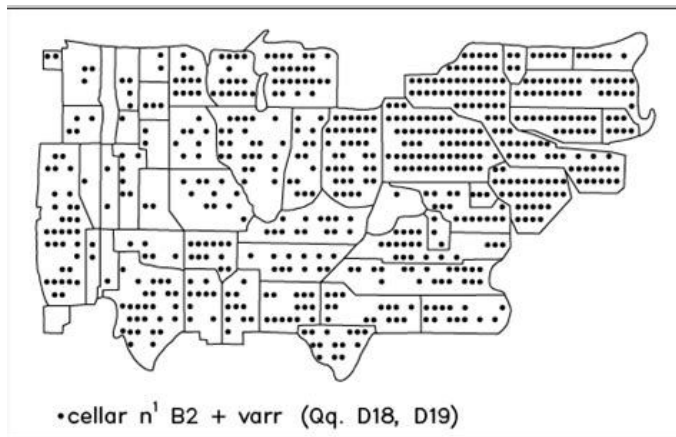
2,777 people

Number of participants

1,843 people

Results

(Distribution of people who use the word “cellar”)

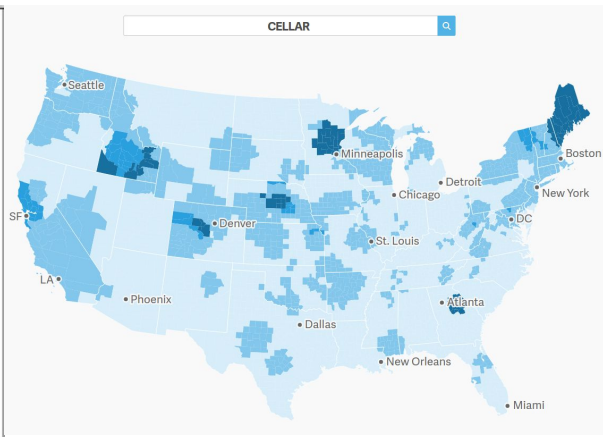


## Word Mapper App

<1 year

4 people

20 million



- Large scale social network analysis
  - Simple (and fast!) to quantify strength & direction of bonds
- Avoid common data collection problems:
  - **Bradley effect:** People tend to tell researcher what they think they want to hear
  - **Response bias:** The sample of people willing to do an experiment/survey differ in a meaningful way from the population as a whole
  - **Observer's paradox/Hawthorne effect:** People change their behaviour when they know they're being observed
- Computer mediated communication without speech analogs
  - Emoticons/Emoji
  - Reaction gif's
  - Memes(but see Silvia Sierra's work)

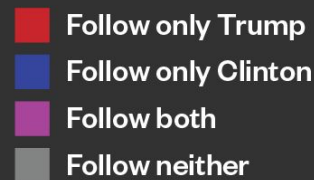
# Clinton and Trump supporters live in their own Twitter worlds

**Clinton Supporters**

Hillary Clinton supporters in this user group are not as cohesive as Trump supporters and they interact more frequently with users who follow both or neither candidate. They have few mutual follower networks in common with the far-right conservative cluster.

This large cluster of Trump supporters on Twitter have little mutual follower overlap with other users and are a remarkably cohesive group. They exist in their own information bubble.

**Trump Supporters**





# Self-Representation on Twitter Using Emoji Skin Color Modifiers

Alexander Robertson, Walid Magdy, Sharon Goldwater

Institute for Language, Computation and Cognition

School of Informatics, University of Edinburgh

alexander.robertson@ed.ac.uk, {wmagdy, sgwater}@inf.ed.ac.uk

## Abstract

Since 2015, it has been possible to modify certain emoji with a skin tone. The five different skin tones were introduced with the aim of representing more human diversity, but some commentators feared they might be used as a way to negatively represent other users/groups. This paper presents a quantitative analysis of the use of skin tone modifiers on emoji on Twitter, showing that users with darker-skinned profile photos employ them more often than users with lighter-skinned profile photos, and the vast majority of skin tone usage matches the color of a user's profile photo—i.e., tones represent the self, rather than the other. In the few cases where users do use opposite-toned emoji, we find no evidence of negative racial sentiment. Thus, the introduction of skin tones seems to have met the goal of better representing human diversity.

## Introduction

Emoji—icons used to represent emotions, ideas, or objects—became a formally recognized component of the Unicode Standard in 2010. At that time, all emoji depicting humans were rendered with the same yellow skin tone. However, in 2015, special modifier codes were introduced to allow users to change the skin tone of certain emoji. Unicode Technical



Figure 1: The current 102 tone-modifiable emoji. The first five rows are modified by a skin tone (shown in the left-most column), and the final row contains unmodified emoji.

TME+ were light colors. By annotating the profile photos of 4,099 users and analyzing their tweet histories, we show that the overall prevalence of lighter skin tone modifiers is likely due to the prevalence of lighter-skinned users. Indeed, users with darker-skinned profile photos are more likely to use TME, and more likely to modify them. Moreover, the vast majority of TME+ are similar in tone to the user's profile photo, suggesting that TME+ are used overwhelmingly for

## Can

- Increase the size and speed of research projects
- Ask new types of questions

## Can't

- **Know for sure who you're looking at**
- **Violate developer's agreements**

## Shouldn't

- Violate reasonable expectations of privacy
- Enable potentially unethical projects

Working with social media data, you...

# Who are your "participants"?

- Social media users aren't an unbiased sample
- Demographics vary dramatically between platforms
  - Snapchat and Instagram users are younger
  - Pinterest has more female users
  - LinkedIn users are predominantly middle/upper middle class
- You can't be sure...
  - That the data in your sample was produced by humans (as opposed to bots)
  - Of the identity of individuals in your data (even human-run accounts may be sock puppets)
  - Of the validity of social ties

# Developer Agreements

# Example: Twitter's Developer Agreement

PSA: as of July 24, 2018 there are [NEW requirements](#)

- You now need to apply for a developer account in order to get API access
- “When applying, all developers will be required to provide detailed information about how they use or intend to use Twitter’s APIs so that we can better ensure compliance with our policies”
- More details: <https://developer.twitter.com/>

Let's look at the parts of agreement  
most relevant to researchers

**Do not** (and do not permit others to) **associate the Twitter Content with any** person, household, device, browser, or other **individual identifier, unless** you or the entity on whose sole behalf you make such an association do so (a) with the express opt-in consent of the applicable individual; or (b) **based solely on publicly available data and/or data provided directly by the applicable individual that the individual would reasonably expect to be used for that purpose.**

My take:

- Try to avoiding storing PII (personally identifiable information) if possible
- You're probably in the clear with truly public figures (like celebrities or politicians)

If **Twitter Content** is deleted, gains protected status, or is otherwise suspended, withheld, modified, or **removed from the Twitter Service** (including removal of location information), you will make all reasonable efforts to **delete** or modify such Twitter Content (as applicable) **as soon as reasonably possible**, and in any case within 24 hours after a request to do so by Twitter or by a Twitter user with regard to their Twitter Content, unless otherwise prohibited by applicable law or regulation, and with the express written permission of Twitter.

My take:

- Don't store tweets for longer than you need for your analysis, instead store stats/relevant metrics

If your Service will **display Twitter Content to the public** or to end users of your Service, and you do not use Twitter Kit or Twitter for Websites to do so, then **you must use the Twitter API to retrieve the most current version of the Twitter Content for such display**. If Twitter Content ceases to be available through the Twitter API, you may not display such Twitter Content and must remove it from non-display portions of your Service as soon as reasonably possible.

My take:

- If you're going to share data, you can only share share Twitter ID's (or, for non-commercial research purposes only, less than 50,000 tweets via manual download)
- Use the API or a GUI wrapper of it like [DocNow's hydrator](#) to re-collect tweets as needed



If you provide Twitter Content to third parties, including downloadable datasets of Twitter Content or an API that returns Twitter Content, **you will only distribute or allow download of Tweet IDs**, Direct Message IDs, and/or User IDs.

1. **You may, however, provide export via non-automated means** (e.g., download of spreadsheets or PDF files, or use of a “save as” button) of **up to 50,000 public Tweet Objects** and/or User Objects per user of your Service, per day.
2. Any Twitter Content provided to third parties remains subject to this Policy, and those third parties must agree to the Twitter [Terms of Service](#), [Privacy Policy](#), [Developer Agreement](#), and [Developer Policy](#) before receiving such downloads.
  1. You may not distribute more than 1,500,000 Tweet IDs to any entity (inclusive of multiple individual users associated with a single entity) within any given 30 day period, unless you are doing so on behalf of an academic institution and for the sole purpose of non-commercial research or you have received the express written permission of Twitter.
  2. **You may not distribute Tweet IDs for the purposes of (a) enabling any entity to store and analyze Tweets for a period exceeding 30 days unless you are doing so on behalf of an academic institution and for the sole purpose of non-commercial research** or you have received the express written permission of Twitter, or (b) enabling any entity to circumvent any other limitations or restrictions on the distribution of Twitter Content as contained in this Policy, the Twitter Developer Agreement, or any other agreement with Twitter.

Working with social media data, you...

## Can

- Increase the size and speed of research projects
- Ask new types of questions

## Can't

- Know for sure who you're looking at
- Violate developer's agreements

## Shouldn't

- **Violate reasonable expectations of privacy**
- **Enable potentially unethical projects**

# Research Ethics in 60 Seconds

The Belmont Report lays out three primary ethical principles:

1. **Respect for Persons:** Individuals should be treated as autonomous agents and persons with diminished autonomy are entitled to protection.
2. **Beneficence:** 1) Do not harm and 2) maximize possible benefits and minimize possible harms.
3. **Justice:** Both the risks and benefits of research should be distributed equally.

Social media research may not count as human subjects research... but (in my opinion) we should be holding ourselves to the same ethical standard.

**Table 4.** “How Would You Feel If a Tweet of Yours Was Used in a Research Study and . . .” (n = 268).

	Very uncomfortable	Somewhat uncomfortable	Neither uncomfortable nor comfortable	Somewhat comfortable	Very comfortable
. . . you were not informed at all?	35.1%	31.7%	16.4%	13.4%	3.4%
. . . you were informed about the use after the fact?	21.3%	29.1%	20.5%	22.0%	7.1%
. . . it was analyzed along with millions of other tweets?	2.6%	18.7%	25.5%	30.0%	23.2%
. . . it was analyzed along with only a few dozen tweets?	16.5%	30.3%	24.0%	20.2%	9.0%
. . . it was from your “protected” account?	54.9%	20.5%	13.8%	6.0%	4.9%
. . . it was a public tweet you had later deleted?	31.3%	32.5%	20.5%	10.4%	5.2%
. . . no human researchers read it, but it was analyzed by a computer program?	2.6%	14.3%	30.5%	32.3%	20.3%
. . . the human researchers read your tweet to analyze it?	9.7%	27.6%	25.0%	25.4%	12.3%
. . . the researchers also analyzed your public profile information, such as location and username?	32.2%	23.2%	21.0%	13.9%	9.7%
. . . the researchers did not have any of your additional profile information?	4.9%	15.4%	25.1%	34.1%	20.6%
. . . your tweet was quoted in a published research paper, attributed to your Twitter handle?	34.3%	21.6%	21.6%	13.1%	9.3%
. . . your tweet was quoted in a published research paper, attributed anonymously?	9.0%	16.8%	26.5%	28.4%	19.4%

Note. The shading was used to provide a visual cue about higher percentages.

Fiesler, C., & Proferes, N. (2018). “Participant” Perceptions of Twitter Research Ethics. *Social Media+ Society*, 4(1), 2056305118763366.

## More acceptable:

- Large datasets
- Analyzed automatically
- Social media users informed about research
- Anonymized tweets being quoted (note that enough words should be changed that a reverse search isn't possible)

## Less acceptable:

- Small datasets
- Analysis done by hand (presumably including analysis by Mechanical Turk workers)
- Tweets from protected accounts or deleted tweets analyzed
- Quoting with citation (very different from academic norms!)

# Unethical projects?

- Being involved in research (**especially** without their consent) should not result in harm to users
- Need to consider not only the project for which the data was originally collected but also potential unintended consequences
- Examples:
  - Collecting tweets discussing recreational marijuana use in the US in order to infer demographics and compare to current patterns of arrest
    - **Problem:** Research data isn't protected from subpoena; data you collect could be used as evidence against them in court.
  - Create a database of the hometowns of users based in New York to study migration patterns
    - **Problem:** Might contain information incriminating information on immigration status of users

“The same-sex dating app Grindr says it will stop sharing its users' HIV status with other companies, after it was discovered the app was allowing third parties to access encrypted forms of the sensitive data.”

“Grindr Admits It Shared HIV Status Of Users”  
Scott Neuman for NPR, April 2018



“In 2016, [Optum] filed a patent application to gather what people share on platforms like Facebook and Twitter, and link this material to the person’s clinical and payment information. A company spokesman said in an email that the patent application never went anywhere. But the company’s current marketing materials say it combines claims and clinical information with social media interactions.”

“Health Insurers Are Vacuuming Up Details About You  
— And It Could Raise Your Rates”

Marshall Allen for ProPublica, July 2018

We need to be cautious  
and courteous when  
using social media data.

## Working with social media data, you...

### Can

- Increase the size and speed of research projects
- Ask new types of questions

### Can't

- Know for sure who you're looking at
- Violate developer's agreements

### Shouldn't

- Violate reasonable expectations of privacy
- Enable potentially unethical projects

Thank you!

Questions?