



Potentials and Pitfalls of Social Media Data

Indira Sen & Katrin Weller
GESIS workshop - December 2022

0. Welcome and logistics



Photo: Katrin Weller

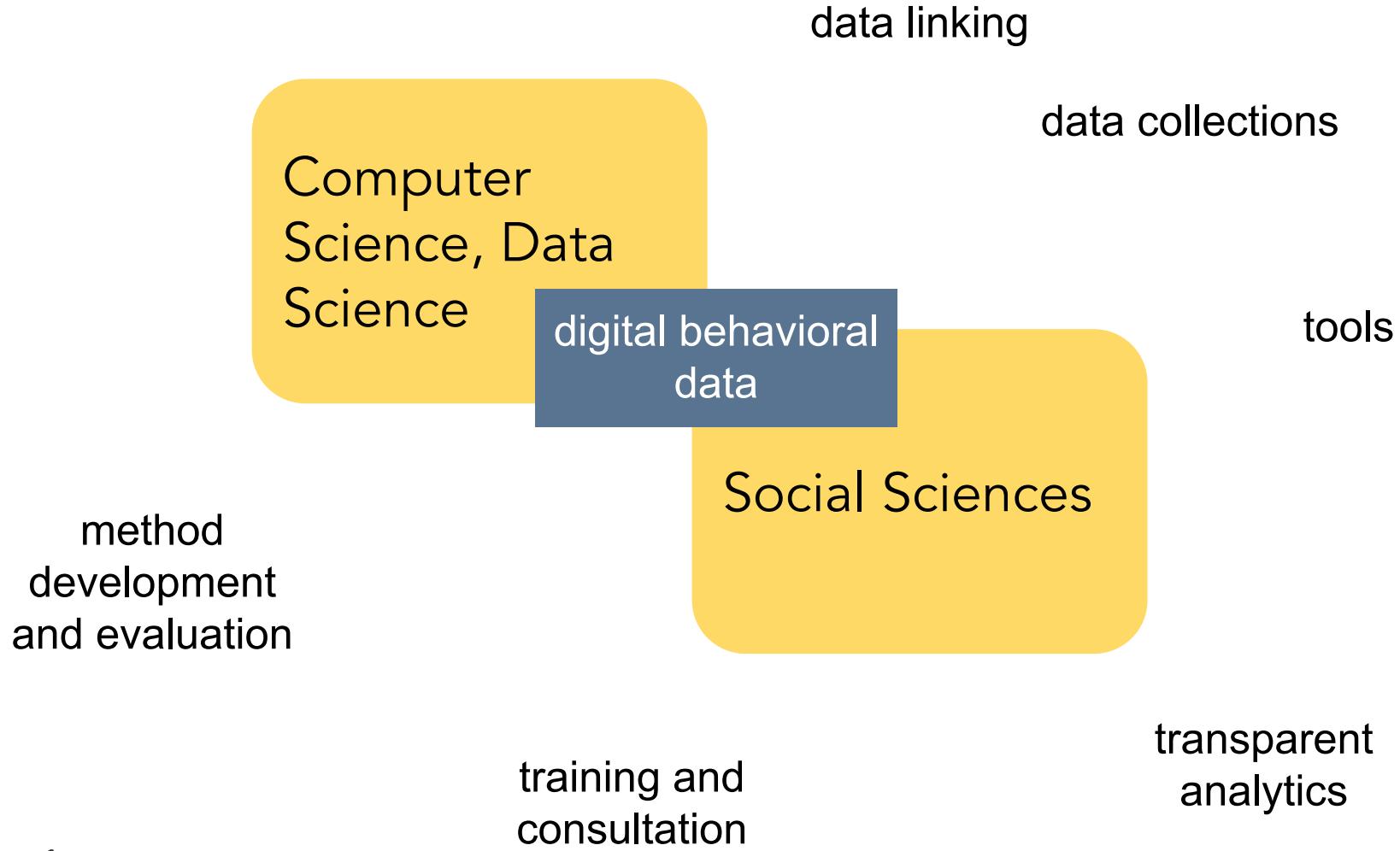
About us

Contact

- Indira: indira.sen@gesis.org, @indiigosky
- Katrin: katrin.weller@gesis.org, @kwelle



Computational Social Science at GESIS



Acknowledgements

We thank our colleagues at GESIS, especially:

Claudia Wagner

Fabian Flöck

Mattia Samory

Arnim Bleier

Bernd Weiß

Leon Fröhling

Felix Soldner

... and the entire CSS department

About the workshop

Zoom Etiquette

- Microphones off when not speaking
- Camera preferably on
- Asking questions or making comments (both welcome!)
 - raise hand on zoom
 - type in zoom chat
 - => **We will try to get to you ASAP**
- We have planned for breaks in the schedule, but feel free to request for additional short breaks through chat

Resources

Materials will be shared via:

https://github.com/Indiiigo/social_media_data_research_2022

This includes:

- **Slides**
- List of **references** at the end of the last slide deck
- **Notebooks** with examples of code that can be executed
- **Discussion Board:**

<https://tinyurl.com/socialmediaworkshop22>

Workshop goals

- provide an overview on current approaches in research **based on digital traces from social media**
- outline different steps in the research process when working with social media data, and provide **practical examples** for data collection, cleaning and analysis
- offer a **structured approach** to think about potential pitfalls and error sources in social media research, that can help to design, present, talk about research approaches

Motivation & Background



SERIOUSLY? DO THEY NOT REALIZE THAT 99% OF TWEETS ARE WORTHLESS BABBLE THAT READ SOMETHING LIKE 'JUST WOKE UP. GOING TO STARBUCKS NOW. GETTING LATTE.'

READERS' COMMENT FOUND IN THE COMMENT SECTION FOR GROSS, D. (2010, APRIL 14). LIBRARY OF CONGRESS TO ARCHIVE YOUR TWEETS. CNN. RETRIEVED FROM [HTTP://EDITION.CNN.COM/2010/TECH/04/14/LIBRARY.CONGRESS.TWITTER/](http://EDITION.CNN.COM/2010/TECH/04/14/LIBRARY.CONGRESS.TWITTER/),
RETRIEVED NOVEMBER 19.
PHOTOS: [HTTPS://WWW.FLICKR.COM/SEARCH/?TEXT=COFFEE&LICENSE=4%2C5%206%2C9%2C10](https://www.flickr.com/search/?text=coffee&license=4%2C5%206%2C9%2C10)



Digital Traces: So much data!

(...about people's behavior and attitudes)

2019 *This Is What Happens In An Internet Minute*



And so much research!

Computer Science

Social Sciences

2010 From tweets to polls
[O'Connor et al.]

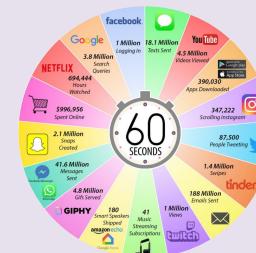
2014 Predicting tie strength
with social media
[Gilbert et al.]

2019 Investigating
commentator bias in
football broadcasts
[Merullo et al.]

2013 Text as data
[Grimmer et al.]

2013 Big Data in Survey
Research: AAPOR Task Force
Report [Japec et al.]

2019 Combining surveys and
digital traces
[Stier et al., Pasek et al.]



No formal research field - no standard methods

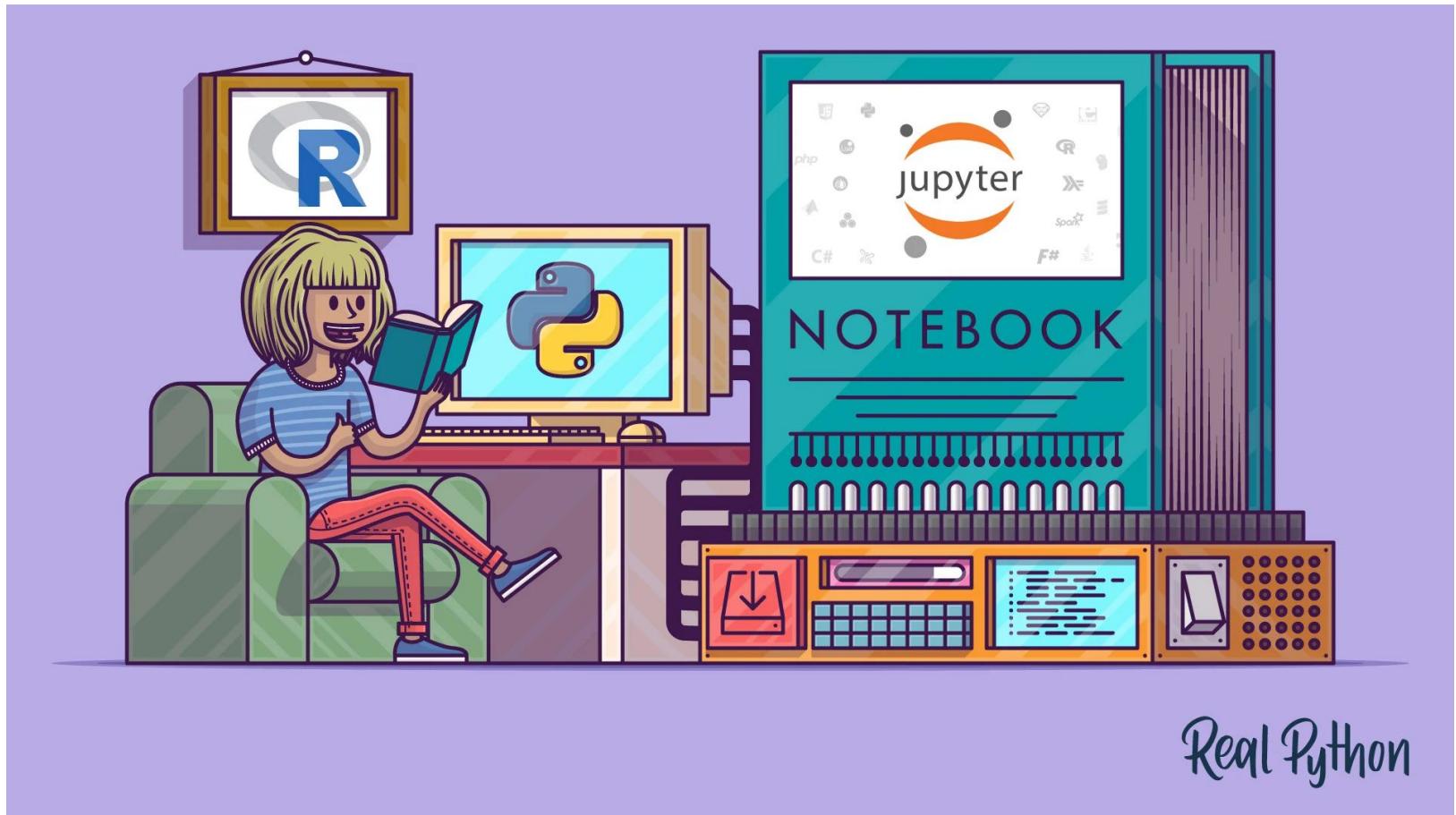
- multidisciplinary entry points
- diverse range of approaches (also mixed methods)
- many exploratory studies

how to get started?

Workshop goal

- Combine critical reflections and hand-on exercises
- get some first idea on how approaches work in practice
- provide reusable material

We will be using Notebooks for Coding



We will be using Notebooks for Coding

- Python + Google Colab:
<https://colab.research.google.com/>
→ you need to login via a google account
- light-weight, easy to share, good for reproducing results
- but of course there are also other options

*different exercises
throughout the workshop*

Agenda

- Session 1: Introduction to Research with Social Media Data (SMD)
- Session 2: SM Data Collection
- Session 3: SMD Preprocessing and Analysis
- Session 4: Potential Pitfalls of SMD
- Session 5: Identifying Pitfalls with help from surveys
- Session 6: Identifying Pitfalls in SMD
- Session 7: Mitigating Pitfalls
- Session 8: Documenting Pitfalls
- Session 9: Recap and Conclusions

Today's Schedule

Monday, 05.12.	
9:30-11:00	Introduction to Social Media Research
11:00-11:15	<i>Break</i>
11:15-12:00	Data Collection
12:00-12:30	<i>Break</i>
12:30-14:00	Hands on data collection

About you

Introductory Round

What brought you to this course?

We would like to learn more about

- Your interest in social media data (specific types of data or methods you might want to learn more about)
- Prior experience
- Your disciplinary background
- <https://tinyurl.com/socialmediaworkshop22>

Agenda

- Session 1: Introduction to Research with Social Media Data (SMD)
- Session 2: SM Data Collection
- Session 3: SMD Preprocessing and Analysis
- Session 4: Potential Pitfalls of SMD
- Session 5: Identifying Pitfalls with help from surveys
- Session 6: Identifying Pitfalls in SMD
- Session 7: Mitigating Pitfalls
- Session 8: Documenting Pitfalls
- Session 9: Recap and Conclusions

1. Introduction and Data Collection

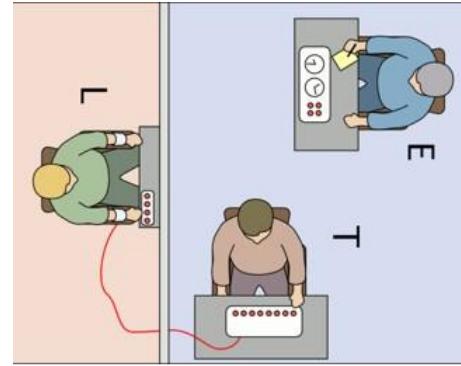


Introduction: digital traces from social media

Sources of data for SocSci research



Surveys



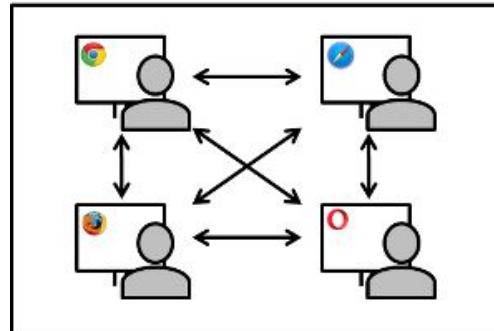
Experiments



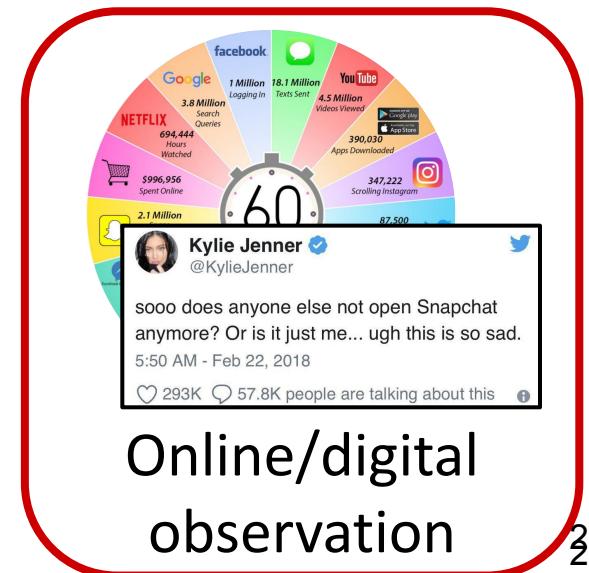
Observation



Online/digital
surveys



Online/digital
experiments



Online/digital
observation

Two main sources of data

Digital traces of humans

“Data documenting the interactions of users with digital devices or services”
[Howison et al., 2011]



Surveys (via Questionnaires)

“The collection of information from a **sample of individuals** through their **responses to questions**”
[Check & Schutt, 2012]

Two main sources of data

Digital traces of humans

“Data documenting the
interactions of users with
digital devices or services”
[Howison et al., 2011]

Surveys (via Questionnaires)

“The collection of information
from a **sample of individuals**
through their **responses to**
questions”
[Check & Schutt, 2012]



Two main sources of data

Properties of...

Digital traces of humans

- non-designed (“found”) & not reactive to sp. stimulus
- self-selected & n=all per system
- very heterogeneous signals
- large volume, cheap
- fine-grained (longitudinal, space)
- relational

Surveys

- designed in setup & specific stimulus
- usually probabilistically sampled f. target pop.
- reaction space constrained
- small scale, costly
- limited set of collections
- unlinked units

Computer Science

Social Sciences

2010 From tweets to polls
[O'Connor et al.]

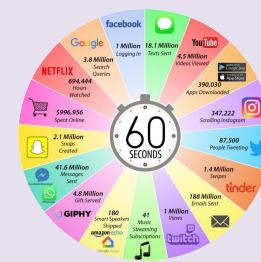
2014 Predicting tie strength
with social media
[Gilbert et al.]

2019 Investigating
commentator bias in
football broadcasts

2013 Text as data
[Grimmer et al.]

2013 Big Data in Survey
Research: AAPOR Task Force
Report [Japec et al.]

2019 Combining surveys and
digital traces
[Stier et al., Pasek et al.]



United by data: Social scientists and computer scientists are both interested in how digital traces can be used for inferring human behaviour & attitudes

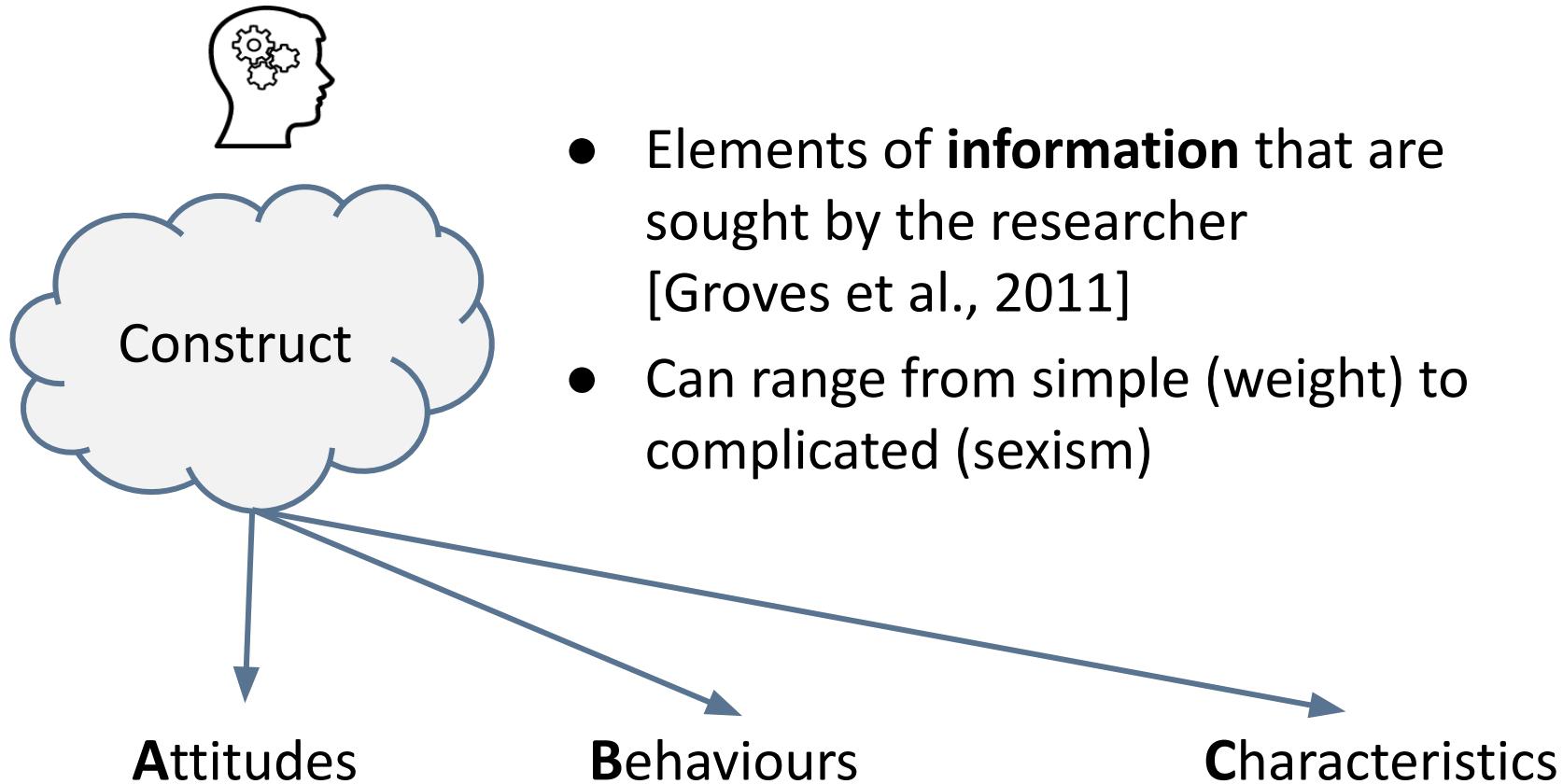
Let's unite on Methods and Theory

- **Common vocabulary** that would improve communication and cross-development of measurement theories
- **Knowledge transfer re: overlapping approaches**, especially in addressing challenges and errors
- **Standards in documenting pipelines and reporting results**

Understanding Social Phenomena



Understanding Social Phenomena



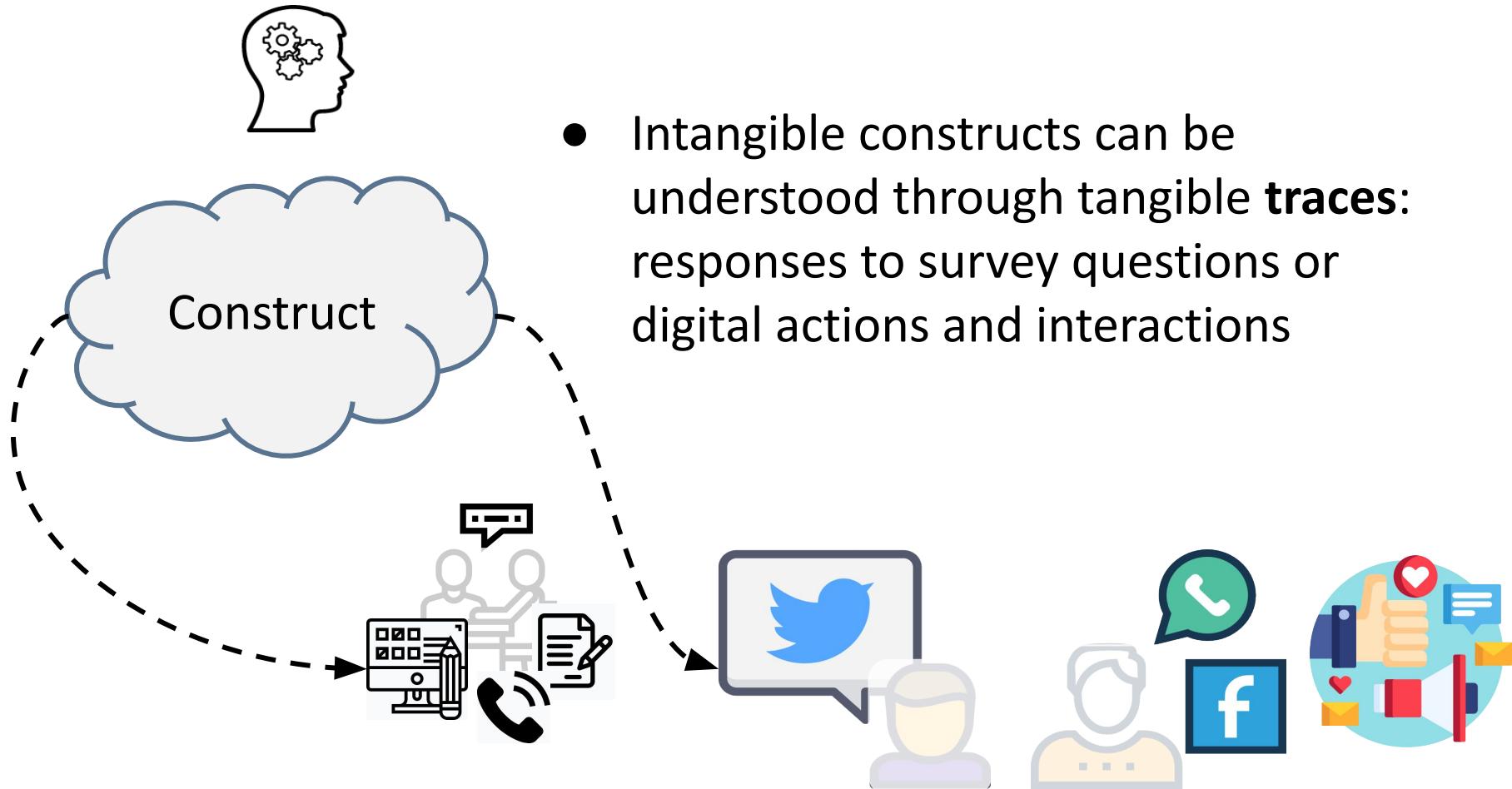
Understanding Social Phenomena



- The set of **people** to be studied
[Groves et al., 2011]
- Can be easily defined (all the preschool students in a city) to more difficult (all refugees)
- Usually a national population, can also be any “system population”
(-> platform study vs. “Social Sensing”)

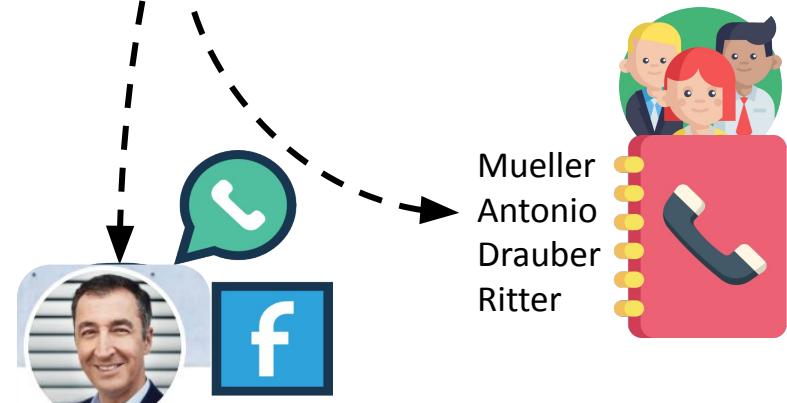
Target Population

Understanding Social Phenomena



Understanding Social Phenomena

- **Users** of the target population (TP) are the *representation of TP elements* in a record. This can be offline records (census, phonebook), online platforms (profiles, IPs), or other digital forms (bank accounts, fitness trackers)



Studying social media as digital traces

- Researchers value social media as a new type of data
- Previously „ephemeral data“ become visible
- Immediate – quick reaction to events
- Structured
- „natural“ data

“What I find really interesting is that structure becomes manifest in internet communication. So it’s the first time in history actually that we can, that social structures between people become manifest within a technology. (...) They become visible, they become crawlable, they become analyzable.”

Kinder-Kurlanda, Katharina E., and Katrin Weller. 2014. "I always feel it must be great to be a hacker!": The role of interdisciplinary work in social media research." In Proceedings of the 2014 ACM conference on Web Science, 91-98. New York: ACM.

Studying social media as digital traces

Data on social media platforms is not created for research purposes.

This leads to challenges in

- Accessibility
- Quality
- Interpretation
- Ethics

2006

DECEMBER 25, 2006 / JANUARY 1, 2007

www.time.com

TIME

PERSON OF THE YEAR

You.



Yes, you.
You control the Information Age.
Welcome to your world.

Social Media

Online platforms, that are typically build around user activities and connections/interactions between users.

Examples:

*Social-Networking-Sites, Microblogs, Blogs (& podcasts),
Media Sharing Sites (video sharing, photo sharing ...),
Social games, Social news, Q&A communities, Wikis,
Product-based communities*

Digital traces may also come from:

search engines, shopping portals, dating websites...

Categorizing social media platforms?

Web 2.0 (O'Reilly 2005)

Social Web

Social Networking Sites SNS

Social Media

Virtual Communities

Prosumerism (Buzzetto-More 2013) [based on Toffler, 1980]

Produsage (Bruns 2008)

User-generated content

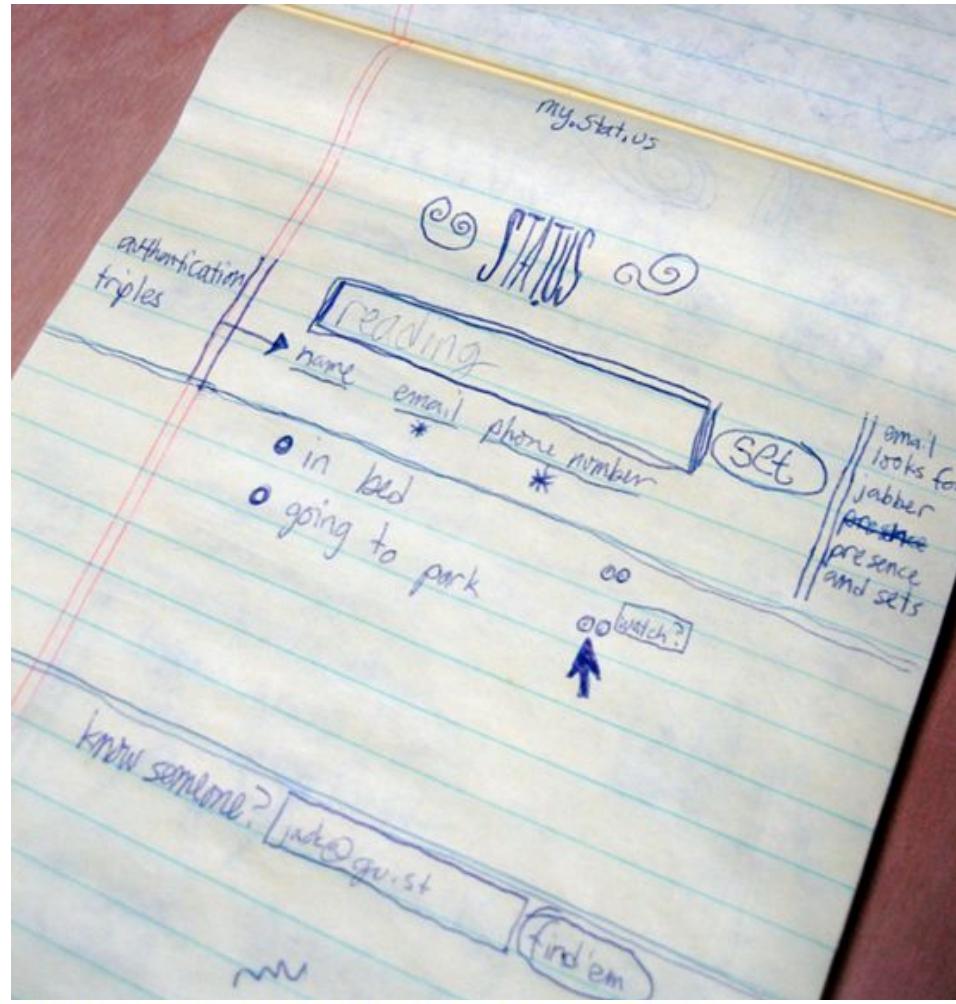
Online data / Internet data / Web data

Digital behavioural data

Categorizing social media platforms?

- What is the main type of „content“? (photos, texts, videos, users...)
- What are the types of interactions?
- What are the main activities and use cases?
(stay informed, connect, bookmark, share, play, sell...)
- Openness of the platform?
- Professional or private use?
- International or local?

Platforms and affordances



<https://www.flickr.com/photos/jackdorsey/182613360>

Cross-platform connections

- Facebook post linking to a blog post
- Instagram photo shared through Twitter
- Tweet consisting of just a link to a Facebook post
- Newspaper article shared and commented through Facebook
- ...

Platform affordances / functionalities

- Different interpretations of functionalities („Retweets are not endorsements“)
- Platform affordances can influence user behaviour - and the other way round
- Users creating new practices / standards

Chris Messina @chrismessina

Folgen

how do you feel about using # (pound) for groups. As in [#barcamp](#) [msg]?

RETWEETS GEFÄLLT
1.304 2.392

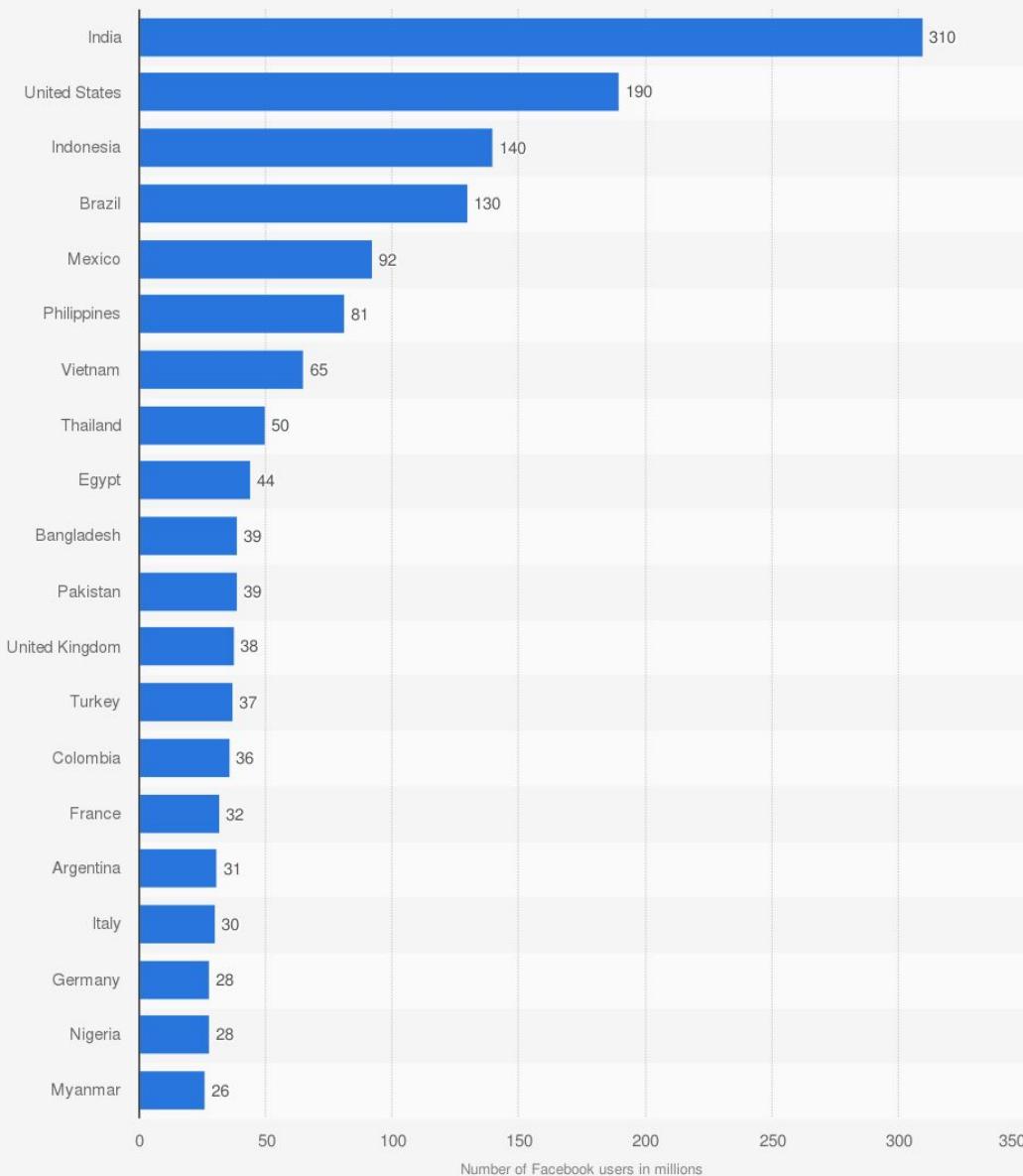
12:25 - 23. Aug. 2007

1,3 Tsd. 2,4 Tsd. ***

Social media users and usage

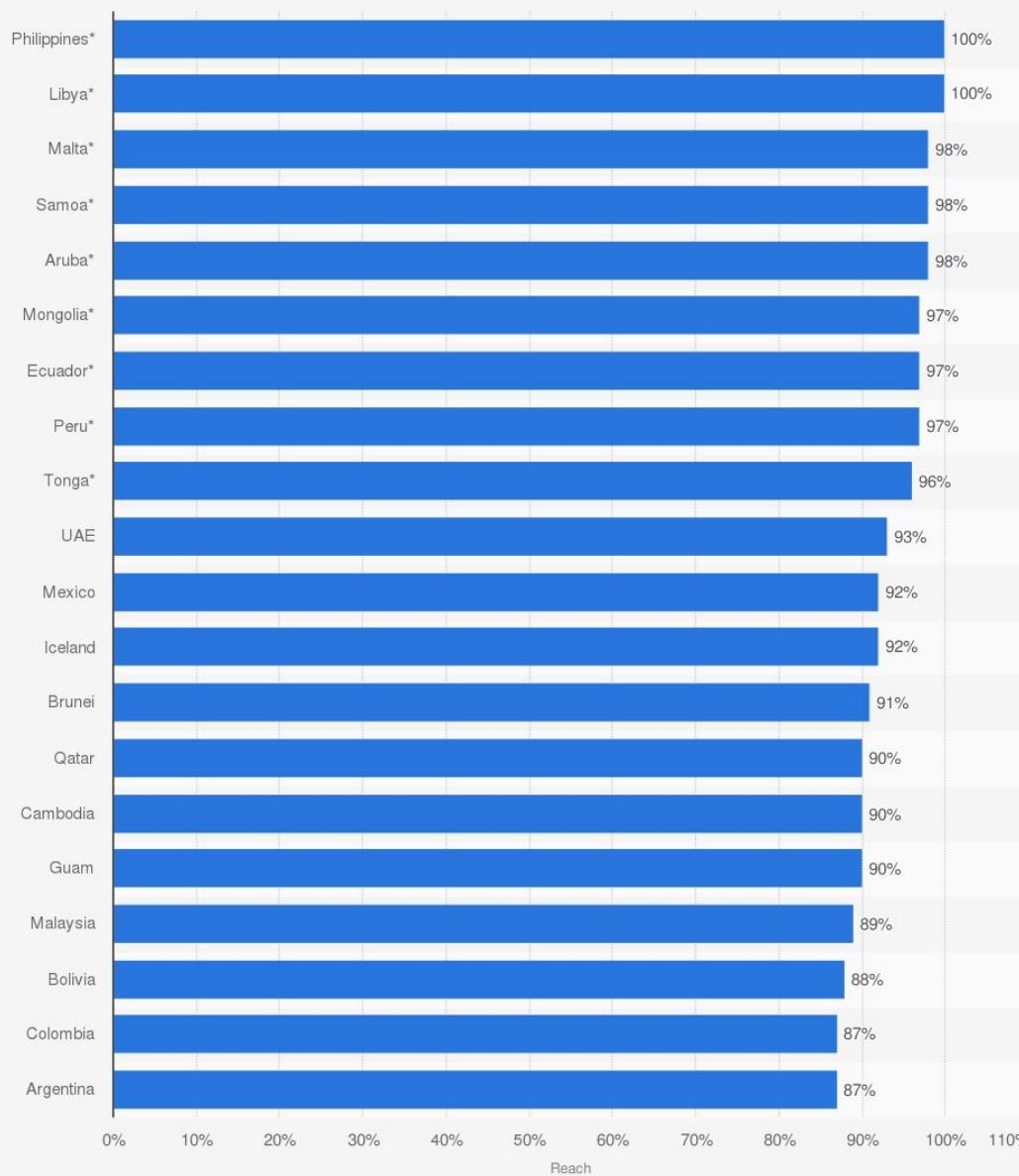
- Different user demographics per platform (often limited profile information)
- Different popular platforms per country
- Changes over time

Leading countries based on Facebook audience size as of October 2020 (in millions)



<https://www.statista.com/statistics/268136/top-15-countries-based-on-number-of-facebook-users/>

Countries with the highest Facebook audience reach as of October 2020



<https://www.statista.com/statistics/278435/percentage-of-selected-countries-internet-users-on-facebook/>

User numbers provided by platforms

What is being measured?

- Daily active users
- Monthly active users
- „active“ users
- Everyone with an account

Other challenges when working with usage numbers

- Different contexts for usage (e.g. mobile)
- Different approaches to capture demographics
- “Active” and “passive” users
- User behaviour is constantly changing, people are exploring new practices.

non-human social media users

Curiosity Rover [@MarsCuriosity](#)

NASA's latest mission to explore the surface of Mars. Roving the Red Planet since Aug. 5, 2012 (PDT) (Aug 6 UTC)

Gale Crater, Mars
mars.jpl.nasa.gov/msl/
Joined July 2008

2,803 Tweets 161 Following 1.74M Followers 317 Favorites 6 Lists

[Follow](#)

Tweets Tweets & replies Photos & videos

Curiosity Rover [@MarsCuriosity](#) Dec 3

Pics I take on Mars get posted online. Some high-res, some thumbnail. Some B&W, others color:
mars.jpl.nasa.gov/msl/multimedia...

View more photos and videos

Worldwide Trends Change

#Gala14GH15
#DPDA
#FitInHarmonyAtTheWhiteHouse
#MTVIsHere
#trainwreckLDC
Eric Garner
Attali
Alper Taş
Jessica Jones
SeçimBarajını Dağızı AYM

<https://twitter.com/MarsCuriosity>

32

2011



Research ethics

- Research ethics practices are also still evolving
- Main focus on privacy
- Lack of informed consent
- “participants” not expecting research activities on their data [[Fiesler & Proferes, 2018](#)]
- Different assumptions for different types of user groups (e.g. vulnerable groups)
- Ethically problematic research questions
- Potential starting point: [AoIR ethics guidelines](#)

Next week, we will discuss examples of research designs with social media data. Especially social science questions that have been investigated with other types of data, but can now be explored with social media data.

Email us until Tuesday afternoon if you want your specific examples to be included. Vague ideas are totally okay!