

Data Ethics: Understanding Big Data, Algorithmic Bias, and Research Ethics

SOCL 1101 Introduction to Sociology
Ineke Marshall
Fall 2021



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Workshop Agenda

- Objectives
- Introduce 'Big Data' Concepts
- Discuss data, privacy, and algorithms
- Activity: Adopt or Not?
- Discuss ethical implications of big data and lessons for (digital) research

Slides, handouts, and data available at

<https://bit.ly/diti-fa21-marshall-data-ethics>



Workshop Goals

- Understand the ways data are being used in society as well as how algorithms impact and shape our daily lives
- Explore the ways in which privacy and security are being reshaped and redefined through the use of big data, algorithms, and policy
- Understand the ways in which technology reflects cultural, social, and political biases.
- Explore the ways in which these questions and methods are influencing how social scientists do research and practice their craft



What is “Big Data”?



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Big Data is here (and it's getting *bigger*)

1 How much data is generated every minute?

Source: Domo

 **41,666,667**

messages shared
by WhatsApp users

 **1,388,889**

video / voice calls made
by people worldwide

 **404,444**

hours of video streamed
by Netflix users



2.1 Million



3.8 Million



4.5 Million



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

What is “Big Data”?

Companies, governments, and other groups **collect vast amounts of data from vast numbers of users** and analyze that data quickly for a variety of purposes, including advertising, marketing, surveillance, building profiles, etc.

The goal of big data is to predict individual user behavior based on patterns from the user as well as patterns from “similar” users (based on demographic information, behavioral patterns, etc).

If we’re living in an era of “surveillance capitalism,” **our information can be considered to be a valuable *product*.**



40 ZETTABYTES

[43 TRILLION GIGABYTES]

of data will be created by 2020, an increase of 300 times from 2005



6 BILLION PEOPLE have cell phones

Volume SCALE OF DATA

It's estimated that **2.5 QUINTILLION BYTES**

[2.3 TRILLION GIGABYTES] of data are created each day

Most companies in the U.S. have at least **100 TERABYTES** [100,000 GIGABYTES] of data stored

The New York Stock Exchange captures

1 TB OF TRADE INFORMATION

during each trading session



Velocity ANALYSIS OF STREAMING DATA

By 2016, it is projected there will be

18.9 BILLION NETWORK CONNECTIONS

— almost 2.5 connections per person on earth



Modern cars have close to **100 SENSORS** that monitor items such as fuel level and tire pressure



The FOUR V's of Big Data

From traffic patterns and music downloads to web history and medical records, data is recorded, stored, and analyzed to enable the technology and services that the world relies on every day. But what exactly is big data, and how can these massive amounts of data be used?

As a leader in the sector, IBM data scientists break big data into four dimensions: **Volume, Velocity, Variety and Veracity**

Depending on the industry and organization, big data encompasses information from multiple internal and external sources such as transactions, social media, enterprise content, sensors and mobile devices. Companies can leverage data to adapt their products and services to better meet customer needs, optimize operations and infrastructure, and find new sources of revenue.

By 2015 **4.4 MILLION IT JOBS** will be created globally to support big data, with 1.9 million in the United States



As of 2011, the global size of data in healthcare was estimated to be

150 EXABYTES

[161 BILLION GIGABYTES]



30 BILLION PIECES OF CONTENT are shared on Facebook every month



Variety DIFFERENT FORMS OF DATA

By 2014, it's anticipated there will be

420 MILLION WEARABLE, WIRELESS HEALTH MONITORS



4 BILLION+ HOURS OF VIDEO are watched on YouTube each month



400 MILLION TWEETS are sent per day by about 200 million monthly active users



1 IN 3 BUSINESS LEADERS

don't trust the information they use to make decisions



27% OF RESPONDENTS

in one survey were unsure of how much of their data was inaccurate



Veracity UNCERTAINTY OF DATA

Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



Why should we care?

- Big data is characterized by its **scale**
- Big data **sources** include: digitized records, social media/internet activity, or sensors from the physical environment.
- Big data is often **privately owned**
 - Example: an insurance company purchasing social media activity from Facebook in order to make insurance sales decisions.



Online Presence & Data Privacy



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

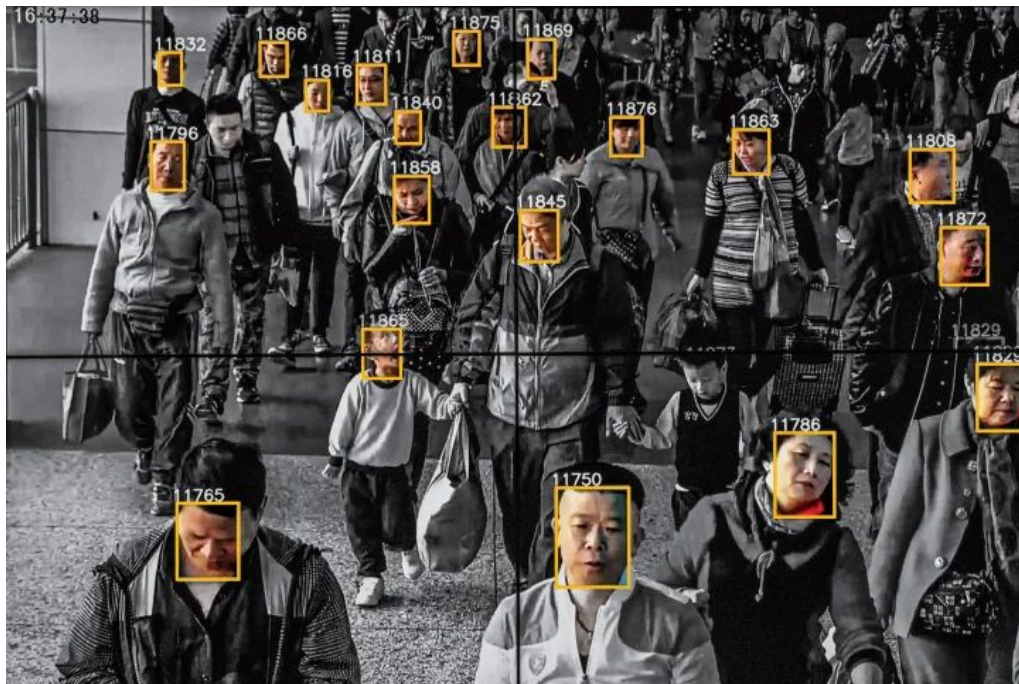
Questions to consider

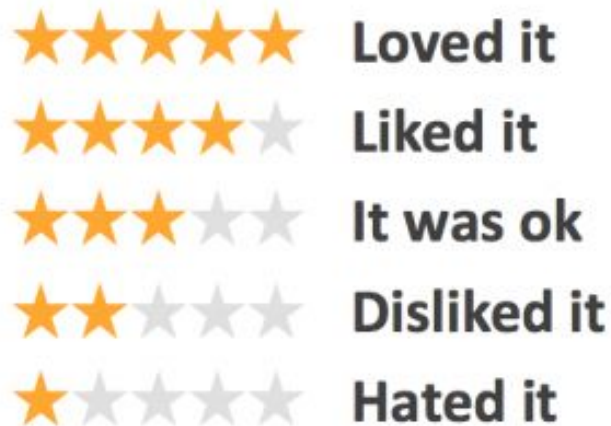
- How are we being represented online?
- How are our data being used?
- Who is using our data and for what purposes?
- How might our data be used in the future?



An Example: China's Social Credit System

- What is China's Social Credit system?
How does it work?



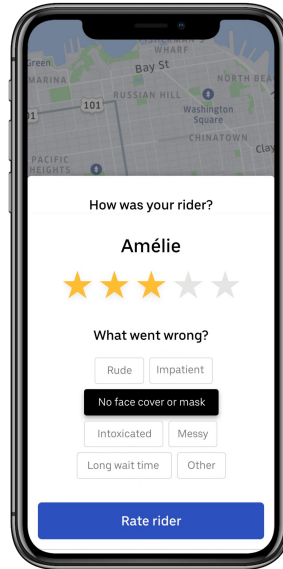
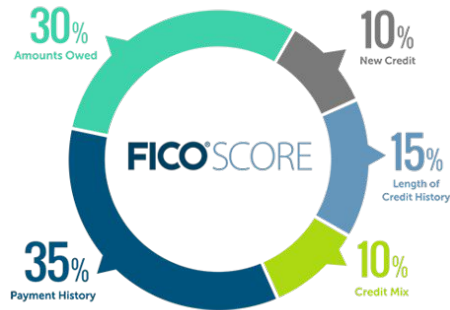


Black Mirror: Nosedive (2016)



Discussion: America's Social Credit System

In what ways might America have similar or different technological infrastructures when compared with China?



The bouncer that never forgets a face

Spot trouble from 50,000+ individuals known for assaults, chargebacks, drugs and property damage.

Reduce nightlife incidents by as much as 97% by spotting trouble before it becomes a problem. Receive alerts when troublemakers scan their ID including details on why they've been flagged.

[Book Demo](#)



Northeastern University
NULab for Texts, Maps, and Networks

Feel free to ask questions at any point during the presentation!

What about Social Media?

Facebook collects, stores, and sells information about you so you get more targeted ads and your newsfeed is tailored to your categories.

Other social media sites that do this:

- Instagram (owned by Facebook)
- Google
- YouTube (owned by Google)
- Twitter





AWARENESS | SCIENCE & TECH | AUG 3, 2019 AT 11:08 AM.

Google's File on You is 10 Times Bigger Than Facebook's — Here's How to View It

Google, Amazon, Apple, and Microsoft are all central players in “surveillance capitalism” and prey on our data.



Example: If you have **location services** turned on for Google (like if you use Google maps), Google can track your every move.

Go to:

<https://www.google.com/maps/timeline>



Northeastern University
NULab for Texts, Maps, and Networks

Feel free to ask questions at any point during the presentation!



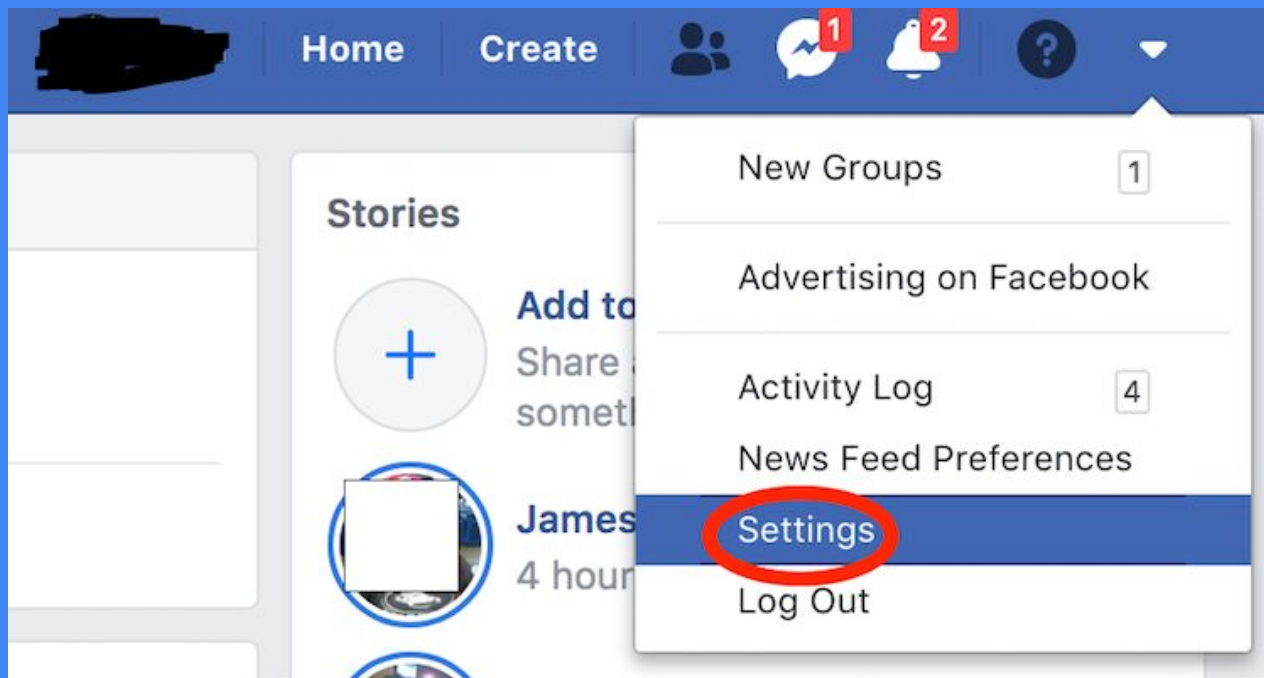
Image and Audio Information

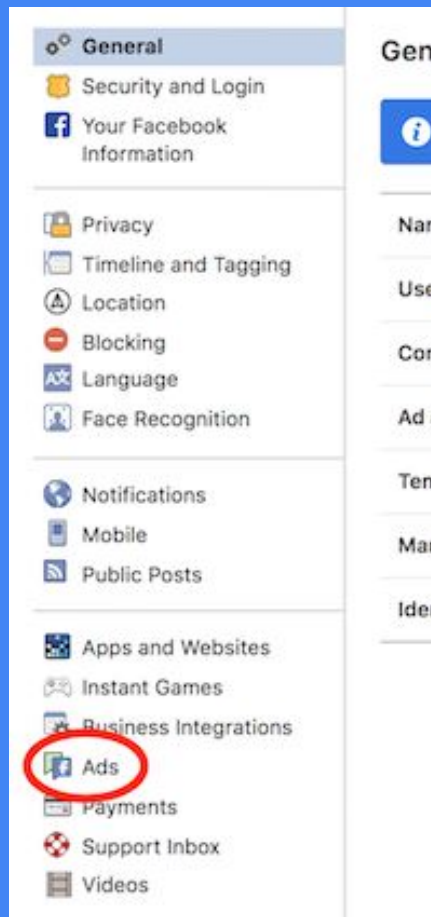
We may collect information about the images and audio that are a part of your User Content, such as identifying the objects and scenery that appear, the existence and location within an image of face and body features and attributes, the nature of the audio, and the text of the words spoken in your User Content. We may collect this information to enable special video effects, for content moderation, for demographic classification, for content and ad recommendations, and for other non-personally-identifying operations. We may collect biometric identifiers and biometric information as defined under US laws, such as faceprints and voiceprints, from your User Content. Where required by law, we will seek any required permissions from you prior to any such collection.



Your Privacy on Facebook

Find Out:
Settings > Ads >
Your information
> Categories





Your ad preferences

Learn what influences the ads you see and take control over your ad experience.

[Learn about Facebook Ads](#)



Your interests



Advertisers you've interacted with



Your information





Your information

Close ^

About you

Your categories

The categories in this section help advertisers reach people who are most likely to be interested in their products, services, and causes. We've added you to these categories based on information you've provided on Facebook and other activity.

Away from family

Close Friends of Men with a Birthday in 0-7 days

Away from hometown

Birthday in March

Close friends of people with birthdays in a month

US politics (very liberal)

Sales

Education and Libraries

Administrative Services

Facebook access (mobile): smartphones and tablets

Frequent Travelers

Technology early adopters



Downloading Your Data

Facebook: Settings > Your Facebook Information > Download your Information

Google:

<https://support.google.com/accounts/answer/3024190?hl=en>

Instagram: Settings > Privacy and Security > Data download/Request Download



DIY Cybersecurity and Tightening your Privacy

Want to make your life more private?

Follow this “DIY Guide to Feminist Cybersecurity”

<https://hackblossom.org/cybersecurity/>



Issues in Big Data: Ethics and Algorithmic Bias



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Algorithms

- An algorithm is a process of instructions provided, usually for computers to interpret and follow.
 - There is usually an **input**, which is determined by the programmer; then there is a **set of rules** (the algorithm) that help lead to the **output**, or the results.
 - Algorithms can be fairly simple, but they can also be much more complex.
- "**Machine learning**" happens when an algorithm tells a computer to make decisions based on a set of patterns derived from data, instead of following specific predetermined instructions.



“Big Data” Unbounded — Ethical Issues

Some recent controversies:

- [Cambridge Analytica controversy](#): psychological profiles of American voters
- [Racial bias in health algorithms](#): results in reduced access to care for Black people
- [Use of facial recognition](#)
 - [Clearview AI](#): sells facial recognition “services”
 - [Case of Robert Williams](#): wrongfully arrested
 - [Machine Bias](#): Software used to predict future criminals, biased against Black men
 - Stanford study creates AI that can [predict sexual orientation based on a photo](#) with up to 91% accuracy



It's not all bad . . .

- Prof. Lazar and NetSI researchers, at Northeastern, [working on COVID-19](#)
- Algorithms predicting the likelihood of cancer ([Breast cancer](#), [Prostate cancer](#))
- [Allegheny County PA “family screening tool”](#) to support human screeners in the Department of Children, Youth, and Families



Identifying Bias: Some Guiding Questions

- In what way(s) is the software used in each scenario biased?
- Do technology and big data-driven solutions **eliminate** human bias or **amplify** it?
- What can be done to decrease bias and improve data-driven decision-making software?



The takeaway?

Algorithms are NOT neutral!



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Algorithms and Bias Activity



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Activity: Data deciding dog adoption

You will be assigned into small groups. You work for an adoption agency and have to decide if someone can adopt a dog. On your handouts, please read the four previous adoption applications and decide if the new adoption applicant can adopt or not.

Do you think this new applicant should be allowed to adopt a dog? Why or why not?



Let's Discuss

Please elect one representative from your group to explain your group's responses to the following questions:

- Would you ACCEPT or REJECT their application? Why?
- What questions from the application did you weigh more? Why?
- What might be some implicit biases in this application form, the process, and in your choices?



Moving Forward

How can we be cognizant
of 'big data' & algorithms
in our research?



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*

Questions Researchers Must Ask

- What **information** is being collected and from where? To whom does this data **belong**?
- How is it being **collected**? Do **participants** know that it is collected, how it will be collected, and how will it be used?
- **How** will the data be analyzed? What **biases** and **ideologies** may be implicit in this analysis?
- Who will this research impact? Who will it **benefit**? Who will it potentially **harm**?



Thank you!

If you have any questions, contact DITI at nulab.info@gmail.com

Developed by DITI Research Fellows:

**Tieanna Graphenreed, Vaishali Kushwaha, Cara Messina, Yana Mommadova,
Garrett Morrow, and Claire Tratnyek**

Slides, handouts, and data available at

<https://bit.ly/diti-fa21-marshall-data-ethics>

Schedule an appointment with us! <https://calendly.com/diti-nu>



Northeastern University
NULab for Texts, Maps, and Networks

*Feel free to ask questions at any point
during the presentation!*