

Introduction to Digital Trace Data: Quality, ethics, and analysis

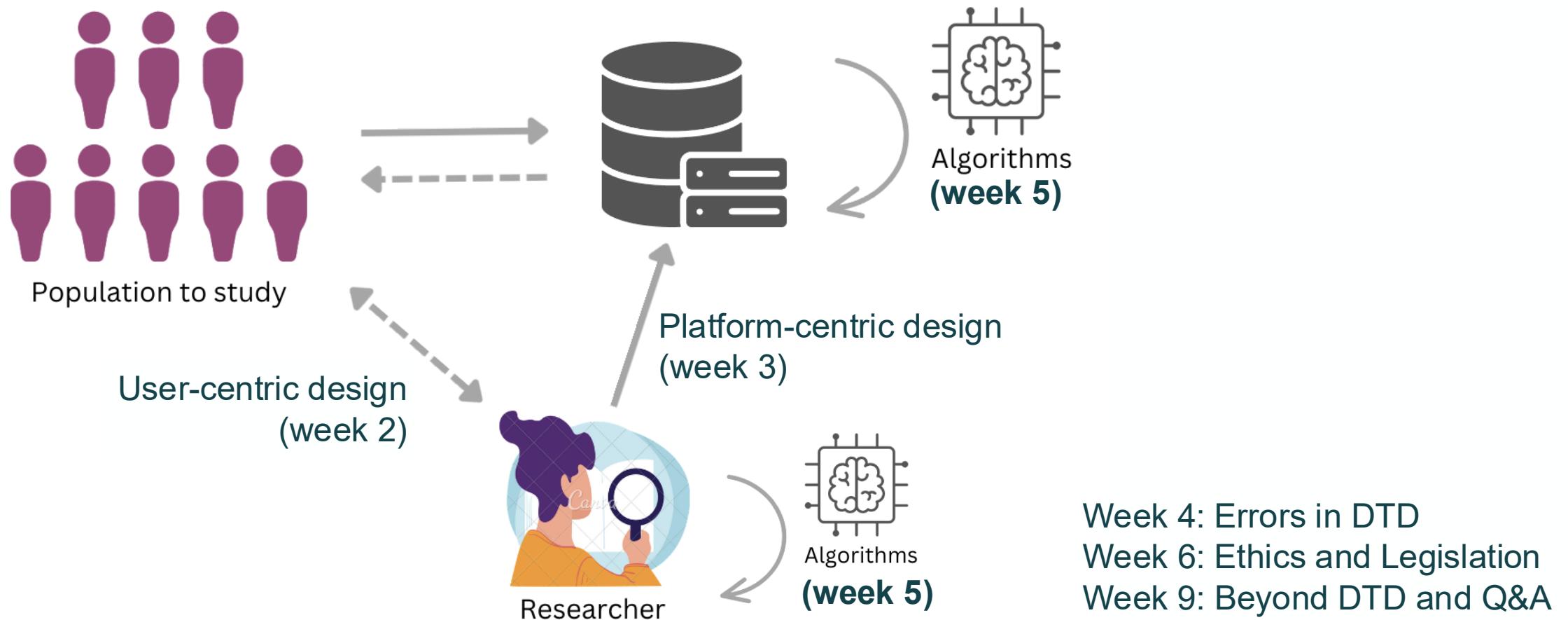
Lecture 5: The role of AI in DTD

Javier Garcia-Bernardo

Assistant Professor

Department of Methodology and Statistics

Where are we?



Today's material

1. What are Algorithms/AI/Machine Learning?
2. Using ML to study societies
3. The impact of biased ML on Digital Trace Data
4. The impact of ML on societies
5. Dealing with bias in ML

TODAY

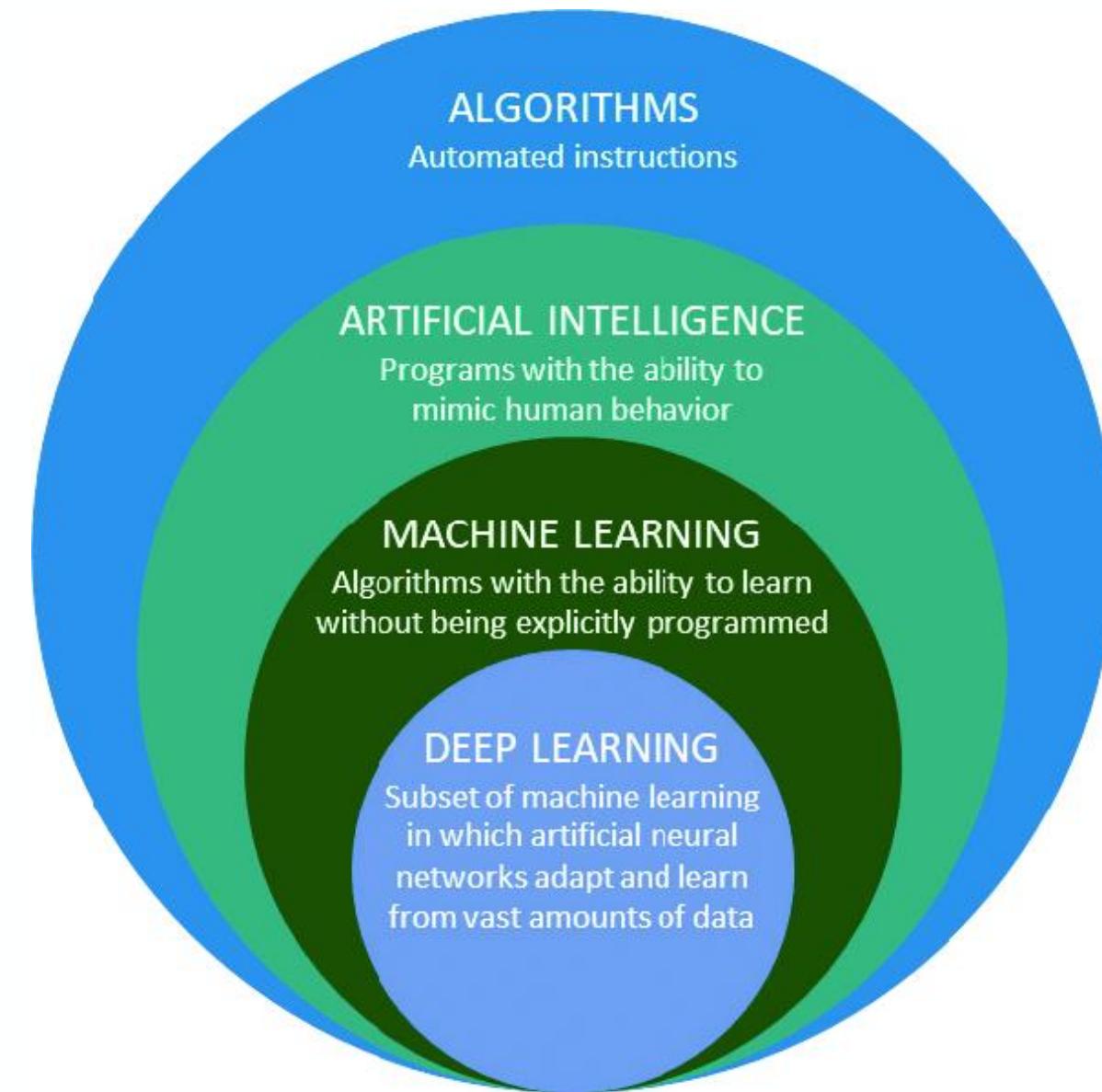
Lecture

1. Explain machine learning in your own words
2. Explain *why* machine learning models may be biased.
3. Understand the effects of ML on DTD and society.
4. Assess bias in ML models

Lab

- Apply a ML model to text data
- Audit a ML model

1. What is machine learning?



Machine Learning

"A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P** if its performance at task T , as measured by P , improves with experience E ." (Samuel/Mitchell, 1959/1997)

- **Experience: Data** (e.g. comments from TikTok)
- **Task: Goal of the model** (e.g. predict hate speech)
- **Performance measure: Accuracy, R^2 , etc**

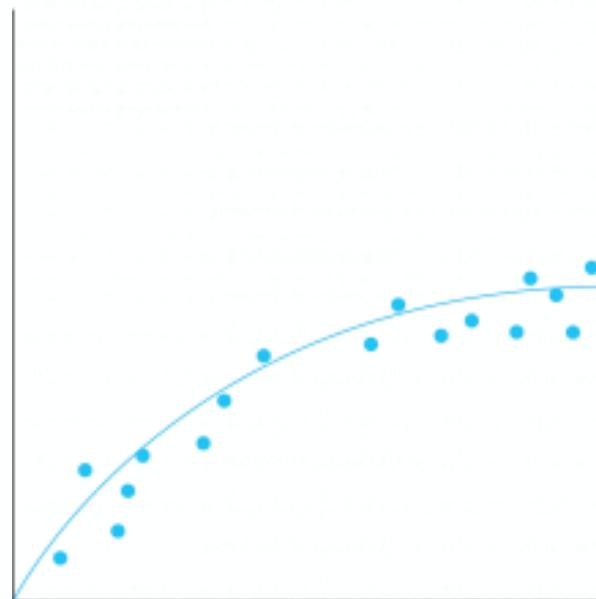
Is a linear regression a machine learning model?

Supervised vs unsupervised ML

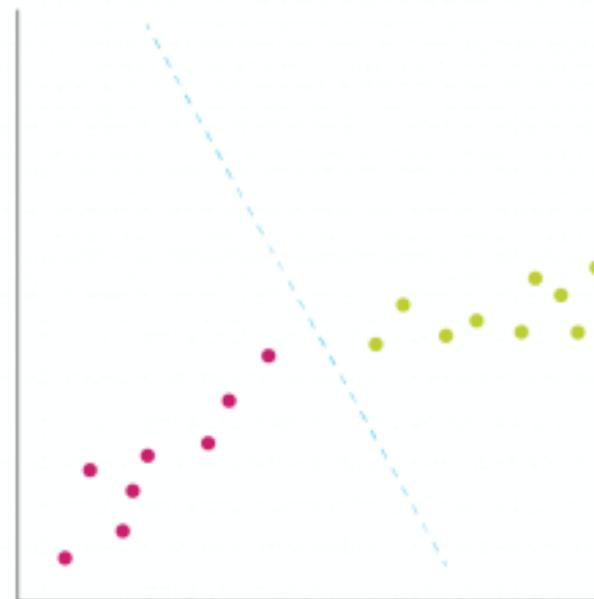
Supervised ML: We have inputs (features, independent variables) and an output (target, dependent variable)

Unsupervised ML: We have inputs and (mostly) try to find groups

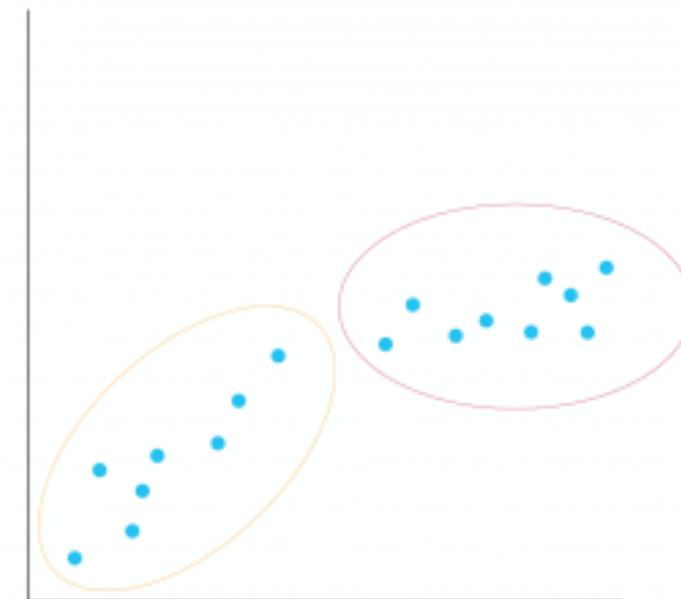
Regression



Classification



Clustering



Examples

(A) You use at an advertisement company and want to customers into segments based on their purchasing behaviours, age, maximal educational degree attained, etc.

→ Unsupervised

(B) You work at a bank and want to develop a model that helps them predicting which loan applicants will default (not be able to pay the loan) based on their financial transactions.

→ Supervised (classification)

(C) You use news and social media analytics to predict changes in the stock market. You will use a small number of stocks as target indicators, and web-scraped text from social media and the news. You have access to previous instances of this data, and you want to predict the values for your indicators in the near future.

→ Supervised (regression)

Using ML to understand societies

- **Description:** The goal is to describe patterns or groupings in historical data.
- **Prediction:** The goal here is to predict outcomes. For example, predict missing data, predict risk of developing diseases, or label data.
- **Explain:** The goal here is to understand the causal relationships in the data to influence the outcomes we care about.

In data analysis: Descriptive (unsupervised ML)

RESEARCH ARTICLE

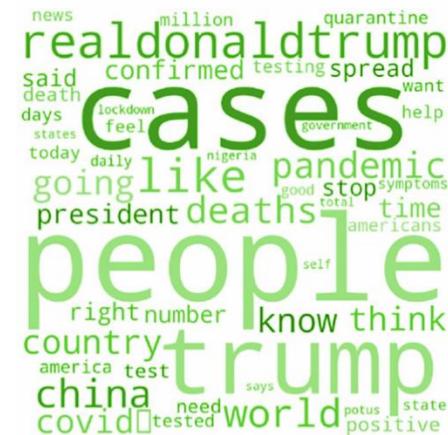
Framing COVID-19: How we conceptualize and discuss the pandemic on Twitter

Philipp Wicke¹*, Marianna M. Bolognesi²

- **Question:** What is the framing of the COVID pandemic? Framing of WAR (fight, combat, battle), STORM (wave, storm, cloud), MONSTER (evil, horror, killer) or TSUNAMI (wave, tragedy, catastrophe).
 - **Data:** Twitter around #Covid-19 (80 hashtags)



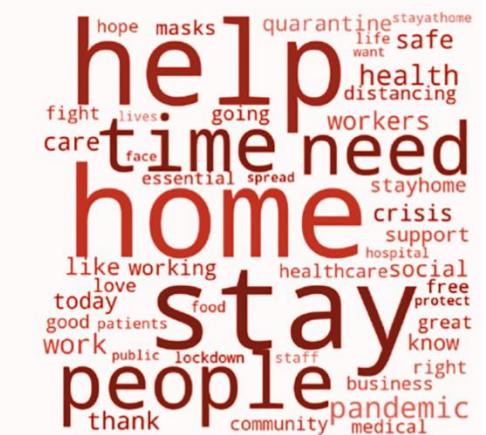
Topic #I: Communications and Reporting



Topic #III: Politics



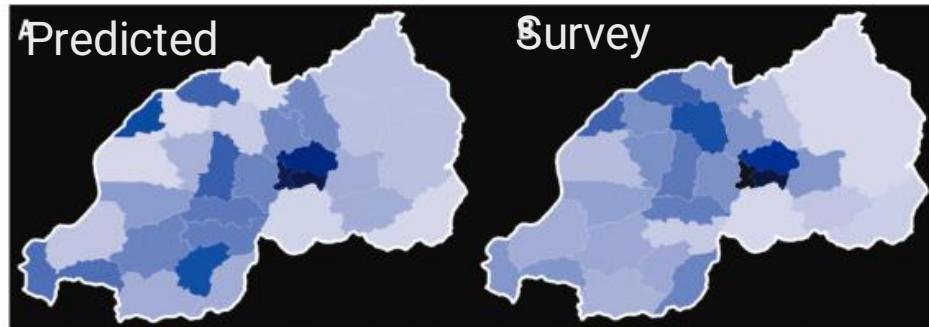
Topic #II: Community and Social Compassion



Topic #IV: Reacting to the epidemic

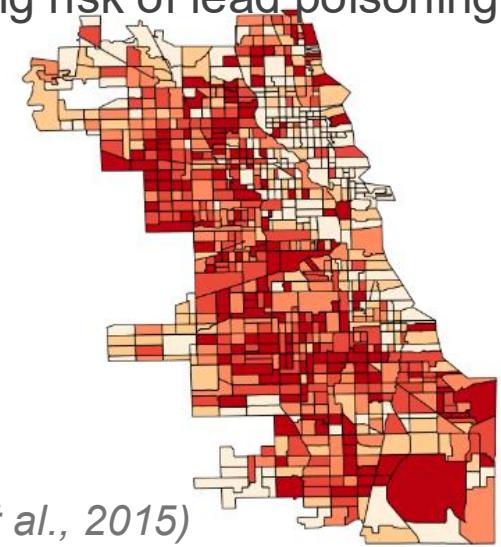
In data analysis: Prediction (supervised ML)

Predicting wealth/SES:

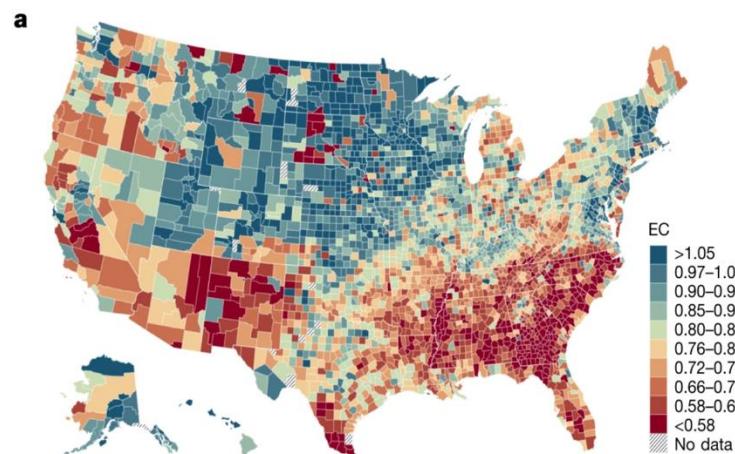


(Blumenstock et al., 2015)

Predicting risk of lead poisoning:

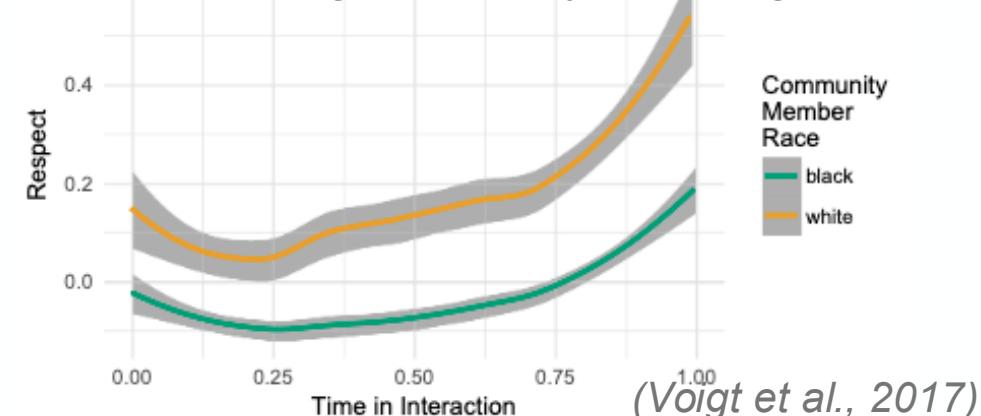


(Potash et al., 2015)



(Chetty et al., 2022)

Predicting respect by police agents



(Voigt et al., 2017)

2. Why are ML models often biased*?

*The model performance is different for different subgroups

Errors in algorithms

We need to consider the errors of every algorithm:

- How often they fail?
- For whom do they fail? (*bias*)

No model will be perfect, but we need to understand when and for whom they fail.

Fitting and using ML models (supervised ML)

1) Training data



2) Model



3) New data

X: Flights
to Russia

Y: Criminal



3



12



1



80



If $X > 10$: criminal
If $X \leq 10$: not criminal

X: Flights
to Russia



20



6



3

How often do they fail? The confusion matrix

	Predicted criminal	Predicted not criminal
Criminal	True positive	False negative
Not criminal	False positive	True negative

For whom do they fail?

Group A	Predicted criminal	Predicted not criminal
Criminal	10	10
Not criminal	1	100

Group B	Predicted criminal	Predicted not criminal
Criminal	10	1
Not criminal	10	100

Exercise (in pairs)

You work at a Dutch bank and you want to develop a model that will help them predict which loan applicants will default (not be able to pay the loan) based on their financial transactions. You have labeled data from customers in Utrecht for the last three years. You want to predict default for all new loan applicants.

In which parts of the process may bias be introduced? Think about

- The representativeness of the training data
- The quality of the outcome (default/non-default) in the training data for different subpopulations
- The quality of the features (financial transactions) in the training data for different subpopulations
- The machine learning pipeline

Five main sources of bias in a ML model

Sample bias: the training data does not generalize to the prediction data.

X: Flights to Russia	Y: Criminal
3	
12	
1	
80	

Your sample may be only Dutch people, who don't go to Russia often (unless they are criminals in this example). This could happen because:

- you are getting the data from a Dutch governmental body only tracking Dutch people
- you are using a script that cannot handle Russian names, and they are not merged properly in the process
- you collected data during COVID times, or 10 years ago (drift)

Five main sources of bias in a ML model

ML pipeline bias: e.g., you are using a model or performance metric that focuses mostly on the majority class; there are errors in your code; you used a wrong model.

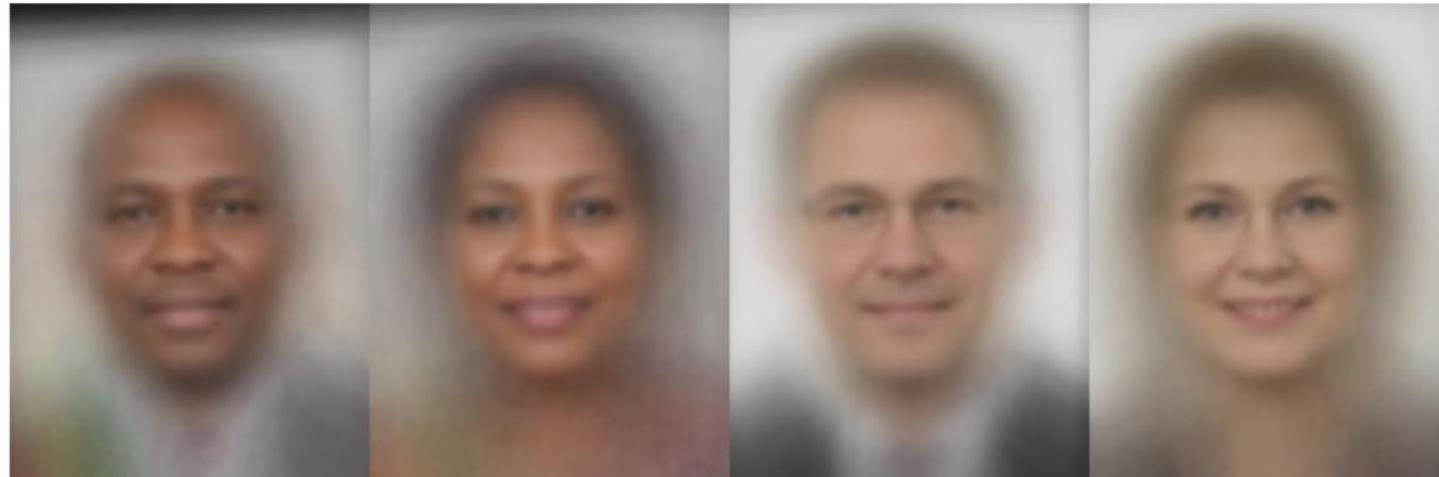
Features bias: the features (flights to Russia) have different meaning from some subpopulation.

Outcome (label) bias: the label (criminal/not criminal) has different meaning from some subpopulation. For example if police only looks for criminals among non-Dutch, the label “non-criminal” does not for Dutch people mean something else (lack of policing, not lack of crime). Or if humans creating the labels are themselves biased.

Application bias: the model and data is right, but it is applied in a bias way. E.g. a manager always trust a ML system for firing foreigners, but ignore it for Dutch.



Gender Classifier	Darker Male	Darker Female	Lighter Male	Lighter Female	Largest Gap
Microsoft	94.0%	79.2%	100%	98.3%	20.8%
FACE++	99.3%	65.5%	99.2%	94.0%	33.8%
IBM	88.0%	65.3%	99.7%	92.9%	34.4%



What type of bias?

Dissecting racial bias in an algorithm used to manage the health of populations

ZIAD OBERMEYER  , BRIAN POWERS, CHRISTINE VOGELI, AND SENDHIL MULLAINATHAN  [Authors Info & Affiliations](#)

SCIENCE • 25 Oct 2019 • Vol 366, Issue 6464 • pp. 447-453 • DOI: 10.1126/science.aax2342

 139,001  1,266



Racial bias in health algorithms

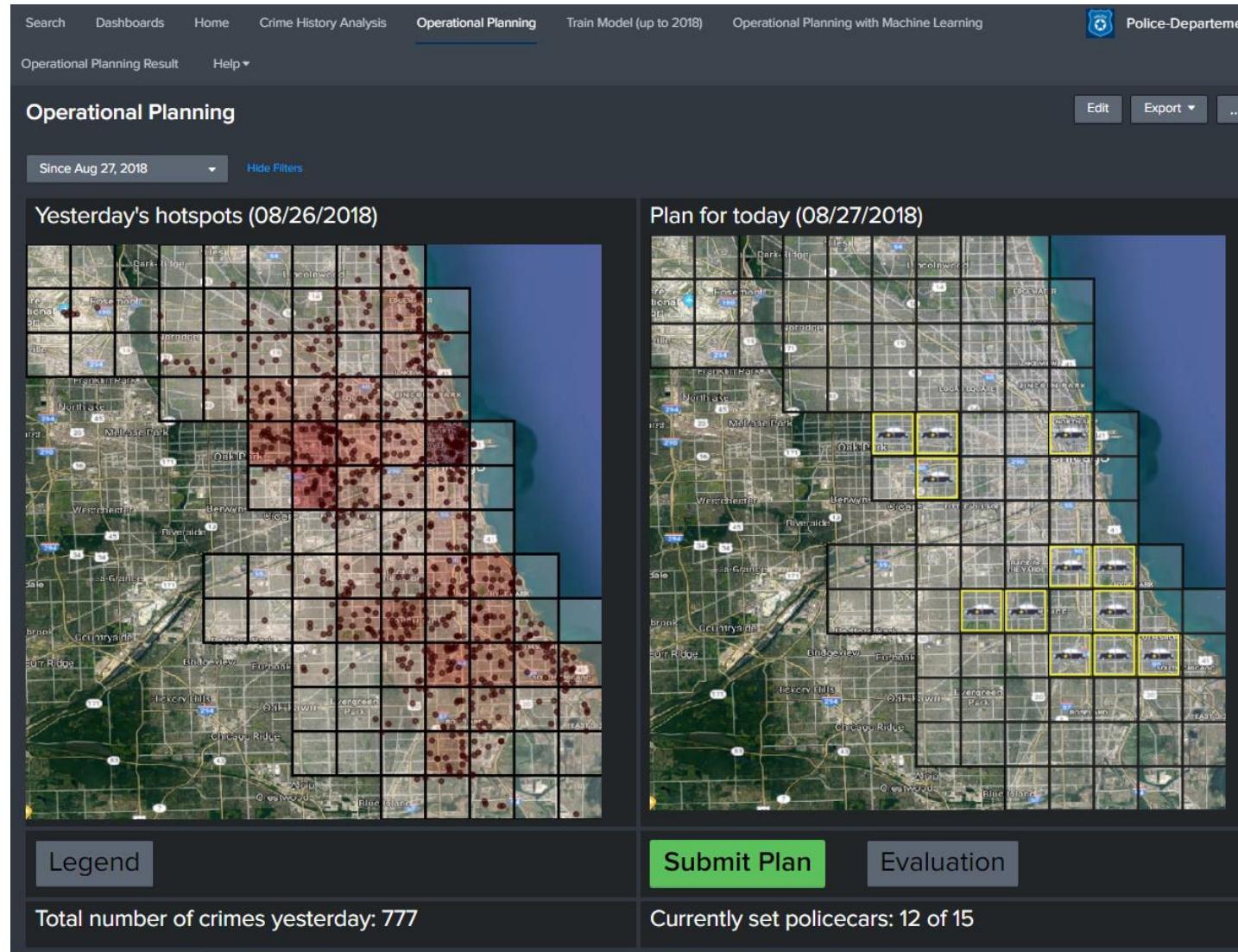
The U.S. health care system uses commercial algorithms to guide health decisions. Obermeyer *et al.* find evidence of racial bias in one widely used algorithm, such that Black patients assigned the same level of risk by the algorithm are sicker than White patients (see the Perspective by Benjamin). The authors estimated that this racial bias reduces the number of Black patients identified for extra care by more than half. Bias occurs because the algorithm uses health costs as a proxy for health needs. Less money is spent on Black patients who have the same level of need, and the algorithm thus falsely concludes that Black patients are healthier than equally sick White patients. Reformulating the algorithm so that it no longer uses costs as a proxy for needs eliminates the racial bias in predicting who needs extra care.



What type of bias?

Predictive policing

“Quickly connect the dots, identify hidden patterns and discover trends” (Splunk)



What type of bias?

Judge Rules \$400 Million Algorithmic System Illegally Denied Thousands of People's Medicaid Benefits

Thousands of children and adults were automatically terminated from Medicaid and disability benefits programs by a computer system that was supposed to make applying for and receiving health coverage easier.

“The system often doesn’t load the appropriate data, assigns beneficiaries to the wrong households, and makes incorrect eligibility determinations”

What type of bias?

Humans are biased too

Case Study: Blind Auditions & Women in Orchestras (U.S.)

Goldin & Rouse, American Economic Review (2000); NBER Working Paper (1997)

Policy shift: Most major U.S. orchestras adopted 'blind' auditions (with a screen) in the 1970s–1980s; Boston Symphony used screened prelims as early as 1952.

Measured effect: Using actual audition records, a screen increased the probability a woman advanced from prelim rounds by ~50%, and raised the chance she won the final round by severalfold.

Contribution to change: Blind auditions explain ~30–55% of the increase in the female share among new hires and ~25–46% of the increase in the overall female share of orchestras since 1970.

Contextual outcome: In the top five U.S. orchestras, women were <5% of players in 1970 versus ~25% by the mid-1990s.

Sources: Goldin & Rouse (1997) *NBER Working Paper 5903*; Goldin & Rouse (2000) *AER, "Orchestrating Impartiality."*

- *IN this case, this creates label bias (success=being hired) does not reflect the same for different groups*

Causal loop diagram

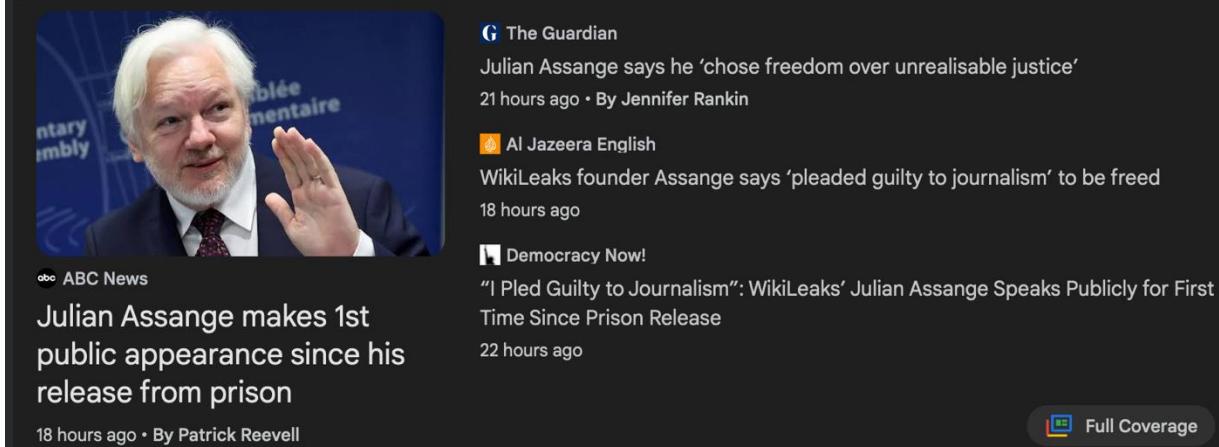
How different complex, interdependent issues are causally interrelated

3a. How does machine learning affect DTD?

From the point of view of data quality and data analysis

ML is used in the collection, processing and interpretation of DTC

Main uses of ML in Digital Trade Data: Recommendation algorithms



The screenshot shows a news feed with several articles:

- G The Guardian**: Julian Assange says he 'chose freedom over unrealisable justice' (21 hours ago, By Jennifer Rankin)
- Al Jazeera English**: WikiLeaks founder Assange says 'pledged guilty to journalism' to be freed (18 hours ago)
- Democracy Now!**: "I Pled Guilty to Journalism": WikiLeaks' Julian Assange Speaks Publicly for First Time Since Prison Release (22 hours ago)

ABC News: Julian Assange makes 1st public appearance since his release from prison (18 hours ago, By Patrick Reevell) - This article includes a photo of Assange waving.

Full Coverage

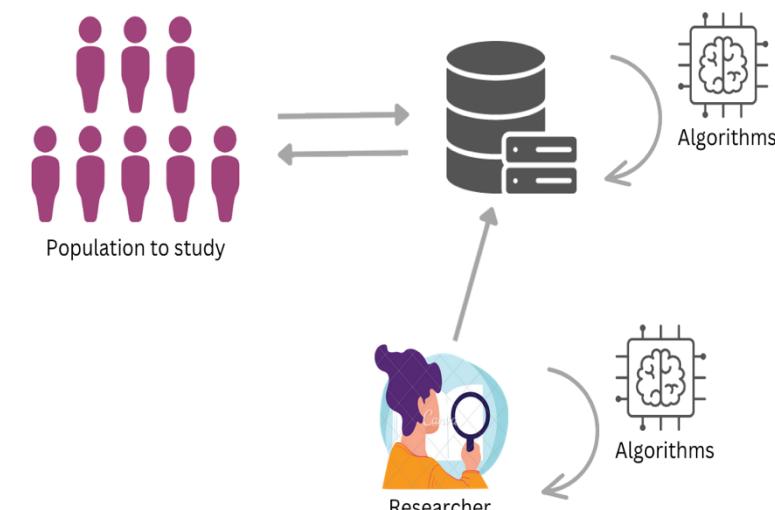
Data interpretation Data augmentation

```
1,  
  "genderInfo" : {  
    "gender" : "male"  
  }  
  
"inferredAgeInfo" : {  
  "age" : [  
    ">50"  
  ],  
  "birthDate" : ""  
}
```

```
[  
  {  
    "name" : "Rap",  
    "isDisabled" : false  
  },  
  {  
    "name" : "Retired life",  
    "isDisabled" : false  
  },  
  {  
    "name" : "Rom-com films",  
    "isDisabled" : false  
  },  
  {  
    "name" : "Drama",  
    "isDisabled" : false  
  }]
```

Errors introduced by ML

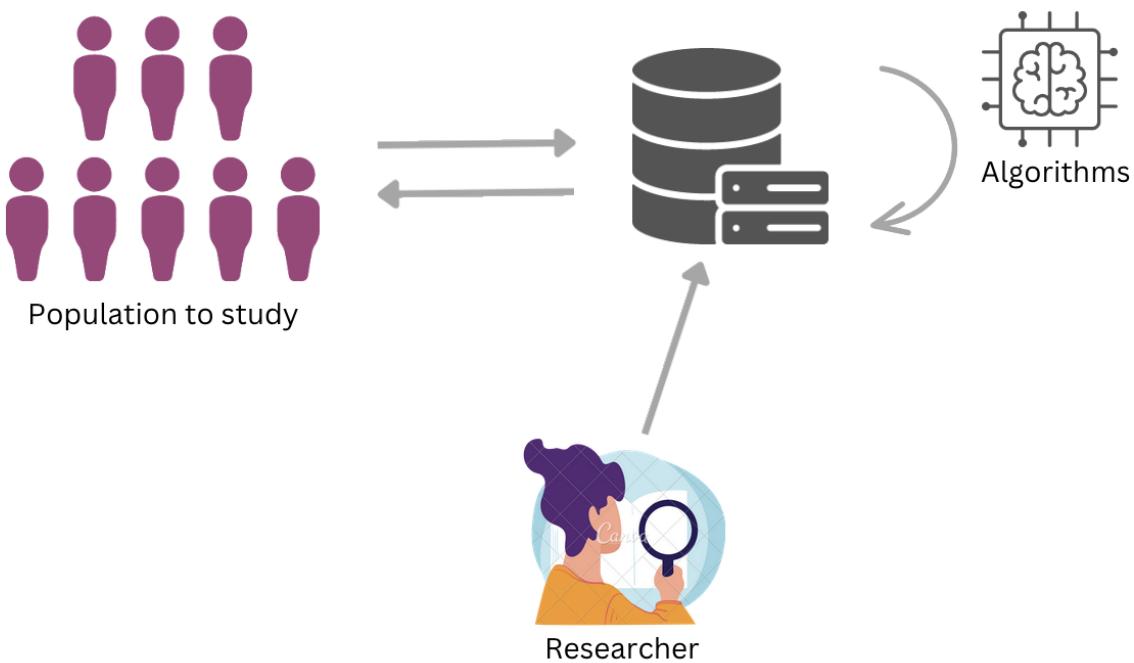
- Total error = measurement error + representation error
- ML impacts mostly measurement error (from lecture 4):
 - On a conceptual level: **validity**
 - Data can be algorithmically confounded (recommendation algorithms)
 - Our measurement will include both human behavior and the influence of the algorithm
 - **Processing error:**
 - When using ML to process data
 - A ML trained using biased (human) data will typically be biased



Problems with validity

Facebook uses the “clustering coefficient” to recommend friends: e.g., if you have two friends, Sanne and Joep, that are not Facebook friends, Facebook will suggest Sanne and Joep to add each other as friends.

Your measurement of social closure (clustering coefficient) is measuring *both* social closure and the effect of the algorithm.

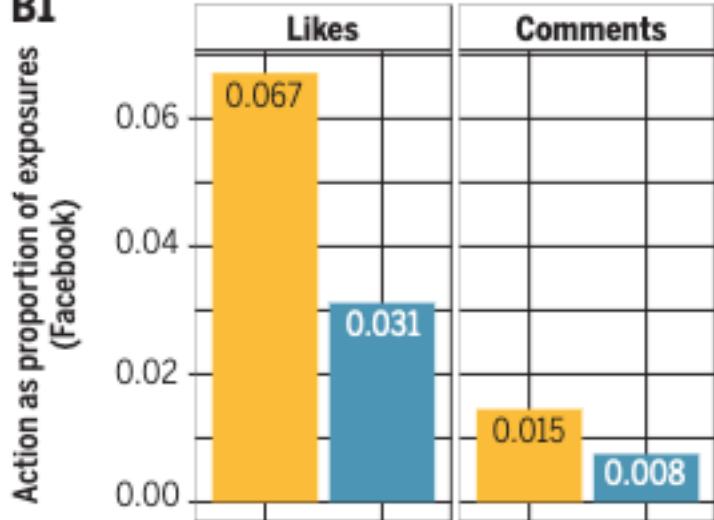


Problems with validity

More diverse sources!

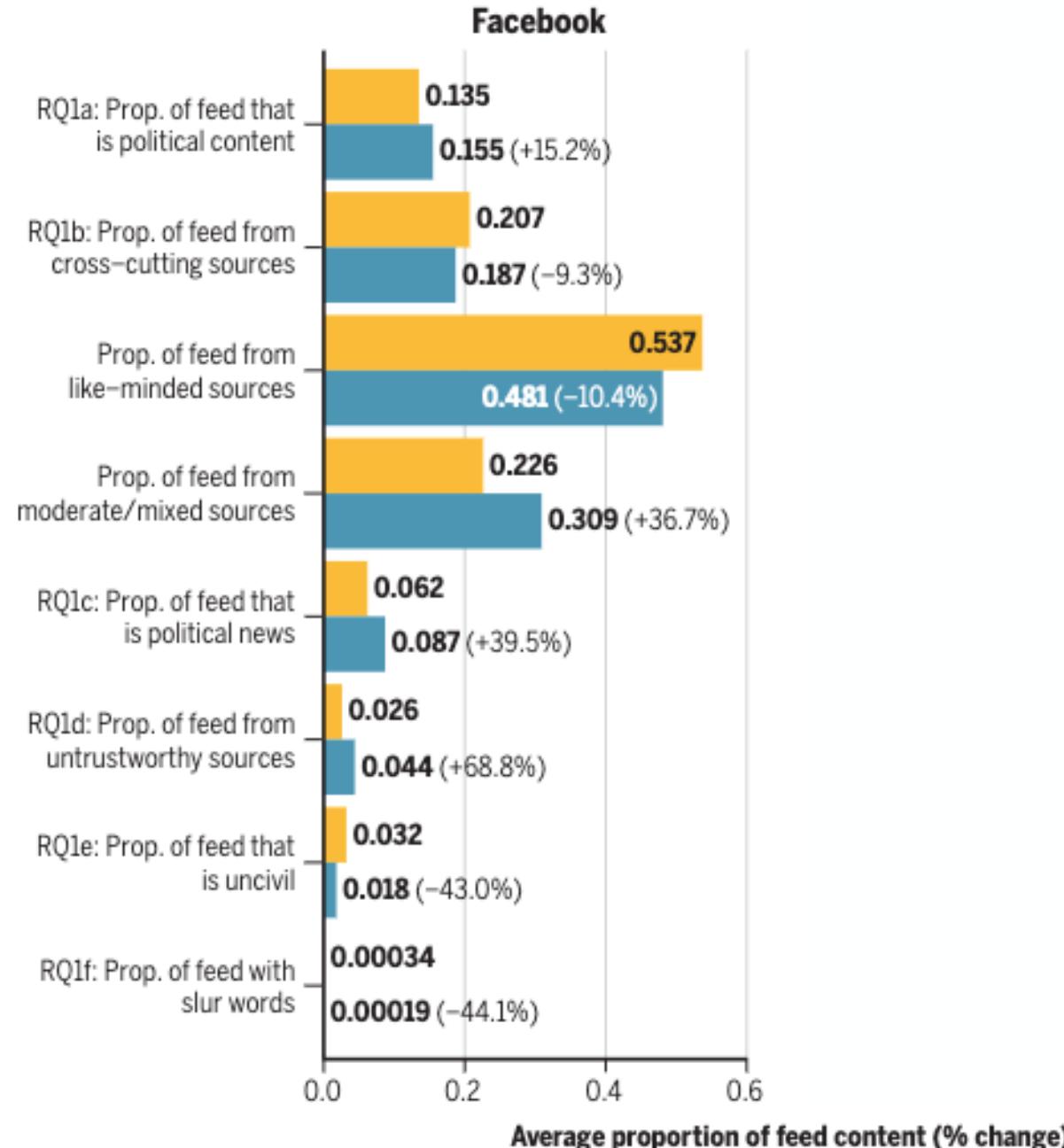
Algorithmic feed Chronological feed

B1



People less engaged!

How do social media feed algorithms affect attitudes and behavior in an election campaign?
Guess et al., 2023, Science



In your group project

Think about how ML may be affecting:

- The text data you collected (validity)
- The labels you infer from the text (processing error)

3b. How does machine learning affect societies?

When machine learning models are being used to make decisions, they cannot be separated from the social and ethical context in which they are applied (*Kit T. Rodolfa, Pedro Saleiro, and Rayid Ghani, Big Data book*)

- a) Through errors in algorithms
- b) Through feedback loops
- c) By reinforcing power structures

A) Errors in algorithms

We need to consider the errors of every algorithm:

- How often they fail?
- For whom do they fail? (**bias**)

Remember there are people behind the data:

- What are the costs of those failures?
- What are the long-term effects? (feedback effects)

Things can go horribly wrong

ML is used in many crucial areas for human wellbeing:

- Who to hire – CV screening
- Who to promote – performance reviews
- Who to jail – predictive policing
- Who to kill – “we kill people based on metadata” (based on similarity with somebody who was labeled as an enemy)

Those algorithms are often (1) “opaque,” (2) “beyond dispute or appeal,” and (3) disproportionately impact the underprivileged (Cathy O’Neal)

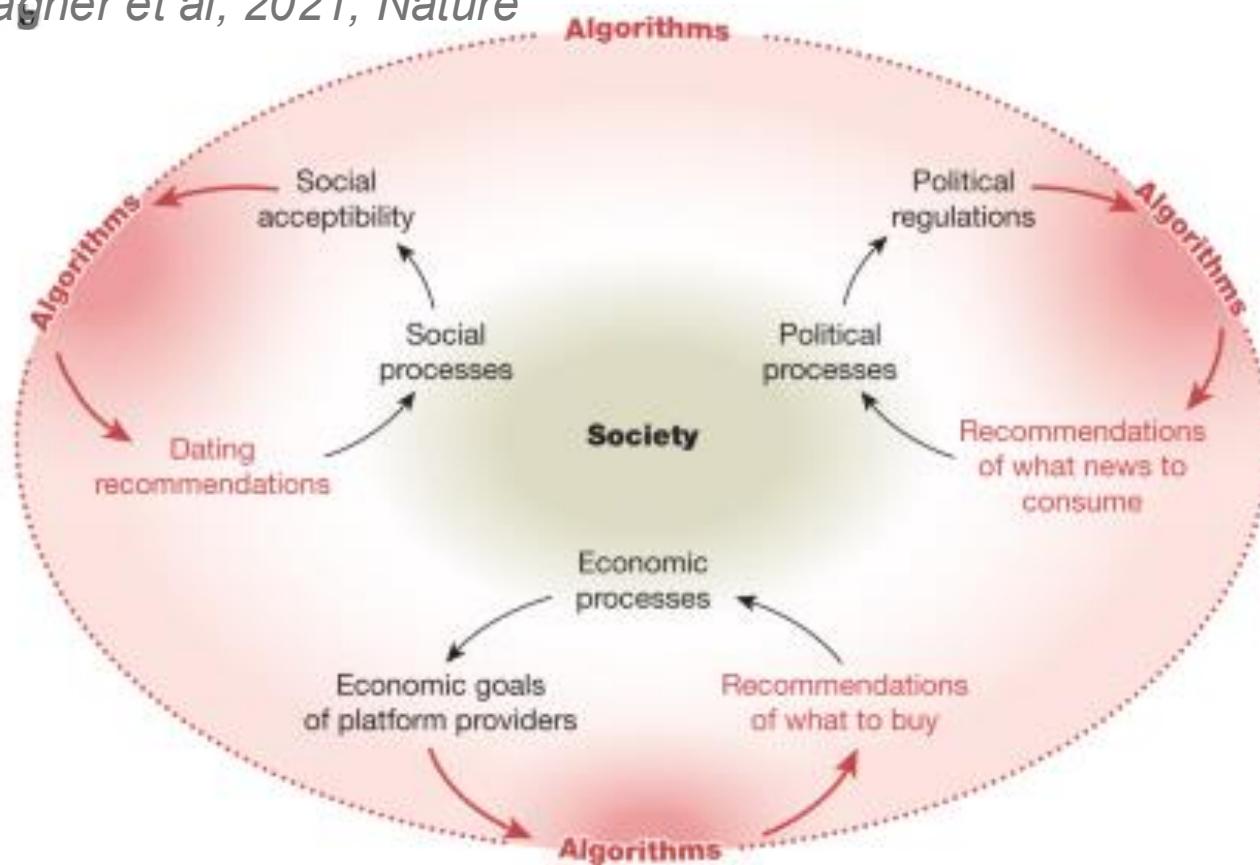


Cathy O’Neil

Mathematician (Harvard/MIT/Barnard College)
Worked four years in finance and advertisement

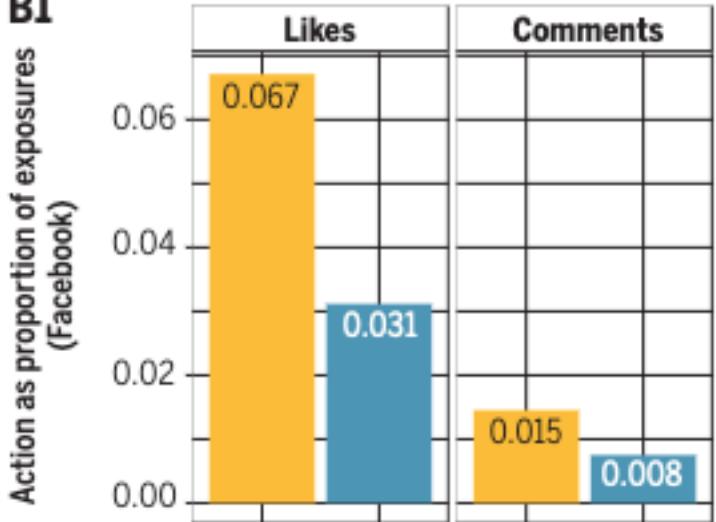
B) Algorithms create feedback effects

Measuring algorithmically infused societies
Wagner et al, 2021, Nature



How do social media feed algorithms affect attitudes and behavior in an election campaign?
Guess et al., 2023, Science

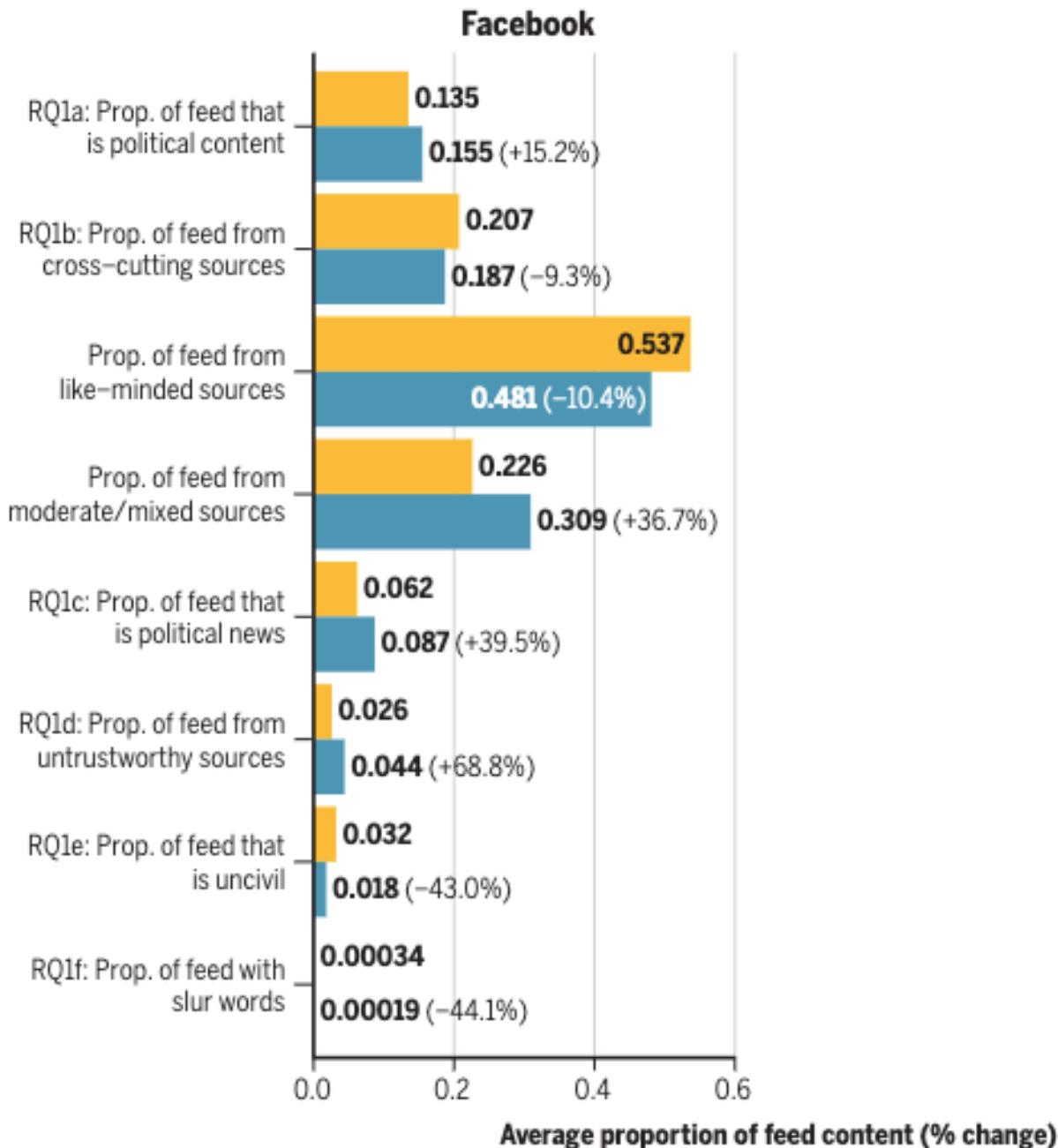
B1



People less engaged!

More diverse sources!

● Algorithmic feed ● Chronological feed

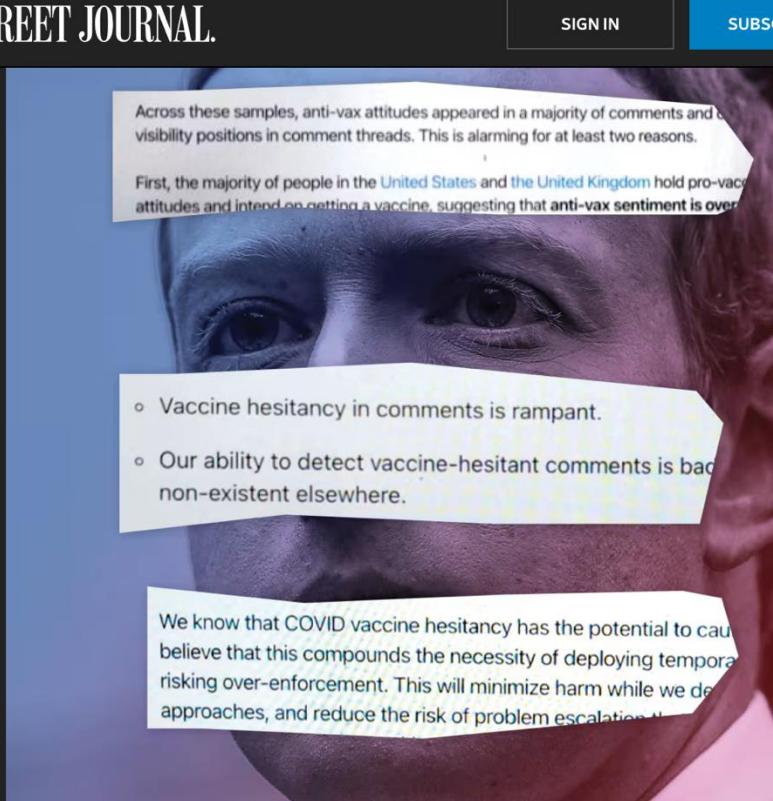




the facebook files □

How Facebook Hobbled Mark Zuckerberg's Bid to Get America Vaccinated

Company documents show antivaccine activists undermined the CEO's ambition to support the rollout by flooding the site and using Facebook's own tools to sow doubt about the Covid-19 vaccine



Facebook told the White House to focus on the ‘facts’ about vaccine misinformation. Internal documents show it wasn’t sharing key data.

The tech giant meticulously tracked how misleading medical information spread — but didn’t tell policymakers, even as they demanded it do so.

<https://www.washingtonpost.com/technology/2021/10/28/facebook-covid-misinformation>

https://www.wsj.com/articles/facebook-mark-zuckerberg-vaccinated-11631880296?mod=article_inline

Upvoting extremism: Collective identity formation and the extreme right on Reddit

Tiana Gaudette  , Ryan Scrivens , [...], and Richard Frank  [View all authors and affiliations](#)

Volume 23, Issue 12 | <https://doi.org/10.1177/1461444820958123>

 [Contents](#)

 [Get access](#)

Abstract

Since the advent of the Internet, right-wing extremists and those who subscribe to extreme right views have exploited online platforms to build a collective identity among the like-minded. Research in this area has largely focused on extremists' use of websites, forums, and mainstream social media sites, but overlooked in this research has been an exploration of the popular social news aggregation site Reddit. The current study explores the role of Reddit's unique voting algorithm in facilitating "othering" discourse and, by extension, collective identity formation among members of a notoriously hateful subreddit community, r/The_Donald. The results of the thematic analysis indicate that those who post extreme-right content on r/The_Donald use Reddit's voting algorithm as a tool to mobilize like-minded members by promoting extreme discourses against two prominent out-groups: Muslims and the Left. Overall, r/The_Donald's "sense of community" facilitates identity work among its members by creating an environment wherein extreme right views are continuously validated.

Simil



C) ML reinforces power structures



Meredith Whittaker

- Employed by Google for 13 years
- Research Professor at New York University
- Co-founder and faculty director of the AI Now Institute.
- President of Signal

*"Private computational systems marketed as artificial intelligence (AI) are threading through our public life and institutions, **concentrating industrial power, compounding marginalization, and quietly shaping access to resources and information**" (in *The Steep Cost of Capture*, 2021)*

"The commoditization of our data enables an asymmetric redistribution of power that is weighted toward the actors who have access and the capability to make sense of information" (Sarah Myers West, 2017)

Even if the algorithms are unbiased!

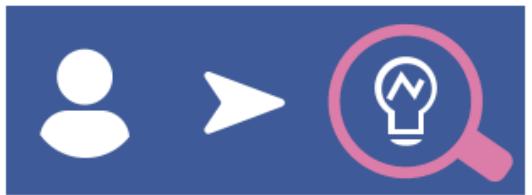
We understand social media platforms as ways to share and see content.

Private companies influence societal outcomes by controlling information flows and target ads and services.

How was Facebook users' data misused?

1

In 2014 a Facebook quiz invited users to find out their personality type



2

The app collected the data of those taking the quiz, but also recorded the public data of their friends



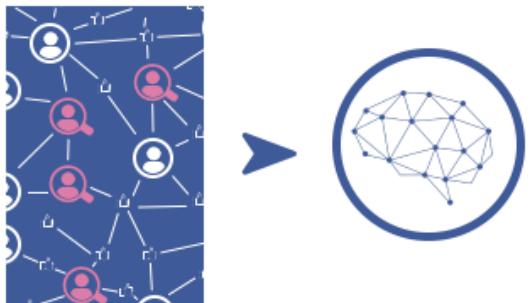
3

About 305,000 people installed the app, but it gathered information on up to 87 million people, according to Facebook



4

It is claimed the data was sold to Cambridge Analytica (CA), which used it to psychologically profile voters in the US



Data feminism

Data is never neutral: And data are never neutral; they are always the biased output of unequal social, historical, and economic conditions:
Data feminism focused on power: who has it, who doesn't, how power shapes data collection, interpretation, and use.

Who is doing the work of data science (and who is not)?

- AI labs dominated by white men from elite institutions → limited perspectives and privilege hazard: those with privilege often don't see how data systems disadvantage others.

- e.g. *It was black people who realized that facial recognition did not work on black people. Why? They were not designing it*

"The biggest threat from artificial intelligence systems is not that they will become smarter than humans, but that they will hard-code sexism, racism, and other forms of discrimination into the digital infrastructure of our societies." — Kate Crawford

Whose goals are prioritized in data science (and whose are not)?

Why do we collect data on how to predict crime, instead of how to support those in need?

e.g., until recently, crash test dummies were designed in the size and shape of men, an oversight that meant that women had a 47 percent higher chance of car injury than men.

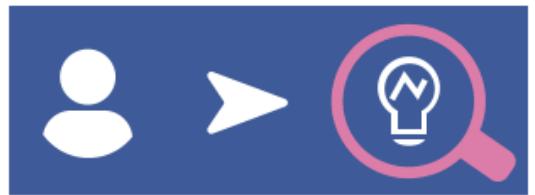
And who benefits from data science (and who is either overlooked or actively harmed)?

- Algorithms often surveil the poor more than the rich (Allegheny County child welfare risk scores).
- Target's pregnancy model exposed a teenager's private health information → profit over privacy.
- Facebook–Cambridge Analytica: corporate/political gain vs. public trust.

How was Facebook users' data misused?

1

In 2014 a Facebook quiz invited users to find out their personality type



2

The app collected the data of those taking the quiz, but also recorded the public data of their friends



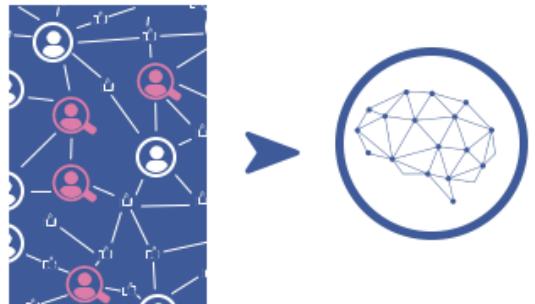
3

About 305,000 people installed the app, but it gathered information on up to 87 million people, according to Facebook



4

It is claimed the data was sold to Cambridge Analytica (CA), which used it to psychologically profile voters in the US



Whose goals are prioritized (and whose are not)?

And who benefits (and who is either overlooked or actively harmed)?

Cambridge Analytica infographic

Harm to society is often accepted as part of the business model

THE WALL STREET JOURNAL.

the facebook files



Facebook Knows Instagram Is Toxic for Teen Girls, Company Documents Show

Its own in-depth research shows a significant teen mental-health issue that Facebook plays down in public



For years they had focus groups, online surveys, diary studies - so this was not one chance finding.

- *A 2019 presentation slide said: "We make body-image issues worse for one in three teenage girls"*
- *Another slide said teenagers blamed Instagram for increased levels of anxiety and depression*
- *Some 13% of UK teenagers and 6% of US users surveyed traced a desire to kill themselves to Instagram*

Instagram response: "Based on our research and feedback from experts, we've developed features so people can protect themselves from bullying, we've given everyone the option to hide 'like' counts and we've continued to connect people who may be struggling with local support organisations."

<https://www.bbc.com/news/technology-58570353>

Performing Platform Governance: Facebook and the Stage Management of Data Relations

Karen Huang¹ · P. M. Krafft²

Received: 2 April 2021 / Accepted: 12 February 2024 / Published online: 4 April 2024

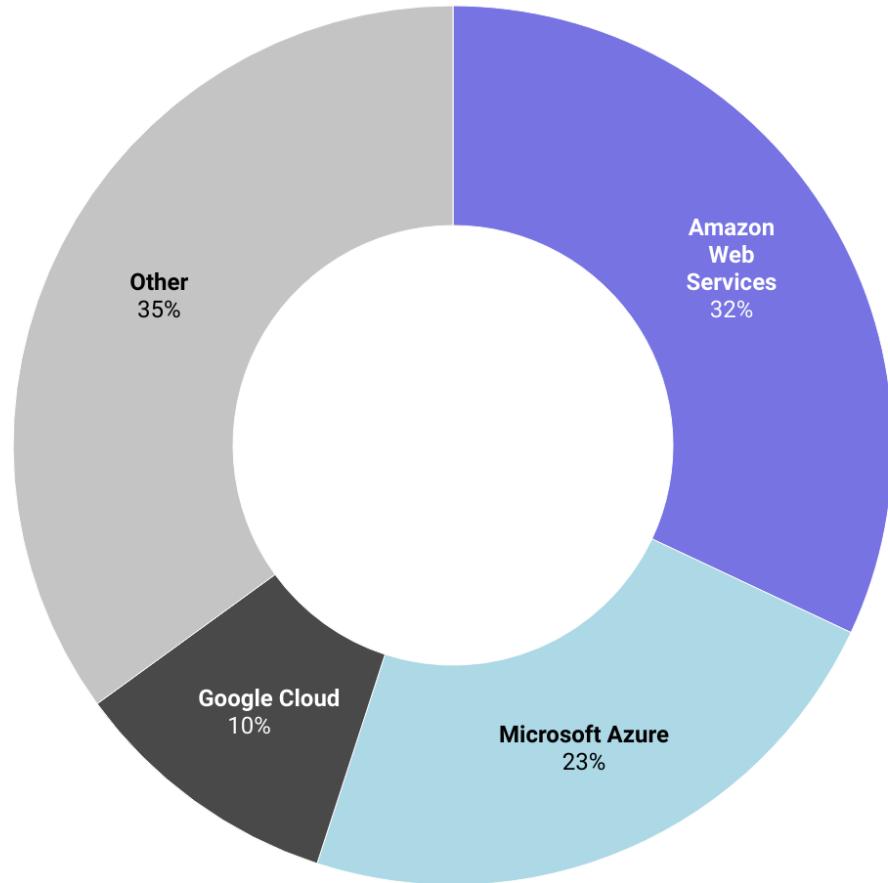
Abstract

Controversies surrounding social media platforms have provided opportunities for institutional reflexivity amongst users and regulators on how to understand and govern platforms. Amidst contestation, platform companies have continued to enact projects that draw upon existing modes of privatized governance. We investigate how social media companies have attempted to achieve closure by continuing to set the terms around platform governance. We investigate two projects implemented by Facebook (Meta)—authenticity regulation and privacy controls—in response to the Russian Interference and Cambridge Analytica controversies surrounding the 2016 U.S. Presidential Election. Drawing on Goffman's metaphor of stage management, we analyze the techniques deployed by Facebook to reinforce a division between what is visible and invisible to the user experience. These platform governance projects propose to act upon *front-stage data relations*: information that users can see from other users—whether that is content that users can see from “bad actors”, or information that other users can see about oneself. At the same time, these projects relegate *back-stage data relations*—information flows between users constituted by recommendation and targeted advertising systems—to invisibility and inaction. As such, Facebook renders the user experience actionable for governance, while foreclosing governance of back-stage data relations central to the economic value of the platform. As social media companies continue to perform platform governance projects following controversies, our paper invites reflection on the politics of these projects. By destabilizing the boundaries drawn by platform companies, we

Concentration of control

Market share of cloud computing providers

■ Amazon Web Services ■ Microsoft Azure ■ Google Cloud ■ Other



Only hegemonic companies have the capital and power to thrive in the new era, reinforcing their power and dominance.

Revenue in 2024:

Amazon: \$638 bn.

Netherlands: \$440 bn. (407 bn. eur)

Apple: \$391 bn.

Alphabet: \$350 bn.

Microsoft: \$262 bn.

Meta: \$164 bn.

(updated from Babic et al., 2017)

Exercise for the practical

Predictive policing is increasingly used (also by the Netherlands). Let's imagine these systems are trained in historical and personal data. The use of personal data for predictive policing has recently been banned in the EU.

Where can biases enter those models? (sample/outcome/features/pipeline/application)

Who are the actors involved? Whose interests are prioritized (and who's not)?

Do you think the benefits outweigh the biases?



What are some things that we can do as individuals?

WHITE COLLAR CRIME RISK ZONES

White Collar Crime Risk Zones uses machine learning to predict where financial crimes are mostly likely to occur across the US. To learn about our methodology, read our white paper.

By Brian Clifton, Sam Lavigne and Francis Tseng for *The New Inquiry Magazine*, Vol. 59: ABOLISH.

Enter a Location

Search

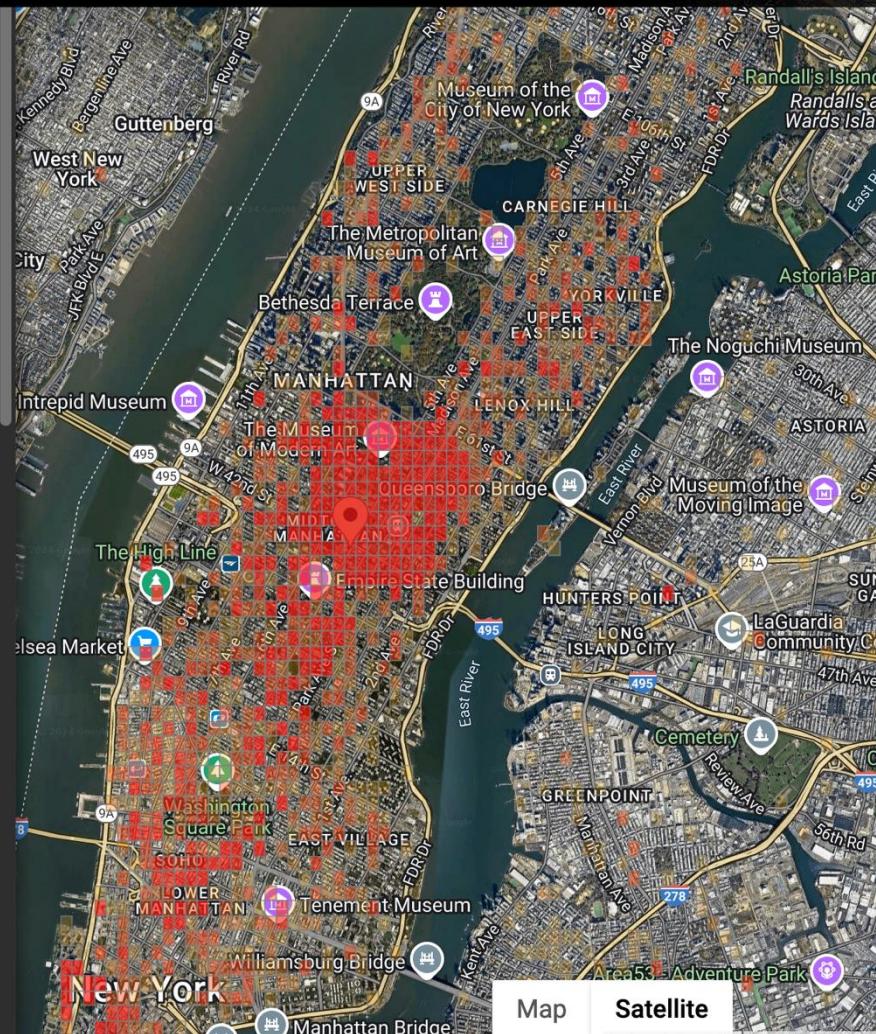
Most Likely Suspect



Top Risk Likelihoods

DEFAMATION
(13.23%)

THE NEW INQUIRY



<https://whitecollar.thenewinquiry.com/#dr5rudb>

Respect ethical principles (next week!): do no harm / do good.

Reduce tracking:

- Think if the app you are using contributes to your wellbeing, or if there are better ways for it.
- Use privacy-friendly software when possible (e.g., Signal, Proton, Cryptpad)
- Use plugins to reduce tracking

Don't take services for granted. Examine power structures:

- Who benefits?
- Who gets harm?
- Who gets to decide?

Require model transparency and accountability.

4. Dealing with bias in ML

What information do we need to understand bias?

Information on *protected attributes* (gender, sex, class, etc)

Information on *true labels* to create the confusion matrix (per attribute) → Often we need to manually label data

A definition of bias/fairness → There is no universally-accepted definition of what it means for a model to be fair. This is not an excuse for ignoring fairness!

Group A		Predicted criminal	Predicted not criminal	Group B	
Criminal	Predicted	10	10	Predicted criminal	Predicted not criminal
	criminal	100	1		
Not criminal	Predicted	1	10	Predicted not criminal	Predicted criminal
	not criminal	100	10		

Errors in algorithms

We need to consider the errors of every algorithm:

- How often they fail?
- For whom do they fail? (**bias**)

Remember there are people behind the data:

- **What are the costs of those failures?**

- *Assistive intervention:* “Individuals may be harmed by being incorrectly included in the “low need” population that does not receive an intervention” (**harm = false negatives**, e.g., not supporting somebody in need)

- *Punitive intervention:* “Individuals may be harmed by being incorrectly included in the “high risk” population that receives an intervention” (**harm = false positives**, e.g. jailing an innocent)

- What are the long-term effects? (feedback effects)

	Assistive	Predicted in need		Predicted not in need	
		In need	Not in need	In need	Not in need
In need	True positive	10		False negative	10
	False positive	1		True negative	100
Not in need	Punitive				
	True positive	10		False negative	10
Criminal	Predicted criminal				
	True positive	10		False negative	10
Not criminal	True negative				
	False positive	1		True negative	100

Assessing bias, punitive example

- Among the general population (T), the probability of *being wrongly jailed* is independent of race
- Among the jailed population (those predicted criminals), the probability of being wrongly jailed is independent of race: Parity in **False Discovery Rate**.
→ Focuses on those affected by the intervention
- Among innocents (the actual non-criminals), the probability of being wrongly jailed is independent of race: Parity in **False Positive Rate** (Predictive Equality)
→ Focuses on those who should not be affected by the intervention

False positive
All

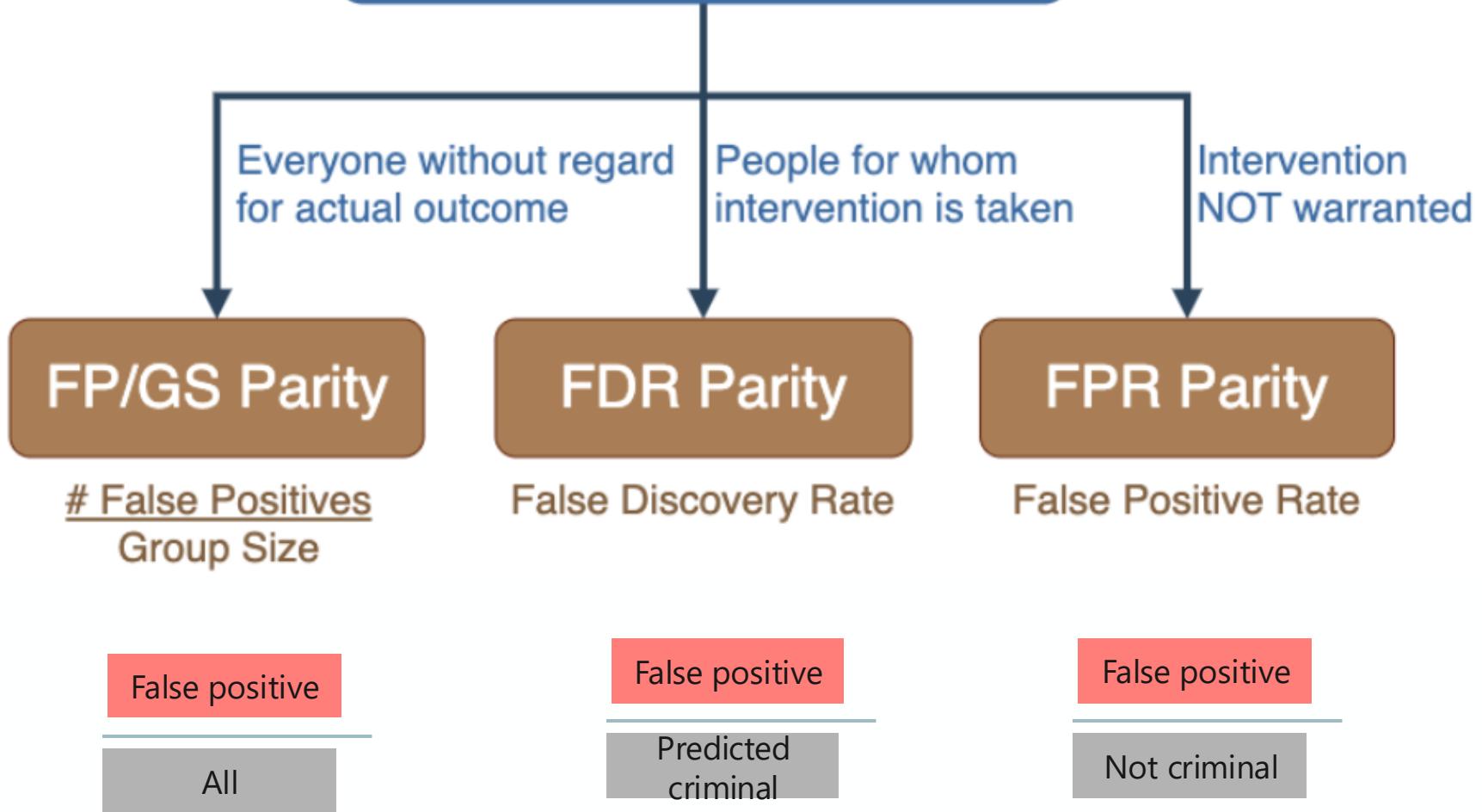
False positive
Predicted criminal

False positive
Not criminal

	Predicted criminal	Predicted not criminal
Criminal	True positive	False negative
Not criminal	False positive	True negative

No universal definition of fairness!

Among which group are you most concerned with ensuring predictive equity?



Models can be fair and unfair at the same time

	Criminal	Predicted criminal	Predicted not criminal
	Criminal	True positive	False negative
	Not criminal	False positive	True negative

False positive

Not criminal

False positive

Predicted
criminal

Table 11.1: COMPAS Fairness Metrics

Metric	Caucasian	African American
False Positive Rate (<i>FPR</i>)	23%	45%
False Discovery Rate (<i>FDR</i>)	41%	37%

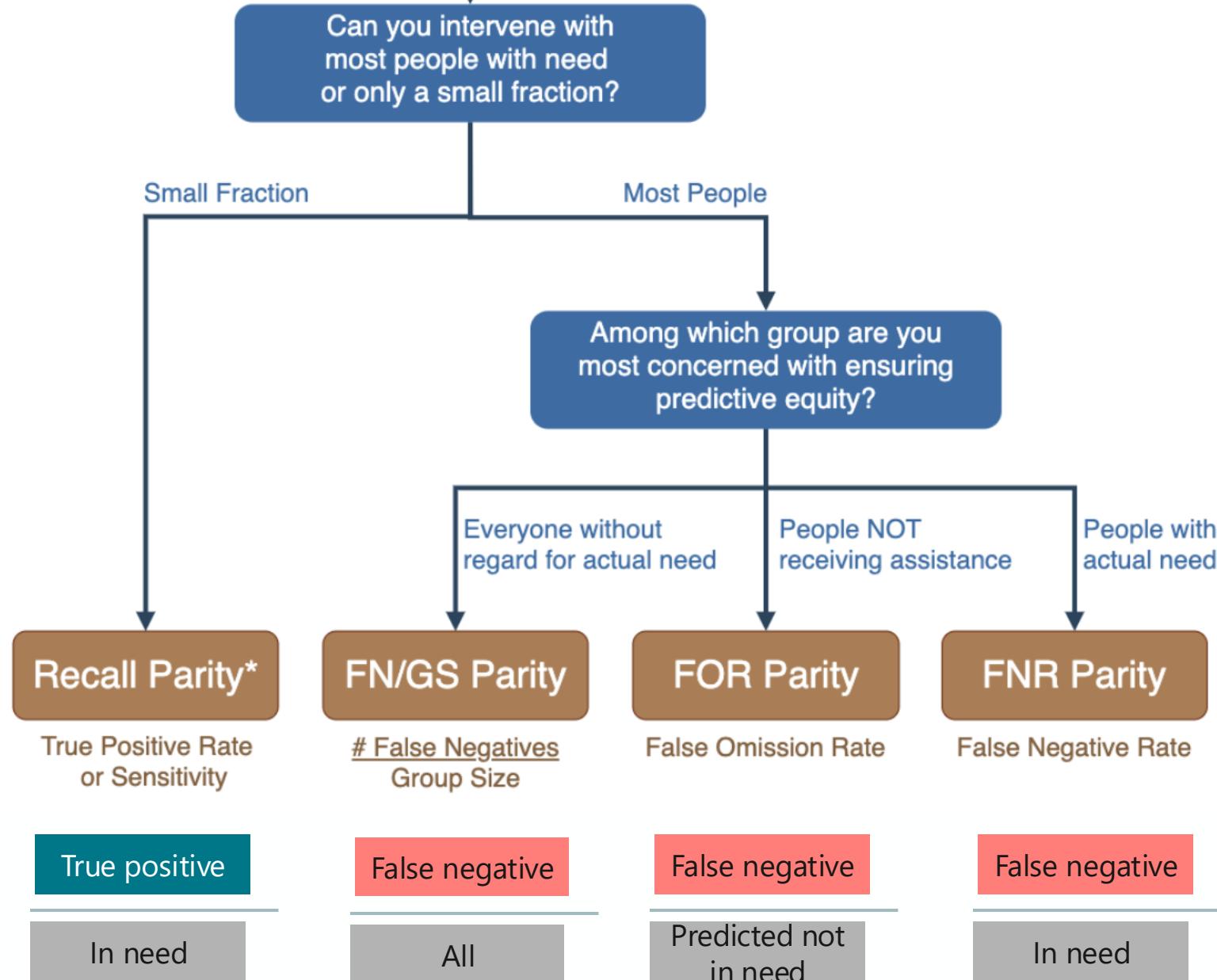
Correctional Offender Management Profiling for Alternative Sanctions (COMPAS): Evaluates the likelihood of an offender committing another crime in the future.

- *FPR*: Among black defendants who did not end up with another arrest, 45% were labeled by the system as high risk, almost twice the rate for whites (23%).
- *FDR*: Among individuals labelled as high risk, a similar fraction of black and white defendants were arrested again.

It's generally impossible for a model to maximize both fairness criteria at the same time

Assessing bias, assistive example

	Predicted in need	Predicted not in need
In need	True positive	False negative
Not in need	False positive	True negative



The application of the model may also introduce bias

“Perhaps a model developed to screen out *unqualified job candidates* is only “trusted” by a *hiring manager for female candidates* but often ignored or overridden for men. In a perverse way, applying an unbiased model in such a context might serve to increase inequities by giving bad actors more information with which to (wrongly) justify their discriminatory practices.” (Big Data and Social Science)

Perhaps a model developed to *identify criminal behavior* is only *deployed in low-income areas*. In a perverse way, applying an unbiased model in such a context might serve to increase inequities by giving bad actors more information with which to (wrongly) justify their discriminatory practices.

Dealing with bias in ML

Audit the model to understand bias

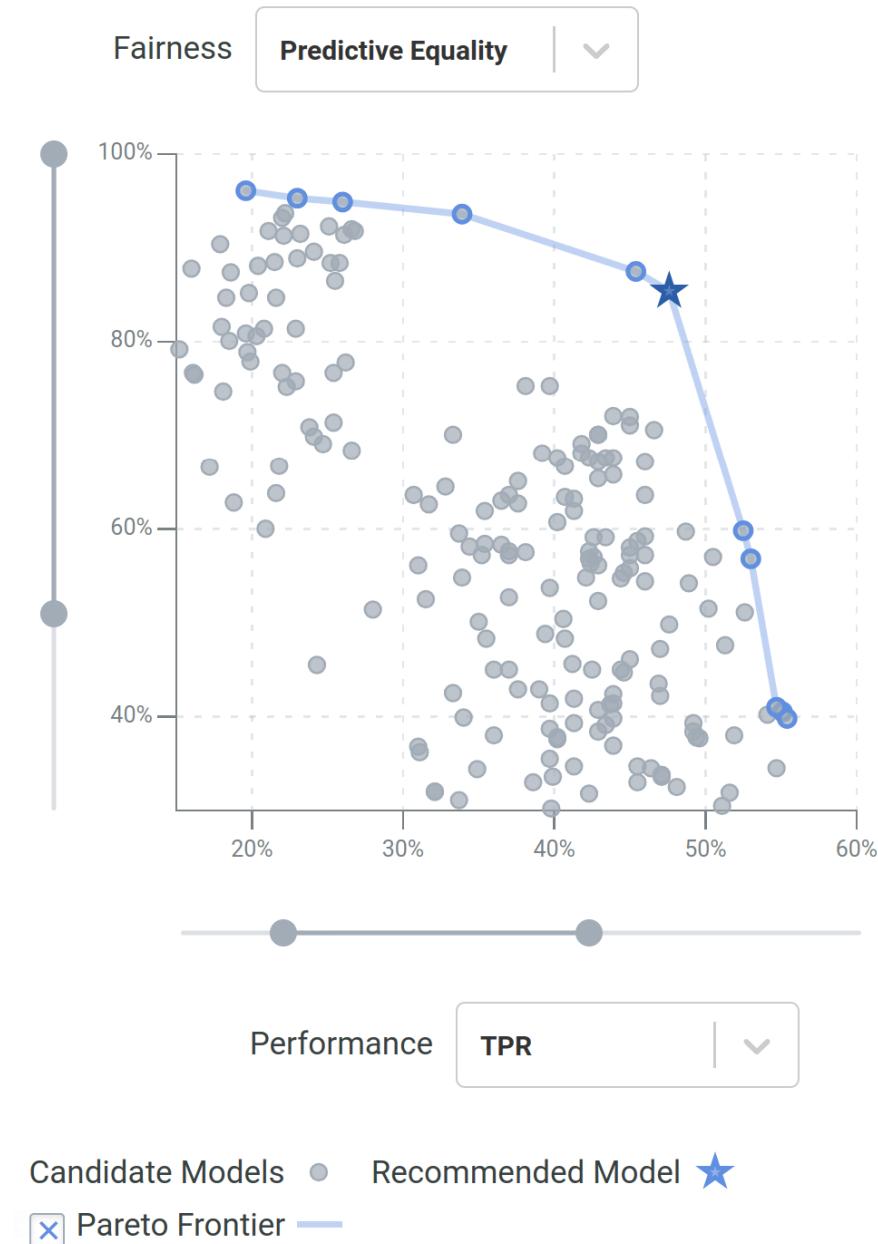
Mitigate bias:

- Removing "protected attributes" (gender/race) generally does *not* mitigate bias.
- Test different models and select one with strong performance across fairness and accuracy (Pareto optimality).
- Adjust thresholds to increase/decrease FP or FN. Example: Offer a subsidy to Group A if the model predicts a need with over 50% probability, and to Group B if the need is predicted with over 25% probability.

Consider Intersectionality: Optimizing for one factor (e.g., gender) may introduce bias in another (e.g., class).

Regularly Test for Bias: Monitor for concept drift and ensure ongoing fairness.

Consider if Bias May Be Acceptable: For example, if the intervention is most useful to a specific subpopulation.



TODAY

Lecture

1. Explain machine learning in your own words

2. Explain why machine learning models may be biased.

3. Understand the effects of ML on DTD and society.

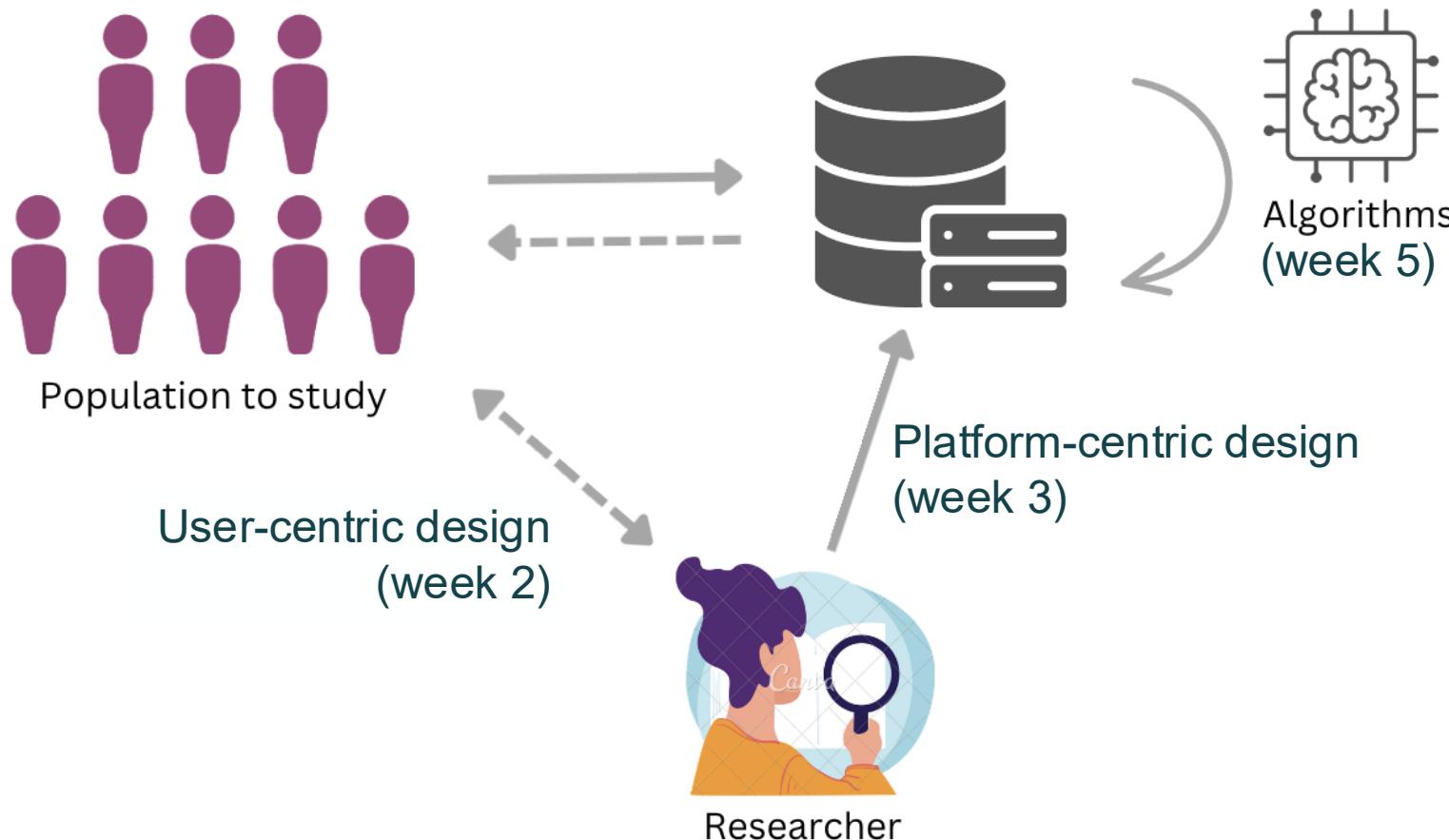
4. Assess bias in ML models

Lab

Apply a ML model to text data

Audit a ML model

Summary of the course



Week 4: Errors in DTD
Week 6: Ethics and Legislation
Week 9: Recap and Q&A