

# Design & Governance of Data-Driven Algorithms

Sarah H. Cen  
MIT EECS

Advised by Aleksander Mądry and Devavrat Shah

October 2, 2023  
Young Researchers Workshop, Cornell

# The Role of Data-Driven Algorithms



Should data-driven algorithms intervene on high-stakes decisions?



Should social media platforms choose how election-related content is disseminated?



When should the government regulate the AI industry, and when should it abstain?

Normative Factors

Practical Considerations

## Data-Driven Algorithms

### Design

What is feasible from a design perspective?  
(e.g., when do trade-offs exist)

[**CS'22**, **CIM'23**, **ZCS'23**,  
**ACSY'23**, **CS'21**, **CMS'23**]

### Governance

What is feasible from a governance perspective?  
(e.g., what can be regulated)

[**CS'21**, **CIM'22**, **CR'23**,  
**CSFMM'23**, **CHIMSV'23**, **CMS'23**]

# Today: Auditing

**Goal:** Verify that a data-driven algorithm satisfies some criteria.

**Input:** Criteria (e.g., given by law)

**Constraint:** Black-box access

## **Questions:**

1. To what extent does the audit test for the given criteria?
2. How much data does the audit need?
3. Are there side effects?

# Case Study: Auditing Social Media

Cen and Shah, 2021



NEWS CULTURE MUSIC PODCASTS & SHOWS SEARCH

POLITICS

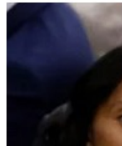
Did Fake News On Facebook  
Trump? Here's What We Know

April 11, 2018 ·



DANIEL

ACLU



NEWS & COMMENTARY

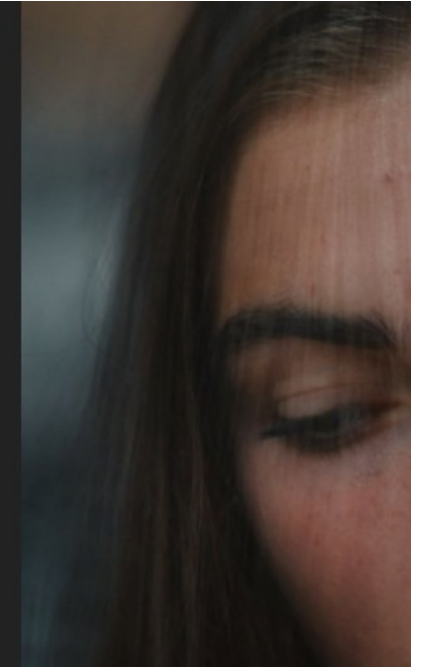
## Holding Facebook Accountable for Digital Redlining

Online ad-targeting practices often reflect and replicate existing disparities, effectively locking out marginalized groups from housing, job, and credit opportunities.

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



There are rising calls for social media regulations.

# Calls to Regulate

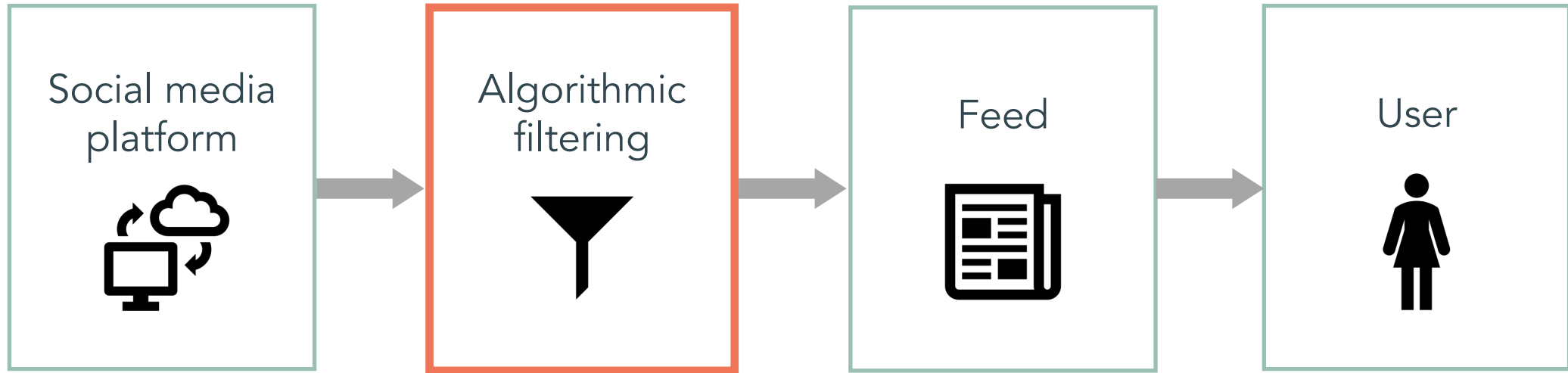
Ex 1: Ads not be based on user's sexual orientation.

Ex 2: Info on public health (e.g., COVID-19) not reflect political affiliation.

Ex 3: Not sway voting preferences beyond serving as a social network.

Translating desiderata → audit is difficult

- Performance cost
- Censorship
- User privacy
- Trade secrets



## **Main contribution: data-driven auditing procedure**

Strong statistical guarantees

Not necessarily a performance-audit trade-off

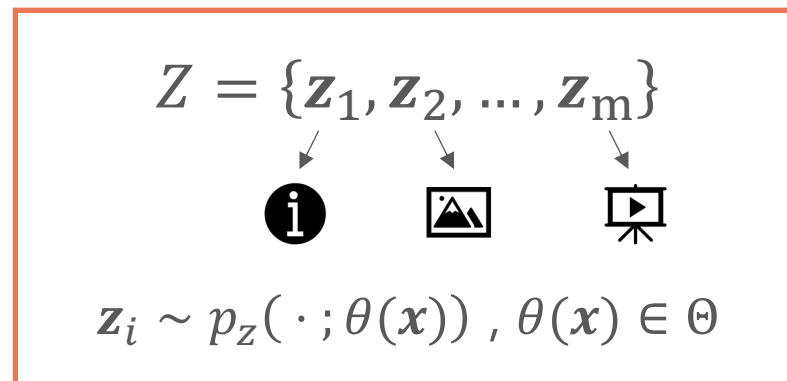
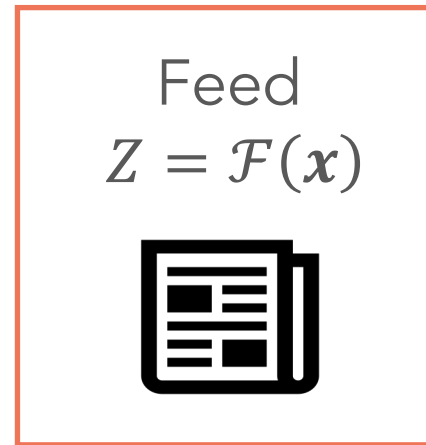
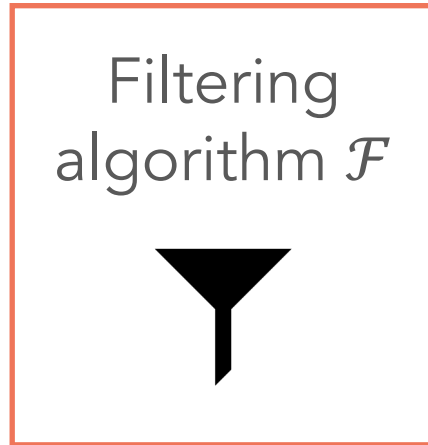
Incentivizes platform to inject content diversity

Requires only black-box access

Does not remove content or require user data



# Problem setup



Black-box access:  
Run  $\mathcal{F}$  on  $\{x_j\}$   
and observe  $\{Z_j\}$ .

## Auditor's task

Given a CF desideratum  
& black-box access to  $\mathcal{F}$ ,  
check if platform complies.

# Counterfactual desideratum

hypothetical!

"Algorithm  $\mathcal{F}$  must behave similarly under  $x$  and  $x'$  for all  $(x, x') \in S$ ."

Articles with  
medical advice on  
COVID-19 must be  
robust to user's  
political affiliation.



Articles shown by  $\mathcal{F}$  that  
have medical advice on  
COVID-19 should be  
**similar** whether a user is  
left- or right-leaning.

# Counterfactual desideratum

hypothetical!

“Algorithm  $\mathcal{F}$  must behave similarly under  $x$  and  $x'$  for all  $(x, x') \in S$ .”

The platform  
should not sway  
voting beyond  
serving as a  
social network.



Election-related posts  
that  $\mathcal{F}$  injects should be  
**similar** to those a user  
would see from its social  
network (without filtering).

What is an appropriate notion of “similarity” ?

$$\|Z_i - Z'_i\|_q < \delta ?$$

**Observation:** Algorithmic filtering is powerful (sometimes harmful) because information influences decisions.

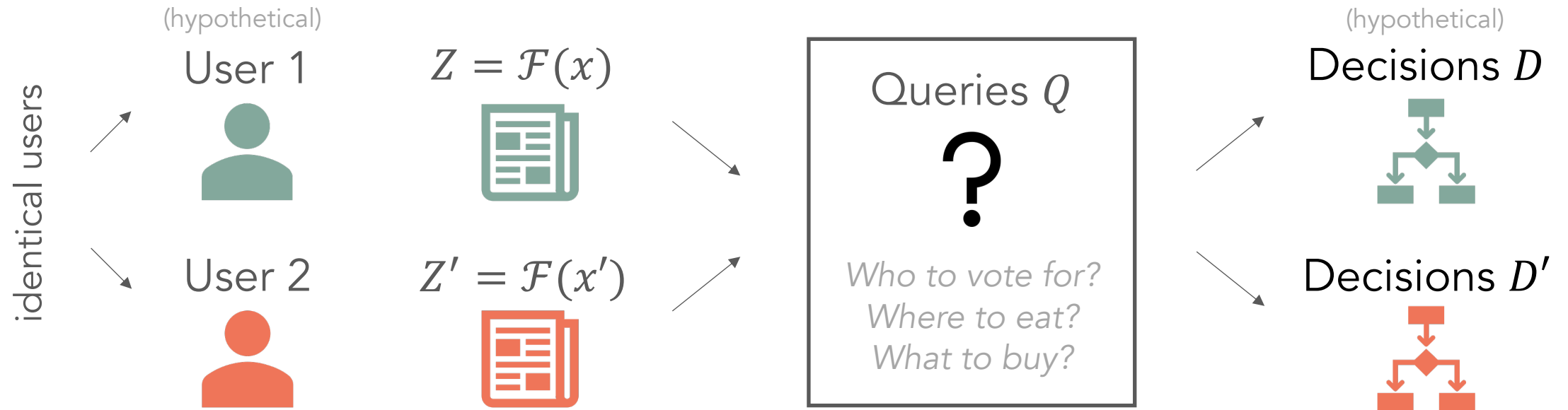
**Examples.** The content that  $\mathcal{F}$  filters affects ...

- What a user eats
- What they buy
- How the user votes

**Implication.** Should enforce **similarity** between  $\mathcal{F}(x)$  and  $\mathcal{F}(x')$  w.r.t. the outcome of interest: **the users' decisions**.

# Decision robustness

CF Reg: " $\mathcal{F}$  must behave similarly under  $x$  and  $x'$  for all  $(x, x') \in S$ ."



$\mathcal{F}$  is decision-robust to  $(x, x')$  if and only if, for any  $Q$ , one cannot confidently determine that  $x \neq x'$  from  $D$  and  $D'$ .

↳ can formalize as hypothesis test

# Auditing procedure

"Algorithm  $\mathcal{F}$  must behave similarly under  $\mathbf{x}$  and  $\mathbf{x}'$  for all  $(\mathbf{x}, \mathbf{x}') \in S$ ."

Inputs:

$\mathcal{F}$

$\mathbf{x}$

$\mathbf{x}'$

$\Theta$

$\epsilon$

- 1  $\tilde{\theta} \leftarrow \mathcal{L}^+(\mathcal{F}(\mathbf{x}));$
  - 2  $\tilde{\theta}' \leftarrow \mathcal{L}^+(\mathcal{F}(\mathbf{x}'));$
  - 3 **if**  $(\tilde{\theta} - \tilde{\theta}')^\top I(\tilde{\theta})(\tilde{\theta} - \tilde{\theta}') \geq \frac{2}{m} \chi_r^2(1 - \epsilon)$  **then**
  - 4 | Does not pass the test for  $(\mathbf{x}, \mathbf{x}')$ ;
  - 5 **end**
  - 6 Passes the test for  $(\mathbf{x}, \mathbf{x}')$ ;
- Minimum-variance unbiased estimator (MVUE)

# Advantages

- Only needs black-box access to  $\mathcal{F}$ .
- Does not require access to users or their personal data.
- Modular. Can scale up for any  $(\mathbf{x}, \mathbf{x}')$  pairs.
- Intuitive tunable parameter.  $\epsilon$  is false positive rate.
- No content removal.



## Guarantee on how well the audit enforces the regulation.

**Theorem (informal).** If the filtering algorithm  $\mathcal{F}$  passes the audit, then  $\mathcal{F}$  is guaranteed to be approximately asymptotically decision-robust.

### Alternative statement

If  $\mathcal{F}$  does not pass the audit, then the auditor is  $(1 - \epsilon)$ -confident that  $\mathcal{F}$  is not decision-robust as  $m \rightarrow \infty$ .

### Takeaways

- The audit enforces strong similarity between  $\mathcal{F}(\mathbf{x})$  and  $\mathcal{F}(\mathbf{x}')$ .
- $\epsilon$  is the allowable false positive rate: increasing  $\epsilon$  increases strictness.

## Why the MVUE?

**Proposition (informal).** Faced with a decision between a finite number of options, the decision of the hypothetical user whose belief after viewing content  $Z$  is given by the MVUE is more sensitive to  $Z$  than any other user.

### Takeaway

MVUE = user whose decisions are **most sensitive** to the content they see.

The MVUE allows us to reason about how content affects users *without access to users' decisions* → expensive or unethical to obtain.

There isn't a regulation-performance trade-off.

**Theorem (informal).** Consider a finite feed. If performance is independent of elements in  $\theta$  that can increase the Fisher information and the available content is diverse, then there is no regulation-performance trade-off.

Takeaway

Platform can pass audit without sacrificing performance.

Content diversity can reduce the cost of regulation

The lower the content diversity of  $Z$  and  $Z'$ , the more easily an auditor can distinguish between how  $\mathcal{F}$  behaves under  $x$  and  $x'$ .

# Design & Governance of Human-Facing Algorithms

## Case Study: Auditing Social Media

Black-box auditing procedure

Audit is consistent with existing laws

## Extensions

Instantiated a viable regulation [CMS '23]

Running a live audit this month

Auditing from dataset [ACSY '23]

Here, we asked: What is feasible from an auditing perspective?

Interplay between design & governance is going to be important

Thank you!

[shcen@mit.edu](mailto:shcen@mit.edu)