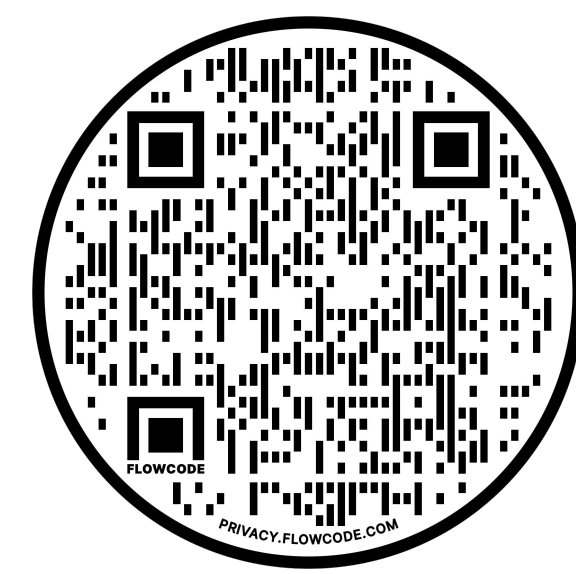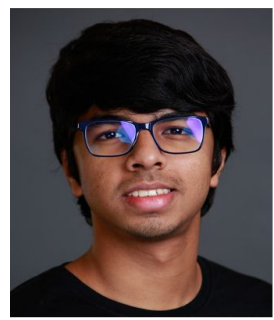# Finding Valid Adjustments under Non-ignorability with Minimal DAG Knowledge

Abhin Shah
MIT
abhin@mit.edu

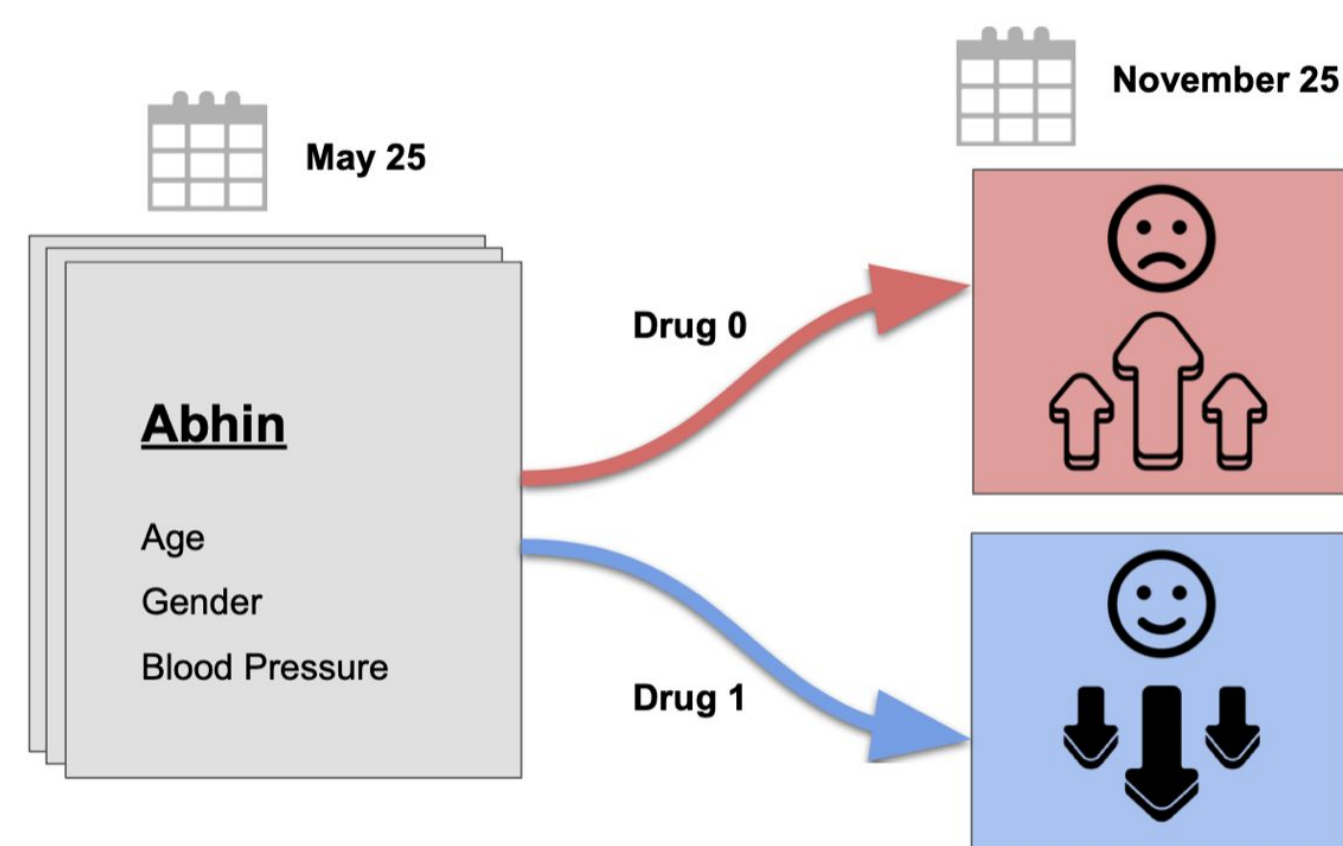Karthikeyan Shanmugam
Google Research

Kartik Ahuja
Meta AI

arXiv Link

## Causal Effect Estimation

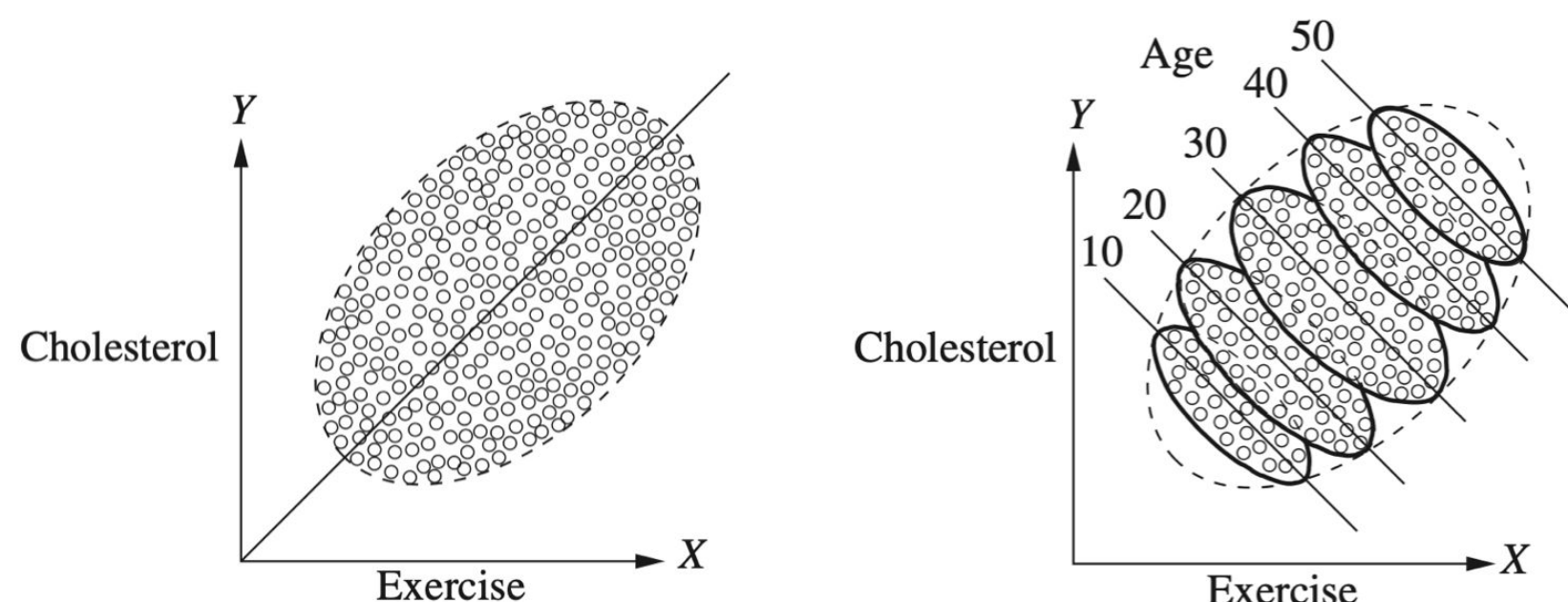Causal effect of a drug on cholesterol level from observational data

May 25

November 25

**Abhin**

Age
Gender
Blood Pressure

Drug 0

Drug 1

$\mathbb{P}(\text{cholesterol} \mid do(\text{drug}))$ ?

### Observational Data

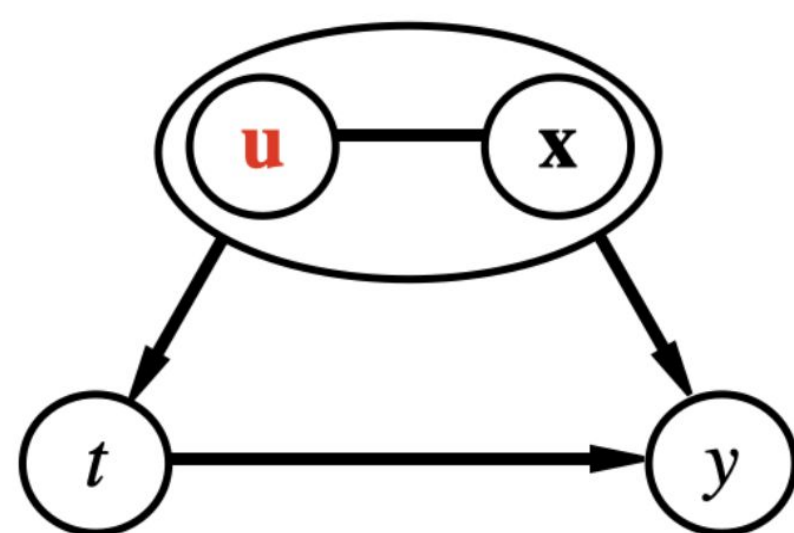| Age | Gender | Blood Pressure | Drug | Cholesterol (0) | Cholesterol (1) |
|-----|--------|----------------|------|-----------------|-----------------|
| 22 | Male | 145/95 | 0 | ⇑⇑ | ? |
| 26 | Female | 135/80 | 0 | ⬇⬇⬇ | ? |
| 58 | Female | 130/70 | 1 | ? | ⇑⇑ |
| 50 | Male | 145/80 | 1 | ? | ⬇⬇ |
| 24 | Female | 150/85 | 1 | ? | ⇑⇑ |

## Challenge — Unobserved Confounding

**Simpson's paradox:** Which subsets of the observed features should be used?

## Problem Formulation

- **u** : unobserved exogenous variables
- **x** : observed features
- $t$ : observed binary treatment variable
- $y$ : observed outcome
- $\mathscr{G}$ : DAG over the set of vertices $\{\mathbf{u}, \mathbf{x}, t, y\}$

## Valid adjustments

$\mathbf{z}$ is a valid adjustment set if $\mathbb{P}(y \mid do(t = t)) = \mathbb{E}_{\mathbf{z}}[\mathbb{P}(y \mid \mathbf{z} = z, t = t))]$

### Pearlian Framework

DAG knowledge

Given the complete knowledge of the DAG, graphical criteria could be used to check whether $\mathbf{z}$ is valid for adjustment
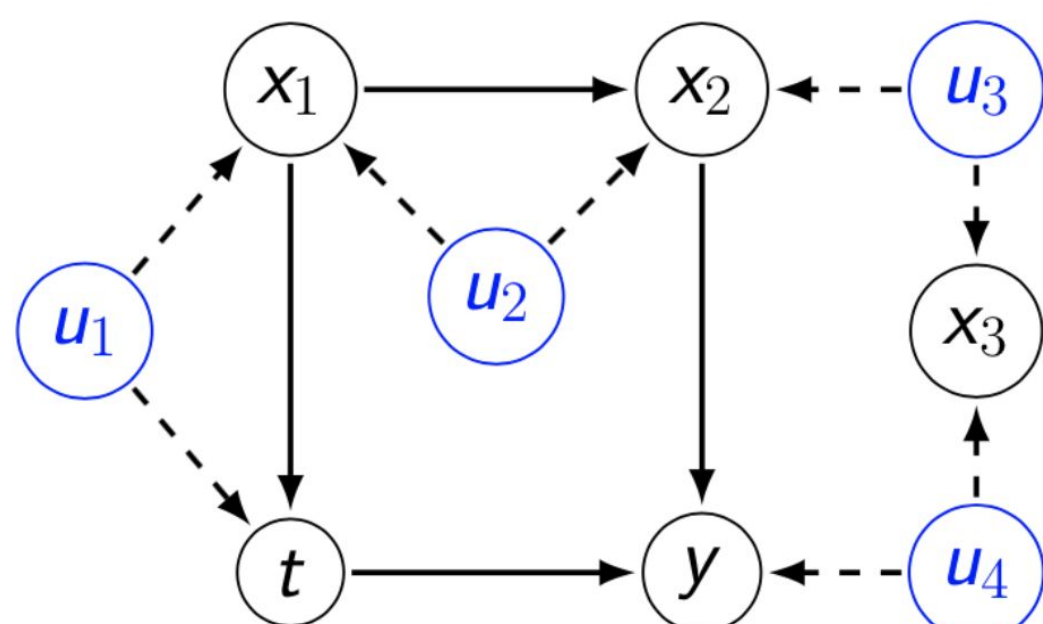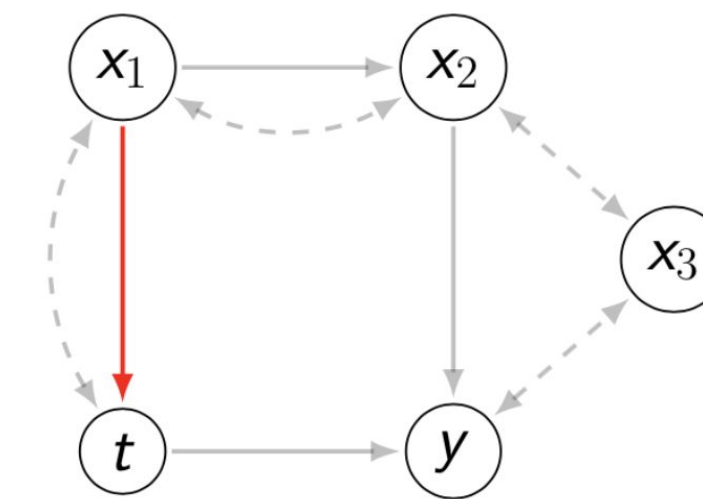
### Potential Outcomes

Ignorability

**x** satisfies ignorability

⇓

**x** is a valid adjustment.

## How much of the DAG do we need to know?

To find the causal effect of $t$ on $y$, i.e., $\mathbb{P}(y \mid do(t = t))$

---

Can we significantly reduce the structural knowledge required about the DAG and yet find valid adjustment sets?
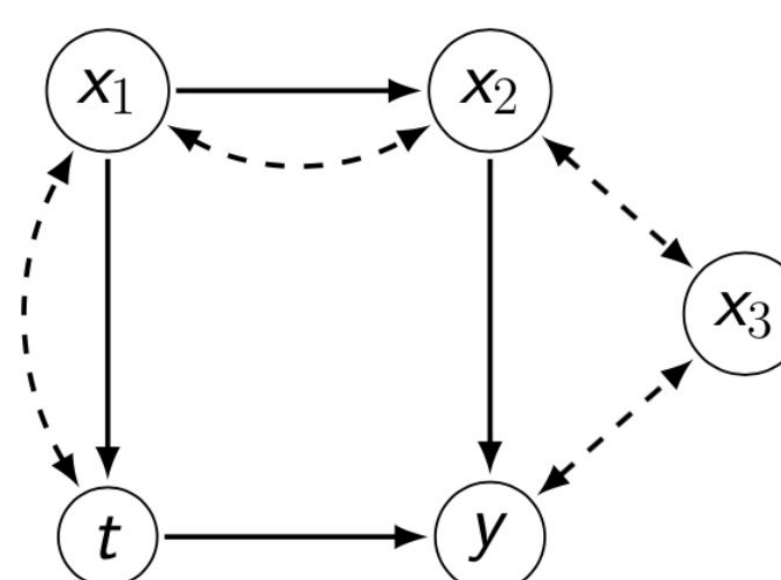
The knowledge of one causal parent of the treatment is sufficient to find a class of valid adjustment sets!

## Assumptions

**Semi-Markovian model**

1. The treatment $t$ has the outcome $y$ as its only child.

2. The outcome $y$ has no child.

## Back-door Criterion

**A popular sufficient graphical criterion for finding valid adjustments**

Under our assumptions, a set $\mathbf{z}$ satisfies the back-door criterion in $\mathscr{G}$ if

1. $\mathbf{z}$ blocks every path between $t$ and $y$ in $\mathscr{G}$ that contains an arrow into $t$.

Sets satisfying back-door: $\{x_1, x_2\}$ and $\{x_2\}$

## Conditional Independence ⟺ Back-door

- $x_t$ : an observed feature that is a direct causal parent of $t$.

- Consider any subset of the remaining observed features i.e., $\mathbf{z} \subseteq \mathbf{x} \setminus \{x_t\}$.

- $\mathbf{z}$ satisfies the back-door criterion if and only if $x_t \perp y \mid \mathbf{z}, t$.

## Algorithms

- Subset Search:
  ➡ Use a subset based search procedure that exploits conditional independence (CI) testing to check our invariance criterion.

- IRM-based:
  ➡ Use a sub-sampling trick to leverage Invariant Risk Minimization (IRM) as a scalable approximation for CI testing.

## IHDP

**A RCT studying cognitive test score of low-birth-weight, premature infants.**



## Cattaneo

**Studies the effect of maternal smoking on babies' birth weight.**

# On Counterfactual Inference with Unobserved Confounding

Abhin Shah — abhin@mit.edu — MIT

Raaz Dwivedi — raaz@mit.edu — MIT, HARVARD UNIVERSITY

Devavrat Shah — devavrat@mit.edu — MIT

Greg Wornell — gww@mit.edu — MIT

arXiv Link

## Observational Setting

**a** — treatment/intervention
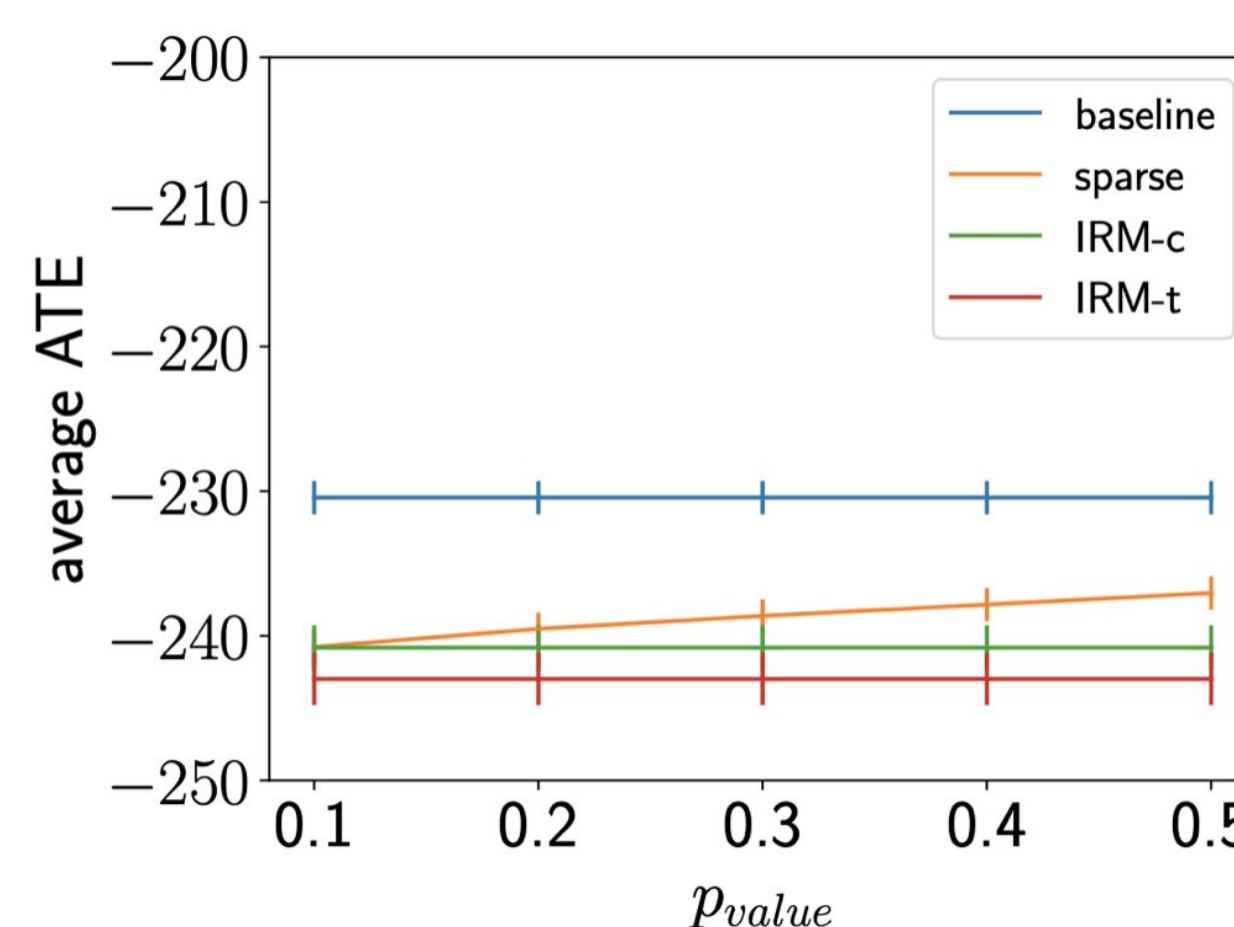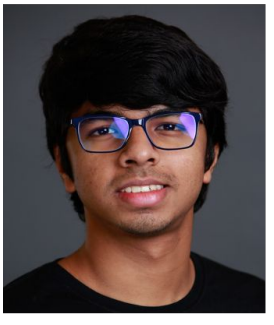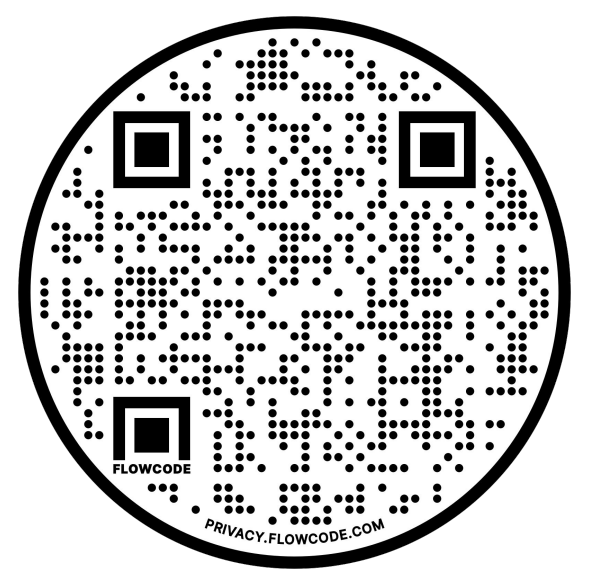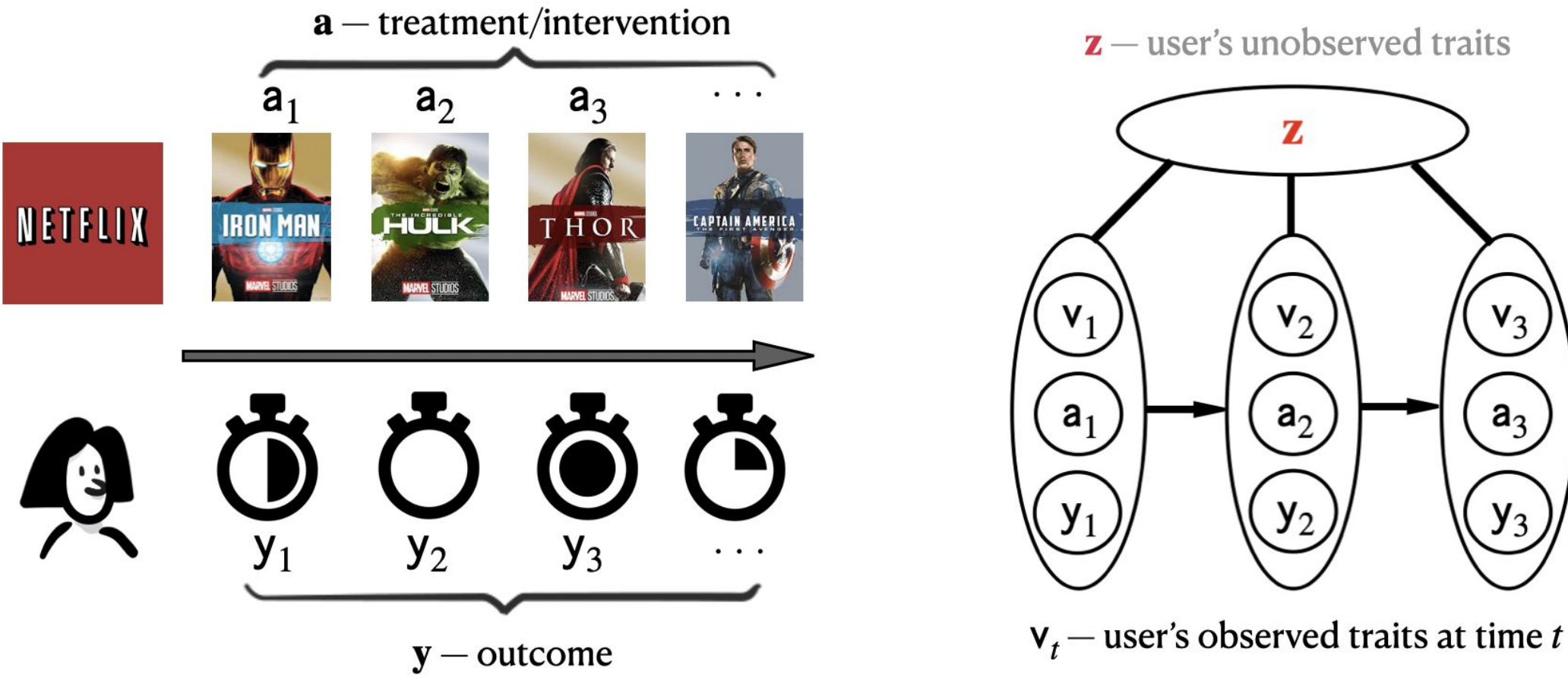
$a_1$ $a_2$ $a_3$ $\cdots$

NETFLIX

$y_1$ $y_2$ $y_3$ $\cdots$

**y** — outcome

**z** — user's unobserved traits

$\mathbf{z}$

$v_1$ $v_2$ $v_3$
$a_1$ $a_2$ $a_3$
$y_1$ $y_2$ $y_3$

$v_t$ — user's observed traits at time $t$

## Panel Data

**Potential Outcomes**

$\{\mathbf{y}^{(1)}(\mathbf{a})\}_{\mathbf{a}\in\mathscr{A}}$

$\mathbf{a}^{(1)}$    $\mathbf{y}^{(1)}$    $\mathbf{y}^{(1)} = \mathbf{y}^{(1)}(\mathbf{a}^{(1)})$

$\vdots$

$\{\mathbf{y}^{(n)}(\mathbf{a})\}_{\mathbf{a}\in\mathscr{A}}$

$\mathbf{a}^{(n)}$    $\mathbf{y}^{(n)}$    $\mathbf{y}^{(n)} = \mathbf{y}^{(n)}(\mathbf{a}^{(n)})$

## Goal

$\widetilde{\mathbf{a}}^{(1)}$    $\mathbf{y}^{(1)}(\widetilde{\mathbf{a}}^{(1)})$ ?

$\vdots$

$\widetilde{\mathbf{a}}^{(n)}$    $\mathbf{y}^{(n)}(\widetilde{\mathbf{a}}^{(n)})$ ?

## Challenges

1. unobserved factors $\rightarrow$ spurious associations
2. users $\rightarrow$ heterogeneous
3. each user $\rightarrow$ a single interaction trajectory

engagement level — comedy — Simpson's paradox

engagement level — comedy — 10 20 30 40 50 age

## Problem Setup

unobserved covariates $\mathbf{z}$ — $\mathbf{v}$ observed covariates

observed interventions $\mathbf{a}$ $\rightarrow$ $\mathbf{y}$ observed outcomes

$n$ heterogenous and independent users with one observation each - $\underbrace{\{\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^n}_{p-\text{dimensional}}$

## Goal : Counterfactual Questions

For user $i \in [n]$, what would have happened if alternative treatments were assigned?

$\equiv$

Estimate $\mathbf{y}^{(i)}(\widetilde{\mathbf{a}}^{(i)})$ for $\widetilde{\mathbf{a}}^{(i)} \in \mathscr{A}$?

Suffices to learn $f(\mathbf{y} = \cdot \mid \mathbf{a} = \cdot, \mathbf{z}^{(i)}, \mathbf{v}^{(i)})$ for all $i \in [n]$, but each user may have *different* $\mathbf{z}$

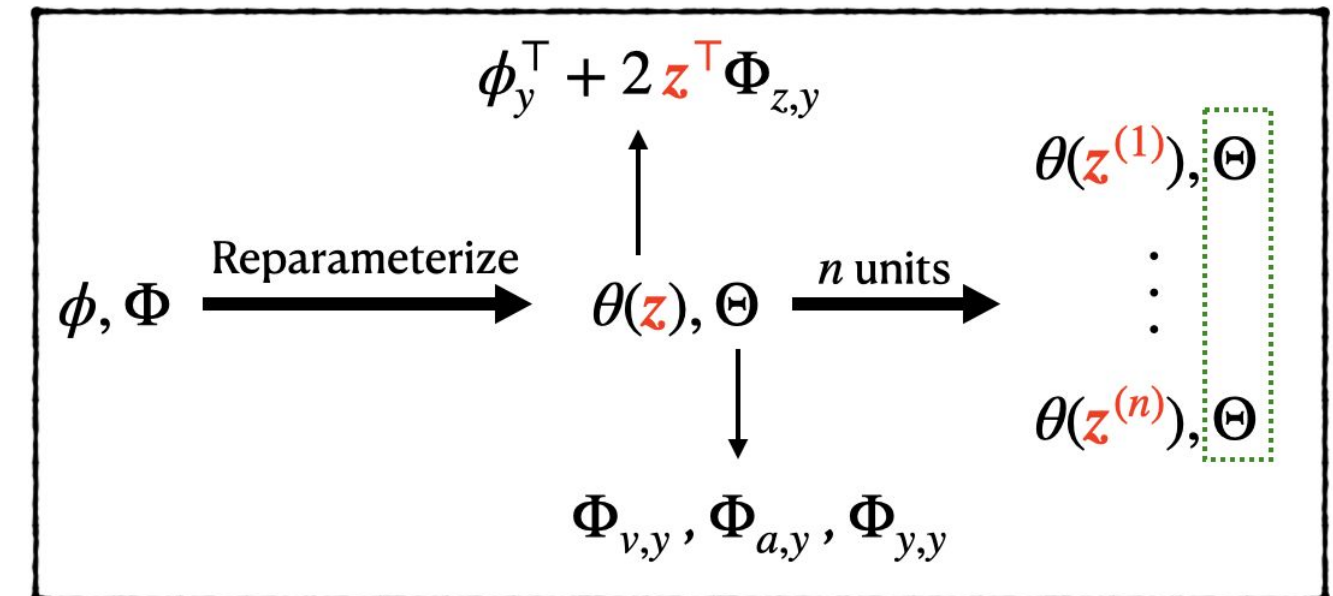Can we learn $n$ different distributions with *one* sample per distribution ?

## Our Approach

We posit a joint exponential family distribution for $\mathbf{w} \triangleq (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y})$

$$f(\mathbf{w}) \propto \exp\Big(\boldsymbol{\phi}^\top \mathbf{w} + \mathbf{w}^\top \boldsymbol{\Phi} \mathbf{w}\Big)$$

$$f(\mathbf{y} \mid \mathbf{a}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}) \propto \exp\Big(\Big[\ \boldsymbol{\phi}_y^\top + 2\mathbf{z}^{(i)\top}\boldsymbol{\Phi}_{z,y} + 2\mathbf{v}^{(i)\top}\boldsymbol{\Phi}_{v,y} + 2\mathbf{a}^\top\boldsymbol{\Phi}_{a,y}\ \Big]\mathbf{y} + \mathbf{y}^\top\boldsymbol{\Phi}_{y,y}\mathbf{y}\Big)$$

different for different users

$n$ heterogeneous conditional distributions $\Rightarrow$ same exp. family but with diff. parameters

$\boldsymbol{\phi}_y^\top + 2\mathbf{z}^\top\boldsymbol{\Phi}_{z,y}$

$\boldsymbol{\phi}, \boldsymbol{\Phi}$ $\xrightarrow{\text{Reparameterize}}$ $\theta(\mathbf{z}), \Theta$ $\xrightarrow{n \text{ units}}$ $\theta(\mathbf{z}^{(1)}), \widehat{\Theta}$ $\vdots$ $\theta(\mathbf{z}^{(n)}), \widehat{\Theta}$

$\boldsymbol{\Phi}_{v,y}, \boldsymbol{\Phi}_{a,y}, \boldsymbol{\Phi}_{y,y}$

## Inference Tasks

**1. Parameters:**   User-level — $\theta^\star(\mathbf{z}^{(i)})$ for all $i \in [n]$ $\rightarrow$ counterfactual distribution

Population-level — $\Theta^\star$

**2. Potential Outcomes:**   $\mu^{(i)} \triangleq \mathbb{E}\Big[\mathbf{y}^{(i)}(\widetilde{\mathbf{a}}^{(i)}) \mid \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}\Big]$ $\rightarrow$ counterfactual mean

## Parameter Estimation

$\{\mathbf{v}^{(i)}, \mathbf{a}^{(i)}, \mathbf{y}^{(i)}\}_{i=1}^n$ pool all $n$ samples $\longrightarrow$ $\widehat{\theta}^{(1)}, \cdots, \widehat{\theta}^{(n)}, \widehat{\Theta}$ estimates

$$\min_{\theta^{(1)}, \cdots, \theta^{(n)}, \Theta} \sum_{t\in[p]} \frac{1}{n} \sum_{i\in[n]} \exp\Big(-\big[\theta_t^{(i)} + 2\Theta_t^\top \mathbf{x}^{(i)}\big]x_t^{(i)}\Big)$$

Assum 1:   $\Theta^\star$ has sparse rows

Assum 2:   $\theta^\star(\mathbf{z}^{(i)}) \in$ set $\mathscr{B}$

$$\|\Theta^\star - \widehat{\Theta}\|_{2,\infty} \leq \epsilon \qquad \text{when } n \geq O\Big(\frac{p^2(p + M_n(\epsilon^2))}{\epsilon^4}\Big)$$
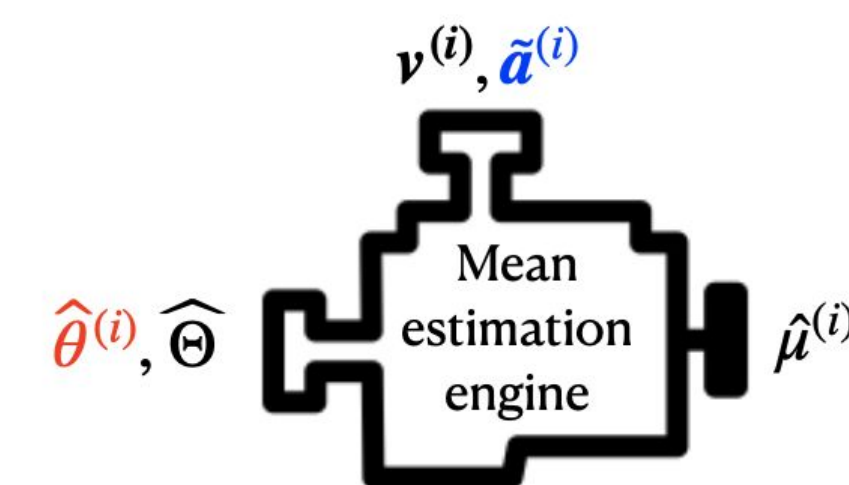
$$\text{For all } i, \text{MSE}\big(\theta^\star(\mathbf{z}^{(i)}), \hat{\theta}^{(i)}\big) \leq \max\Big\{\epsilon^2, \frac{M(c)}{p}\Big\} \quad \text{when } n \geq O\Big(\frac{p^2(pM(c) + M_n(\epsilon^2))}{\epsilon^4}\Big)$$

metric entropy of $\mathscr{B}$     $M_n(\epsilon) = nM(n\epsilon)$

$\star$ When $\mathscr{B} = s-$sparse linear combinations of $k$ known vectors,

$M(c) = O\big(s\log(k)\big)$ and $M_n(\epsilon) = O\Big(\frac{s\log k}{\epsilon}\Big)$

## Outcome Estimation

$\mathbf{v}^{(i)}, \widetilde{\mathbf{a}}^{(i)}$

$\widehat{\theta}^{(i)}, \widehat{\Theta}$ $\rightarrow$ Mean estimation engine $\rightarrow$ $\hat{\mu}^{(i)}$

$$\hat{f}(\mathbf{y} \mid \mathbf{a} = \widetilde{\mathbf{a}}^{(i)}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}) \propto \exp\Big(\Big[\ \widehat{\theta}(\mathbf{z}^{(i)}) + 2\mathbf{v}^{(i)\top}\widehat{\boldsymbol{\Phi}}_{v,y} + 2\widetilde{\mathbf{a}}^{(i)\top}\widehat{\boldsymbol{\Phi}}_{a,y}\ \Big]\mathbf{y} + \mathbf{y}^\top\widehat{\boldsymbol{\Phi}}_{y,y}\mathbf{y}\Big)$$

For all $i$ and any $\widetilde{\mathbf{a}}^{(i)} \in \mathscr{A}$,

$$\text{MSE}\big(\mu^{(i)}, \hat{\mu}^{(i)}\big) \leq \epsilon^2 + \frac{M(c)}{p} \quad \text{when } n \geq O\Big(\frac{p^2(pM(c) + M_n(\epsilon^2))}{\epsilon^4}\Big)$$

## Application: Denoise User-wise Data

No systematically unobserved covariates

Noisy observed data = true data + measurement error

$\overline{\mathbf{X}}$     $\mathbf{X}$     $\Delta\mathbf{x}$

Assum 1: Only half users have error: $\Delta\mathbf{x}^{(i)} = \mathbf{0}$ for $i \in \{n/2, \cdots, n\}$

Assum 2: Data has sparse error: $\|\Delta\mathbf{x}^{(i)}\|_0 \leq s$ for $i \in \{1, \cdots, n/2\}$

Goal: Estimate the true data

$$\text{For all } i, \|\Delta\mathbf{x}^{(i)}, \widehat{\Delta\mathbf{x}^{(i)}}\|^2 \leq \max\Big\{\frac{\epsilon^2}{s}, \frac{s}{p}\Big\} + \epsilon^2 \quad \text{when } n \geq O\Big(\frac{s^2 p}{\epsilon^4}\Big)$$