

On Counterfactual Inference with Unobserved Confounding

Abhin Shah

MIT



Raaz Dwivedi
Harvard & MIT



Devavrat Shah
MIT

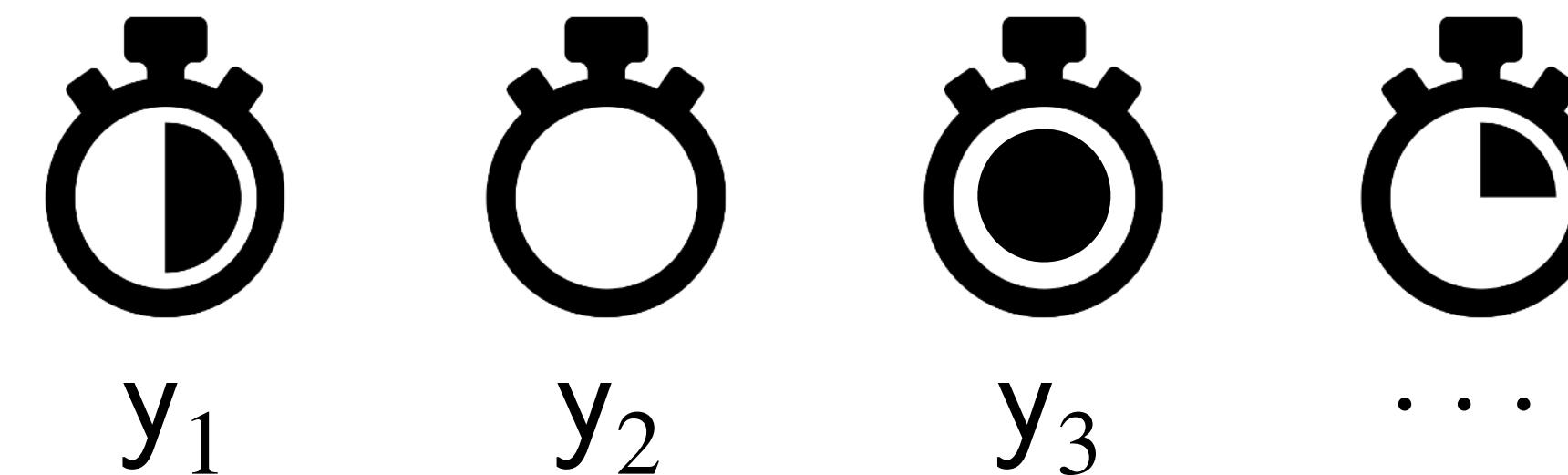


Greg Wornell
MIT

Observational Setting



a — treatment/intervention

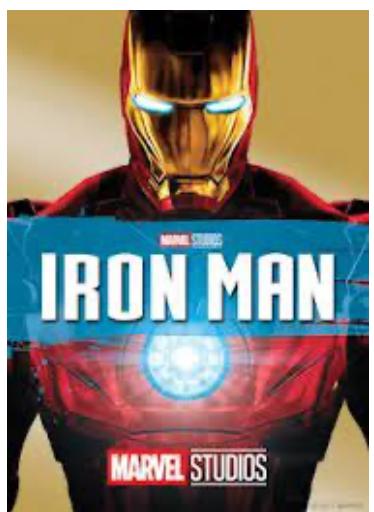


y — outcome

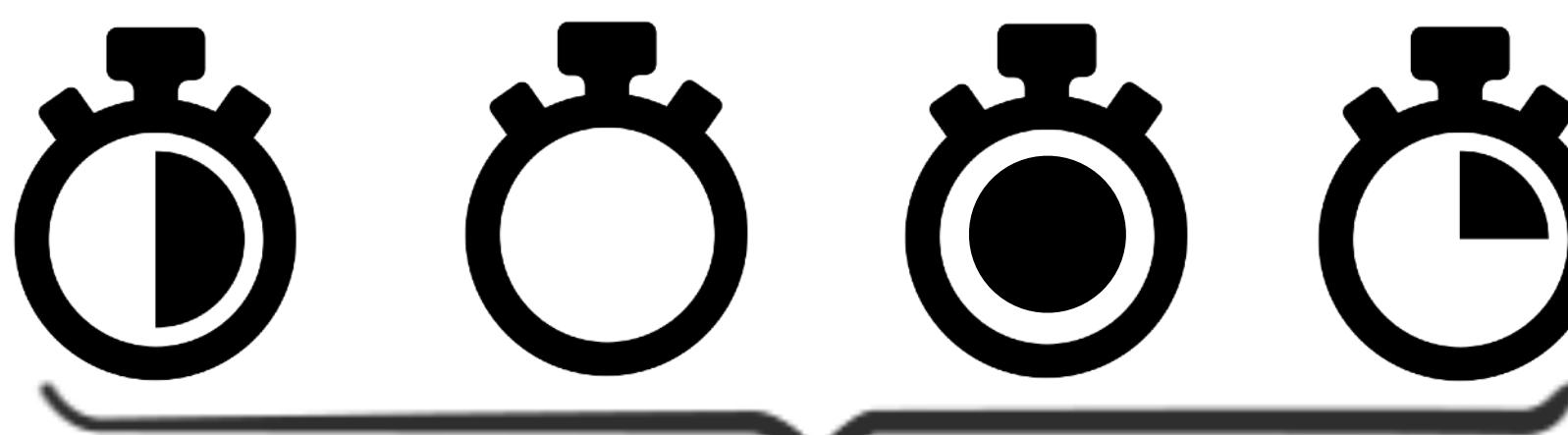


Panel Data

Potential
Outcomes



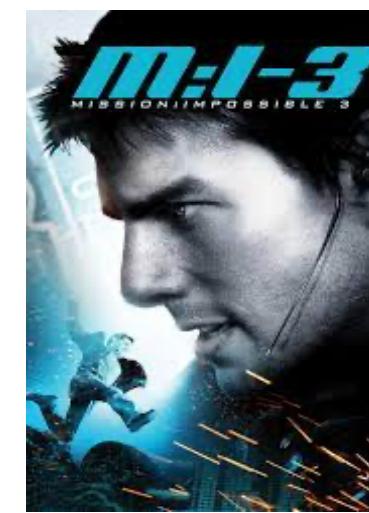
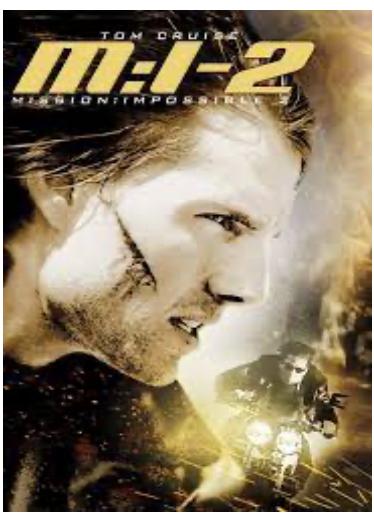
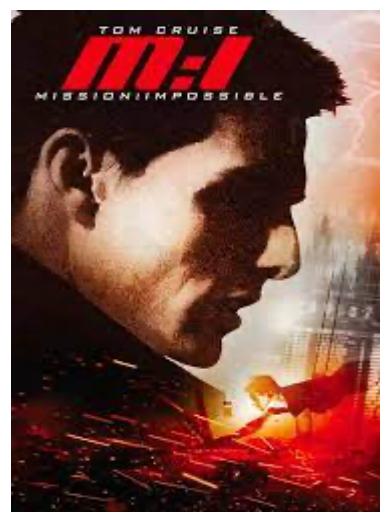
$a^{(1)}$



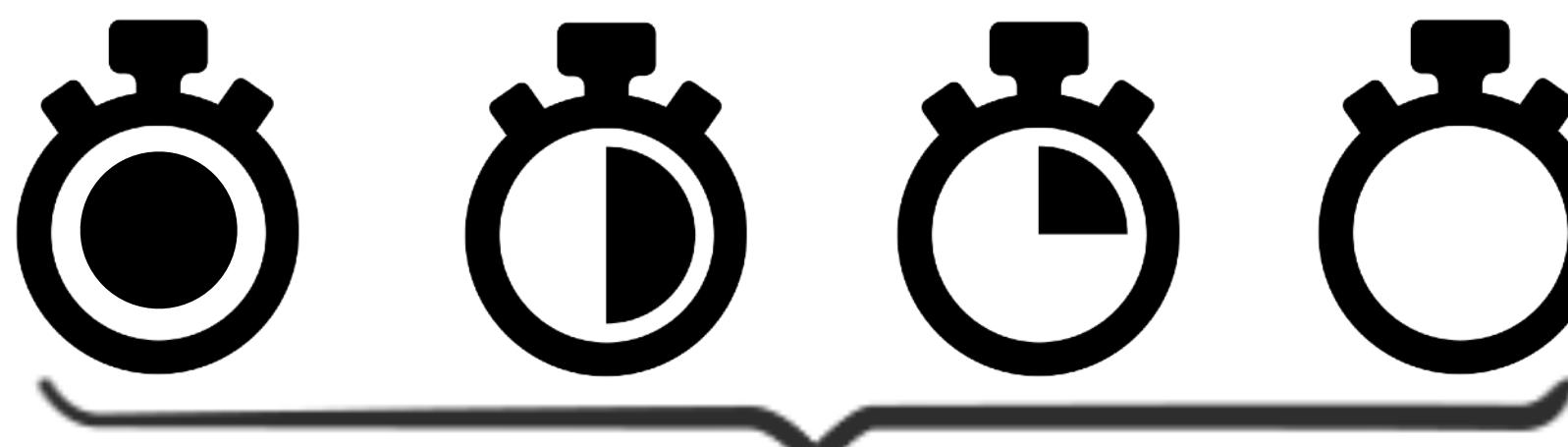
$y^{(1)}$

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

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



$a^{(n)}$



$y^{(n)}$

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

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

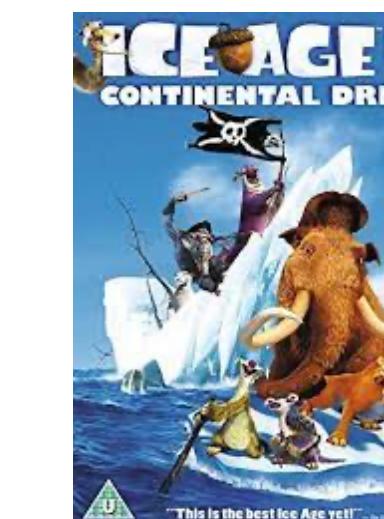
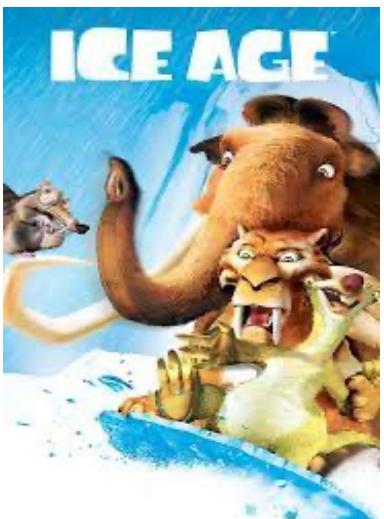
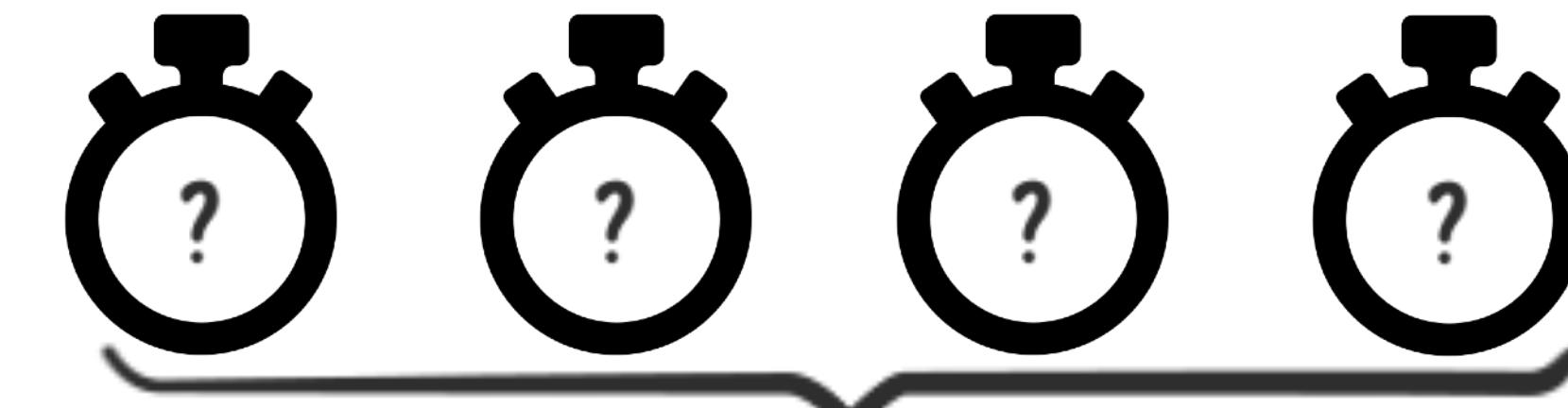
Goal



⋮ ⋮

$\tilde{\mathbf{a}}^{(1)}$

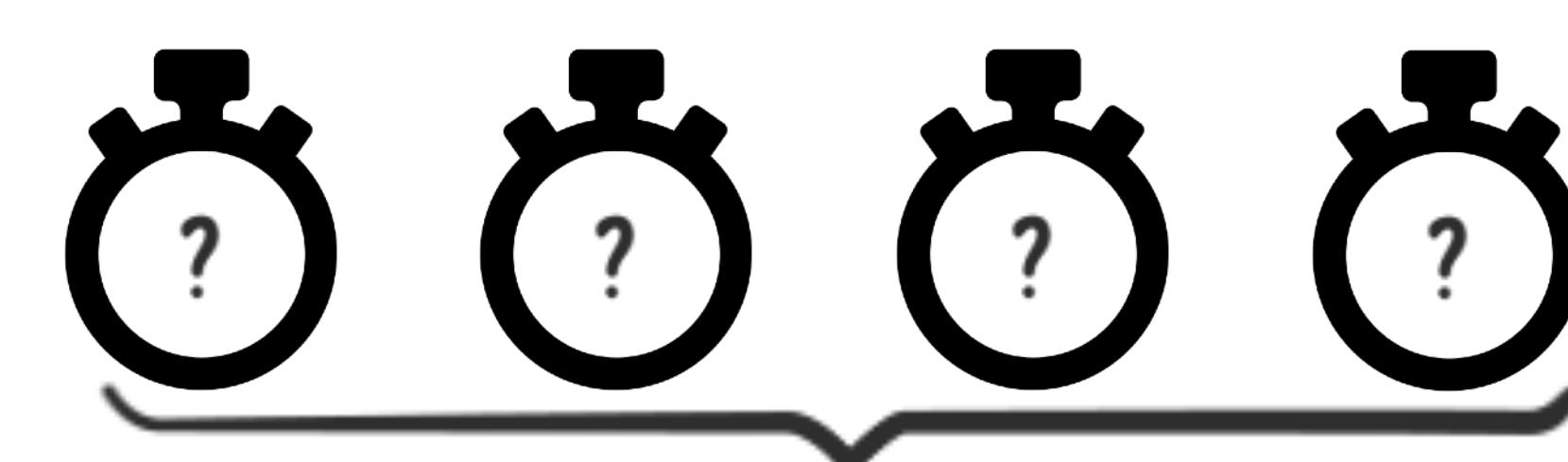
$\mathbf{y}^{(1)}(\tilde{\mathbf{a}}^{(1)}) ?$



⋮ ⋮

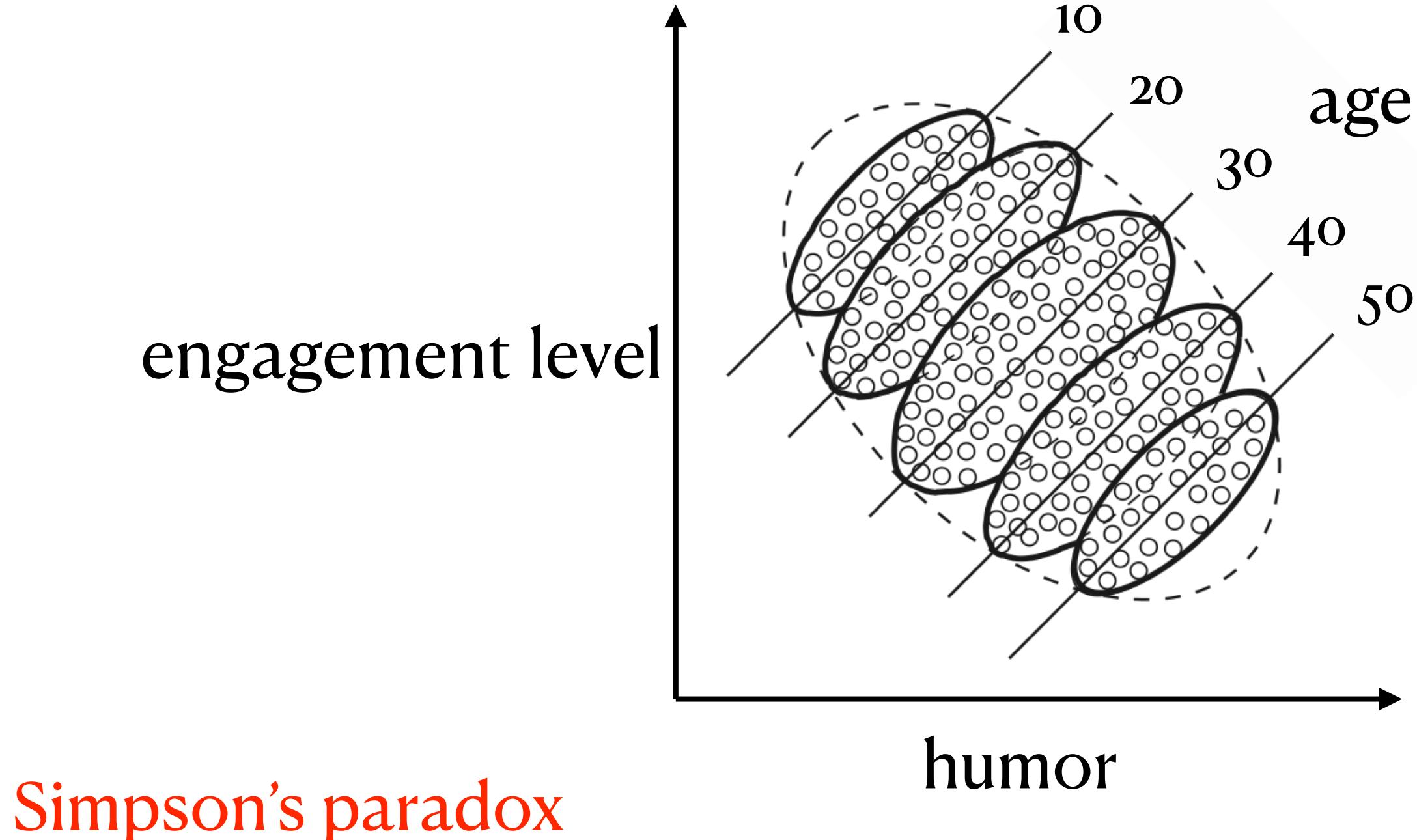
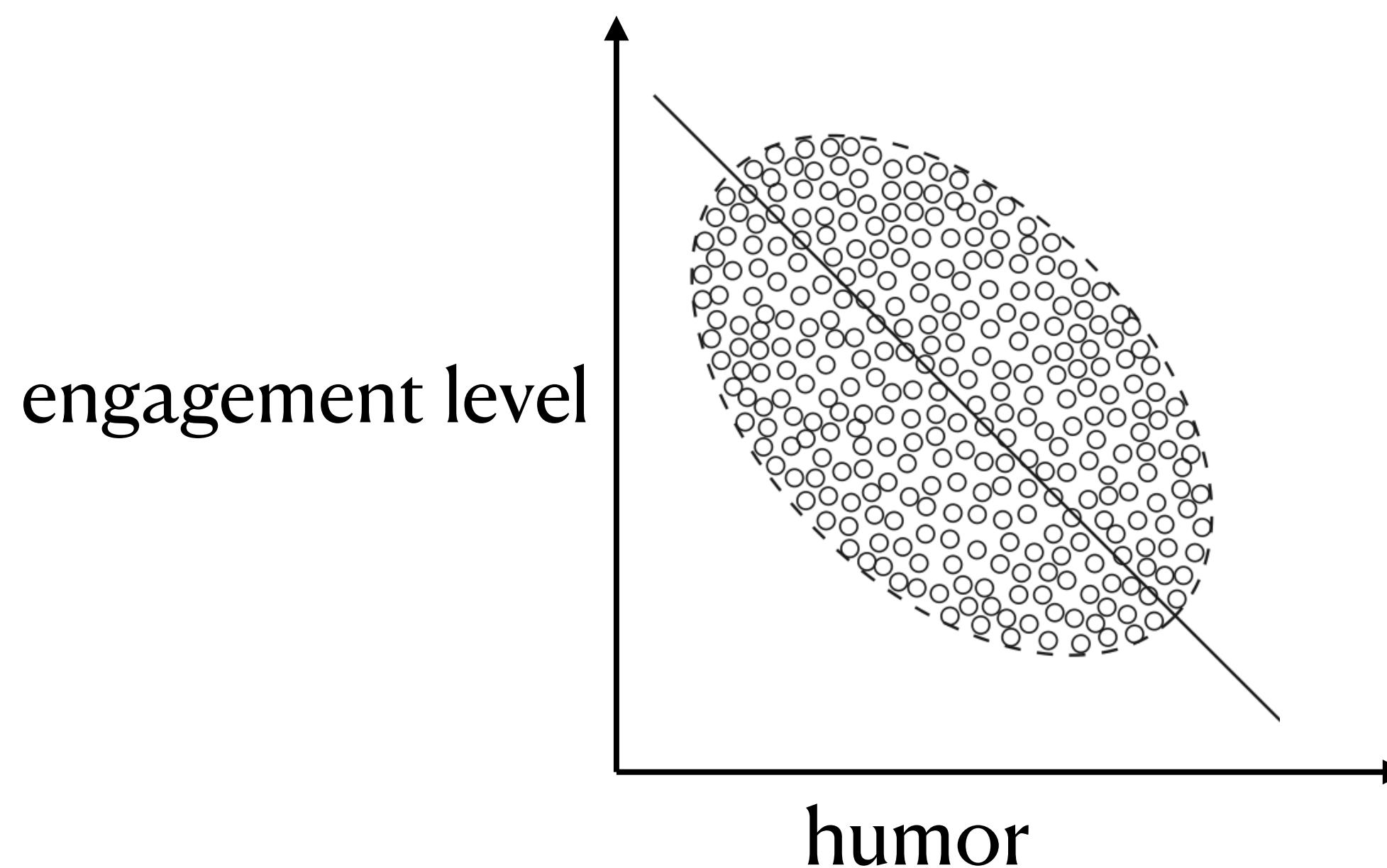
$\tilde{\mathbf{a}}^{(n)}$

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

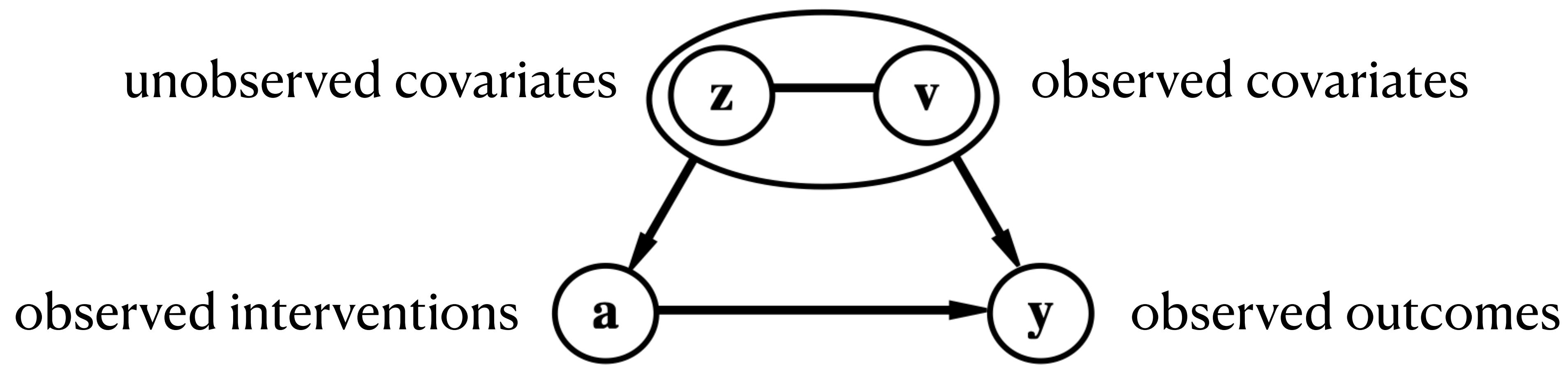


Challenges

1. unobserved factors → **spurious associations**
2. users → **heterogeneous**
3. each user → a **single** interaction trajectory



Problem Setup

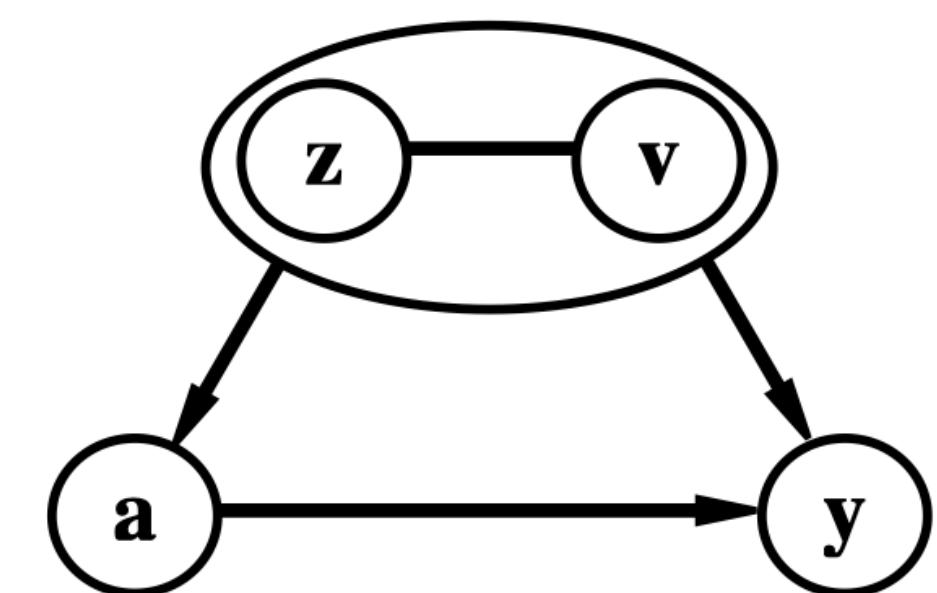
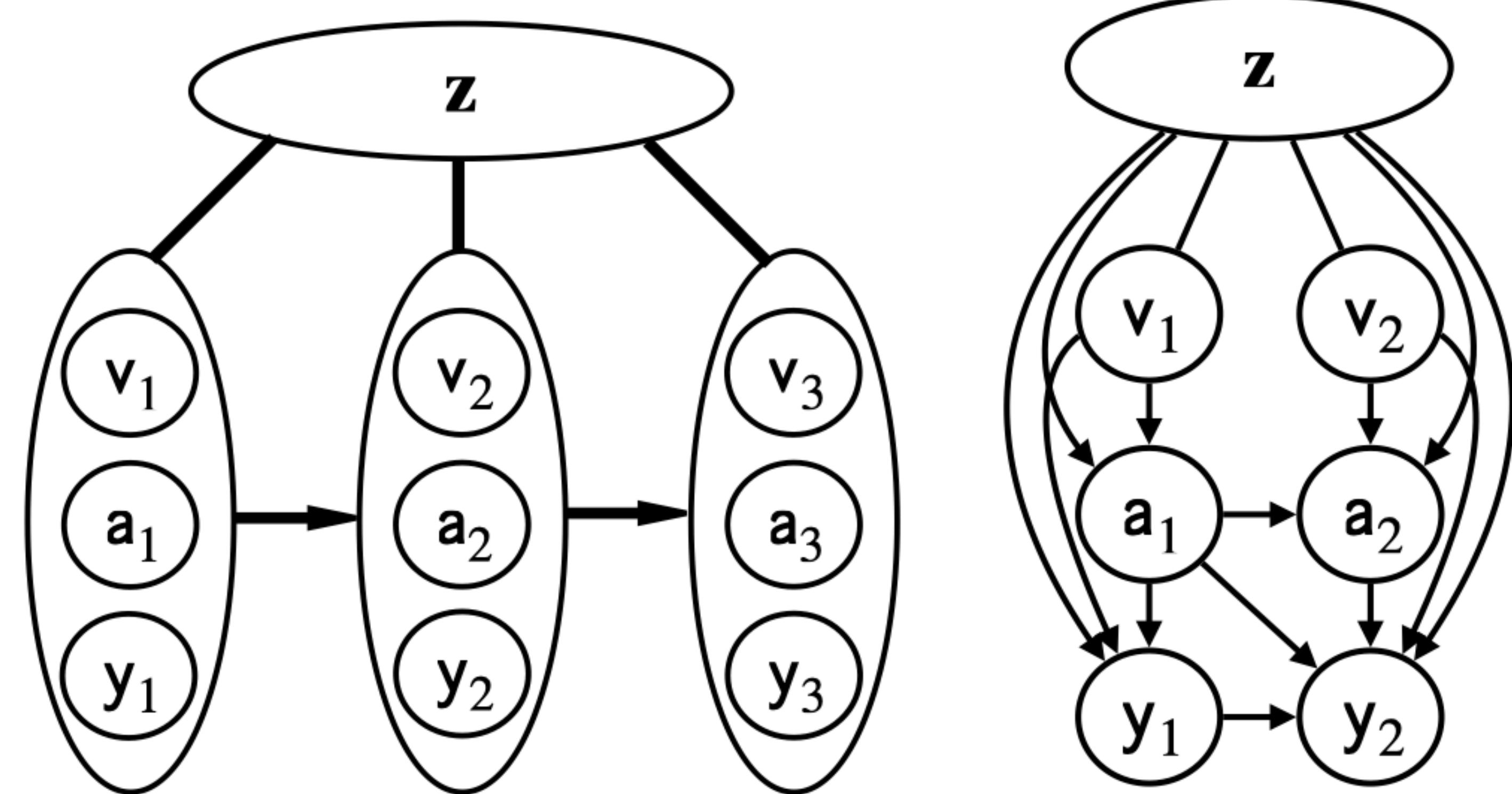


n heterogenous and independent (SUTVA) users

one observation per user - $\{v^{(i)}, a^{(i)}, y^{(i)}\}_{i=1}^n$ (high-dim)

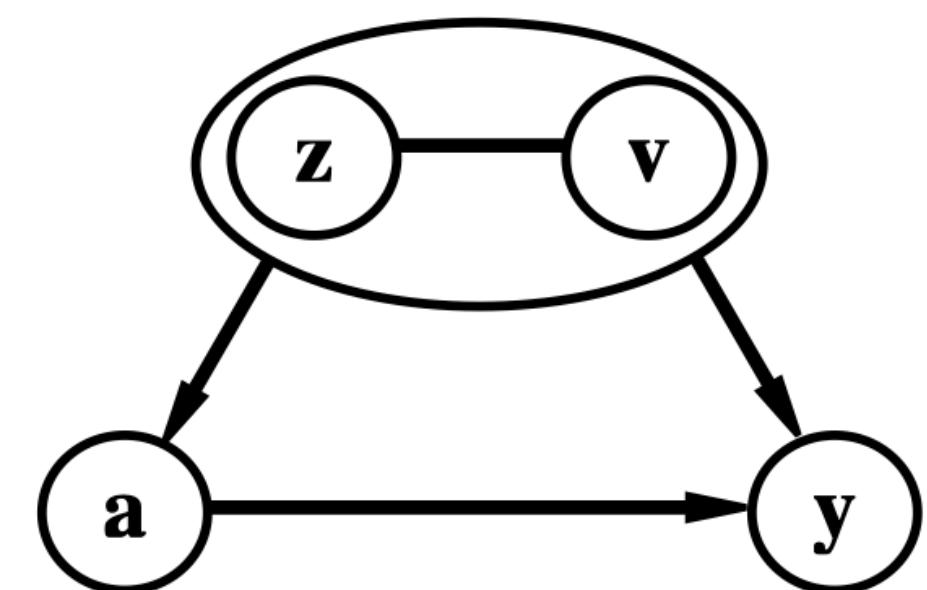
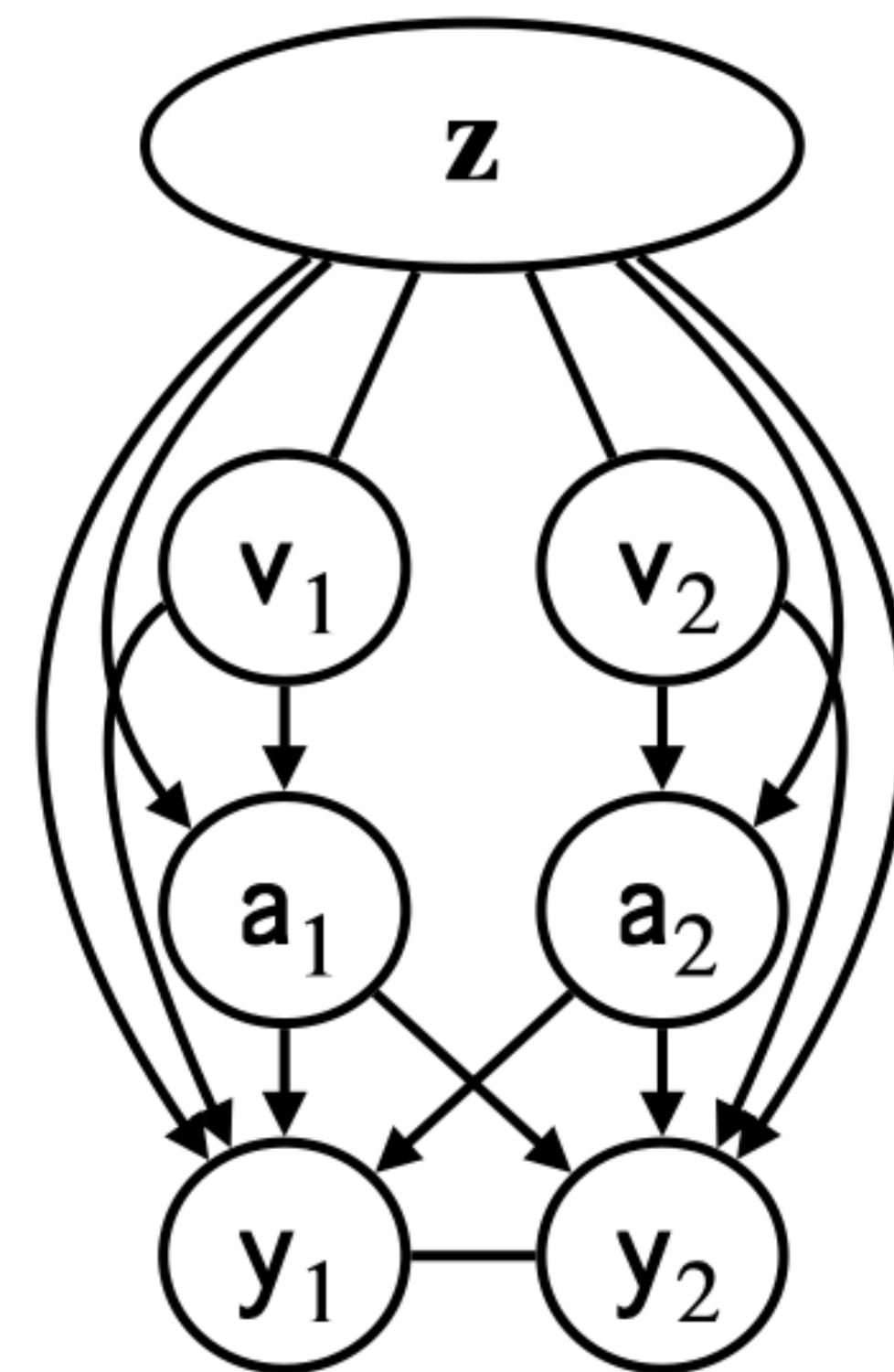
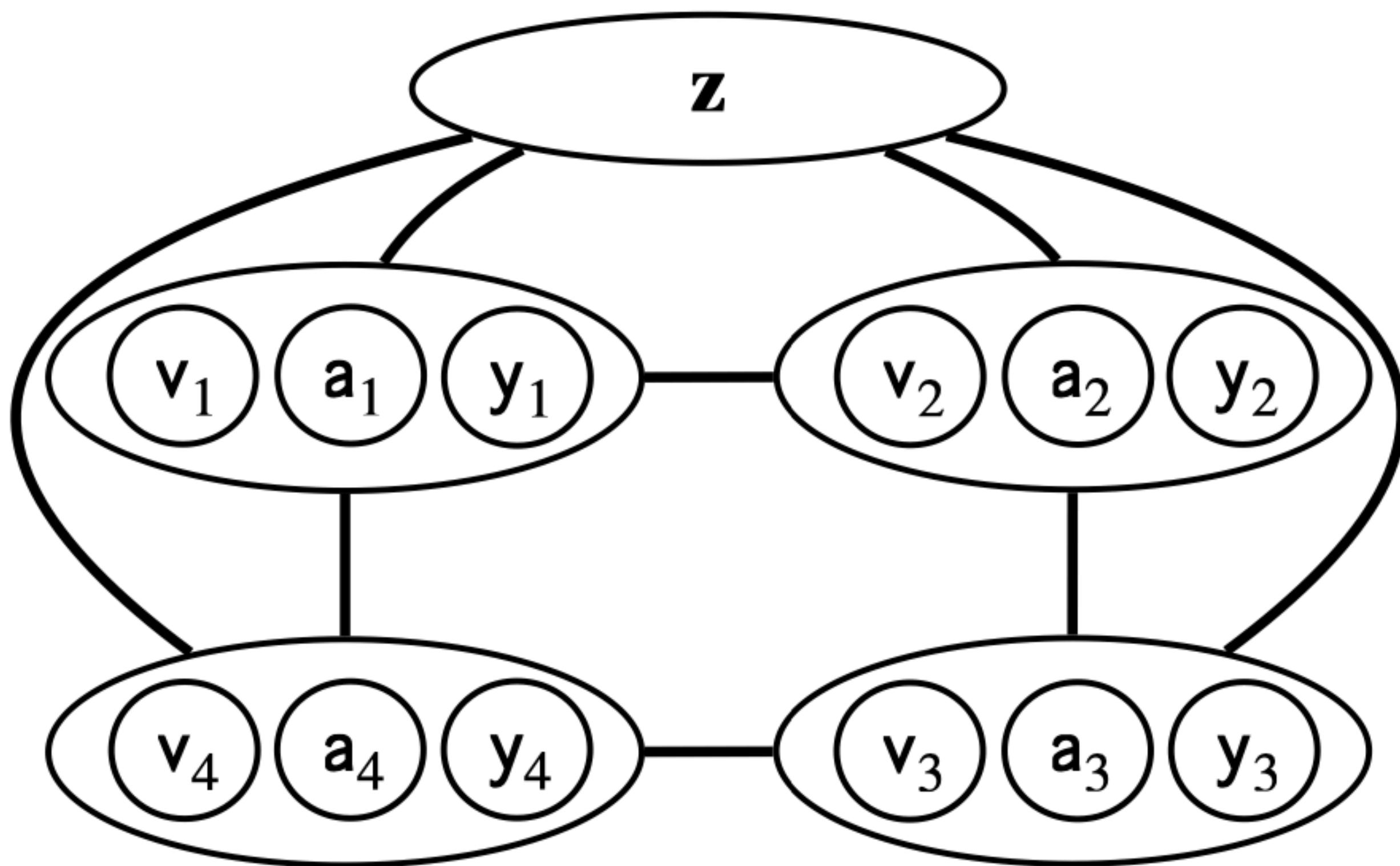
Examples

Sequential setting



Examples

Network setting



Goal: Counterfactual Questions

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

≡

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

suffices to learn $p(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \mathbf{z}^{(i)}, \mathbf{v}^{(i)})$ for all $i \in [n]$

Heterogeneity: each user may have different *unobserved* \mathbf{z}

Can we learn n different distributions

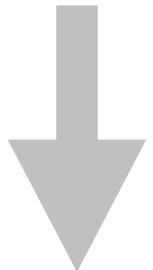
$p(\mathbf{y} = \cdot | \mathbf{a} = \cdot, \mathbf{z}^{(i)}, \mathbf{v}^{(i)}) \quad i \in [n]$

with *one* sample per distribution?

Our Approach

Model the joint distribution of $\mathbf{w} \triangleq (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y})$ as a particular **exponential family**

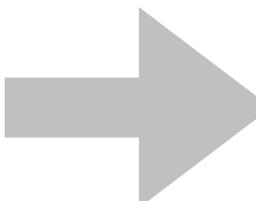
$$p_{\phi, \Phi}(\mathbf{w}) \propto \exp\left(\phi^\top \mathbf{w} + \mathbf{w}^\top \Phi \mathbf{w}\right)$$



$$p(y | \mathbf{a}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}) \propto \exp\left(\left[\phi_y^\top + 2\mathbf{z}^{(i)\top} \Phi_{z,y} + 2\mathbf{v}^{(i)\top} \Phi_{v,y} + 2\mathbf{a}^\top \Phi_{a,y} \right] \mathbf{y} + \mathbf{y}^\top \Phi_{y,y} \mathbf{y}\right)$$

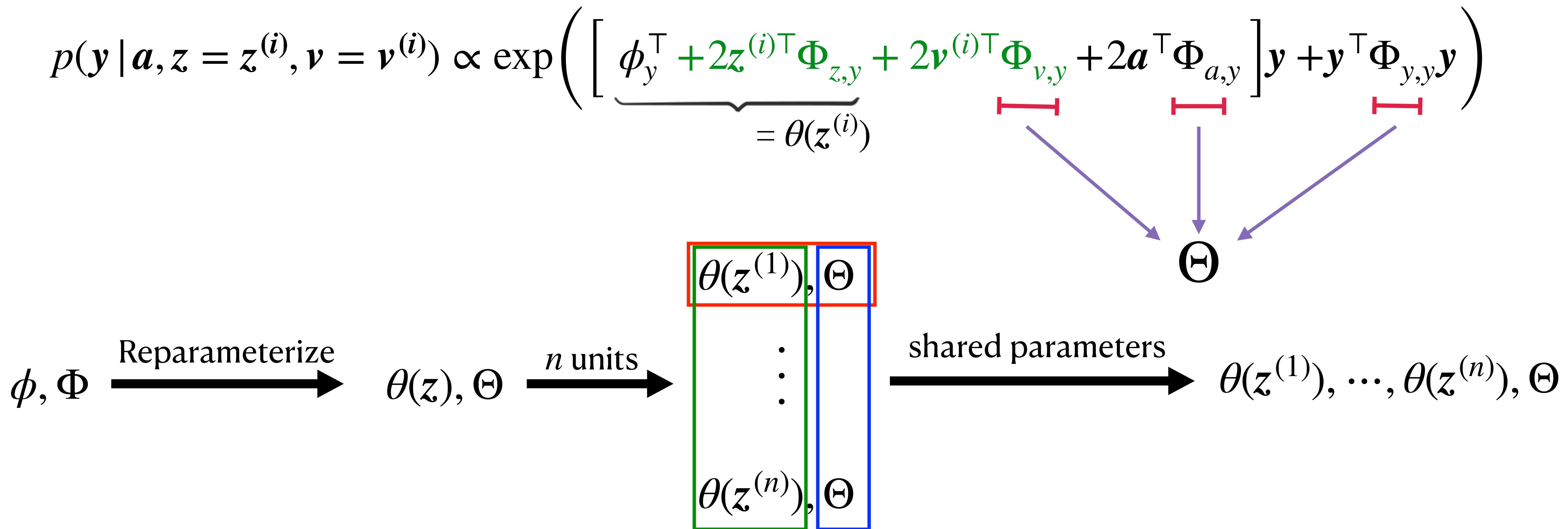
different for different users

n heterogeneous conditional distributions



same exp. family but with diff. parameters

Our Approach



1. If $\mathbf{z}^{(1)} = \dots = \mathbf{z}^{(n)} \rightarrow$ a single exponential family with n samples [B '15, KM '17, SSW '21, VML '22]
2. If $n = 1 \rightarrow$ a single exponential family with one sample (assume Θ is known) [KDDGD '21, MHBM' 21]

Inference Tasks

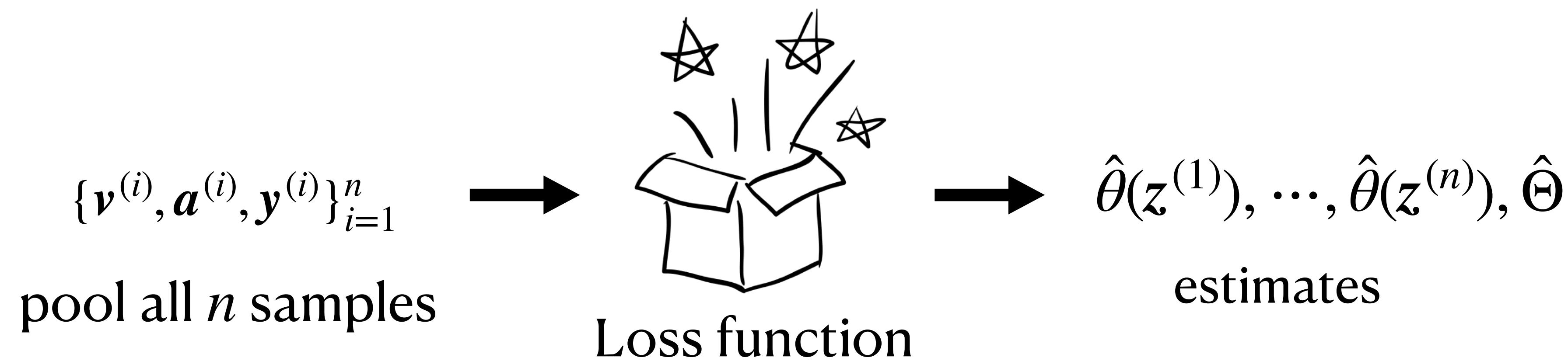
1. Parameters
 - A. User-level – $\theta^\star(z^{(i)})$ for all $i \in [n]$ → counterfactual distribution
 - B. Population-level – Θ^\star
2. Expected potential outcomes – $E[y^{(i)}(\tilde{a}^{(i)}) | z = z^{(i)}, v = v^{(i)}]$ → counterfactual mean

Complexity of Parameter Space

Λ_θ – the set containing the true model parameters $\{\theta^\star(z^{(i)})\}_{i \in [n]}$

	Linear combination of k -known vectors	s -sparse linear combination of k -known vectors
$M(\epsilon) = \log N(\Lambda_\theta, \epsilon)$	$O\left(k \cdot \log\left(1 + \frac{1}{\epsilon}\right)\right)$	$O\left(s \log k \cdot \log\left(1 + \frac{1}{\epsilon}\right)\right)$
$M_n(\epsilon) = nM(n\epsilon)$	$O\left(\frac{k}{\epsilon}\right)$	$O\left(\frac{s \log k}{\epsilon}\right)$

Parameter Estimation



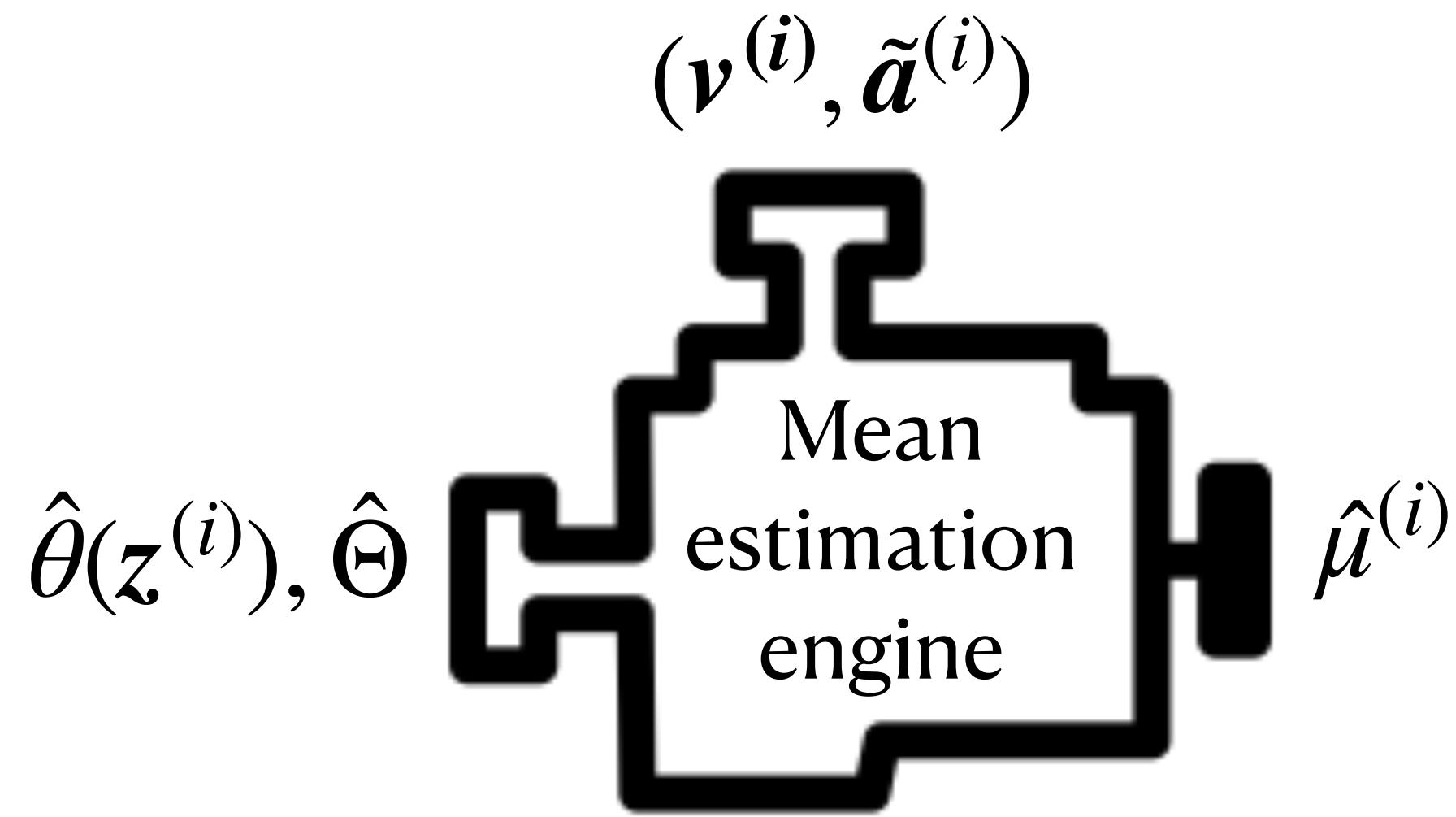
$$\|\Theta^\star - \hat{\Theta}\|_{2,\infty} \leq \epsilon$$

$$\text{For all } i, \text{MSE}\left(\theta^\star(z^{(i)}), \hat{\theta}(z^{(i)})\right) \leq \max \left\{ \epsilon^2, \frac{M(c)}{p} \right\}$$

$$\text{when } n \geq O\left(\frac{p^2(p + M_n(\epsilon^2))}{\epsilon^4}\right)$$
$$\text{when } n \geq O\left(\frac{p^2(pM(c) + M_n(\epsilon^2))}{\epsilon^4}\right)$$

Outcome Estimation

Expected potential outcomes – $\mu^{(i)} \triangleq \mathbf{E}\left[y^{(i)}(\tilde{a}^{(i)}) \mid \mathbf{z} = z^{(i)}, \mathbf{v} = v^{(i)}\right]$



For all i , $MSE\left(\mu^{(i)}, \hat{\mu}^{(i)}\right) \leq \epsilon^2 + \frac{M(c)}{p}$ when $n \geq O\left(\frac{p^2(pM(c) + M_n(\epsilon^2))}{\epsilon^4}\right)$

$$\mathbf{x} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$$

An Application

Denoise user-wise data

No systematically unobserved covariates \mathbf{z}

Noisy observed data = true data + measurement error

$$\bar{\mathbf{x}}$$

$$\mathbf{x}$$

$$\Delta\mathbf{x}$$

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

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

Goal: Estimate the true data

$$\mathbf{x} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$$

An Application

Denoise user-wise data

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

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

Our results can recover data such that

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

Remainder Of The Talk

The Loss Function, And Why It Works

Log likelihood as the loss function = Max. Likelihood Est.

not the right answer (computational challenge)

An alternative

a *proper* loss function, computationally efficient

ingredients of analysis

population-level parameter est through concentration due to many samples

user-level parameter est through *single* sample concentration due to

Log Sobolev + Dobrushin's criteria

Maximum Likelihood Estimation

Consider an exponential family $p_\theta(\mathbf{x}) = \frac{\exp(\theta^\top f(\mathbf{x}))}{Z(\theta)}$

Given independent $\mathbf{x}_1, \dots, \mathbf{x}_n \sim p_{\theta^*}(\mathbf{x}), \exp(\theta^* \mathbf{w} + \mathbf{w}^\top \Phi \mathbf{w})$

$$\hat{\theta}_{MLE} \triangleq \operatorname{argmax}_\theta \text{Log-Likelihood}(\theta)$$

$$\text{Log-likelihood}(\theta) = \frac{1}{n} \sum_{i \in [n]} \theta^\top f(\mathbf{x}_i) - \log Z(\theta)$$

- A. Consistency ?
- B. Asymptotic normality ?
- C. Asymptotic efficiency ?
- D. Computational tractability ? **How to compute $Z(\theta)$?**

An Alternative

Consider an exponential family $p_\theta(\mathbf{x}) = \frac{\exp(\theta^\top f(\mathbf{x}))}{Z(\theta)}$

Given independent $\mathbf{x}_1, \dots, \mathbf{x}_n \sim p_{\theta^*}(\mathbf{x})$, learn $\hat{\theta}$

$$\hat{\theta} \triangleq \operatorname{argmin}_{\theta} \ell(\theta) \quad \text{likelihood}$$

$$\ell(\theta) = -\frac{1}{n} \sum_{i \in [n]} \exp(-\theta^\top f(\mathbf{x}_i)) \quad \prod_{i \in [n]} \frac{1}{Z(\theta)} \exp(\theta^\top f(\mathbf{x}_i))$$

Convex constraint on θ \longrightarrow convex optimization problem

avoids $Z(\theta)$

but is it any good?

A Proper Loss Function

$$\ell(\theta) = \frac{1}{n} \sum_{i \in [n]} \exp(-\theta^\top f(x_i))$$

$$\mathbb{E}_{\theta^*}[\ell(\theta)] = \mathbb{E}_{\theta^*}[\exp(-\theta^\top f(x))]$$

Theorem [SSW '21]. $\operatorname{argmin}_\theta \mathbb{E}_{\theta^*}[\ell(\theta)] = \operatorname{argmin}_\theta KL(\mathcal{U} \| p_{\theta^* - \theta})$

$\mathcal{U} \rightarrow$ uniform distribution



$\mathbb{E}_{\theta^*}[\ell(\theta)]$ is minimized uniquely when $\theta = \theta^*$

Properties

[SSW'21]

$$\hat{\theta} \in \arg \min_{\theta} \ell(\theta) = \frac{1}{n} \sum_{i \in [n]} \exp(-\theta^\top f(x_i))$$

- A. Consistency ? \$\hat{\theta}\$ is MLE w.r.t. \$p_{\theta^* - \theta}\$ (not \$p_\theta\$)
- B. Asymptotic normality ?
- C. Asymptotic efficiency ?
- D. Computational tractability ? No need to compute \$Z(\theta)\$!
- E. Finite sample guarantees ?

Finite Sample Guarantees

[ssw'21]

$$\hat{\theta} \in \arg \min_{\theta} \ell(\theta) = \frac{1}{n} \sum_{i \in [n]} \exp(-\theta^\top f(x_i))$$

$$\|\theta^* - \hat{\theta}\|_2 \sim n^{-1/4}$$

Proof ingredients —

- A. Concentration of gradient
- B. Anti-concentration of Hessian (Restricted strong convexity)

Hoeffding's
inequality

Back To Our Setting

Condition on \mathbf{z}

Recall the joint distribution of $\mathbf{w} = (\mathbf{z}, \mathbf{v}, \mathbf{a}, \mathbf{y})$

$$p_{\phi, \Phi}(\mathbf{w}) \propto \exp\left(\phi^\top \mathbf{w} + \mathbf{w}^\top \Phi \mathbf{w}\right)$$

Letting $\mathbf{x} \triangleq (\mathbf{v}, \mathbf{a}, \mathbf{y})$, the conditional distribution of \mathbf{x} given \mathbf{z} can be written as

$$p_{\theta(z), \Theta}(\mathbf{x} | \mathbf{z}) \propto \exp\left(\left[\theta(\mathbf{z})\right]^\top \mathbf{x} + \mathbf{x}^\top \Theta \mathbf{x}\right)$$
$$p(y | \mathbf{a}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}) \propto \exp\left(\left[\phi_y + 2\mathbf{z}^{(i)\top} \Phi_{z,y} + 2\mathbf{v}^{(i)\top} \Phi_{v,y} + 2\mathbf{a}^\top \Phi_{a,y} \right] y + y^\top \Phi_{y,y} y \right)$$

The diagram illustrates the relationship between the two equations. A purple arrow points from the first equation down to the second. Another purple arrow points from the right side of the first equation to the right side of the second equation. Three red horizontal bars with brackets underneath the terms in the second equation group the terms by color: green for the first term, blue for the second, and red for the third.

Structure On Parameters

$$p_{\theta(z), \Theta}(x | z) \propto \exp\left(\left[\theta(z)\right]^\top x + x^\top \Theta x\right)$$

- (A) Every element of $\theta^*(z^{(1)}), \dots, \theta^*(z^{(n)})$, and Θ^* are bounded
- (B) Every row of Θ^* is sparse

$$\Lambda_\theta \triangleq \{\theta : \theta \text{ is consistent with (A)} + \boxed{\text{low complexity}}\} \rightarrow \begin{array}{l} \text{required to} \\ \text{provide} \\ \text{meaningful} \\ \text{guarantees} \end{array}$$
$$\Lambda_\Theta \triangleq \{\Theta : \Theta \text{ is consistent with (A) and (B)}\}$$

Learning Population-level Parameter

$$\mathcal{L}(\underline{\Theta}) = \sum_{t \in [p]} \frac{1}{n} \sum_{i \in [n]} \exp\left(-[\theta_t^{(i)} + 2\Theta_t^\top \mathbf{x}^{(i)}] x_t^{(i)}\right) \text{ where } \underline{\Theta} \triangleq [\theta^{(1)}, \dots, \theta^{(n)}, \Theta]$$

$\Lambda_\Theta \rightarrow$ (A) bounded elements, (B) sparse rows

Λ_Θ places independent constraints on the rows of Θ

p independent optimization problems

$$\mathcal{L}_t(\underline{\Theta}_t) = \frac{1}{n} \sum_{i \in [n]} \exp\left(-[\theta_t^{(i)} + 2\Theta_t^\top \mathbf{x}^{(i)}] x_t^{(i)}\right) \text{ for all } t \in [p] \longrightarrow \widehat{\Theta}_t$$

$$\|\Theta_t^* - \widehat{\Theta}_t\|_{2,\infty} \sim n^{-1/4}$$

- A. Concentration of gradient
B. Anti-concentration of Hessian
- $\left. \right\} \text{Hoeffding's inequality}$

Learning Unit-level Parameter

$$\mathcal{L}(\theta^{(1)}, \dots, \theta^{(n)}) = \sum_{t \in [p]} \frac{1}{n} \sum_{i \in [n]} \exp\left(-[\theta_t^{(i)} + 2\widehat{\Theta}_t^\top \mathbf{x}^{(i)}]x_t^{(i)}\right)$$

$\Lambda_\theta \rightarrow$ (A) bounded elements, (B) low complexity

$\theta^{(1)}, \dots, \theta^{(n)} \in \Lambda_\theta^n$ places independent constraints on units, i.e., $\theta^{(i)} \in \Lambda_\theta$ for all $i \in [n]$

n independent optimization problems

$$\mathcal{L}^{(i)}(\theta^{(i)}) = \sum_{t \in [p]} \exp\left(-[\theta_t^{(i)} + 2\widehat{\Theta}_t^\top \mathbf{x}^{(i)}]x_t^{(i)}\right) \text{ for all } i \in [n] \quad \longrightarrow \quad \widehat{\theta}^{(i)}$$
$$\|\theta^\star(z^{(i)}) - \widehat{\theta}^{(i)}\|_2 \sim \max\{n^{-1/4}, M\}$$

- A. Concentration of gradient
B. Anti-concentration of Hessian
- $\left. \begin{array}{l} \\ \end{array} \right\}$ Logarithmic
Sobolev
inequality