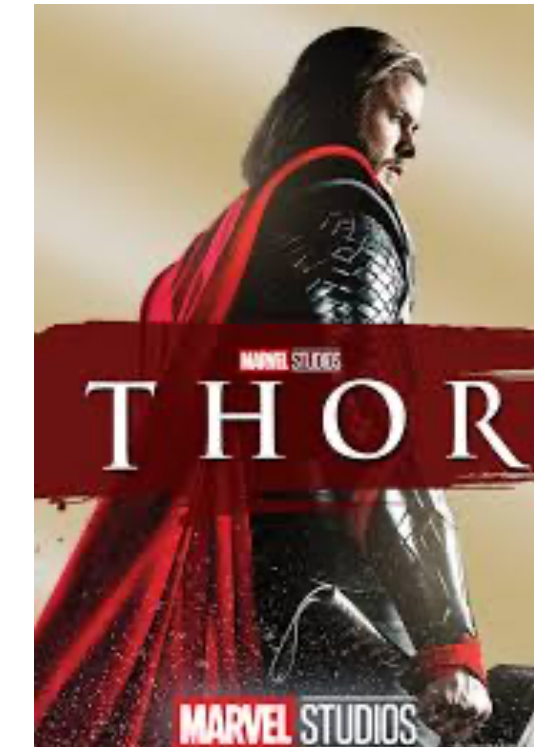


Can we do counterfactual inference in the presence of unobserved confounding with one sample?

with Raaz Dwivedi, Devavrat Shah, and Greg Wornell

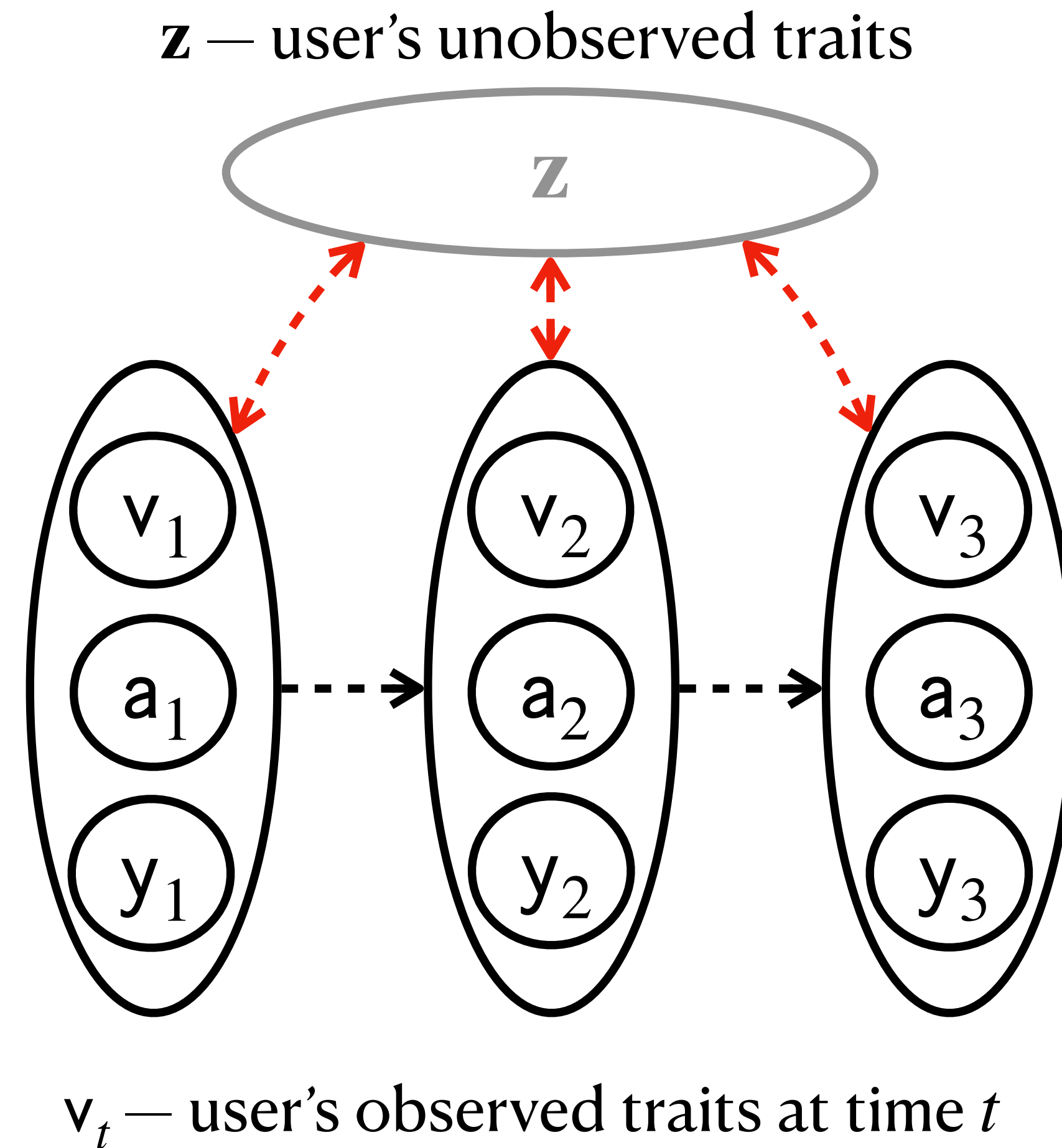
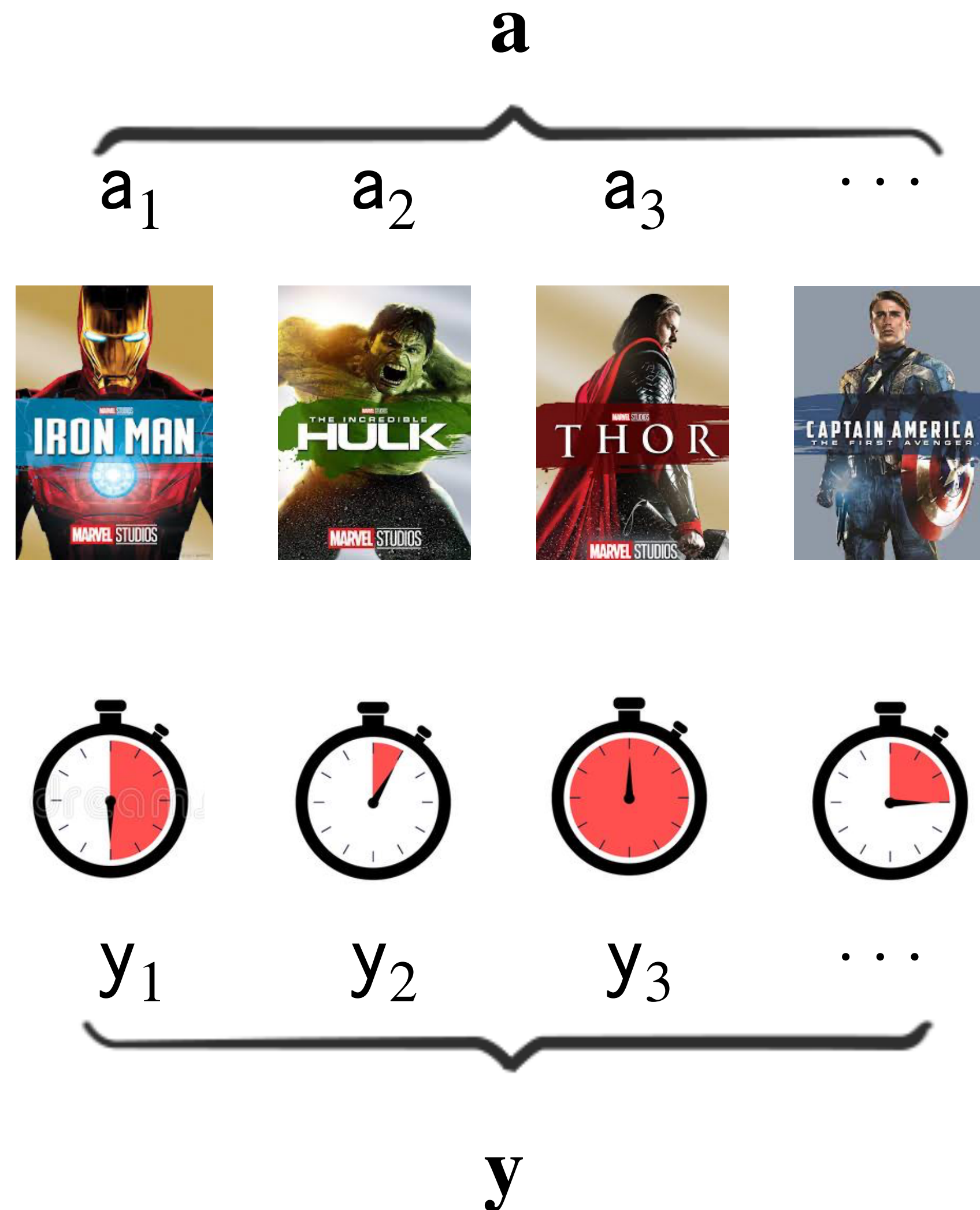


Sequential Recommender System



Sequential Recommender System

A graphical model

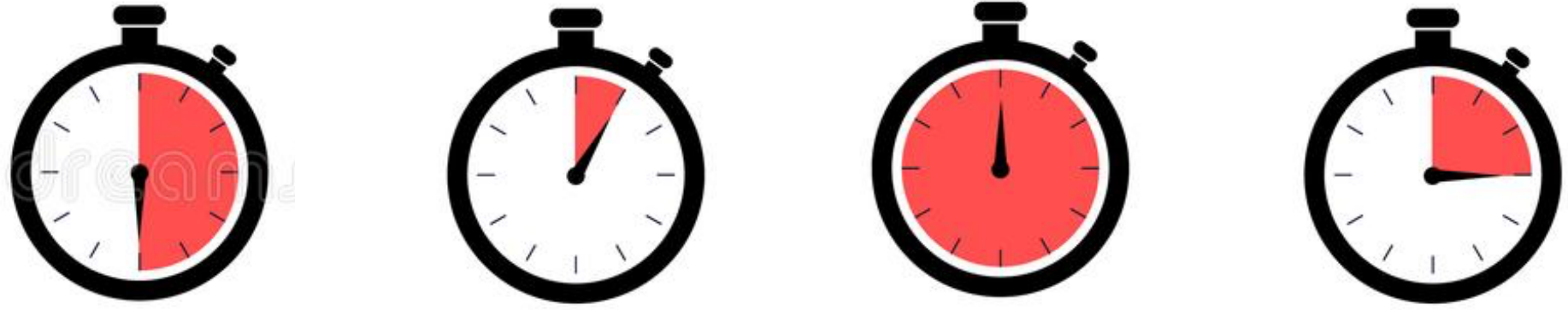


Sequential Recommender System

Data



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



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

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

⋮

⋮

⋮



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



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

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

Sequential Recommender System

Goal



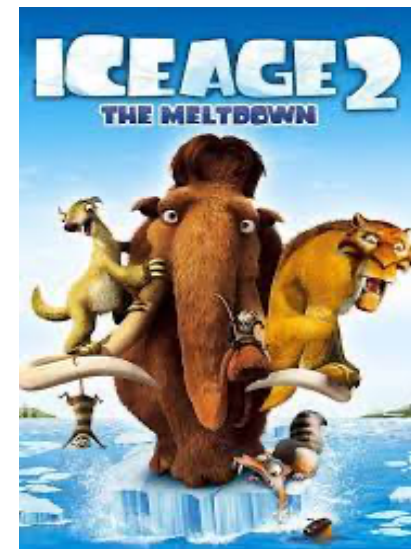
$$y^{(1)}(\tilde{a}^{(1)})$$



⋮

⋮

⋮



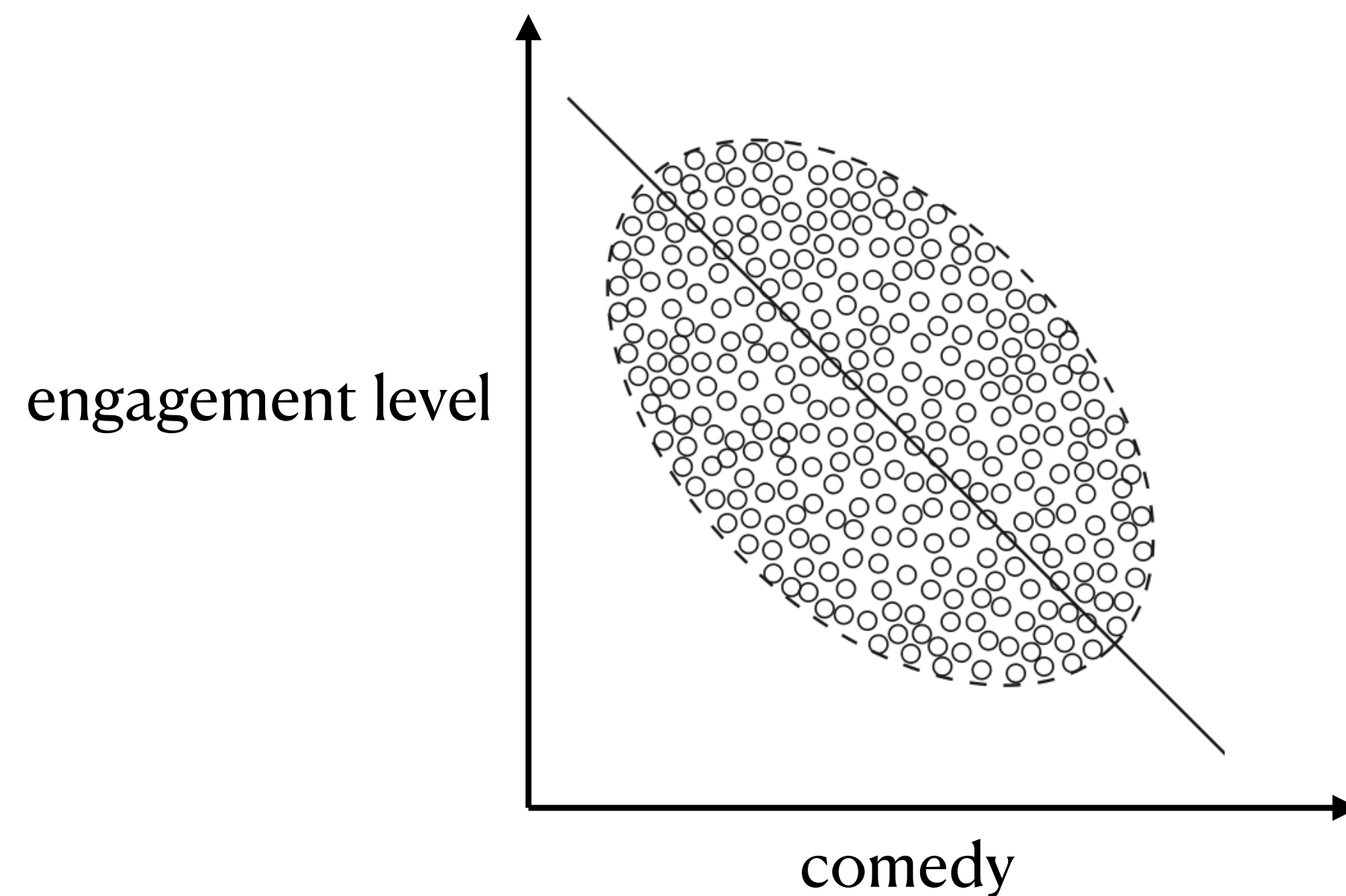
$$y^{(n)}(\tilde{a}^{(n)})$$



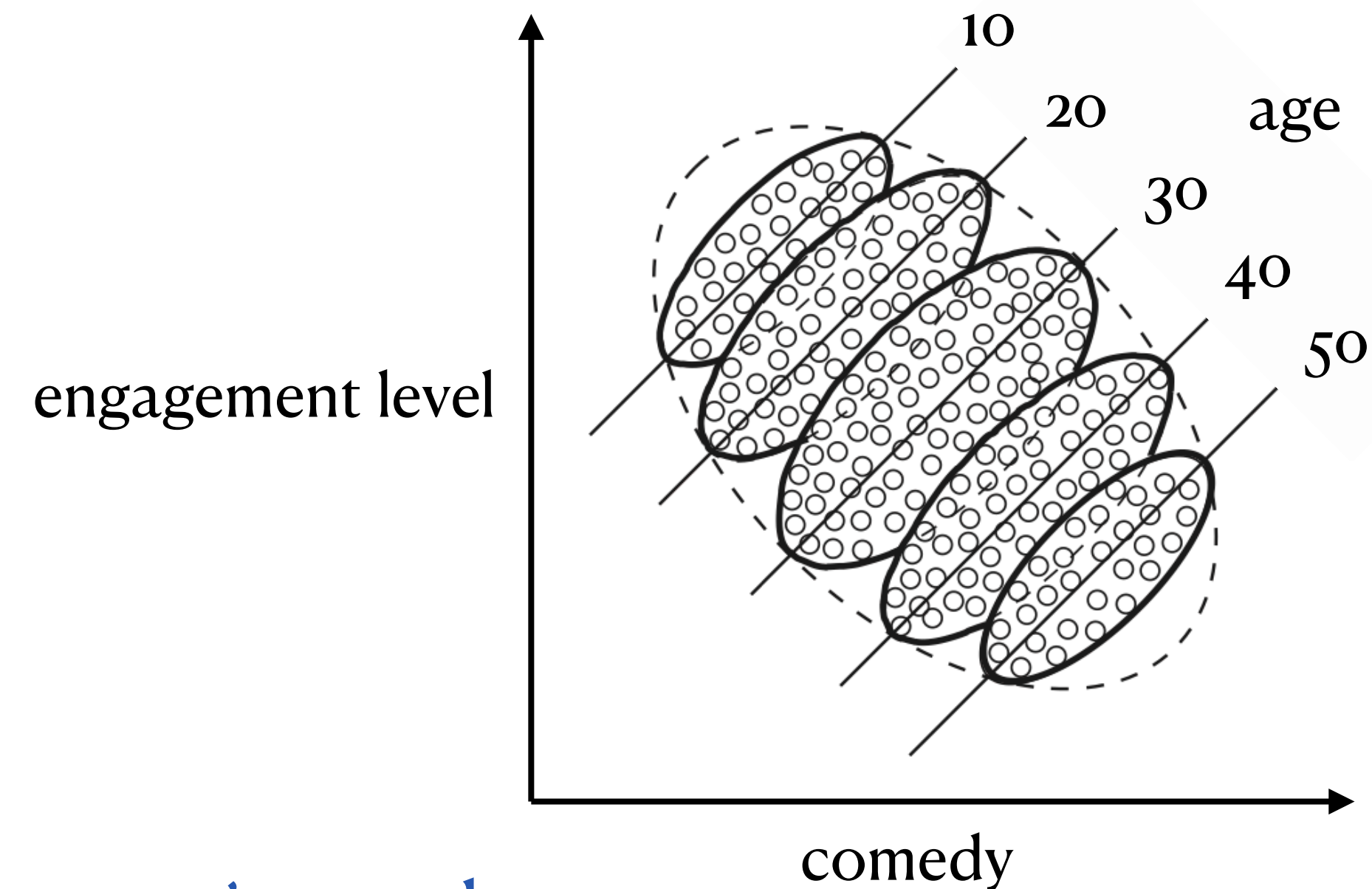
Sequential Recommender System

Challenges

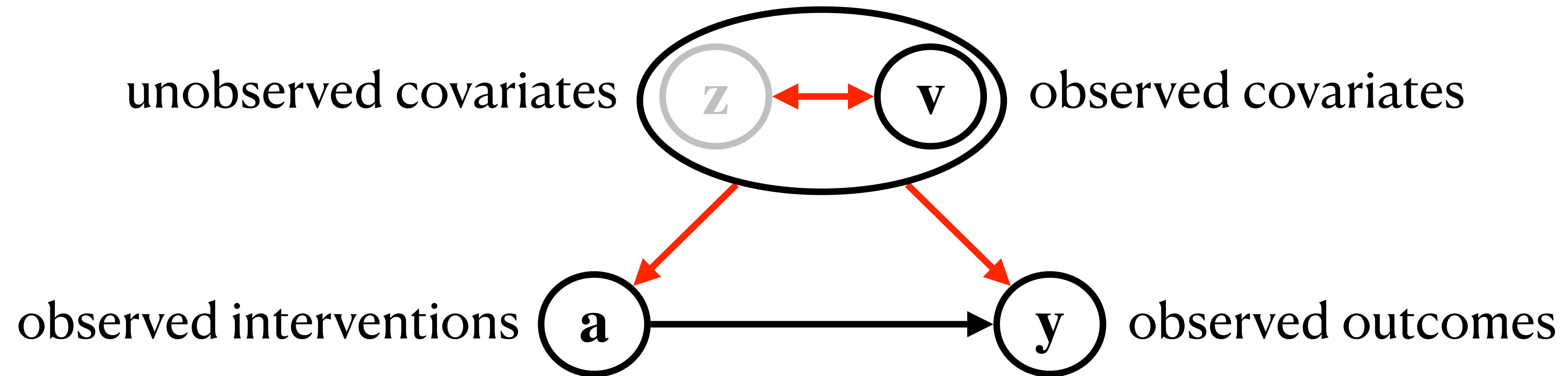
- **unobserved factors** → **spurious associations**
- users could be **heterogeneous**
- each user → a **single** interaction trajectory



Simpson's paradox



Problem Setup



The micro-level graphical is consistent with this macro-level graphical model

- n heterogenous and independent units
- only one observation per unit - $\{v^{(i)}, a^{(i)}, y^{(i)}\}_{i=1}^n$

Goal

Counterfactual questions for these n units

For every unit $i \in [n]$, what would the potential outcome $y^{(i)}(\tilde{\mathbf{a}}^{(i)})$ be while $\mathbf{z} = \mathbf{z}^{(i)}$ and $\mathbf{v} = \mathbf{v}^{(i)}$?

Under SUTVA, learning unit-level counterfactual distributions is equivalent to learning unit-level conditional distributions:

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

Challenges

1. Unobserved confounding — \mathbf{z} introduces statistical dependence between \mathbf{a} and \mathbf{y}
 2. Heterogeneity — $(\mathbf{z}^{(i)}, \mathbf{v}^{(i)})$ could be different for different units
- ➡ we only observe one realization that is consistent with $p(\mathbf{y} = \cdot \mid \mathbf{a}, \mathbf{z}^{(i)}, \mathbf{v}^{(i)})$

Is it possible to learn n heterogeneous distributions
with only one sample per distribution?

Our approach

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

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

- The conditional distribution of \mathbf{y} given $\mathbf{a}, \mathbf{z}, \mathbf{v}$ can be written as

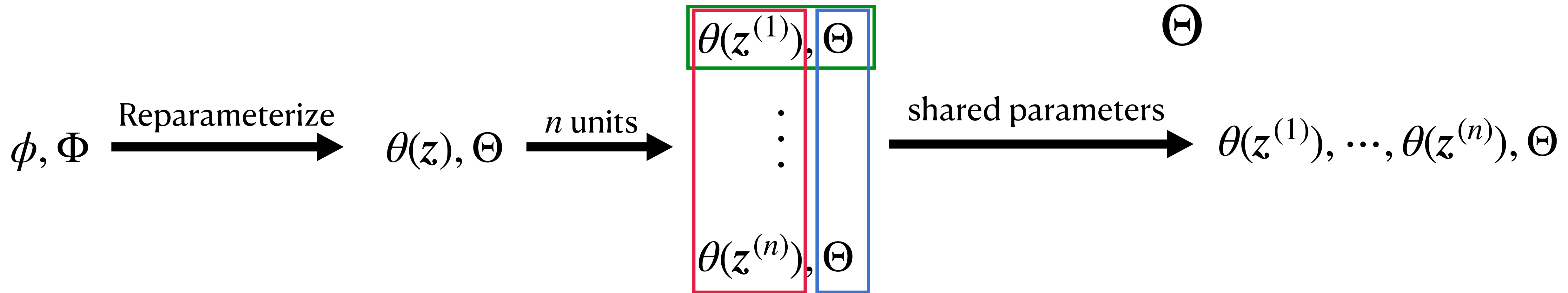
$$p(\mathbf{y} | \mathbf{a}, \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}) \propto \exp\left(\left[\underbrace{\phi_y^\top + 2\mathbf{z}^{(i)\top} \Phi_{z,y} + 2\mathbf{v}^{(i)\top} \Phi_{v,y}}_{\text{different for different units}} + 2\mathbf{a}^\top \Phi_{a,y} \right] \mathbf{y} + \underbrace{\mathbf{y}^\top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

n heterogeneous conditional distributions 

n distributions from the same exponential family
but with parameters that vary across units

Our approach

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

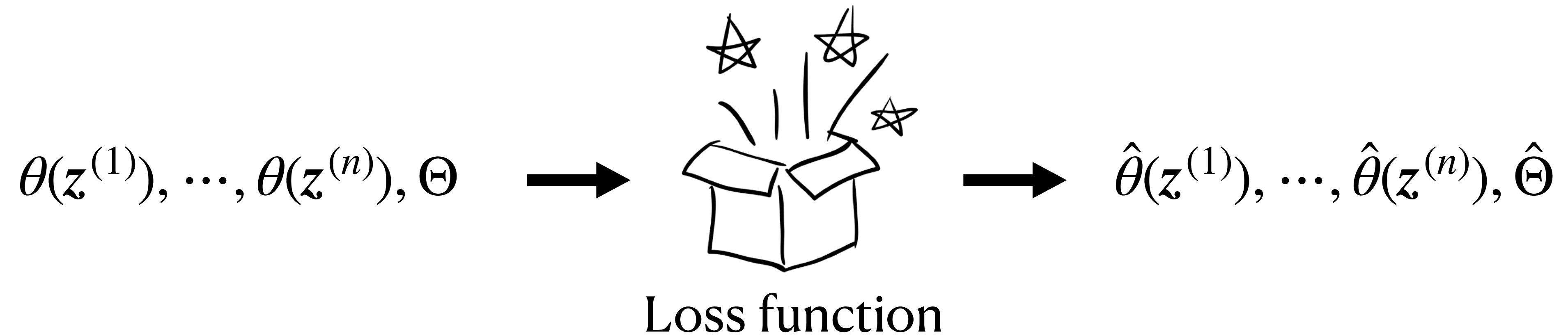


1. If $\mathbf{z}^{(1)} = \dots = \mathbf{z}^{(n)} \rightarrow$ a single exponential family with n samples
2. If $n = 1 \rightarrow$ a single exponential family with one sample (assume Θ is known)

Inference Tasks

1. Parameters
 - A. Unit-level — $\theta^*(\mathbf{z}^{(i)})$ for all $i \in [n]$
 - B. Population-level — Θ^*
2. Expected potential outcomes — $\mathbf{E}[\mathbf{y}^{(i)}(\tilde{\mathbf{a}}^{(i)}) \mid \mathbf{z} = \mathbf{z}^{(i)}, \mathbf{v} = \mathbf{v}^{(i)}]$

Parameter estimation



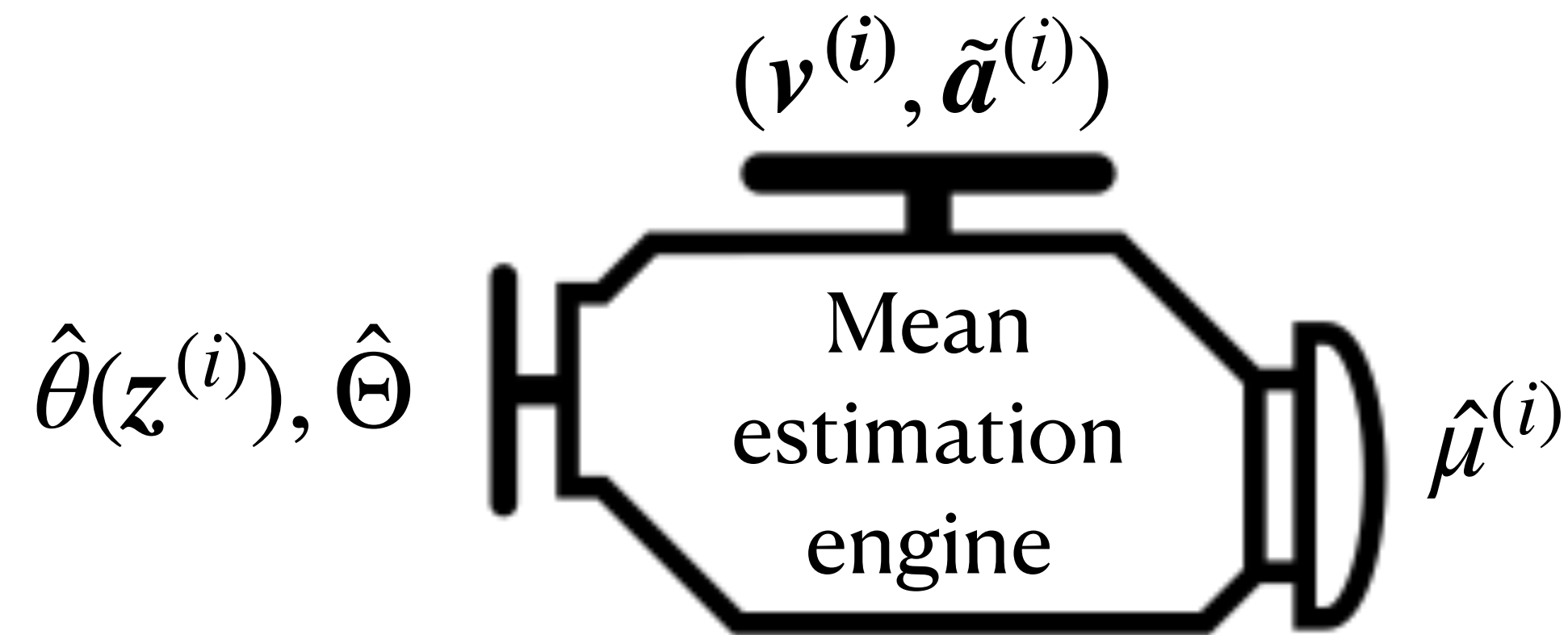
For all row j , $\|\Theta_j^\star - \hat{\Theta}_j\|_2 \leq \epsilon$ whenever $n \geq O\left(\frac{\log(\text{dim})}{\epsilon^4}\right)$

For all unit i , $\|\theta^\star(\mathbf{z}^{(i)}) - \hat{\theta}(\mathbf{z}^{(i)})\|_2 \leq \max\{\epsilon, \text{Comp}\}$ whenever $n \geq O\left(\frac{\text{dim}^2 \text{Comp}^2}{\epsilon^4}\right)$

When the true parameters are s -sparse linear combination of k known vectors, $\text{Comp} = O(s \log(k \cdot \text{dim}))$

Outcome estimation

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



When the true parameters are s -sparse linear combination of k known vectors, for any $\{\tilde{\mathbf{a}}^{(i)} \in \mathcal{A}\}_{i=1}^n$

$$\text{For all unit } i, \text{MSE}(\mu^{(i)}, \hat{\mu}^{(i)}) \leq \frac{s \log(k \cdot \text{dim}) + \epsilon^2}{\text{dim}} \quad \text{whenever } n \geq O\left(\frac{sp^2 \log(k \cdot \text{dim})}{\epsilon^4}\right)$$

Loss function

Condition on \mathbf{z}

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

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

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

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

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

Loss function

Assumptions

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

- (A) Every element of $\theta^\star(\mathbf{z}^{(1)}), \dots, \theta^\star(\mathbf{z}^{(n)})$, and Θ^\star is bounded
- (B) Every row of Θ^\star is sparse
- (C) Every diagonal entry of Θ^\star is zero

- $\Lambda_\theta \triangleq \{ \theta : \theta \text{ is consistent with (A) + low complexity} \}$
- $\Lambda_\Theta \triangleq \{ \Theta : \Theta \text{ is consistent with (A), (B), and (C)} \}$

Loss function



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

- inspired by the conditional distribution of x_t given \mathbf{x}_{-t} and \mathbf{z} —

$$p(x_t | \mathbf{x}_{-t}, \mathbf{z}; \theta_t(\mathbf{z}), \Theta_t) \propto \exp\left([\theta_t(\mathbf{z}) + 2\Theta_t^\top \mathbf{x}]x_t\right)$$

- maps the parameters $\theta^{(1)}, \dots, \theta^{(n)}$, and Θ to \mathcal{L}

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

Loss function

Estimate

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

$$\hat{\underline{\Theta}} = \arg \min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \mathcal{L}(\underline{\Theta})$$

- convex optimization problem
- strictly proper loss function — $\underline{\Theta}^\star$ uniquely maximizes $\mathbb{E}[\mathcal{L}(\cdot)]$

Loss function

Decomposition

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

$$\hat{\underline{\Theta}} = \arg \min_{\underline{\Theta} \in \Lambda_\theta^n \times \Lambda_\Theta} \mathcal{L}(\underline{\Theta})$$

$$\min_{a,b} f(a,b) = \min_a \min_b f(a,b)$$

1. minimize w.r.t Θ
2. minimize w.r.t $\theta^{(1)}, \dots, \theta^{(z)}$

Loss function

Learning population-level parameter

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

$\Lambda_{\Theta} \rightarrow$ (A) bounded elements, (B) sparse rows (C) zero diagonals

Λ_{Θ} places independent constraints on the rows of Θ

p independent convex optimization problems

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

Loss function

Learning unit-level parameter

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

$\Lambda_\theta \rightarrow$ (A) bounded elements, (B) 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 convex optimization problems

$$\mathcal{L}^{(i)}(\theta^{(i)}) = \sum_{t \in [\text{dim}]} \exp\left(-[\theta_t^{(i)} + 2\hat{\Theta}_t^\top \mathbf{x}^{(i)}]x_t^{(i)}\right) \text{ for all } t \in [\text{dim}] \longrightarrow \hat{\theta}^{(i)}$$

Can we do counterfactual inference in the presence of unobserved confounding with one sample?



**For every user, Netflix can estimate
the expected potential outcomes
with MSE scaling as $1/\text{dimension}$**

Exponential family to the rescue!



Social network setting

