# On Counterfactual Inference with Unobserved Confounding

Abhin Shah

MIT



Raaz Dwivedi

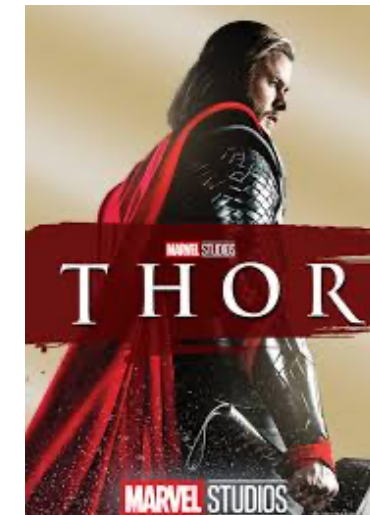Harvard & MIT

Devavrat Shah
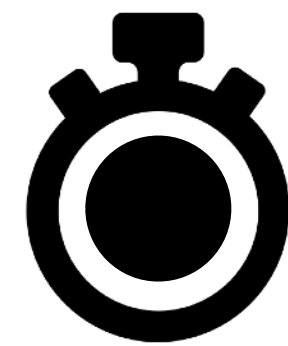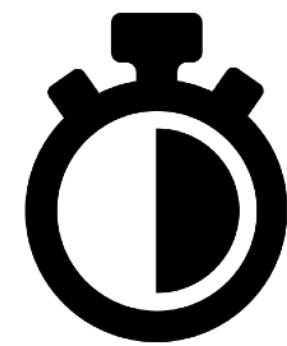
MIT

Greg Wornell

MIT

# Observational Setting



$\mathbf{a}$ – action/intervention

$\mathbf{y}$ – outcome

# **Panel Data**



**Potential Outcomes**

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

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

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

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

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

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

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

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

# Goal : What-if?



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

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

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

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

# Challenges

1. unobserved factors → spurious associations

2. users → heterogeneous

3. each user → a single interaction trajectory



engagement level

humor

Simpson's paradox

engagement level

humor

age

10

20

30

40

50

# Problem Setup



unobserved covariates    z — v    observed covariates

observed interventions    a → y    observed outcomes

$n$ heterogenous and independent users

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

# Examples

## Sequential setting

# Examples

## Network setting

# Goal : Counterfactual Questions

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

$$\equiv$$

$$\boxed{\text{Estimate } \boldsymbol{y}^{(i)}(\tilde{\boldsymbol{a}}^{(i)}) \text{ for } \tilde{\boldsymbol{a}}^{(i)} \in \mathscr{A}?}$$

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

**Heterogeneity:** each user may have different *unobserved* z

$$\boxed{\begin{array}{c} \text{Can we learn } n \text{ different distributions} \\ p(\mathbf{y} = \cdot \mid \mathbf{a} = \cdot, \boldsymbol{z}^{(i)}, \boldsymbol{v}^{(i)}) \quad i \in [n] \\ \text{with } one \text{ sample per distribution?} \end{array}}$$

# Related Work

## Linear panel data

- Examples: difference in differences, synthetic control + variants, synthetic interventions + variants

- Only finitely many interventions

- Special structure on intervention assignment

# Our Approach

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

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

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

different for different users

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

# Our Approach

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

$\Theta$

$\theta(z^{(1)}), \Theta$

$\vdots$

$\theta(z^{(n)}), \Theta$

$\phi, \Phi \xrightarrow{\text{Reparameterize}} \theta(z), \Theta \xrightarrow{n \text{ units}} \xrightarrow{\text{shared parameters}} \theta(z^{(1)}), \cdots, \theta(z^{(n)}), \Theta$

1.  If $z^{(1)} = \cdots = z^{(n)} \to$ a single exponential family with $n$ samples **[B '15, KM '17, SSW '21, VML '22]**

2.  If $n = 1 \to$ a single exponential family with one sample (assume $\Theta$ is known) **[KDDGD '21, MHBM' 21]**

# Inference Tasks

1. Parameters

   A. User-level $-\theta^{\star}(z^{(i)})$ for all $i \in [n]$ $\quad\longrightarrow$ <span style="color:blue">counterfactual distribution</span>

   B. Population-level $-\Theta^{\star}$

2. Expected potential outcomes $-\mathbf{E}\left[y^{(i)}(\tilde{a}^{(i)}) \mid \mathbf{z} = z^{(i)}, \mathbf{v} = v^{(i)}\right]$ $\quad\longrightarrow$ <span style="color:blue">counterfactual mean</span>

# Complexity of Parameter Space

$\Lambda_\theta$ — 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 |
|---|---|---|
| $M(\epsilon) = logN(\Lambda_\theta, \epsilon)$ | $O\left(k \cdot log\left(1 + \dfrac{1}{\epsilon}\right)\right)$ | $O\left(s \log k \cdot log\left(1 + \dfrac{1}{\epsilon}\right)\right)$ |
| $M_n(\epsilon) = nM(n\epsilon)$ | $O\left(\dfrac{k}{\epsilon}\right)$ | $O\left(\dfrac{s \log k}{\epsilon}\right)$ |

# Parameter Estimation

$$\{\boldsymbol{v}^{(i)}, \boldsymbol{a}^{(i)}, \boldsymbol{y}^{(i)}\}_{i=1}^{n} \longrightarrow \qquad \longrightarrow \hat{\theta}(\boldsymbol{z}^{(1)}), \cdots, \hat{\theta}(\boldsymbol{z}^{(n)}), \hat{\Theta}$$

pool all $n$ samples $\qquad\qquad$ Loss function $\qquad\qquad$ estimates

$$\|\Theta^{\star} - \hat{\Theta}\|_{2,\infty} \le \epsilon \qquad \text{when } n \ge O\left(\frac{p^2\big(p + M_n(\epsilon^2)\big)}{\epsilon^4}\right)$$

$$\text{For all } i, \text{MSE}\Big(\theta^{\star}(\boldsymbol{z}^{(i)}), \hat{\theta}(\boldsymbol{z}^{(i)})\Big) \le \max\left\{\epsilon^2, \frac{M(c)}{p}\right\} \quad \text{when } n \ge O\left(\frac{p^2\big(pM(c) + M_n(\epsilon^2)\big)}{\epsilon^4}\right)$$

# Outcome Estimation

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

$(v^{(i)}, \tilde{a}^{(i)})$

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

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

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

# An Application

**Denoise user-wise data**

No systematically unobserved covariates $\mathbf{z}$

Noisy observed data = true data + measurement error

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

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

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

Goal: Estimate the true data

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

# 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 $p_\theta(\boldsymbol{x}) = \dfrac{\exp(\theta^\top f(\boldsymbol{x}))}{Z(\theta)}$

Given independent $\boldsymbol{x}_1, \cdots, \boldsymbol{x}_n \sim p_{\theta^\star}(\boldsymbol{x})$, learn $\theta^\star$

$\widehat{\theta}_{MLE} \triangleq \text{argmax}_\theta \text{ Log-Likehood}(\theta)$

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

A. ✓ Consistency ?

B. ✓ Asymptotic normality ?

C. ✓ Asymptotic efficiency ?

D. ✗ Computational tractability ? How to compute $Z(\theta)$?

# An Alternative

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

Given independent $\boldsymbol{x}_1, \cdots, \boldsymbol{x}_n \sim p_{\theta^\star}(\boldsymbol{x})$, learn $\theta^\star$

$$\widehat{\theta} \triangleq \operatorname{argmin}_\theta \ell(\theta)$$

likelihood

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

Convex constraint on $\theta$ $\longrightarrow$ convex optimization problem

avoids $Z(\theta)$

but is it any good?

# Properties

**[SSW'21]**

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

✓ A. Strictly proper loss function

✓ B. Consistency ?  $\qquad$  $\hat{\theta}$ is MLE w.r.t. $p_{\theta^{\star}-\theta}$ (not $p_{\theta}$)

✗ C. Asymptotic normality ?

✓ D. Asymptotic efficiency ?

✓ E. Computational tractability ?  $\qquad$  No need to compute $Z(\theta)$ !

✓ F. Finite sample guarantees ?

# Finite Sample Guarantees

## [SSW'21]

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

$$\|\theta^\star - \widehat{\theta}\|_2 \sim n^{-1/4}$$

Proof ingredients —

    A.  Concentration of gradient

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

} Hoeffding's inequality

# Back To Our Setting

**Condition on z**

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

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

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

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

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

$$p(x_t | \boldsymbol{x}_{-t}, \boldsymbol{z}) \propto \exp\Big( \big[\theta_t(\boldsymbol{z}) + 2\Theta_t^\top \boldsymbol{x}\big] x_t\Big)$$

# Structure On Parameters

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

(A) Every element of $\theta^\star(\boldsymbol{z}^{(1)}), \cdots, \theta^\star(\boldsymbol{z}^{(n)}),$ and $\Theta^\star$ are bounded

(B) Every row of $\Theta^\star$ is sparse

$$\Lambda_\theta \triangleq \{\theta : \theta \text{ is consistent with (A)} + \boxed{\text{low complexity}}\} \longrightarrow \text{required to provide meaningful guarantees}$$

$$\Lambda_\Theta \triangleq \{\Theta : \Theta \text{ is consistent with (A) and (B)}\}$$

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

# Learning Population-level Parameter

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

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

$\Lambda_\Theta$ places independent constraints on the rows of $\Theta$

$p$ independent optimization problems

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

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

A. Concentration of gradient

B. Anti-concentration of Hessian

Hoeffding's inequality

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

# Learning Unit-level Parameter

$$\mathscr{L}(\theta^{(1)}, \cdots, \theta^{(n)}) = \sum_{t \in [p]} \frac{1}{n} \sum_{i \in [n]} \exp\Big( -\big[ \theta_t^{(i)} + 2\widehat{\Theta}_t^\top \boldsymbol{x}^{(i)} \big] x_t^{(i)} \Big)$$

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

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

$n$ independent optimization problems

$$\mathscr{L}^{(i)}(\theta^{(i)}) = \sum_{t \in [p]} \exp\Big( -\big[ \theta_t^{(i)} + 2\widehat{\Theta}_t^\top \boldsymbol{x}^{(i)} \big] x_t^{(i)} \Big) \text{ for all } i \in [n] \quad \longrightarrow \quad \widehat{\theta}^{(i)}$$

$$\|\theta^\star(\boldsymbol{z}^{(i)}) - \widehat{\theta}^{(i)}\|_2 \sim \max\{n^{-1/4}, M\}$$

A.  Concentration of gradient  } Logarithmic
                                   Sobolev
B.  Anti-concentration of Hessian  inequality