

Causal Tensor Estimation

Devavrat Shah

Alberto Abadie Anish Agarwal Dennis Shen

Massachusetts Institute of Technology

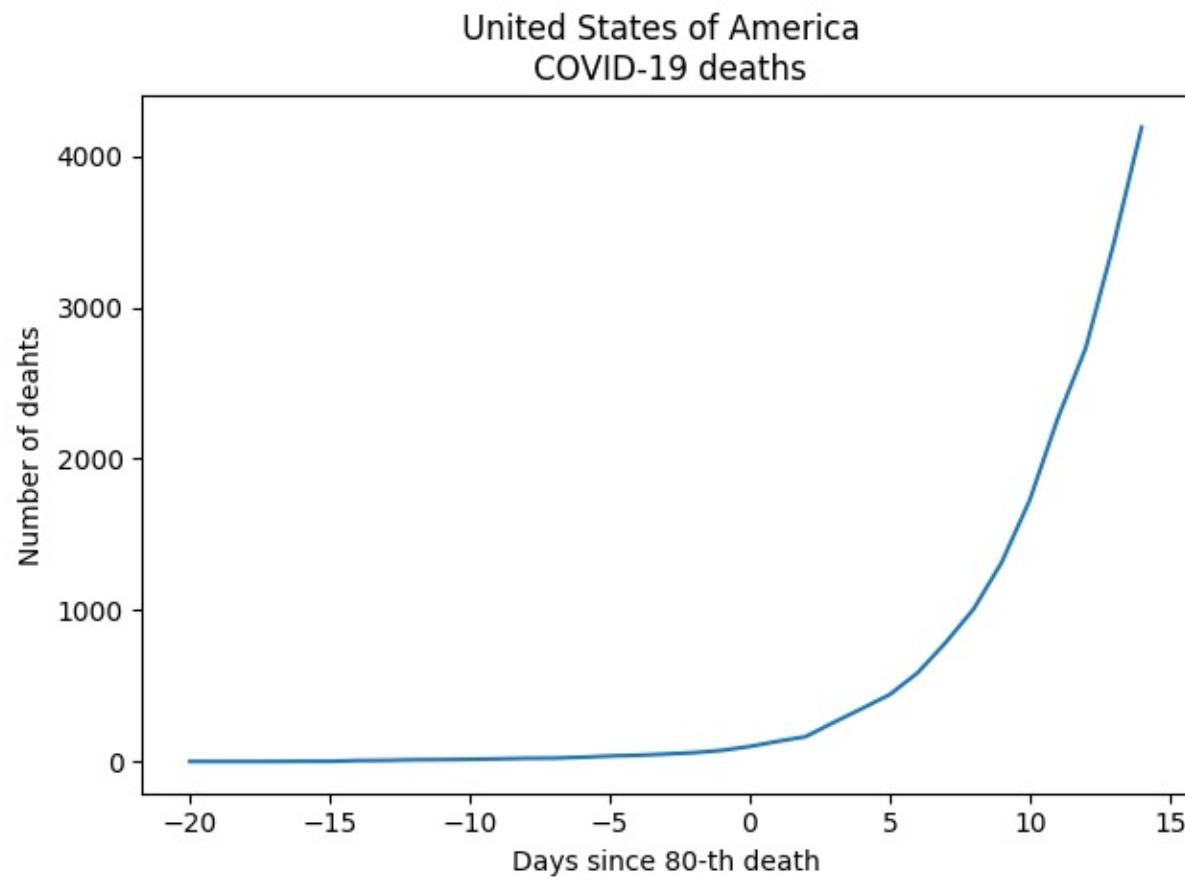
Synthetic Interventions: <https://arxiv.org/abs/2006.07691>

Causal Tensor Estimation: working paper

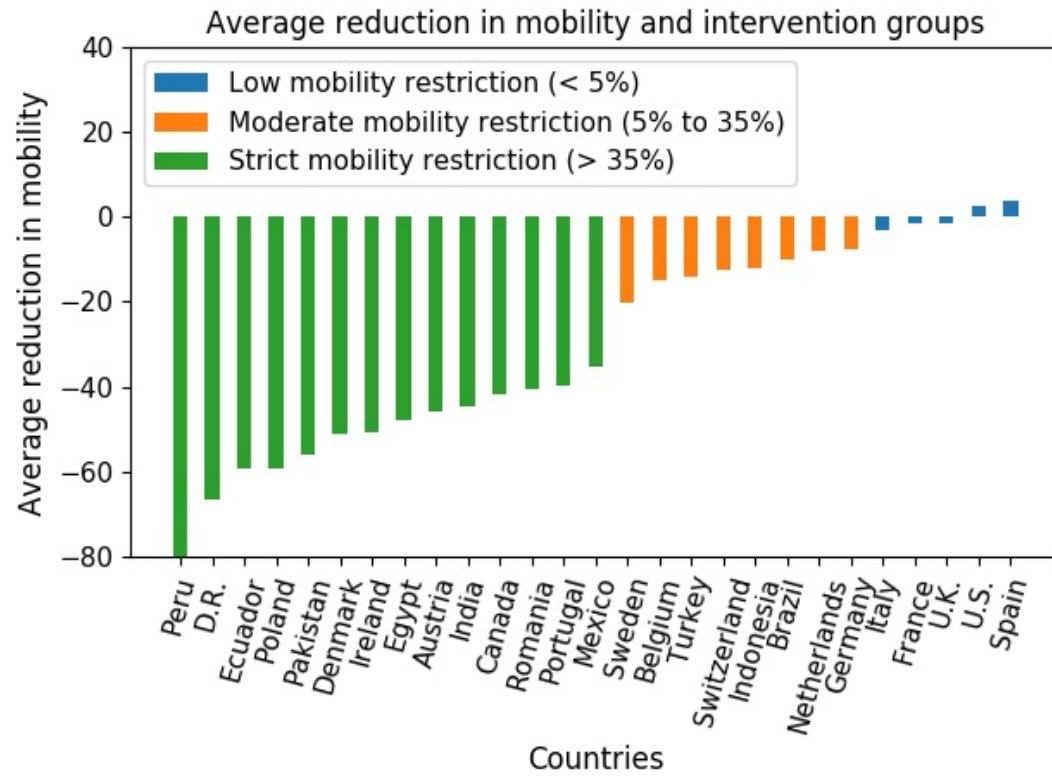
Policy Evaluation

Anish Agarwal Abdullah Alomar Arnab Sarkar Dennis Shen

United States



Looking Across the Globe



Low mobility restriction

< 5% reduction

Moderate mobility restriction

5 - 35% reduction

Severe mobility restriction

> 35% reduction

What would have happened to United States if...

- United States had experienced (through appropriate policy)
 - Low mobility restriction
 - < 5% reduction [the reality]
 - Moderate mobility restriction
 - 5-35% reduction
 - Severe mobility restriction
 - > 35% reduction
- In terms of
 - Trajectory of death counts in region of interest

It's Causal Inference

Potential Outcomes Framework [Newman '23, Rubin '74]

An individual contains many latent selves



Low



Moderate



High

Y = observed outcome (health outcome)

$M^{(d)}$ = potential outcome under intervention d (health outcome under policy d)

d^* = observed intervention (low restrictions), that is $Y \xrightarrow{\mathbb{E}} M^{d^*}$

Goal: estimate $M^{(d)}$, $d \neq d^*$

Fundamental Question

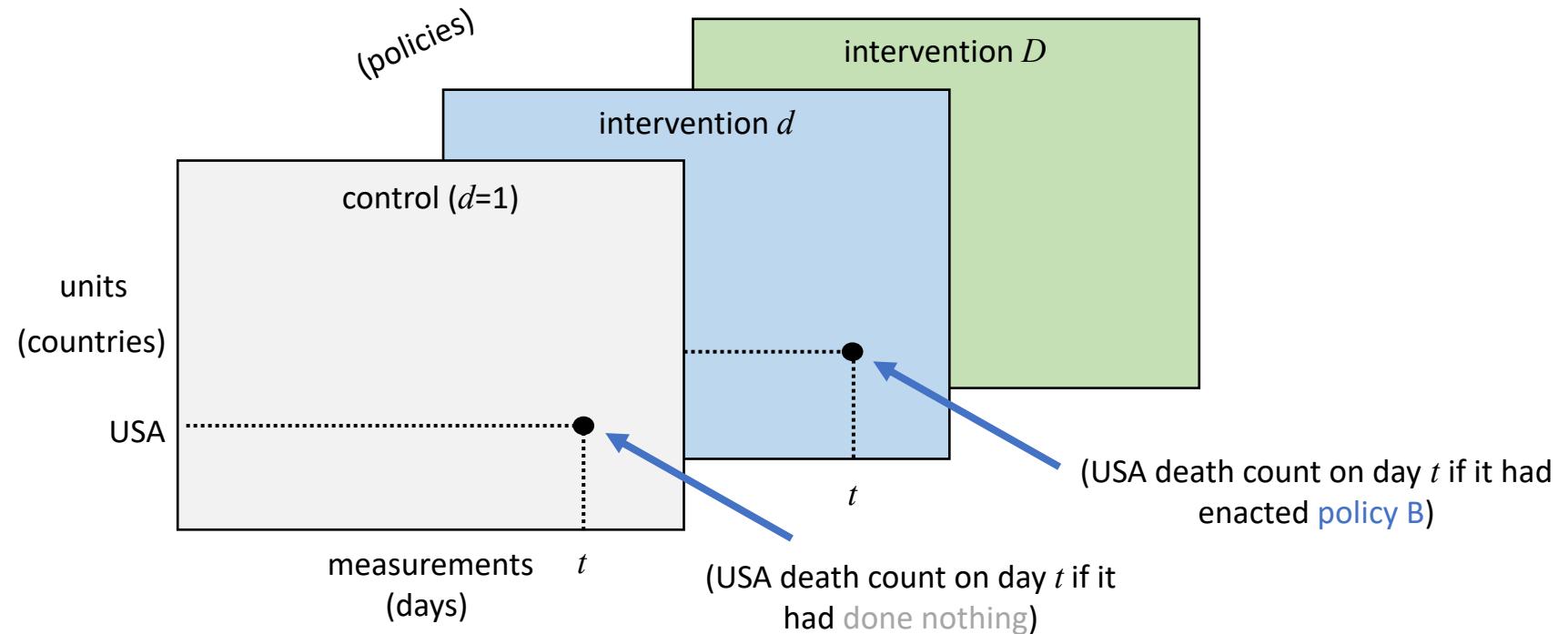
Only one outcome can be revealed
But want to know *all possible* outcomes

Let's Look At An Alternative Representation: Tensor

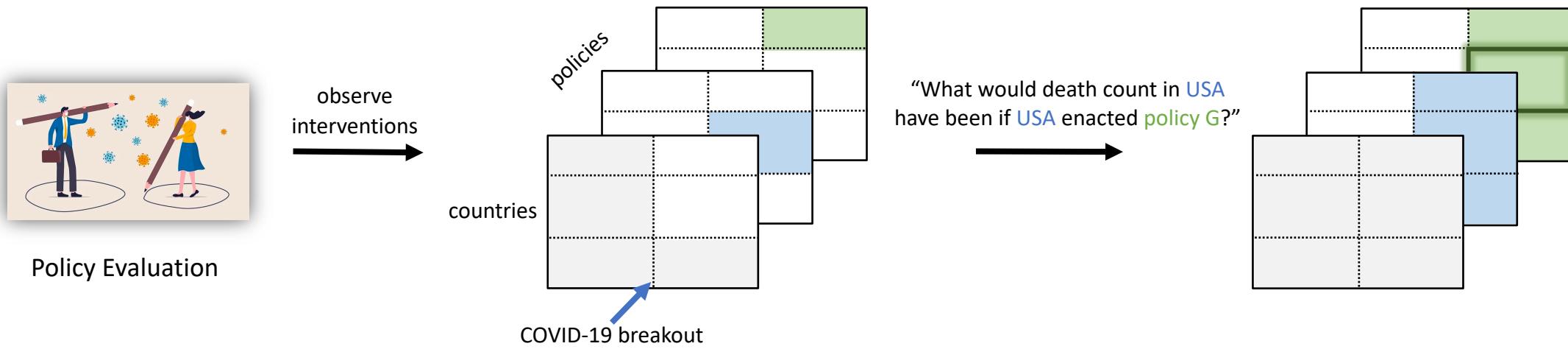
Encode potential outcomes into a [Tensor](#)

N units, T measurements, D interventions

(n, t, d)



Causal Inference = Causal Tensor Estimation



- Causal Tensor Estimation
 - “Imputing” missing values in a Tensor
 - Potentially “confounded” observations (e.g. not missing at random)
 - The policy implemented in a country depends on the “characteristics” of the country!

What is Confounding, Why Is it a Problem

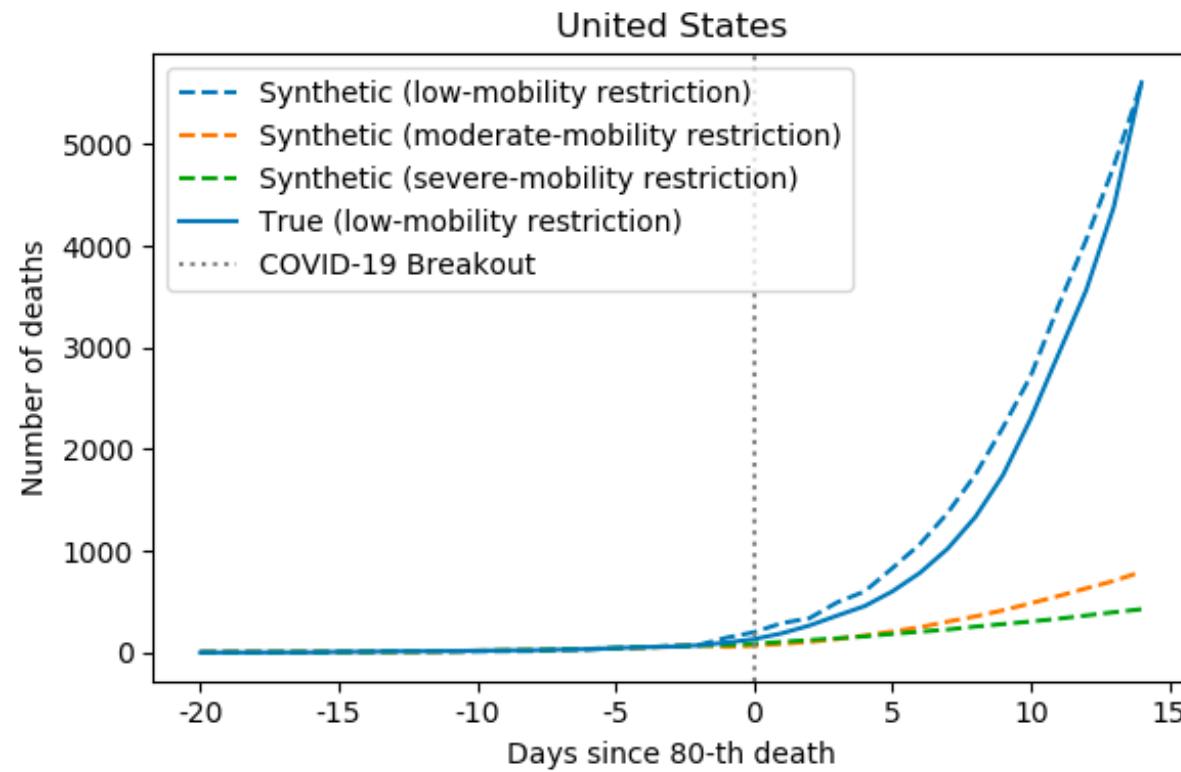
- To determine A vs B:
 - Access to 100 M + 100 W patients
- Randomized Trial
 - 50 M, 50 W receive A (similarly B)
 - Average efficacy: **5 for A** and **10 for B**
 - Conclusion: B is better than A
- Observational data (“confounded” selection)
 - 100 M get A, 100 W get B
 - Average efficacy: **10 for A** and **0 for B**
 - Conclusion: A is better than B

Ground Truth or Potential Outcome

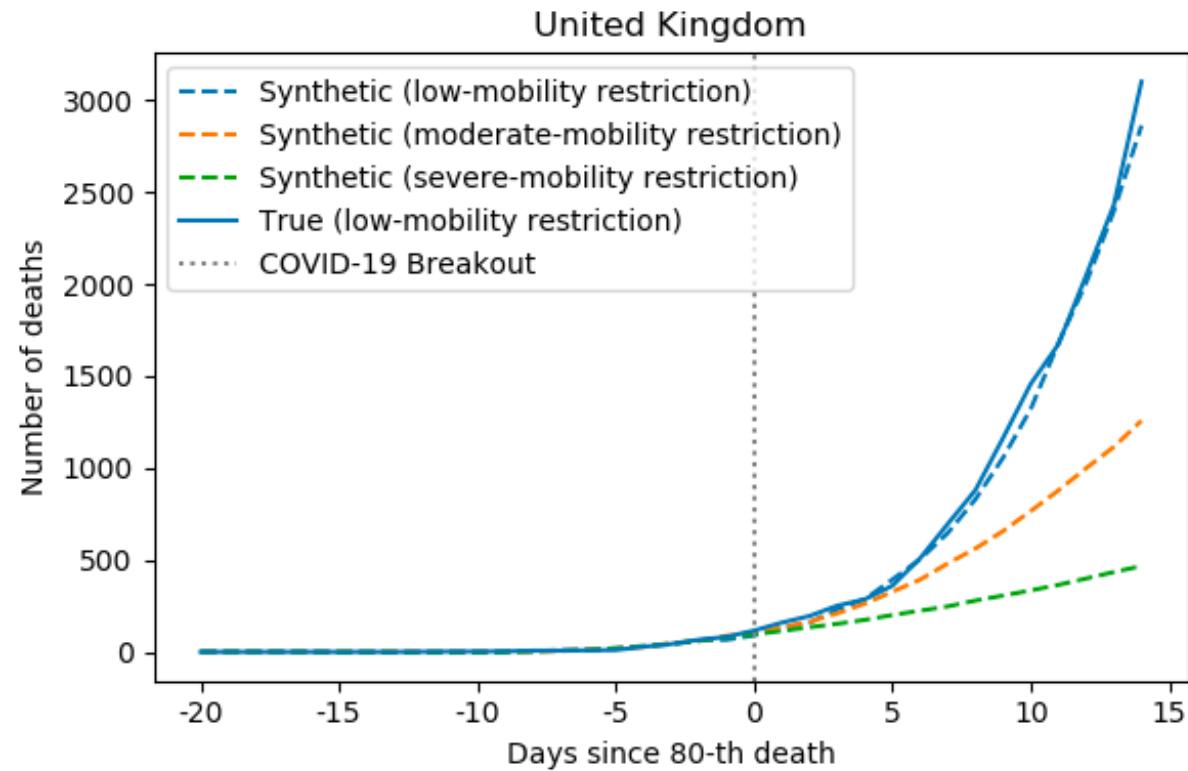
Efficacy of Two Drugs Across Gender

	Drug A	Drug B
Men	10	20
Women	0	0

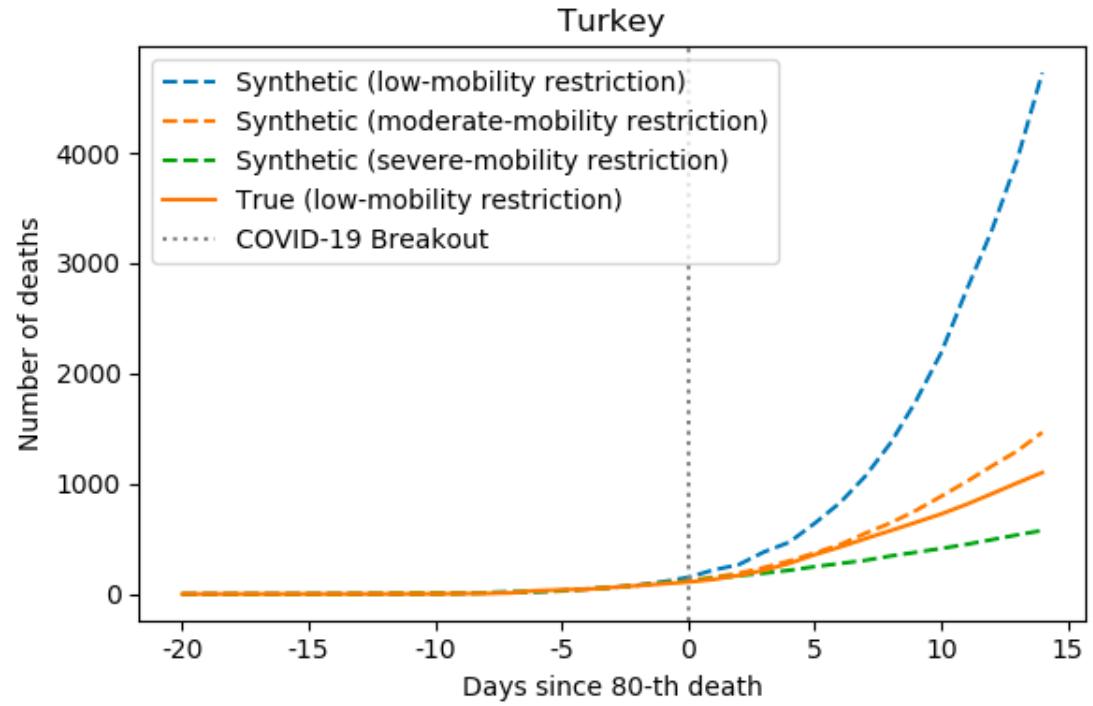
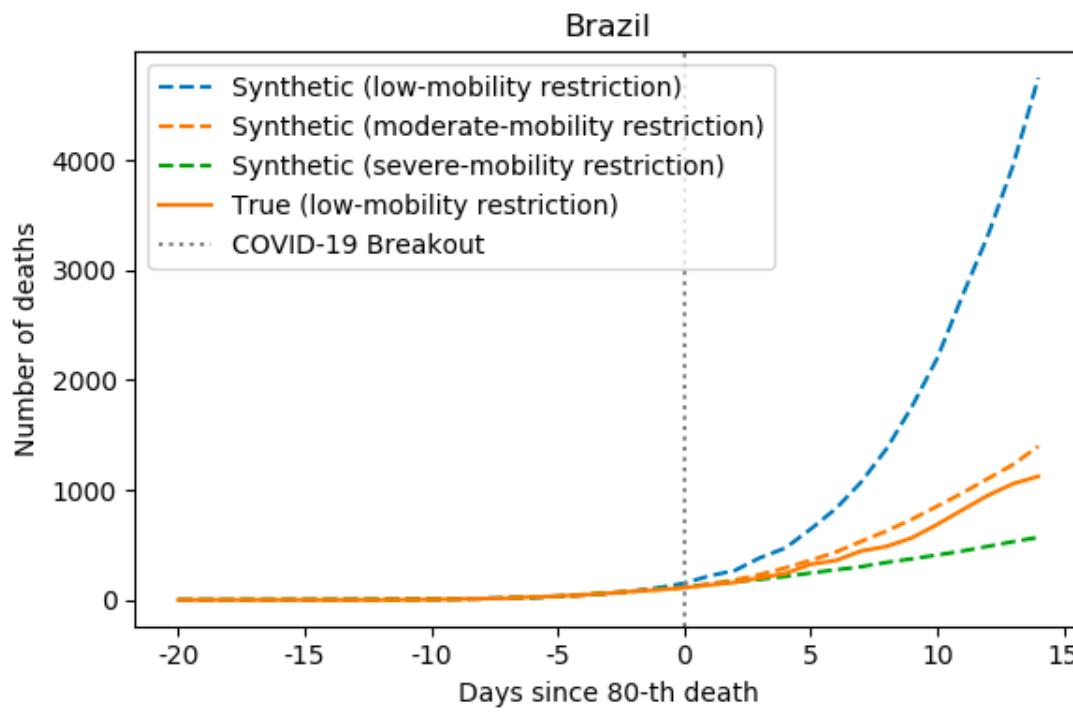
United States: Causal Tensor Estimation w Synthetic Interventions



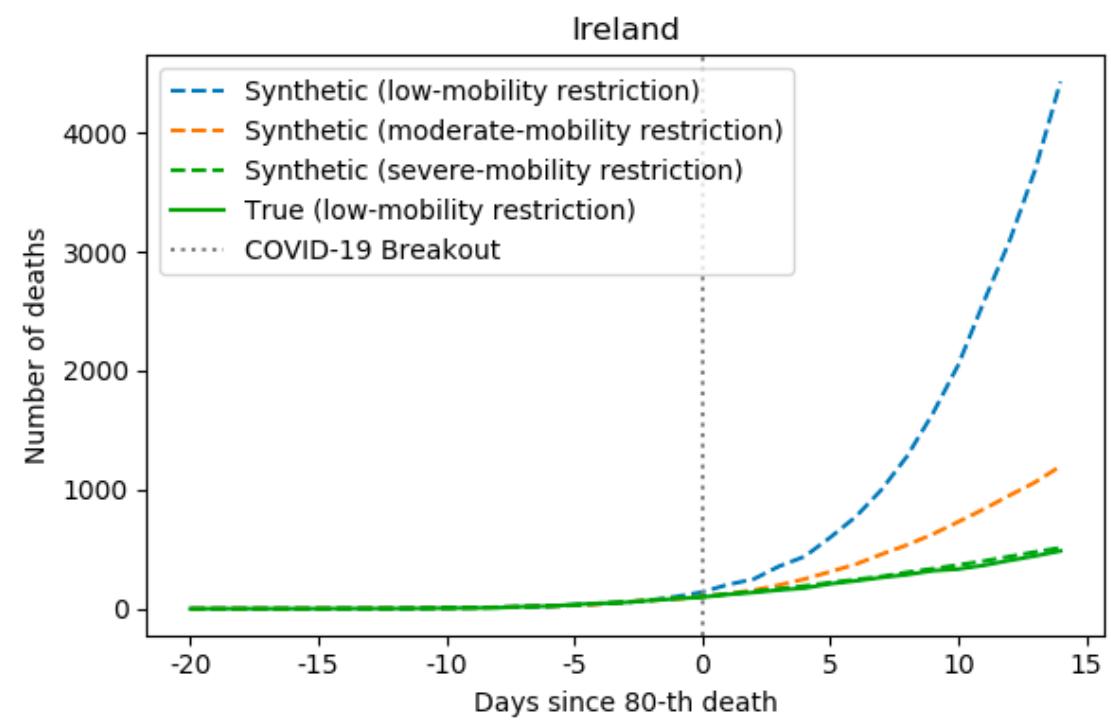
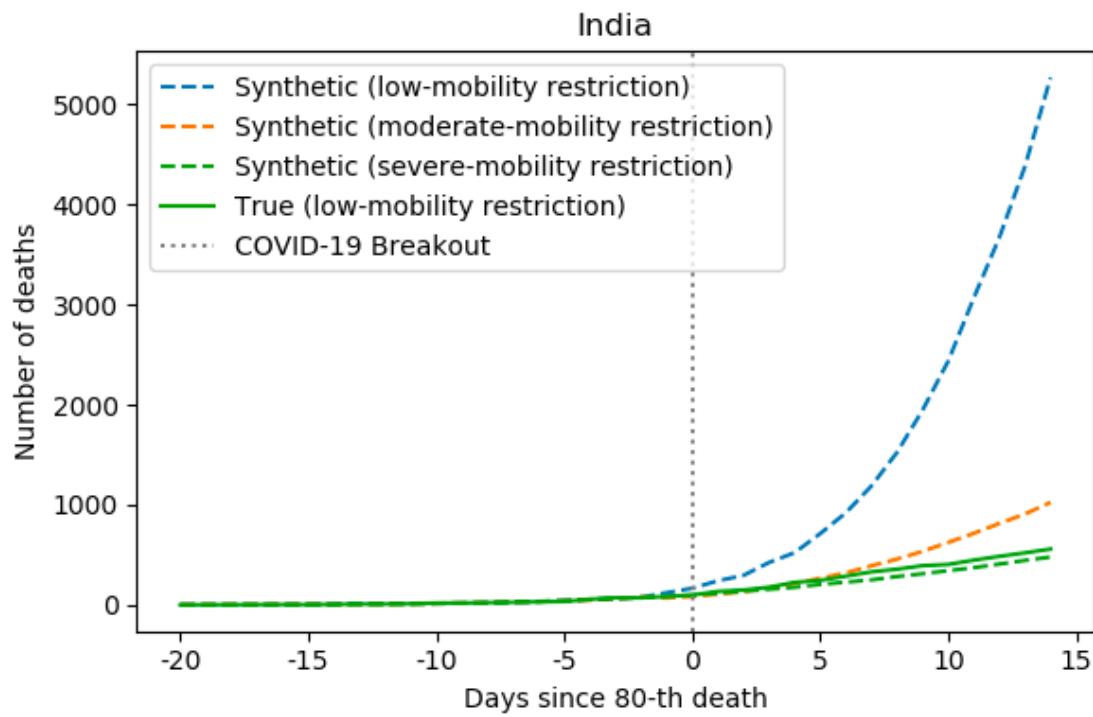
United Kingdom: Causal Tensor Estimation w Synthetic Interventions



Brazil, Turkey: Causal Tensor Estimation w Synthetic Interventions



India, Ireland: Causal Tensor Estimation w Synthetic Interventions



Data Efficient Randomized Control

Anish Agarwal Vishal Misra Dennis Shen

Clinical Trial For Personalized Treatment

6 patient types
(N=6)



3 drugs
(1 placebo) (D=4)



Intervention	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
placebo	✓	✓	✓	✓	✓	✓
drug 1	✓	✓	✓	✓	✓	✓
drug 2	✓	✓	✓	✓	✓	✓
drug 3	✓	✓	✓	✓	✓	✓

Clinical Trial For Personalized Treatment

Real clinical trial: $D \times N$

Intervention	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
placebo	✓	✓	✓	✓	✓	✓
drug 1	✓	✓	✓	✓	✓	✓
drug 2	✓	✓	✓	✓	✓	✓
drug 3	✓	✓	✓	✓	✓	✓

[the reality]

Data-efficient clinical trial: $2 \times N$

Intervention	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
placebo	✓	✓	✓	✓	✓	✓
drug 1	✓	✓				
drug 2			✓	✓		
drug 3					✓	✓

[our proposal]

Clinical Trial For Personalized Treatment = Tensor Estimation

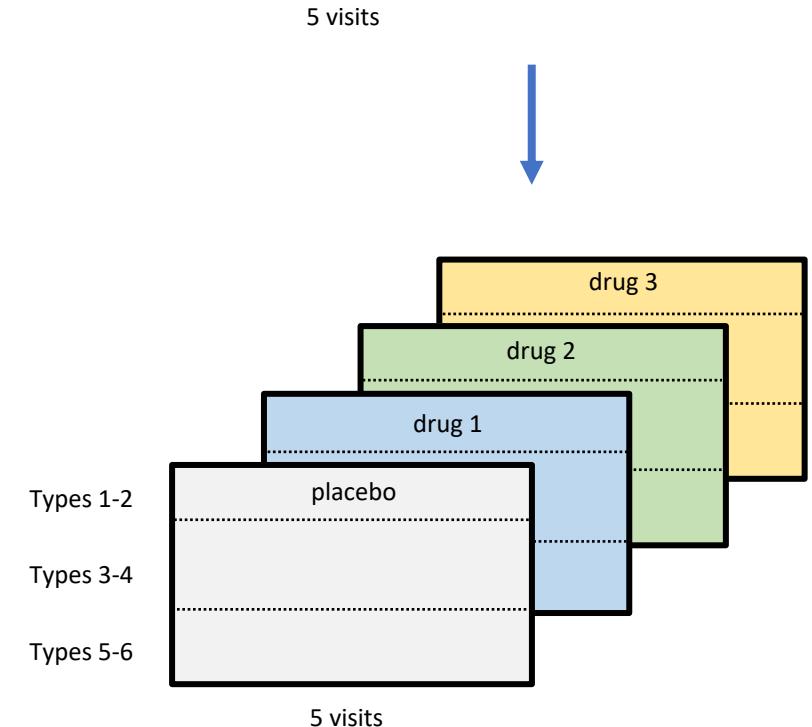
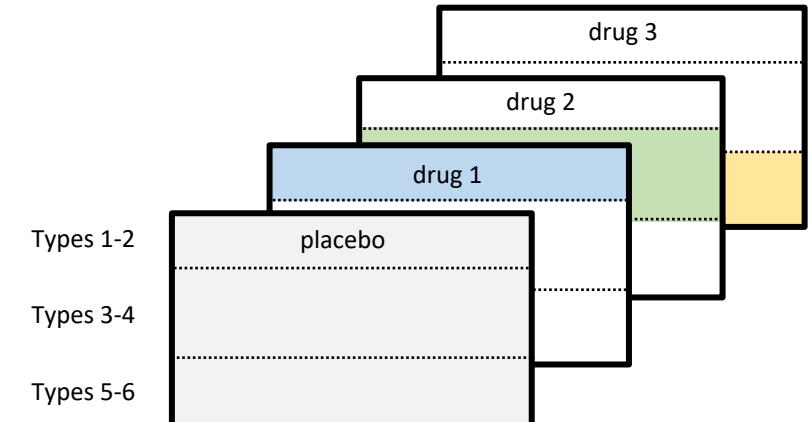
Real clinical trial: $D \times N$

Intervention	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
placebo	✓	✓	✓	✓	✓	✓
drug 1	✓	✓	✓	✓	✓	✓
drug 2	✓	✓	✓	✓	✓	✓
drug 3	✓	✓	✓	✓	✓	✓

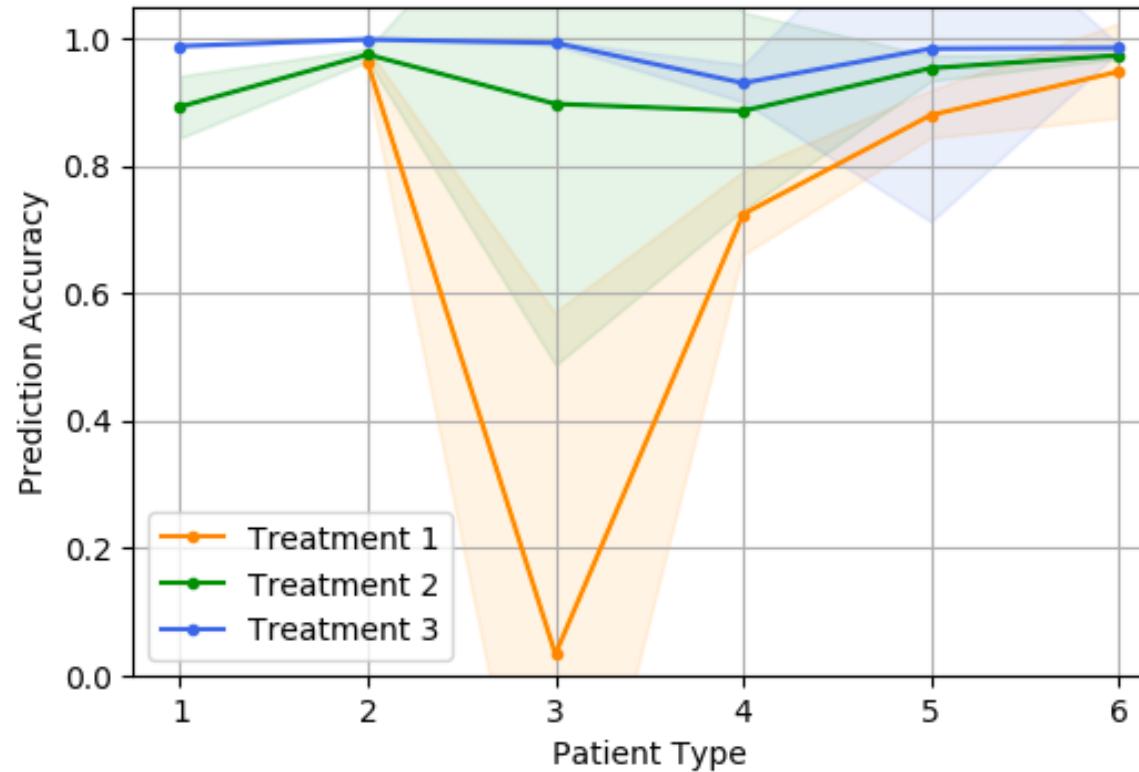
Data-efficient clinical trial: $2 \times N$

Intervention	Type 1	Type 2	Type 3	Type 4	Type 5	Type 6
placebo	✓	✓	✓	✓	✓	✓
drug 1	✓	✓				
drug 2			✓	✓		
drug 3					✓	✓

- Causal Tensor Estimation
 - Estimate outcomes for every (patient type, drug)
 - Using partial observations (no confounding)



Tensor Estimation Using Synthetic Interventions

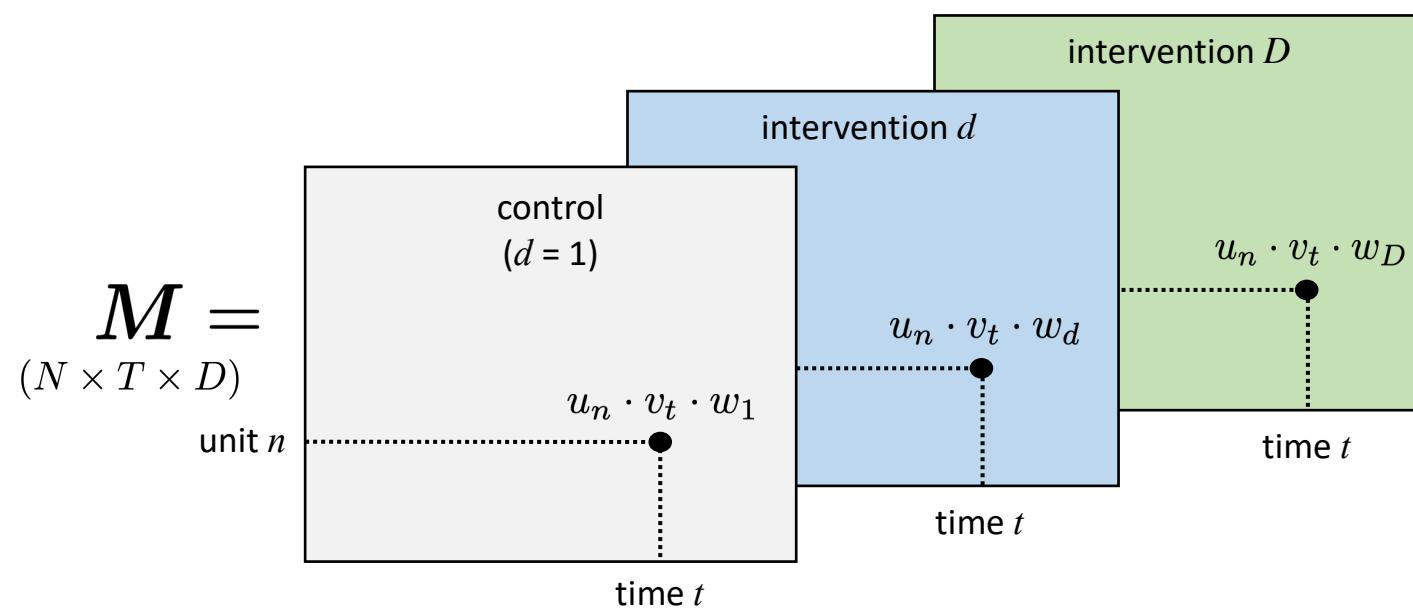


Accurately predicts outcome of 6×4 trials using only 6×2 trials

Framework: Causal Tensor Estimation

Alberto Abadie Anish Agarwal Dennis Shen

Potential Outcomes Tensor



Assumption (low-rank): r is small

$$= \sum_{\ell=1}^r u_\ell \otimes v_\ell \otimes w_\ell$$

intervention

unit

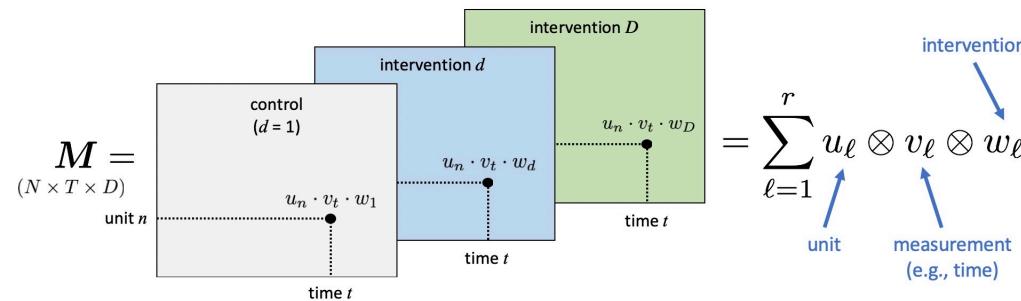
measurement
(e.g., time)

The Model

1. Sample (or given) latent unit, time, intervention factors

$$(u_n, v_t, w_d)$$

2. Sample potential outcomes tensor



3. Sample treatment assignment (determines sparsity pattern of observed tensor)

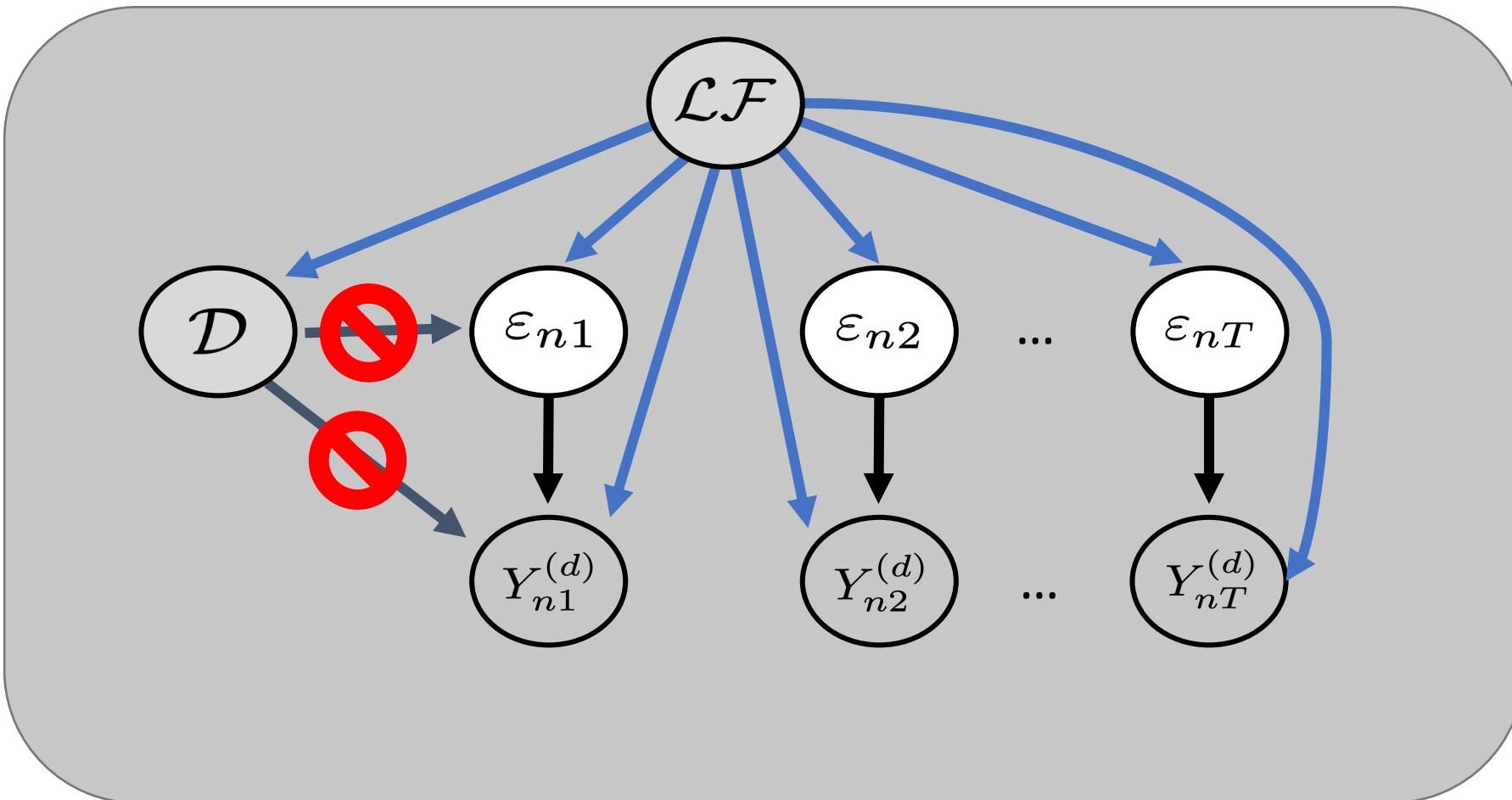
$$D(n, t) : [N] \times [T] \rightarrow 2^{[D]}$$

4. Observe noisy measurements: sampled entries of

$$Y = M + \varepsilon$$

What Type of Confounding is Allowed?

The joint distribution of latent factors (confounders, covariates), treatment assignment and observations satisfy the following Causal Structure



What Type of Confounding is Allowed?

- Recall

$$Y_{nt}^{(d)} = \sum_{\ell=1}^r u_{n\ell} \cdot v_{t\ell} \cdot w_{d\ell} + \varepsilon_{nt}$$

- Why is there confounding?

$$\mathcal{D} \not\perp\!\!\!\perp Y_{nt}^{(d)}$$

- Treatment assignments correlated with latent factors (i.e. unmeasured confounders)
- Selection on Latent Factors

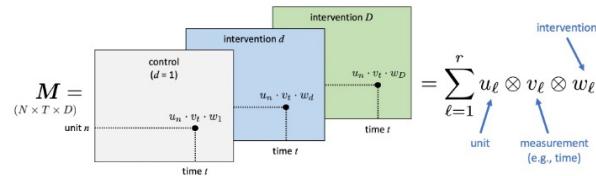
$$\mathcal{D} \perp\!\!\!\perp Y_{nt}^{(d)} \mid \mathcal{LF}$$

Causal Tensor Estimation

1. Sample (or given) latent unit, time, intervention factors

$$(u_n, v_t, w_d)$$

2. Sample potential outcomes tensor



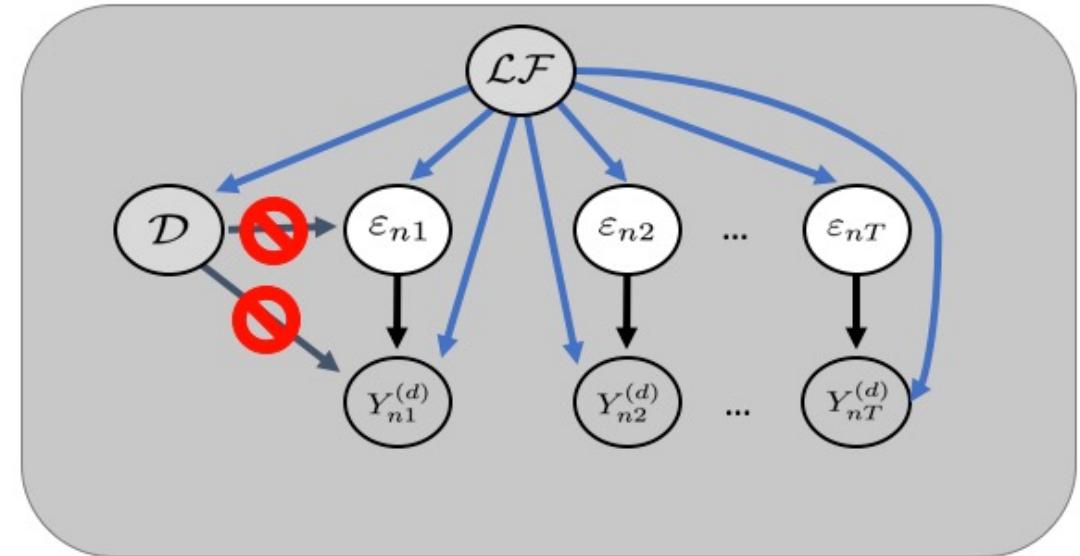
+

3. Sample treatment assignment (determines sparsity pattern of observed tensor)

$$D(n, t) : [N] \times [T] \rightarrow 2^{[D]}$$

4. Observe noisy measurements: sampled entries of

$$Y = M + \varepsilon$$



Produce an estimate \widehat{M} of M so that

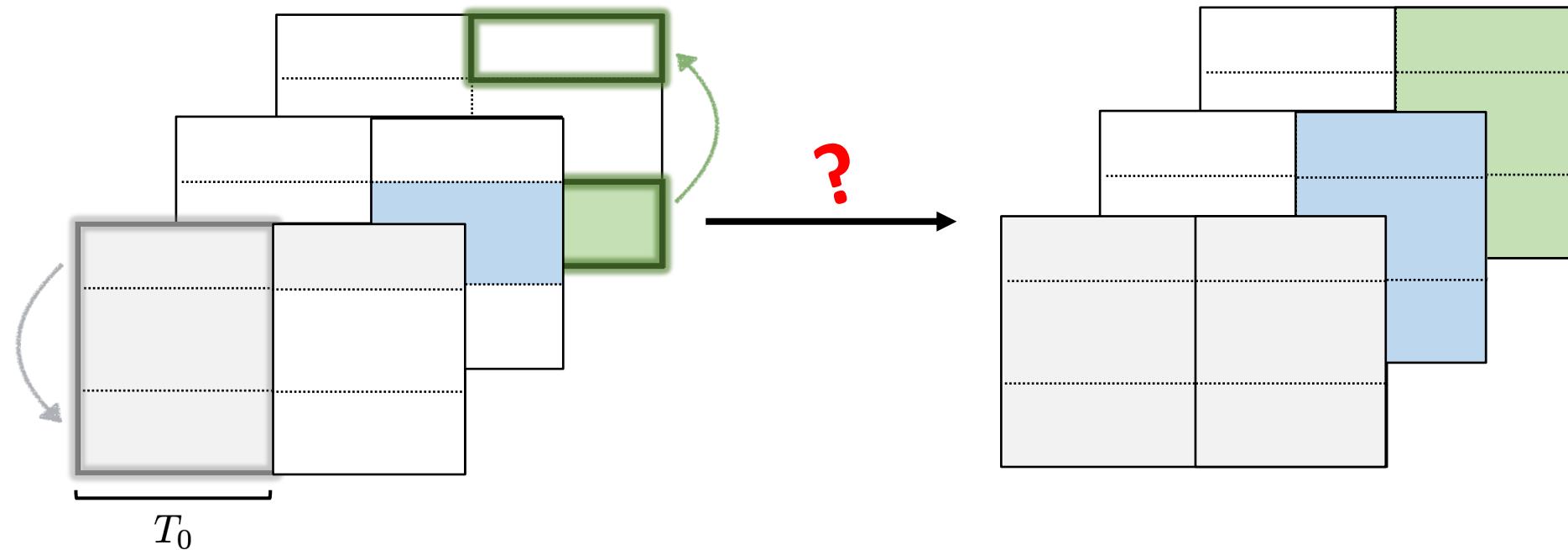
$$\widehat{M} \approx M$$

A Method: Synthetic Interventions

Anish Agarwal Dennis Shen

Key Insight

leverage data from other units
learn relationships between units

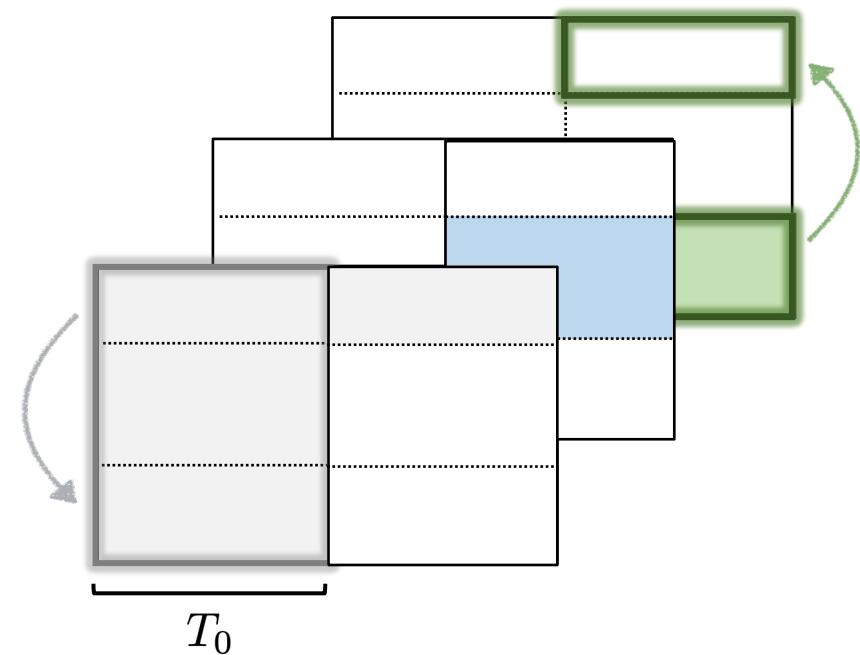


T_0
(# measurements under
common intervention)

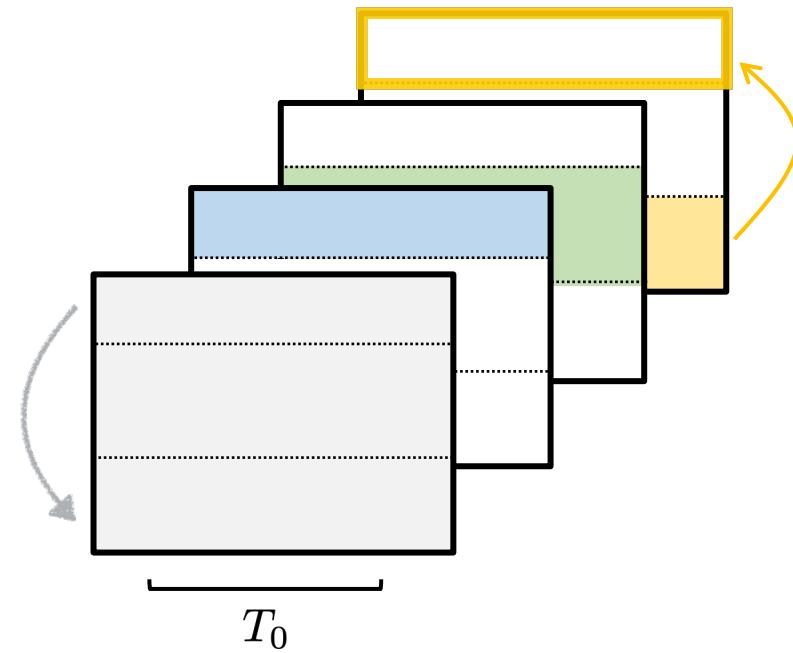
require all units to be under same intervention for
some number of measurements

Key Insight

leverage data from other units
learn relationships between units



T_0
(# measurements under
common intervention)



T_0
(# measurements under
common intervention)

require all units to be under same intervention for
some number of measurements

Synthetic Control (SC)

Estimates counterfactuals in absence of intervention

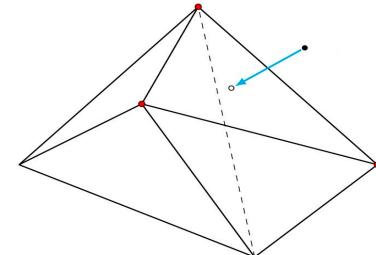
1. Learn Model
under control

$$\hat{\beta} = \underset{w \in \mathcal{W}}{\operatorname{argmin}} \|y - Z_{\text{pre}}w\|_2^2$$

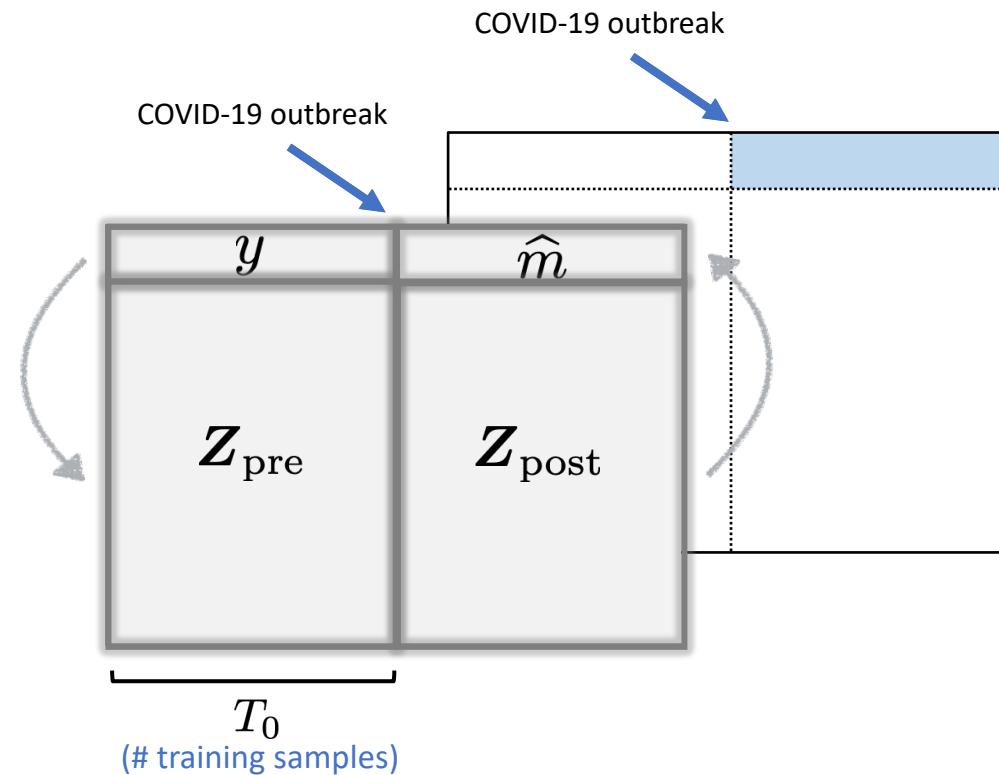
“synthetic” control

2. Predict
under control

$$\hat{m} = Z_{\text{post}}\hat{\beta}$$



“What would US’s death count have been if
US did nothing? ”



“And Other Countries did nothing? ”

Synthetic Interventions (SI)

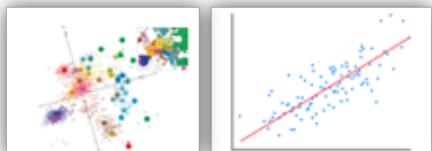
Estimates counterfactuals in absence and presence of interventions

"What would US's death count have been if it enacted de-mobility restricting policies?"

1. Learn Model
under control

$$\hat{\beta} = \underset{w \in \mathcal{W}}{\operatorname{argmin}} \|y - Z_{\text{pre}} w\|_2^2$$

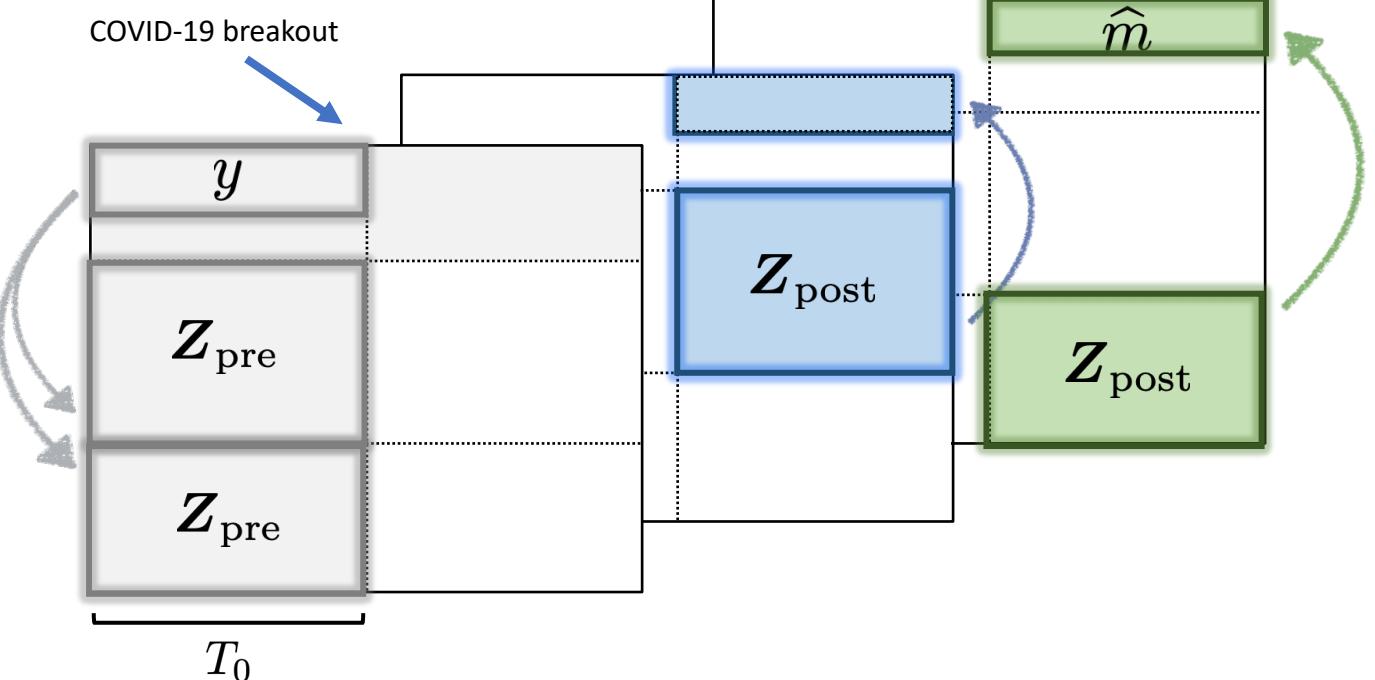
"synthetic" intervention



2. Predict
under intervention

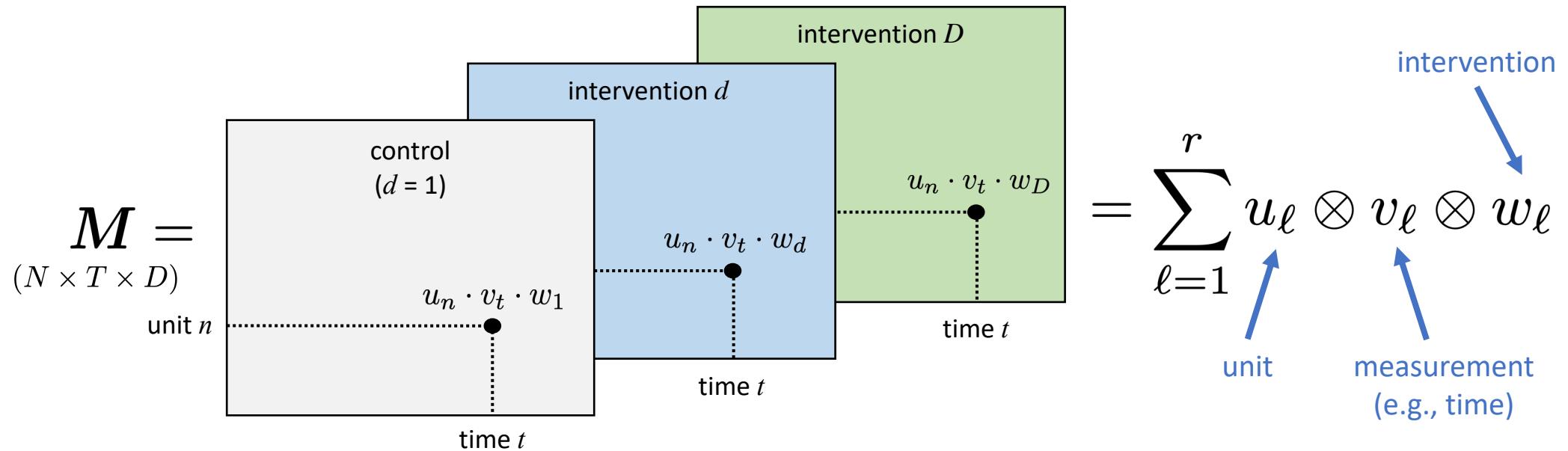
$$\hat{m} = Z_{\text{post}} \hat{\beta}$$

(de-noise via matrix estimation)



Why can we **transfer** learned model between
different interventional frameworks?

Why does SI Work?



Under intervention d

$$M^{(d)} = \sum_{\ell=1}^r u_\ell \otimes (w_{d\ell} \cdot v_\ell) = U(V^{(d)})^T$$

- U describes an **invariant** relationship between units across interventions
- **Each** intervention d is a linear transformation of U
- **SI** learns linear relationship between rows of U

Why does SI Work?

(WLOG) suppose unit 1 satisfies:

$$u_{1\ell} = \sum_{n>1} \beta_n^* \cdot u_{n\ell}$$

(occurs w.h.p.)

(low rank = few canonical unit profiles)

$$M_{1t}^{(d)} = \sum_{\ell=1}^r u_{1\ell} \cdot v_{t\ell} \cdot w_{d\ell} \quad \text{for any } (t,d)$$

(via tensor factor model)

$$= \sum_{\ell=1}^r \sum_{n>1} \beta_n^* \cdot u_{n\ell} \cdot v_{t\ell} \cdot w_{d\ell}$$

(via assumption)

$$= \sum_{n>1} \beta_n^* \cdot M_{nt}^{(d)}$$

invariant across
time, interventions

SI (and thus SC) exists

Identification, Consistency

$\theta_n^{(d)}$ = individual potential outcome under every intervention
averaged over post-intervention period

$$= \frac{1}{T_1} \sum_{t>T_0} \mathbb{E}[Y_{nt}^{(d)} \mid \{u_t^{(d)}, v_n : t > T_0\}]$$

$$\widehat{\theta}_n^{(d)} - \theta_n^{(d)} = \mathcal{O}_p \left(\frac{1}{T_0^{1/4}} + \frac{\|\tilde{w}^{(n,d)}\|_2}{\sqrt{T_1}} + \frac{\|\tilde{w}^{(n,d)}\|_1}{\min\{\sqrt{T_0}, \sqrt{N_d}\}} \right)$$

Normality

$$\sqrt{T_1}(\widehat{\theta}_n^{(d)} - \theta_n^{(d)}) \xrightarrow{d} \mathcal{N}\left(0, \sigma^2 \text{plim} \|\tilde{w}^{(n,d)}\|_2^2\right)$$

95% confidence interval

$$\theta_n^{(d)} \in \left[\widehat{\theta}_n^{(d)} \pm \frac{1.96 \widehat{\sigma} \|\widehat{w}^{(n,d)}\|_2}{\sqrt{T_1}} \right]$$

Computable quantities
(with provable guarantees)

Subspace Inclusion: Hypothesis Test

$$H_0 : \text{span}(\mathbf{V}_{\text{post}}) \subseteq \text{span}(\mathbf{V}_{\text{pre}})$$

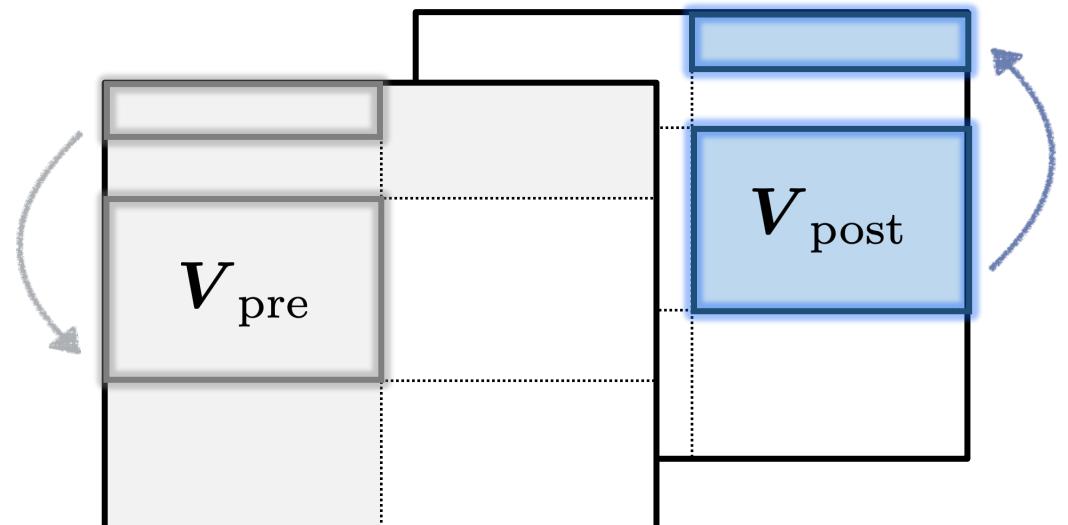
$$H_1 : \text{span}(\mathbf{V}_{\text{post}}) \not\subseteq \text{span}(\mathbf{V}_{\text{pre}})$$

If H_0 holds:

$$\|(\mathbf{I} - \mathbf{V}_{\text{pre}}\mathbf{V}_{\text{pre}}^T)\mathbf{V}_{\text{post}}\|_F^2 = 0$$

If H_1 holds:

$$\|(\mathbf{I} - \mathbf{V}_{\text{pre}}\mathbf{V}_{\text{pre}}^T)\mathbf{V}_{\text{post}}\|_F^2 > 0$$



Subspace Inclusion: Hypothesis Test

$$H_0 : \text{span}(\mathbf{V}_{\text{post}}) \subseteq \text{span}(\mathbf{V}_{\text{pre}})$$

$$H_1 : \text{span}(\mathbf{V}_{\text{post}}) \not\subseteq \text{span}(\mathbf{V}_{\text{pre}})$$

Test statistic

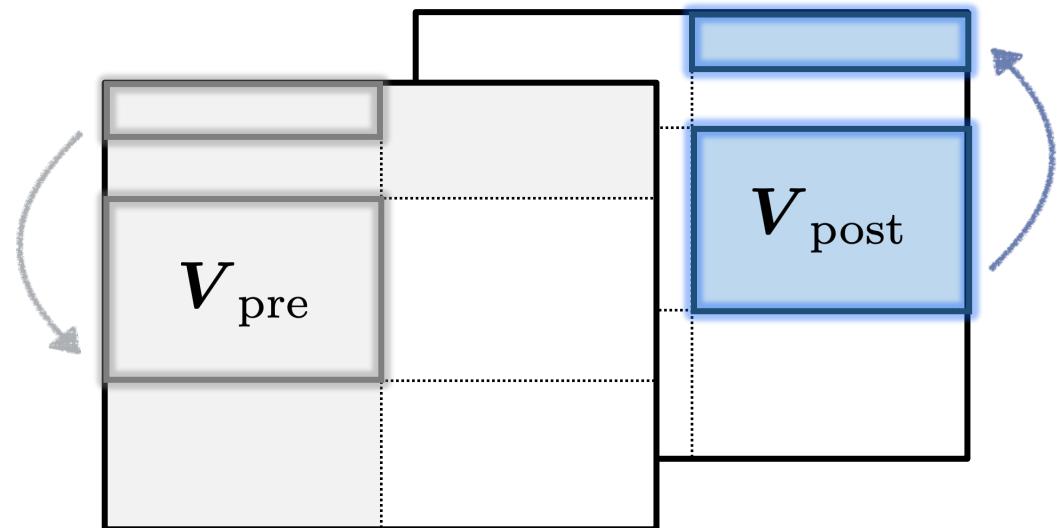
$$\hat{\tau} = \|(\mathbf{I} - \hat{\mathbf{V}}_{\text{pre}} \hat{\mathbf{V}}_{\text{pre}}^T) \hat{\mathbf{V}}_{\text{post}}\|_F^2$$

Test

For any significance level $\alpha \in (0, 1)$

Retain H_0 if $\hat{\tau} \leq \tau(\alpha)$

Reject H_0 if $\hat{\tau} > \tau(\alpha)$



Subspace Inclusion: Type I & Type II Error Guarantees

Fix any $\alpha \in (0, 1)$

Type I error:

$$\mathbb{P}(\hat{\tau} > \tau(\alpha) | H_0) \leq \alpha$$

Type II error:

$$\mathbb{P}(\hat{\tau} \leq \tau(\alpha) | H_1) \leq \alpha$$

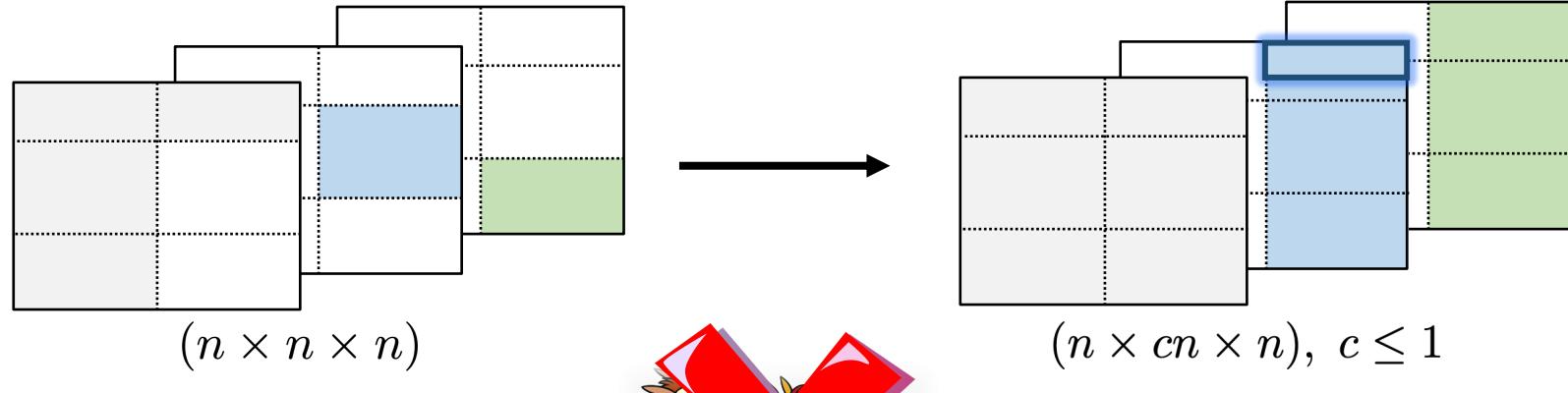
$$\text{where } \tau(\alpha) = \mathcal{O} \left(\frac{\sqrt{\log(1/\alpha)}}{\min\{\sqrt{T_0}, \sqrt{T_1}, \sqrt{N_d}\}} \right)$$

Parting Remarks

Statistical & Computational Tradeoffs in Causal Inference

Block Sparsity

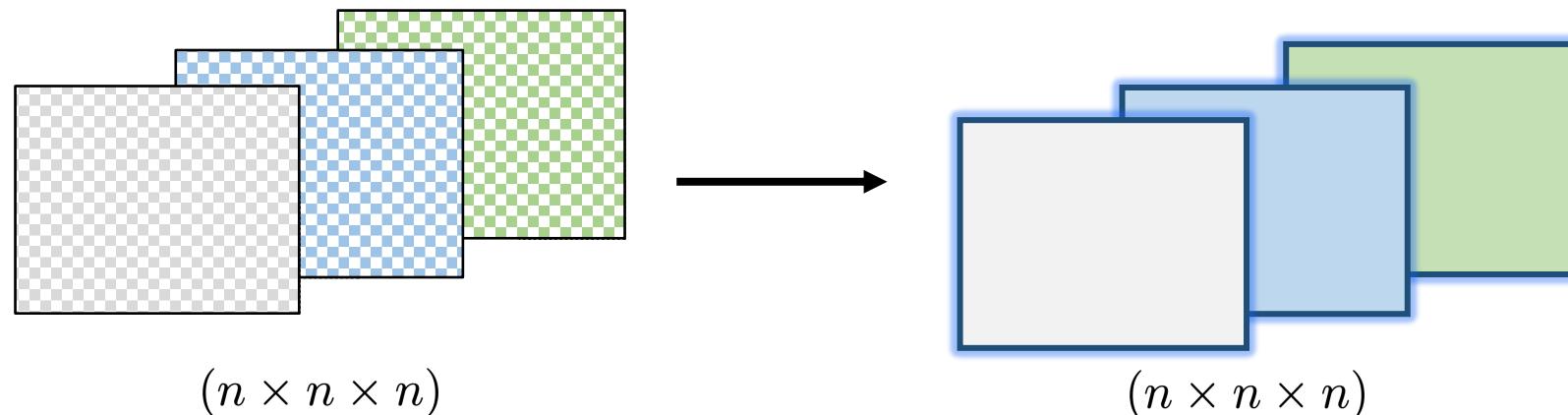
- Sample complexity:
 $\Omega(\sqrt{r} \cdot n^{3/2})$
- Computational complexity:
 $\mathcal{O}(n^3)$
(PCR)



Start discussion about computational/statistical trade-offs
New experimental designs via sampling?

Uniform Sparsity

- Sample complexity:
 $\Omega(r \cdot n^{3/2})$
- Computational complexity:
 $\mathcal{O}(\text{poly}(n))$



Causal Tensor Estimation: A Generic Framework

- Enables novel estimation
 - Regression discontinuity design in the panel data setting
- Experiment design
 - Observational pattern in tensor to enable identification
- Computational and statistical tradeoff
 - A missing discussion in Causal inference
- Role of error metric for tensor estimation
 - What causal quantities can be identified (or not)
- Causal estimation methods
 - SI is one such method, but more is needed

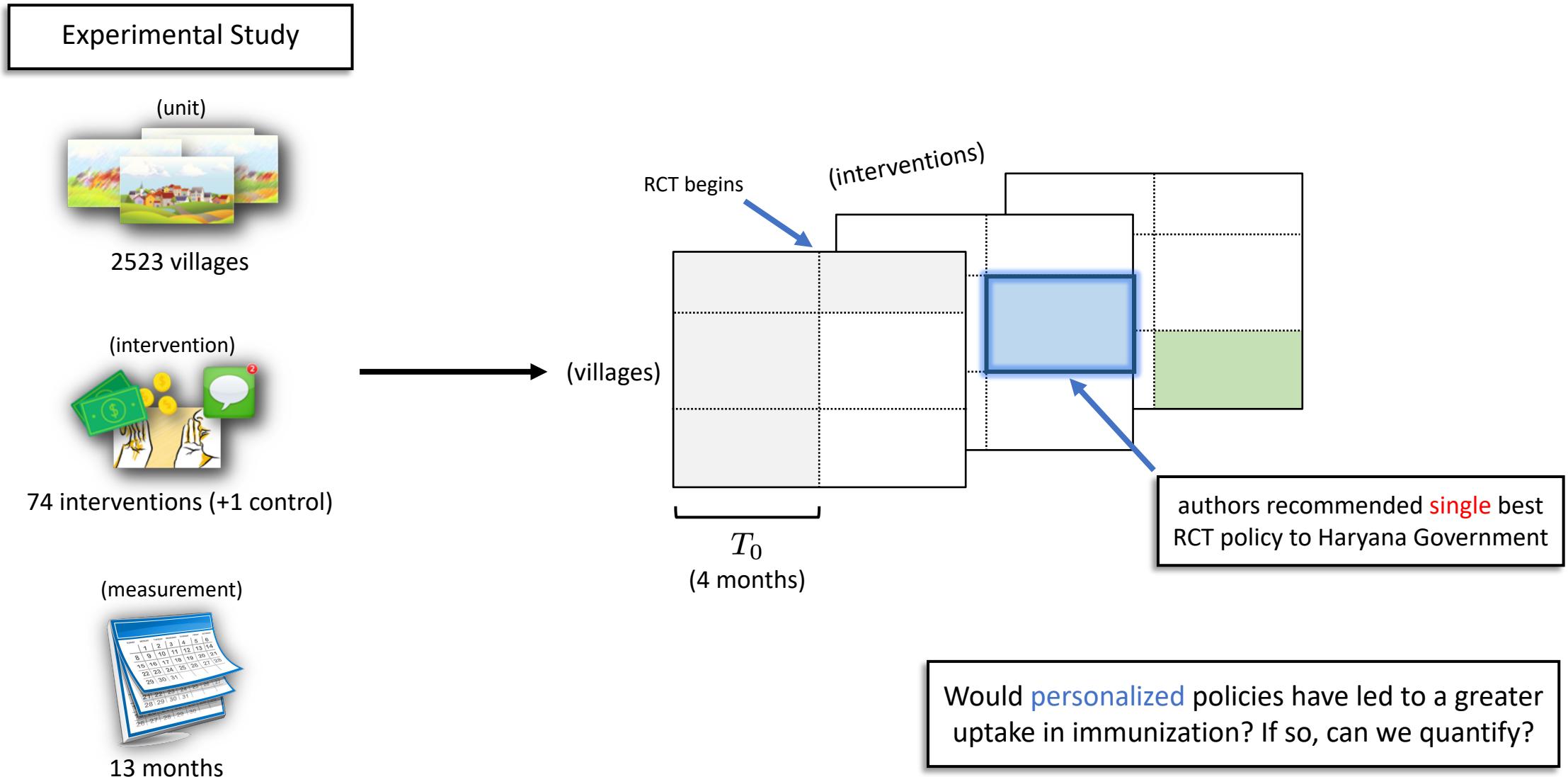
Questions

+ please feel free to contact at:
devavrat@mit.edu

Appendix

Development Economics

Development Economics [Banerjee et al 2019]



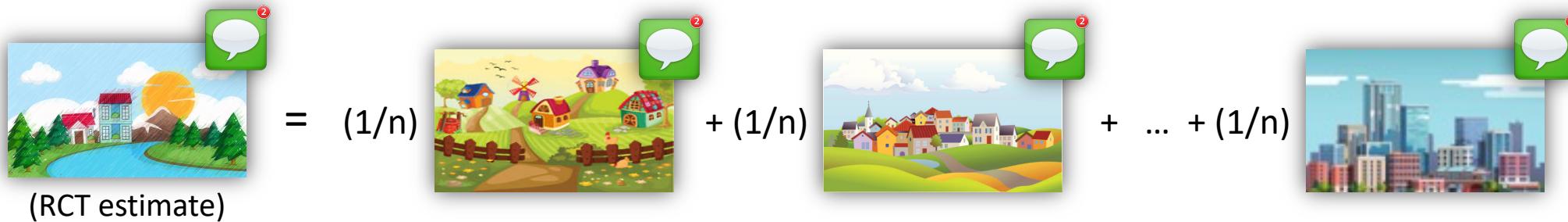
Recreating Observed Immunization Rates



Intervention	000	001	002	010	031	032	040	050	100	101	102	200	201	202	300	301	302	400	401	402
Hyp. Test ($\alpha = 0.05$)	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✗	✓	✓	✗	✗	✓
R^2_{RCT} (vs. RCT)	0.55	0.50	0.48	0.73	0.62	0.73	0.57	0.75	0.50	0.68	0.48	0.70	0.66	0.45	0.34	0.46	0.60	0.29	0.29	0.42

$$R^2_{\text{RCT}} = 1 - \frac{\text{SI error}}{\text{RCT error}}$$

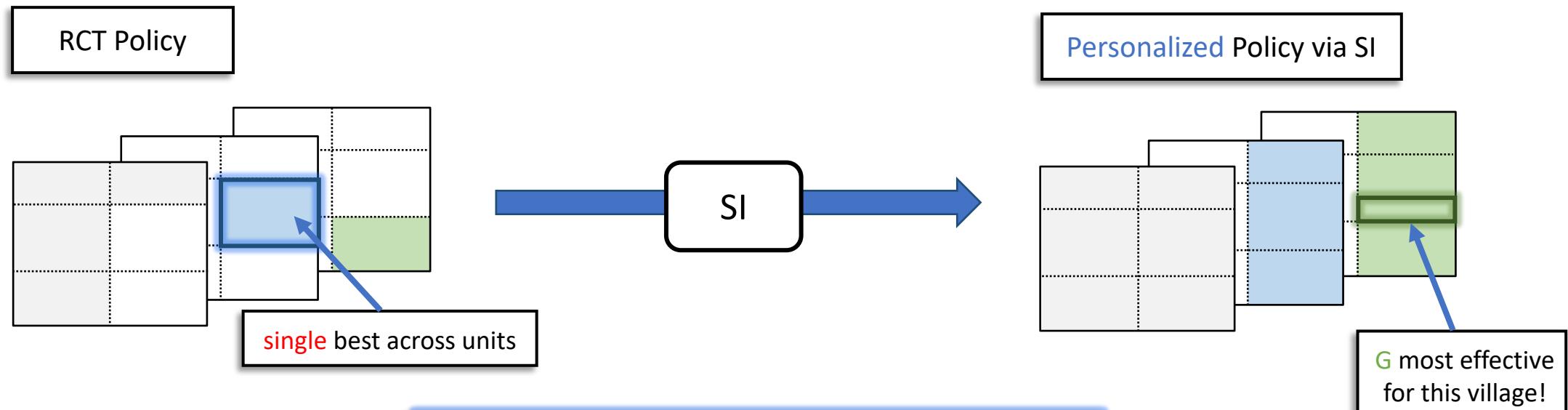
heterogenous villages

Heterogenous villages → SI is a stronger predictor than RCT estimator

Policy Recommendations

Policy Recommendation Method	Avg. net increase in immunization rates (estimated)
Random policy (per village)	1.0
Best RCT policy (031)	1.3
SI's personalized policy (per village)	2.8



SI enables **personalized** policies
from **same** RCT data → significant gains

A/B Testing in E-commerce

Data Efficient RCTs: Web A/B Testing in Ecommerce

25 user groups



(>10k users per group)

3 interventions



(+1 control intervention)

8 days



(measuring customer engagement)

We get access to customer engagement
trajectories of all 25 user groups under all interventions

Data Efficient RCTs: Web A/B Testing in Ecommerce

- Ideal RCT setting – experiments run

Intervention	Groups 1-8	Groups 9-16	Groups 17-25
Control	✓	✓	✓
10% Discount	✓	✓	✓
30% Discount	✓	✓	✓
50% Discount	✓	✓	✓

- Synthetic Interventions – experiments run

Intervention	Groups 1-8	Groups 9-16	Groups 17-25
Control	✓	✓	✓
10% Discount	✓		
30% Discount		✓	
50% Discount			✓

Data Efficient RCTs: Web A/B Testing in Ecommerce

- Hypothesis Test

Intervention	Metric	Projection Test
10% Discount	Subscription rate	(Pass, $\alpha = 0.05$)
30% Discount	Subscription rate	(Pass, $\alpha = 0.05$)
50% Discount	Subscription rate	(Pass, $\alpha = 0.05$)

Data Efficient RCTs: Web A/B Testing in Ecommerce

- Quantifying prediction accuracy
 - R^2 score (access to true counterfactual)

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{reg}}}$$
$$SS_{\text{reg}} = \sum_i \left(Y_{ni}^{(d)} - \bar{Y}_n^{(d)} \right)^2$$
$$SS_{\text{res}} = \sum_i \left(Y_{ni}^{(d)} - \hat{Y}_{ni}^{(d)} \right)^2$$

- Quantifying utility over standard RCTs
 - R^2 score (using RCT as a predictor)

$$R_{\text{rct}}^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{rct}}}$$
$$SS_{\text{rct}} = \sum_i \left(Y_{ni}^{(d)} - \frac{1}{|\mathcal{I}^{(d)}|} \sum_{m \in \mathcal{I}^{(d)}} Y_{mi}^{(d)} \right)^2$$

Intervention	R^2 score (True Counterfactual)	R^2 score (RCT Baseline)
10% Discount	0.76	0.98
30% Discount	0.56	0.99
50% Discount	0.75	0.98

accurate recovery of true customer engagements

heterogeneous user groups

Synthetic Interventions simulates ideal RCT experiment