

# Proximal Causal Inference with Proxies: Valid and Invalid

ACIC 2025

## Instructors



Chan Park  
Assistant Professor at UIUC



Prabirsha Rakshit  
Postdoc at UPenn



Eric Tchetgen Tchetgen  
University Professor at UPenn

## Outline

- 1st half: 8:30am - 10:20am

Instructor: Chan Park

Proximal causal inference with valid proxy variables

- Break: 10:20am - 10:30am

- 2nd half: 10:30am - 12:30am

Instructor: Prabirsha Rakshit

Proximal causal inference with invalid proxy variables

# Outline of the First Half

- Negative Controls

- Brief review of standard causal inference

- Concepts of negative control exposures and negative control outcomes

- Usage of negative controls in the literature

- Proximal Causal Inference

- Identification & estimation of causal effects in the presence of unmeasured confounding

- Application to the SUPPORT study

- Extensions of Proximal Causal Inference to Various Settings

- Longitudinal settings

- Mediation, front-door formula

- Network data

- Time series data

Available at

[github.com/qkrcks0218/ACIC2025](https://github.com/qkrcks0218/ACIC2025)



- Demonstration using R

# Negative Controls

# Causal Inference

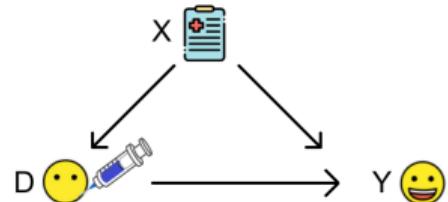
- Elements

Treatment  $D$

Observed outcome  $Y$

Potential outcome  $Y^{(d)}$

(measured) Baseline covariate  $X$



# Causal Inference

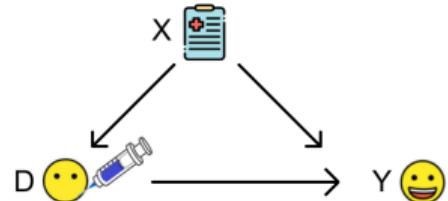
- Elements

Treatment  $D$

Observed outcome  $Y$

Potential outcome  $Y^{(d)}$

(measured) Baseline covariate  $X$



- Flow of Causal Inference

Estimand

$$\tau^* = E[Y^{(1)} - Y^{(0)}]$$

Define an estimand

(eg. Average Treatment Effect (ATE); involves counterfactual, unobserved variables)

# Causal Inference

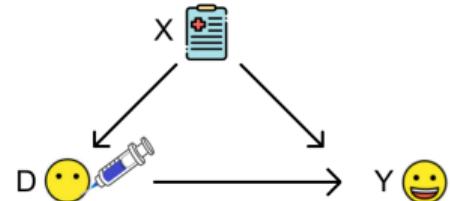
- Elements

Treatment  $D$

Observed outcome  $Y$

Potential outcome  $Y^{(d)}$

(measured) Baseline covariate  $X$



- Flow of Causal Inference

Estimand

$$\tau^* = E[Y^{(1)} - Y^{(0)}]$$

Define an estimand

(eg. Average Treatment Effect (ATE); involves counterfactual, unobserved variables)

- Fundamental problem of causal inference

We only observe either  $Y^{(1)}$  or  $Y^{(0)}$ , but not both

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob	😊	1	30	😊	?
Sally	😢	0	25	?	😢

# Causal Inference

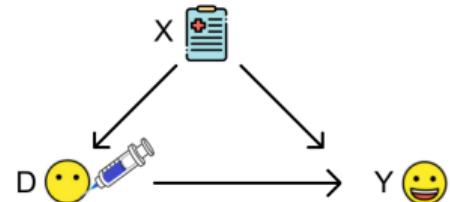
- Elements

Treatment  $D$

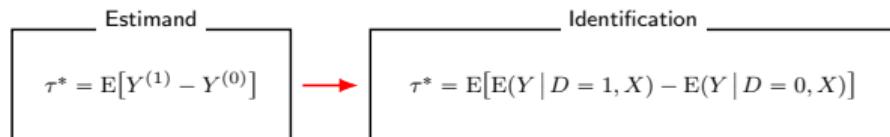
Observed outcome  $Y$

Potential outcome  $Y^{(d)}$

(measured) Baseline covariate  $X$



- Flow of Causal Inference



Define an estimand

(eg. Average Treatment Effect (ATE); involves counterfactual, unobserved variables)

Establish Identification

(eg. g-formula; find a representation of the ATE in terms of the observed data)

- Fundamental problem of causal inference

We only observe either  $Y^{(1)}$  or  $Y^{(0)}$ , but not both

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob	😊	1	30	😊	?
Sally	😢	0	25	?	😢

# Causal Inference

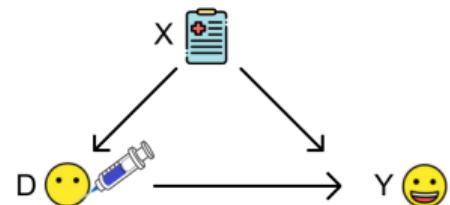
- Elements

Treatment  $D$

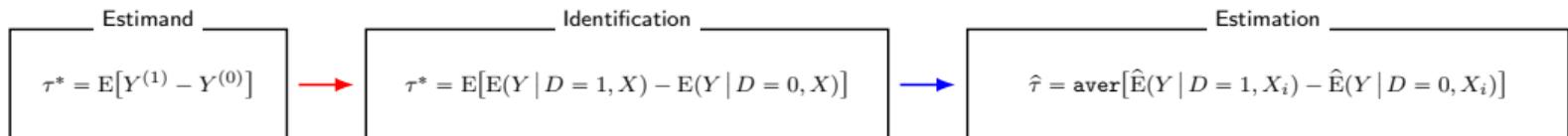
Observed outcome  $Y$

Potential outcome  $Y^{(d)}$

(measured) Baseline covariate  $X$



- Flow of Causal Inference



Define an estimand

(eg. Average Treatment Effect (ATE); involves counterfactual, unobserved variables)

Establish Identification

(eg. g-formula; find a representation of the ATE in terms of the observed data)

Construct an estimator

(eg. OLS; apply statistical methods to the observed data)

- Fundamental problem of causal inference

We only observe either  $Y^{(1)}$  or  $Y^{(0)}$ , but not both

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob	😊	1	30	😊	?
Sally	😢	0	25	?	😢

## Standard Assumptions

- Consistency:  $Y = Y^{(D)}$

The observed outcome = one of the potential outcomes associated with the actually assigned treatment

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob		1	30		?
Sally		0	25	?	

## Standard Assumptions

- Consistency:  $Y = Y^{(D)}$

The observed outcome = one of the potential outcomes associated with the actually assigned treatment

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob		1	30		?
Sally		0	25	?	

- Positivity:  $\Pr(D = d \mid X) > 0$

Given  $X$ , each unit can be assigned to either treatment or control

## Standard Assumptions

- Consistency:  $Y = Y^{(D)}$

The observed outcome = one of the potential outcomes associated with the actually assigned treatment

- **Positivity:**  $\Pr(D = d \mid X) > 0$

Given  $X$ , each unit can be assigned to either treatment or control

- Most importantly, the assumption most relevant to this course:

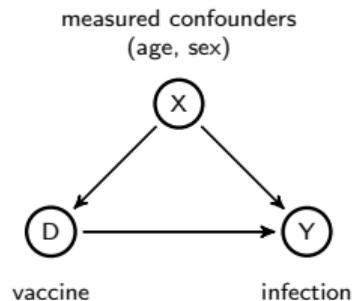
- Ignorability / No Unmeasured Confounding:  $Y^{(d)} \perp\!\!\!\perp D | X$

All common causes of  $Y$  and  $D$  are measured

Treatment is randomized within each stratum of  $X$

Not empirically verifiable

Unit	$Y$	$D$	$X$ (age)	$Y^{(1)}$	$Y^{(0)}$
Bob		1	30		?
Sally		0	25	?	



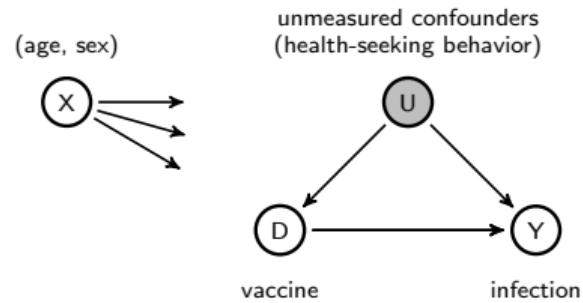
## Unmeasured Confounding

- Hereafter all arguments are made implicitly conditional on  $X$

- Unmeasured confounder  $U$

An **unmeasured** common cause of  $D$  and  $Y$

At the center of much skepticism about observational studies

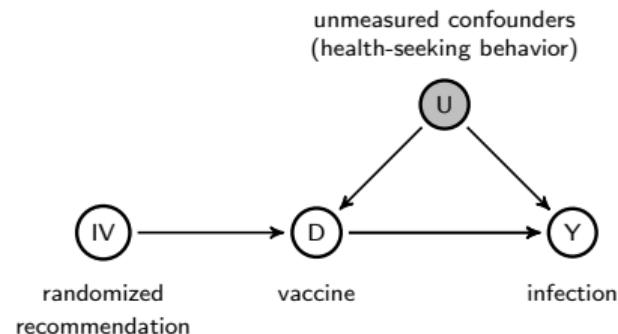


## A Well-known Approach: Instrumental Variables

- An Instrumental Variable (IV) satisfies:

1. Relevance:  $\text{IV} \perp\!\!\!\perp D$
2. Exclusion Restriction:  $Y^{(\text{iv}, d)} = Y^{(d)}$
3. Independence:  $\text{IV} \perp\!\!\!\perp U$

- An ideal IV can be obtained from a randomized assignment



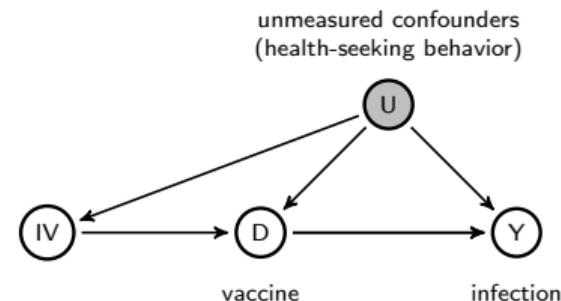
## A Well-known Approach: Instrumental Variables

- An Instrumental Variable (IV) satisfies:

1. Relevance:  $\text{IV} \perp\!\!\!\perp D$
2. Exclusion Restriction:  $Y^{(\text{iv}, d)} = Y^{(d)}$
3. Independence:  $\text{IV} \perp\!\!\!\perp U$

- An ideal IV can be obtained from a randomized assignment

- In observational studies, the independence assumption is often difficult to justify, as it is generally difficult to rule out the possibility that a given variable is correlated with  $U$
- A hidden treasure: negative control variables and proxy variables



## Motivating Example: Does stress during pregnancy affect birth weight? [6, 7]<sup>1</sup>

- Observational study on effect of mother's stress on birth weight

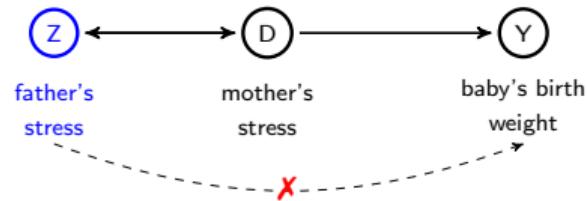


<sup>1</sup> Davey Smith (2008). Assessing intrauterine influences on offspring health outcomes: can epidemiological studies yield robust findings? *Basic & Clinical Pharmacology & Toxicology*

Davey Smith (2012). Negative control exposures in epidemiologic studies. Comments on "Negative controls: a tool for detecting confounding and bias in observational studies." *Epidemiology*

## Motivating Example: Does stress during pregnancy affect birth weight? [6, 7]<sup>1</sup>

- Observational study on effect of mother's stress on birth weight
- No causal effect of father's stress after adjusting for mother's stress
- $\text{lm}(Y \sim D + Z)$  and study the effect of  $Z$



<sup>1</sup> Davey Smith (2008). Assessing intrauterine influences on offspring health outcomes: can epidemiological studies yield robust findings? *Basic & Clinical Pharmacology & Toxicology*

Davey Smith (2012). Negative control exposures in epidemiologic studies. Comments on "Negative controls: a tool for detecting confounding and bias in observational studies." *Epidemiology*

# Motivating Example: Does stress during pregnancy affect birth weight? [6, 7]<sup>1</sup>

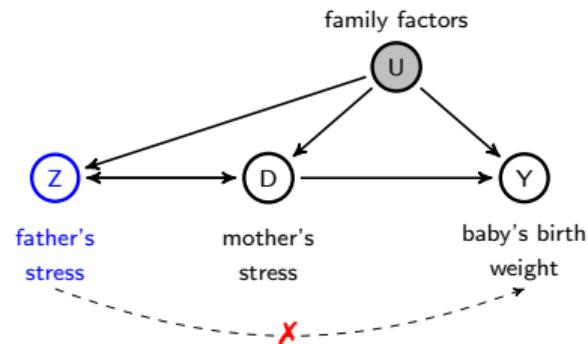
- Observational study on effect of mother's stress on birth weight
- No causal effect of father's stress after adjusting for mother's stress



- $\text{lm}(Y \sim D + Z)$  and study the effect of  $Z$

Nonzero effect of father's stress indicates hidden bias

- Family factors could be an unmeasured confounder



<sup>1</sup> Davey Smith (2008). Assessing intrauterine influences on offspring health outcomes: can epidemiological studies yield robust findings? *Basic & Clinical Pharmacology & Toxicology*

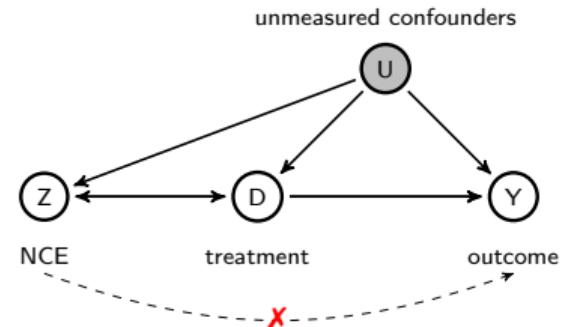
Davey Smith (2012). Negative control exposures in epidemiologic studies. Comments on "Negative controls: a tool for detecting confounding and bias in observational studies." *Epidemiology*

## Negative Control Exposure (NCE)

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$

$Z$  does not causally affect  $Y$

$Z$  is associated with  $Y$  only through  $U$  conditional on  $D$

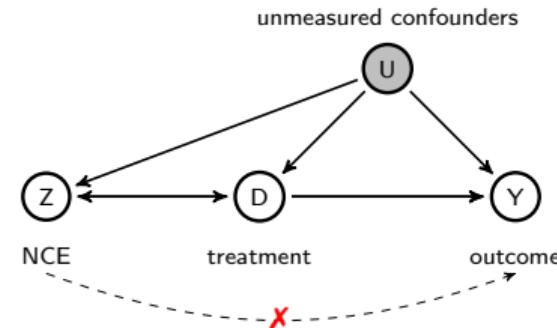


## Negative Control Exposure (NCE)

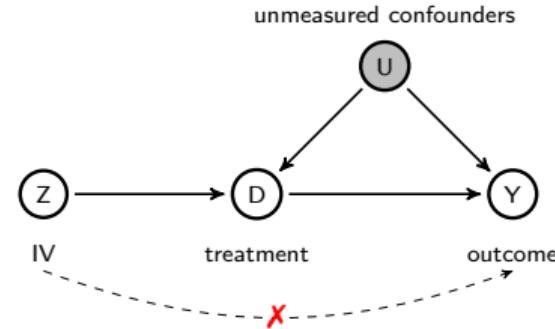
- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$

$Z$  does not causally affect  $Y$

$Z$  is associated with  $Y$  only through  $U$  conditional on  $D$



- A valid IV is a valid NCE

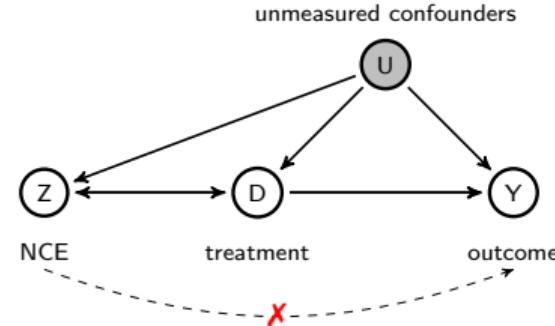


## Negative Control Exposure (NCE)

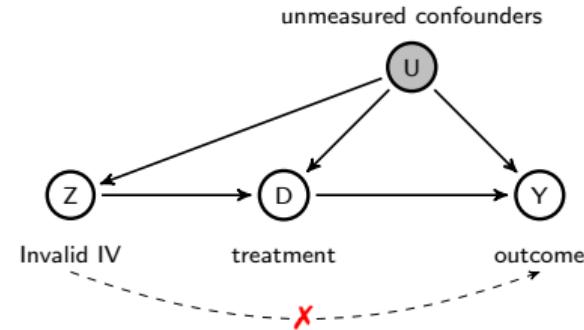
- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$

$Z$  does not causally affect  $Y$

$Z$  is associated with  $Y$  only through  $U$  conditional on  $D$

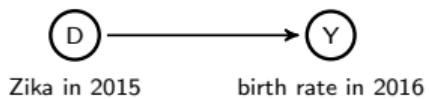
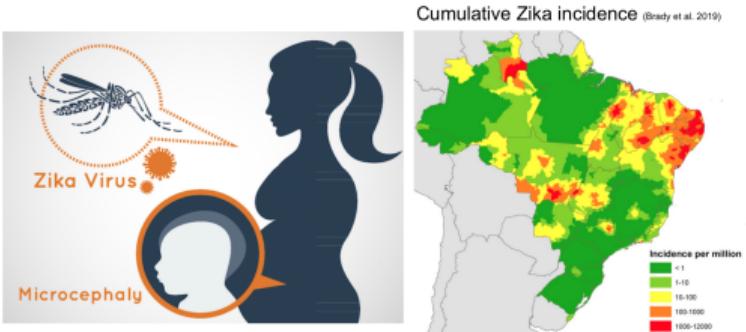


- A valid IV is a valid NCE
- An invalid IV also can be a valid NCE



# Motivating Example: Did the 2015 Zika Virus Outbreak in Brazil Lead to a Decrease in Birth Rate? [23]<sup>1</sup>

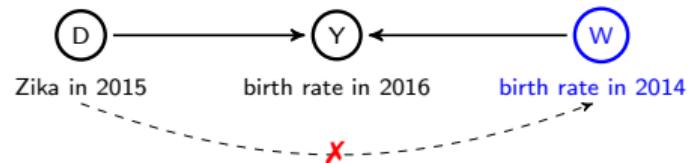
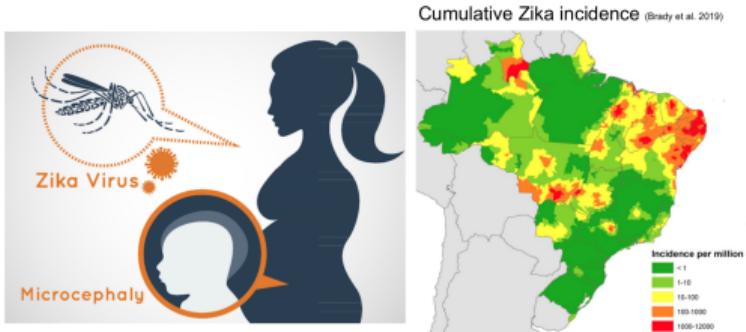
- Observational study on 2015 Zika virus outbreak in Brazil



<sup>1</sup> Park, Richardson, TT (2024). Single proxy control. *Biometrics*

## Motivating Example: Did the 2015 Zika Virus Outbreak in Brazil Lead to a Decrease in Birth Rate? [23]<sup>1</sup>

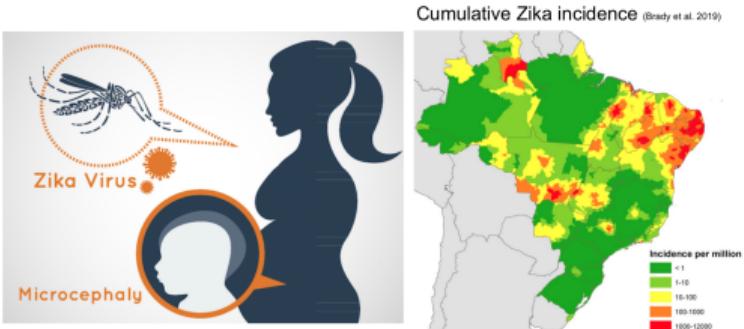
- Observational study on 2015 Zika virus outbreak in Brazil
- No causal effect of 2015 Zika on 2014 birth rate
- $\text{Im}(W \sim D)$  and study the effect of  $D$



<sup>1</sup> Park, Richardson, TT (2024). Single proxy control. *Biometrics*

# Motivating Example: Did the 2015 Zika Virus Outbreak in Brazil Lead to a Decrease in Birth Rate? [23]<sup>1</sup>

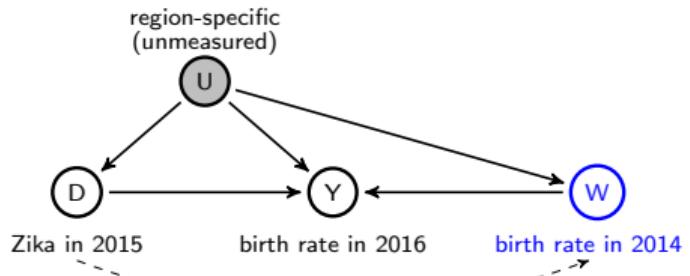
- Observational study on 2015 Zika virus outbreak in Brazil



- No causal effect of 2015 Zika on 2014 birth rate

- $\text{Im}(W \sim D)$  and study the effect of  $D$   
nonzero effect of Zika indicates hidden bias

- Potential unmeasured confounding by region-specific variables  
healthcare infrastructure, socioeconomic status,  
cultural perspectives on childbirth



<sup>1</sup>Park, Richardson, TT (2024). Single proxy control. *Biometrics*

## Negative Control Exposure (NCE) and Negative Control Outcome (NCO)

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$

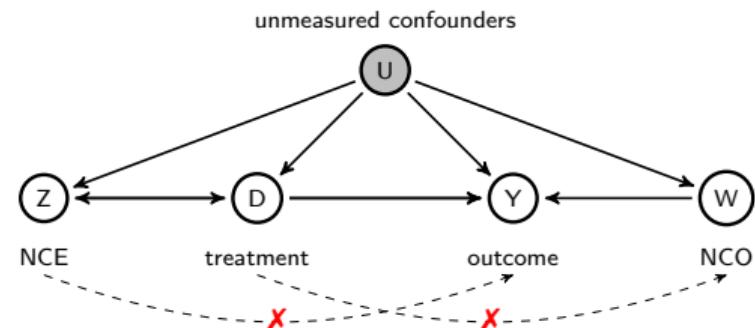
$Z$  does not causally affect  $Y$

$Z$  is associated with  $Y$  only through  $U$  conditional on  $D$

- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | U$

$W$  is not causally affected by  $D$

$W$  is associated with  $(D, Z)$  only through  $U$



## Negative Control Exposure (NCE) and Negative Control Outcome (NCO)

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$

$Z$  does not causally affect  $Y$

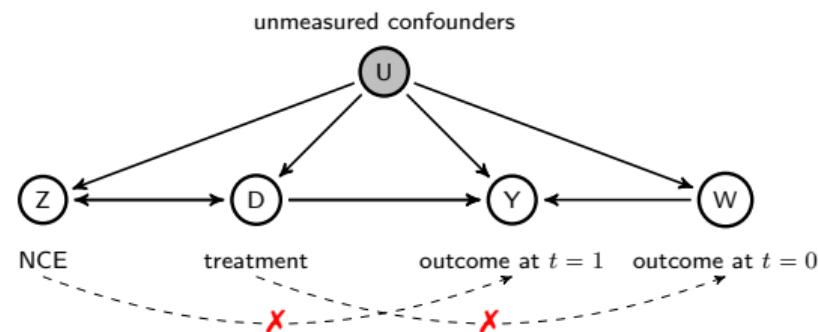
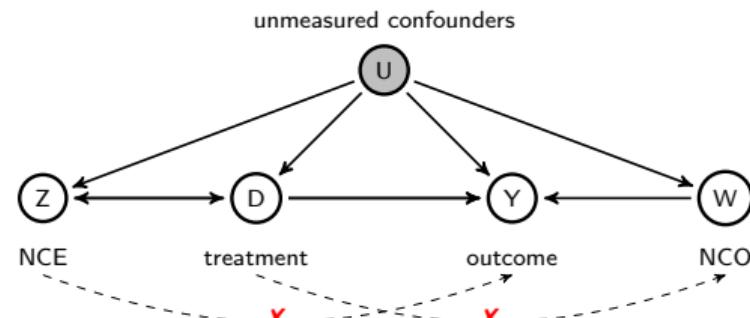
$Z$  is associated with  $Y$  only through  $U$  conditional on  $D$

- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | U$

$W$  is not causally affected by  $D$

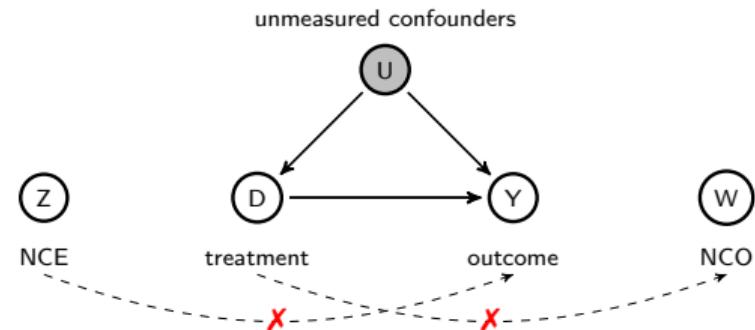
$W$  is associated with  $(D, Z)$  only through  $U$

- A pre-treatment outcome is a good NCO candidate



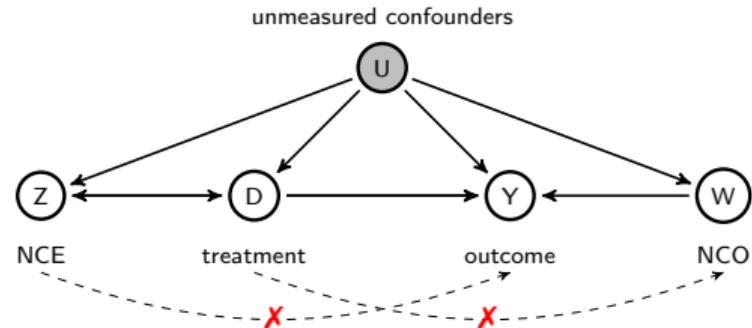
## U-Comparability

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y \mid (U, D)$   
 $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) \mid U$
- Toss a coin, and use the result as  $Z, W$  (?)



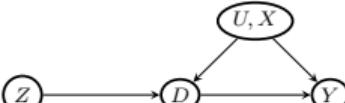
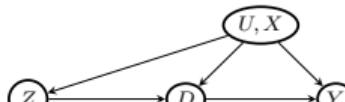
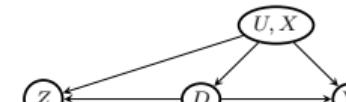
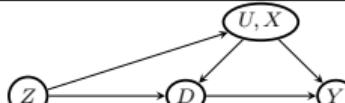
## U-Comparability

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D)$   
 $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | U$
- Toss a coin, and use the result as  $Z, W$  (?)
- A variable completely irrelevant to the problem would not provide any useful information
- Unmeasured confounding mechanism of NCs should be comparable to that of  $D$  and  $Y$
- **U-comparability:**  $Z \not\perp\!\!\!\perp U | D$  and  $W \not\perp\!\!\!\perp U$



## Examples Encoding NC Assumptions [31]<sup>1</sup>

- Gray indicates violation of NC assumptions  $Z \perp\!\!\!\perp Y | (U, D)$  and  $Z \not\perp\!\!\!\perp U | D$

Examples of NCE			
	$Z \rightarrow D$ (pre-treatment)	$D \rightarrow Z$ (post-treatment)	$Z \perp\!\!\!\perp D   (U, X)$
No arrow between $U$ and $Z$ (may violate U-comparability)	Instrumental variable (IV) 		
$U \rightarrow Z$	Invalid IV 		
$Z \rightarrow U$	May violate Assumptions if there is $W \rightarrow U$ (collider bias)		
			

<sup>1</sup> Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

## Examples Encoding NC Assumptions [31]<sup>1</sup>

- Gray indicates violation of NC assumptions  $W \perp\!\!\!\perp (D, Z) \mid U$  and  $W \not\perp\!\!\!\perp U$

Examples of NCO			
	$W \rightarrow Y$ (pre-treatment)	$Y \rightarrow W$ (post-treatment)	$Y \perp\!\!\!\perp W \mid (D, U, X)$
No arrow between $U$ and $W$ (may violate U-comparability)			
$U \rightarrow W$	<b>Pre-trt Outcome</b> 		
$W \rightarrow U$	May violate Assumption if there is $Z \rightarrow U$ (collider bias)		

<sup>1</sup> Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

# Negative Controls Are Widely Available

- Air pollution and health outcomes: the future  $\Rightarrow$  the past [10]<sup>1</sup>

NCE = future exposure; NCO = past outcome

- Genetics research and batch effect [16]<sup>2</sup>

Use control genes ( $W$ ) to remove unwanted variation

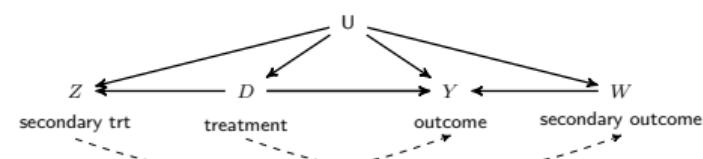
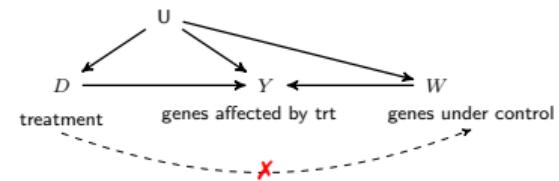
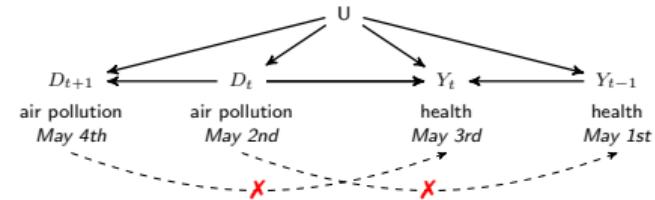
$$W = \gamma U + \epsilon_W, Y = \beta D + \gamma U + \epsilon_Y$$

- Drug/vaccine comparative effectiveness and safety [29]<sup>3</sup>

Use secondary treatments or outcomes in electronic health records

Obtain the p-values from  $\text{lm}(Y \sim Z)$  and  $\text{lm}(W \sim D)$

Calibrate the p-values from  $\text{lm}(Y \sim D)$  to correct for the deflation



<sup>1</sup> Flanders et al. (2011). A method for detection of residual confounding in time-series and other observational studies. *Epidemiology*

<sup>2</sup> Jacob et al. (2016) Correcting gene expression data when neither the unwanted variation nor the factor of interest are observed. *Biostatistics*

<sup>3</sup> Schuemie et al. (2014) Interpreting observational studies: why empirical calibration is needed to correct p-values. *Statistics in Medicine*

# Using Negative Controls to Detect, Reduce, and Correct Bias

Detect	[10] <sup>1</sup>	Time-series study
	[7, 39] <sup>2</sup>	invalid NCE
Reduce	[11, 22] <sup>3</sup>	Time-series study
	[26] <sup>4</sup>	Standardized mortality ratio in occupational cohort study
	[28, 29] <sup>5</sup>	Drug–outcome pairs with no plausible causal effect
Debias	[27, 36] <sup>6</sup>	Time-to-event outcome
	[15, 34] <sup>7</sup>	Generalized difference-in-differences using NCO
	[23, 35] <sup>8</sup>	Calibration using NCO
	[12, 16, 38] <sup>9</sup>	Removing unwanted variation in gene-expression analysis
	[5, 20, 21] <sup>10</sup>	Nonparametric identification using double negative control/proximal causal inference

1 Flanders et al. (2011). A method for detection of residual confounding in time-series and other observational studies. *Epidemiology*

2 Davey Smith (2012). Negative control exposures in epidemiologic studies. Comments on “Negative controls: a tool for detecting confounding and bias in observational studies.” *Epidemiology* Weisskopf, TT, Raz (2016). Commentary: on the use of imperfect negative control exposures in epidemiologic studies. *Epidemiology*

3 Flanders, Strickland, Klein (2017). A new method for partial correction of residual confounding in time-series and other observational studies. *Am J Epi* Miao, TT (2017). Invited commentary: bias attenuation and identification of causal effects with multiple negative controls. *Am J Epi*

4 Richardson et al. (2015). Negative control outcomes and the analysis of standardized mortality ratios. *Epidemiology*

5 Schuemie et al. (2014) Interpreting observational studies: why empirical calibration is needed to correct p-values. *Statistics in Medicine* Schuemie et al. (2018). Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data. *PNAS*

6 Richardson et al. (2014). Assessment and indirect adjustment for confounding by smoking in cohort studies using relative hazards models. *Am J Epi* TT, Sofer, Richardson (2015). NCO for unobserved confounding under a Cox proportional hazards model. *preprint*

7 Sofer et al. (2016). On negative outcome control of unobserved confounding as a generalization of DiD. *Statistical Science*, Glynn, Ichino (2019). Generalized Nonlinear DiDiD. *preprint*

8 TT (2014). The control outcome calibration approach for causal inference with unobserved confounding. *Am J Epi*, Park, Richardson, TT (2024). Single proxy control. *Biometrics*

9 Gagnon-Bartsch, Speed (2012). Using control genes to correct for unwanted variation in microarray data. *Biostatistics* Wang et al. (2017). Confounder adjustment in multiple hypothesis testing. *Annals of Statistics* Jacob et al. (2016) Correcting gene expression data when neither the unwanted variation nor the factor of interest are observed. *Biostatistics*

10 Miao, Geng, TT (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika* Cui, Pu, Shi, Miao, TT (2024). Semiparametric proximal causal inference. *JASA*

Miao, Shi, Li, TT (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*

# Proximal Causal Inference

## Can Imperfect Measurements Still Be Useful?

- Survey Studies [3]<sup>1</sup>
- Obtaining a variable **without measurement error** often requires significant time and resources

This can be an extremely challenging task—perhaps impossible
- **Noisy measures** can still be highly informative about the quantity of interest, provided the measurement errors have the **right structure**
- NC/Proximal Causal Inference
- The **no unmeasured confounding** assumption depends on investigator's ability to accurately measure covariates capturing all potential sources of confounding

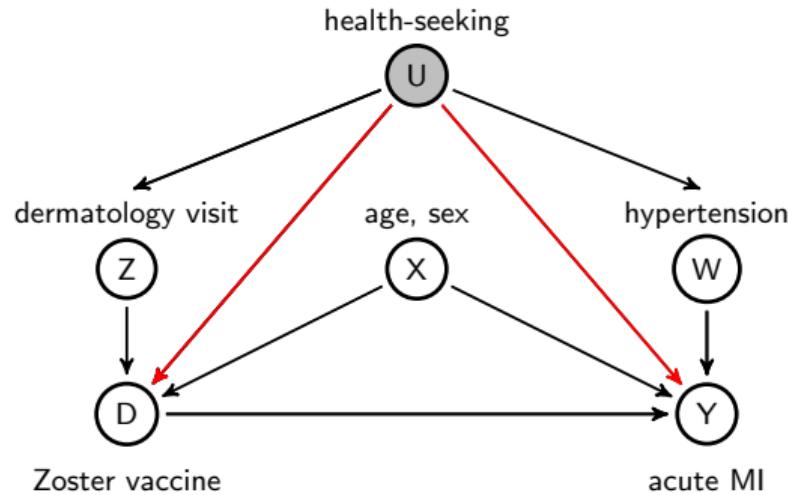
This assumption can be unrealistic  
In practice, the most one can hope for is that measured covariates are at best proxies of the true underlying confounding mechanism
- Acknowledge that measured covariates are **imperfect proxies** of confounders

If these proxies satisfying **certain assumptions**, it may still be possible to infer causal quantities of interest

<sup>1</sup> Browning, Crossley (2009). Are two cheap, noisy measures better than one expensive, accurate one? *American Economic Review*

## An Example in Vaccine Safety Study [17]<sup>1</sup>

- Adverse effect of a new Zoster vaccine on acute Myocardial Infarction (MI; heart attack)
- Plan to adjust for the following baseline variables
  - age, sex ( $X$ )
  - dermatology visit ( $Z$ ; check skin)
  - hypertension ( $W$ ; high blood pressure)
- Naively adjusting ( $X, W, Z$ ) does not eliminate confounding bias due to  $U$
- Need to use  $X, W, Z$  differently!



<sup>1</sup> Li et al. (2024). Using Double Negative Controls to Adjust for Healthy User Bias in a Recombinant Zoster Vaccine Safety Study. *Am J Epi*

# An Example in Vaccine Safety Study [17]<sup>1</sup>

- Three types of measured covariates

- Confounder

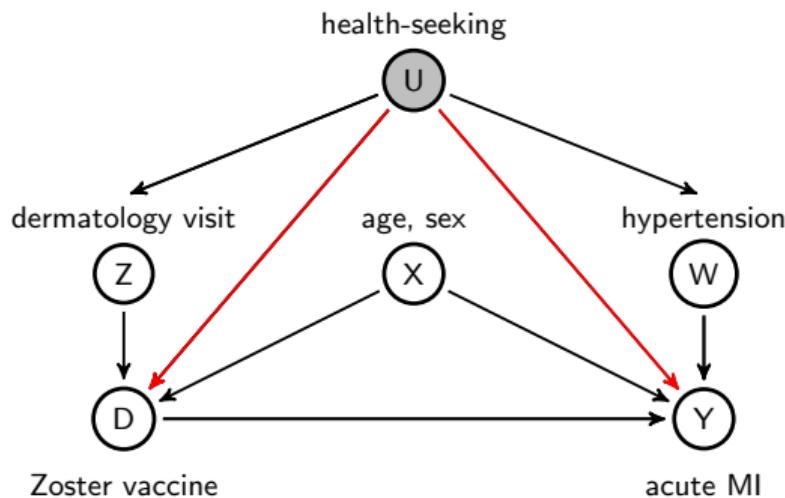
Common causes of the treatment and outcome  
age, sex ( $X$ )

- Treatment confounding proxy

$Y \perp\!\!\!\perp Z | (D, U, X)$  (NCE)  
dermatology visit ( $Z$ )

- Outcome confounding proxy

$W \perp\!\!\!\perp (D, Z) | (U, X)$  (NCO)  
hypertension ( $W$ )



<sup>1</sup> Li et al. (2024). Using Double Negative Controls to Adjust for Healthy User Bias in a Recombinant Zoster Vaccine Safety Study. *Am J Epi*

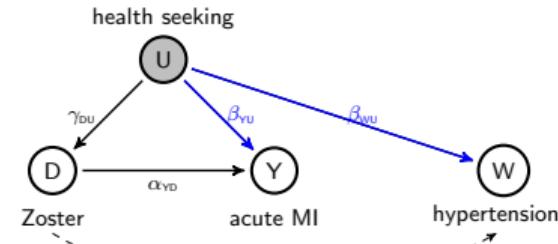
## Intuitions for Identification

- $[V_1 \sim V_2] = \text{coef}(\text{lm}(V_1 \sim V_2))$   
= direct effect of  $V_2$  on  $V_1$   
+ confounding bias between  $V_2$  and  $V_1$
- $[Y \sim D] = \alpha_{YD} + [D \sim U] \times [Y \sim U] = \alpha_{YD} + \gamma_{DU} \beta_{YU}$
- $[W \sim D] = \alpha_{WD} + [D \sim U] \times [W \sim U] = \gamma_{DU} \beta_{WU}$

$$\begin{array}{lcl} [Y \sim D] & = & \alpha_{YD} + \gamma_{DU} \beta_{YU} \\ [W \sim D] & = & \gamma_{DU} \beta_{WU} \\ \hline \text{diff in coeffs of } D & = & \alpha_{YD} \end{array}$$

- A Special Case: Difference-in-Differences (DiD)

$$[Y_0 \sim U] = [Y_1 \sim U] \quad [27, 34, 36]^1$$



<sup>1</sup> Richardson et al. (2014). Assessment and indirect adjustment for confounding by smoking in cohort studies using relative hazards models. *Am J Epi*  
 TTT, Sofer, Richardson (2015). NCO for unobserved confounding under a Cox proportional hazards model. *preprint*  
 Sofer et al. (2016). On negative outcome control of unobserved confounding as a generalization of DiD. *Statistical Science*

## Intuitions for Identification

- $[V_1 \sim V_2] = \text{coef}(\text{lmm}(V_1 \sim V_2))$   
= direct effect of  $V_2$  on  $V_1$   
+ confounding bias between  $V_2$  and  $V_1$

- $[Y \sim D] = \alpha_{YD} + [D \sim U] \times [Y \sim U] = \alpha_{YD} + \gamma_{DU} \beta_{YU}$

- $[Y \sim Z] = \alpha_{YZ} + [Z \sim U] \times [Y \sim U] = \gamma_{ZU} \beta_{YU}$

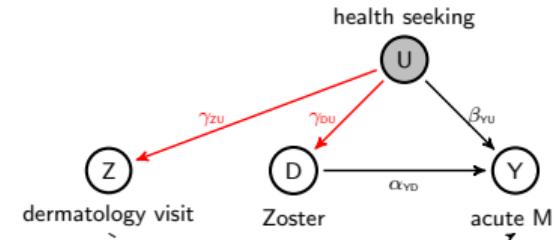
- If  $[D \sim U] = [Z \sim U] \Leftrightarrow \gamma_{DU} = \gamma_{ZU}$

$$[Y \sim D] = \alpha_{YD} + \gamma_{DU} \beta_{YU}$$

$$[Y \sim Z] = \gamma_{ZU} \beta_{YU}$$

---


$$\text{diff in coeffs} = \alpha_{YD}$$

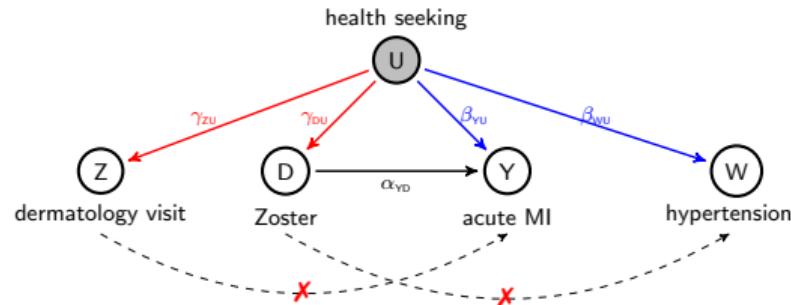


- A Special Case: air pollution studies [10, 11, 22]<sup>1</sup>

<sup>1</sup> Flanders et al. (2011). A method for detection of residual confounding in time-series and other observational studies. *Epidemiology*  
 Flanders, Strickland, Klein (2017). A new method for partial correction of residual confounding in time-series and other observational studies. *Am J Epi*  
 Miao, TT (2017). Invited commentary: bias attenuation and identification of causal effects with multiple negative controls. *Am J Epi*

## Identification Strategy

- What if  $[Y \sim U] \neq [W \sim U]$  and  $[D \sim U] \neq [Z \sim U]$ ?
- Solution: Use both proxies [20, 21, 31, 32]<sup>1</sup>

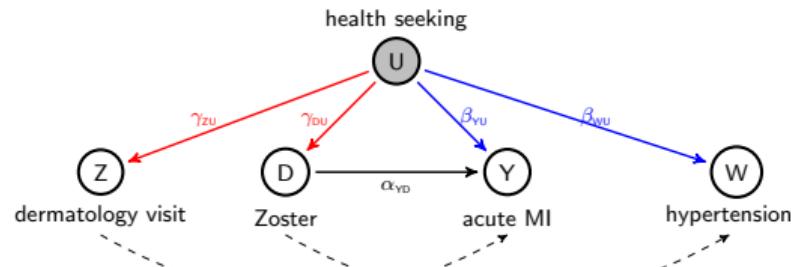


<sup>1</sup> Miao, Shi, Li, TT (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*  
Miao, Geng, TT (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*  
Shi, Miao, TT (2020). Multiply robust causal inference with double negative control adjustment for categorical unmeasured confounding. *JRSSB*  
Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

## Identification Strategy

- What if  $[Y \sim U] \neq [W \sim U]$  and  $[D \sim U] \neq [Z \sim U]$ ?
- Solution: Use both proxies  $[20, 21, 31, 32]$ <sup>1</sup>
- Note that

$$\begin{aligned} [Y \sim D] &= \alpha_{YD} + \gamma_{DU}\beta_{YU} \\ [W \sim D] &= \gamma_{DU}\beta_{WU} \\ \hline \Rightarrow [Y \sim D] &= \alpha_{YD} + \underbrace{\frac{\beta_{YU}}{\beta_{WU}} [W \sim D]}_{\text{bias}} \end{aligned}$$



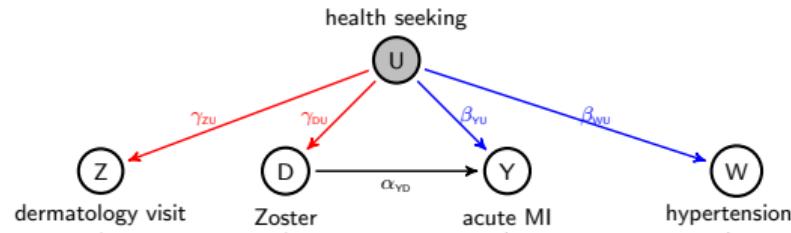
<sup>1</sup> Miao, Shi, Li, TT (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*  
 Miao, Geng, TT (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*  
 Shi, Miao, TT (2020). Multiply robust causal inference with double negative control adjustment for categorical unmeasured confounding. *JRSSB*  
 Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

## Identification Strategy

- What if  $[Y \sim U] \neq [W \sim U]$  and  $[D \sim U] \neq [Z \sim U]$ ?
- Solution: Use both proxies [20, 21, 31, 32]<sup>1</sup>
- Note that

$$\begin{aligned} [Y \sim D] &= \alpha_{YD} + \gamma_{DU}\beta_{YU} \\ [W \sim D] &= \gamma_{DU}\beta_{WU} \\ \hline \Rightarrow [Y \sim D] &= \alpha_{YD} + \underbrace{\frac{\beta_{YU}}{\beta_{WU}} [W \sim D]}_{=bias} \end{aligned}$$

$$\begin{aligned} [Y \sim Z] &= \gamma_{ZU}\beta_{YU} \\ [W \sim Z] &= \gamma_{ZU}\beta_{WU} \end{aligned} \Rightarrow \frac{[Y \sim Z]}{[W \sim Z]} = \underbrace{\frac{\beta_{YU}}{\beta_{WU}}}_{=scale}$$



- $W$  recovers bias up to a scale
- $Z$  recovers that scale

<sup>1</sup> Miao, Shi, Li, TT (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*  
 Miao, Geng, TT (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*  
 Shi, Miao, TT (2020). Multiply robust causal inference with double negative control adjustment for categorical unmeasured confounding. *JRSSB*  
 Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

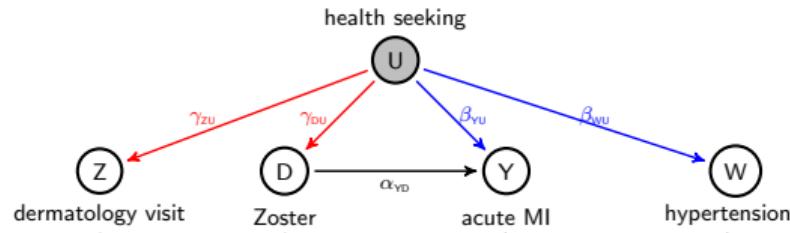
## Identification Strategy

- What if  $[Y \sim U] \neq [W \sim U]$  and  $[D \sim U] \neq [Z \sim U]$ ?
- Solution: Use both proxies [20, 21, 31, 32]<sup>1</sup>
- Note that

$$\begin{aligned} [Y \sim D] &= \alpha_{YD} + \gamma_{DU}\beta_{YU} \\ [W \sim D] &= \gamma_{DU}\beta_{WU} \\ \hline \Rightarrow [Y \sim D] &= \alpha_{YD} + \underbrace{\frac{\beta_{YU}}{\beta_{WU}} [W \sim D]}_{=bias} \end{aligned}$$

$$\begin{aligned} [Y \sim Z] &= \gamma_{ZU}\beta_{YU} \Rightarrow \frac{[Y \sim Z]}{[W \sim Z]} = \frac{\beta_{YU}}{\beta_{WU}} \\ [W \sim Z] &= \gamma_{ZU}\beta_{WU} \quad = scale \end{aligned}$$

$$\Rightarrow \alpha_{YD} = [Y \sim D] - \frac{[Y \sim Z]}{[W \sim Z]} [W \sim D]$$



- $W$  recovers bias up to a scale
- $Z$  recovers that scale

<sup>1</sup> Miao, Shi, Li, TT (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*  
 Miao, Geng, TT (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*  
 Shi, Miao, TT (2020). Multiply robust causal inference with double negative control adjustment for categorical unmeasured confounding. *JRSSB*  
 Shi, Miao, TT (2020). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*

## Linear Additive Model

- Consider

$$E(Y | D, Z, U) = \beta_0 + \beta_D D + \beta_U U$$

$$E(W | D, Z, U) = \gamma_0 + \gamma_U U$$

$$\beta_D = \text{ATE} = E[Y^{(1)} - Y^{(0)}] = E[Y^{(1)} - Y^{(0)} | U]$$

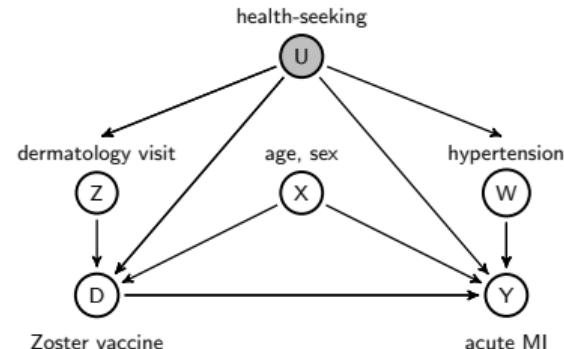
- One can establish that

$$E(Y | D, Z) = \tilde{\beta}_0 + \beta_D D + \tilde{\beta}_U E(W | D, Z)$$

- Two-stage Regression Procedure [37]<sup>1</sup>

$$1. \quad \widehat{W} \leftarrow \text{lm}(W \sim D + Z)$$

$$2. \quad \beta_D \leftarrow \text{coef}(\text{lm}(Y \sim D + \widehat{W}))$$



<sup>1</sup>TT et al. (2024). An Introduction to Proximal Causal Inference. *Statistical Science*

## Generalized Linear Additive Models

- Liu et al. [18]<sup>1</sup>: Generalization of the two-stage regression procedure to count/binary/categorical  $Y$  and  $W$
- The coefficient of  $D$  in the 2nd stage is consistent for the causal effect (on the scale defined by the corresponding  $Y$  link function)

		$Y$		
		Continuous (Identity Link)	Count (Log Link)	Binary/Categorical (Logit Link)
$W$	Continuous (Identity Link)	$\text{Linear } W \sim D + Z$ $S = E(W D, Z)$ $\text{Linear } Y \sim D + S$	$\text{Linear } W \sim D + Z$ $S = E(W D, Z)$ $\text{Poisson } Y \sim D + S$	$\text{Linear } W \sim D + Z + Y$ $S = E(W D, Z, Y = 0)$ $\text{Logistic } \mathbb{1}(Y = t) \sim D + S$
	Count (Log Link)	$\text{Poisson } W \sim D + Z$ $S = \log(E(W D, Z))$ $\text{Linear } Y \sim D + S$	$\text{Poisson } W \sim D + Z$ $S = \log(E(W D, Z))$ $\text{Poisson } Y \sim D + S$	$\text{Poisson } W \sim D + Z + Y$ $S = \log(E(W D, Z, Y = 0))$ $\text{Logistic } \mathbb{1}(Y = t) \sim D + S$
	Binary/Categorical (Logit Link)	$\text{Logistic } \mathbb{1}(W = k) \sim D + Z$ $S_k = \text{logit}(P(W = k D, Z))$ $\text{Linear } Y \sim D + S + W$	$\text{Logistic } \mathbb{1}(W = k) \sim D + Z$ $S_k = \text{logit}(P(W = k D, Z))$ $\text{Poisson } Y \sim D + S + W$	$\text{Logistic } \mathbb{1}(W = k) \sim D + Z + Y$ $S_k = \text{logit}(P(W = k D, Z, Y = 0))$ $\text{Logistic } \mathbb{1}(Y = t) \sim D + S + W$

<sup>1</sup>Liu, Park, Li, TT (2024). Regression-Based Proximal Causal Inference. *Am J Epi*

## Review: Nonparametric Identification Under No Unmeasured Confounding

- Assuming no unmeasured confounding  $Y^{(d)} \perp\!\!\!\perp D \mid L$  where  $L$  = pre-treatment covariates

$$\begin{aligned} E[Y^{(d)}] &= E[p(d, L) \mathbb{1}(D = d) Y] && \text{IPW} \\ &= E[\mu(d, L)] && \text{g-formula} \\ &= E[p(d, L) \mathbb{1}(D = d) \{Y - \mu(d, L)\} + \mu(d, X)] && \text{AIPW} \end{aligned}$$

where

$$p(d, L) = \frac{1}{\Pr(D = d \mid L)} = \text{inverse propensity score}$$

$$\mu(d, L) = E(Y \mid D = d, L) = \text{outcome regression}$$

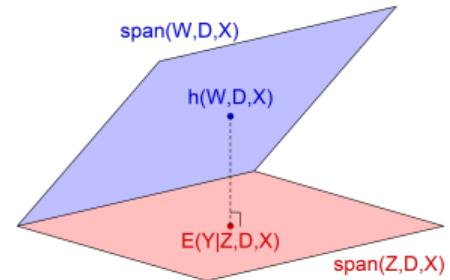
- These identifying formulas are nonparametric
- Question:** Can we establish similar identifying formulas in the presence of  $U$  using proximal causal inference?

## Bridge Function

- Bridge function:

Find  $h$  such that  $E[h(W, D, X)|Z, D, X] = E(Y|Z, D, X)$

Find  $q$  such that  $E[q(Z, D, X)|W, D, X] = \frac{1}{\Pr(D|W, X)}$



## Bridge Function

- Bridge function:

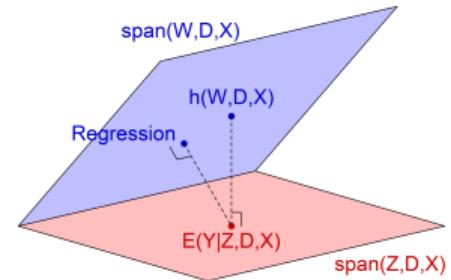
Find  $h$  such that  $E[h(W, D, X)|Z, D, X] = E(Y|Z, D, X)$

Find  $q$  such that  $E[q(Z, D, X)|W, D, X] = \frac{1}{\Pr(D|W, X)}$

- $h$  and  $q$  cannot be obtained from standard regression!

$h \neq$  projection of  $E(Y|Z, D, X)$  onto  $(W, D, X)$

$h =$  a function whose projection onto  $(Z, D, X)$  is  $E(Y|Z, D, X)$



## Bridge Function

- Bridge function:

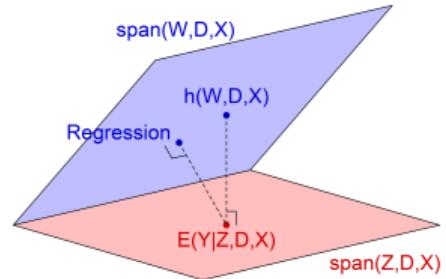
Find  $h$  such that  $E[h(W, D, X) | Z, D, X] = E(Y | Z, D, X)$

Find  $q$  such that  $E[q(Z, D, X) | W, D, X] = \frac{1}{\Pr(D | W, X)}$

- $h$  and  $q$  cannot be obtained from standard regression!

$h \neq$  projection of  $E(Y | Z, D, X)$  onto  $(W, D, X)$

$h =$  a function whose projection onto  $(Z, D, X)$  is  $E(Y | Z, D, X)$



- Ill-posed inverse problems, Fredholm integral equations of the first kind

## Bridge Function

- Bridge function:

Find  $h$  such that  $E[h(W, D, X)|Z, D, X] = E(Y|Z, D, X)$

Find  $q$  such that  $E[q(Z, D, X)|W, D, X] = \frac{1}{\Pr(D|W, X)}$

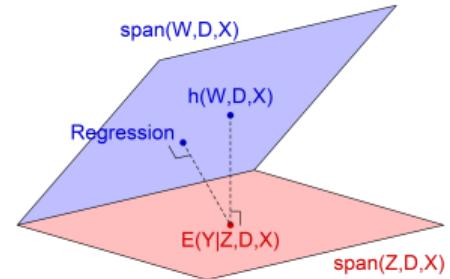
- $h$  and  $q$  cannot be obtained from standard regression!

$h \neq$  projection of  $E(Y|Z, D, X)$  onto  $(W, D, X)$

$h =$  a function whose projection onto  $(Z, D, X)$  is  $E(Y|Z, D, X)$

- Ill-posed inverse problems, Fredholm integral equations of the first kind

- $h =$  outcome confounding bridge function,     $q =$  treatment confounding bridge function



## Nonparametric Identification [5, 37]<sup>1</sup>

- If  $h$  and  $q$  exist,

$$\begin{aligned} \mathbb{E}[Y^{(d)}] &= \mathbb{E}[q(Z, d, X)\mathbb{1}(D = d)Y] && \text{Proximal IPW} \\ &= \mathbb{E}[h(W, d, X)] && \text{Proximal g-formula} \\ &= \mathbb{E}[q(Z, d, X)\mathbb{1}(D = d)\{Y - h(W, d, X)\} + h(W, d, X)] && \text{Proximal AIPW} \end{aligned}$$

- Under no unmeasured confounding:

$$\begin{aligned} \mathbb{E}[Y^{(d)}] &= \mathbb{E}[p(d, L)\mathbb{1}(D = d)Y] && \text{IPW} \\ &= \mathbb{E}[\mu(d, L)] && \text{g-formula} \\ &= \mathbb{E}[p(d, L)\mathbb{1}(D = d)\{Y - \mu(d, L)\} + \mu(d, X)] && \text{AIPW} \end{aligned}$$

<sup>1</sup>Cui, Pu, Shi, Miao, TT (2024). Semiparametric proximal causal inference. *JASA*, TT et al. (2024). An Introduction to Proximal Causal Inference. *Statistical Science*

## Example: Binary $Z, W$

- Suppressing  $X$ ,  $h(w, d)$  solves

$$\mathbb{E}[h(W, D)|Z, D] = \mathbb{E}(Y|Z, D)$$

$$\Leftrightarrow \underbrace{\begin{bmatrix} \Pr(W=1|Z=1, D) & \Pr(W=0|Z=1, D) \\ \Pr(W=1|Z=0, D) & \Pr(W=0|Z=0, D) \end{bmatrix}}_{=B(D)} \begin{bmatrix} h(W=1, D) \\ h(W=0, D) \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbb{E}(Y|Z=1, D) \\ \mathbb{E}(Y|Z=0, D) \end{bmatrix}}_{=C(D)}$$

if the inverse exists  
 $\Rightarrow \begin{bmatrix} h(W=1, D) \\ h(W=0, D) \end{bmatrix} = B^{-1}(D)C(D)$

- $\mathbb{E}[Y^{(d)}]$  is identified by

$$\mathbb{E}[Y^{(d)}] = \mathbb{E}[h(W, d)] = \Pr(W=1)h(W=1, d) + \Pr(W=0)h(W=0, d)$$

## Review: Estimation Under No Unmeasured Confounding

- Posit parametric models for  $p$  and  $\mu$ <sup>1</sup>

$$p(D, X; \theta_p) = 1/\Pr(D | X; \theta_p); \quad \mu(D, X; \theta_\mu) = E(Y | D, X; \theta_\mu)$$

- Estimate  $\theta_p$  and  $\theta_\mu$  from regression models

`glm(D ~ X, family="binomial"); lm(Y ~ D + X)`

- Obtain estimators using the ID formula

$$\text{AIPW estimator} = \text{aver} \left[ \hat{p}(D, X) \mathbb{1}(D = d) \{Y - \hat{\mu}(d, X)\} + \hat{\mu}(d, X) \right]$$

- Asymptotically normal for  $E[Y^{(d)}]$  under certain conditions

<sup>1</sup>One can use nonparametric/machine learning approaches with cross-fitting

## Estimation for Proximal Causal Inference

- Outcome/Treatment confounding bridge function:

$$E[h(W, D, X) | Z, D, X] = E(Y | Z, D, X) \quad \Rightarrow \quad E[g_h(Z, D, X)\{Y - h(W, D, X)\}] = 0 \text{ for any } g_h$$

- Posit parametric models for  $h = h(w, d, x; \theta_h)$  and  $q = q(z, d, x; \theta_q)$
- Estimate  $h$  and  $q$  from generalized method of moments (GMM)<sup>1</sup>

$$\hat{\theta}_h \leftarrow \text{aver} \left[ g_h(Z, D, X) \{Y - h(W, D, X; \theta_h)\} \right] = 0,$$

`ivreg:::ivreg(Y~W+D+X | Z+D+X)`  
`gmm::gmm(g=Y~W+D+X, x=~Z+D+X)`

- Obtain estimators using the ID formula

$$\text{Proximal AIPW estimator} = \text{aver} \left[ \hat{q}(Z, d, X) \mathbb{1}(D = d) \{Y - \hat{h}(W, d, X)\} + \hat{h}(W, d, X) \right]$$

- Asymptotically normal for  $E[Y^{(d)}]$  under certain conditions
- One can use nonparametric approaches [14, 19, 33]<sup>2</sup>

<sup>1</sup>Can use existing instrumental variable software packages; R: gmm, ivreg, sem, ivpack, AER; SAS: SYSLIN; Stata: ivregress, ivreg, ivreg2

<sup>2</sup>Singh, Sahani, Gretton (2019). Kernel IV regression. *NeurIPS*, Mastouri et al. (2021). Proximal causal learning with kernels: Two-stage estimation and moment restriction. *ICML*. Ghassami, Ying, Shpitser, TT (2022). Minimax kernel machine learning for a class of doubly robust functionals with application to proximal causal inference. *AISTAT*

# Comparison to Standard Causal Inference

- Standard Causal Inference

- Usage of covariates  $L$

Treat equally as confounders

$$L = \{\text{age, sex, hypertension, dermatology visit}\}$$

- Proximal Causal Inference

Divide into three buckets

$$X = \{\text{age, sex}\}, W = \{\text{hypertension}\}, Z = \{\text{dermatology visit}\}$$

- Identifying formula for  $E[Y^{(d)}]$

$$\text{IPW} \quad E[p(d, L) \mathbb{1}(D = d) Y]$$

$$\text{g-formula} \quad E[\mu(d, L)]$$

$$\text{AIPW} \quad E \left[ \begin{array}{l} p(d, L) \mathbb{1}(D = d) \{Y - \mu(d, L)\} \\ + \mu(d, X) \end{array} \right]$$

$$\text{Proximal IPW} \quad E[q(Z, d, X) \mathbb{1}(D = d) Y]$$

$$\text{Proximal g-formula} \quad E[h(W, d, X)]$$

$$\text{Proximal AIPW} \quad E \left[ \begin{array}{l} q(Z, d, X) \mathbb{1}(D = d) \{Y - h(W, d, X)\} \\ + h(W, d, X) \end{array} \right]$$

- Estimators for  $E[Y^{(d)}]$

$$\text{OLS} \quad \text{lm}(Y \sim D + L)$$

$$\begin{aligned} \text{2 Stage} \quad & \text{lm}(Y \sim D + X + \widehat{W}) \\ & \widehat{W} \leftarrow \text{lm}(W \sim D + X + Z) \end{aligned}$$

$$\text{AIPW} \quad \text{aver} \left[ \begin{array}{l} \hat{p}(d, L) \mathbb{1}(D = d) \{Y - \hat{\mu}(d, L)\} \\ + \hat{\mu}(a, X) \end{array} \right]$$

$$\text{Proximal AIPW} \quad \text{aver} \left[ \begin{array}{l} \hat{q}(Z, d, X) \mathbb{1}(D = d) \{Y - \hat{h}(W, d, X)\} \\ + \hat{h}(W, d, X) \end{array} \right]$$

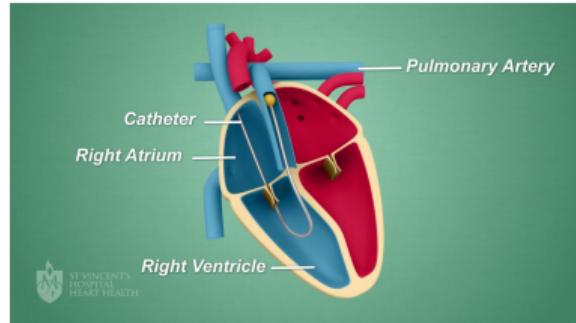
## Application to the SUPPORT study

- Right heart catheterization (RHC) procedure

A catheter (a thin, flexible tube) is inserted into a vein and guided through the blood vessels to reach the right side of the heart and pulmonary artery

Performed to measure blood flow and heart pressure

Common belief: measurements from the RHC can help doctors assess heart function and diagnose conditions



source: St Vincent's Hospital

- The Study to Understand Prognoses and Preferences for Outcomes and Risks of Treatments (SUPPORT) [4]<sup>1</sup>

Evaluate the effectiveness of RHC among patients admitted to the intensive care unit (ICU)

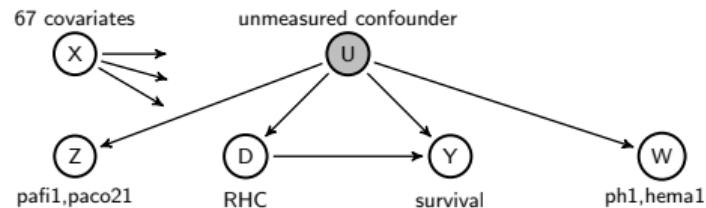
2184 patients with RHC, 3551 without RHC

- The SUPPORT study found that RHC is harmful (associated with a lower chance of survival)

<sup>1</sup> Connors et al. (1996). The effectiveness of right heart catheterization in the initial care of critically ill patients. *Jama*

## Application to the SUPPORT study

- $Y$  = the number of days between admission and death or censoring at 30 days
- 71 covariates
  - demographics, comorbidity, vital signs, functional status
  - Physiological status measured from a blood test during the first day in the ICU (proxy candidates)
- Determine  $Z$  and  $W$ 
  - (pafi1, paco21, ph1, hema1) are strongly correlated with  $D$  and  $Y$
  - Variables most strongly associated with  $D$  and  $Y$  were selected as  $Z$  and  $W$ , respectively
  - $Z = (\text{pafi1}, \text{paco21})$ ;  $W = (\text{ph1}, \text{hema1})$



## Application to the SUPPORT study

		Standard Causal Inference	Proximal Causal Inference
Estimator	OLS / 2 Stage [37] <sup>2</sup>	-1.25 (-1.80,-0.70)	-1.80 (-2.64,-0.96)
	AIPW [5] <sup>2</sup>	-1.17 (-1.79,-0.55)	-1.66 (-2.50,-0.83)

- Estimation methods show little difference
- Standard and proximal causal inference yield substantially different results
- RHC consistently exhibits a negative impact on survival

<sup>2</sup>TT et al. (2024). An Introduction to Proximal Causal Inference. *Statistical Science*

<sup>2</sup>Cui, Pu, Shi, Miao, TT (2024). Semiparametric proximal causal inference. *JASA*

## Extensions of Proximal Causal Inference to Various Settings

## Extensions

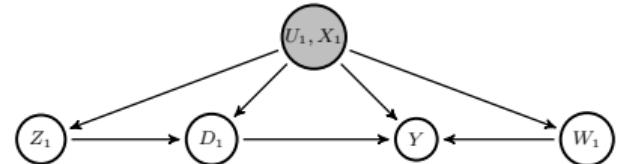
- Longitudinal Settings
- Mediation analysis with hidden confounders
- Mediation analysis and front-door model with hidden mediators
- Interference with homophily driven by hidden factors
- Proximal synthetic controls

## Longitudinal Setting

- Discussed the point exposure setting

$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$



## Longitudinal Setting

- Discussed the point exposure setting

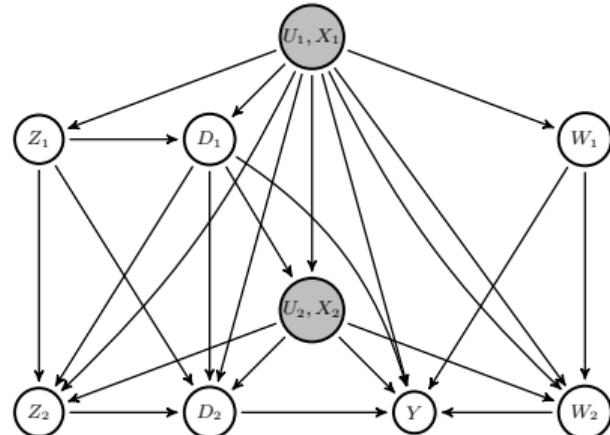
$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$

- Extension to longitudinal settings

$$(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(\vec{d})}) \mid (U_1, X_1)$$

$$(\vec{Z}_2, \vec{D}_2) \perp\!\!\!\perp (\vec{W}_2, Y^{(\vec{d})}) \mid (D_1 = d_1, \vec{U}_2, \vec{X}_2)$$



## Longitudinal Setting

- Discussed the point exposure setting

$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

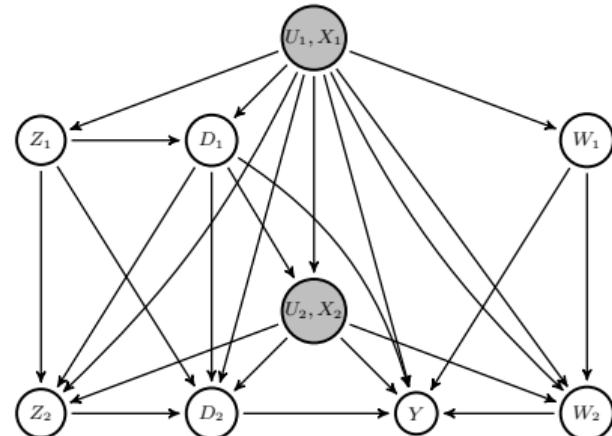
- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$

- Extension to longitudinal settings

$$(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(\vec{d})}) \mid (U_1, X_1)$$

$$(\vec{Z}_2, \vec{D}_2) \perp\!\!\!\perp (\vec{W}_2, Y^{(\vec{d})}) \mid (D_1 = d_1, \vec{U}_2, \vec{X}_2)$$

- Estimand:  $E[Y^{(\vec{d})}]$



## Longitudinal Setting

- Discussed the point exposure setting

$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$

- Extension to longitudinal settings

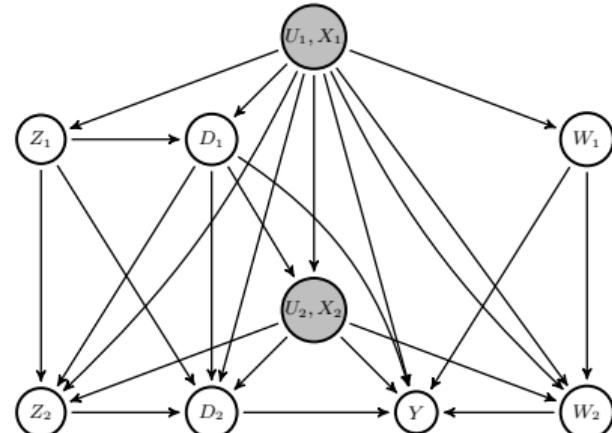
$$(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(\vec{d})}) \mid (U_1, X_1)$$

$$(\vec{Z}_2, \vec{D}_2) \perp\!\!\!\perp (\vec{W}_2, Y^{(\vec{d})}) \mid (D_1 = d_1, \vec{U}_2, \vec{X}_2)$$

- Estimand:  $E[Y^{(\vec{d})}]$

- Nested bridge functions:

$$E[h_2(\vec{D}_2, \vec{W}_2, \vec{X}_2) \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2] = E(Y \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2)$$



## Longitudinal Setting

- Discussed the point exposure setting

$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$

- Extension to longitudinal settings

$$(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(\vec{d})}) \mid (U_1, X_1)$$

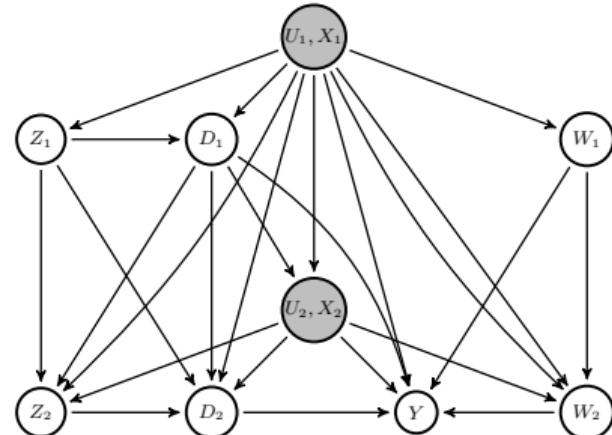
$$(\vec{Z}_2, \vec{D}_2) \perp\!\!\!\perp (\vec{W}_2, Y^{(\vec{d})}) \mid (D_1 = d_1, \vec{U}_2, \vec{X}_2)$$

- Estimand:  $E[Y^{(\vec{d})}]$

- Nested bridge functions:

$$E[h_2(\vec{D}_2, \vec{W}_2, \vec{X}_2) \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2] = E(Y \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2)$$

$$E[h_1(\vec{D}_2, W_1, X_1) \mid \vec{D}_2, Z_1, X_1] = E[h_2(\vec{D}_2, \vec{W}_2, \vec{X}_2) \mid \vec{D}_2, Z_1, X_1]$$



<sup>1</sup>Ying, Miao, Shi, TT (2023). Proximal causal inference for complex longitudinal studies. *JRSSB*

## Longitudinal Setting

- Discussed the point exposure setting

$$Z_1 \perp\!\!\!\perp Y \mid (U_1, D_1, X_1) \text{ and } W_1 \perp\!\!\!\perp (D_1, Z_1) \mid (U_1, X_1)$$

- Suffices to write as  $(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(d)}) \mid (U_1, X_1)$

- Extension to longitudinal settings

$$(Z_1, D_1) \perp\!\!\!\perp (W_1, Y^{(\vec{d})}) \mid (U_1, X_1)$$

$$(\vec{Z}_2, \vec{D}_2) \perp\!\!\!\perp (\vec{W}_2, Y^{(\vec{d})}) \mid (D_1 = d_1, \vec{U}_2, \vec{X}_2)$$

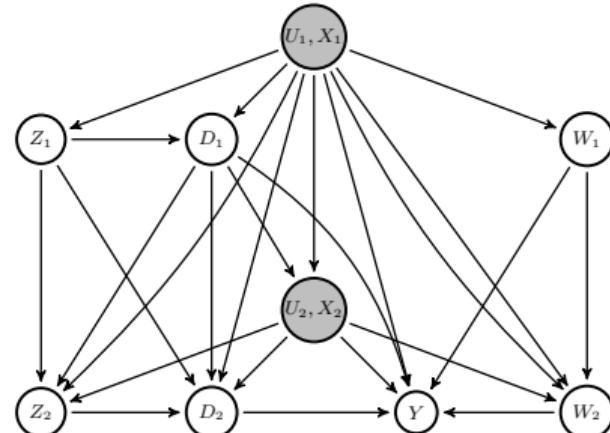
- Estimand:  $E[Y^{(\vec{d})}]$

- Nested bridge functions:

$$E[h_2(\vec{D}_2, \vec{W}_2, \vec{X}_2) \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2] = E(Y \mid \vec{D}_2, \vec{Z}_2, \vec{X}_2)$$

$$E[h_1(\vec{D}_2, W_1, X_1) \mid \vec{D}_2, Z_1, X_1] = E[h_2(\vec{D}_2, \vec{W}_2, \vec{X}_2) \mid \vec{D}_2, Z_1, X_1]$$

- Identification:**  $E[Y^{(\vec{d})}] = E[h_1(\vec{d}, W_1, X_1)]$ ; Proximal AIPW identification is also available [40]<sup>1</sup>



<sup>1</sup>Ying, Miao, Shi, TT (2023). Proximal causal inference for complex longitudinal studies. *JRSSB*

## Longitudinal Setting

- Suppressing  $X$

- Nested bridge functions:

$$E[h_2(\vec{D}_2, \vec{W}_2) | \vec{D}_2, \vec{Z}_2] = E(Y | \vec{D}_2, \vec{Z}_2) \quad \Rightarrow \quad E[g_2(\vec{D}_2, \vec{Z}_2)\{Y - h_2(\vec{D}_2, \vec{W}_2)\}] = 0$$

$$E[h_1(\vec{D}_2, W_1) | \vec{D}_2, Z_1] = E[h_2(\vec{D}_2, \vec{W}_2) | \vec{D}_2, Z_1] \quad \Rightarrow \quad E[g_1(\vec{D}_2, Z_1)\{h_2(\vec{D}_2, \vec{W}_2) - h_1(\vec{D}_2, W_1)\}] = 0$$

- $E[Y^{(\vec{d})}] = E[h_1(\vec{d}, W_1)]$

- Moment equation:

$$(\hat{\theta}_1, \hat{\theta}_2) \leftarrow \text{aver} \begin{bmatrix} g_2(\vec{D}_2, \vec{Z}_2)\{Y - h_2(\vec{D}_2, \vec{W}_2; \theta_2)\} \\ g_1(\vec{D}_2, Z_1)\{h_2(\vec{D}_2, \vec{W}_2; \theta_2) - h_1(\vec{D}_2, W_1; \theta_1)\} \end{bmatrix} = 0$$

- $\hat{E}[Y^{(\vec{d})}] = \text{aver}[h_1(\vec{d}, W_1; \hat{\theta}_1)]$

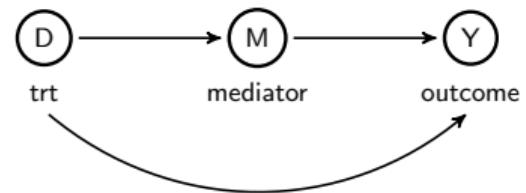
- Asymptotically normal under certain conditions

## Extensions

- Longitudinal Settings
- Mediation analysis with hidden confounders
- Mediation analysis and front-door model with hidden mediators
- Interference with homophily driven by hidden factors
- Proximal synthetic controls

## Mediation

- Mediator ( $M$ ) explains part of the relationship between  $D$  and  $Y$



## Mediation

- Mediator ( $M$ ) explains part of the relationship between  $D$  and  $Y$

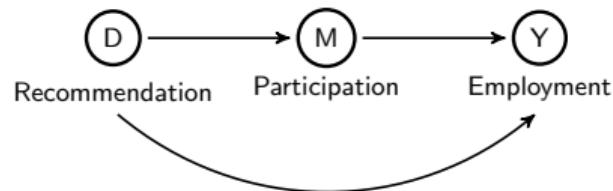
- Example: Job Corps program

Education and career training program administered by the U.S. Department of Labor

$D$ : Randomized, recommendation to participate in Job Corps

$M$ : Actual participation

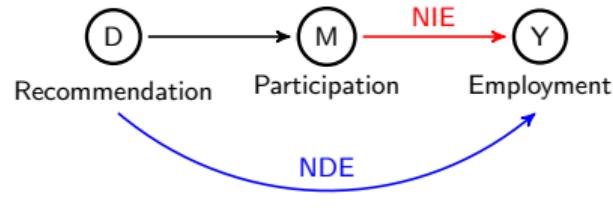
$Y$ : Employment



## Natural Direct/Indirect Effects

- The average treatment effect of  $D$  on  $Y$

$$\begin{aligned} & E[Y^{(1)} - Y^{(0)}] \\ &= E[Y^{(1,M^{(1)})} - Y^{(1,M^{(0)})}] \\ &= \underbrace{E[Y^{(1,M^{(1)})} - Y^{(1,M^{(0)})}]}_{\text{Natural Indirect Eff.}} + \underbrace{E[Y^{(1,M^{(0)})} - Y^{(0,M^{(0)})}]}_{\text{Natural Direct Eff.}} \end{aligned}$$



- Interpretation:

**NDE:** Effect of the JC recommendation ( $D$ ) on employment ( $Y$ ) that is not mediated by the participation ( $M$ )

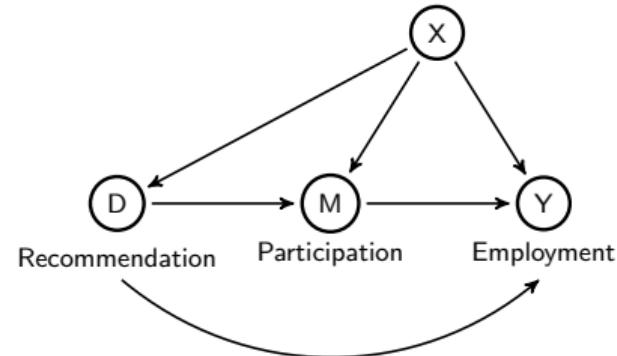
**NIE:** Effect of the JC recommendation ( $D$ ) on employment ( $Y$ ) that is mediated by the participation ( $M$ )

- $E[Y^{(1,M^{(0)})}] = \text{mediation functional}$

## Review: Identification of the Mediation Functional

- Under no unmeasured confounding (and certain assumptions),  
the mediation functional is identified by the **mediation formula** [25]<sup>1</sup>:

$$E[Y^{(1, M^{(0)})}] = E_X \left[ \sum_m E(Y \mid D = 1, M = m, X) \Pr(M = m \mid D = 0, X) \right]$$

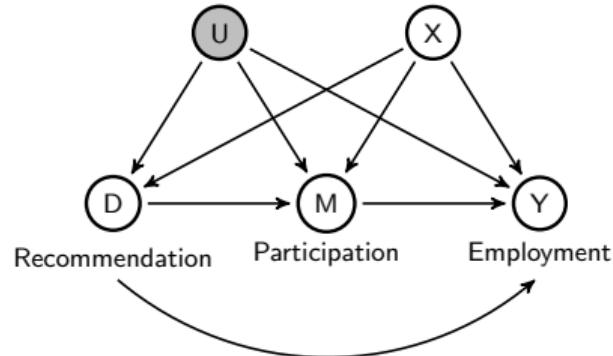


<sup>1</sup>Pearl (2001). Direct and indirect effects. *UAI*

## Review: Identification of the Mediation Functional

- Under no unmeasured confounding (and certain assumptions),  
the mediation functional is identified by the **mediation formula** [25]<sup>1</sup>:

$$E[Y^{(1, M^{(0)})}] = E_X \left[ \sum_m E(Y | D = 1, M = m, X) \Pr(M = m | D = 0, X) \right]$$



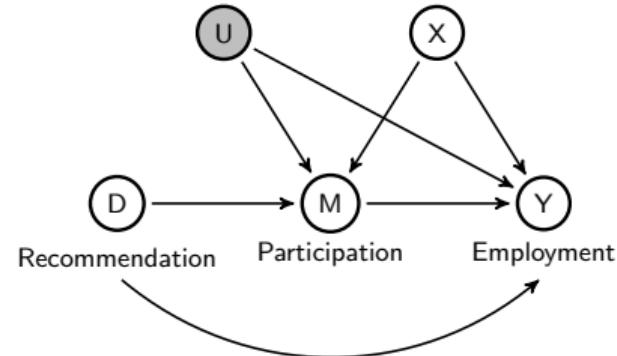
- In the presence of  $U$ , the mediation formula is no longer valid

<sup>1</sup>Pearl (2001). Direct and indirect effects. *UAI*

## Review: Identification of the Mediation Functional

- Under no unmeasured confounding (and certain assumptions),  
the mediation functional is identified by the **mediation formula** [25]<sup>1</sup>:

$$E[Y^{(1, M^{(0)})}] = E_X \left[ \sum_m E(Y | D = 1, M = m, X) \Pr(M = m | D = 0, X) \right]$$



- In the presence of  $U$ , the mediation formula is no longer valid
- Even in randomized trials, common causes of  $M$  and  $Y$  may exist

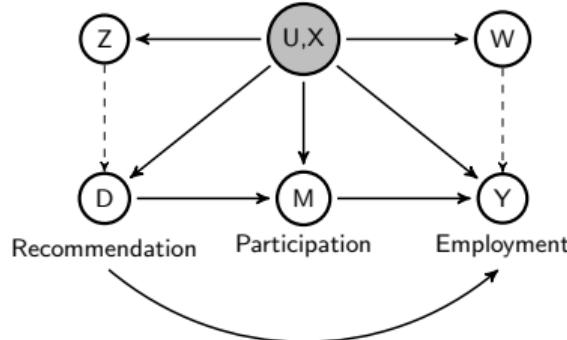
<sup>1</sup>Pearl (2001). Direct and indirect effects. *UAI*

## Proximal Identification of the Mediation Functional

- Suppose that there exist  $Z$  and  $W$  satisfying

$$Z \perp\!\!\!\perp (Y, M) \mid (U, X, D)$$

$$W \perp\!\!\!\perp (D, M, Z) \mid (U, X)$$



- Define two outcome confounding bridge functions

$$\mathbb{E}[h_1(W, M, X) \mid Z, D = 1, M, X] = \mathbb{E}(Y \mid Z, D = 1, M, X)$$

$$\mathbb{E}[h_0(W, X) \mid Z, D = 0, X] = \mathbb{E}[h_1(W, M, X) \mid Z, D = 0, X]$$

- Proximal mediation formula [8]<sup>1</sup>

$$\mathbb{E}[Y^{(1, M^{(0)})}] = \mathbb{E}[h_0(W, X)]$$

- An asymptotically normal estimator can be constructed

<sup>1</sup>Dukes, Shpitser, TT (2023). Proximal mediation analysis. *Biometrika*

## Proximal Identification of the Mediation Functional

- Suppose that there exist  $Z$  and  $W$  satisfying

$$Z \perp\!\!\!\perp Y \mid (U, X, D)$$

$$W \perp\!\!\!\perp (D, M, Z) \mid (U, X)$$

$D$  is randomized

- Define two outcome confounding bridge functions

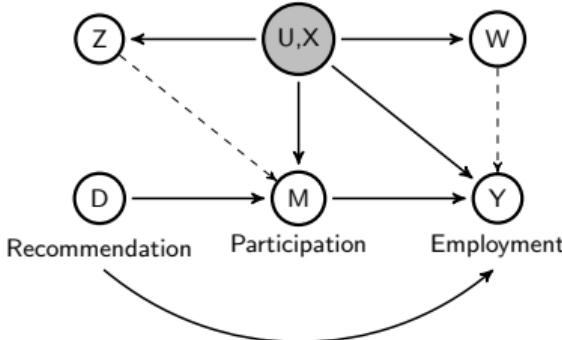
$$\mathbb{E}[h_1(W, M, X) \mid Z, D = 1, M, X] = \mathbb{E}(Y \mid Z, D = 1, M, X)$$

$$\mathbb{E}[h_0(W, X) \mid Z, D = 0, X] = \mathbb{E}[h_1(W, M, X) \mid Z, D = 0, X]$$

- Proximal mediation formula [8]<sup>1</sup>

$$\mathbb{E}[Y^{(1, M^{(0)})}] = \mathbb{E}[h_0(W, X)]$$

- An asymptotically normal estimator can be constructed



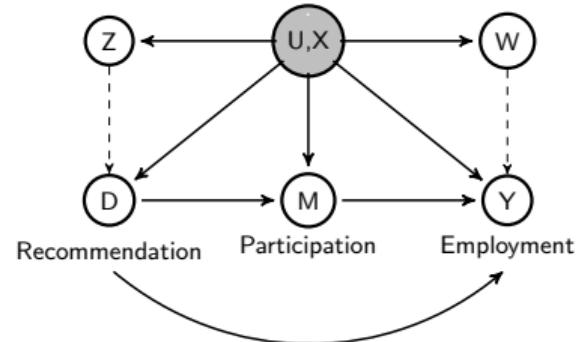
<sup>1</sup>Dukes, Shpitser, TT (2023). Proximal mediation analysis. *Biometrika*

## Extensions

- Longitudinal Settings
- Mediation analysis with hidden confounders
- Mediation analysis and front-door model with hidden mediators
- Interference with homophily driven by hidden factors
- Proximal synthetic controls

## Hidden Mediators

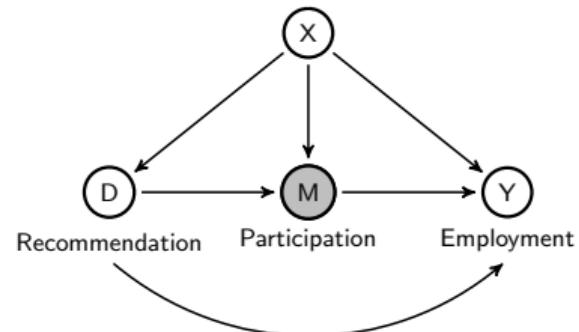
- So far, we allow for  $U$  as long as certain proxies are available



## Hidden Mediators

- So far, we allow for  $U$  as long as certain proxies are available
- What if there is no  $U$ , but  $M$  is not observed?

Respondents may inaccurately report their true participation  $M$  due to pressure to comply, confusion with other programs, concerns about not receiving an incentive, etc



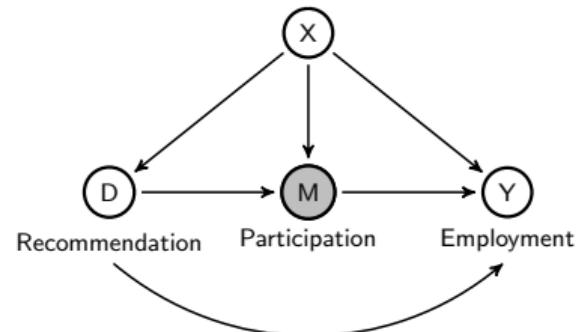
## Hidden Mediators

- So far, we allow for  $U$  as long as certain proxies are available

- What if there is no  $U$ , but  $M$  is not observed?

Respondents may inaccurately report their true participation  $M$  due to pressure to comply, confusion with other programs, concerns about not receiving an incentive, etc

- Idea: if we have proxies of  $M$ , we can identify  $E[Y^{(1, M^{(0)})}]$



# Identification of the Mediation Functional with Hidden Mediators

- Suppose that there exist  $Z$  and  $W$  satisfying

$$Z \perp\!\!\!\perp Y \mid (X, D, M) \quad W \perp\!\!\!\perp (D, Z) \mid (X, M)$$

- Define an outcome confounding bridge functions

$$\text{E}[h(W, X) \mid Z, D = 1, X] = \text{E}(Y \mid Z, D = 1, X)$$

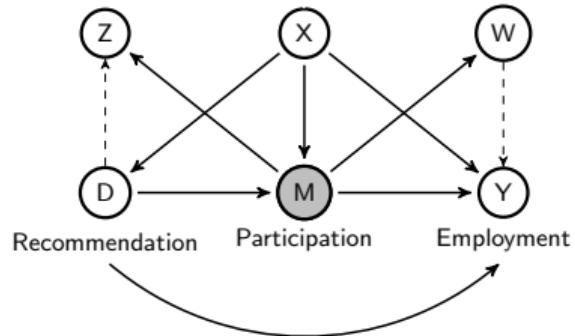
- The mediation functional is identified [13]<sup>1</sup>

$$\text{E}[Y^{(1, M^{(0)})}] = \text{E}_X \left[ \sum_w h(w, X) \Pr(W = w \mid D = 0, X) \right]$$

Original mediation formula:

$$\text{E}[Y^{(1, M^{(0)})}] = \text{E}_X \left[ \sum_m \text{E}(Y \mid D = 1, M = m, X) \Pr(M = m \mid D = 0, X) \right]$$

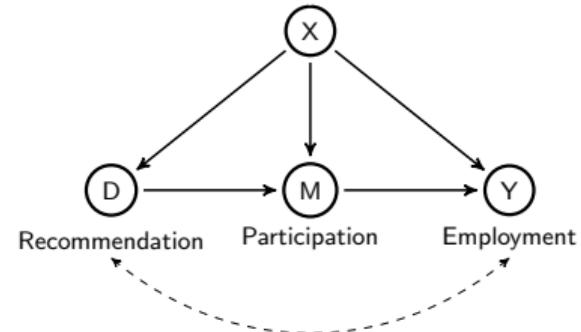
- An asymptotically normal estimator can be constructed



<sup>1</sup> Ghassami, Yang, Shpitser, TT (2024). Causal inference with hidden mediators. *Biometrika*

## Front-Door Model

- Suppose  $Y^{(d,m)} = Y^{(m)}$  (exclusion restriction)  
and allow for unmeasured confounding between  $D$  and  $Y$



- Front-door Model [24]<sup>1</sup>:  
the ATE can be identified even though  $D-Y$  are confounded

$$E[Y^{(0)}] = E[(1 - D)Y] + E_X \left[ \sum_m E(Y | D = 1, M = m, X) \Pr(M = m | D = 0, X) \Pr(D = 1 | X) \right]$$

- What if  $M$  is hidden?

<sup>1</sup>Pearl (2000). *Causality: Models, Reasoning, and Inference*

## Identification of the Front-Door Model with Hidden Mediators

- Suppose that there exist  $Z$  and  $W$  satisfying

$$Z \perp\!\!\!\perp Y \mid (X, D, M) \quad W \perp\!\!\!\perp (D, Z) \mid (X, M)$$

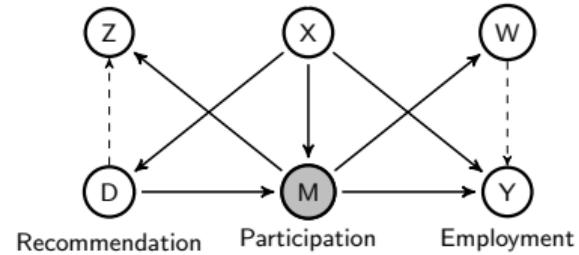
- Define an outcome confounding bridge functions

$$\text{E}[h(W, X) \mid Z, D = 0, X] = \text{E}(Y \mid Z, D = 0, X)$$

- The ATE is identified [13]<sup>1</sup>

$$\text{E}[Y^{(0)}] = \text{E}[(1 - D)Y] + \text{E}_X \left[ \sum_w h(w, X) \Pr(W = w \mid D = 1, X) \Pr(D = 0 \mid X) \right]$$

- An asymptotically normal estimator can be constructed



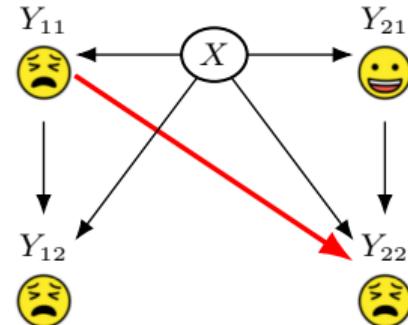
<sup>1</sup> Ghassami, Yang, Shpitser, TT (2024). Causal inference with hidden mediators. *Biometrika*

## Extensions

- Longitudinal Settings
- Mediation analysis with hidden confounders
- Mediation analysis and front-door model with hidden mediators
- Interference with homophily driven by hidden factors
- Proximal synthetic controls

## Causal Inference with Dependence/Interference

- Consider a married couple & flu transmission
- $Y_{jt}$  = unit  $j$  in time  $t$
- $Y_{22}^{(y_{11})}$  = pot. outcome of unit 2 in time 2 provided that  $Y_{11} = y_{11}$
- Consistency:  $Y_{22} = Y_{22}^{(Y_{11})}$
- $X$  = Couple's measured covariate



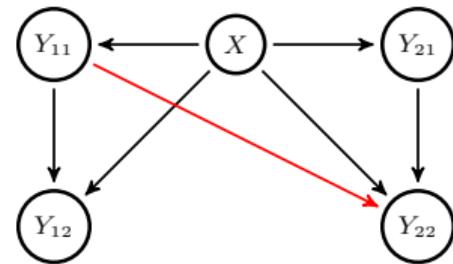
- **Average Causal Peer Effect** =  $E[Y_{22}^{(y_{11})} - Y_{22}^{(y'_{11})}]$
- Can be identified under no unmeasured confounding assumption  $Y_{22}^{(y_{11})} \perp\!\!\!\perp Y_{11} \mid (Y_{21}, X)$

$$\text{ACPE} = E[Y_{22} \mid Y_{11} = y_{11}, Y_{21}, X] - E[Y_{22} \mid Y_{11} = y'_{11}, Y_{21}, X]$$

- What are we missing?

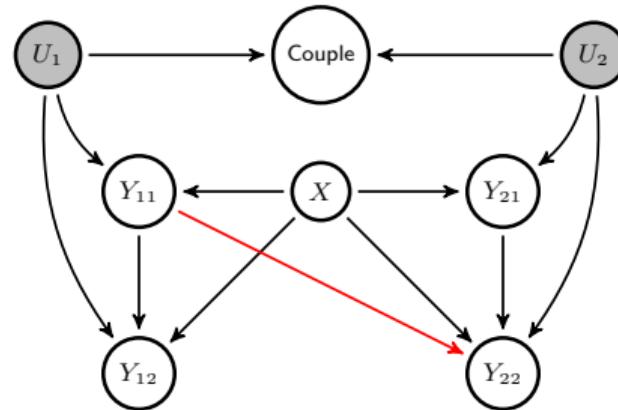
## Homophily

- Why did they become a couple/dyad?
- Measured variables  $X$  (age, gender) may not contain the key reason



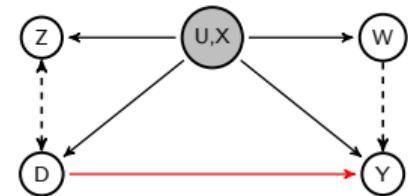
## Homophily

- Why did they become a couple/dyad?
- Measured variables  $X$  (age, gender) may not contain the key reason
- Instead, it may be driven by unmeasured characteristics  $U$   
(cat/dog person, hometown, hobbies, overlapping schools, use of dating apps, etc)
- **Homophily Bias** arises when people become connected due to unobserved characteristics
- Proximal causal inference can be used to address homophily bias



## Connecting the Dyad Example to Proximal Causal Inference

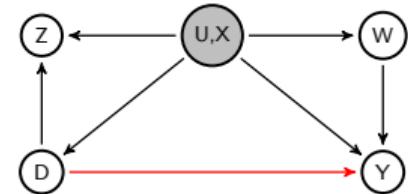
- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D, X)$
- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | (U, X)$



<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Connecting the Dyad Example to Proximal Causal Inference

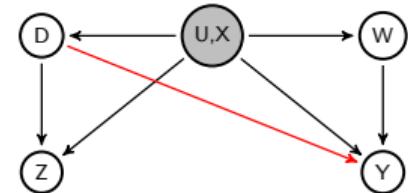
- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D, X)$
- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | (U, X)$



<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSS:B*

## Connecting the Dyad Example to Proximal Causal Inference

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y | (U, D, X)$
- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) | (U, X)$



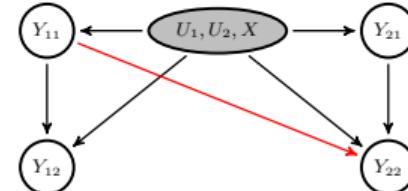
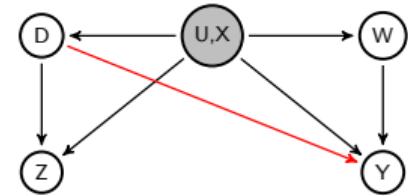
<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Connecting the Dyad Example to Proximal Causal Inference

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y \mid (U, D, X)$
- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) \mid (U, X)$
- Consistent with the dyad model under homophily bias

NCE:  $Y_{12} = \text{peer's outcome at time 2}$

NCO:  $Y_{21} = \text{focal unit's outcome at time 1}$



- Suppose  $h(Y_{21}, Y_{11}, X)$  satisfies

$$E[h(Y_{22}, Y_{11}, X) \mid Y_{12}, Y_{11}, X] = E(Y_{22} \mid Y_{12}, Y_{11}, X)$$

- [9]<sup>1</sup>: The ACPE is identified by

$$\text{ACPE} = E[Y_{22}^{(y_{11})} - Y_{22}^{(y'_{11})}] = E[h(Y_{22}, y_{11}, X) - h(Y_{22}, y'_{11}, X)]$$

- An asymptotically normal estimator can be constructed

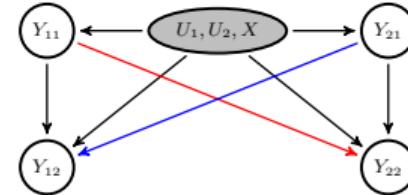
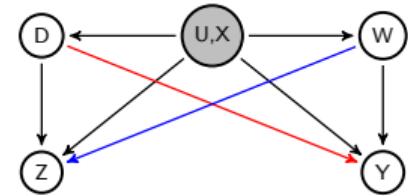
<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Connecting the Dyad Example to Proximal Causal Inference

- $Z$  is an NCE if  $Z \perp\!\!\!\perp Y \mid (U, D, X)$
- $W$  is an NCO if  $W \perp\!\!\!\perp (D, Z) \mid (U, X)$
- Consistent with the dyad model under homophily bias

NCE:  $Y_{12} = \text{peer's outcome at time 2}$

NCO:  $Y_{21} = \text{focal unit's outcome at time 1}$



- Suppose  $h(Y_{21}, Y_{11}, X)$  satisfies

$$E[h(Y_{22}, Y_{11}, X) \mid Y_{12}, Y_{11}, X] = E(Y_{22} \mid Y_{12}, Y_{11}, X)$$

- [9]<sup>1</sup>: The ACPE is identified by

$$\text{ACPE} = E[Y_{22}^{(y_{11})} - Y_{22}^{(y'_{11})}] = E[h(Y_{22}, y_{11}, X) - h(Y_{22}, y'_{11}, X)]$$

- An asymptotically normal estimator can be constructed
- If  $Y_{21} \rightarrow Y_{12}$ , we need another NCE or NCO

<sup>1</sup> Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

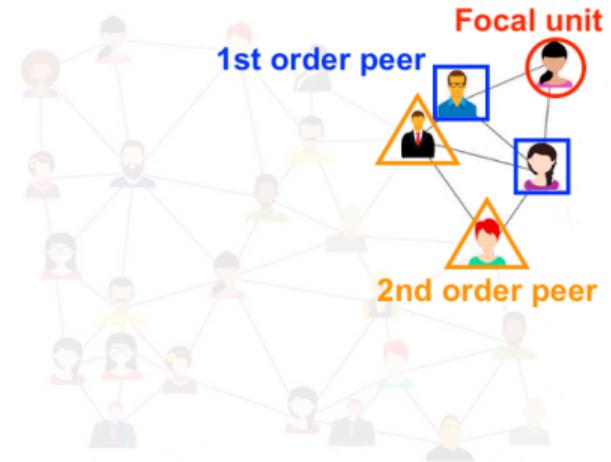
- Consider a social network



<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers

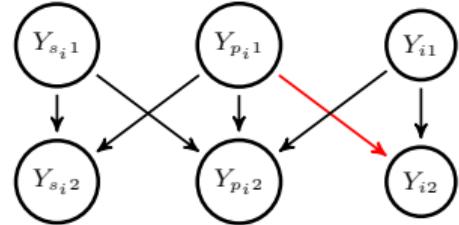


<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers
- Average Causal Peer Effect

$$\text{ACPE} = E \left[ \frac{1}{N} \sum_{i=1}^N \left\{ Y_{i2}^{(y_{p_i 1})} - Y_{i2}^{(y'_{p_i 1})} \right\} \right]$$



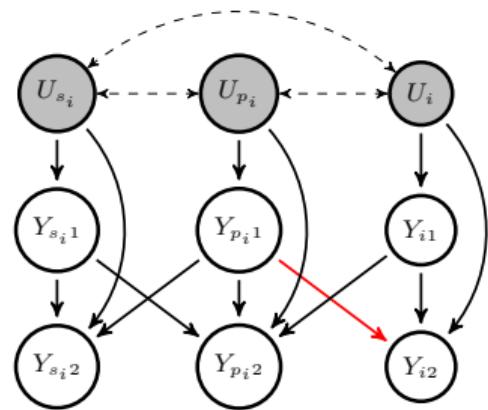
<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers
- **Average Causal Peer Effect**

$$\text{ACPE} = E \left[ \frac{1}{N} \sum_{i=1}^N \left\{ Y_{i2}^{(y_{p_i1})} - Y_{i2}^{(y'_{p_i1})} \right\} \right]$$

- Unmeasured characteristics may exist



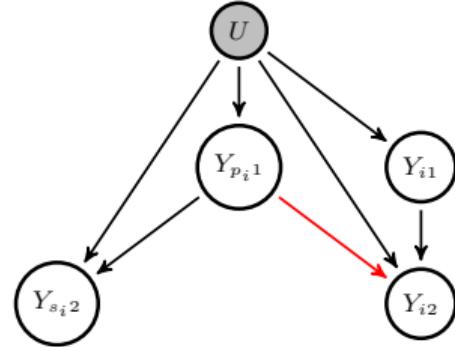
<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers
- Average Causal Peer Effect

$$\text{ACPE} = E \left[ \frac{1}{N} \sum_{i=1}^N \left\{ Y_{i2}^{(y_{p_i1})} - Y_{i2}^{(y'_{p_i1})} \right\} \right]$$

- Unmeasured characteristics may exist



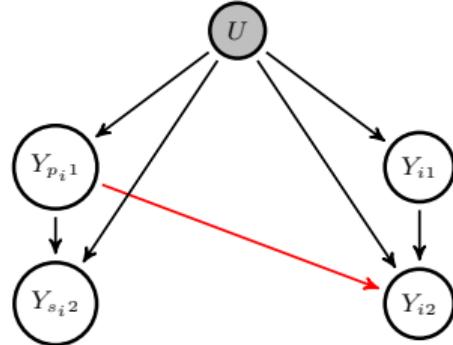
<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers
- **Average Causal Peer Effect**

$$\text{ACPE} = E \left[ \frac{1}{N} \sum_{i=1}^N \left\{ Y_{i2}^{(y_{p_i1})} - Y_{i2}^{(y'_{p_i1})} \right\} \right]$$

- Unmeasured characteristics may exist

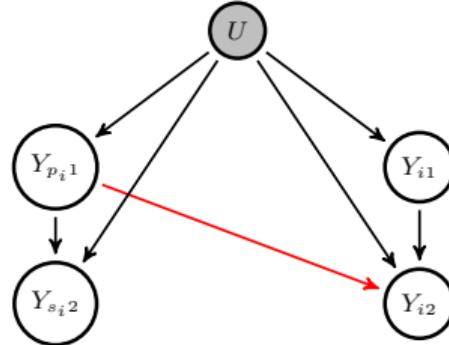


<sup>1</sup>Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

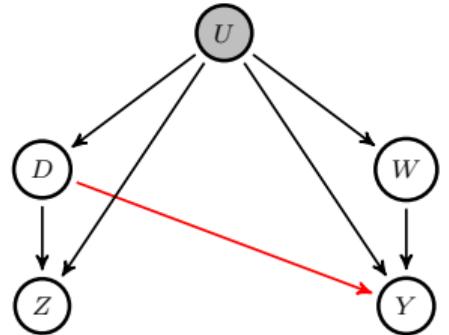
## Extension to Network Data

- Consider a social network
- Given a focal unit, find its 1st and 2nd-order peers
- **Average Causal Peer Effect**

$$\text{ACPE} = E \left[ \frac{1}{N} \sum_{i=1}^N \left\{ Y_{i2}^{(y_{p_i1})} - Y_{i2}^{(y'_{p_i1})} \right\} \right]$$



- Unmeasured characteristics may exist
- Proximal causal inference for network data:
  - NCE =  $Y_{s_i2}$  = 2nd-order peers' outcome at time 2
  - NCO =  $Y_{i1}$  = focal unit's outcome at time 1
- Identification & estimation can be established [9]<sup>1</sup>



<sup>1</sup> Egami, TT (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *JRSSB*

## Extensions

- Longitudinal Settings
- Mediation analysis with hidden confounders
- Mediation analysis and front-door model with hidden mediators
- Interference with homophily driven by hidden factors
- Proximal synthetic controls

## Synthetic Control Approach

- California tobacco control program in Abadie et al. [1]<sup>1</sup>

Treatment:

Proposition 99 (a large-scale tobacco control program) in 1988

Pre- and post-intervention periods:

$T_0 = 19$  (1970-1988) and  $T_1 = 12$  (1989-2000) years

Outcome  $Y_t$ :

Per-capita cigarette sales of California in year  $t$

Outcome  $W_{it}$ :

Per-capita cigarette sales of other 38 states in year  $t$

<sup>1</sup> Abadie, Diamond, Hainmueller (2010). Synthetic control methods for comparative case studies: Estimating the effect of California's tobacco control program. *JASA*

## Synthetic Control Approach

- California tobacco control program in Abadie et al. [1]<sup>1</sup>

Treatment:

Proposition 99 (a large-scale tobacco control program) in 1988

Pre- and post-intervention periods:

$T_0 = 19$  (1970-1988) and  $T_1 = 12$  (1989-2000) years

Outcome  $Y_t$ :

Per-capita cigarette sales of California in year  $t$

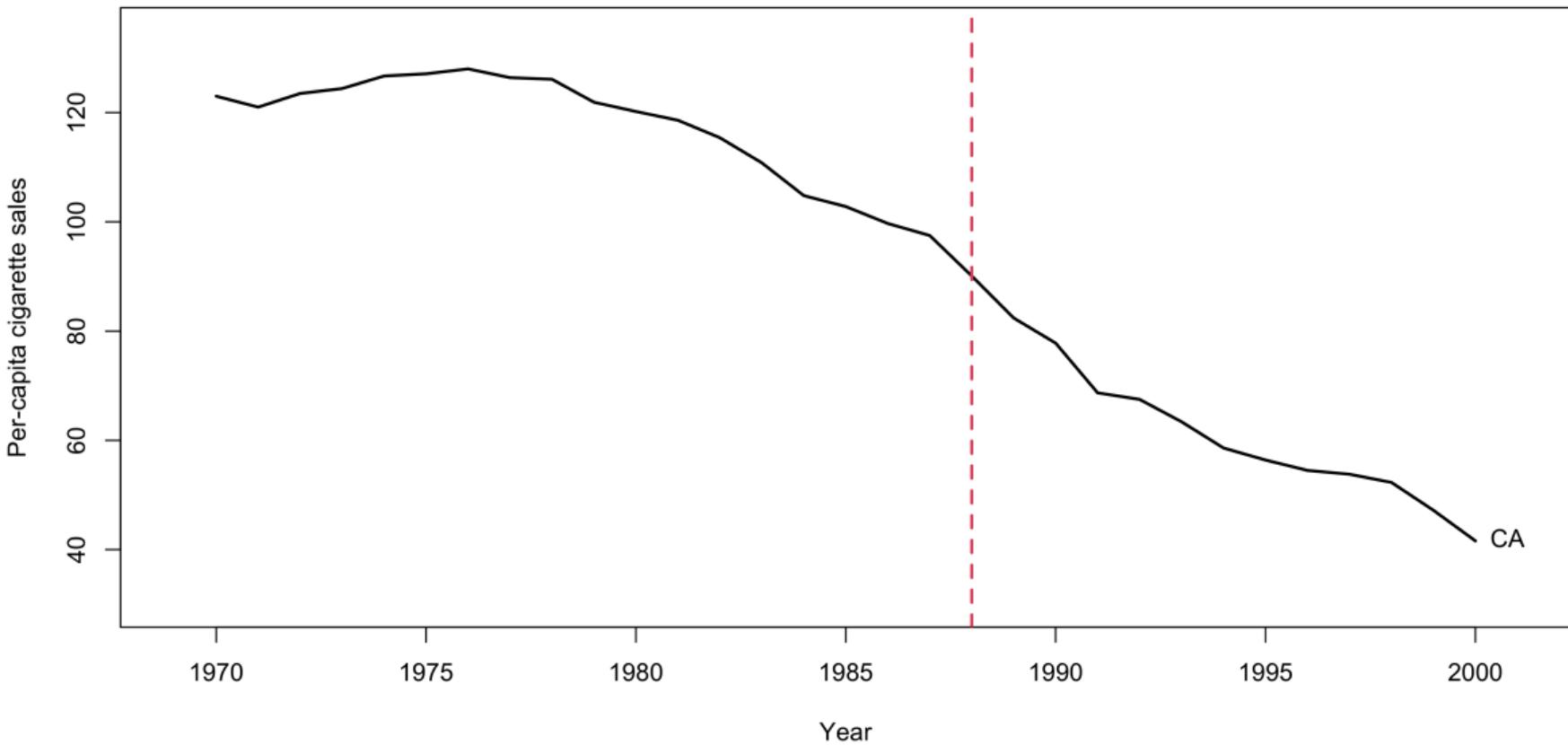
Outcome  $W_{it}$ :

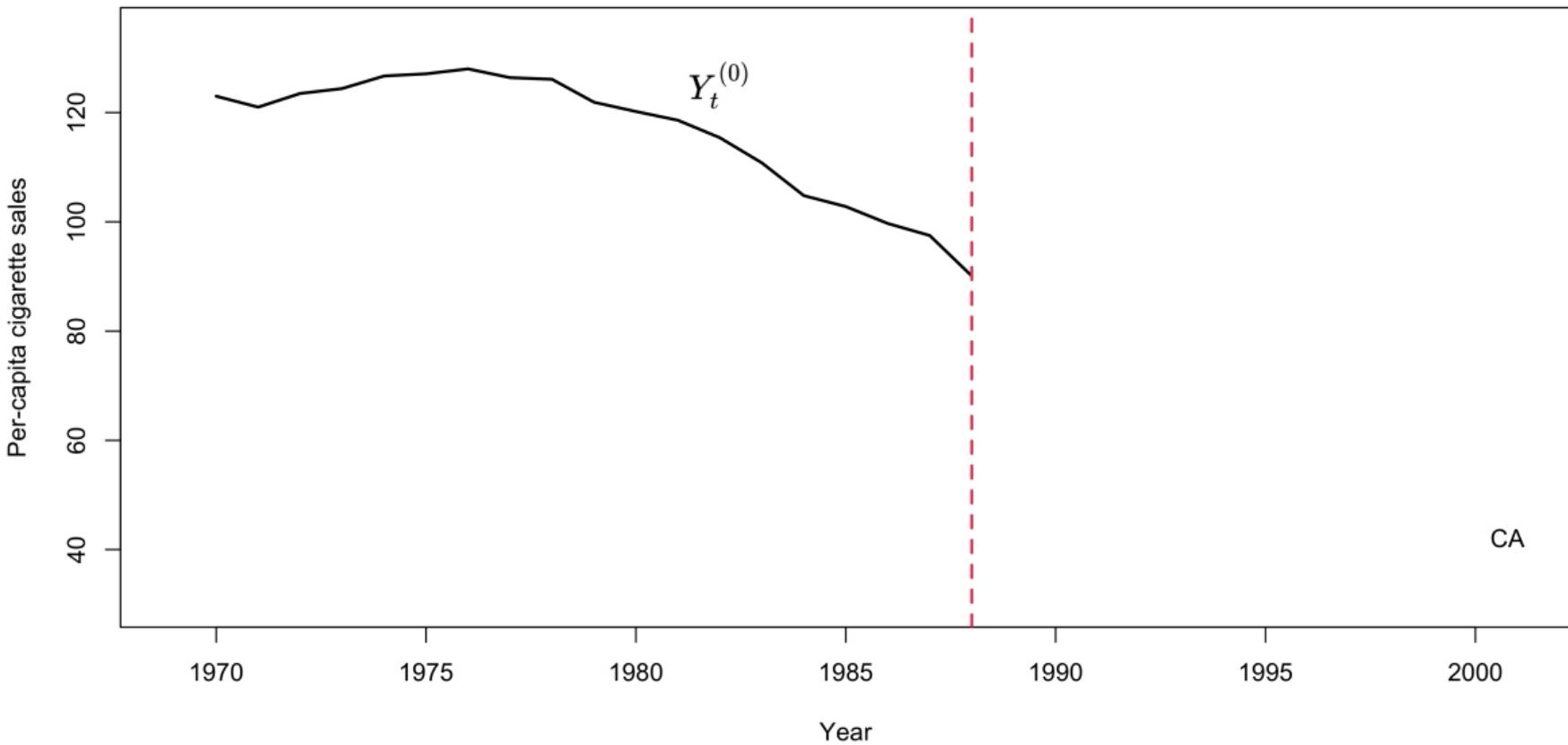
Per-capita cigarette sales of other 38 states in year  $t$

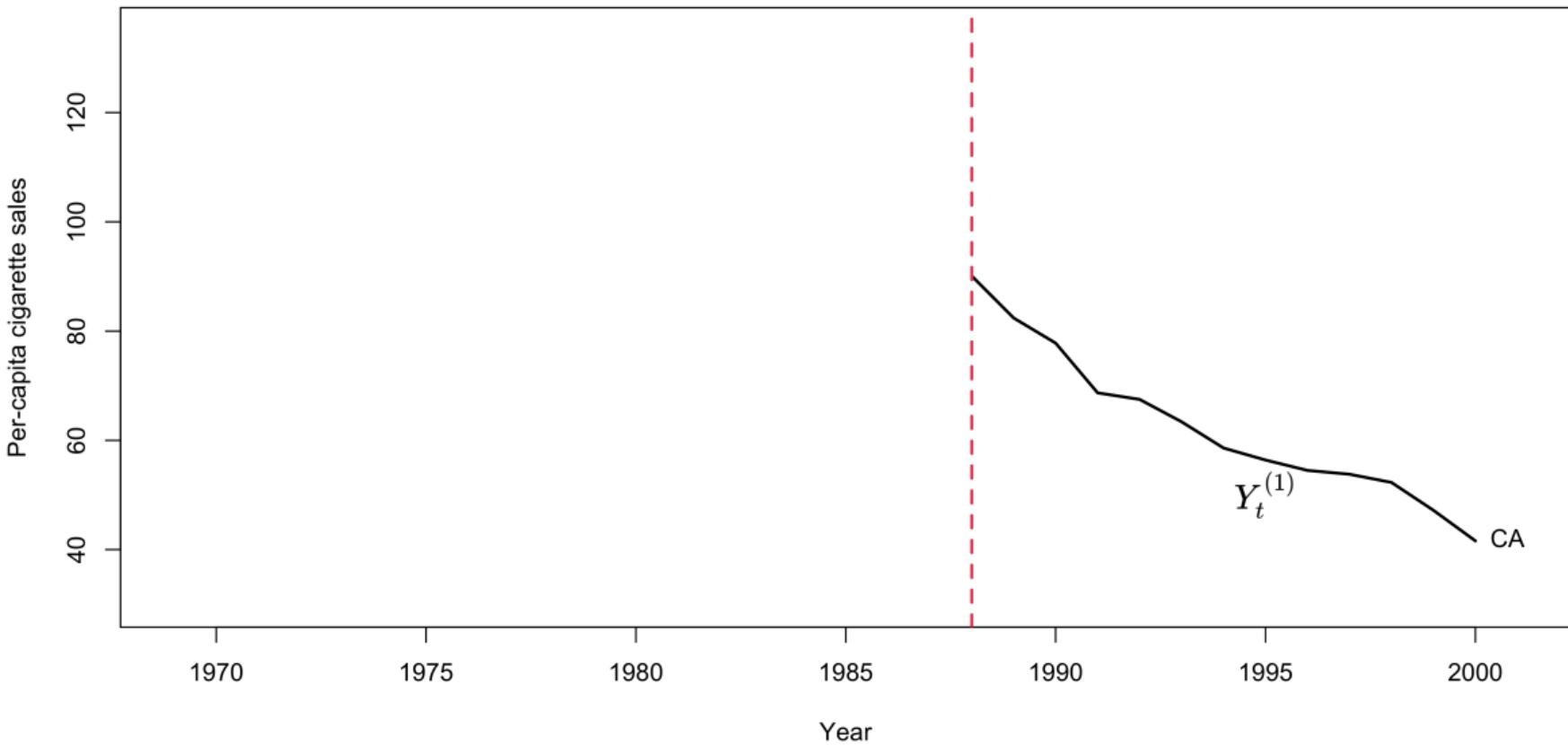
- Estimand: Effect of California's Proposition 99 on California's per-capita cigarette sales

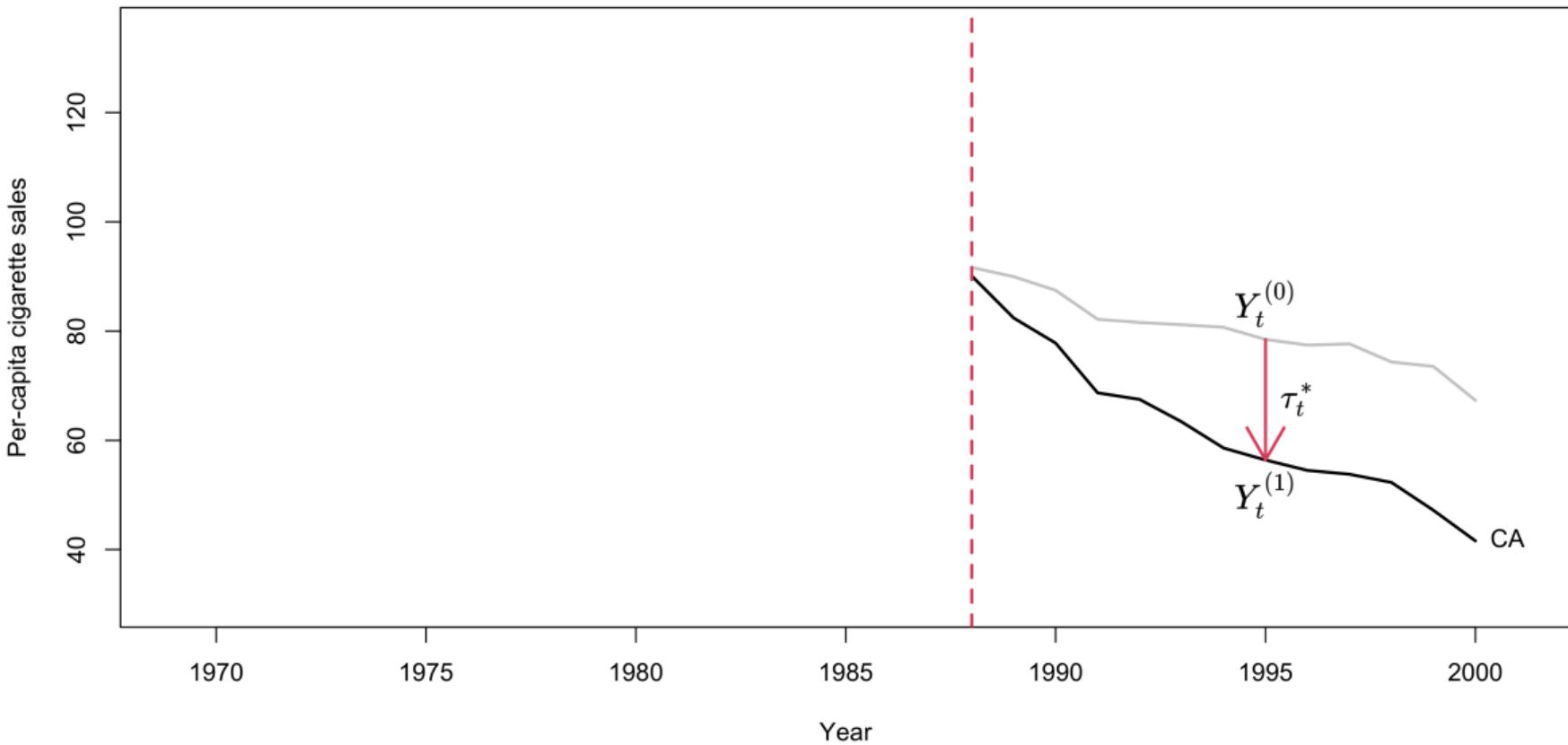
$$\tau_t^* = E[Y_t^{(1)} - Y_t^{(0)}] , \quad 1988 < t$$

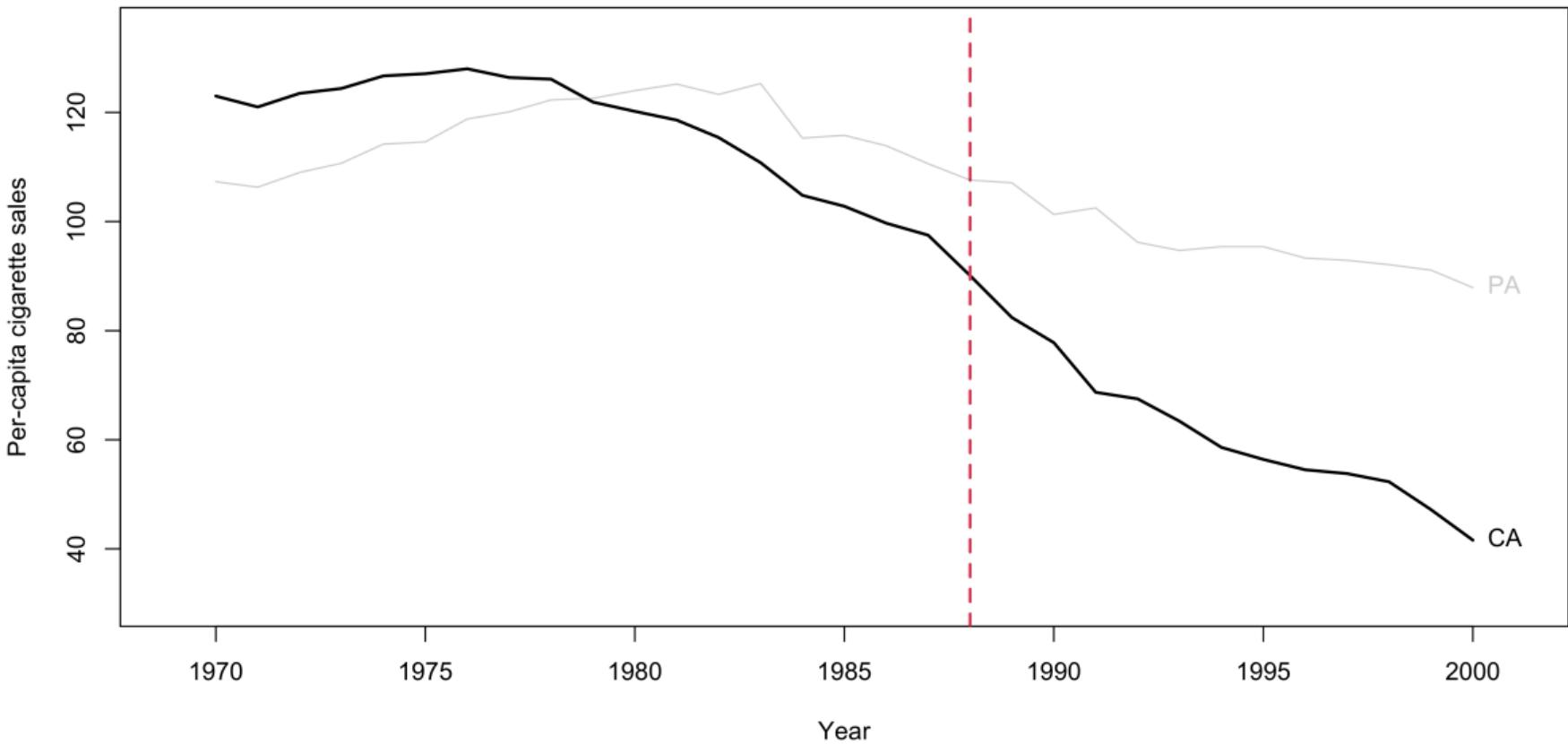
<sup>1</sup> Abadie, Diamond, Hainmueller (2010). Synthetic control methods for comparative case studies: Estimating the effect of California's tobacco control program. *JASA*

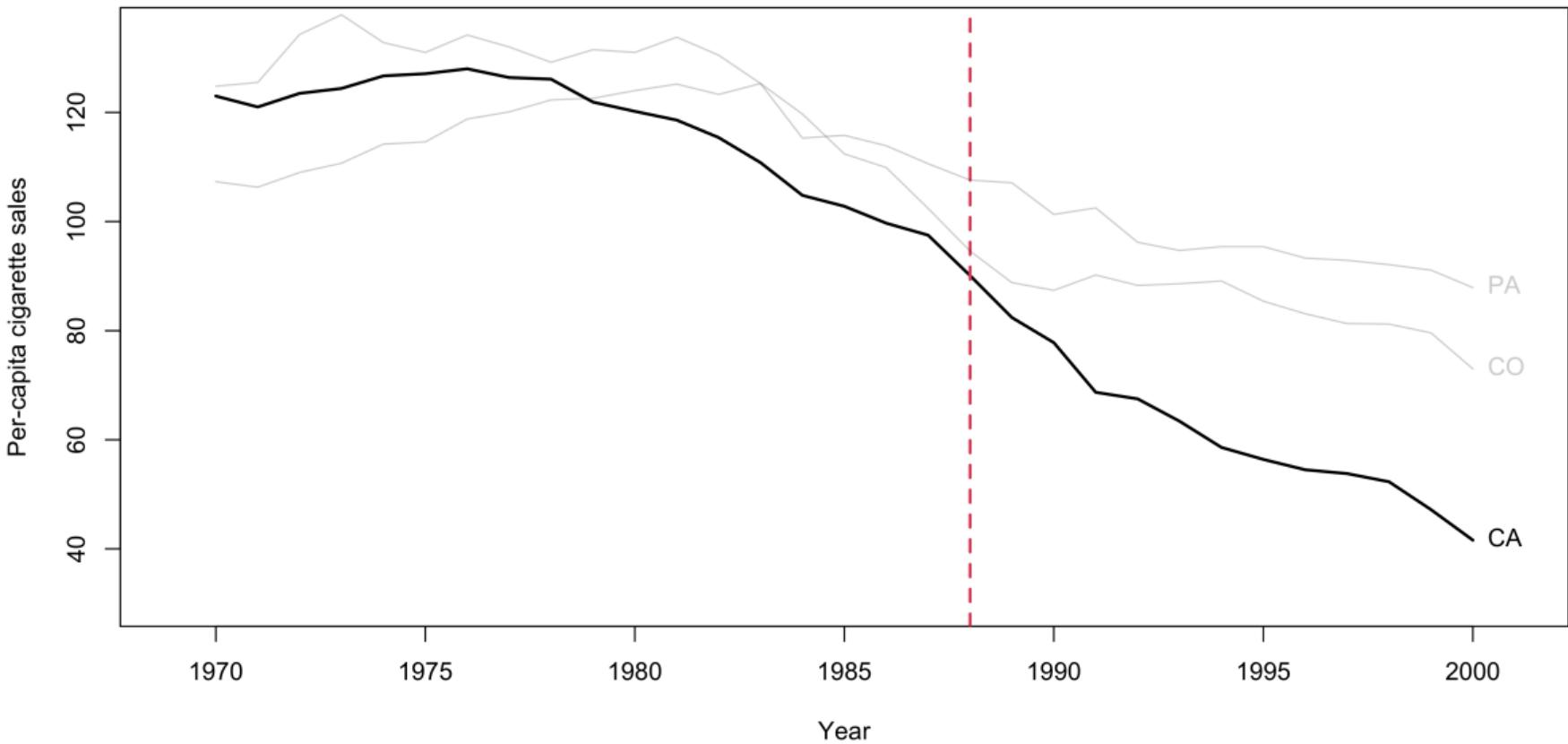


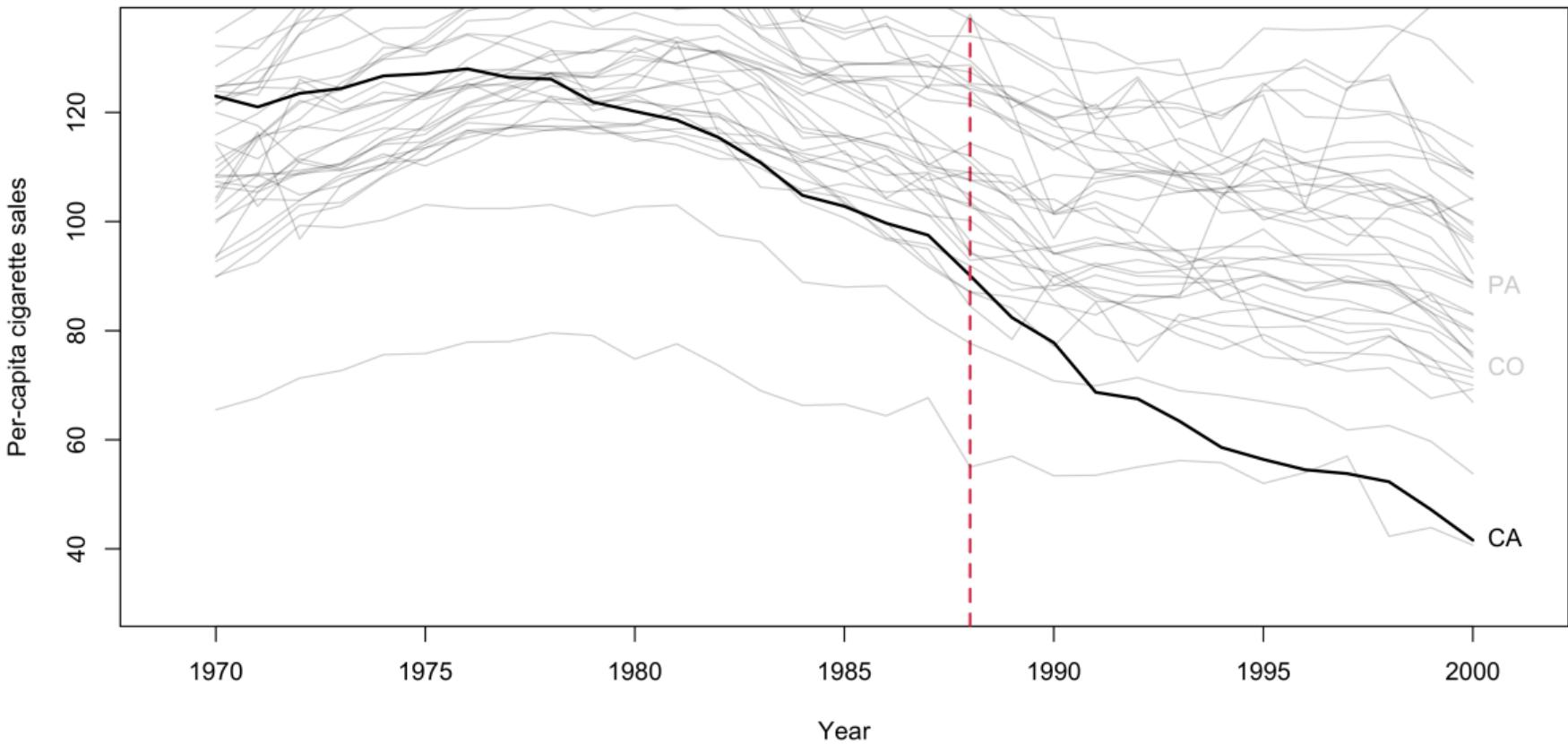


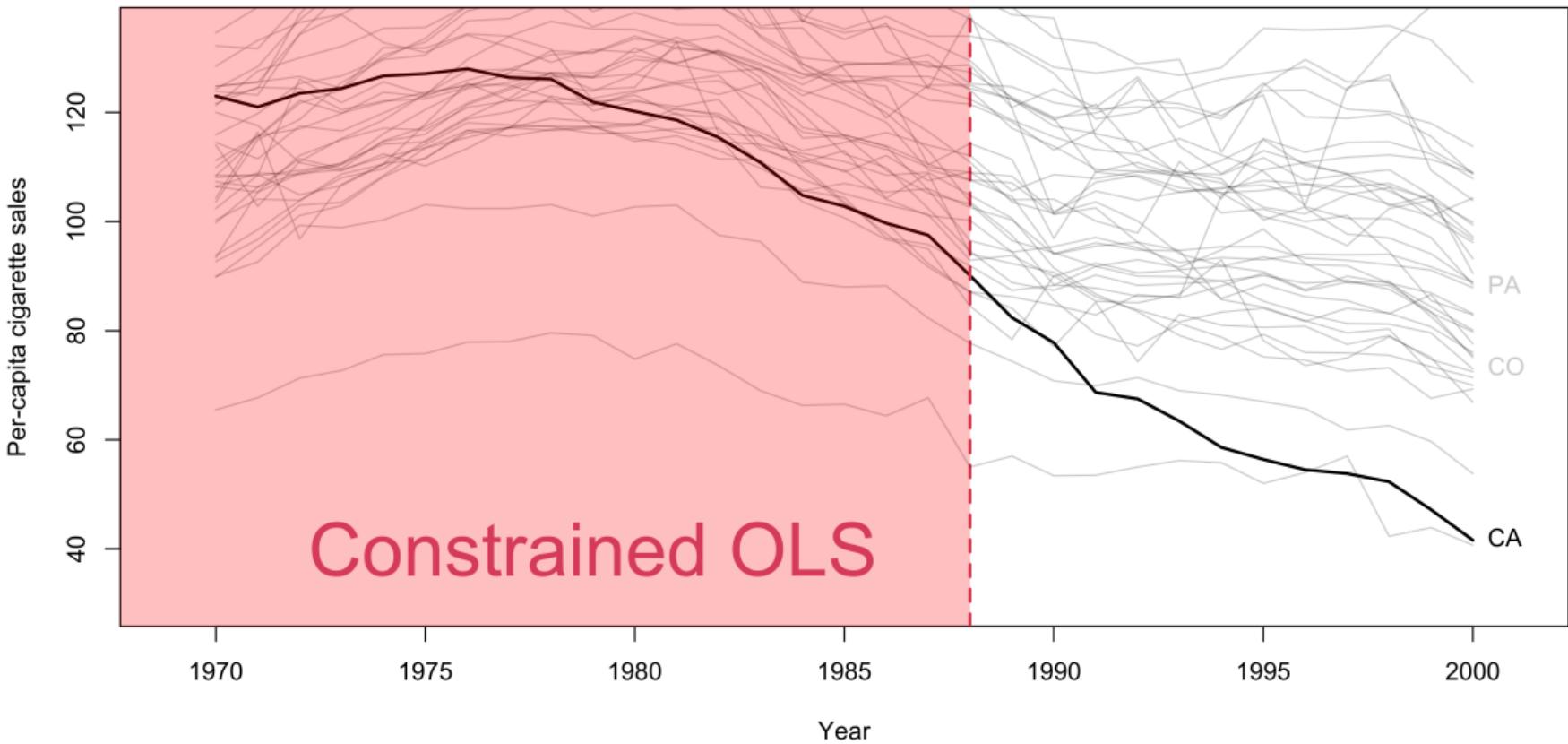


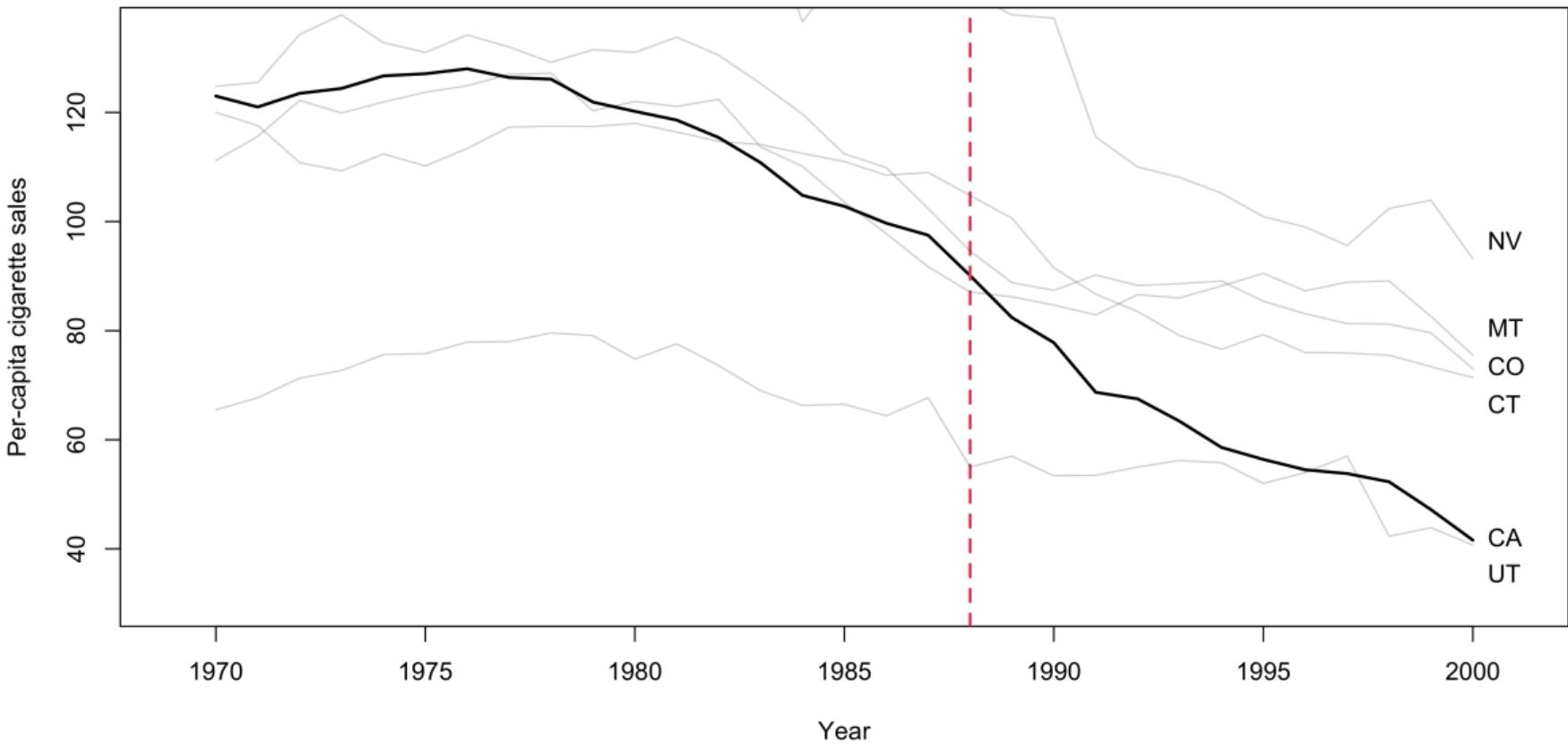


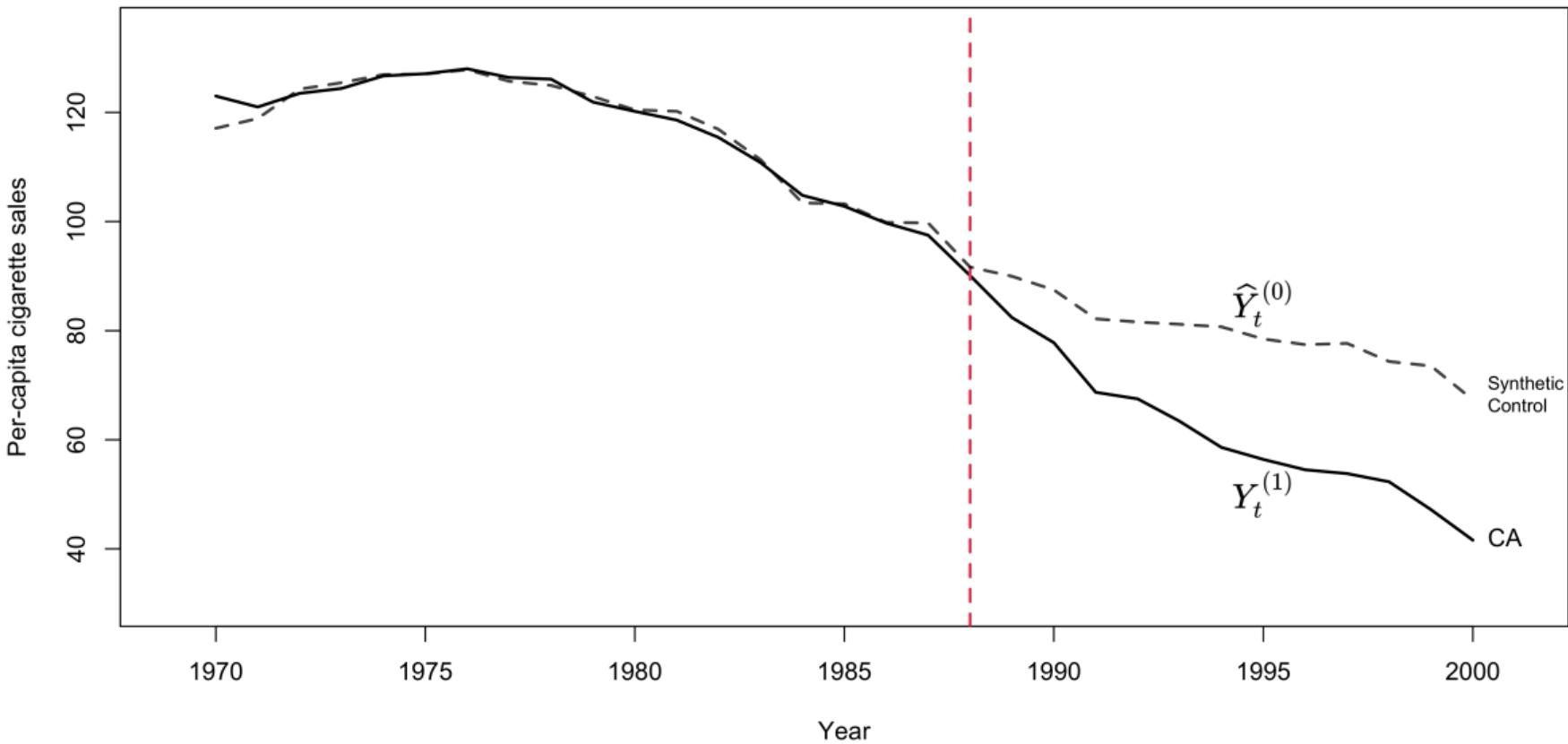


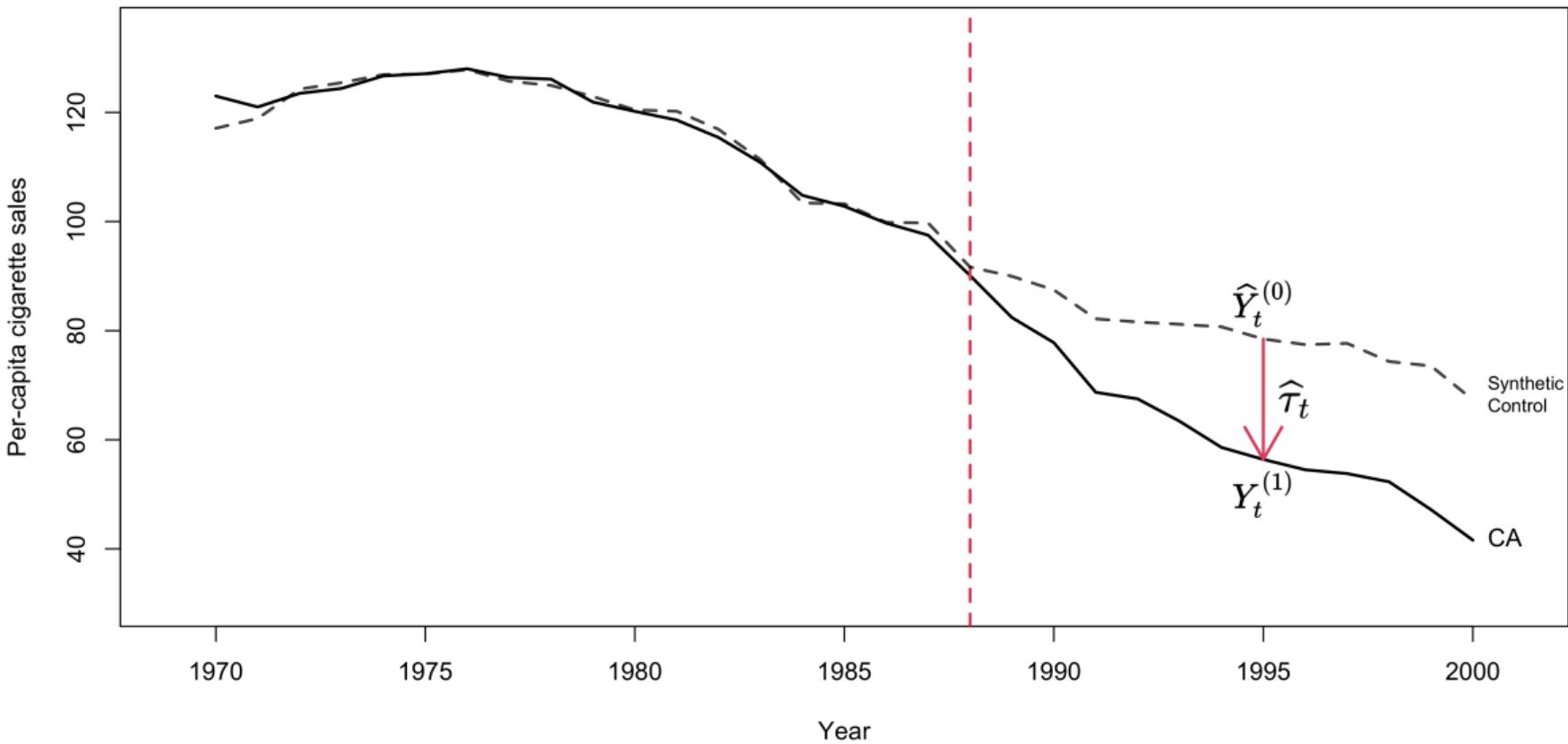












## Synthetic Control Approach

- California tobacco control program in Abadie et al. [1]<sup>1</sup>

Treatment:

Proposition 99 (a large-scale tobacco control program) in 1988

Pre- and post-intervention periods:

$T_0 = 19$  (1970-1988) and  $T_1 = 12$  (1989-2000) years

Outcome  $Y_t$ :

Per-capita cigarette sales of California in year  $t$

Outcome  $W_{it}$ :

Per-capita cigarette sales of other 38 states in year  $t$

- Estimand: Effect of California's Proposition 99 on California's per-capita cigarette sales

$$\tau_t^* = E[Y_t^{(1)} - Y_t^{(0)}], \quad 1988 < t$$

- Synthetic control: find a linear combination of other states

$$\hat{Y}_t^{(0)} \leftarrow \text{constrained OLS}(Y_t \sim W_t, \text{data} = \text{pre-trt data}) \Rightarrow \hat{\tau}_t = Y_t - \hat{Y}_t^{(0)}$$

<sup>1</sup> Abadie, Diamond, Hainmueller (2010). Synthetic control methods for comparative case studies: Estimating the effect of California's tobacco control program. *JASA*

## Interactive Fixed Effects Model

- Consider an IFEM [2]<sup>1</sup>

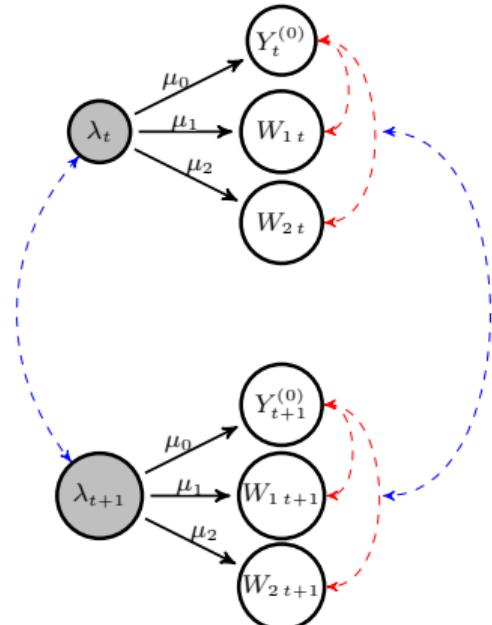
$$Y_t^{(0)} = \mu_0^\top \lambda_t + \epsilon_{0t}$$

$$W_{it} = \mu_i^\top \lambda_t + \epsilon_{it}$$

$\lambda_t$ : latent factors

$\mu_i$ : factor loadings

$\epsilon_{it}$ : error term with  $E(\epsilon_{it} | \lambda_t) = 0$



- Suppose there exists  $\gamma$  such that  $\mu_0 = \sum_i \gamma_i \mu_i$

$$E[Y_t^{(0)} - \sum_i \gamma_i W_{it}] = E[\underbrace{(\mu_0 - \sum_i \gamma_i \mu_i)^\top \lambda_t}_{=0} + \underbrace{\epsilon_{0t} - \sum_i \gamma_i \epsilon_{it}}_{\text{error}}] = 0$$

- $W_t^\top \gamma = \sum_i \gamma_i W_{it}$  is a **synthetic control** with  $E[Y_t^{(0)}] = E[W_t^\top \gamma]$

- $\tau_t = E[Y_t^{(1)} - Y_t^{(0)}] = E[Y_t - W_t^\top \gamma]$  for  $t \in \{\text{post-trt}\}$

<sup>1</sup> Bai (2009). Panel Data Models with Interactive Fixed Effects. *Econometrica*

## Naive OLS

- IFEM with a SC

$$\begin{aligned} Y_t^{(0)} &= \mu_0^\top \lambda_t + \epsilon_{0t} \\ W_{it} &= \mu_i^\top \lambda_t + \epsilon_{it} \end{aligned}, \quad \mu_0 = \sum_i \gamma_i \mu_i$$

- A linear model over the pre-treatment periods

$$Y_t = \sum_i \gamma_i W_{it} + \underbrace{\left( \epsilon_{0t} - \sum_i \gamma_i \epsilon_{it} \right)}_{\text{error}} \Rightarrow Y_t = \sum_i \gamma_i W_{it} + \eta_t$$

- Naive OLS:  $\text{lm}(Y_t \sim 0 + W_{it})$  to estimate  $\gamma$  using  $t \in \{\text{pre-trt}\}$
- The OLS estimator is **inconsistent!** because the regressors  $W_{it}$  and the error are correlated
- An exception: noise-less setting ( $\epsilon_{it} = 0$ ), but nearly impossible in real-world applications

## Proximal Synthetic Control Approach [30]<sup>1</sup>

- IFEM with a SC

$$\begin{aligned} Y_t^{(0)} &= \mu_0^\top \lambda_t + \epsilon_{0t} \\ W_{it} &= \mu_i^\top \lambda_t + \epsilon_{it} \end{aligned}, \quad \mu_0 = \sum_i \gamma_i \mu_i$$

- Suppose we have  $Z_t$  satisfying  $Z_t \perp\!\!\!\perp (Y_t^{(0)}, W_t) | \lambda_t$

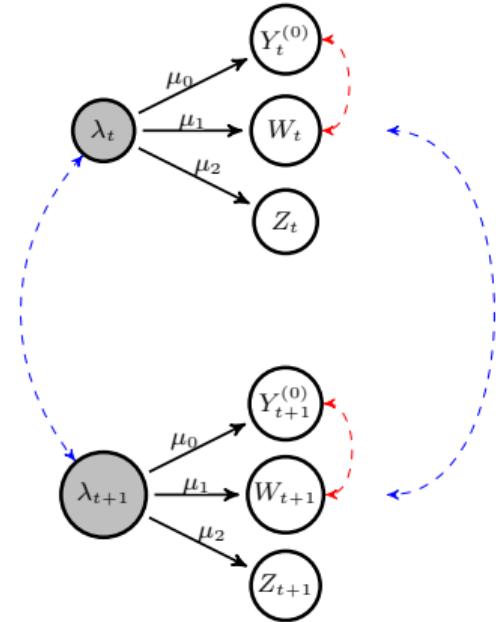
- The SC weight  $\gamma$  can be identified via a moment condition:

$$E[g(Z_t) \times (Y_t - W_t^\top \gamma)] = 0, \quad \forall g(Z_t), \quad t \in \{\text{pre-trt}\}$$

- Estimate  $\gamma$ , and construct a consistent SC estimator:

$$\hat{\gamma} \leftarrow \text{aver}[g(Z_t) \times (Y_t - W_t^\top \gamma)] = 0, \quad t \in \{\text{pre-trt}\}$$

$$\hat{\tau}_t \leftarrow Y_t - W_t^\top \hat{\gamma}, \quad t \in \{\text{post-trt}\}$$



<sup>1</sup>Shi, Li, Miao, Hu, TT (2023). Theory for identification and inference with synthetic controls: A proximal causal inference framework. arXiv

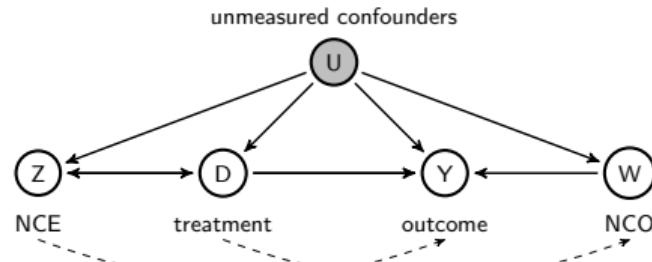
# Summary

- Negative Controls

NCE:  $Z \perp\!\!\!\perp Y | (U, D)$  and  $Z \not\perp\!\!\!\perp U | D$

NCO:  $W \perp\!\!\!\perp (Z, D) | U$  and  $W \not\perp\!\!\!\perp U$

Can be used to detect, reduce, and correct confounding bias



- Proximal Causal Inference

Divide covariates into three buckets

confounder, treatment confounding proxy (NCE), outcome confounding proxy (NCO)

Identification and estimation

GLM → 2-stage regression procedure

bridge functions → proximal IPW, proximal g-formula, proximal AIPW

- Extensions of Proximal Causal Inference to Various Settings

Longitudinal settings

Mediation, front-door formula

Network data

Time series data

# Demonstration using R

[github.com/qkrcks0218/ACIC2025](https://github.com/qkrcks0218/ACIC2025)



## Reference

- [1] Abadie, A., Diamond, A., and Hainmueller, J. (2010). Synthetic control methods for comparative case studies: Estimating the effect of California's tobacco control program. *Journal of the American Statistical Association*, 105(490):493–505.
- [2] Bai, J. (2009). Panel data models with interactive fixed effects. *Econometrica*, 77(4):1229–1279.
- [3] Browning, M. and Crossley, T. (2009). Are two cheap, noisy measures better than one expensive, accurate one? *American Economic Review*, 99(2):99–103.
- [4] Connors, A. F., Speroff, T., Dawson, N. V., Thomas, C., Harrell, F. E., Wagner, D., Desbiens, N., Goldman, L., Wu, A. W., Califf, R. M., et al. (1996). The effectiveness of right heart catheterization in the initial care of critically ill patients. *Jama*, 276(11):889–897.
- [5] Cui, Y., Pu, H., Shi, X., Miao, W., and Tchetgen Tchetgen, E. (2024). Semiparametric proximal causal inference. *Journal of the American Statistical Association*, 119(546):1348–1359.
- [6] Davey Smith, G. (2008). Assessing intrauterine influences on offspring health outcomes: can epidemiological studies yield robust findings? *Basic & Clinical Pharmacology & Toxicology*, 102(2):245–256.
- [7] Davey Smith, G. (2012). Negative control exposures in epidemiologic studies. Comments on “Negative controls: a tool for detecting confounding and bias in observational studies”. *Epidemiology*, 23(2):350–351.
- [8] Dukes, O., Shpitser, I., and Tchetgen Tchetgen, E. J. (2023). Proximal mediation analysis. *Biometrika*. asad015.
- [9] Egami, N. and Tchetgen Tchetgen, E. J. (2023). Identification and estimation of causal peer effects using double negative controls for unmeasured network confounding. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 86(2):487–511.
- [10] Flanders, W. D., Klein, M., Darrow, L. A., Strickland, M. J., Sarnat, S. E., Sarnat, J. A., Waller, L. A., Winquist, A., and Tolbert, P. E. (2011). A method for detection of residual confounding in time-series and other observational studies. *Epidemiology*, 22(1):59.
- [11] Flanders, W. D., Strickland, M. J., and Klein, M. (2017). A new method for partial correction of residual confounding in time-series and other observational studies. *American Journal of Epidemiology*, 185(10):941–949.
- [12] Gagnon-Bartsch, J. A. and Speed, T. P. (2012). Using control genes to correct for unwanted variation in microarray data. *Biostatistics*, 13(3):539–552.

## Reference

- [13] Ghassami, A., Yang, A., Shpitser, I., and Tchetgen Tchetgen, E. (2024). Causal inference with hidden mediators. *Biometrika*, 112(1):asae037.
- [14] Ghassami, A., Ying, A., Shpitser, I., and Tchetgen Tchetgen, E. (2022). Minimax kernel machine learning for a class of doubly robust functionals with application to proximal causal inference. In Camps-Valls, G., Ruiz, F. J. R., and Valera, I., editors, *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151 of *Proceedings of Machine Learning Research*, pages 7210–7239. PMLR.
- [15] Glynn, A. and Ichino, N. (2019). Generalized nonlinear difference-in-difference-in-differences. *preprint*, 90.
- [16] Jacob, L., Gagnon-Bartsch, J. A., and Speed, T. P. (2016). Correcting gene expression data when neither the unwanted variation nor the factor of interest are observed. *Biostatistics*, 17(1):16–28.
- [17] Li, K., Emerman, I., Cook, A. J., Fireman, B. H., Sundaram, M., Tseng, H.-F. X., Weintraub, E. S., Yu, O., Nelson, J. L., and Shi, X. (2024). Using double negative controls to adjust for healthy user bias in a recombinant zoster vaccine safety study. *American Journal of Epidemiology*, page kwaе439.
- [18] Liu, J., Park, C., Li, K., and Tchetgen Tchetgen, E. J. (2024). Regression-based proximal causal inference. *American Journal of Epidemiology*, page kwaе370.
- [19] Mastouri, A., Zhu, Y., Gultchin, L., Korba, A., Silva, R., Kusner, M., Gretton, A., and Muandet, K. (2021). Proximal causal learning with kernels: Two-stage estimation and moment restriction. In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 7512–7523. PMLR.
- [20] Miao, W., Geng, Z., and Tchetgen Tchetgen, E. J. (2018). Identifying causal effects with proxy variables of an unmeasured confounder. *Biometrika*, 105(4):987–993.
- [21] Miao, W., Shi, X., Li, Y., and Tchetgen Tchetgen, E. J. (2024). A confounding bridge approach for double negative control inference on causal effects. *Statistical Theory and Related Fields*, 8(4):262–273.
- [22] Miao, W. and Tchetgen Tchetgen, E. J. (2017). Invited commentary: bias attenuation and identification of causal effects with multiple negative controls. *American Journal of Epidemiology*, 185(10):950–953.
- [23] Park, C., Richardson, D. B., and Tchetgen Tchetgen, E. J. (2024). Single proxy control. *Biometrics*, 80(2):ujae027.
- [24] Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press, Cambridge.

## Reference

- [25] Pearl, J. (2001). Direct and indirect effects. In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*, UAI'01, page 411–420. Morgan Kaufmann Publishers Inc.
- [26] Richardson, D. B., Keil, A., Tchetgen Tchetgen, E. J., and Cooper, G. S. (2015). Negative control outcomes and the analysis of standardized mortality ratios. *Epidemiology*, 26(5):727–732.
- [27] Richardson, D. B., Laurier, D., Schubauer-Berigan, M. K., Tchetgen Tchetgen, E. J., and Cole, S. R. (2014). Assessment and indirect adjustment for confounding by smoking in cohort studies using relative hazards models. *American Journal of Epidemiology*, 180(9):933–940.
- [28] Schuemie, M. J., Hripcsak, G., Ryan, P. B., Madigan, D., and Suchard, M. A. (2018). Empirical confidence interval calibration for population-level effect estimation studies in observational healthcare data. *Proceedings of the National Academy of Sciences*, 115(11):2571–2577.
- [29] Schuemie, M. J., Ryan, P. B., DuMouchel, W., Suchard, M. A., and Madigan, D. (2014). Interpreting observational studies: why empirical calibration is needed to correct p-values. *Statistics in Medicine*, 33(2):209–218.
- [30] Shi, X., Li, K., Miao, W., Hu, M., and Tchetgen Tchetgen, E. (2023). Theory for identification and inference with synthetic controls: A proximal causal inference framework. *Preprint arXiv:2108.13935*.
- [31] Shi, X., Miao, W., and Tchetgen Tchetgen, E. (2020a). A selective review of negative control methods in epidemiology. *Current Epidemiology Reports*, 7:190–202.
- [32] Shi, X., Miao, W., and Tchetgen Tchetgen, E. J. (2020b). Multiply robust causal inference with double negative control adjustment for categorical unmeasured confounding. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(2):521–540.
- [33] Singh, R., Sahani, M., and Gretton, A. (2019). Kernel instrumental variable regression. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 32.
- [34] Sofer, T., Richardson, D. B., Colicino, E., Schwartz, J., and Tchetgen Tchetgen, E. J. (2016). On negative outcome control of unobserved confounding as a generalization of difference-in-differences. *Statistical Science*, 31(3):348–361.

## Reference

- [35] Tchetgen Tchetgen, E. (2014). The control outcome calibration approach for causal inference with unobserved confounding. *American Journal of Epidemiology*, 179(5):633–640.
- [36] Tchetgen Tchetgen, E. J., Sofer, T., and Richardson, D. (2015). Negative outcome control for unobserved confounding under a cox proportional hazards model. *preprint*.
- [37] Tchetgen Tchetgen, E. J., Ying, A., Cui, Y., Shi, X., and Miao, W. (2024). An introduction to proximal causal inference. *Statistical Science*, 39(3):375 – 390.
- [38] Wang, J., Zhao, Q., Hastie, T., and Owen, A. B. (2017). Confounder adjustment in multiple hypothesis testing. *Annals of Statistics*, 45(5):1863–1894.
- [39] Weisskopf, M. G., Tchetgen Tchetgen, E. J., and Raz, R. (2016). Commentary: on the use of imperfect negative control exposures in epidemiologic studies. *Epidemiology*, 27(3):365–367.
- [40] Ying, A., Miao, W., Shi, X., and Tchetgen Tchetgen, E. J. (2023). Proximal causal inference for complex longitudinal studies. *Journal of the Royal Statistical Society Series B: Statistical Methodology*. qkad020.