

# Causal Inference: **Fundamentals of Partial Identification**

YES Workshop 2023  
Eindhoven, 15th March 2023

**Jakob Zeitler,**  
*PhD Candidate in Foundational Artificial Intelligence*  
Centre for Doctoral Training in Foundational AI at UCL  
**mail@jakob-zeitler.de**



**UCL**

# The basis for this talk

## **The Causal Marginal Polytope for Bounding Treatment Effects**

Jakob Zeitler, Ricardo Silva

<https://arxiv.org/abs/2202.13851>

## **Stochastic Causal Programming for Bounding Treatment Effects**

Kirtan Padh, Jakob Zeitler, David Watson, Matt Kusner, Ricardo Silva, Niki Kilbertus

<https://arxiv.org/abs/2202.10806>

**Learn more @ CLeaR (Causal Learning and Reasoning) 2023!**

**Literature Review:** <https://tinyurl.com/partial-identification>

# You have seen this before this week

Mats  
Stensrud,  
Tuesday,  
9:30am

Slide 30  
onwards

## Example on Vitamin A supplementation in northern Sumatra

$Z = 0$			$Z = 1$		
	$Y = 1$	$Y = 0$		$Y = 1$	$Y = 0$
$A = 0$	0.0064	0.9936	$A = 0$	0.0028	0.1972
$A = 1$	0.0000	0.0000	$A = 1$	0.0010	0.7990

$$-0.0054 \leq \mathbb{E}(Y^{a=1} - Y^{a=0}) \leq 0.1946$$

Consider now the effect among those who intend to be treated:

$$\mathbb{E}(Y^{a=1} - Y^{a=0} \mid A = 1) = -0.0032.$$

and among those who intend to be untreated

$$-0.007 \leq \mathbb{E}(Y^{a=1} - Y^{a=0} \mid A = 0) \leq 0.331.$$

# Program

- 1. Causal Inference**
- 2. Partial Identification**
  - a. Problem
  - b. Challenges
  - c. Solutions
- 3. Q&A**

---

# Causal Inference: **Identifiability**

# Identifiability

Can I express the effect of interest from the data?

# Identifiability

Can I express the effect of interest from the data?



**YES**

The effect is  
**identifiable**

# Identifiability

Can I express the effect of interest from the data?

**YES**

The effect is  
**identifiable**

**NO**

The effect is  
**not identifiable**



We usually ask:



*Is there  
confounding?*

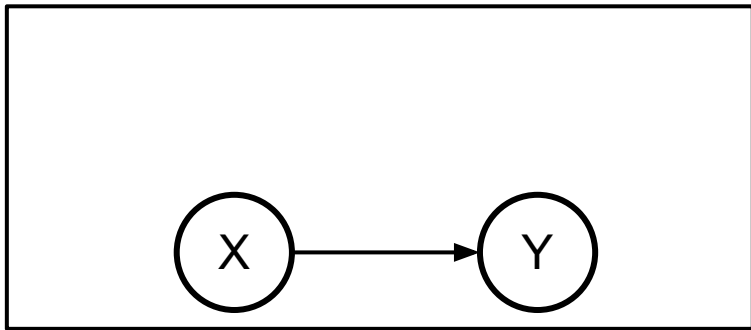
If yes, apply “**Back-Door Adjustment**”

# Example: **Confounding backdoor adjustment**

*If we want to adjust for confounding, we need to identify the confounders!*

## Scenario 1:

“No Confounding”

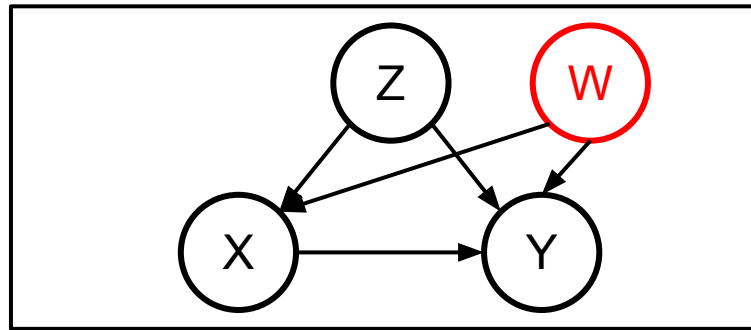


Average Treatment Effect (ATE):

$$P(Y = y | do(X = x))$$

## Scenario 2:

“Confounding” **Problem?**



Average Treatment Effect (ATE):

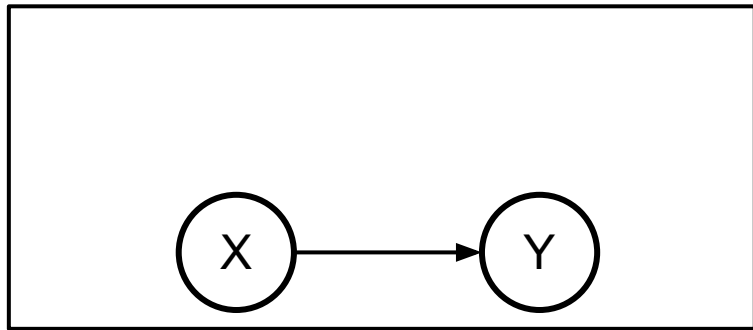
$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Example: **Confounding backdoor adjustment**

*If we want to adjust for confounding, we need to identify the confounders!*

## Scenario 1:

“No Confounding”

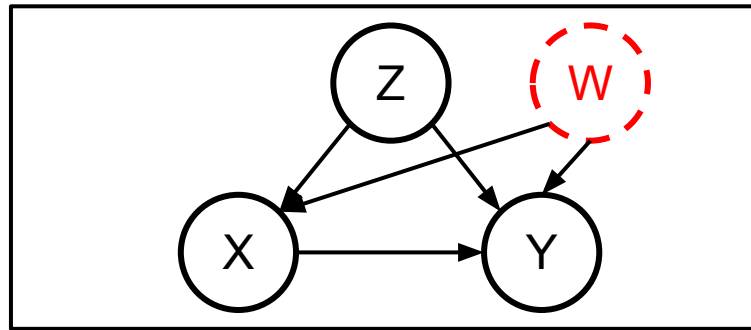


Average Treatment Effect (ATE):

$$P(Y = y | do(X = x))$$

## Scenario 2:

“Confounding” **Problem?**



Average Treatment Effect (ATE):

$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Example: Confounding backdoor adjustment

If we were able to identify the confounders!

Scenario 1:  
Success

## Problem!

### Assumption:

“No unmeasured confounding”

### Violated?: Yes!

Because we don't observe

**W, and cannot adjust for it.**

### Consequence:

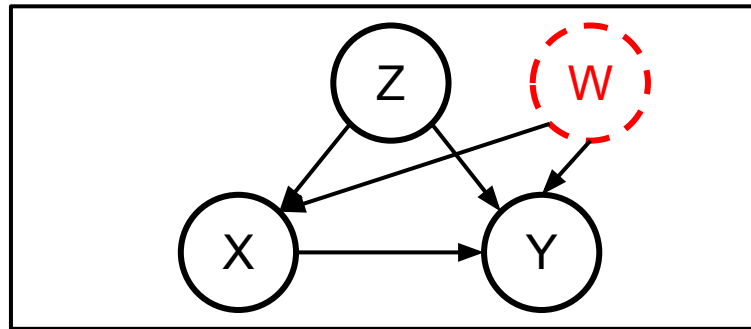
**Biased effect estimate**

Average

## Scenario 2:

“Confounding”

*Problem!*



Average Treatment Effect (ATE):

$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Assumptions are essential for **Identification**

Causal Inference is  
a language for encoding  
your causal assumptions

It's all about the **assumptions** you  
make:

**“No Causes in, No Causes out”**

**No assumptions, no identification**

- Nancy Cartwright,  
<https://doi.org/10.1093/0198235070.003.0003>

---

# Partial Identification: **Why**

*Why is partial identification **relevant**?*



# Identifiability

Can I express the effect of interest from the data?

**YES**

The effect is  
**identifiable**

**NO**

The effect is  
**not identifiable**

# Identifiability

Can I express the effect of interest from the data?

**YES**

The effect is  
**identifiable**

**NO**

The effect is  
**not identifiable**

**MAYBE**

The effect might  
be partially  
identifiable



# Identifiability

Can I express the effect of interest from the data?

**YES**

The effect is **identifiable**

**NO**

The effect is **not identifiable**

1. No **confidence in strong assumptions**, e.g. “no unmeasured confounding”
2. We want to **compare and report models** with weaker and stronger assumptions

**MAYBE**

The effect might be **partially identifiable**

# Bounds on effects of interest

Instead of doing our usual *backdoor identification strategy* of ...

$$P(Y = y|do(X = x)) = \sum_z P(Y = y|X = x, Z = z)P(Z = z) = 0.34$$

# Bounds on effects of interest

Instead of doing our usual *backdoor identification strategy* of ...

$$P(Y = y|do(X = x)) = \sum_z P(Y = y|X = x, Z = z)P(Z = z) = 0.34$$

... we try to calculate **Bounds** on the effect of interest:

i.e. **Lower Bound** < **ATE(X->Y)** < **Upper Bound**

$$\text{e.g. } 0.1 < P(Y = y|do(X = x)) < 0.4$$



# Let's review

So far, we have learned:

1. Causal Inference is all about **assumptions**
2. If we **assume a certain structure**, we can **apply the backdoor** adjustment for **full identification**
3. If we do not want to assume we know all confounders, we can calculate **bounds**, i.e. **partial identification**

Questions?

Explaining partial identification  
with the help of

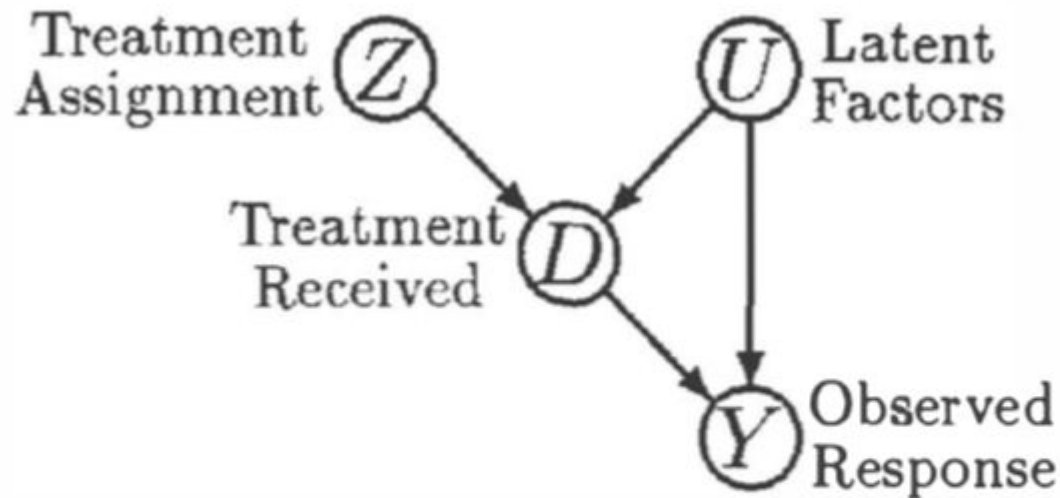
# Instrumental Variable Models

*A classical workhorse method of causal  
inference, especially econometrics.*

---

# Instrumental Variable Model

**Practically:** “Imperfect Compliance” Model



**Assume:**

All binary, i.e.  $Z, D, Y = \{0, 1\}$



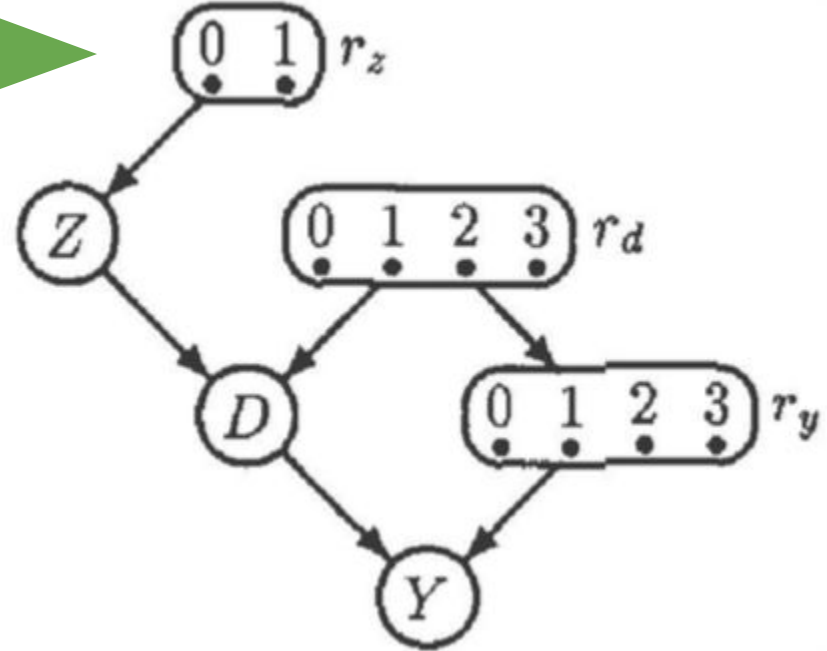
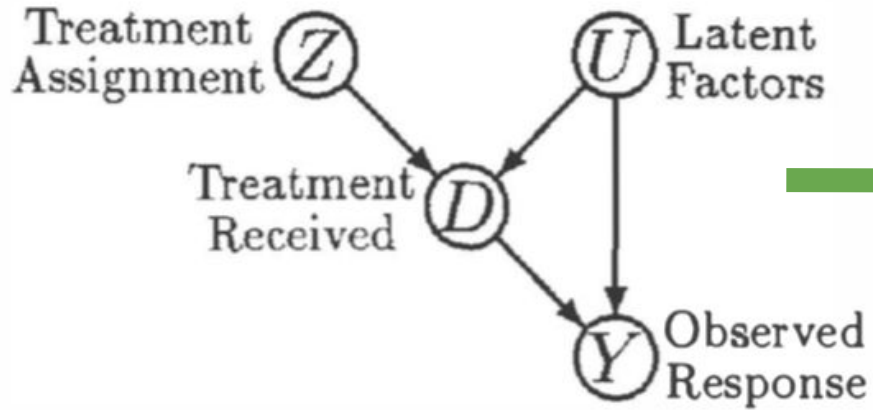
Introducing

# Response Functions Variables

*A very simple idea*

---

# Imperfect Compliance *with Response Functions*





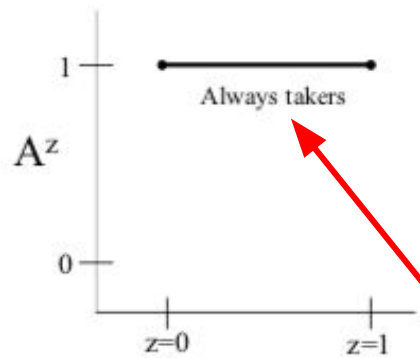
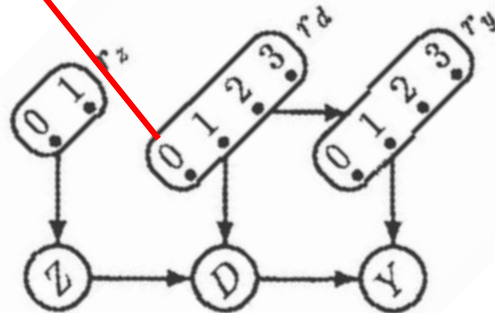


Figure 16.4



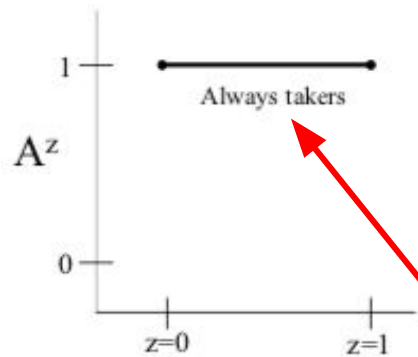


Figure 16.4

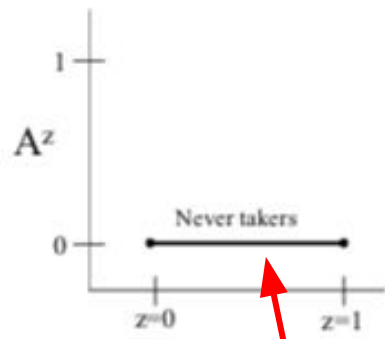
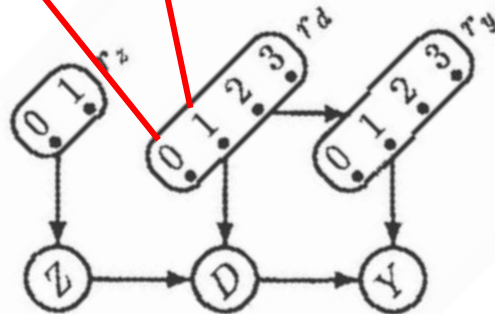


Figure 16.5



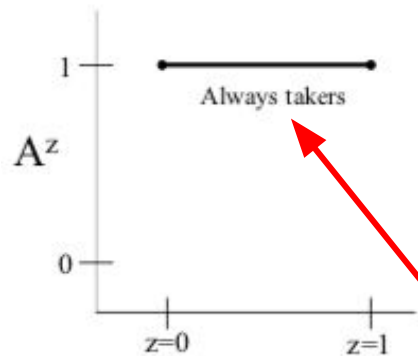


Figure 16.4

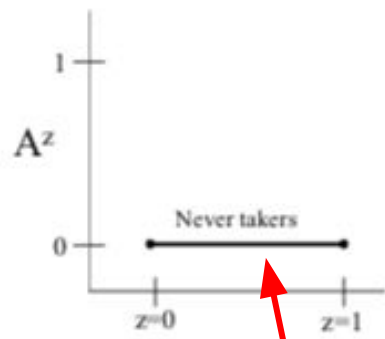


Figure 16.5

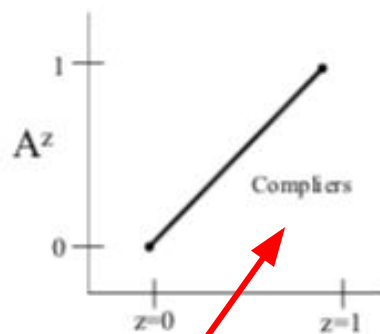
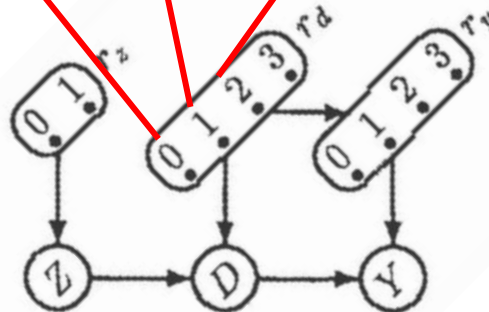


Figure 16.6



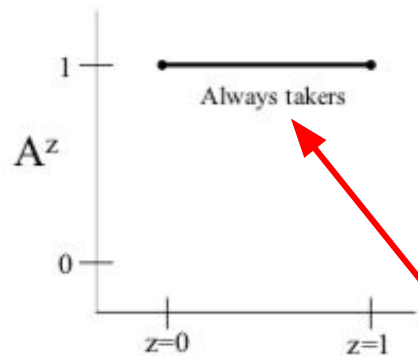


Figure 16.4

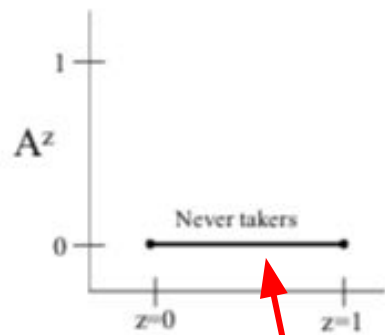


Figure 16.5



Figure 16.6

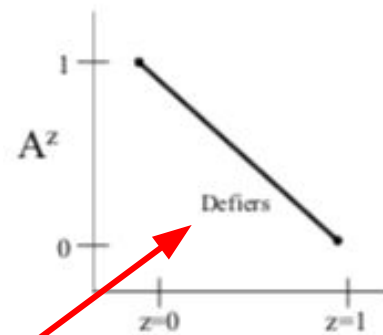
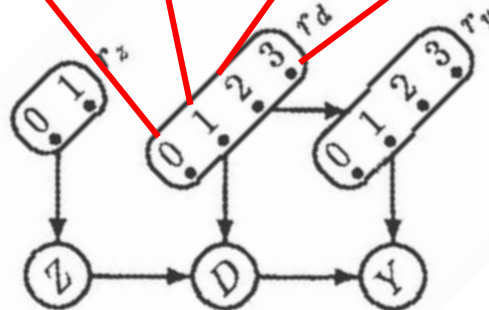


Figure 16.7



# Expressing $P(V=\{Z,D,Y\})$ with $P(R=\{R_z,R_d,R_y\})$

$$P(\mathbf{v}) = \sum_{\mathbf{r}} P(\mathbf{r}) \prod_{V \in \mathbf{V}} \mathbb{I}(v; \text{pa}_V, r_V),$$

**Observed** data  $P(\mathbf{v})$

**Latent** model  $P(\mathbf{r})$

# Calculating Upper and Lower Bounds

Using Linear Programming

LPs are convex, i.e. solutions are exact

---

# Linear Programming

$$\min/\max_{x_1, x_0} \quad \text{TCE}(x_1, x_0) = \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_1, r_Y) - \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_0, r_Y),$$

# Linear Programming

$$\begin{aligned} \min/\max \quad & \text{TCE}(x_1, x_0) = \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_1, r_Y) - \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_0, r_Y), \\ \text{s.t.} \quad & P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0, \end{aligned}$$



# Linear Programming

$$\min/\max \quad \text{TCE}(x_1, x_0) = \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_1, r_Y) - \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_0, r_Y),$$

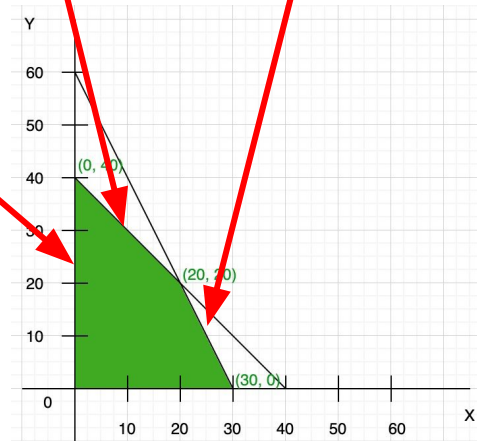
s.t.

$$P(\mathbf{V}) = P(\mathcal{D}),$$

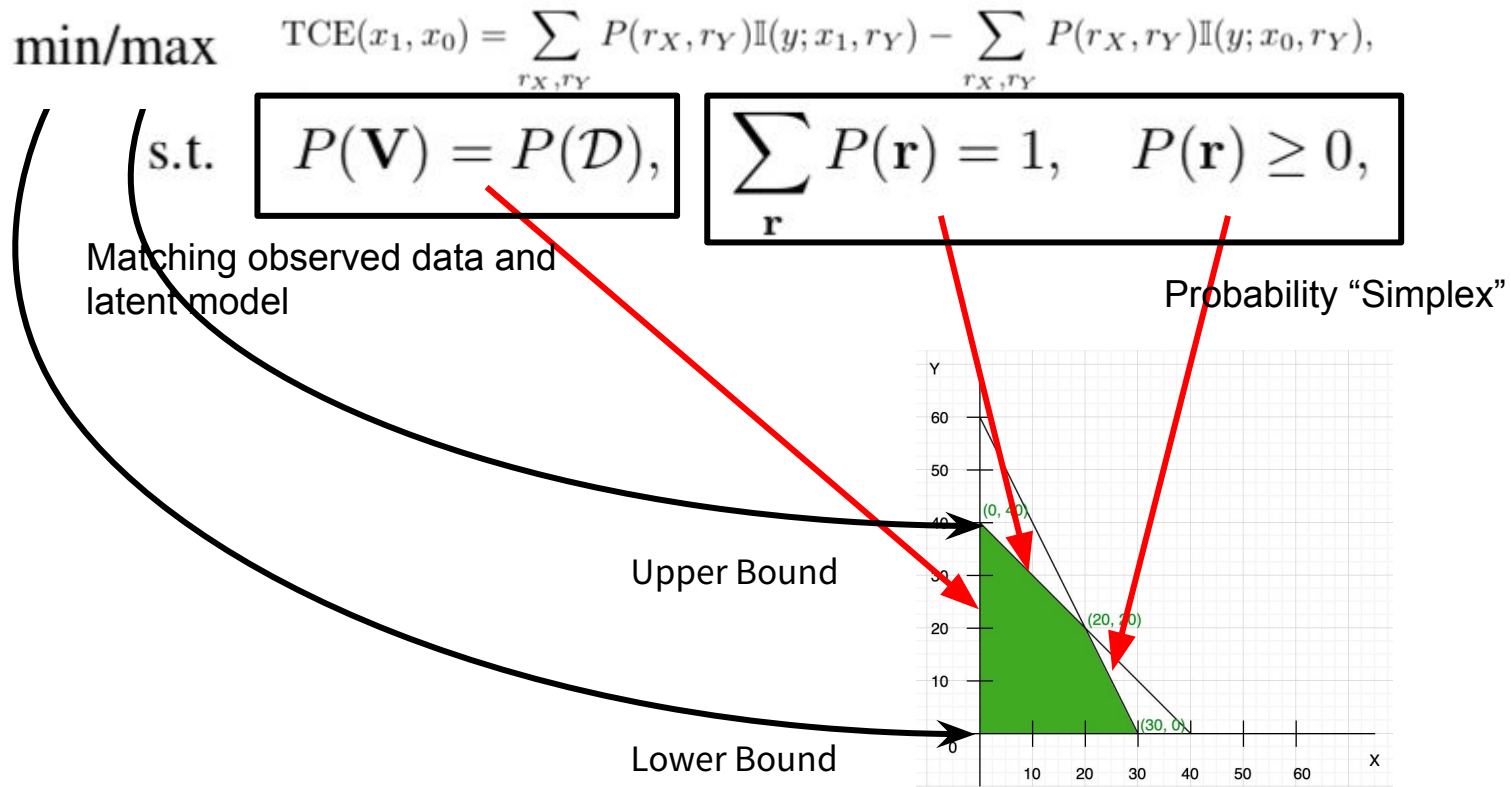
Matching observed data and  
latent model

$$\sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$$

Probability “Simplex”



# Linear Programming



# Partial Identification: Challenges

*Why is partial identification **hard**?*

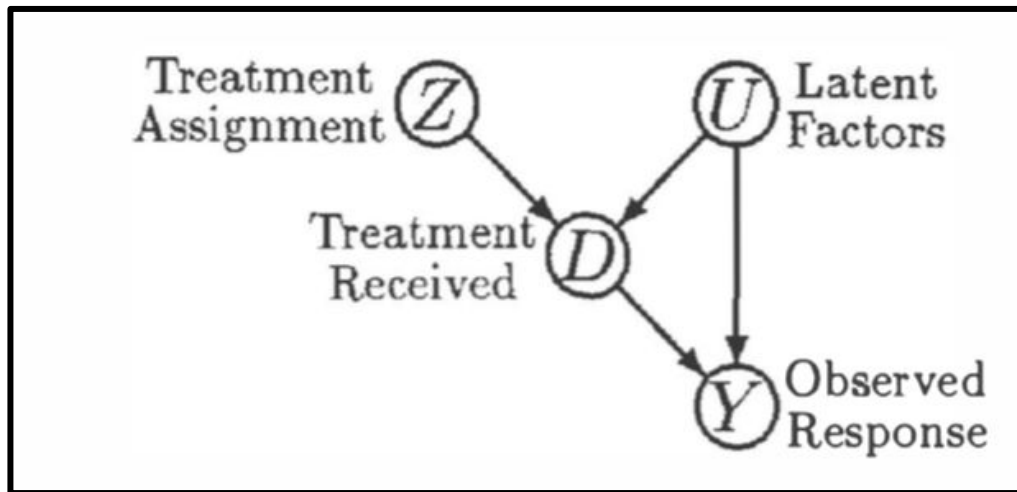


# Beyond Instrumental Variable Models

How to generalise response  
functions variables

---

So far ...



# Response Functions Variables are *universal*

Response Function Variables:

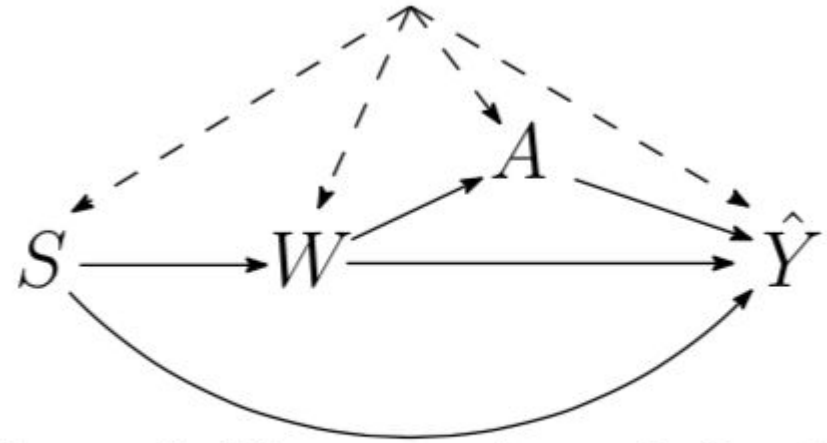
1. Number of states “coming in” from parents
2. Numbers of “states” of node (the domain)

$$N_V = |V|^{|\text{PA}_V|}$$

How many Response Function states **per variable**?

$$N_V = |V| |PA_V|$$

S =  
W =  
A =  
Y =



How many Response Function states **per variable**?

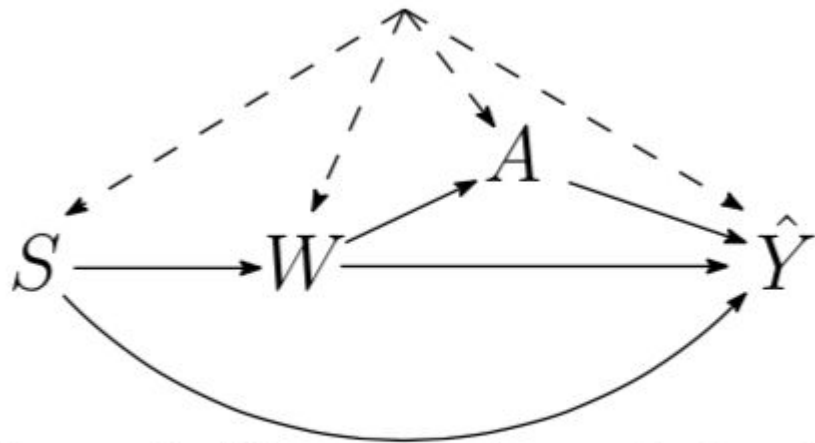
$$N_V = |V| |PA_V|$$

$$S = 2$$

$$W =$$

$$A =$$

$$Y =$$





How many Response Function states **per variable**?

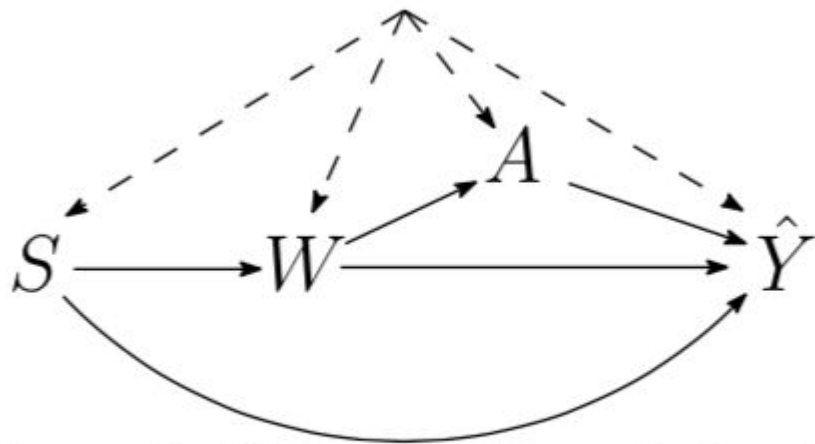
$$N_V = |V| |PA_V|$$

$$S = 2$$

$$W = 2 ** 2$$

$$A =$$

$$Y =$$



How many Response Function states **per variable**?

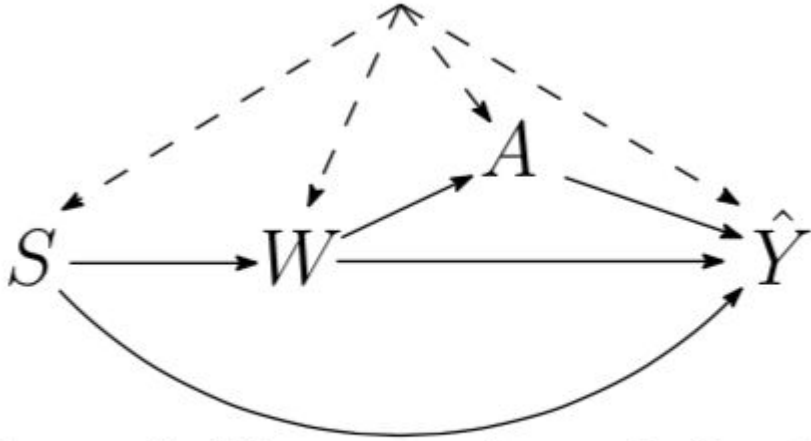
$$N_V = |V| |PA_V|$$

$$S = 2$$

$$W = 2 ** 2$$

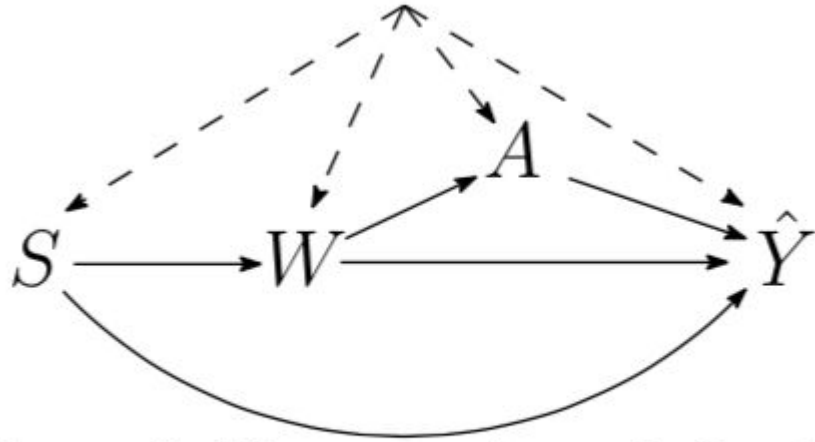
$$A = 2 ** 2$$

$$Y =$$



How many Response Function states **per variable**?

$$N_V = |V| |PA_V|$$



$$S = 2$$

$$W = 2 ** 2$$

$$A = 2 ** 2$$

$$Y = 2 ** (2 * 2 * 2) = 256 (!!!)$$

How many Response Function states **per variable**?

$$N_V = |V| |PA_V|$$

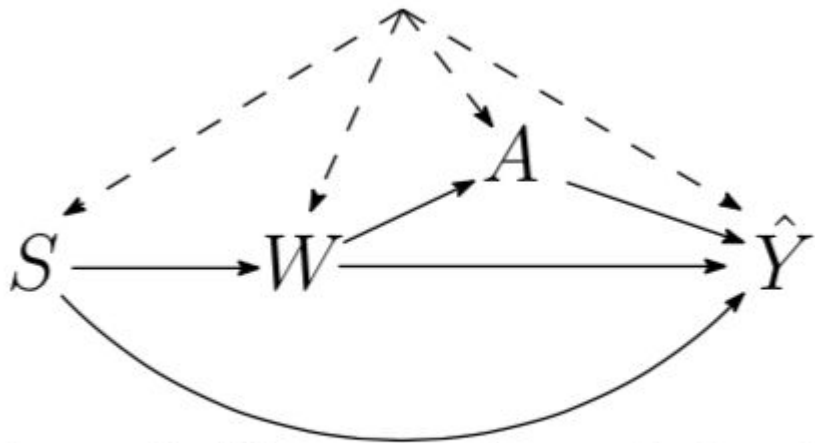
$$S = 2$$

$$W = 2 ** 2$$

$$A = 2 ** 2$$

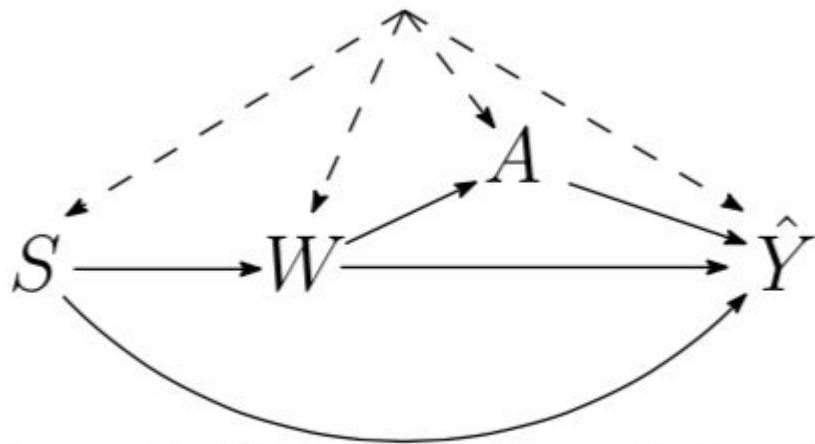
$$Y = 2 ** (2 * 2 * 2) = 256 (!!!)$$

How many Response Function states **for the whole model**?



How many Response Function states **per variable**?

$$N_V = |V| |PA_V|$$



$$S = 2$$

$$W = 2 ** 2$$

$$A = 2 ** 2$$

$$Y = 2 ** (2 * 2 * 2) = 256 (!!!)$$

How many Response Function states **for the whole model**?

$$2 * 4 * 4 * 256 = \underline{8192}$$

How many Response Function states **per variable**?

$$N_V = |V| |PA_V|$$

$$S = 2$$

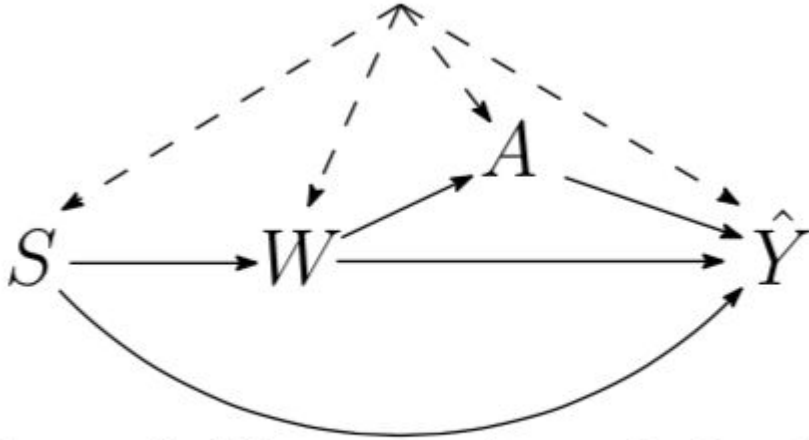
$$W = 2 ** 2$$

$$A = 2 ** 2$$

$$Y = 2 ** (2 * 2 * 2) = 256 (!!!)$$

How many Response Function states **for the whole model**?

$$2 * 4 * 4 * 256 = \underline{8192}$$

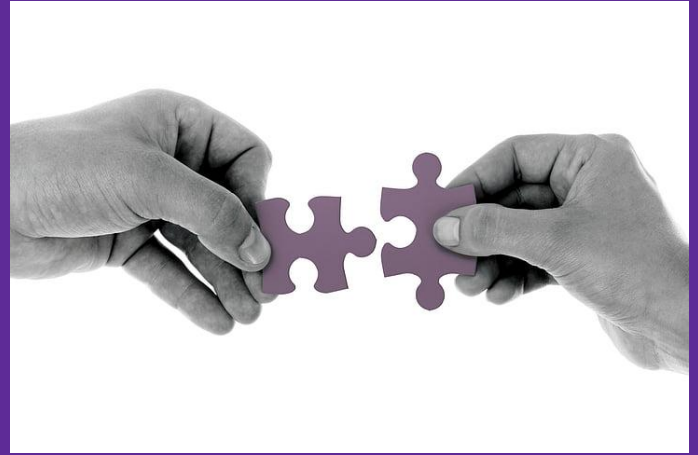


**Problem:** RFV space increases *super-exponentially*.

How can we deal with that?

# Partial Identification: **Solutions**

*How is partial identification **feasible**?*

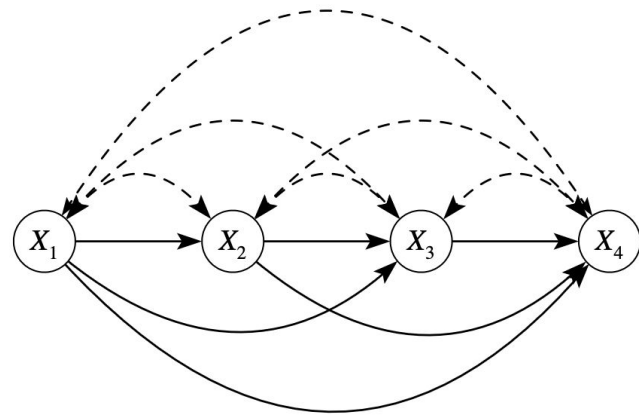


# Solution: **Causal Marginal Polytope**

**Idea:** Consider smaller parts of the graph, aka. “marginals”

## **Example:**

- 4 Variables
- Fully connected
- Fully confounded



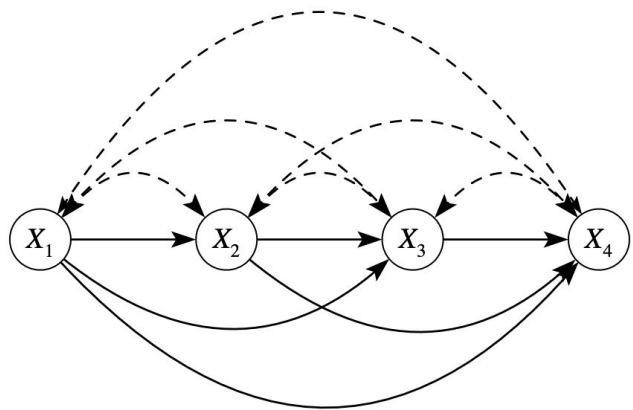
## **The Causal Marginal Polytope for Bounding Treatment Effects**

Jakob Zeitler, Ricardo Silva

<https://arxiv.org/abs/2202.13851>



# “Complete” model

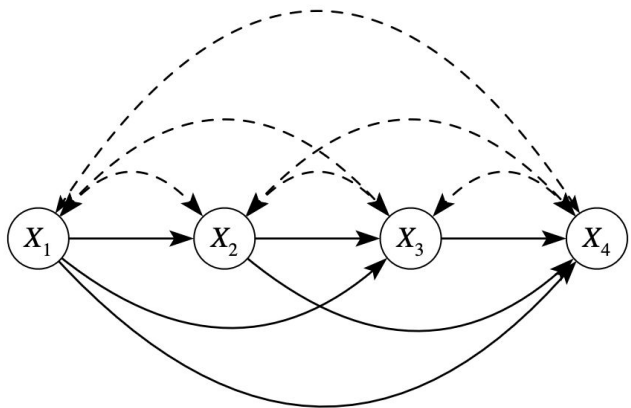


**Marginals:** 1, i.e. the complete 4 variable model

**Number of Parameters:** 32,768

**Tightness of bounds:** Sharp

# “Complete” model

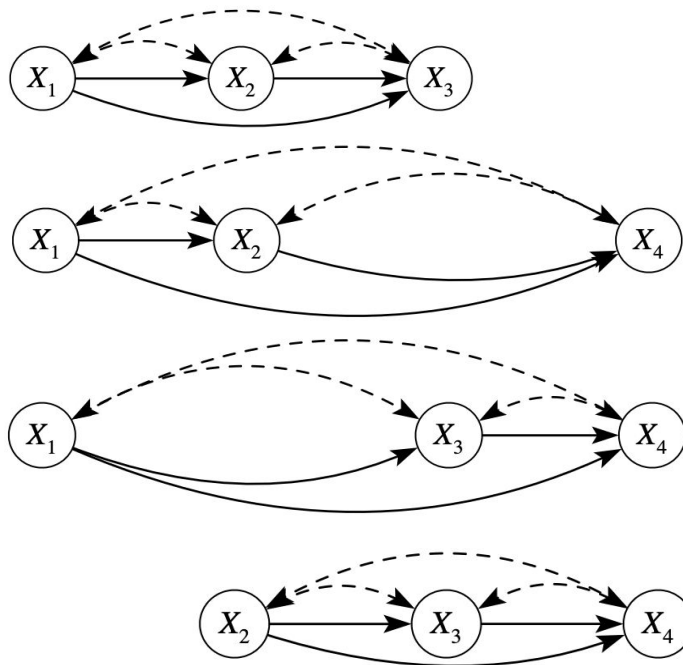


**Marginals:** 1, i.e. the complete 4 variable model

**Number of Parameters:** 32,768

**Tightness of bounds:** Sharp

# Causal Marginal Polytope

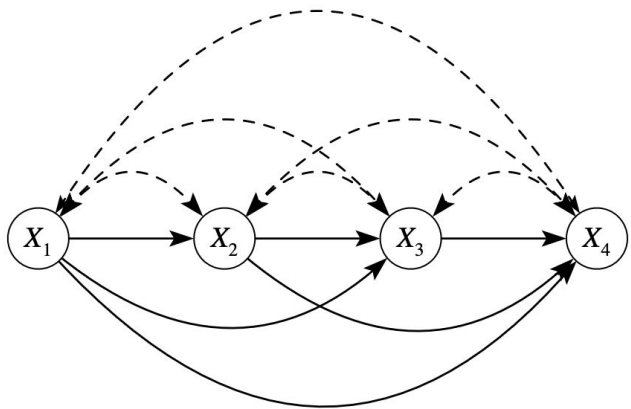


**Marginals:** 4, each with 3 variables

**Number of Parameters:**  $128 \times 4 = 512$

**Tightness of bounds:** More loose

# “Complete” model

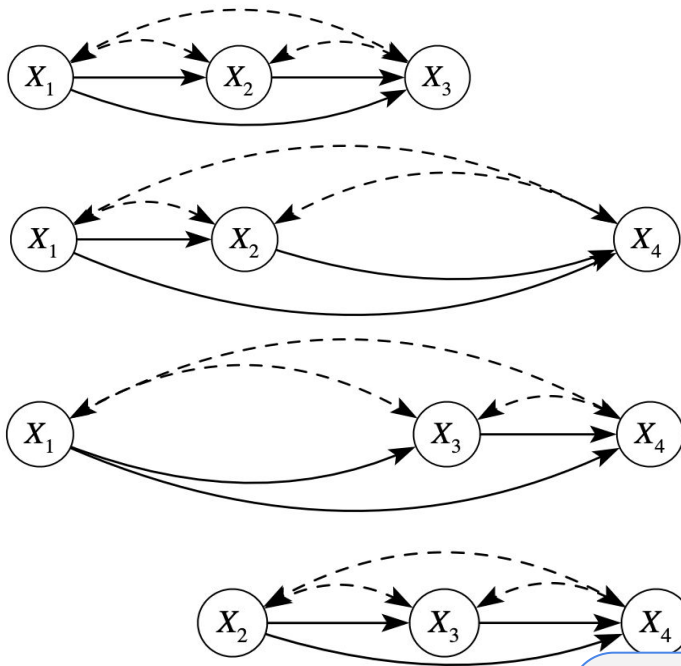


**Marginals:** 1, i.e. the complete 4 variable model

**Number of Parameters:** 32,768

**Tightness of bounds:** Sharp

# Causal Marginal Polytope



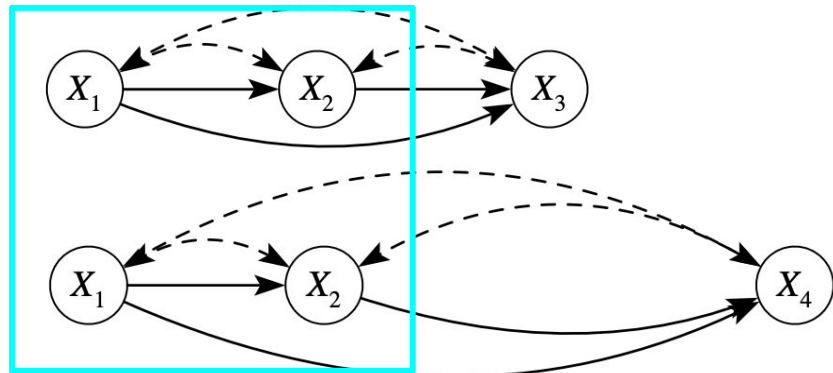
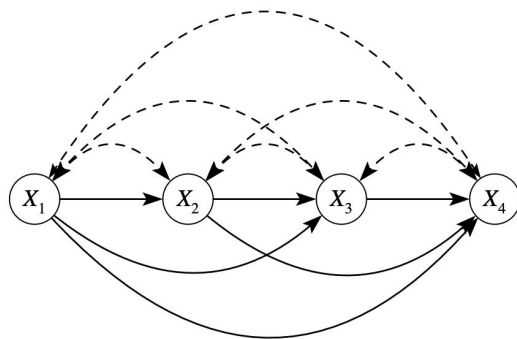
**Marginals:** 4, each with 3 variables

**Number of Parameters:**  $128 \times 4 = 512$

**Tightness of bounds:** More loose

Let's add constraints to tighten the bounds

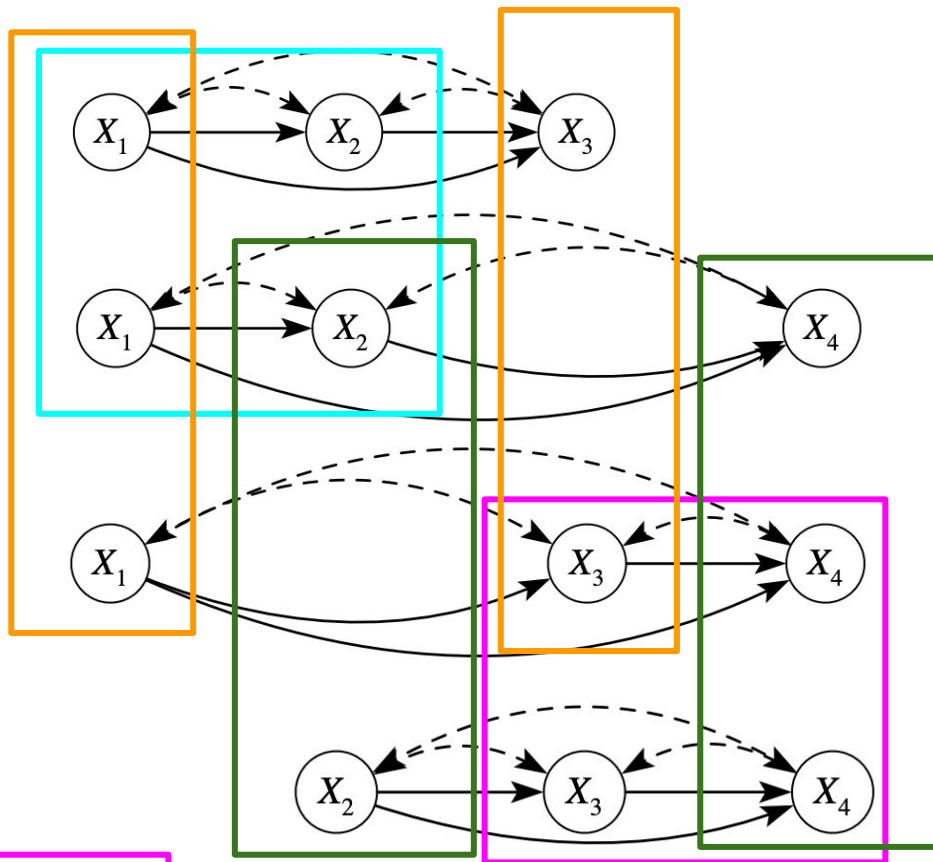
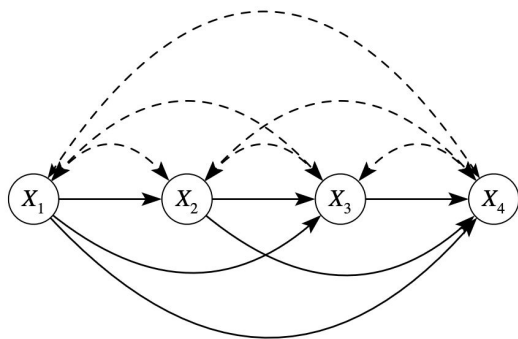
# Constraint A: Overlaps



$$\sum_{X_3} P(X_1, X_2, X_3) = \sum_{X_4} P(X_1, X_2, X_4)$$

$\{\{X_1, X_2\},$

# Constraint A: **Overlaps**



$\{\{X_1, X_2\}, \{X_2, X_4\}, \{X_1, X_3\}, \{X_3, X_4\}\}$

min/max

Objective: e.g. TCE

s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$

Constraint A:  
**Overlaps**

min/max

Objective: e.g. TCE

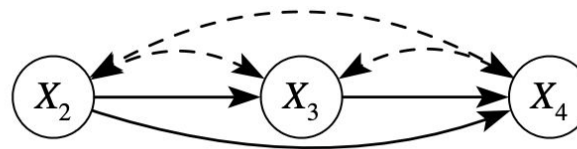
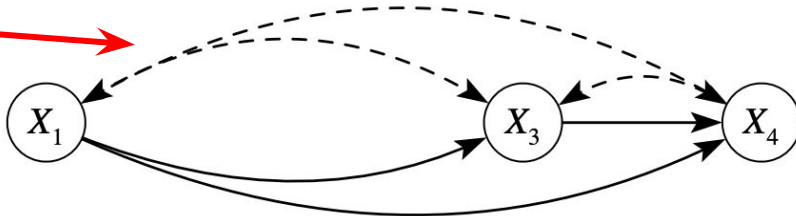
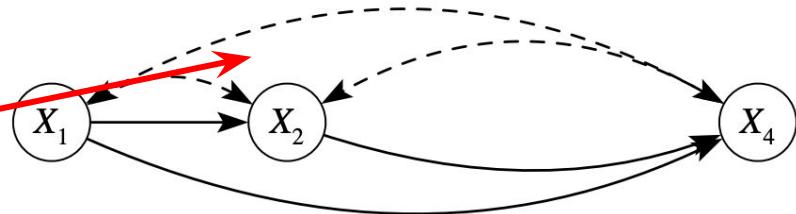
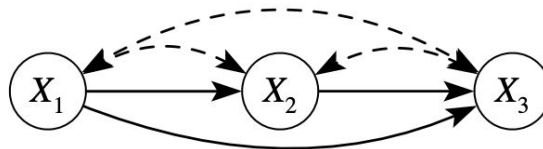
s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$

Constraint A:  
**Overlaps**

$$\sum_{X_3} P(X_1, X_2, X_3) = \sum_{X_4} P(X_1, X_2, X_4)$$

# Constraint B: **Expert knowledge**

1. Bidirected edges
2. Directed edges





min/max

Objective: e.g. TCE

s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$

Constraint A:  
**Overlaps**

$$\sum_{X_3} P(X_1, X_2, X_3) = \sum_{X_4} P(X_1, X_2, X_4)$$

Constraint B:  
**Expert  
knowledge**

min/max

Objective: e.g. TCE

s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$

Constraint A:  
**Overlaps**

$$\sum_{X_3} P(X_1, X_2, X_3) = \sum_{X_4} P(X_1, X_2, X_4)$$

Constraint B:  
**Expert  
knowledge**

Directed edges

$$|P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_i \setminus j}), do(V_j = 1), v_{\mathcal{M}'}) - \\ P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_i \setminus j}), do(V_j = 0), v_{\mathcal{M}'})| \leq \epsilon_{ij},$$

min/max

Objective: e.g. TCE

s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$

Constraint A:  
**Overlaps**

$$\sum_{X_3} P(X_1, X_2, X_3) = \sum_{X_4} P(X_1, X_2, X_4)$$

Constraint B:  
**Expert  
knowledge**

Directed edges

$$|P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_i \setminus j}), do(V_j = 1), v_{\mathcal{M}'}) - \\ P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_i \setminus j}), do(V_j = 0), v_{\mathcal{M}'})| \leq \epsilon_{ij},$$

Bidirected edges

$$|P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_{ij}}), do(V_j = v_j), v_{\mathcal{M}'}) - \\ P_{\mathcal{M}}(V_i = 1 \mid do(v_{pa_{ij}}), V_j = v_j, v_{\mathcal{M}'})| \leq \epsilon_{ij}^c,$$

# Example: **Simulation**

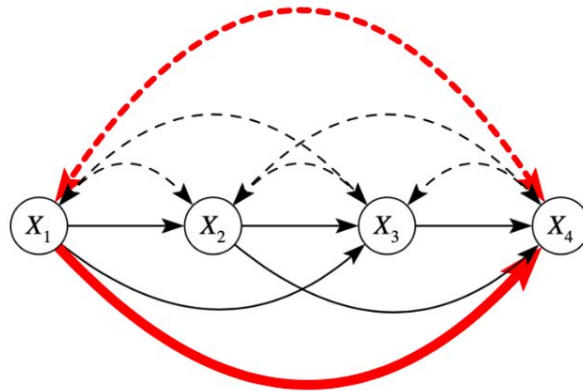
- 4 Variables
- We constrain the red directed and bidirected edges
- We supply the following data regimes

$do(\emptyset)$

$do(x_2), do(x_3)$

$do(x_2, x_3)$

$do(x_1 = 0, x_3 = 0), do(x_1 = 0, x_3 = 1)$



# Example: **Simulation**

- 4 Variables
- We constrain the red directed and bidirected edges
- We supply the following data regimes

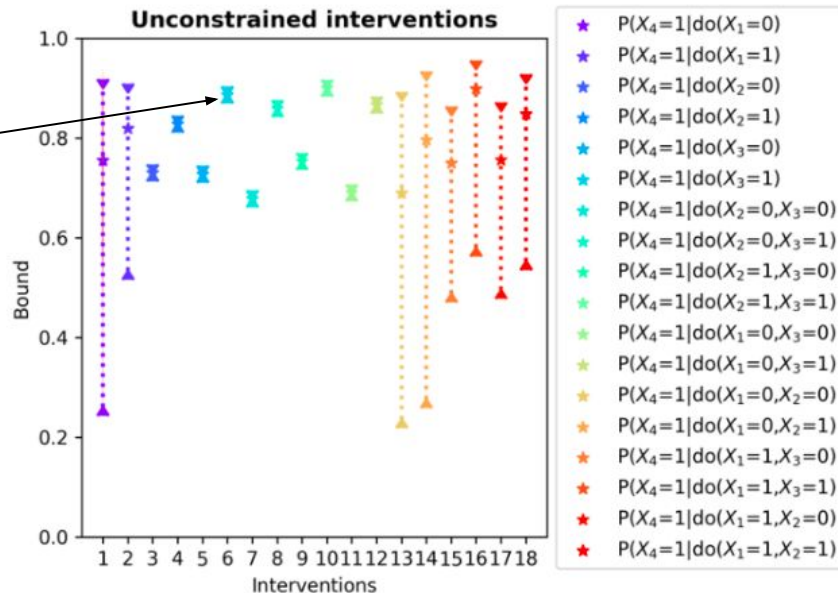
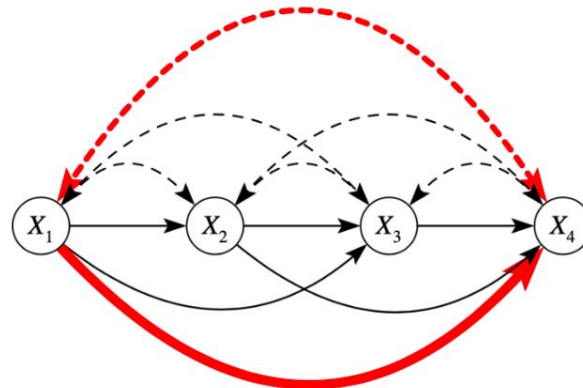
$do(\emptyset)$

$do(x_2), do(x_3)$

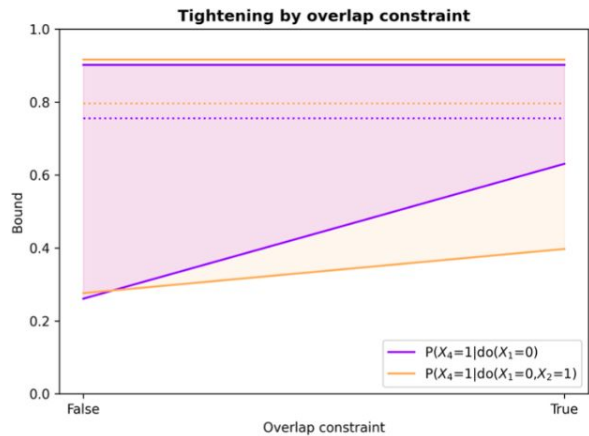
$do(x_2, x_3)$

$do(x_1 = 0, x_3 = 0), do(x_1 = 0, x_3 = 1)$

Identification from data



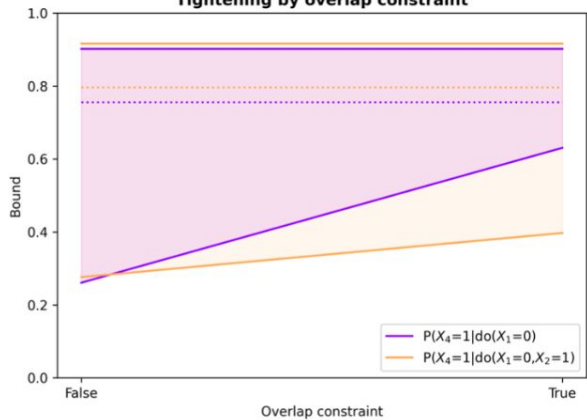
# Results: Each constraint by itself



Tightening of bounds by  
constraining the **overlap**  
of the margins.

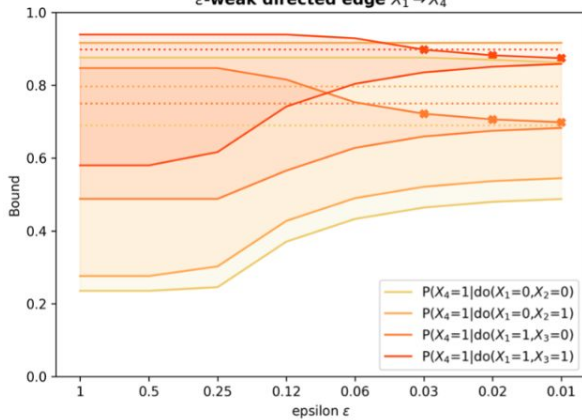
# Results: Each constraint by itself

Tightening by overlap constraint



Tightening of bounds by  
constraining the **overlap**  
**of the margins.**

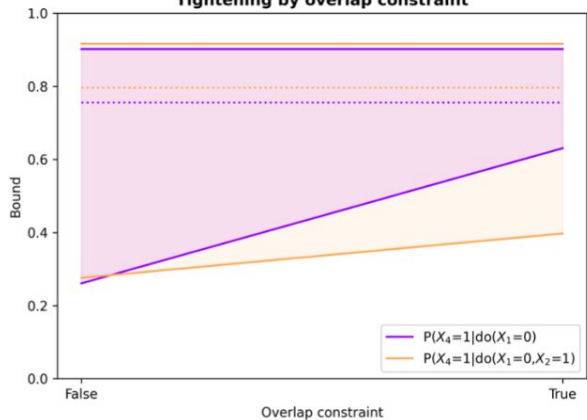
$\epsilon$ -weak directed edge  $X_1 \rightarrow X_4$



A **weak directed edge**  
tightens bounds. Invalid  
bounds are marked with  
a cross.

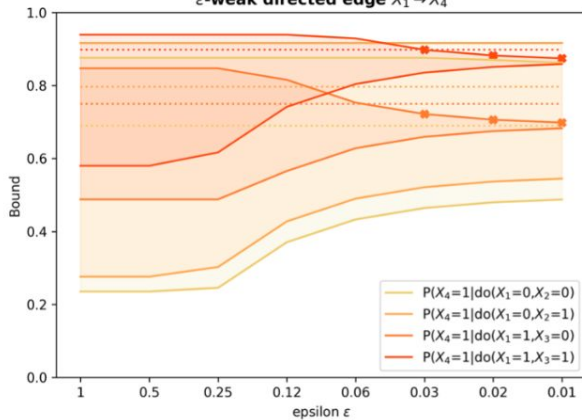
# Results: Each constraint by itself

Tightening by overlap constraint



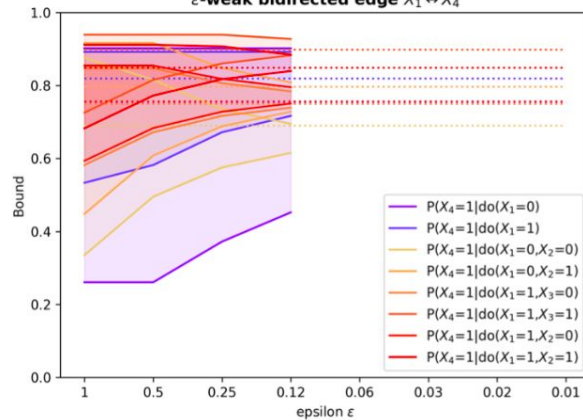
Tightening of bounds by constraining the **overlap of the margins**.

$\epsilon$ -weak directed edge  $X_1 \rightarrow X_4$



A **weak directed edge** tightens bounds. Invalid bounds are marked with a cross.

$\epsilon$ -weak bidirected edge  $X_1 \leftrightarrow X_4$



A **weak bidirected edge** tightens bounds. Lower epsilon values are falsified by infeasibility.





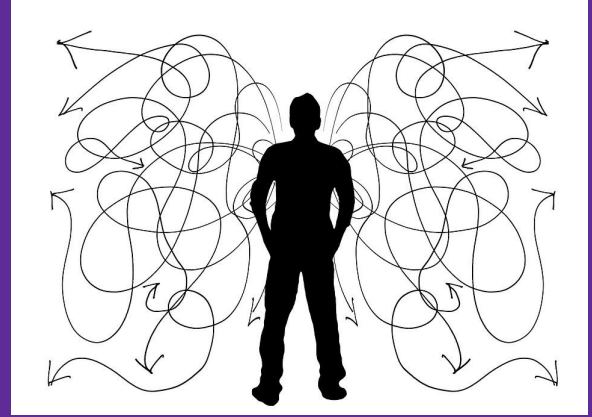
# Let's review, again

We have learned:

1. If we do not want to assume we know all confounders, we can calculate **bounds**, i.e. **partial identification**
2. But *exact* partial identification is hard: **there is no free lunch.**

# Partial Identification: **Alternatives**

*What **alternative solutions** are there?*



# What's really behind partial identification

Take the IV model on the left, with *two possible models for B*.

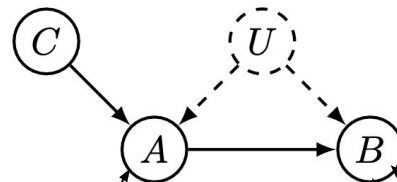


Figure 1: Example:  $A, B, C$  are observed variables and  $U$  is a hidden variable.

$U$	$C$	$A = f_A$
0	0	0
0	1	0
1	0	1
1	1	1

Table 1: Equation  $f_A(c, u)$  for determining values of  $A$ .

$U$	$A$	$B = f_B^1$	$B = f_B^2$
0	0	0	0
0	1	1	0
1	0	0	0
1	1	1	1

Table 2: Equations  $f_B^1(a, u)$  and  $f_B^2(a, u)$  for determining values of  $B$ .

# What's really behind partial identification

Take the IV model on the left, with *two possible models for B*.

1. Setting  $P(U = 1)$  as  $p$
2. The first model implies:
  - a.  $P(B = 1 | \text{do}(A = 1)) = 1$
  - b.  $P(B = 1 | \text{do}(A = 0)) = 0$
  - c. Therefore,  $\text{ACE} = 1 - 0 = 1$

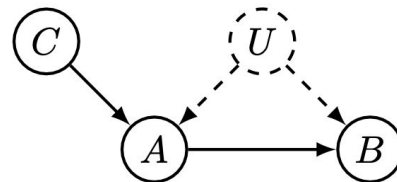


Figure 1: Example:  $A, B, C$  are observed variables and  $U$  is a hidden variable.

$U$	$C$	$A = f_A$
0	0	0
0	1	0
1	0	1
1	1	1

Table 1: Equation  $f_A(c, u)$  for determining values of  $A$ .

$U$	$A$	$B = f_B^1$	$B = f_B^2$
0	0	0	0
0	1	1	0
1	0	0	0
1	1	1	1

Table 2: Equations  $f_B^1(a, u)$  and  $f_B^2(a, u)$  for determining values of  $B$ .

# What's really behind partial identification

Take the IV model on the left, with *two possible models for B*.

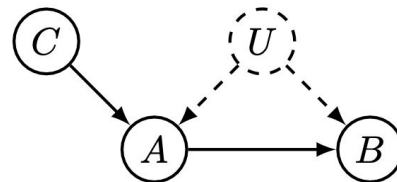


Figure 1: Example:  $A, B, C$  are observed variables and  $U$  is a hidden variable.

1. Setting  $P(U = 1)$  as  $p$
2. The first model implies:
  - a.  $P(B = 1 | \text{do}(A = 1)) = 1$
  - b.  $P(B = 1 | \text{do}(A = 0)) = 0$
  - c. Therefore,  $\text{ACE} = 1 - 0 = 1$

3. The second model implies:
  - a.  $P(B = 1 | \text{do}(A = 1)) = \mathbf{p}$
  - b.  $P(B = 1 | \text{do}(A = 0)) = 0$
  - c. Therefore,  $\text{ACE} = p - 0 = p$

4. We cannot reject either model
5. Therefore, the bounds are:  $[p, 1]$

$U$	$A$	$B = f_B^1$	$B = f_B^2$
0	0	0	0
0	1	1	0
1	0	0	0
1	1	1	1

Table 2: Equations  $f_B^1(a, u)$  and  $f_B^2(a, u)$  for determining values of  $B$ .

# Approximate Partial Identification

1. Define a class of possible models
2. Reject models that disagree with that data
3. Calculate ATEs for all models left
4. The highest and lowest are your bounds

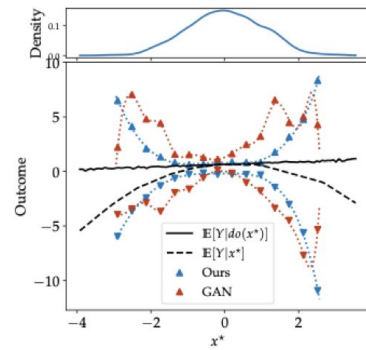
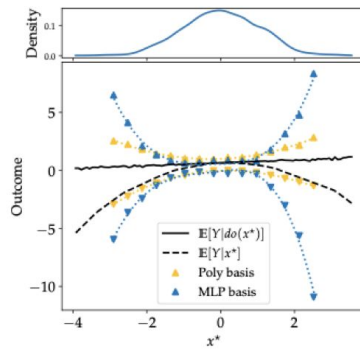
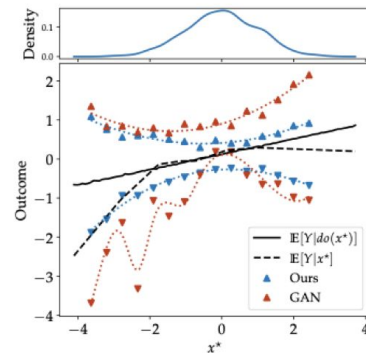
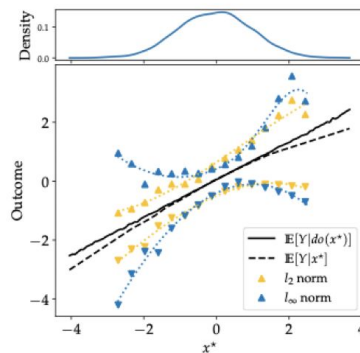
$$\min_{\eta \in \mathbb{R}^d} / \max_{\eta \in \mathbb{R}^d} \quad o_{x^*}(\eta) = \psi(x^*)^\top \mathbb{E}_N[\mu_{\eta_0}(N)]$$

$$\text{subject to} \quad \text{dist}(p_\eta, \hat{p}) \leq \epsilon.$$

## Stochastic Causal Programming for Bounding Treatment Effects

Kirtan Padh, Jakob Zeitler, David Watson, Matt Kusner, Ricardo Silva, Niki Kilbertus

<https://arxiv.org/abs/2202.10806>







## Where we end today:

The goal of causal modelling is **not identification as such**, but to make the best out of the observables that we do have.

Don't say: "*I can't provide a unique solution from your assumptions, go home.*"

Do say: "*These are all the solutions compatible with your assumptions*".

*No academic  
talk without a  
**meme***

how it started

The **core**  
**question** of  
causal  
inference:  
**Identifiability**

*“I can’t provide a  
unique solution from  
your assumptions, go  
home.”*

how it's going

**Make the best  
out of the  
observables**  
that we have

*“These are all the  
solutions compatible  
with your assumptions”.*



# Questions & Answers

## Time for

- Questions
- Discussion
- Feedback
- 40+ slides in the Appendix

You can also talk to me about

1. **Synthetic Control**
2. **(Causal) Bayesian Optimisation**
3. **Topological Perspectives of Causality**

# partial-identification.com

- Book (draft)
- 45 minute version of this talk

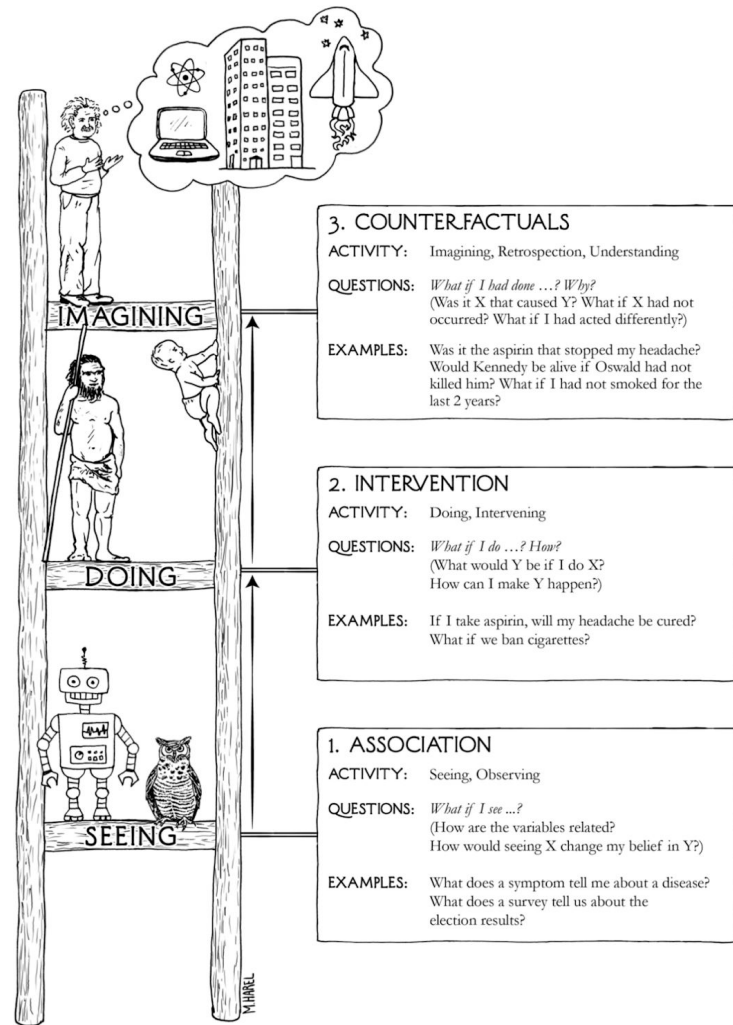
# Contribute to Causal Genealogy: genealogy.causality.link

	A	B	C	D	E	
1	<p><b>External Link to Spreadsheet for Editing:</b> <a href="https://docs.google.com/spreadsheets/d/1rNrb7rvIVnUR0CPBqbvJNNix76wfA5038rJ5zHrSZFk/edit?usp=sharing">https://docs.google.com/spreadsheets/d/1rNrb7rvIVnUR0CPBqbvJNNix76wfA5038rJ5zHrSZFk/edit?usp=sharing</a></p> <p><b>How to use this Genealogy</b> (inspired by <a href="https://www.genealogy.math.ndsu.nodak.edu">https://www.genealogy.math.ndsu.nodak.edu</a>)</p> <p>This spreadsheet is an overview of the researchers in causal inference.</p> <p>For beginners, this is hopefully a useful resource.</p>					
2	<p><b>FAQ</b></p> <p><i>How do I add more people to the sheet?</i></p> <p><i>Simply use the Google Docs commenting function and we will turn your comment into an entry.</i></p>					
3	<b>Name</b>	<b>Institution</b>	<b>Supervisor</b>	<b>Location</b>	<b>Previous Positions</b>	<b>Link</b>
4	<b>Academia</b>					
5	<b>UCLA</b>					
6	Judea Pearl	UCLA	?	US	Rutgers, Technion, Newark College of Engineering	<a href="http://b">http://b</a>
7	Wesley Salmon	UCLA	Hans Reichenbach	US	?	
8	Hans Reichenbach	UCLA	Paul Hensel, Max Noeth	US	Berlin, Istanbul, Erlangen	
9	<b>John Hopkins</b>					
10	Ilya Shpitser	John Hopkins		US	UCLA, Judea Pearl	<a href="https://">https://</a>
11	<b>Oregon State University</b>					
12	Karthika Mohan	Oregon State University	Judea Pearl	US		<a href="http://w">http://w</a>
13	<b>CMU</b>					
14	Kun Zhang	CMU		Pittsburgh, US	MPI Tübingen	
15	Clark Glymour	CMU	Wesley Salmon	Pittsburgh, US		
16	Peter Spirtes	CMU		Pittsburgh, US		
17	<b>ETH Zürich</b>					
18	Peter Bühlmann	ETH		Zürich	?	
19	Marloes Maathuis	ETH		Zürich	?	<a href="https://">https://</a>
20	Nicolai Meinshausen	ETH				

# Appendix

# The Causal Ladder

(Pearl, 2018: "Book of Why")



# Two schools of thought: Apples and Oranges?

	Potential Outcomes	Graphical Models (DAGs)
Areas	Epidemiology, Econometrics	Computer Science, Artificial Intelligence
Schools	Harvard, Berkeley and more	UCLA, CMU and more

*Discussion see here:*

<https://www.jstor.org/stable/4616823?seq=1>



# Compare: Identifiability in RCTs and obs. studies

## Randomized experiments

**(Conditional) Exchangeability**

*True by  
design!  
(randomized)*

**Identifiability**

Calculate effect via e.g. IPW.

## Observational studies (*The focus of this talk*)

**Identifiability?**

*True?*

(Alternatively, if these conditions are not met, one can try to use a *instrumental variable model*)

# Compare: Identifiability in RCTs and obs. studies

## Randomized experiments

**(Conditional) Exchangeability**

*True by  
design!  
(randomized)*

**Identifiability**

Calculate effect via e.g. IPW.

## Observational studies (*The focus of this talk*)

**Consistency**  
**(Conditional) Exchangeability**  
**Positivity**

*True?*

**Identifiability**

Observational study acts *as if* it was a conditionally randomized study.

Calculate effect via e.g. IPW.

(Alternatively, if these conditions are not met, one can try to use a *instrumental variable model*)



# What does the causal graph look like?

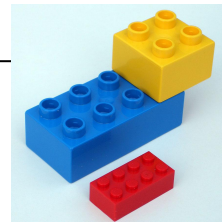
*If we want to adjust for confounding, we need to **identify the confounders!***

Define SCM:  $\langle \mathbf{V}, \mathbf{F}, \mathbf{U} \rangle$

$\mathbf{V}$ : observed variables

$\mathbf{F}$ : structural equations

$\mathbf{U}$ : background variables, with  $P(\mathbf{U})$



*“The building blocks of causal DAGs (directed acyclic graphs)”*

# What does the causal graph look like?

If we want to adjust for confounding, we need to **identify the confounders!**

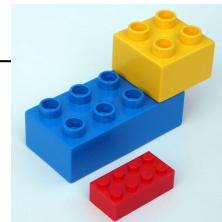
$$V = \{X, Y, Z\}, U = \{U_X, U_Y, U_Z\}, F = \{f_X, f_Y, f_Z\}$$

Define SCM:  $\langle \mathbf{V}, \mathbf{F}, \mathbf{U} \rangle$

$\mathbf{V}$ : observed variables

$\mathbf{F}$ : structural equations

$\mathbf{U}$ : background variables, with  $P(\mathbf{U})$



This SCM implies this DAG:



*“The building blocks of causal DAGs (directed acyclic graphs)”*

# What does the causal graph look like?

If we want to adjust for confounding, we need to **identify the confounders!**

$$V = \{X, Y, Z\}, U = \{U_X, U_Y, U_Z\}, F = \{f_X, f_Y, f_Z\}$$

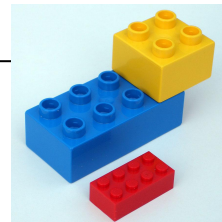
$$f_X : X = U_X$$

Define SCM:  $\langle \mathbf{V}, \mathbf{F}, \mathbf{U} \rangle$

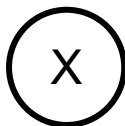
$\mathbf{V}$ : observed variables

$\mathbf{F}$ : structural equations

$\mathbf{U}$ : background variables, with  $P(\mathbf{U})$



This SCM implies this DAG:



*“The building blocks of causal DAGs (directed acyclic graphs)”*

# What does the causal graph look like?

If we want to adjust for confounding, we need to **identify the confounders!**

$$V = \{X, Y, Z\}, U = \{U_X, U_Y, U_Z\}, F = \{f_X, f_Y, f_Z\}$$

$$f_X : X = U_X$$

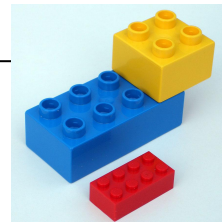
$$f_Y : Y = 4x + U_Y$$

Define SCM:  $\langle \mathbf{V}, \mathbf{F}, \mathbf{U} \rangle$

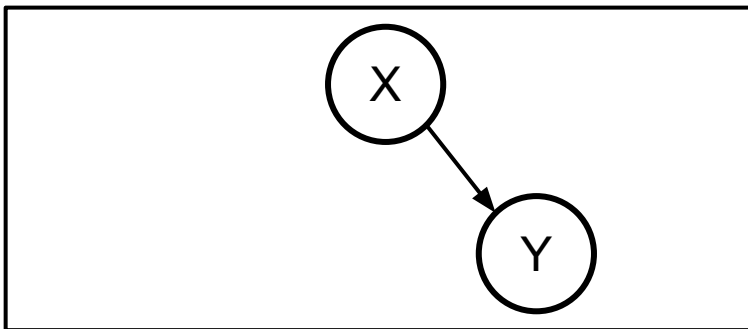
$\mathbf{V}$ : observed variables

$\mathbf{F}$ : structural equations

$\mathbf{U}$ : background variables, with  $P(\mathbf{U})$



This SCM implies this DAG:



*“The building blocks of causal DAGs (directed acyclic graphs)”*

# What does the causal graph look like?

If we want to adjust for confounding, we need to **identify the confounders!**

$$V = \{X, Y, Z\}, U = \{U_X, U_Y, U_Z\}, F = \{f_X, f_Y, f_Z\}$$

$$f_X : X = U_X$$

$$f_Y : Y = 4x + U_Y$$

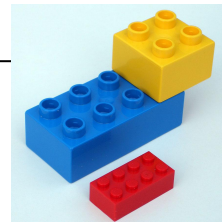
$$f_Z : Z = \frac{x}{10} + U_Z$$

Define SCM:  $\langle \mathbf{V}, \mathbf{F}, \mathbf{U} \rangle$

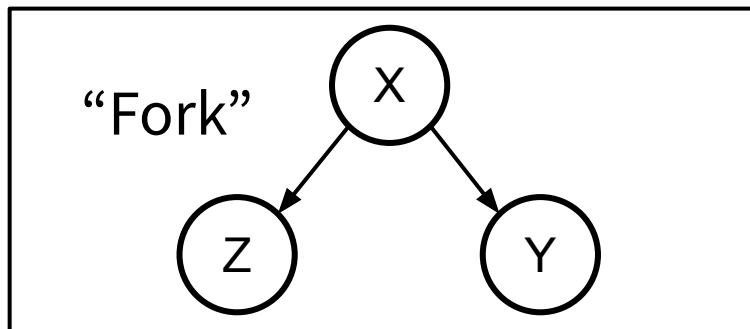
**V**: observed variables

**F**: structural equations

**U**: background variables, with  $P(\mathbf{U})$



This SCM implies this DAG:

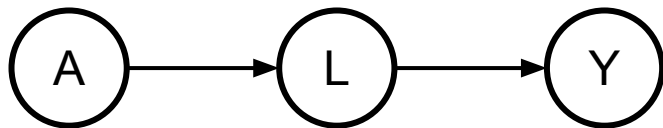


*“The building blocks of causal DAGs (directed acyclic graphs)”*

# What does the causal graph look like?

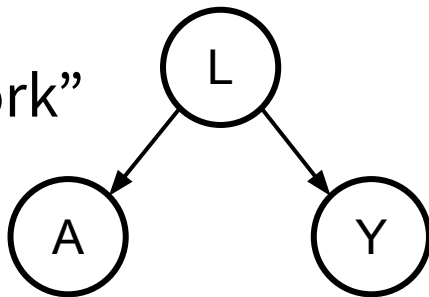
If we want to adjust for confounding, we need to **identify the confounders!**

“Chain”



$V=\{A,L,Y\}$ ,  $F=\{f_A(U_A), f_L(A, U_L), f_Y(L, U_Y)\}$

“Fork”



$V=\{A,L,Y\}$ ,  $F=\{f_A(L, U_A), f_L(U_L), f_Y(L, U_Y)\}$

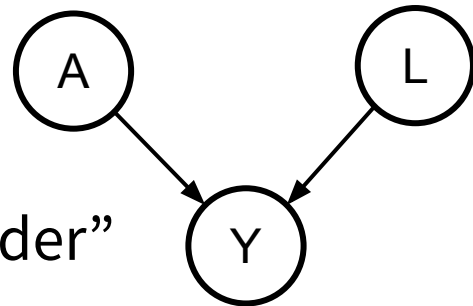
Define SCM:  $\langle V, F, U \rangle$  (which implies a DAG)

**V**: observed variables

**F**: structural equations

**U**: background variables, with  $P(U)$

“Collider”



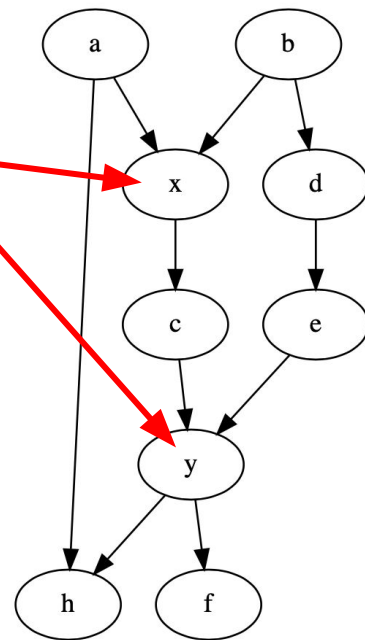
$V=\{A,L,Y\}$ ,  $F=\{f_A(U_A), f_L(U_L), f_Y(L, A, U_Y)\}$

Pearl, Causality, 2009

# Example: Choosing the right adjustment set

Effect of X on Y:

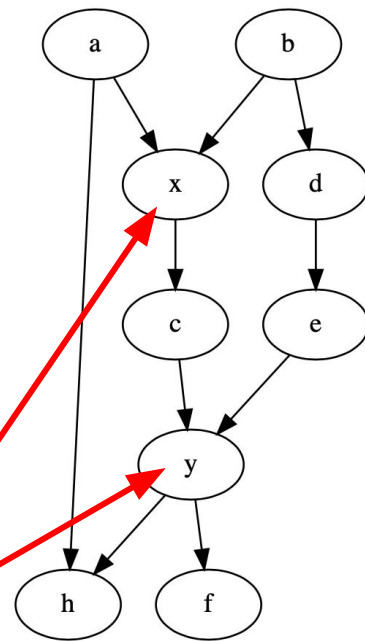
Choosing the  
correct  
**adjustment set Z**  
is critical for  
**unbiased** effect  
estimation.



# Example: Choosing the right adjustment set

**Effect of X on Y:**

Choosing the  
correct  
**adjustment set Z**  
is critical for  
unbiased effect  
estimation.

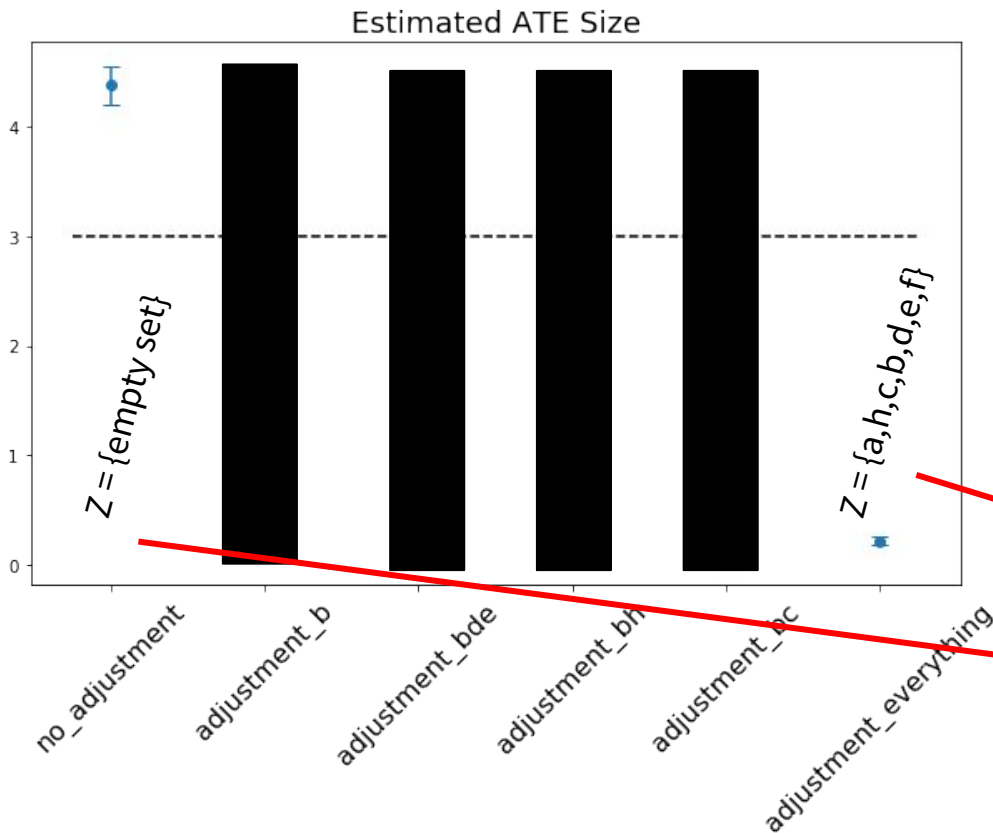


**Backdoor adjusted ATE:**

$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

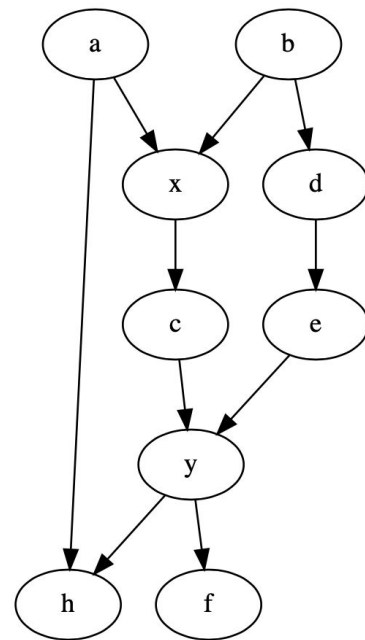


# Example: Choosing the right adjustment set



## Effect of X on Y:

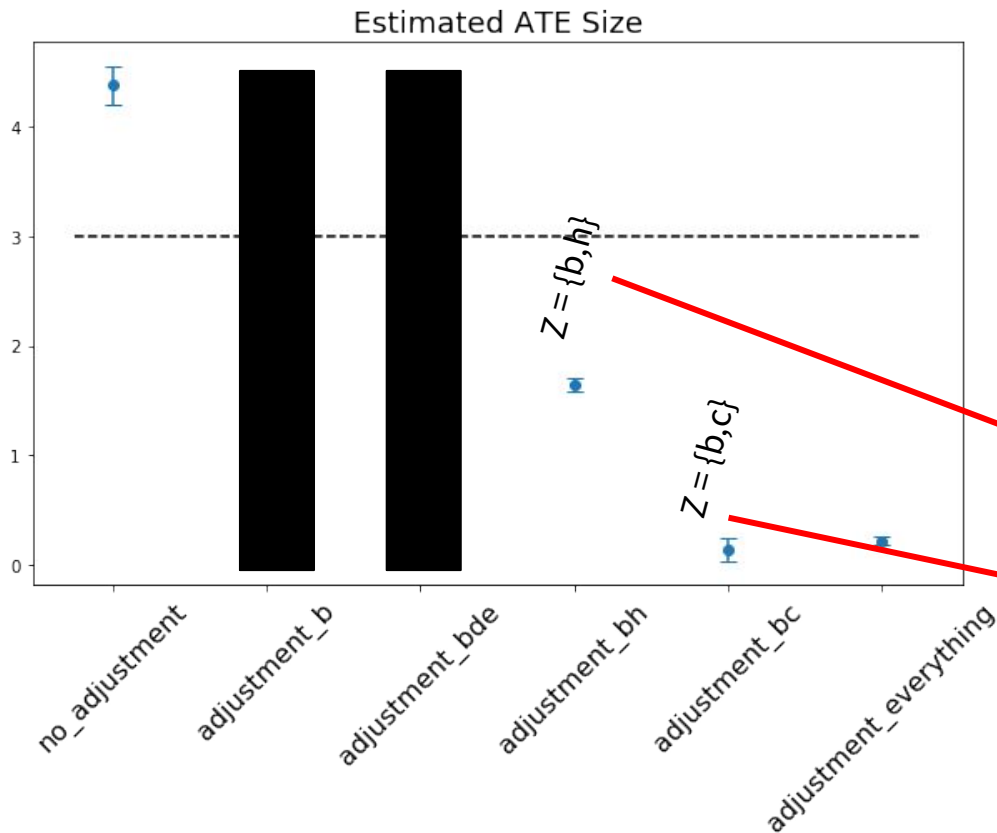
Choosing the correct **adjustment set Z** is critical for unbiased effect estimation.



## Backdoor adjusted ATE:

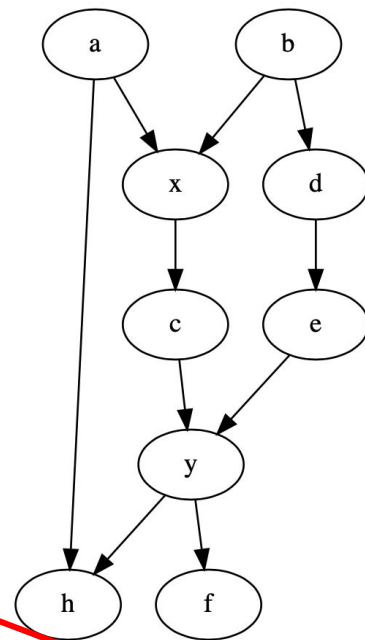
$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Example: Choosing the right adjustment set



## Effect of X on Y:

Choosing the correct **adjustment set Z** is critical for unbiased effect estimation.

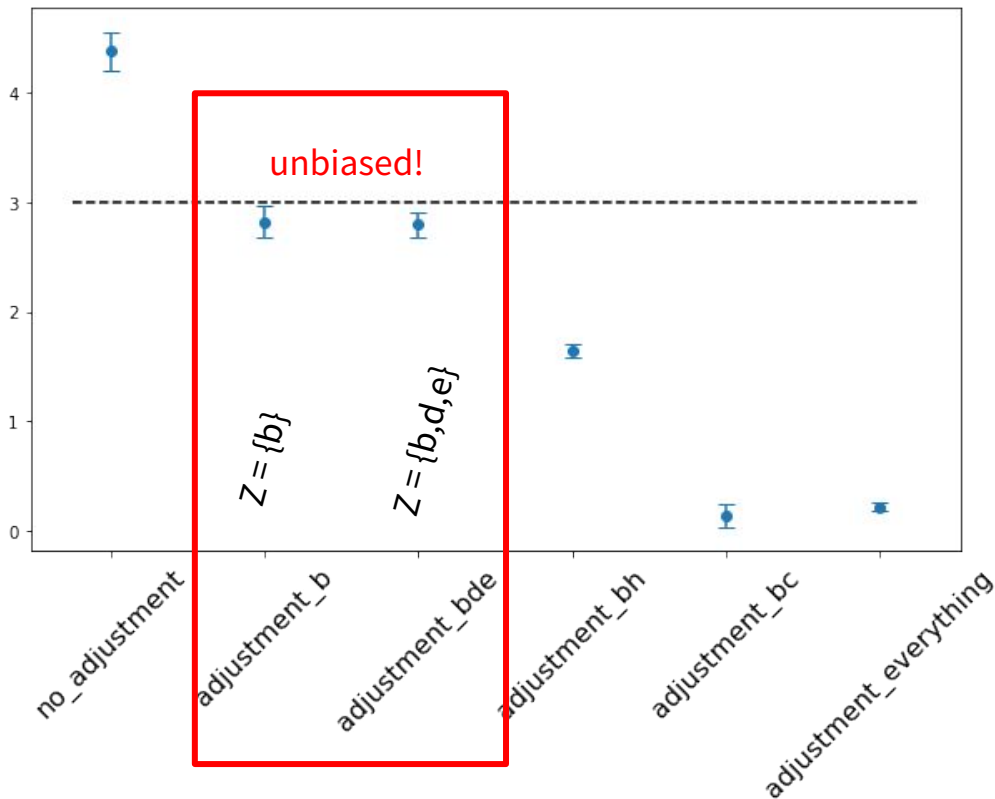


## Backdoor adjusted ATE:

$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

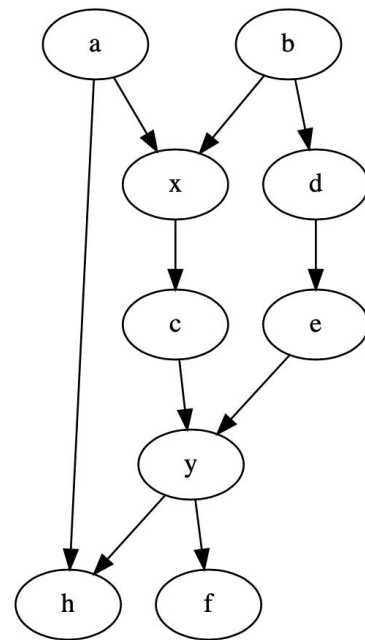
# Example: Choosing the right adjustment set

Estimated ATE Size



## Effect of X on Y:

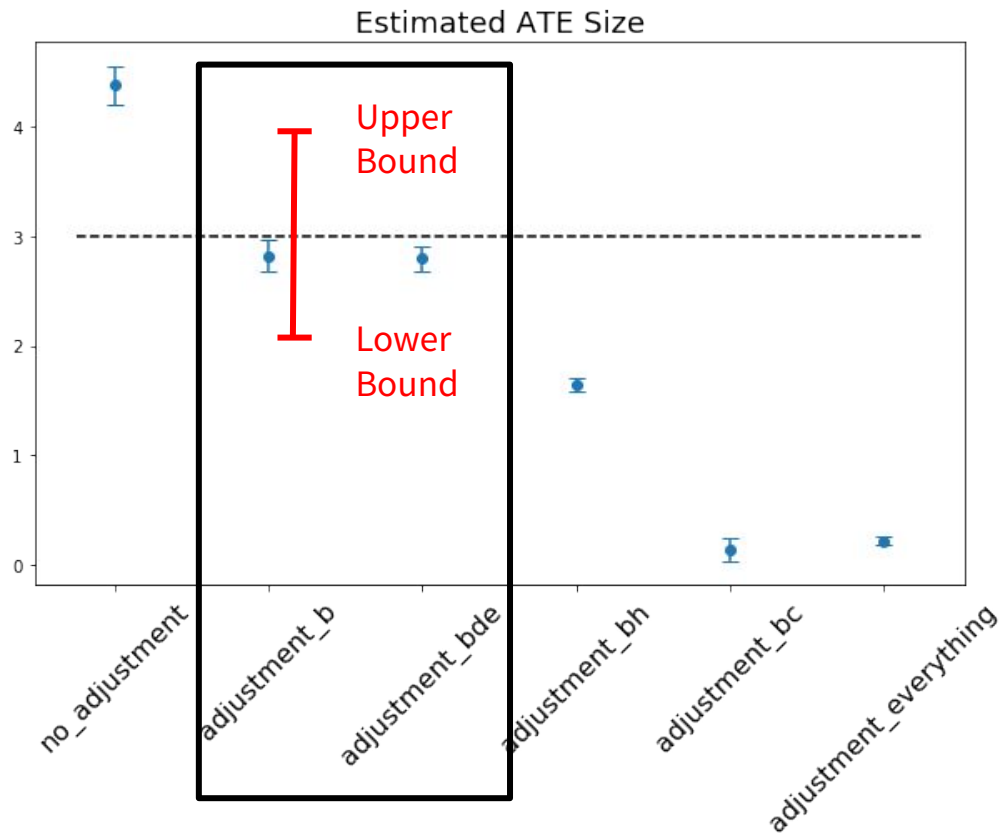
Choosing the correct **adjustment set Z** is critical for unbiased effect estimation.



## Backdoor adjusted ATE:

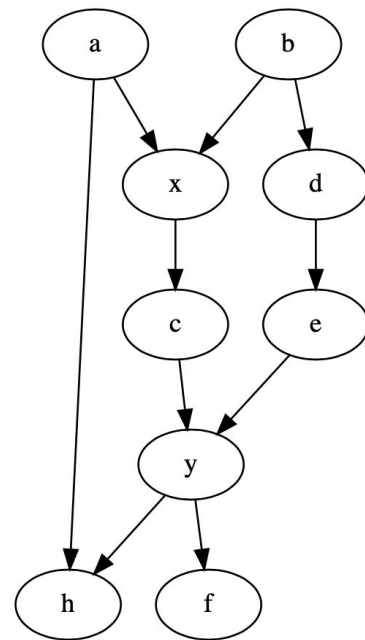
$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Example: Choosing the right adjustment set



## Effect of X on Y:

Choosing the correct **adjustment set Z** is critical for unbiased effect estimation.

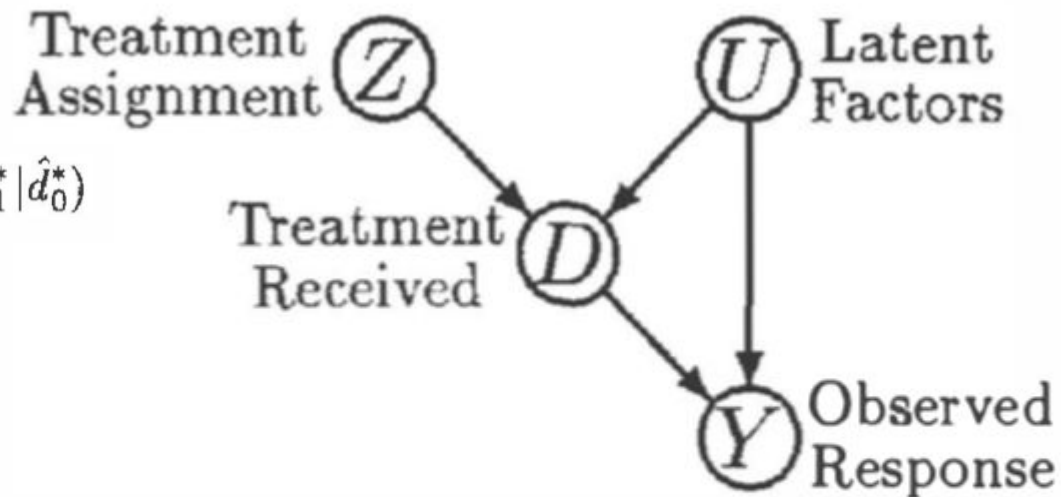


## Backdoor adjusted ATE:

$$P(Y = y | do(X = x)) = \sum_z P(Y = y | X = x, Z = z) P(Z = z)$$

# Effect Estimation

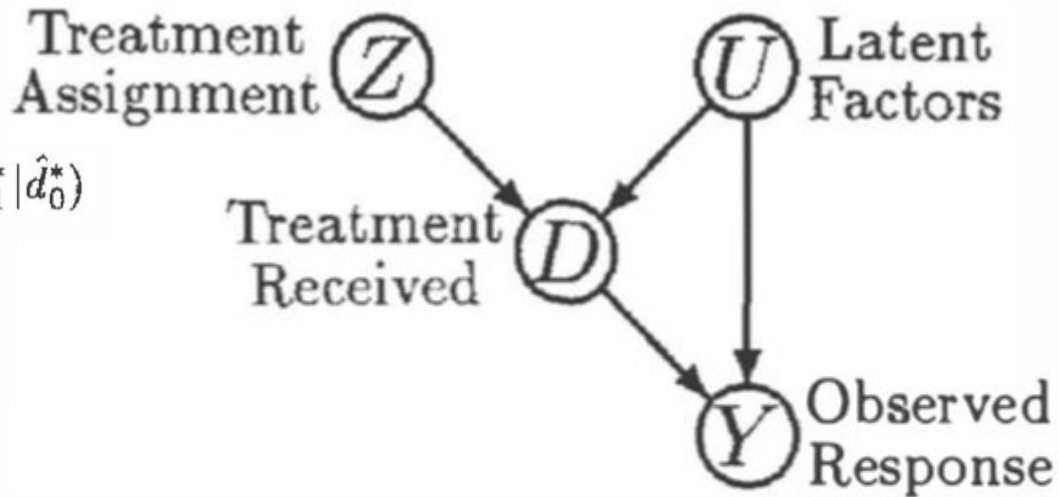
$$\text{ACE}(D \rightarrow Y) = P(y_1^* | \hat{d}_1^*) - P(y_1^* | \hat{d}_0^*)$$



# Effect Estimation

$$\text{ACE}(D \rightarrow Y) = P(y_1^* | \hat{d}_1^*) - P(y_1^* | \hat{d}_0^*)$$

$$P(y^* | \hat{d}^*) = \sum_u P(y | d, u) P(u)$$

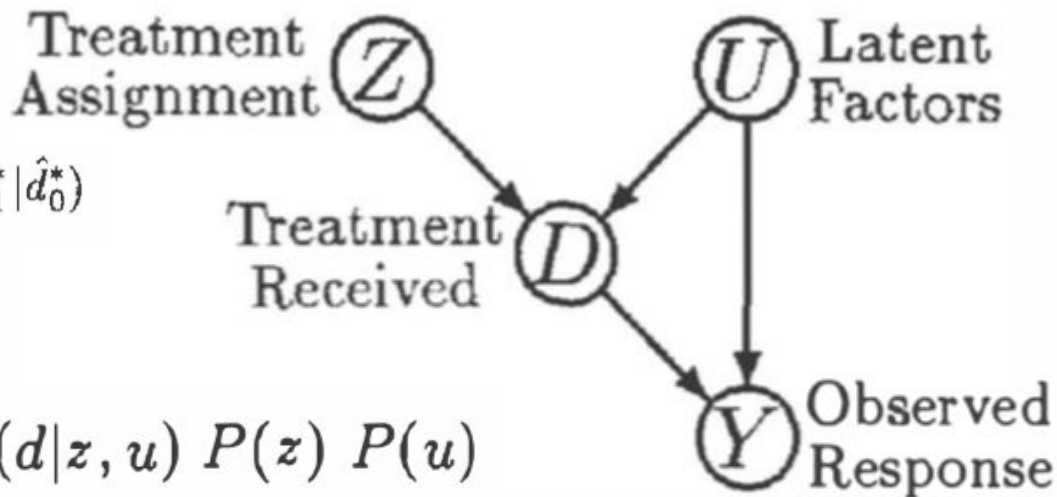


# Effect Estimation

$$\text{ACE}(D \rightarrow Y) = P(y_1^* | \hat{d}_1^*) - P(y_1^* | \hat{d}_0^*)$$

$$P(y^* | \hat{d}^*) = \sum_u P(y | d, u) P(u)$$

$$P(y, d, z, u) = P(y | d, u) P(d | z, u) P(z) P(u)$$

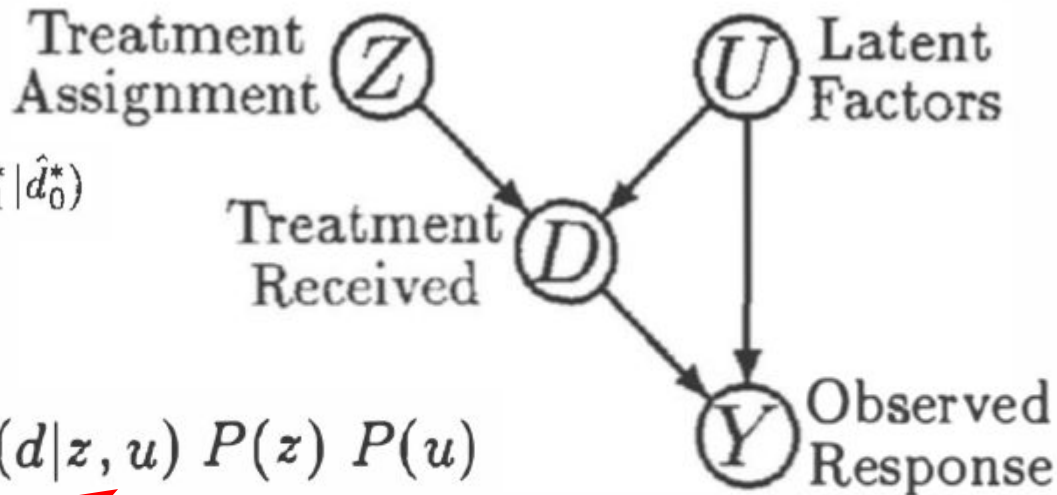


# Effect Estimation

$$\text{ACE}(D \rightarrow Y) = P(y_1^* | \hat{d}_1^*) - P(y_1^* | \hat{d}_0^*)$$

$$P(y^* | \hat{d}^*) = \sum_u P(y | d, u) P(u)$$

$$P(y, d, z, u) = P(y | d, u) P(d | z, u) P(z) P(u)$$



U: unobserved (latent)  
**Unidentifiable!**

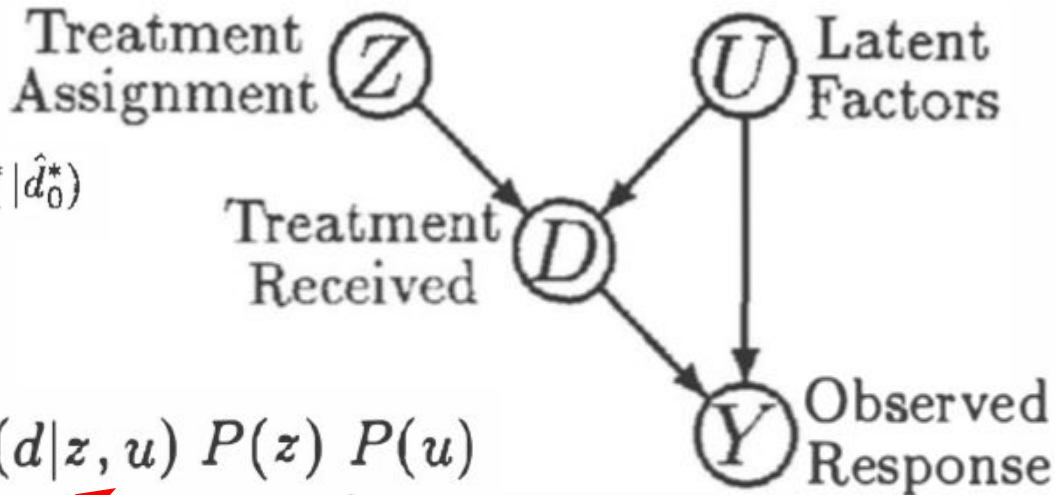


# Effect Estimation

$$\text{ACE}(D \rightarrow Y) = P(y_1^* | \hat{d}_1^*) - P(y_1^* | \hat{d}_0^*)$$

$$P(y^* | \hat{d}^*) = \sum_u P(y | d, u) P(u)$$

$$P(y, d, z, u) = P(y | d, u) P(d | z, u) P(z) P(u)$$



U: unobserved (latent)  
**Unidentifiable!**

**BUT:** Observed marginal

$$P(y, d, z)$$

# Standard Causal Model

$$v = f_V(\text{pa}_V, u_V)$$

**Problem:**  $U_v$  can be anything: “any type with any domain”

Core Message Alert!!1111

# Standard Causal Model

Core Message Alert!!1111

$$v = f_V(\text{pa}_V, u_V)$$

**Problem:**  $U_v$  can be anything: “any type with any domain”

**But:** “For each  $u_v$  of  $U_v$ , the functional mapping from  $\text{PA}_v$  to  $V$  is *particular deterministic response function*.”

# Standard Causal Model

Core Message Alert!!1111

$$v = f_V(\text{pa}_V, u_V)$$

**Problem:**  $U_v$  can be anything: “any type with any domain”

**But:** “For each  $u_v$  of  $U_v$ , the functional mapping from  $PA_v$  to  $V$  is ***particular deterministic response function.***”

**Consequence:** Can map each value of  $U_v$  to a deterministic response function.

**Reason:** Although the domain size of  $U_v$  is unknown, which might be very large or even infinite, the number of different deterministic response functions is **known and limited**, given the domain sizes of  $PA_v$  and  $V$

# Standard Causal Model

Core Message Alert!!1111

$$v = f_V(\text{pa}_V, u_V)$$

**Problem:**  $U_v$  can be anything: “any type with any domain”

**But:** “For each  $u_v$  of  $U_v$ , the functional mapping from  $PA_v$  to  $V$  is ***particular deterministic response function.***”

**Consequence:** Can map each value of  $U_v$  to a deterministic response function.

**Reason:** Although the domain size of  $U_v$  is unknown, which might be very large or even infinite, the number of different deterministic response functions is **known and limited**, given the domain sizes of  $PA_v$  and  $V$

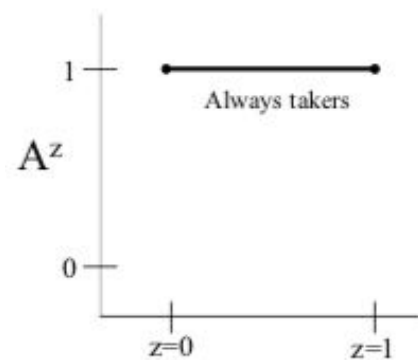


Figure 16.4

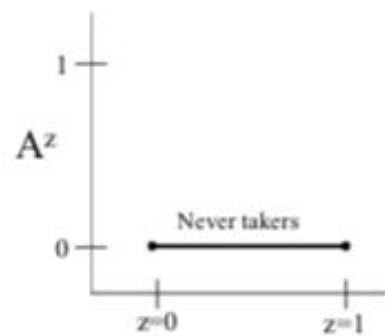


Figure 16.5

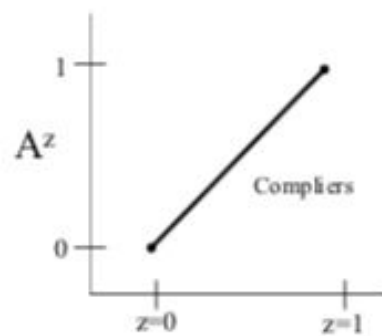


Figure 16.6

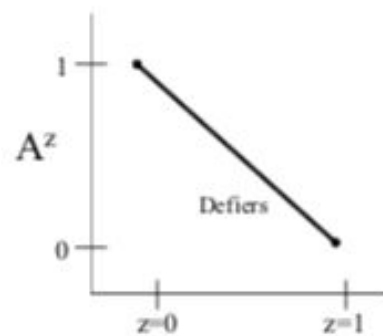


Figure 16.7

$$h_{d,0}(z) = d_0 \quad , \quad h_{d,1}(z) = \begin{cases} d_0 & \text{if } z = z_0 \\ d_1 & \text{if } z = z_1 \end{cases}$$

$$h_{d,3}(z) = d_1 \quad , \quad h_{d,2}(z) = \begin{cases} d_1 & \text{if } z = z_0 \\ d_0 & \text{if } z = z_1 \end{cases}$$



## ***“Replacing the $U$ with the $R$ ”***

$$f_V(\text{pa}_V, u_V) = f_V(\text{pa}_V, \ell_V^{-1}(r_V)) = f_V \circ \ell_V^{-1}(\text{pa}_V, r_V) = g_V(\text{pa}_V, r_V),$$

where  $g_V$  is the composition of  $f_V$  and  $\ell_V^{-1}$

and denotes the response functions represented by  $r_V$ .

$$\mathbb{I}(v; \text{pa}_V, r_V) = \begin{cases} 1 & \text{if } g_V(\text{pa}_V, r_V) = v, \\ 0 & \text{otherwise,} \end{cases}$$



## Different Formulations: *All the same stuff*

$$f_{W_i}(\mathbf{r}) = f_{W_i}(f_{W_{i1}}(\mathbf{r}), \dots, f_{W_{ik_i}}(\mathbf{r}), r_{W_i}).$$

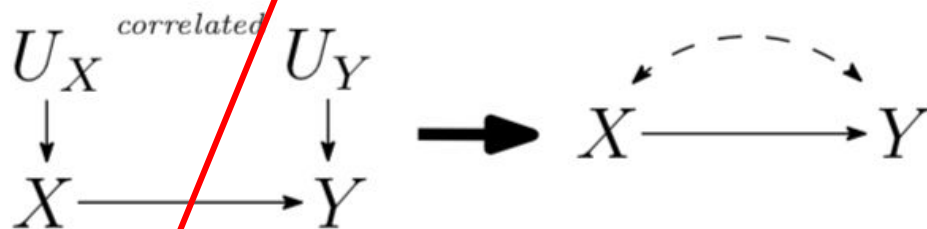
$$\begin{aligned} x_i &= f_{x_i}(\mathbf{r}) \\ &= f_{x_i}(f_{u_1}(\mathbf{r}), f_{u_2}(\mathbf{r}), \dots, f_{u_k}(\mathbf{r}), r_{x_i}) \end{aligned}$$

$$f_V(\mathbf{pa}_V, u_V) = f_V(\mathbf{pa}_V, \ell_V^{-1}(r_V)) = f_V \circ \ell_V^{-1}(\mathbf{pa}_V, r_V) = g_V(\mathbf{pa}_V, r_V),$$

min/max

Objective: e.g. TCE

s.t.  $P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$



**Example:**

$$P(x, y) = \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(x; r_X) \mathbb{I}(y; x, r_Y).$$

$$\text{TCE}(x_1, x_0) = \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_1, r_Y) - \sum_{r_X, r_Y} P(r_X, r_Y) \mathbb{I}(y; x_0, r_Y),$$

# Linear Programming

$$\min/\max \quad P(\hat{y}_{s_1} | \pi, s_0 | \bar{\pi} | \mathbf{o}) - P(\hat{y}_{s_0} | \mathbf{o}),$$

$$\text{s.t.} \quad P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$$

*Wu et al. (2019): The key idea is to*

- 1. parameterize the causal model using so-called response-function variables, whose distribution **captures all randomness encoded in the causal model**,*
- 2. so that we can explicitly **traverse all possible causal models** to find the tightest possible bounds.*

# Linear Programming

$$\min/\max \quad P(\hat{y}_{s_1} | \pi, s_0 | \bar{\pi} | \mathbf{o}) - P(\hat{y}_{s_0} | \mathbf{o}),$$

$$\text{s.t.} \quad P(\mathbf{V}) = P(\mathcal{D}), \quad \sum_{\mathbf{r}} P(\mathbf{r}) = 1, \quad P(\mathbf{r}) \geq 0,$$

*Wu et al. (2019): The key idea is to*

- 1. parameterize the causal model using so-called response-function variables, whose distribution **captures all randomness encoded in the causal model**,*
- 2. so that we can explicitly **traverse all possible causal models** to find the tightest possible bounds.*

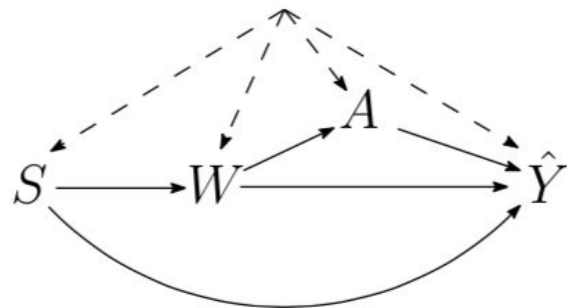
**Recap:** Unidentifiable situation means that there exist two causal models which exactly agree with the same observational distribution.

# Beyond Simple Effects

Everything you can imagine

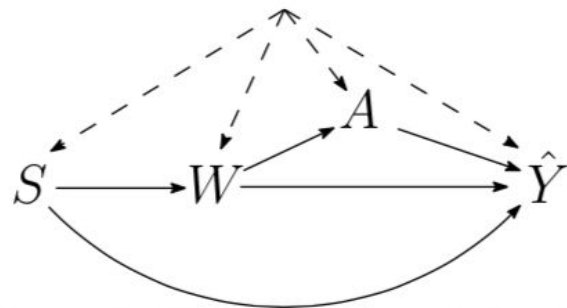
---

$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$



$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$

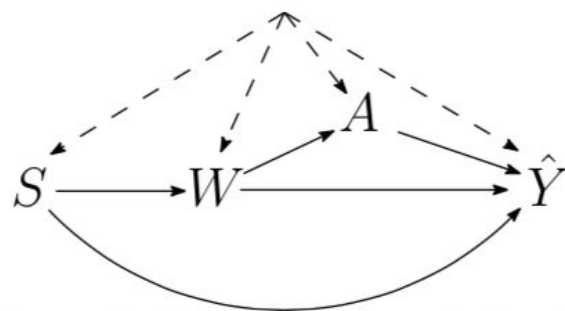
$$\text{PE}_\pi(x_1, x_0) = P(y_{x_1}|\pi, x_0|\bar{\pi}) - P(y_{x_0}),$$



$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$

$$\text{PE}_\pi(x_1, x_0) = P(y_{x_1}|\pi, x_0|\bar{\pi}) - P(y_{x_0}),$$

$$\text{CE}(x_1, x_0|\mathbf{o}) = P(y_{x_1}|\mathbf{o}) - P(y_{x_0}|\mathbf{o}).$$



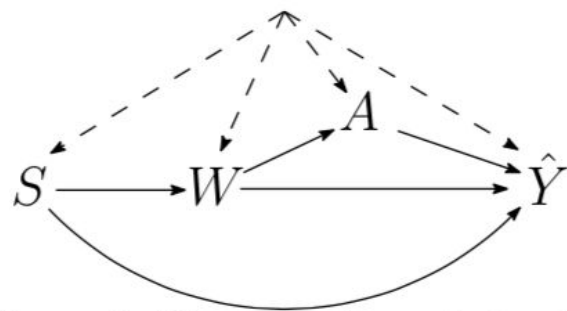


$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$

$$\text{PE}_\pi(x_1, x_0) = P(y_{x_1}|\pi, x_0|\bar{\pi}) - P(y_{x_0}),$$

$$\text{CE}(x_1, x_0|\mathbf{o}) = P(y_{x_1}|\mathbf{o}) - P(y_{x_0}|\mathbf{o}).$$

$$\text{PCE}_\pi(x_1, x_0|\mathbf{o}) = P(y_{x_1}|\pi, x_0|\bar{\pi}|\mathbf{o}) - P(y_{x_0}|\mathbf{o}).$$



$$\pi = \{S \rightarrow W \rightarrow A \rightarrow \hat{Y}, \\ S \rightarrow \hat{Y}\}$$

$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$

$$\text{PE}_\pi(x_1, x_0) = P(y_{x_1} | \pi, x_0 | \bar{\pi}) - P(y_{x_0}),$$

$$\text{CE}(x_1, x_0 | \mathbf{o}) = P(y_{x_1} | \mathbf{o}) - P(y_{x_0} | \mathbf{o}).$$

$$\text{PCE}_\pi(x_1, x_0 | \mathbf{o}) = P(y_{x_1} | \pi, x_0 | \bar{\pi} | \mathbf{o}) - P(y_{x_0} | \mathbf{o}).$$

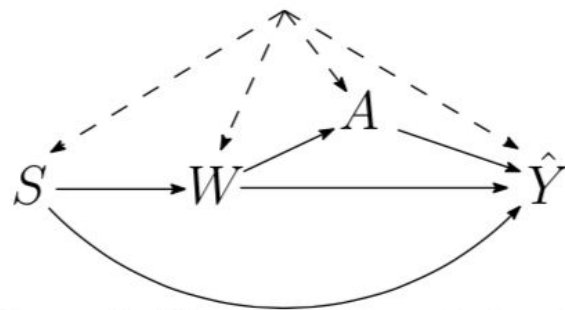


Table 1: Connection between previous fairness notions and PC fairness

Description	References	Relating to PC fairness
<b>Total effect</b>	[19, 16]	$\mathbf{O} = \emptyset$ and $\pi = \Pi$
(System) <b>Direct</b> discrimination	[19, 7, 16]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
(System) <b>Indirect</b> discrimination	[19, 7, 16]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_i \subset \Pi$
Individual <b>direct</b> discrimination	[17]	$\mathbf{O} = \{S, \mathbf{X}\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
Group <b>direct</b> discrimination	[18]	$\mathbf{O} = \mathbf{Q} = \text{PA}_Y \setminus \{S\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
Counterfactual fairness	[5, 9, 14]	$\mathbf{O} = \{S, \mathbf{X}\}$ and $\pi = \Pi$
Counterfactual error rate	[15]	$\mathbf{O} = \{S, Y\}$ and $\pi = \pi_d$ or $\pi_i$

NEW! →

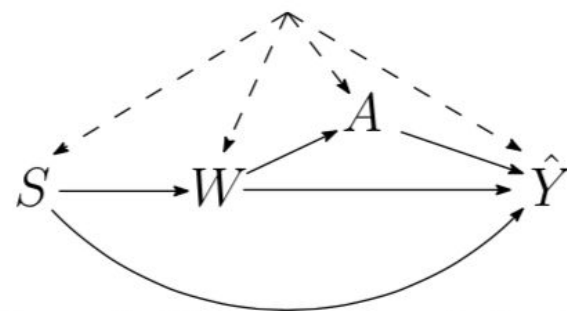


$$\text{TCE}(x_1, x_0) = P(y_{x_1}) - P(y_{x_0}).$$

$$\text{PE}_\pi(x_1, x_0) = P(y_{x_1} | \pi, x_0 | \bar{\pi}) - P(y_{x_0}),$$

$$\text{CE}(x_1, x_0 | \mathbf{o}) = P(y_{x_1} | \mathbf{o}) - P(y_{x_0} | \mathbf{o}).$$

$$\text{PCE}_\pi(x_1, x_0 | \mathbf{o}) = P(y_{x_1} | \pi, x_0 | \bar{\pi} | \mathbf{o}) - P(y_{x_0} | \mathbf{o}).$$



*Choice of  $\mathbf{O}$  and  $\pi$  determines effect! Just plug in into  $\text{PCE}_\pi$*

Table 1: Connection between previous fairness notions and PC fairness

Description	References	Relating to PC fairness
<b>Total</b> effect	[19, 16]	$\mathbf{O} = \emptyset$ and $\pi = \Pi$
(System) <b>Direct</b> discrimination	[19, 7, 16]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
(System) <b>Indirect</b> discrimination	[19, 7, 16]	$\mathbf{O} = \emptyset$ or $\{S\}$ and $\pi = \pi_i \subset \Pi$
Individual <b>direct</b> discrimination	[17]	$\mathbf{O} = \{S, \mathbf{X}\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
Group <b>direct</b> discrimination	[18]	$\mathbf{O} = \mathbf{Q} = \text{PA}_Y \setminus \{S\}$ and $\pi = \pi_d = \{S \rightarrow \hat{Y}\}$
Counterfactual fairness	[5, 9, 14]	$\mathbf{O} = \{S, \mathbf{X}\}$ and $\pi = \Pi$
Counterfactual error rate	[15]	$\mathbf{O} = \{S, Y\}$ and $\pi = \pi_d$ or $\pi_i$

NEW! →



# Partial Identification: **Applications**

*How is partial identification **useful**?*



# Algorithmic Recourse

**Scenario:** Joe wants to know what he needs to change to get his loan approved.



Graph	A	B
	<pre>graph TD; X1((X1)) -.-&gt; X2((X2)); X1((X1)) -.-&gt; X3((X3)); X2((X2)) -.-&gt; X3((X3));</pre>	<pre>graph TD; X1((X1)) -.-&gt; X2((X2)); X1((X1)) --&gt; X3((X3)); X2((X2)) --&gt; X3((X3));</pre>
Confounding	Full	Partial
Recourse	Bounds ( $LB_{FC}, UB_{FC}$ )	Bounds ( $LB_{PC}, UB_{PC}$ )
Example	Joe's loan application will cross the acceptance threshold if his salary increases above €80k.	... if his salary increases above €78k, i.e. $LB_{PB}$ might be tighter due to lack of confounding, 'flipping' earlier.