

COUNTERFACTUAL REASONING

SIMON FERREIRA

simon.ferreira@sorbonne-universite.fr

L'INSTITUT PIERRE LOUIS D'ÉPIDÉMIOLOGIE ET DE SANTÉ
PUBLIQUE, INSERM, SORBONNE UNIVERSITÉ

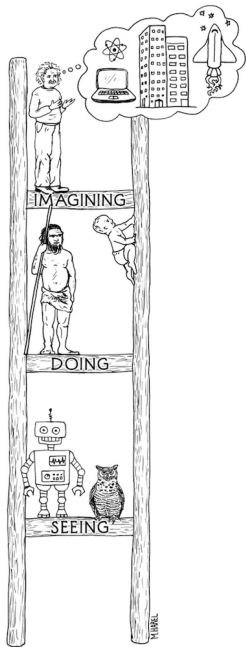
JUNE, 2026



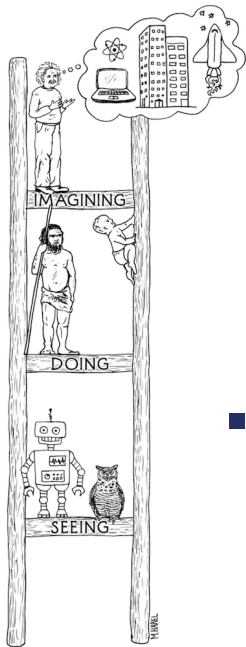
- 1 The Causal Hierarchy
- 2 Examples of Counterfactuals
- 3 Twin Networks
- 4 Single World Intervention Graphs (SWIGs)
- 5 Ancestral Multi World Network (AMWN)
- 6 Conclusion

1

THE CAUSAL HIERARCHY



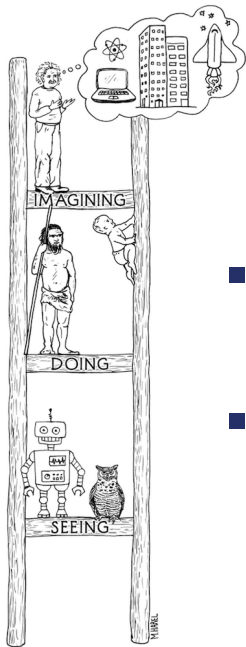
Pearl 2009



■ Associations

- ▶ What if I see...?
- ▶ observational data
- ▶ Tools: statistics, correlation

Pearl 2009



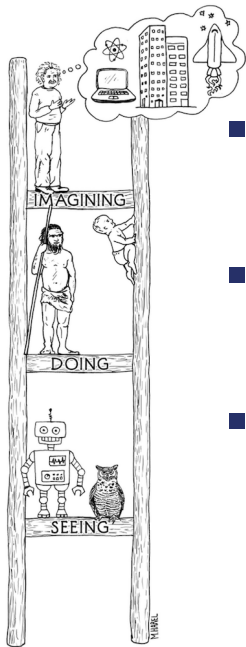
■ Interventions

- ▶ Questions: What if I do...? How?
- ▶ interventional data
- ▶ Tools: RCTs, DAGs, do-calculus

■ Associations

- ▶ What if I see...?
- ▶ observational data
- ▶ Tools: statistics, correlation

Pearl 2009



■ Counterfactuals

- ▶ Questions: What if I had done...? Why?
- ▶ Generating process (SCM)
- ▶ Tools: Twin Networks, SWIGs, AMWN

■ Interventions

- ▶ Questions: What if I do...? How?
- ▶ interventional data
- ▶ Tools: RCTs, DAGs, do-calculus

■ Associations

- ▶ What if I see...?
- ▶ observational data
- ▶ Tools: statistics, correlation

Pearl 2009

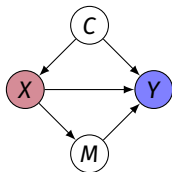
$\xi_C, \xi_X, \xi_M, \xi_Y$

$C := \xi_C$

$X := f_X(C, \xi_X)$

$M := f_M(X, \xi_M)$

$Y := f_Y(C, X, M, \xi_Y)$



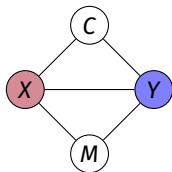
$\xi_C, \xi_X, \xi_M, \xi_Y$

$C := \xi_C$

$X := f_X(C, \xi_X)$

$M := f_M(X, \xi_M)$

$Y := f_Y(C, X, M, \xi_Y)$



■ **Associations:**

$\Pr(Y = y \mid X = x)$

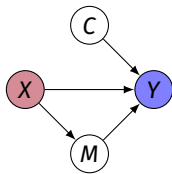
$\xi_C, \xi_X, \xi_M, \xi_Y$

$C := \xi_C$

$X := x$

$M := f_M(X, \xi_M)$

$Y := f_Y(C, X, M, \xi_Y)$



■ **Interventions:**

$\Pr(Y = y \mid \text{do}(X = x))$

■ **Associations:**

$\Pr(Y = y \mid X = x)$

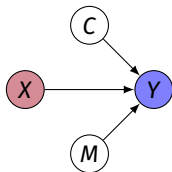
$\xi_C, \xi_X, \xi_M, \xi_Y$

$C := \xi_C$

$X := x$

$M := f_M(x', \xi_M)$

$Y := f_Y(C, X, M, \xi_Y)$



■ Counterfactuals:

$\Pr(Y = y \mid \text{do}(X = x), \text{do}(M = m_{\text{do}(X=x')}))$

■ Interventions:

$\Pr(Y = y \mid \text{do}(X = x))$

■ Associations:

$\Pr(Y = y \mid X = x)$

$$\xi_C, \xi_X, \xi_M, \xi_Y$$

$$C := \xi_C$$

$$X := f_X(C, \xi_X)$$

$$M := f_M(X, \xi_M)$$

$$Y := f_Y(C, X, M, \xi_Y)$$

$$\xi_C, \xi_X, \xi_M, \xi_Y$$

$$C := \xi_C$$

$$X_1 := x'$$

$$M := f_M(X_1, \xi_M)$$

$$X_2 := X$$

$$Y := f_Y(C, X_2, M, \xi_Y)$$

Counterfactuals:

$$\Pr(Y = y \mid \text{do}(X = x), \text{do}(M = m_{\text{do}(X=x')}))$$

2

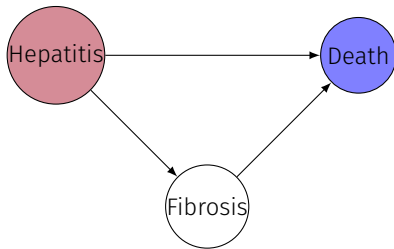
EXAMPLES OF COUNTERFACTUALS

- Effect of the treatment on the treated:

$$\mathbb{E}(Y_{\text{do}(X=1)} \mid X = 1) - \mathbb{E}(Y_{\text{do}(X=0)} \mid X = 1)$$

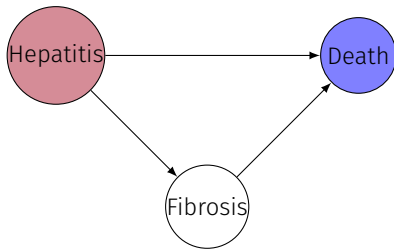
- Effect of the treatment on the treated:
$$\mathbb{E}(Y_{\text{do}(X=1)} \mid X = 1) - \mathbb{E}(Y_{\text{do}(X=0)} \mid X = 1)$$
- Individual effect: “Bob did not take the treatment and died. Would he have lived had he taken the treatment?”

- Effect of the treatment on the treated:
$$\mathbb{E}(Y_{\text{do}(X=1)} | X = 1) - \mathbb{E}(Y_{\text{do}(X=0)} | X = 1)$$
- Individual effect: “Bob did not take the treatment and died. Would he have lived had he taken the treatment?”
- Natural Direct Effect



Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

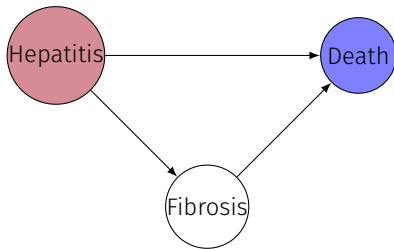
Robins and Greenland 1992; Pearl 2001



Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

$$\begin{aligned}
 CDE(D, H_{0 \rightarrow 1}, F = f) &= \mathbb{E}(D \mid \text{do}(H = 1), \text{do}(F = f)) \\
 &\quad - \mathbb{E}(D \mid \text{do}(H = 0), \text{do}(F = f))
 \end{aligned}$$

This is called the controlled direct effect.



Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

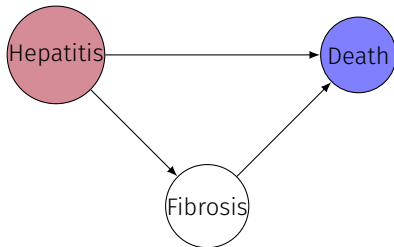
$$\begin{aligned}
 CDE(D, H_{0 \rightarrow 1}, F = f) &= \mathbb{E}(D \mid \text{do}(H = 1), \text{do}(F = f)) \\
 &\quad - \mathbb{E}(D \mid \text{do}(H = 0), \text{do}(F = f))
 \end{aligned}$$

This is called the controlled direct effect.

Problem: This definition depends on the value of f . Thus there exists as many values of CDE as F has possible values.

Robins and Greenland 1992; Pearl 2001

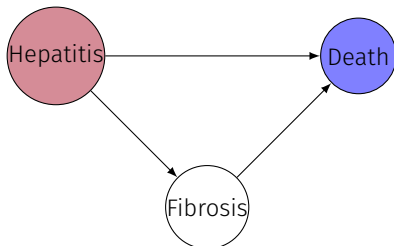
THE NATURAL DIRECT EFFECT



Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

What is the natural value for the mediator Fibrosis?

Robins and Greenland 1992; Pearl 2001

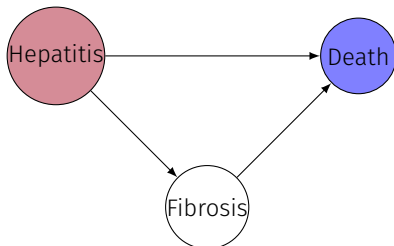


Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

What is the natural value for the mediator Fibrosis?

The natural value for Fibrosis is the value it would have taken had we not modified the exposure (Hepatitis)

Robins and Greenland 1992; Pearl 2001



Direct Effect: How will a change in the exposure (Hepatitis) modify the outcome (Death) if the mediator (Fibrosis) is kept constant?

What is the natural value for the mediator Fibrosis?

The natural value for Fibrosis is the value it would have taken had we not modified the exposure (Hepatitis)

$$\begin{aligned}
 NDE(D, H_{0 \rightarrow 1}, F) &= \mathbb{E} (D \mid \text{do}(H = 1), \text{do}(F = f_{\text{do}(H=0)})) \\
 &\quad - \mathbb{E} (D \mid \text{do}(H = 0), \text{do}(F = f_{\text{do}(H=0)}))
 \end{aligned}$$

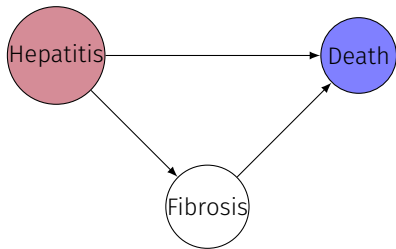
Robins and Greenland 1992; Pearl 2001

$$\xi_H, \xi_F, \xi_D \leftarrow \mathcal{U}(0, 1)$$

$$H := \xi_H$$

$$F := H \vee \xi_F$$

$$D := (H \vee F) \wedge \xi_D$$



observational Data:

Noise			Variables		
ξ_H	ξ_F	ξ_D	H	F	D
0	0	0	0	0	0
0	0	1	0	0	0
0	1	0	0	1	0
0	1	1	0	1	1
1	0	0	1	1	0
1	0	1	1	1	1
1	1	0	1	1	0
1	1	1	1	1	1

$$\begin{aligned} & \mathbb{E}(D \mid H = 1, F = 0) - \mathbb{E}(D \mid H = 0, F = 0) \\ &= 0 - 0 = 0 \end{aligned}$$

TOY EXAMPLE

interventional Data:

Noise			do ($H = 0, F = 0$)			do ($H = 1, F = 0$)		
ξ_H	ξ_F	ξ_D	H	F	D	H	F	D
0	0	0	0	0	0	1	0	0
0	0	1	0	0	0	1	0	1
0	1	0	0	0	0	1	0	0
0	1	1	0	0	0	1	0	1
1	0	0	0	0	0	1	0	0
1	0	1	0	0	0	1	0	1
1	1	0	0	0	0	1	0	0
1	1	1	0	0	0	1	0	1

Controlled direct effect (CDE)

$$\begin{aligned} & \mathbb{E}(D \mid \text{do}(H = 1, F = 0)) - \mathbb{E}(D \mid \text{do}(H = 0, F = 0)) \\ &= \frac{1}{2} - 0 = \frac{1}{2} \end{aligned}$$

Counterfactual Data:

Noise			Counterfactual Data:					
			do ($H = 0, F = f_{do(H=0)}$)			do ($H = 1, F = f_{do(H=0)}$)		
ξ_H	ξ_F	ξ_D	H	F	D	H	F	D
0	0	0	0	0	0	1	0	0
0	0	1	0	0	0	1	0	1
0	1	0	0	1	0	1	1	0
0	1	1	0	1	1	1	1	1
1	0	0	0	0	0	1	0	0
1	0	1	0	0	0	1	0	1
1	1	0	0	1	0	1	1	0
1	1	1	0	1	1	1	1	1

Natural direct effect:

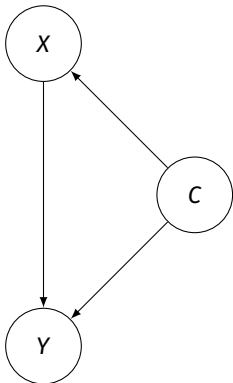
$$\begin{aligned} & \mathbb{E}(D \mid \text{do}(H = 1, F = f_{\text{do}(H=0)})) - \mathbb{E}(D \mid \text{do}(H = 0, F = f_{\text{do}(H=0)})) \\ &= \frac{1}{2} - \frac{1}{4} = \frac{1}{4} \end{aligned}$$

3

TWIN NETWORKS

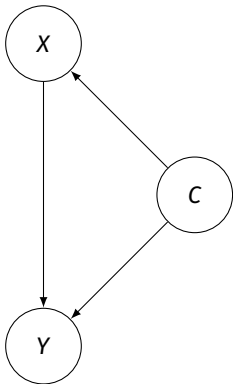
- One graph for the actual world

- One graph for the actual world



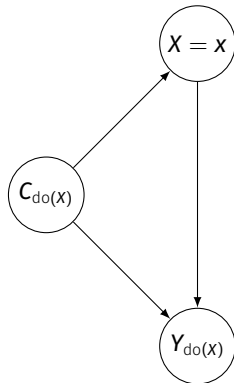
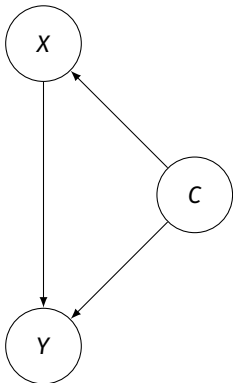
Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world



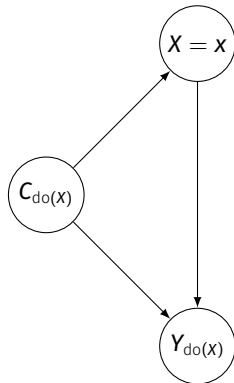
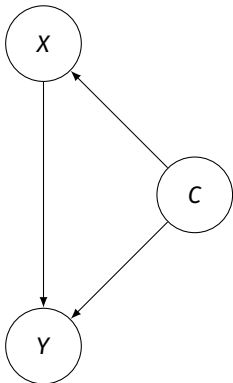
Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world



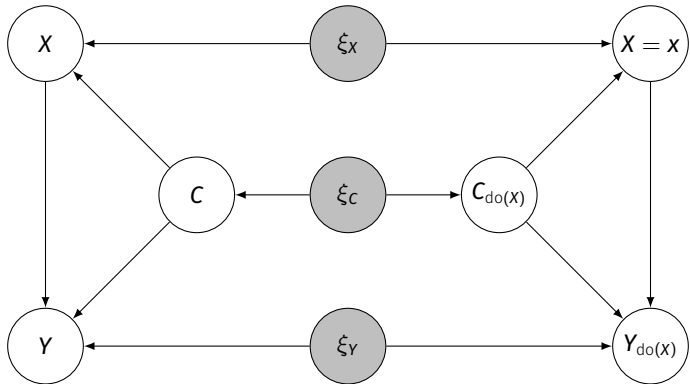
Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world
- Common noise variables



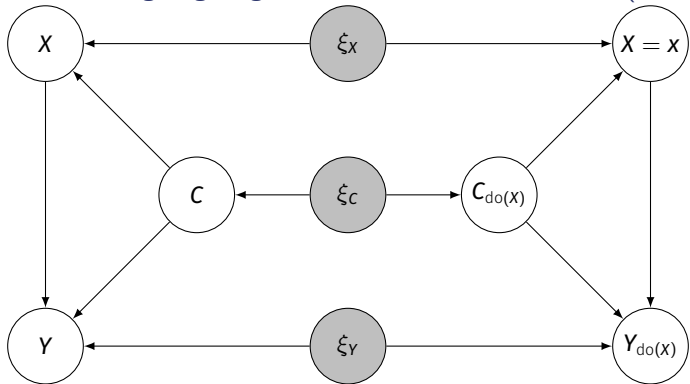
Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world
- Common noise variables



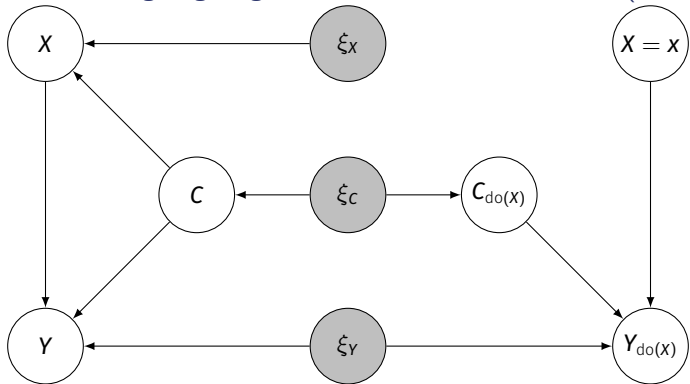
Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world
- Common noise variables
- Remove edges going in the intervened node ($\text{do}(X = x)$)



Balke and Pearl 1994

- One graph for the actual world
- One graph for the counterfactual world
- Common noise variables
- Remove edges going in the intervened node ($\text{do}(X = x)$)



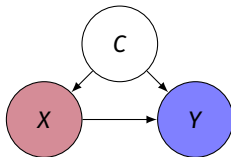
Balke and Pearl 1994

What can Twin Networks be used for?

What can Twin Networks be used for?

SCM framework (Pearl)

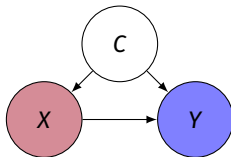
- Assume the graph
- Backdoor criterion: C blocks all backdoor paths from X to Y .



What can Twin Networks be used for?

SCM framework (Pearl)

- Assume the graph
- Backdoor criterion: C blocks all backdoor paths from X to Y .



Potential Outcome framework (Robins)

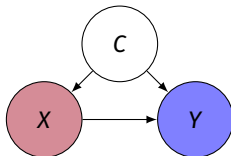
- Assume conditional exchangeability

$$Y_{\text{do}(X)} \perp\!\!\!\perp_{\text{Pr}} X \mid C$$

What can Twin Networks be used for?

SCM framework (Pearl)

- Assume the graph
- Backdoor criterion: C blocks all backdoor paths from X to Y .



Potential Outcome framework (Robins)

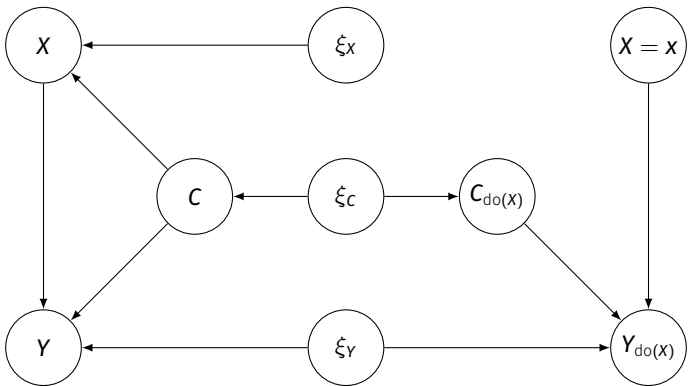
- Assume conditional exchangeability

$$Y_{\text{do}(X)} \perp\!\!\!\perp_{\text{Pr}} X \mid C$$

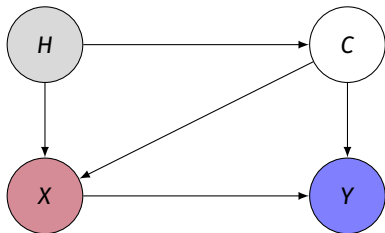
Average Treatment Effect:

$$\sum_c (\mathbb{E}(Y \mid X = 1, C = c) - \mathbb{E}(Y \mid X = 0, C = c)) \Pr(C = c)$$

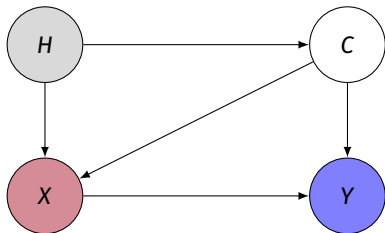
Conditional Exchangeability: $Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$



TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS

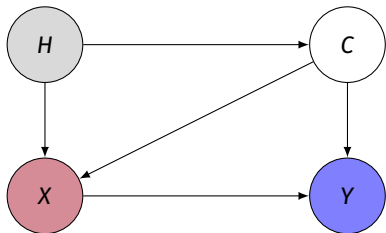


TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS

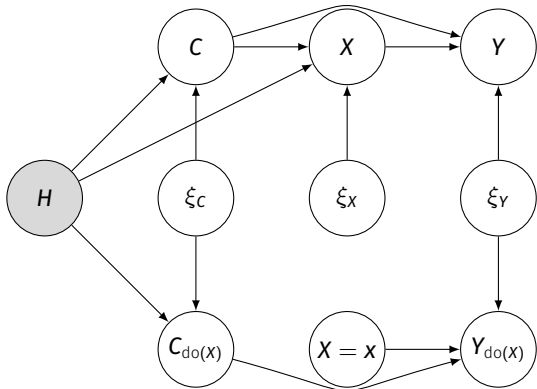


Does C verify the backdoor criterion for X, Y ?

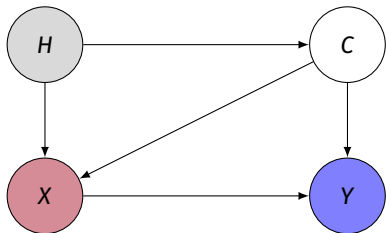
TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS



Does C verify the backdoor criterion for X, Y ?

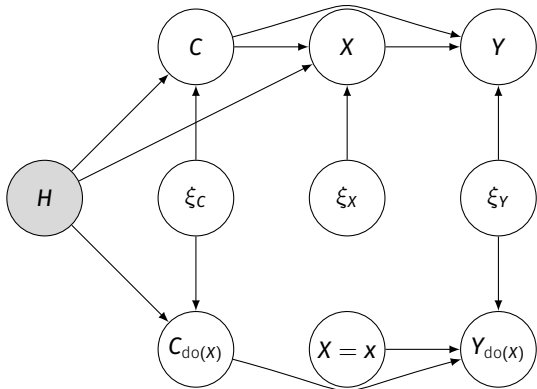


TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS

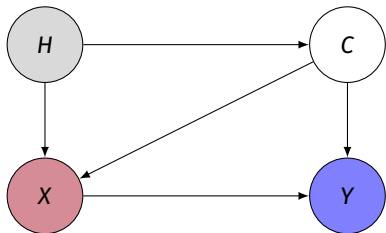


Does C verify the backdoor criterion for X, Y ?

$$Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$$

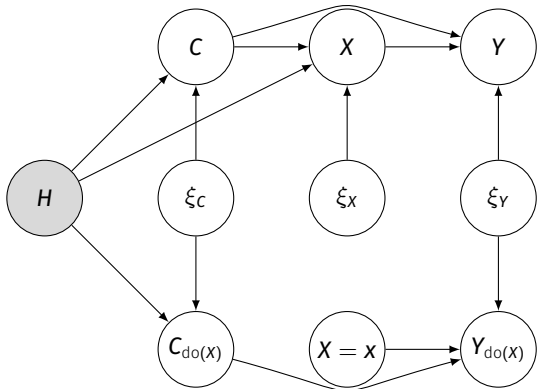


TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS

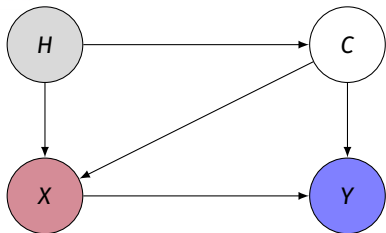


Does C verify the backdoor criterion for X, Y ?

$Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$
No!

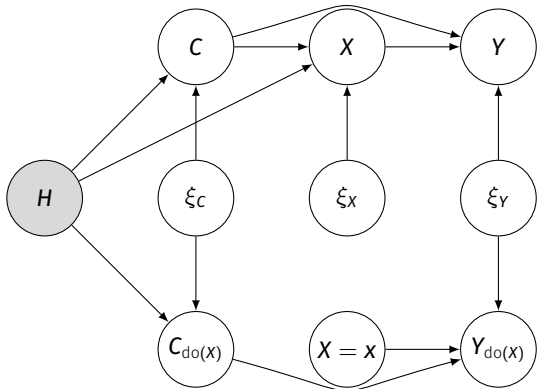


TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS

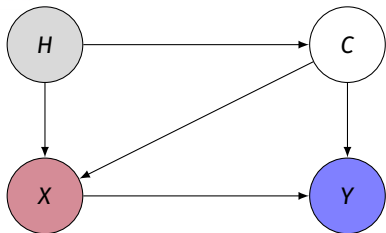


Does C verify the backdoor criterion for X, Y ?

$Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$
 No! $Y_{do(X)} \perp\!\!\!\perp_d X \mid C, C_{do(X)}?$

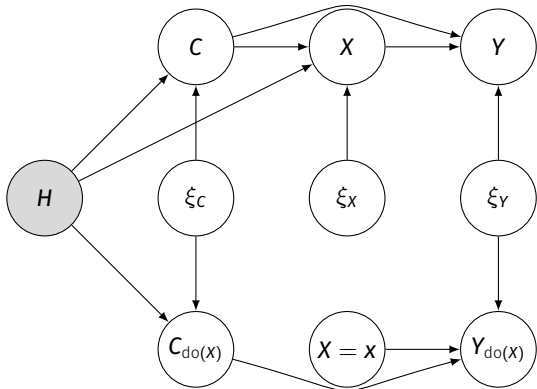


TWIN NETWORKS: AN EXAMPLE OF NON-COMPLETENESS



Does C verify the backdoor criterion for X, Y ?

$Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$
 No! $Y_{do(X)} \perp\!\!\!\perp_d X \mid C, C_{do(X)}?$ Yes!



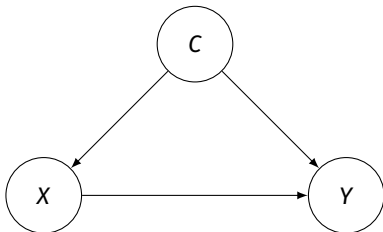
4

SINGLE WORLD
GRAPHS (SWIGs)

INTERVENTION

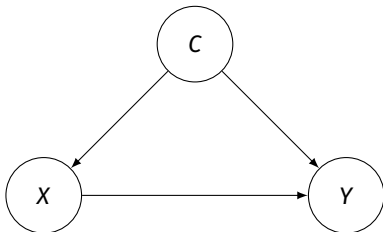
- One graph for both worlds

- One graph for both worlds



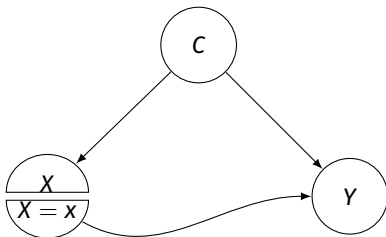
Richardson and Robins 2013

- One graph for both worlds
- Divide the intervened nodes ($\text{do}(X = x)$)



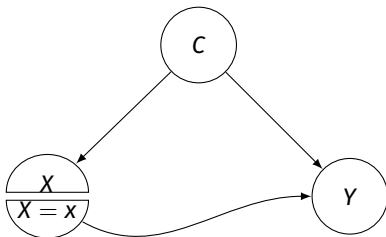
Richardson and Robins 2013

- One graph for both worlds
- Divide the intervened nodes (do ($X = x$))



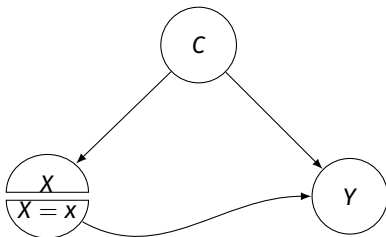
Richardson and Robins 2013

- One graph for both worlds
- Divide the intervened nodes ($\text{do}(X = x)$)
- Notice the edges



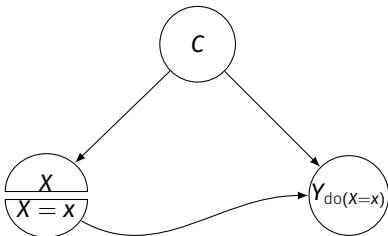
Richardson and Robins 2013

- One graph for both worlds
- Divide the intervened nodes ($\text{do}(X = x)$)
- Notice the edges
- Relabel nodes descendant of interventions

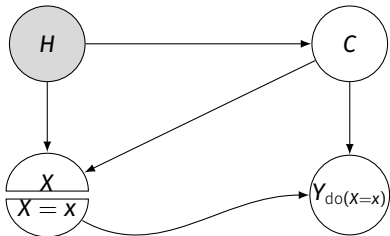
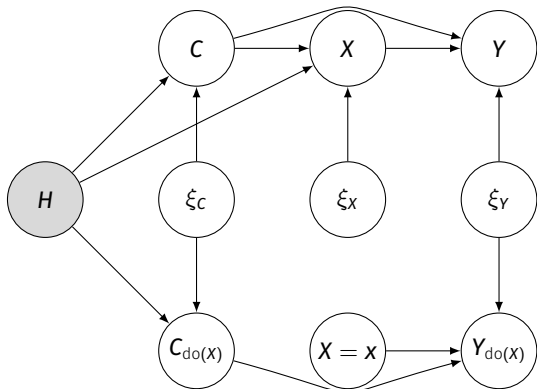


Richardson and Robins 2013

- One graph for both worlds
- Divide the intervened nodes ($\text{do}(X = x)$)
- Notice the edges
- Relabel nodes descendant of interventions



Richardson and Robins 2013



$$Y_{do(X)} \perp\!\!\!\perp_d X \mid C?$$

The do-calculus is sound
 (and complete for interventional queries)

Rule 1 $\Pr(y \mid \text{do}(x), z, w) = \Pr(y \mid \text{do}(x), w)$

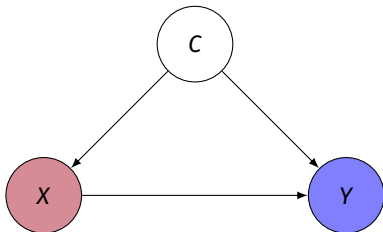
if $(Y \perp\!\!\!\perp_d Z \mid X, W)_{\mathcal{G}_{\overline{X}}}$

Rule 2 $\Pr(y \mid \text{do}(x, z), w) = \Pr(y \mid \text{do}(x), z, w)$

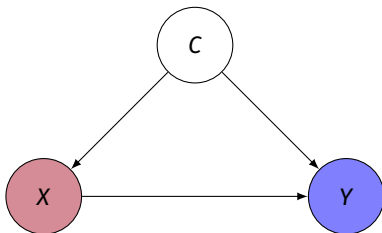
if $(Y \perp\!\!\!\perp_d Z \mid X, W)_{\mathcal{G}_{\overline{XZ}}}$

Rule 3 $\Pr(y \mid \text{do}(x, z), w) = \Pr(y \mid \text{do}(x), w)$

if $(Y \perp\!\!\!\perp_d Z \mid X, W)_{\mathcal{G}_{\overline{XZ}(w)}}$

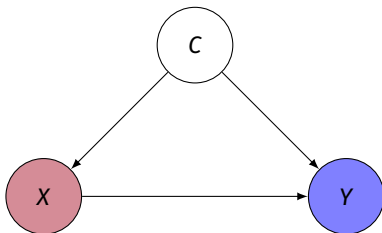


$$\text{ATE: } \mathbb{E}(Y \mid \text{do}(X = x_1)) - \mathbb{E}(Y \mid \text{do}(X = x_0))$$



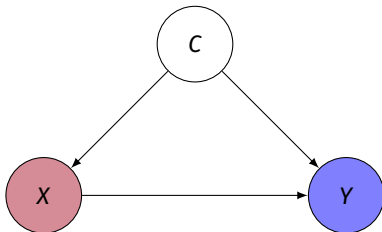
$$\text{ATE: } \mathbb{E}(Y \mid \text{do}(X = x_1)) - \mathbb{E}(Y \mid \text{do}(X = x_0))$$

$$\Pr(Y \mid \text{do}(X))$$



$$\text{ATE: } \mathbb{E}(Y \mid \text{do}(X = x_1)) - \mathbb{E}(Y \mid \text{do}(X = x_0))$$

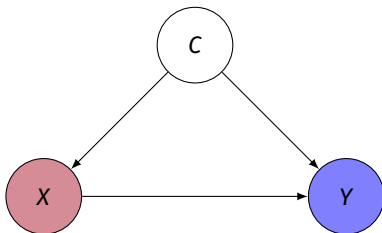
$$\Pr(Y \mid \text{do}(X)) = \sum_c \Pr(Y \mid \text{do}(X), C = c) \Pr(C = c \mid \text{do}(X))$$



$$\text{ATE: } \mathbb{E}(Y \mid \text{do}(X = x_1)) - \mathbb{E}(Y \mid \text{do}(X = x_0))$$

$$\Pr(Y \mid \text{do}(X)) = \sum_c \Pr(Y \mid \text{do}(X), C = c) \Pr(C = c \mid \text{do}(X))$$

$$\text{(Rule 2)} = \sum_c \Pr(Y \mid X, C = c) \Pr(C = c \mid \text{do}(X))$$



$$\text{ATE: } \mathbb{E}(Y \mid \text{do}(X = x_1)) - \mathbb{E}(Y \mid \text{do}(X = x_0))$$

$$\Pr(Y \mid \text{do}(X)) = \sum_c \Pr(Y \mid \text{do}(X), C = c) \Pr(C = c \mid \text{do}(X))$$

$$\text{(Rule 2)} = \sum_c \Pr(Y \mid X, C = c) \Pr(C = c \mid \text{do}(X))$$

$$\text{(Rule 3)} = \sum_c \Pr(Y \mid X, C = c) \Pr(C = c)$$

$$\text{Rule 1 } \Pr(Y_{\text{do}(x)} \mid Z_{\text{do}(x)}, W_{\text{do}(x)}) = \Pr(Y_{\text{do}(x)} \mid W_{\text{do}(x)})$$

$$\text{if } (Y_{\text{do}(x)} \perp\!\!\!\perp_d Z_{\text{do}(x)} \mid W_{\text{do}(x)})_{\mathcal{G}_{\text{do}(x)}}$$

$$\text{Rule 2 } \Pr(Y_{\text{do}(x,z)} \mid W_{\text{do}(x,z)}) = \Pr(Y_{\text{do}(x)} \mid Z_{\text{do}(x)}, W_{\text{do}(x)})$$

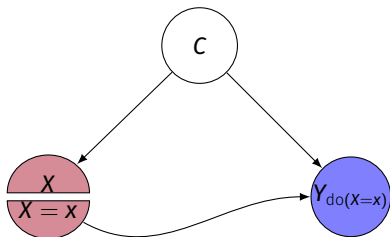
$$\text{if } (Y_{\text{do}(x,z)} \perp\!\!\!\perp_d Z_{\text{do}(x,z)} \mid W_{\text{do}(x,z)})_{\mathcal{G}_{\text{do}(x,z)}}$$

$$\text{Rule 3 } \Pr(Y_{\text{do}(x,z)} \mid W_{\text{do}(x,z)}) = \Pr(Y_{\text{do}(x)} \mid W_{\text{do}(x)})$$

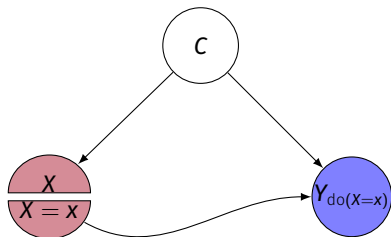
$$\text{if } (Y_{\text{do}(x,z_1)}, W_{\text{do}(x,z_1)} \perp\!\!\!\perp_d Z_1)_{\mathcal{G}_{\text{do}(x,z_1)}}$$

$$\text{and } (Y_{\text{do}(x,z_1)} \perp\!\!\!\perp_d Z_2 \mid W_{\text{do}(x,z_1)})_{\mathcal{G}_{\text{do}(x,z_1)}}$$

$$\text{where } Z_1 = Z \setminus \text{An}(W, \mathcal{G}_{\text{do}(x)}) \text{ and } Z_2 = Z \cap \text{An}(W, \mathcal{G}_{\text{do}(x)})$$

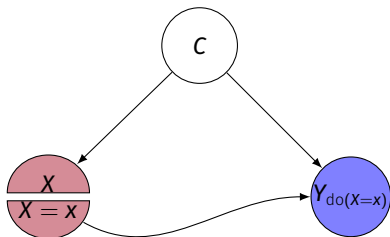


$$\text{ATE: } \mathbb{E} (Y_{do(X=x_1)}) - \mathbb{E} (Y_{do(X=x_0)})$$



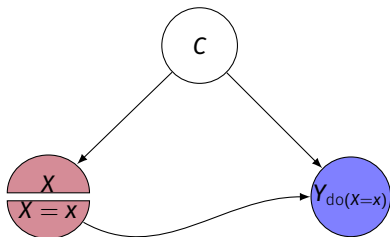
$$\text{ATE: } \mathbb{E} (Y_{\text{do}(X=x_1)}) - \mathbb{E} (Y_{\text{do}(X=x_0)})$$

$$\Pr (Y_{\text{do}(x)})$$



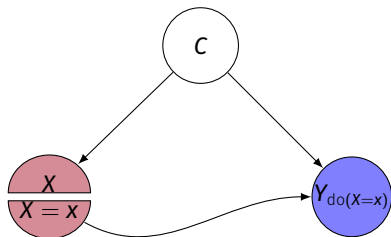
$$\text{ATE: } \mathbb{E} (Y_{\text{do}(X=x_1)}) - \mathbb{E} (Y_{\text{do}(X=x_0)})$$

$$\Pr (Y_{\text{do}(X)}) = \sum_c \Pr (Y_{\text{do}(X)} \mid C_{\text{do}(X)} = c) \Pr (C_{\text{do}(X)} = c)$$



$$\text{ATE: } \mathbb{E} (Y_{do(X=x_1)}) - \mathbb{E} (Y_{do(X=x_0)})$$

$$\begin{aligned} \Pr (Y_{do(X)}) &= \sum_c \Pr (Y_{do(X)} \mid C_{do(X)} = c) \Pr (C_{do(X)} = c) \\ \text{(Rule 2)} &= \sum_c \Pr (Y \mid X, C = c) \Pr (C_{do(X)} = c) \end{aligned}$$



$$\text{ATE: } \mathbb{E} (Y_{do(X=x_1)}) - \mathbb{E} (Y_{do(X=x_0)})$$

$$\Pr (Y_{do(X)}) = \sum_c \Pr (Y_{do(X)} \mid C_{do(X)} = c) \Pr (C_{do(X)} = c)$$

$$\text{(Rule 2)} = \sum_c \Pr (Y \mid X, C = c) \Pr (C_{do(X)} = c)$$

$$\text{(Rule 3)} = \sum_c \Pr (Y \mid X, C = c) \Pr (C = c)$$

SWIGs: Single World Interventional Graphs

SWIGs: **Single** World Interventional Graphs

SWIGs: Single World Interventional Graphs

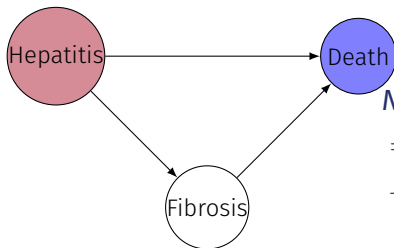
What is the effect of in the population that would die if they do not take it?

$$\mathbb{E} (Y_{\text{do}(X=1)} \mid Y_{\text{do}(X=0)} = 0)$$

SWIGs: Single World Interventional Graphs

What is the effect of in the population that would die if they do not take it?

$$\mathbb{E}(Y_{\text{do}(X=1)} \mid Y_{\text{do}(X=0)} = 0)$$

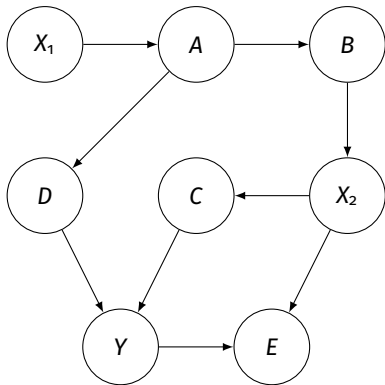


$$\begin{aligned}
 &NDE(D, H_{0 \rightarrow 1}, F) \\
 &= \mathbb{E}(D \mid \text{do}(H = 1, F = f_{\text{do}(H=0)})) \\
 &- \mathbb{E}(D \mid \text{do}(H = 0, F = f_{\text{do}(H=0)}))
 \end{aligned}$$

5

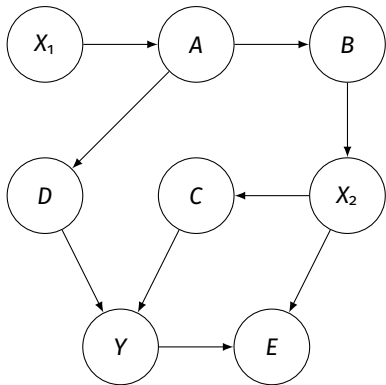
ANCESTRAL MULTI WORLD NETWORK
(AMWN)

ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$\begin{aligned}
 An(Y_{do(x)}) &= \{W_{do(z)} \mid \\
 W &\in An(Y, G_{\bar{X}}) \setminus X, \\
 z &= x \cap An(W, G_{\bar{X}})\}
 \end{aligned}$$

ANCESTOR OF A COUNTERFACTUAL VARIABLE



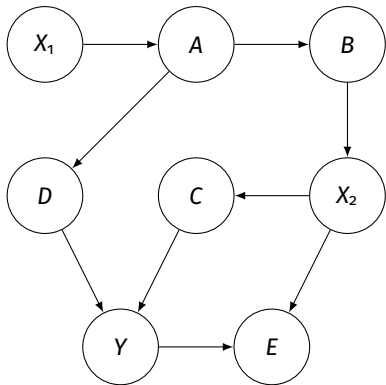
$$An(Y_{do(\mathcal{X})}) = \{W_{do(\mathcal{Z})} \mid$$

$$W \in An(Y, G_{\overline{\mathcal{X}}}) \setminus \mathcal{X},$$

$$\mathcal{Z} = \mathcal{X} \cap An(W, G_{\overline{\mathcal{X}}})\}$$

$$An(Y_{do(X_1, X_2)})?$$

ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(\mathbf{x})}) = \{W_{do(\mathbf{z})} \mid$$

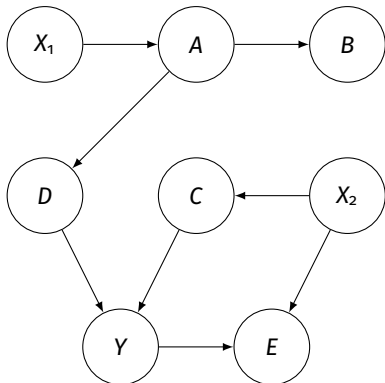
$$W \in An(Y, G_{\overline{\mathbf{X}}}) \setminus \mathbf{X},$$

$$\mathbf{z} = \mathbf{x} \cap An(W, G_{\overline{\mathbf{X}}})\}$$

$An(Y_{do(X_1, X_2)})?$

$\cdot G_{\overline{X_1, X_2}}$

ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(x)}) = \{W_{do(z)} \mid$$

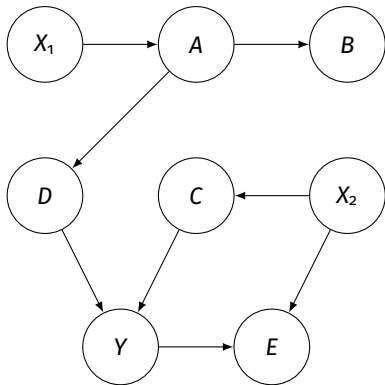
$$W \in An(Y, G_{\overline{X}}) \setminus X,$$

$$z = x \cap An(W, G_{\overline{X}})\}$$

$An(Y_{do(X_1, X_2)})?$

$\cdot G_{\overline{X_1, X_2}}$

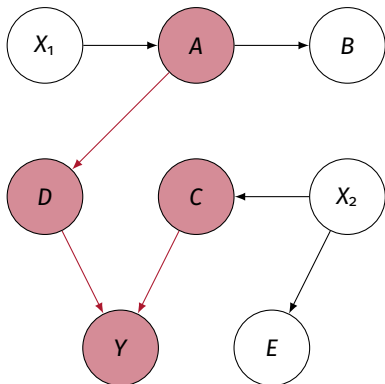
ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(x)}) = \{W_{do(z)} \mid W \in An(Y, G_{\bar{x}}) \setminus \mathcal{X}, z = \mathcal{X} \cap An(W, G_{\bar{x}})\}$$

$$An(Y_{do(X_1, X_2)})? \\ \cdot G_{\overline{X_1, X_2}} \cdot An(Y, G_{\overline{X_1, X_2}}) \setminus \{X_1, X_2\}$$

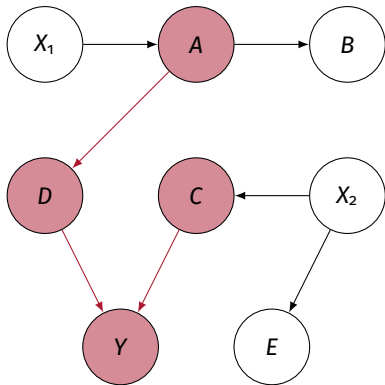
ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(x)}) = \{W_{do(z)} \mid W \in An(Y, G_{\bar{X}}) \setminus X, z = X \cap An(W, G_{\bar{X}})\}$$

$$An(Y_{do(X_1, X_2)})? \\ \cdot G_{\overline{X_1, X_2}} \cdot An(Y, G_{\overline{X_1, X_2}}) \setminus \{X_1, X_2\}$$

ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(x)}) = \{W_{do(z)} \mid$$

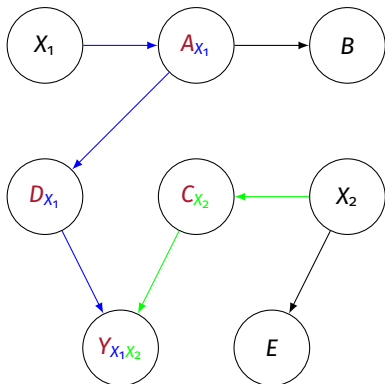
$$W \in An(Y, G_{\overline{X}}) \setminus X,$$

$$z = x \cap An(W, G_{\overline{X}})\}$$

$$An(Y_{do(X_1, X_2)})?$$

$$\cdot G_{\overline{X_1, X_2}} \cdot An(Y, G_{\overline{X_1, X_2}}) \setminus \{X_1, X_2\} \cdot \{x_1, x_2\} \cap An(W, G_{\overline{X_1, X_2}})$$

ANCESTOR OF A COUNTERFACTUAL VARIABLE



$$An(Y_{do(x)}) = \{W_{do(z)} \mid W \in An(Y, G_{\bar{X}}) \setminus X, z = x \cap An(W, G_{\bar{X}})\}$$

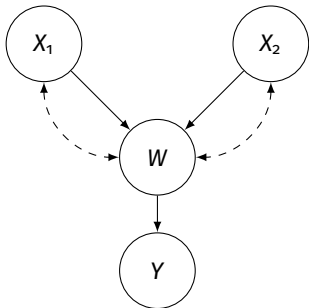
$$An(Y_{do(X_1, X_2)})? \cdot G_{\overline{X_1, X_2}} \cdot An(Y, G_{\overline{X_1, X_2}}) \setminus \{X_1, X_2\} \cdot \{x_1, x_2\} \cap An(W, G_{\overline{X_1, X_2}})$$

For a query on counterfactual variables:

- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables

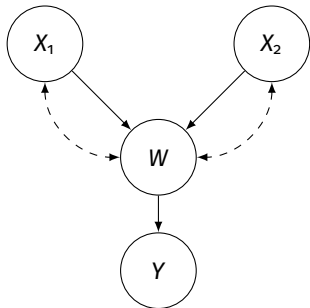
For a query on counterfactual variables:

- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



For a query on counterfactual variables:

- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables

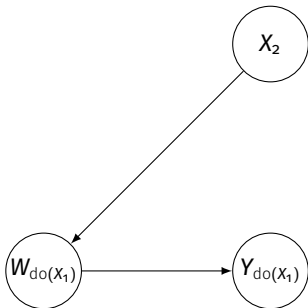
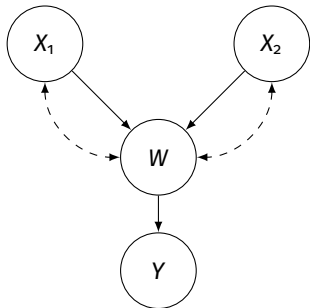


$$Y_{\text{do}(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{\text{do}(X_1)}, W_{\text{do}(X_1, X_2)}$$

Correa and Bareinboim 2025

For a query on counterfactual variables:

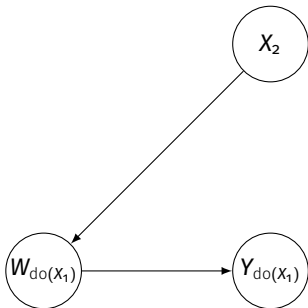
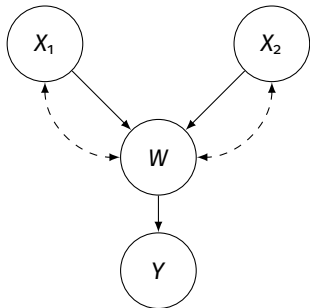
- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

For a query on counterfactual variables:

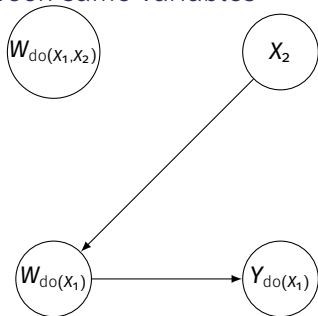
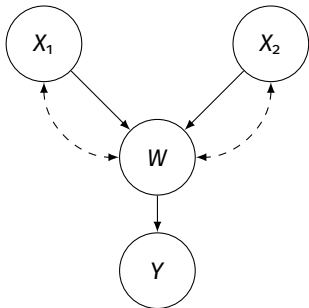
- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

For a query on counterfactual variables:

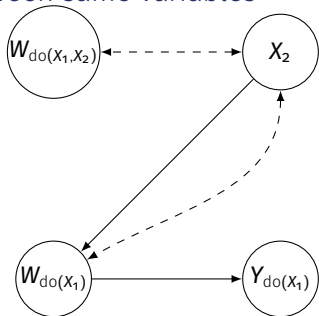
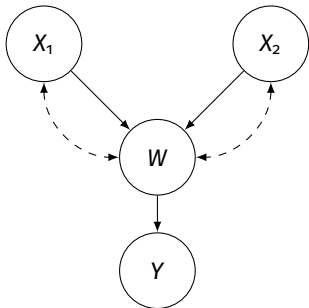
- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

For a query on counterfactual variables:

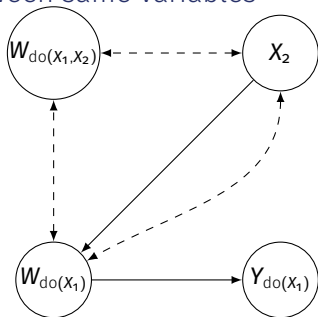
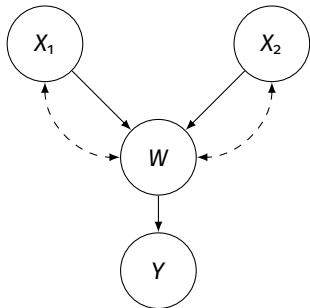
- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestry
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

For a query on counterfactual variables:

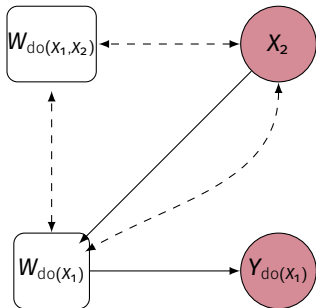
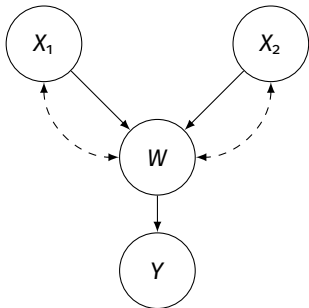
- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

For a query on counterfactual variables:

- Add all ancestors of the counterfactual variables
- Add edges witnessing the ancestrality
- Add the hidden variables
- Add hidden confounders between same variables



$$Y_{do(X_1)} \perp\!\!\!\perp_d X_2 \mid W_{do(X_1)}, W_{do(X_1, X_2)}$$

The ctf-calculus is sound and complete for counterfactual queries.

Rule 1 $\Pr(y_{T^*X}, x_{T^*}, w_*) = \Pr(Y_{T^*}, x_{T^*}, w_*)$

Rule 2 $\Pr(y_r | x_t, w_*) = \Pr(y_r | w_*)$

if $(Y_r \perp\!\!\!\perp_d X_t | W_*)_{\mathcal{G}_A}$

Rule 3 $\Pr(Y_{xz}, w_*) = \Pr(y_z, w_*)$

if $X \cap \text{An}(Y, \mathcal{G}_{\bar{Z}}) = \emptyset$

The ctf-calculus is sound and complete for counterfactual queries.

$$\text{Rule 1 } \Pr(y_{T_*X}, x_{T_*}, w_*) = \Pr(Y_{T_*}, x_{T_*}, w_*)$$

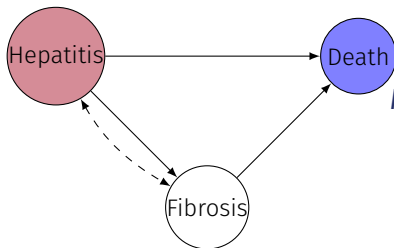
$$\text{Rule 2 } \Pr(y_r | x_t, w_*) = \Pr(y_r | w_*)$$

$$\text{if } (Y_r \perp\!\!\!\perp_d X_t | W_*)_{\mathcal{G}_A}$$

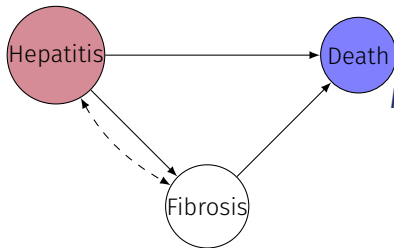
$$\text{Rule 3 } \Pr(Y_{xz}, w_*) = \Pr(y_z, w_*)$$

$$\text{if } X \cap \text{An}(Y, \mathcal{G}_{\bar{Z}}) = \emptyset$$

It transforms counterfactual queries into interventional queries. You might still need to use the do-calculus afterwards to transform the interventional queries into observational queries.



$$\begin{aligned}
 &NDE(D, H_{0 \rightarrow 1}, F) \\
 &= \mathbb{E} (D \mid \text{do} (H = 1, F = f_{\text{do}(H=0)})) \\
 &- \mathbb{E} (D \mid \text{do} (H = 0, F = f_{\text{do}(H=0)}))
 \end{aligned}$$



$$\begin{aligned}
 &NDE(D, H_{0 \rightarrow 1}, F) \\
 &= \mathbb{E}(D \mid \text{do}(H = 1, F = f_{\text{do}(H=0)})) \\
 &- \mathbb{E}(D \mid \text{do}(H = 0, F = f_{\text{do}(H=0)}))
 \end{aligned}$$

How do you identify

$$\Pr(D \mid \text{do}(H = 1, F = f_{\text{do}(H=0)})) = \Pr(d_{H/F_h})?$$

$$\Pr(d_{h'} F_h)$$

$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$(R1) = \sum_f \Pr(d_{h'f}, f_h)$$

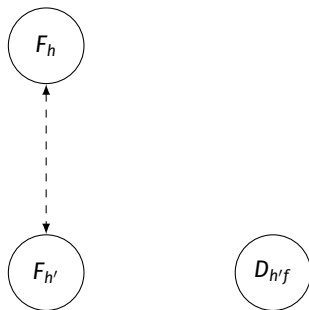
$$\begin{aligned}\Pr(d_{h'F_h}) &= \sum_f \Pr(d_{h'F_h}, f) \\ (R1) &= \sum_f \Pr(d_{h'f}, f_h) \\ &= \sum_f \Pr(d_{h'f} | f_h) \Pr(f_h)\end{aligned}$$

$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$(R1) = \sum_f \Pr(d_{h'f}, f_h)$$

$$= \sum_f \Pr(d_{h'f} | f_h) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'}) \Pr(f_h)$$



$$(D_{h'f} \perp\!\!\!\perp_d F_h, F_{h'})_{\mathcal{G}_A}$$

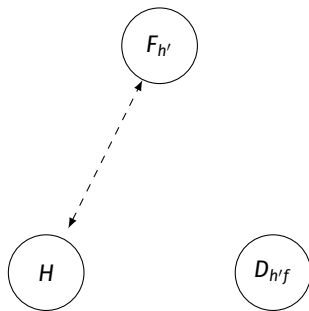
$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$(R1) = \sum_f \Pr(d_{h'f}, f_h)$$

$$= \sum_f \Pr(d_{h'f} | f_h) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'}) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'}h') \Pr(f_h)$$



$$(D_{h'f} \perp\!\!\!\perp_d H | F_{h'})_{G_A}$$

$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$(R1) = \sum_f \Pr(d_{h'f}, f_h)$$

$$= \sum_f \Pr(d_{h'f} | f_h) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'}) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'} h') \Pr(f_h)$$

$$(R1) = \sum_f \Pr(d_{h'f} | f, h') \Pr(f_h)$$

$$\Pr(d_{h'F_h}) = \sum_f \Pr(d_{h'F_h}, f)$$

$$(R1) = \sum_f \Pr(d_{h'f}, f_h)$$

$$= \sum_f \Pr(d_{h'f} | f_h) \Pr(f_h)$$

$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'}) \Pr(f_h)$$

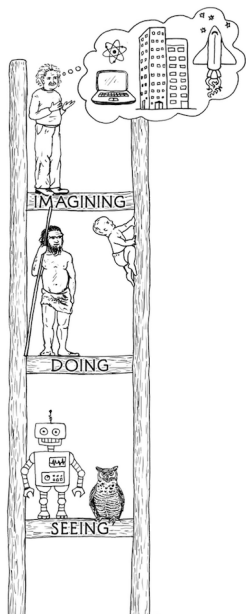
$$(R2) = \sum_f \Pr(d_{h'f} | f_{h'} h') \Pr(f_h)$$

$$(R1) = \sum_f \Pr(d_{h'f} | f, h') \Pr(f_h)$$

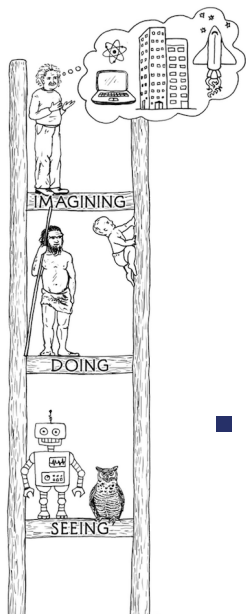
$$(R1) = \sum_f \Pr(d | f, h') \Pr(f_h)$$

6

CONCLUSION

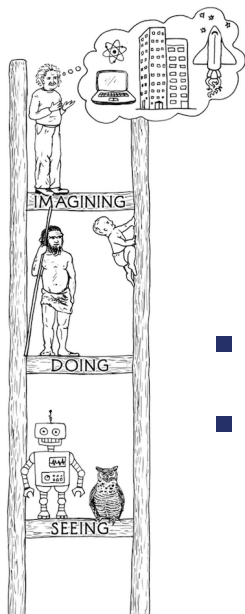


No assumptions



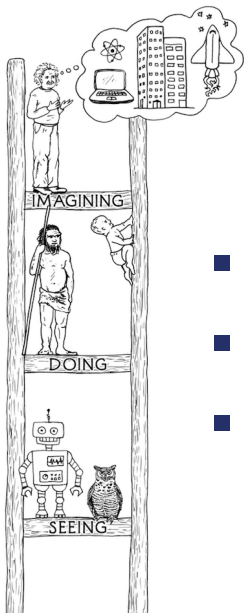
No assumptions

- Associations
 - ▶ observational data



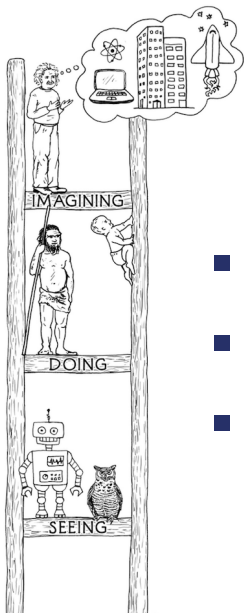
No assumptions

- Interventions
 - ▶ interventional data
- Associations
 - ▶ observational data



No assumptions

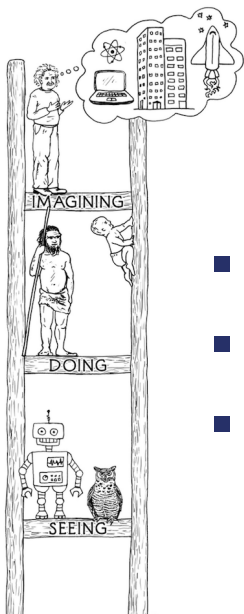
- **Counterfactuals**
 - ▶ Generating process
- **Interventions**
 - ▶ interventional data
- **Associations**
 - ▶ observational data



No assumptions

Assumptions: DAG

- Counterfactuals
 - ▶ Generating process
- Interventions
 - ▶ interventional data
- Associations
 - ▶ observational data



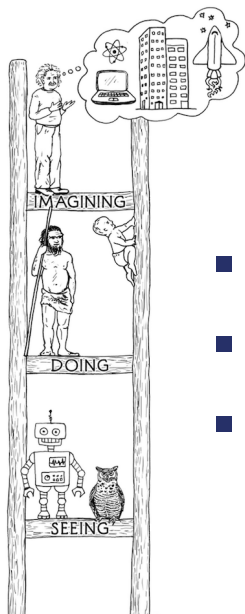
No assumptions

Assumptions: DAG

- Counterfactuals
 - ▶ Generating process
- Interventions
 - ▶ interventional data
- Associations
 - ▶ observational data



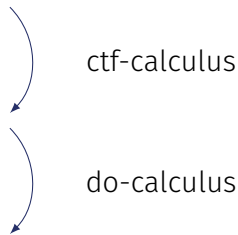
do-calculus



No assumptions

- Counterfactuals
 - ▶ Generating process
- Interventions
 - ▶ interventional data
- Associations
 - ▶ observational data

Assumptions: DAG



THANK YOU FOR YOUR ATTENTION!

ANY QUESTIONS?





Balke, Alexander and Judea Pearl (1994). “Probabilistic evaluation of counterfactual queries”. In: Proceedings of the Twelfth AAAI National Conference on Artificial Intelligence AAAI’94. Seattle, Washington: AAAI Press, 230–237.



Correa, Juan D and Elias Bareinboim (2025). “Counterfactual graphical models: Constraints and inference”. In: Forty-second International Conference on Machine Learning.



Correa, Juan D., Sanghack Lee, and Elias Bareinboim (2021). “Nested counterfactual identification from arbitrary surrogate experiments”. In: Proceedings of the 35th International Conference on Neural Information Processing Systems NIPS ’21. Red Hook, NY, USA: Curran Associates Inc. ISBN: 9781713845393.



Malinsky, Daniel, Ilya Shpitser, and Thomas Richardson (2019). “A potential outcomes calculus for identifying conditional path-specific effects”. In: The 22nd International Conference on Artificial Intelligence and Statistics. PMLR, pp. 3080–3088.



Pearl, Judea (Dec. 1995). “Causal diagrams for empirical research”. In: Biometrika 82.4, pp. 669–688. ISSN: 0006-3444. DOI: [10.1093/biomet/82.4.669](https://doi.org/10.1093/biomet/82.4.669).



— (2001). “Direct and Indirect Effects”. In: Probabilistic and Causal Inference.



— (2009). Causality: Models, Reasoning, and Inference. Cambridge University Press.



Richardson, Thomas S and James M Robins (2013). “Single world intervention graphs: a primer”. In: Second UAI workshop on causal structure learning, Bellevue, Wash



Robins, James M and Sander Greenland (1992).
“Identifiability and exchangeability for direct and indirect effects”. In: Epidemiology 3.2, pp. 143–155.