

# Causal Inference for Health Data

(STATS C160/C260 – Winter 2026)

## Lecture 9: Counterfactuals I

Drago Plečko

# Rules of Do-Calculus

**Theorem.** The following transformations are valid for any do-distribution induced by a SCM  $M$ :

## Rule 1: Adding/removing Observations

$$P(y|do(x), z, w) = P(y|do(x), w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\bar{X}}}$$

## Rule 2: Action/observation exchange

$$P(y|do(x), do(z), w) = P(y|do(x), z, w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\bar{XZ}}}$$

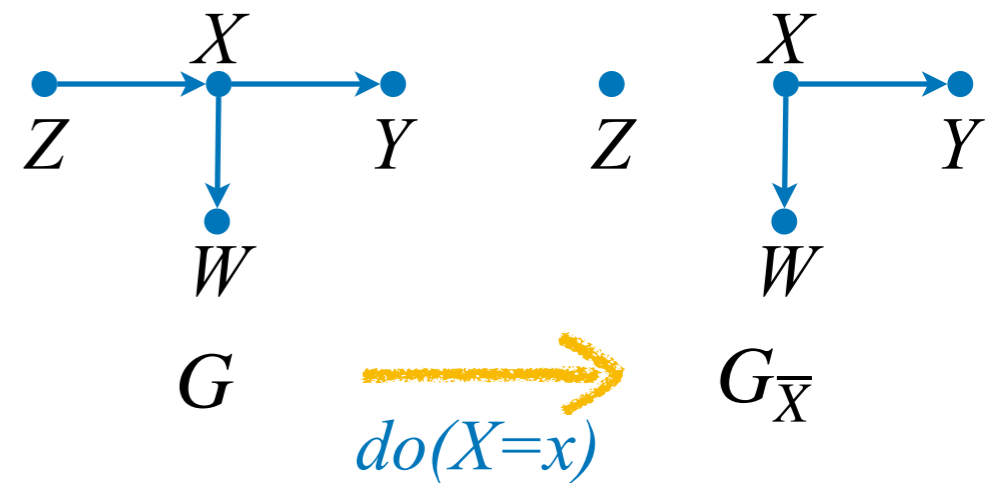
## Rule 3: Adding/removing Actions

$$P(y|do(x), do(z), w) = P(y|do(x), w) \quad \text{if } (Y \perp\!\!\!\perp Z \mid X, W)_{G_{\bar{XZ}(W)}}$$

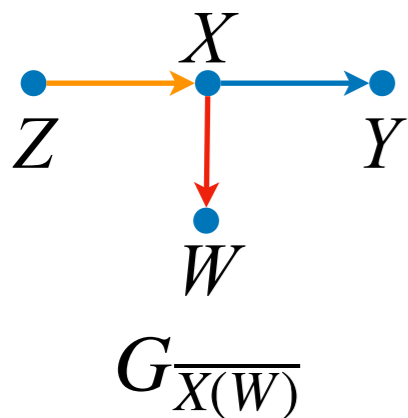
where  $Z(W) = Z \setminus \text{an}(W)_{G_{\bar{X}}}$  is the subset of  $Z$  that are not ancestors of  $W$  in  $G_{\bar{X}}$ .

# A curiosity about R3

- Regarding R3's contingency, note that the probability of  $Z$  given  $W$  may not be the same with or without intervening on  $X$ .



To witness, compare  $G_{\bar{X}}$  and  $G_{\overline{X(W)}}$ :

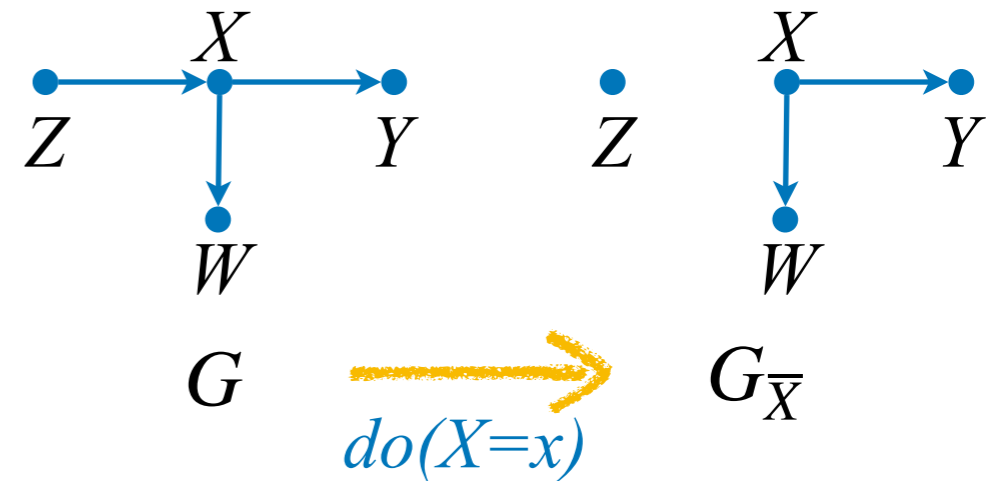


$$(Z \not\perp\!\!\!\perp X \mid W)_{G_{\overline{X(W)}}} \implies \exists_M P(z|do(x), w) \neq P(z|w)$$

What? 🙋

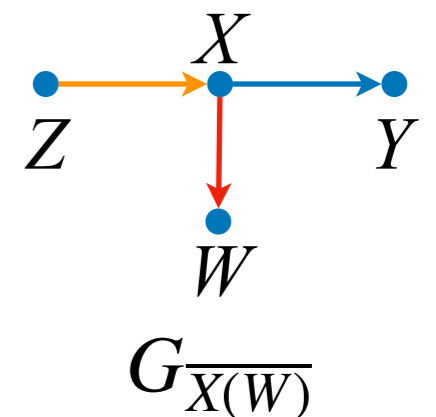
# A curiosity about R3

- Verifying whether  $P(z|do(x), w) \neq P(z|w)$  holds:



$$\begin{aligned}
 P(z|do(x), w) &= \frac{P(z, w | do(x))}{P(w | do(x))} \\
 &= \frac{\sum_y P(z)P(w | x)P(y | x)}{\sum_{y,z} P(z)P(w | x)P(y | x)} \\
 &= \frac{P(z)P(w | x)}{P(w | x)}
 \end{aligned}$$

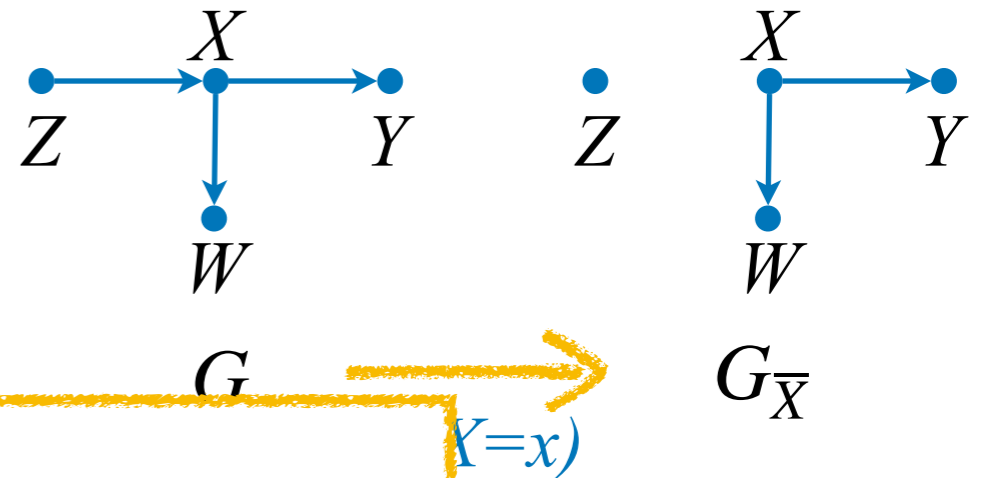
$$= P(z) \neq P(z|w) \quad (Z \not\perp\!\!\!\perp X | W)_{G_{\overline{X(W)}}}$$



in almost any model compatible with G.

# A curiosity about R3

- Verifying whether  $P(z|do(x),w) \neq P(z|w)$  holds:

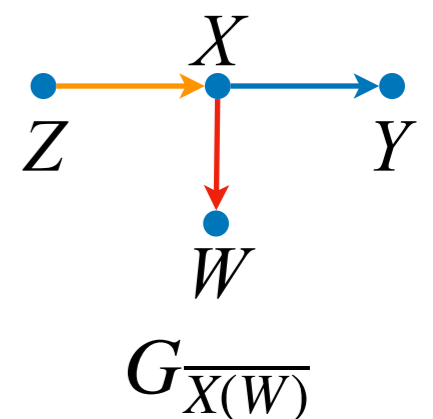


Intuitively:

(i)  $W$  is a descendant of  $X$

(ii) in a  $do(x)$ -world, conditioning on  $W$  does not propagate back to  $an(X)$ , since incoming arrows into  $X$  are removed in the  $do(x)$ -world

$\implies$  (iii) when conditioning on  $W$  in normal world, conditioning on  $W$  **does** propagate to  $an(X)$



**makes the 3rd rule more subtle!**

# Properties of Do-Calculus

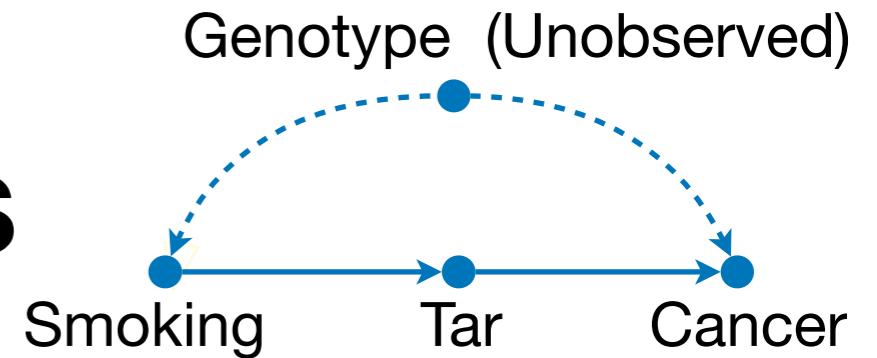
---

Theorem (soundness & completeness of do-calculus for interventional ID from obs. data).

The quantity  $Q = P(y|do(x))$  is identifiable from  $P(v)$  and  $G$  if and only if there exists a sequence of application of the rules of do-calculus and the probability axioms that reduces  $Q$  into a do-free expression.

Syntactic goal: Re-express original  $Q$  without  $do()$ !

# Derivation in Do-Calculus



$$\begin{aligned}
 P(c | do(s)) &= \sum_t P(c | do(s), t)P(t | do(s)) \\
 &= \sum_t P(c | do(s), do(t))P(t | do(s)) \\
 &= \sum_t P(c | do(s), do(t))P(t | s) \\
 &= \sum_t P(c | do(t))P(t | s) \\
 &= \sum_t \sum_{s'} P(c | do(t), s')P(s' | do(t))P(t | s) \\
 &= \sum_t \sum_{s'} P(c | t, s')P(s' | do(t))P(t | s) \\
 &= \sum_t \sum_{s'} P(c | t, s')P(s')P(t | s)
 \end{aligned}$$

Probability Axioms






Probability Axioms



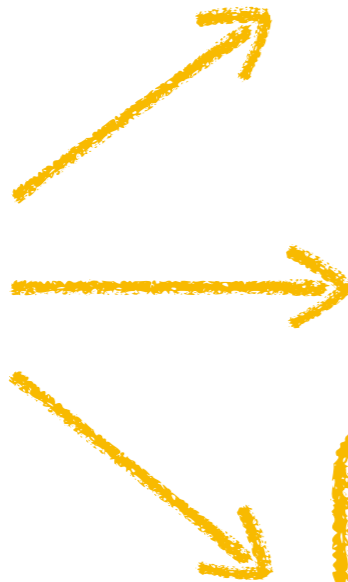
# Pearl's Causal Hierarchy— Layer 3: Counterfactuals

# SCM → Pearl's Causal Hierarchy

	Level (Symbol)	Typical Activity
	1  Association $P(y   x)$	Seeing ML - (Un)Supervised
	2  Intervention $P(y   do(x), c)$	Doing ML - Reinforcement
	3  Counterfactual $P(y_x   x', y')$	Imagining, Retrospection

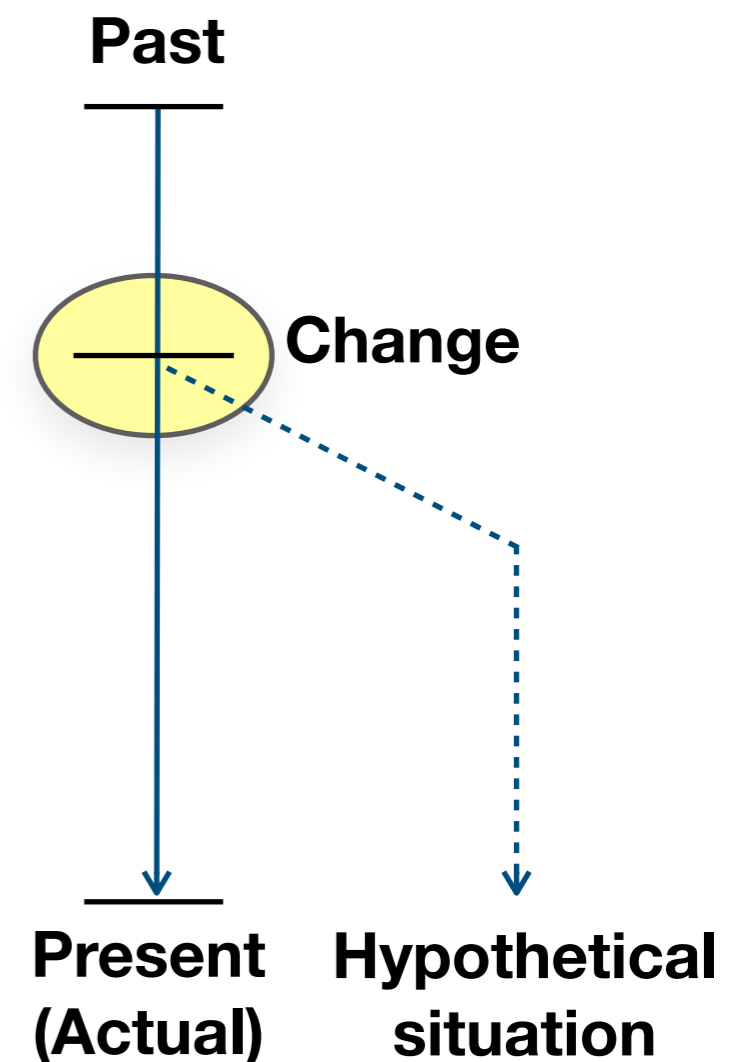


SCM  $M$

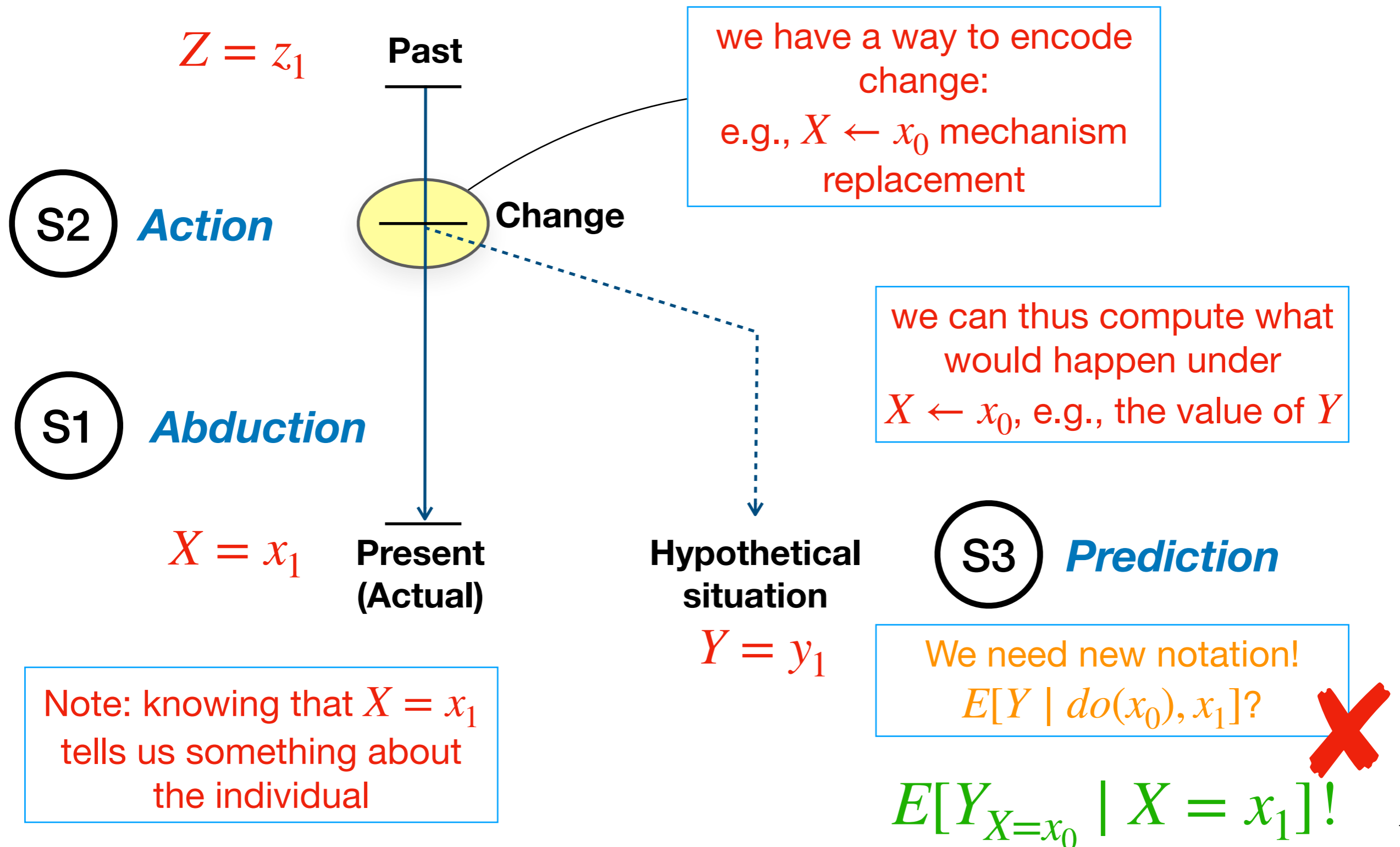


# Layer 3 – Semantics

- Layer 3 of the causal hierarchy allows operationalizing the notion of “imagination,” and the closely related activities of retrospection and introspection.
- This entails thinking about alternative ways the world could be, including ways that might conflict with the current state of the world.



# Layer 3 – Abduction, Action, Prediction (AAP)



# Layer 3 – AAP Example

- Query:  $E[Y_{X=x_0} \mid X = x_1]$

## S1 Abduction

Computing  $P(U \mid \text{evidence})$ :

$$P(U_z, U_x \mid X = 1) = 0.2 \times 0.7 + 0.8 \times 0.1$$

$$= \frac{P(U_z, U_x, X = 1)}{P(X = 1)}$$

$$= \frac{P(U_z)P(U_x)P(X = 1 \mid U_z, U_x)}{P(X = 1)}$$

**Distribution  $P(u \mid x_1)$ !**

$\mathcal{M}$

$$Z \leftarrow U_z$$

$$X \leftarrow (U_{x_1} \wedge \neg Z) \vee (U_{x_2} \wedge Z)$$

$$Y \leftarrow 1(0.8 - 0.3Z + 0.2XZ - 0.1X(1 - Z) > U_y)$$

$U_z \sim B(0.2), U_{x_1} \sim B(0.1),$   
 $U_{x_2} \sim B(0.7), U_Y \sim \mathbf{Unif}[0,1]$

$U_z$	$U_x$	$P(U_z, U_x, x_1)$	$U_z$	$U_x$	$P(U_z, U_x, x_1)$
0	(0, 0)	0	1	(0, 0)	0
0	(0, 1)	0	1	(0, 1)	$0.2 * 0.9 * 0.7$
0	(1, 0)	$0.8 * 0.1 * 0.3$	1	(1, 0)	0
0	(1, 1)	$0.8 * 0.1 * 0.7$	1	(1, 1)	$0.2 * 0.1 * 0.7$

# Layer 3 – AAP Example

Ⓢ2 **Action**

Ⓢ3 **Prediction (in  $M_{x_0}$ )**

$$E[Y_{x_0} | X = 1]$$

$$= \sum_u E[Y_{x_0} | u] P(u | x_1)$$

$$= \frac{32}{220} \times 0.8 + \frac{48}{220} \times 0.8$$

$$U_z = 0, U_x = (1,0) \quad U_z = 0, U_x = (1,1)$$

$$+ \frac{126}{220} \times 0.5 + \frac{14}{220} \times 0.5$$

$$U_z = 1, U_x = (0,1) \quad U_z = 0, U_x = (1,1)$$

$$\approx 0.61$$

$M_{x_0}$

$Z \leftarrow U_z$

$X \leftarrow x_0$

$Y \leftarrow 1(0.8 - 0.3Z + 0.2XZ - 0.1X(1 - Z)) > U_y$

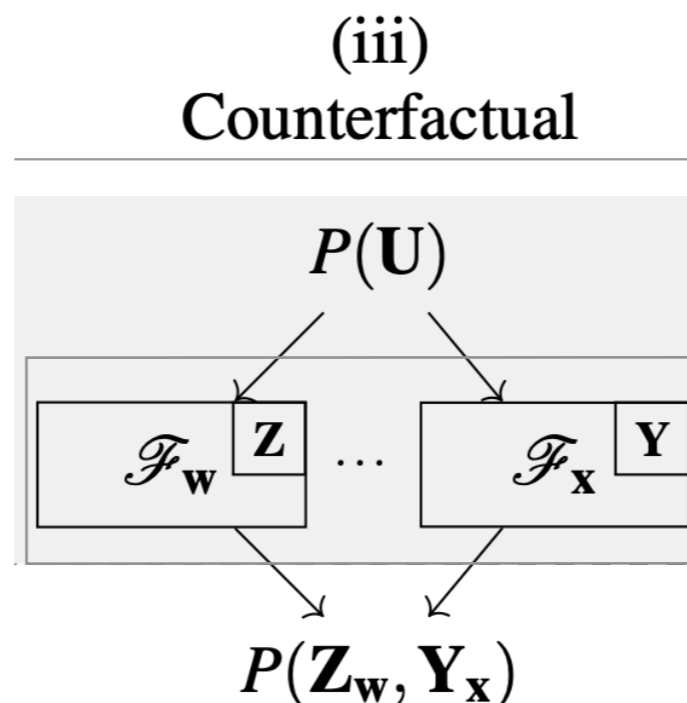
$U_z \sim B(0.2), U_{x_1} \sim B(0.1),$

$U_{x_2} \sim B(0.7), U_Y \sim \mathbf{Unif}[0,1]$

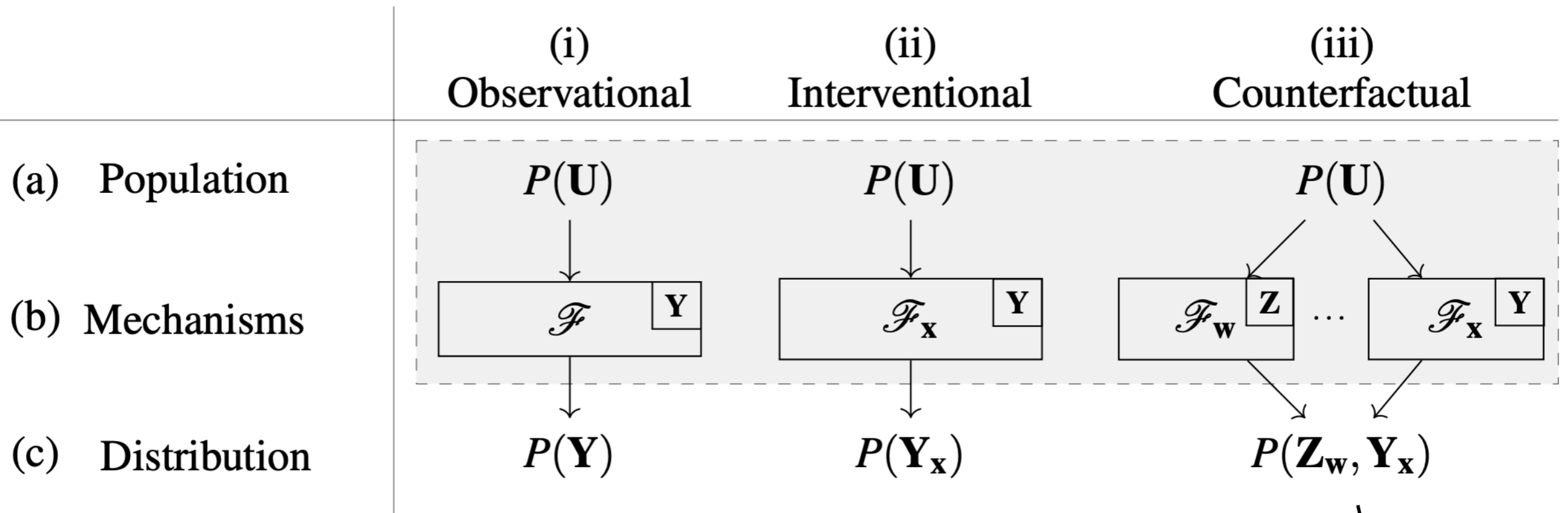
# Layer 3 – Valuation

**Definition 2.2.6 (Layer 3 Valuation).** An SCM  $\mathcal{M} = \langle \mathbf{U}, \mathbf{V}, \mathcal{F}, P(\mathbf{U}) \rangle$  induces a family of joint distributions over counterfactual events  $Y_x, \dots, Z_w$ , for any  $Y, Z, \dots, X, W \subseteq \mathbf{V}$ :

$$P^{\mathcal{M}}(\mathbf{y}_x, \dots, \mathbf{z}_w) = \sum_{\mathbf{u}} \mathbf{1} \left( \mathbf{Y}_x(\mathbf{u}) = \mathbf{y}, \dots, \mathbf{Z}_w(\mathbf{u}) = \mathbf{z} \right) P(\mathbf{u}).$$



# Layer 3 vs. the other layers



Only SCM  $M$   
appears!

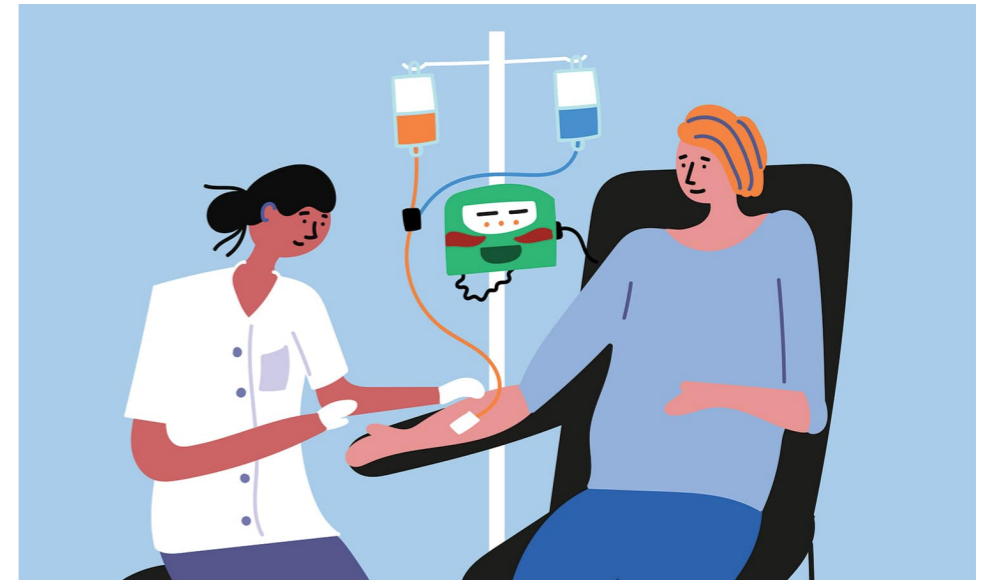
Only SCM  $M_x$   
appears!

Multiple SCMs  
appear ( $M, M_x$ )

# **II. Why Use Counterfactuals (Motivation)**

# Example 5.1. Oncologist and Chemotherapy

- An oncologist needs to decide whether to treat her patients, suffering from breast cancer, with chemotherapy,
- For this purpose, she wants to analyze the data collected in her hospital,
- Let  $X$  represent the treatment, where  $X = 0$  encodes no chemotherapy, and  $X = 1$  for chemotherapy.
- Let  $Y$  represent the patient's health, where  $Y = 0$  means that the patient has died, and  $Y = 1$  that the patient is still alive two years after taking the treatment;  $Z$  represents how advanced the cancer is ( $Z = 0$  for Stage I-II,  $Z = 1$  for Stage III-IV).
- Using the data, how should the physician decide to treat her patients?



# (1) Observational distribution

- She extracts the preliminary data from the hospital EHR, and gets the observational distribution  $P(v)$ ,
- She then decides to compute the metric  $E[Y | X = 1] - E[Y | X = 0]$ ,

- Based on the table she finds that

$$E[Y | X = 1] = 70\%$$

$$E[Y | X = 0] = 77.7\%,$$

$P(v)$	Z	X	Y
0.144	0	0	0
0.03	1	0	0
0.024	0	1	0
0.576	0	0	1
0.042	1	1	0
0.03	1	0	1
0.056	0	1	1
0.098	1	1	1

- The results imply that treated patients are less likely to survive. At first, they seem puzzling, but...
- Recalling the class in causal inference, she realizes that this quantity is not informative, since data it does not adjusted for confounding bias — more sick patients (Stage III-IV) are more likely to receive chemotherapy in the first place.

# (2) Causal Effect Computation

- She then decides to use the adjustment formula based on Z, based on the back-door criterion,  $E[y | do(x)] = \sum_z E[y | x, z]P(z)$

P(v)	Z	X	Y
0.144	0	0	0
0.03	1	0	0
0.024	0	1	0
0.576	0	0	1
0.042	1	1	0
0.03	1	0	1
0.056	0	1	1
0.098	1	1	1

- From BD-adjust, one can compute  $E[Y | do(X = 1)] - E[Y | do(X = 0)]$ ,
- Based on the table she finds that

$$E[Y | do(X = 1)] = 70 \%$$

$$E[Y | do(X = 0)] = 74 \%,$$

- She finds that even after adjustment, the causal effect of chemotherapy results in worse survival. Does this imply the treatment is futile, and should be abandoned entirely?
- How do make sense of this finding?

# (3) Oracle – Effect on Treated

- Let's look at the SCM generating the data, on the left side,
- The physician wants to understand whether the treatment has helped the group of patients who were actually given it,

Would a treated patient ( $X=1$ ) survive without treatment?

$$\begin{aligned} & \mathcal{M} \\ Z & \leftarrow \mathbf{Bern}(0.2) \\ X & \leftarrow \mathbf{Bern}(0.1 + 0.6Z) \\ Y & \leftarrow \mathbf{Bern}(0.8 - 0.3Z + \\ & \quad 0.2XZ - \\ & \quad 0.1X(1 - Z)) \end{aligned}$$

- This query has special meaning in the language of layer 3 ( $L_3$ ) of the PCH, and is written as the counterfactual statement:

$$P(Y_{X=0} = 1 \mid X = 1),$$

This expression is called the *effect of treatment on the treated (ETT)*.

# (3) Oracle – Effect on Treated

---

- What is the probability of survival under a certain treatment, given that the patient's natural treatment regime?
- From the SCM, we can compute:

$$P(Y_{X=1} = 1 \mid X = 1) = 70 \% \quad P(Y_{X=0} = 1 \mid X = 1) = 61 \%$$

- Therefore, among those patients who are treated ( $X = 1$  behind the conditioning bar), being treated ( $X = 1$ ) is actually be better since

$$P(Y_{X=1} = 1 \mid X = 1) > P(Y_{X=0} = 1 \mid X = 1),$$

- For those who are not treated, we can show the opposite:

$$P(Y_{X=1} = 1 \mid X = 0) < P(Y_{X=0} = 1 \mid X = 0),$$

- How is this possible?

# (3) Oracle – Effect on Treated

- We first look at the mechanism for  $Y$ :

$$Y \leftarrow \text{Bern}(0.8 - 0.3Z + 0.2XZ - 0.1X(1 - Z))$$

baseline survival  
of 80%

30% survival decrease  
for late stage

for late stage,  
chemo increases  
survival 20%

for early stage,  
chemo decreases  
survival 10%

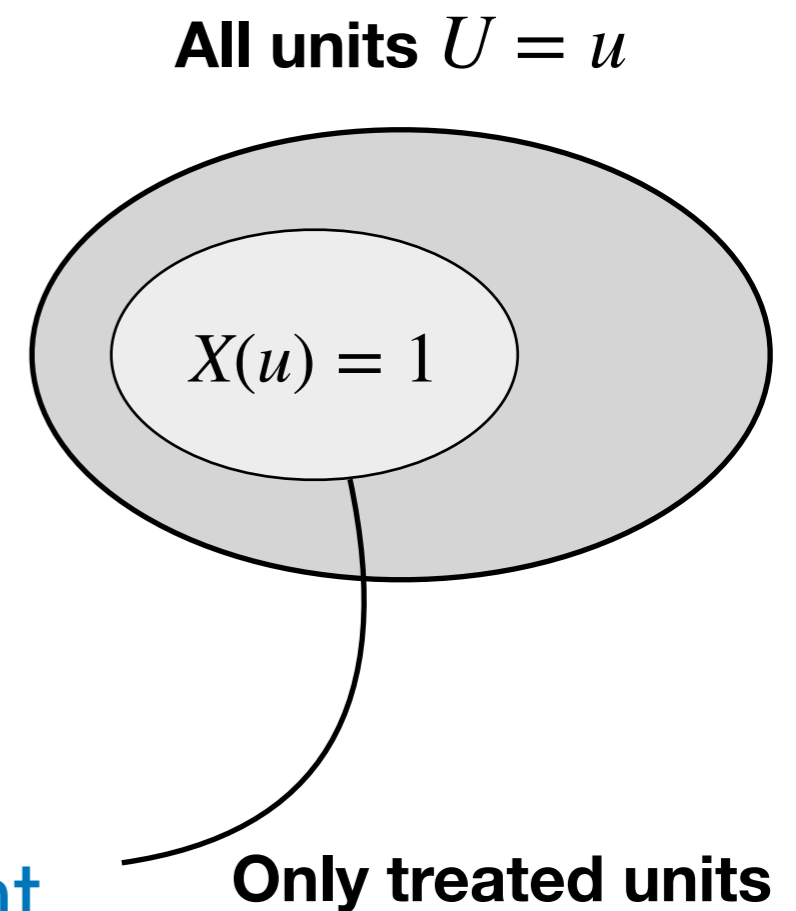
- Therefore, treatment is harmful for early stage, but helpful for late,
- Furthermore, the mechanism for  $X$  shows us that

$$X \leftarrow \text{Bern}(0.1 + 0.6Z)$$

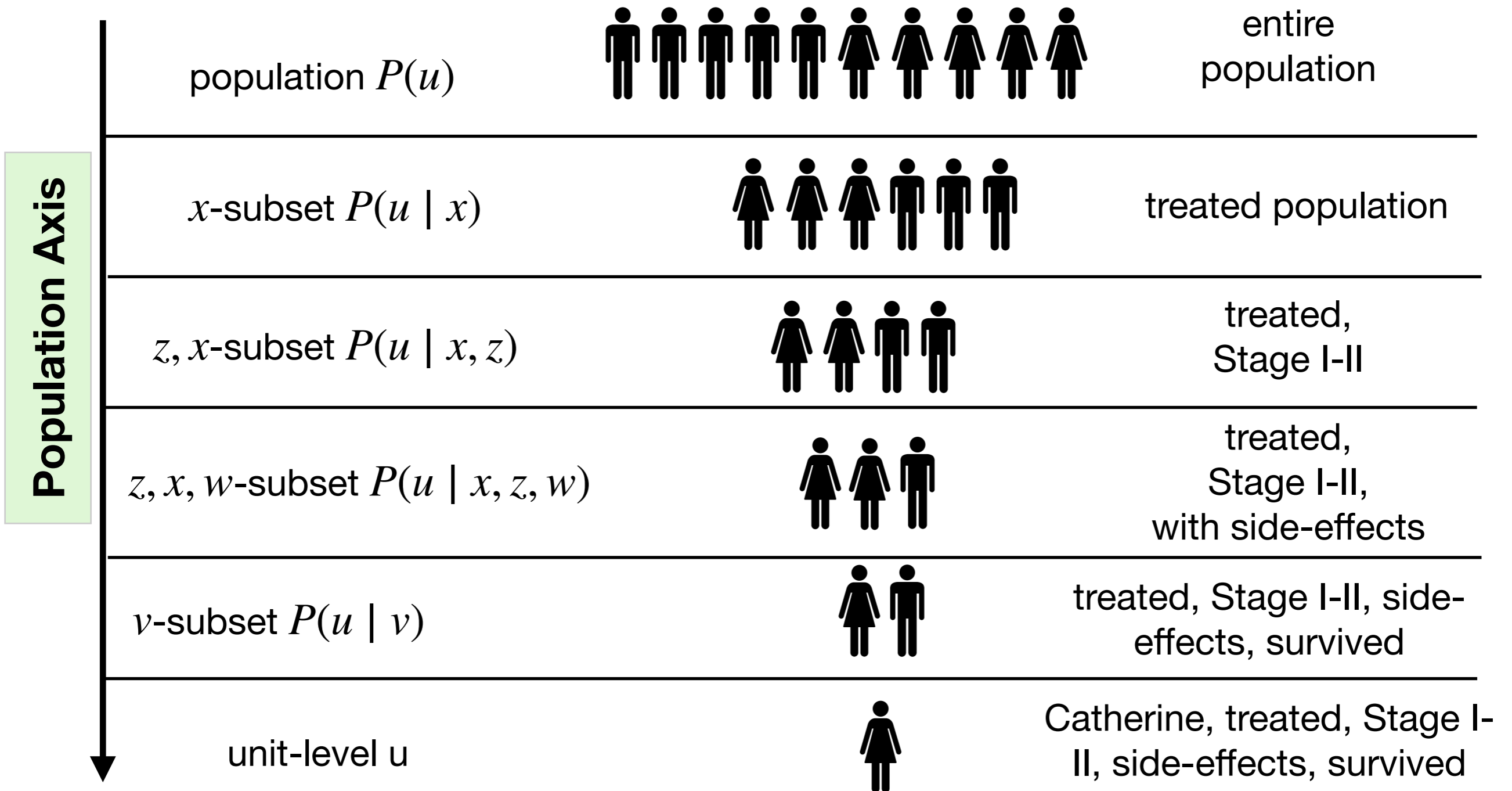
among treated ( $X=1$ ), there are more late-stage patients than usual!

# Example – Takeaways

- Causal effects may exhibit *heterogeneity*, meaning that different parts of the population are affected differently,
- The average causal effect  $E[Y \mid do(X = 1)] - E[Y \mid do(X = 0)]$  is an average over the entire, possibly heterogeneous, population of individuals,
- To better understand heterogeneity, we may use conditional effects, which sometimes require **Layer 3 of PCH**,
- This is the case for the **Effect of Treatment on the Treated (ETT)**.



# (1) Use of Counterfactuals: Granularity



# Example. Semaglutide Effects on Diabetes

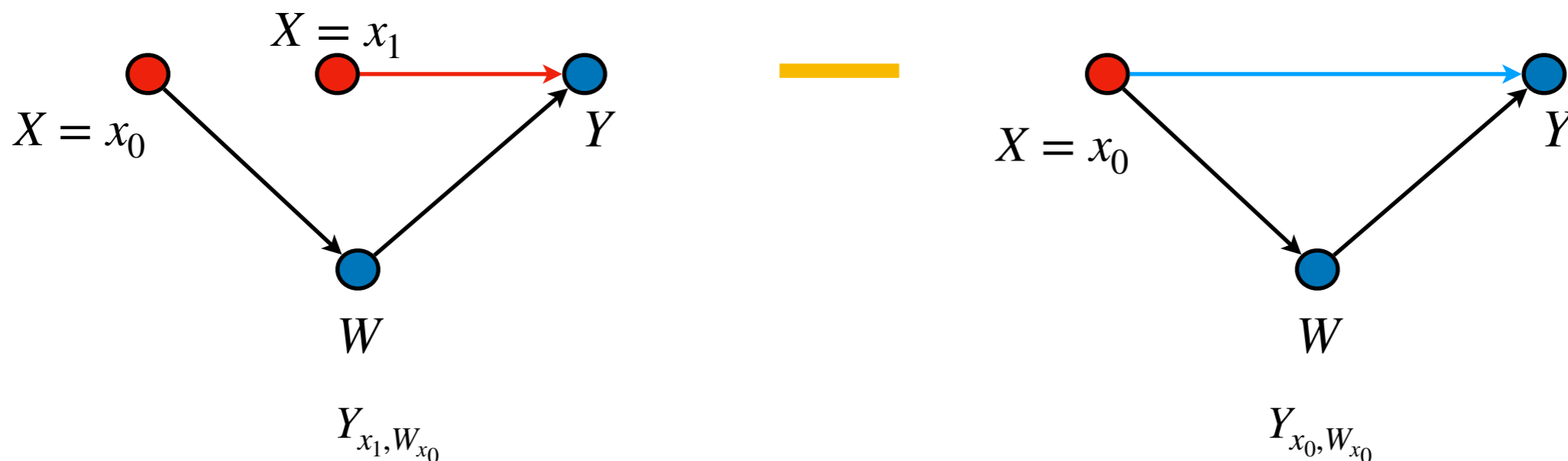
- A diabetologist is treating his patients with semaglutide, to stabilize their HbA1c values,
- However, he realized that there are different causal pathways through which semaglutide may affect HbA1c levels: (1) through weight-loss and associated insulin sensitivity, and (2) direct GLP-1 receptor effects,
- The doctor is interested in how strong each of these pathways is, and he wonders if he can obtain an answer from data:
- $X = 1$  denotes semaglutide treatment,  $X = 0$  standard treatment,  $Y = 1$  denote normal HbA1c levels ( $Y = 0$  otherwise);  $W = 1$  for weight loss ( $W = 0$  otherwise).



# Gedankenexperiment (NDE)

- For investigating this, the doctor performs following thought experiments,
- For an untreated individual ( $X = x_0$ ), how would his HbA1c level ( $Y$ ) change **had he been** treated ( $X = x_1$ ), while keeping the weight loss unchanged (at the natural level  $X = x_0$ )?

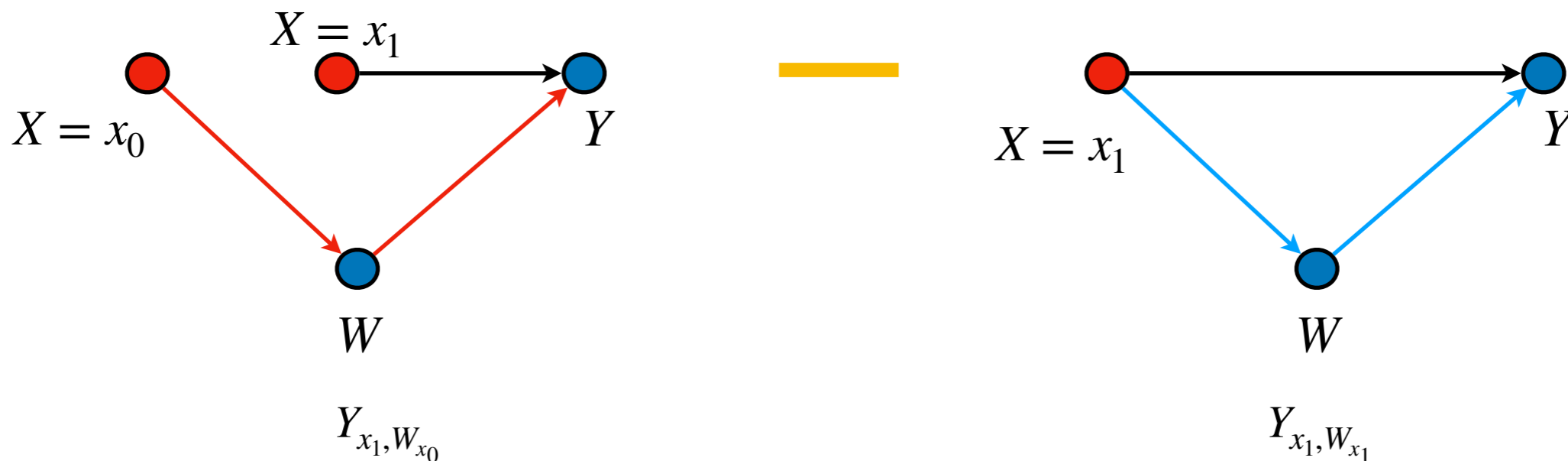
$$\mathbf{NDE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0, W_{x_0}})$$



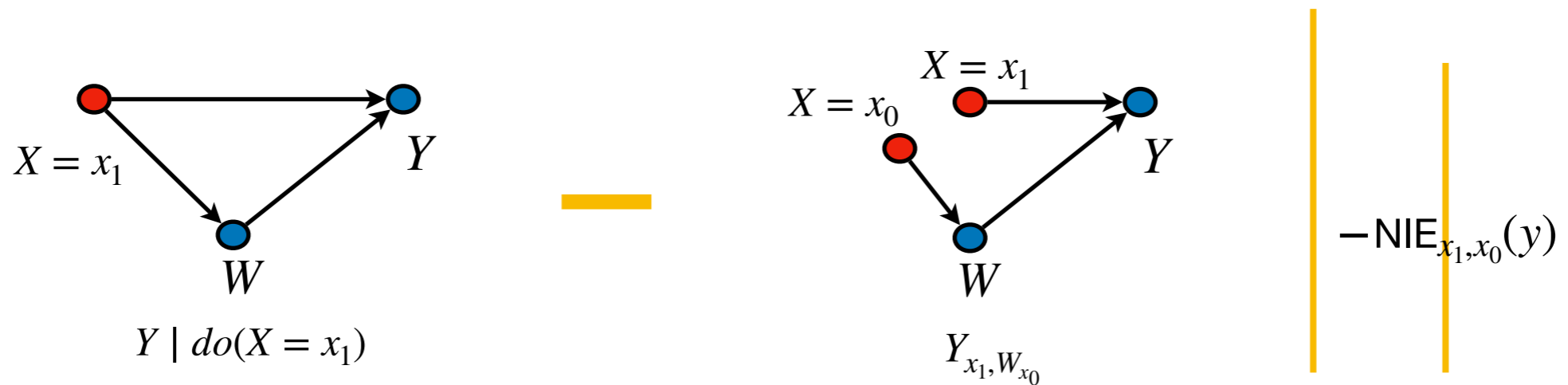
# Gedankenexperiment (NIE)

- For a treated individual ( $X = x_1$ ), how would their HbA1c level ( $Y$ ) change **had they not been** treated ( $X = x_0$ ), while keeping treatment unchanged along the direct causal pathway (at the natural level  $X = x_1$ )?

$$\mathbf{NIE}_{x_1, x_0}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_1, W_{x_1}})$$

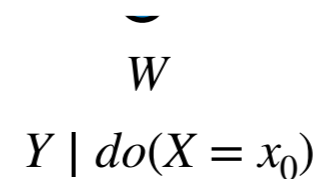
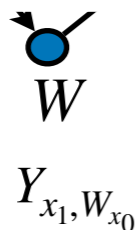


# ATE Decomposition



**Theorem.** The average total effect can be decomposed into its direct and indirect parts:

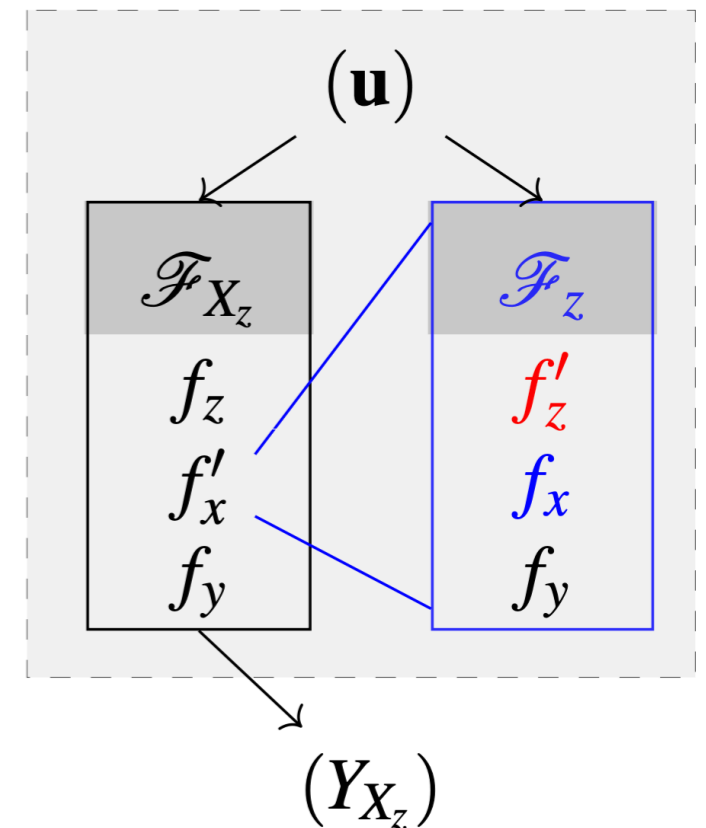
$$ATE_{x_0, x_1}(y) = NDE_{x_0, x_1}(y) - NIE_{x_1, x_0}(y).$$



$NDE_{x_0, x_1}(y)$

# Nested Counterfactuals

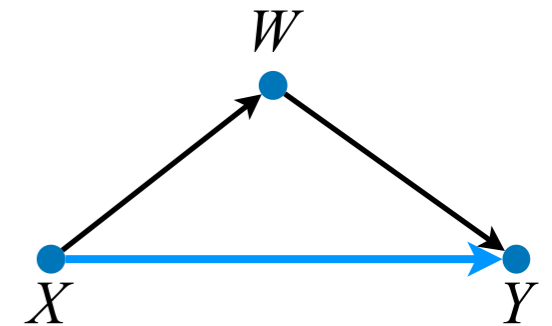
- We used the notion of *nested counterfactuals*.
- $Y_x$  refers to  $Y$  under intervention  $do(X = x)$ , that is, when  $X$  is fixed to a constant value  $x$ .
- Here, we consider an intervention  $do(X \leftarrow X_z)$ , where  $X$  is not fixed to a constant but is supposed to behave as another counterfactual variable  $X_z$ .
- For such intervention, we first need to evaluate  $X_z$  (from  $\mathcal{M}_z$ ), and then consider a model  $\mathcal{M}_{X_z}$  where  $X$  is given the value dictated by  $X_z$ .



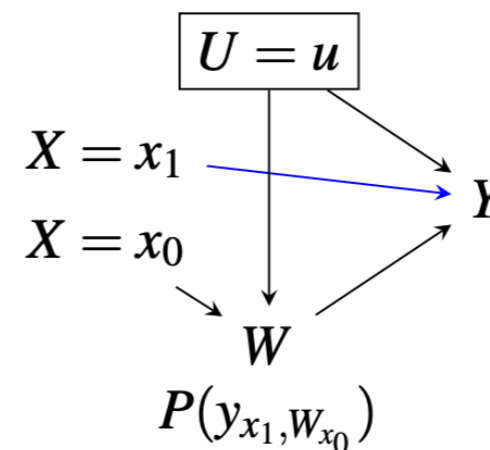
# Natural Direct Effects

- The **Natural Direct Effect (NDE)** is defined as

$$\text{NDE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0, W_{x_0}}) \quad (5.101)$$



- $Y_{x_1, W_{x_0}}$  refers to the outcome  $Y$  under  $X = x_1$  and  $W = W_{x_0}$ , the value that  $W$  would attain had  $X$  been  $x_0$ .
- $Y_{x_0}$  represents a baseline where  $Y$  perceives  $X = x_0$  in all causal paths. It is also equivalent to  $Y_{x_0, W_{x_0}}$ , that is, a situation where  $W$  also gets  $X = x_0$ .
- Taking the difference of those two quantities keeps the path  $X \rightarrow W \rightarrow Y$  constant while changing the level of  $X$  from  $x_0$  to  $x_1$  in the path  $X \rightarrow Y$ , effectively measuring the direct impact of  $X$  on  $Y$ .



# Summary: Using of Counterfactuals

## Total Effect – TE

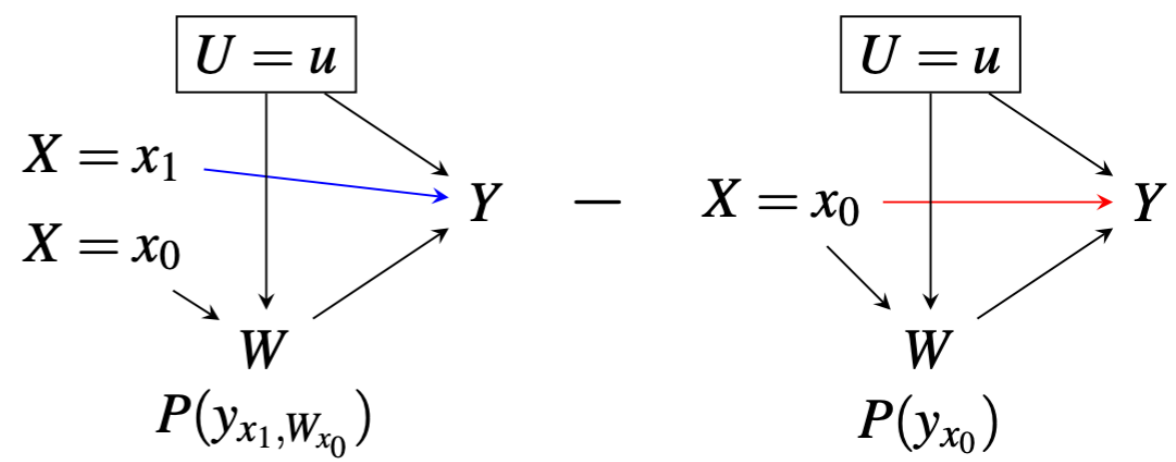
$$E[Y | do(x_1)] - E[Y | do(x_0)]$$

We spent the first 8 lectures building foundations to understand this quantity

Ext. 2: Mechanisms

DE

IE



Ext. 1: Granularity

