

# Causal Inference for Health Data

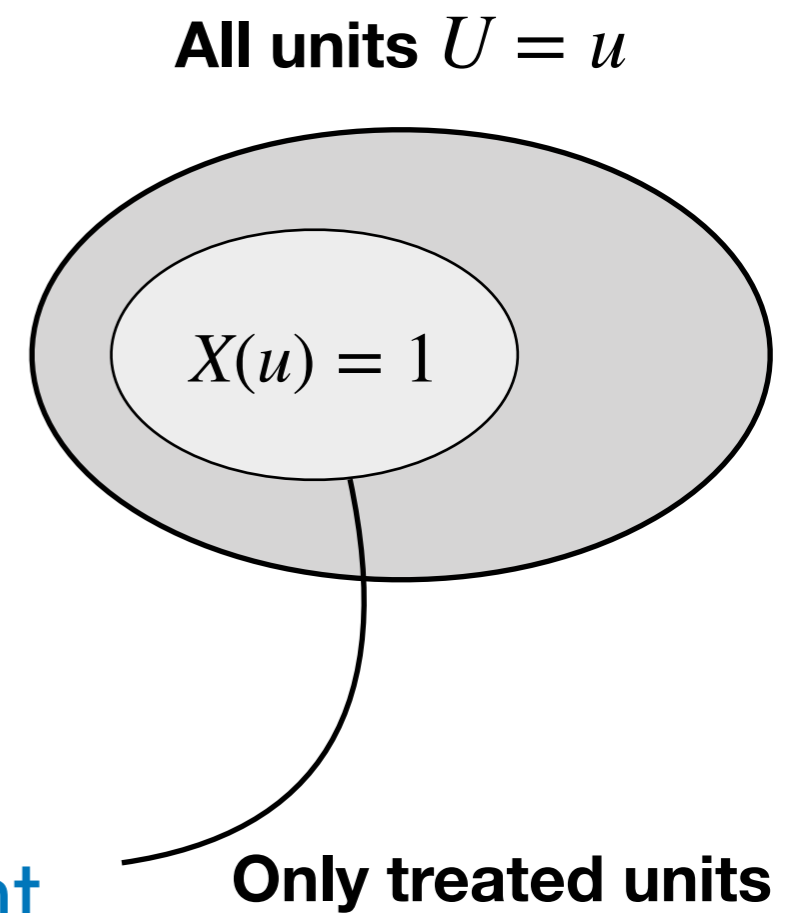
(STATS C160/C260 – Winter 2026)

## Lecture 10: Counterfactuals II

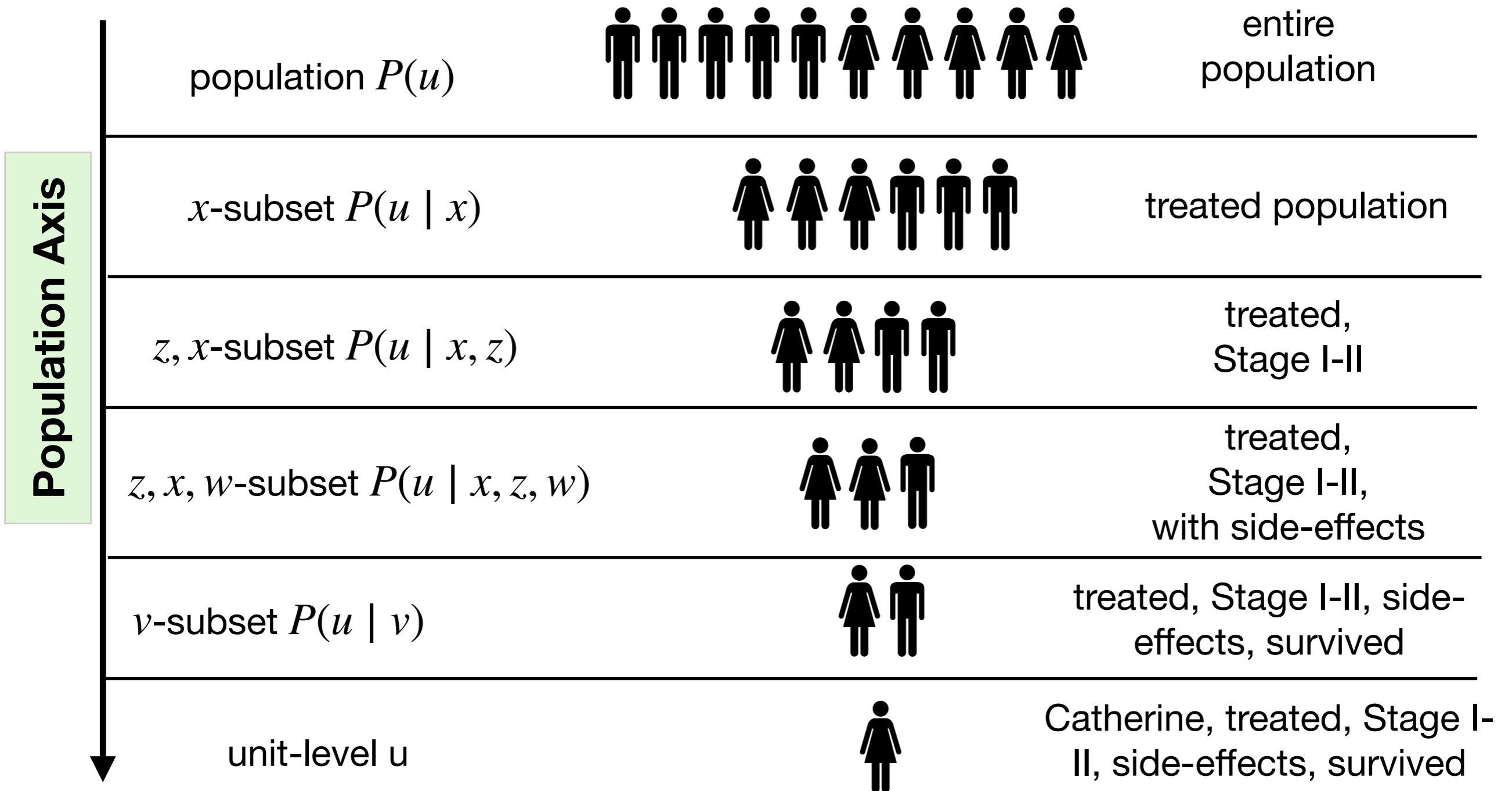
Drago Plečko

# 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

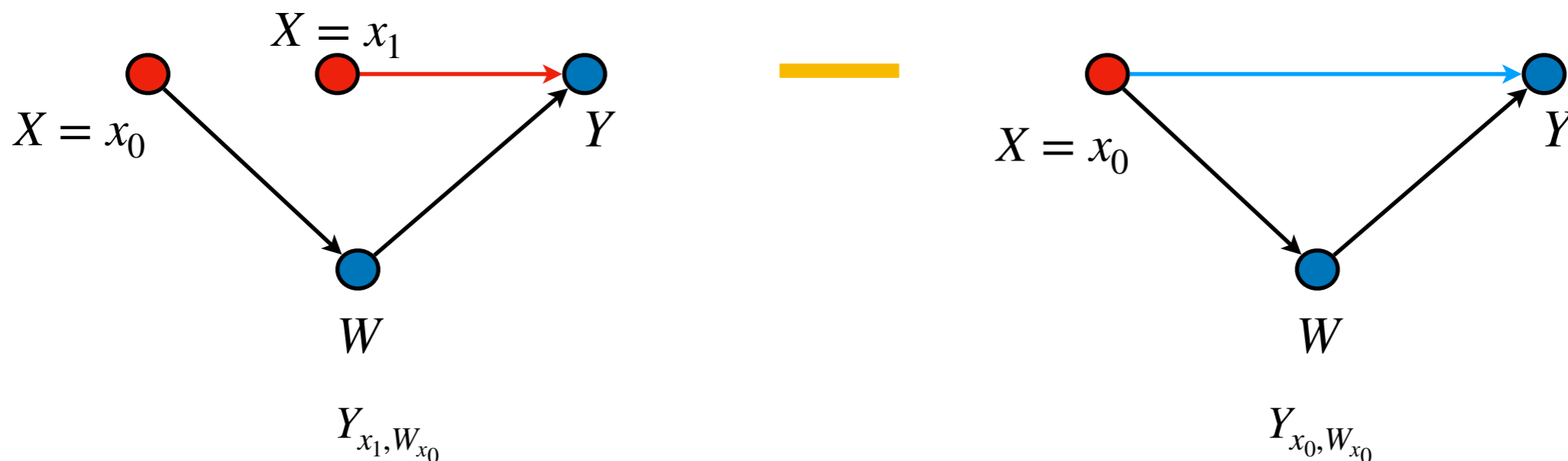
- A diabetologist is treating his patients with semaglutide, to stabilize their HbA1c values,
- However, he realized that there are different causal pathways through which semaglutide may affect HbA1c levels: (1) through weight-loss and associated insulin sensitivity, and (2) direct GLP-1 receptor effects,
- 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  $X = x_0$ )?

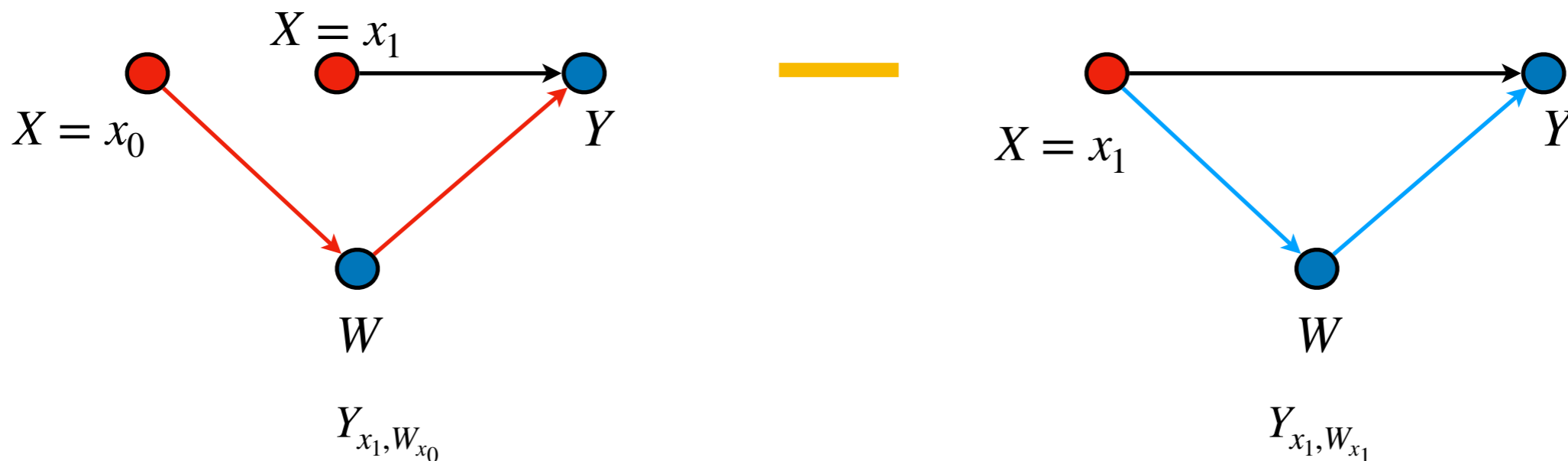
$$\mathbf{NDE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0, W_{x_0}})$$



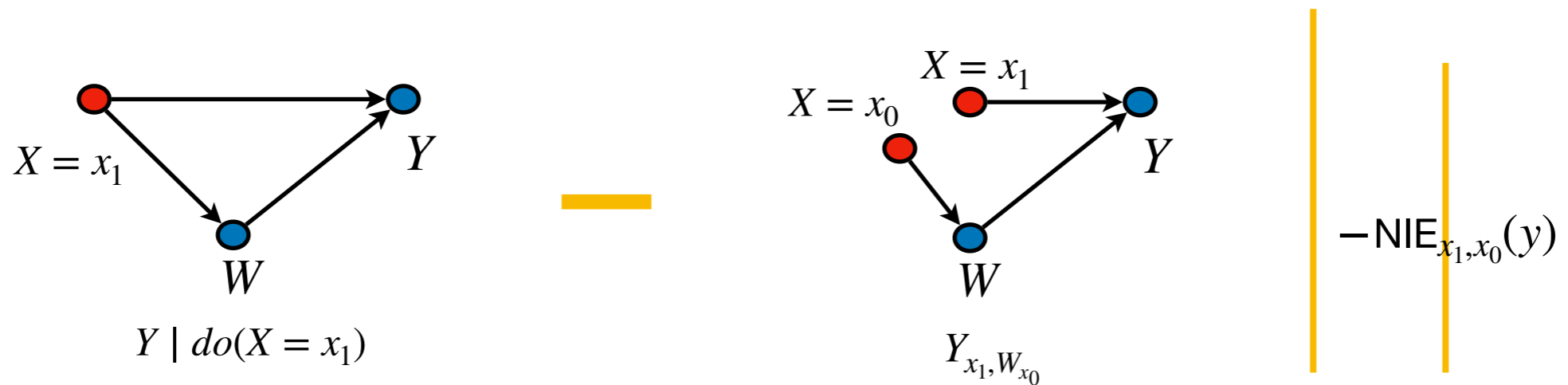
# Gedankenexperiment (NIE)

- 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  $X = x_1$ )?

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

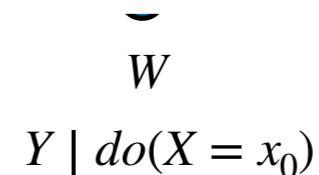
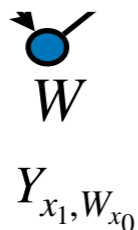


# ATE Decomposition



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

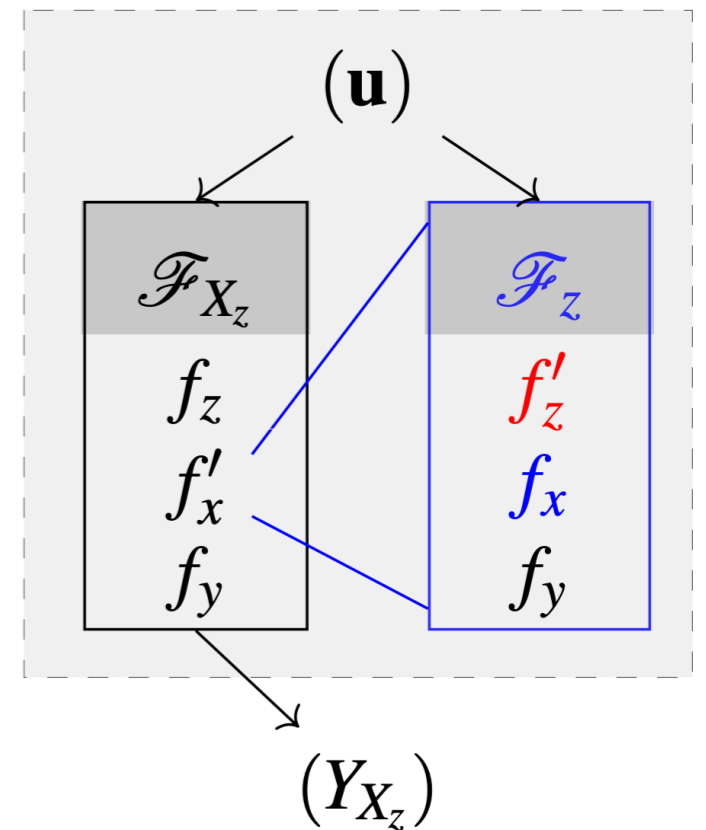
$$ATE_{x_0, x_1}(y) = NDE_{x_0, x_1}(y) - NIE_{x_1, x_0}(y).$$



$NDE_{x_0, x_1}(y)$

# 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  $X_z$ .

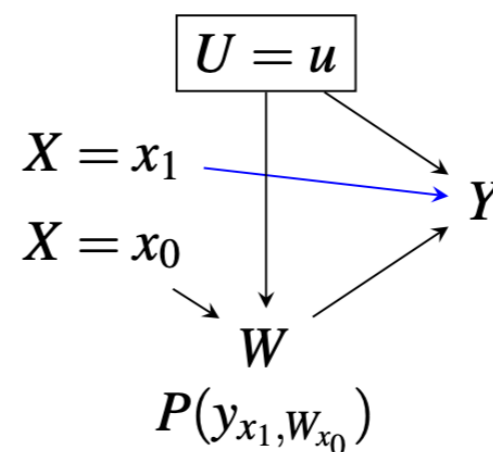
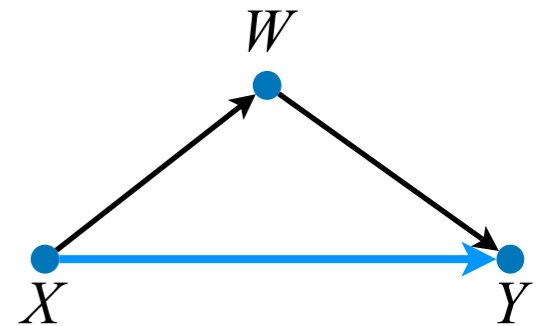


# Natural Direct Effects

- 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

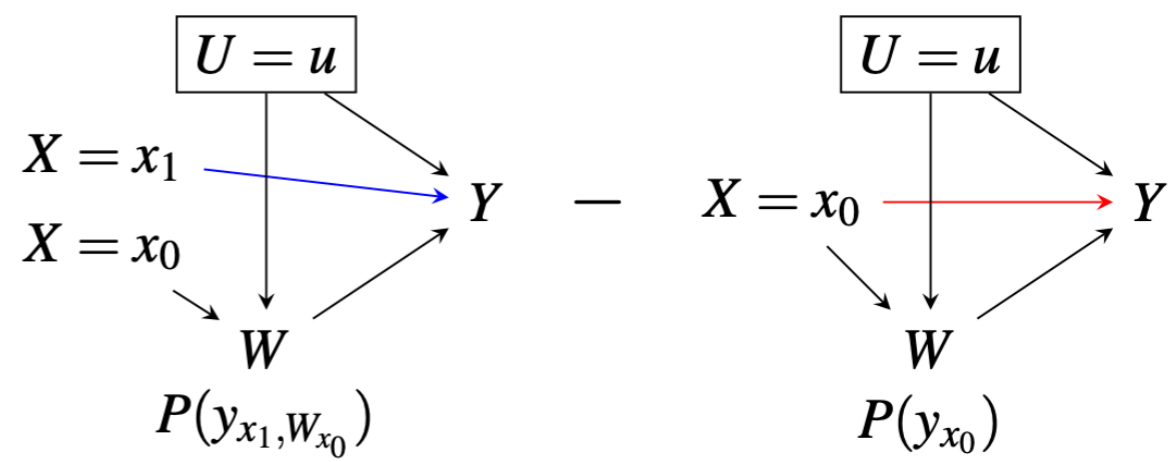
$$E[Y | do(x_1)] - E[Y | do(x_0)]$$

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

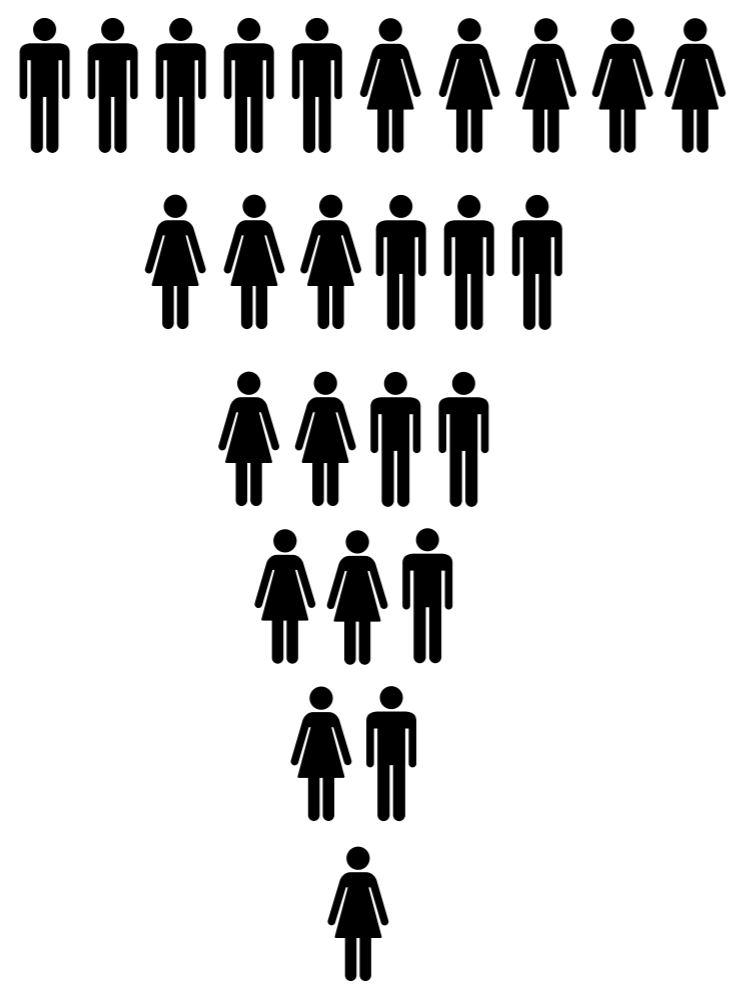
Ext. 2: Mechanisms

DE

IE



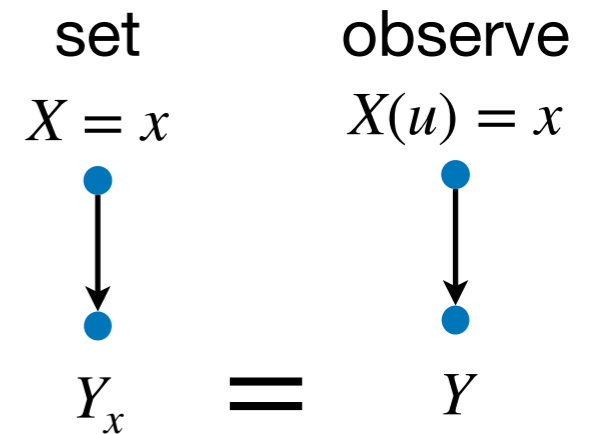
Ext. 1: Granularity



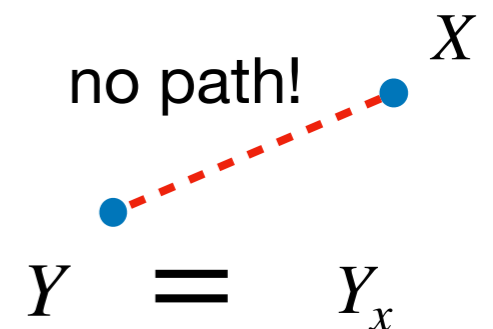
# Counterfactual Constraints

We will study three types of constraints, namely:

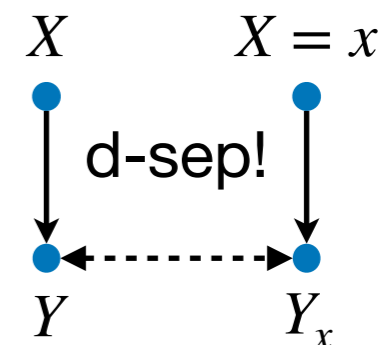
1. **Consistency**: variables behave the same with or without intervention under a certain observational context.



2. **Exclusion**: variables behave the same regardless of the exclusion of one or more interventions under a certain interventional context.



3. **Independence**: (counterfactual) variables are probabilistically independent.



# Consistency

SCM  $M$

unit  $u = (u_1, \dots, u_k) \sim P(u)$

after  $u$  is fixed, the evaluation  
is deterministic

$$V_1 \leftarrow f_1(u_1) \quad X(u) = x$$

$$X \leftarrow f_X(v_1, u_x) \quad \text{same}$$

$\vdots$

$$Y \leftarrow f_Y(v_1, \dots, v_k, u_y)$$

SCM  $M_x$

for the same fixed  $u$ , the  
evaluation is again  
deterministic

$$V_1 \leftarrow f_1(u_1)$$

$$X \leftarrow x$$

$\vdots$

$$Y_x \leftarrow f_Y(v_1, \dots, (v_k)_x, u_y)$$

same  
solution!

# Consistency

---

**Lemma - Consistency.** Given a structural causal model  $M$  and  $X, Y \in \mathbf{V}$ ,  $\mathbf{T} \subseteq \mathbf{V}$ , let  $x$  be a value in the domain of  $X$ . Then,

$$X_{\mathbf{T}} = x \Rightarrow Y_{\mathbf{T}_x} = Y_{\mathbf{T}}.$$

Whenever  $\mathbf{T}$  is empty, consistency can be written as:

$$X = x \Rightarrow Y_x = Y.$$

(In some settings, consistency is used when  $T$  is empty, and is called *composition* otherwise. )

# Consistency Example

- Recall the ETT quantity:

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

$\underbrace{\hspace{10em}}_{\text{consistency?}} \quad \checkmark$        $\underbrace{\hspace{10em}}_{\text{consistency?}} \quad \times$

$= E(Y \mid X=1)$       we compute it later!  
(need other constraints)

# Counterfactual Unnesting Theorem (CUT)

**Theorem – (CUT)** Let  $Y, X \in \mathbf{V}$ ,  $\mathbf{T}, \mathbf{Z} \subseteq \mathbf{V}$ , and let  $\mathbf{z}$  be a set of values for  $\mathbf{Z}$ . Then, the nested counterfactual  $P(Y_{\mathbf{T}X_{\mathbf{z}}} = y)$  can be written as a non-nested counterfactual, as follows:

$$P(Y_{\mathbf{T}X_{\mathbf{z}}} = y) = \sum_x P(Y_{\mathbf{T}x} = y, X_{\mathbf{Tz}} = x).$$

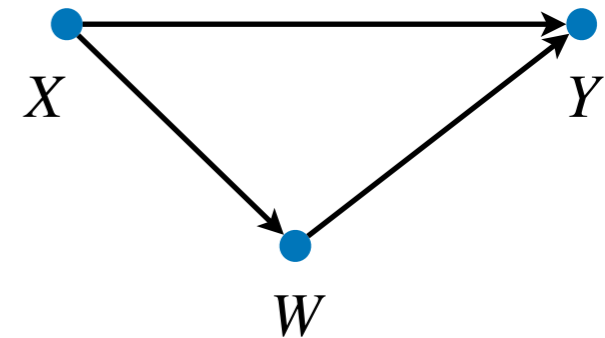
This statement's proof only requires two steps, the law of total probability and consistency itself, i.e.:

$$\begin{aligned} P(Y_{\mathbf{T}X_{\mathbf{z}}} = y) &= \sum_x P(Y_{\mathbf{T}X_{\mathbf{z}}} = y, X_{\mathbf{Tz}} = x) && \text{(sum over } X_{\mathbf{Tz}}) \\ &= \sum_x P(Y_{\mathbf{T}x} = y, X_{\mathbf{Tz}} = x) && \text{(consistency nested ctf.).} \end{aligned}$$

# Un-nesting Example: NDE

- Recall the notion of natural direct effect (NDE) that we defined

$$\text{NDE}_{x_0, x_1}(y) \stackrel{\delta}{=} E[Y_{x_1, W_{x_0}} - Y_{x_0}],$$



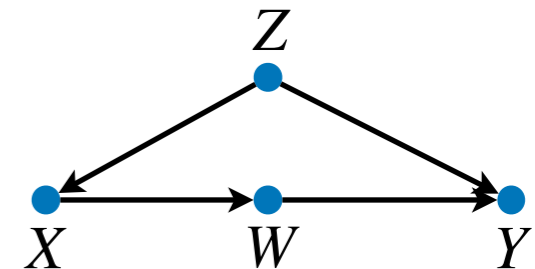
- This quantity requires the nested  $Y_{x_0, W_{x_1}}$ , and applying the CUT we get:

$$P(Y_{x_0, W_{x_1}} = y) = \sum_w P(Y_{x_1 w} = y, W_{x_0} = w) \quad (\text{CUT})$$

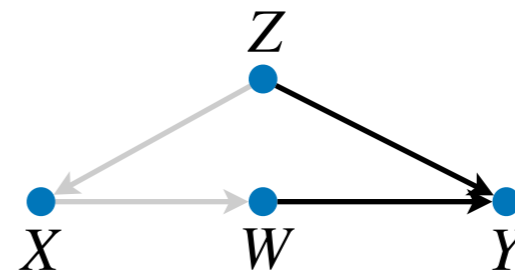
# Counterfactual Ancestors

**Definition – Ancestors (of a counterfactual).** Let  $Y_{\mathbf{x}}$  be such that  $Y \in \mathbf{V}$ ,  $\mathbf{X} \subseteq \mathbf{V}$ . Then, the set of (counterfactual) ancestors of  $Y_{\mathbf{x}}$ , denoted  $An(Y_{\mathbf{x}})$ , consist of each  $W_{\mathbf{z}}$ , such that  $W \in An_{\mathcal{G}_{\bar{\mathbf{X}}}}(Y) \setminus \mathbf{X}$  (which includes  $Y$  itself), and  $\mathbf{z} = \mathbf{x} \cap An_{\mathcal{G}_{\bar{\mathbf{X}}}}(W)$ .

Example. Consider the causal diagram:



- $An(Y_{xw}) = \{Z, Y_w\}$



$$An_{\mathcal{G}_{\bar{XW}}}(Y) \setminus \{X, W\} = \{Z, Y\}$$

*Get ancestors*

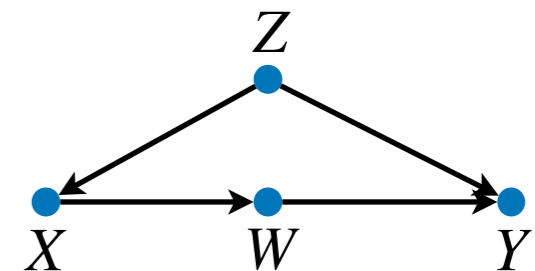
$$\begin{aligned} \emptyset &= \{x, w\} \cap An_{\mathcal{G}_{\bar{XW}}}(Z) \\ \{w\} &= \{x, w\} \cap An_{\mathcal{G}_{\bar{XW}}}(Y) \end{aligned}$$

*Compute subscripts*

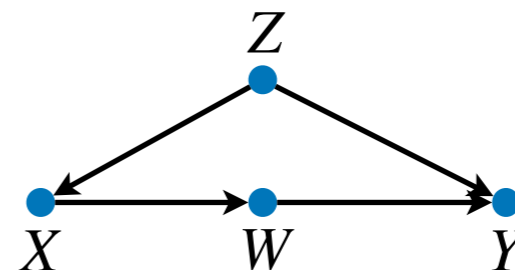
# Counterfactual Ancestors (cont)

**Definition – Ancestors (of a counterfactual).** Let  $Y_{\mathbf{x}}$  be such that  $Y \in \mathbf{V}$ ,  $\mathbf{X} \subseteq \mathbf{V}$ . Then, the set of (counterfactual) ancestors of  $Y_{\mathbf{x}}$ , denoted  $An(Y_{\mathbf{x}})$ , consist of each  $W_{\mathbf{z}}$ , such that  $W \in An_{\mathcal{G}_{\bar{\mathbf{x}}}}(Y) \setminus \mathbf{X}$  (which includes  $Y$  itself), and  $\mathbf{z} = \mathbf{x} \cap An_{\mathcal{G}_{\bar{\mathbf{x}}}}(W)$ .

Example (2). Consider the causal diagram from before:



- $An(Y_{\mathbf{z}}) = \{X_{\mathbf{z}}, W_{\mathbf{z}}, Y_{\mathbf{z}}\}$



$$An_{\mathcal{G}_{\bar{\mathbf{z}}}}(Y) \setminus \{Z\} = \{Y, W, X\}$$

Get ancestors

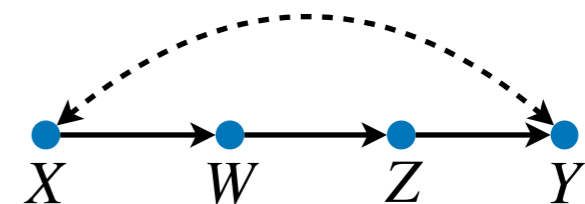
$$\begin{aligned} \{z\} &= \{z\} \cap An_{\mathcal{G}_{\bar{\mathbf{z}}}}(X) \\ \{z\} &= \{z\} \cap An_{\mathcal{G}_{\bar{\mathbf{z}}}}(W) \\ \{z\} &= \{z\} \cap An_{\mathcal{G}_{\bar{\mathbf{z}}}}(Y) \end{aligned}$$

Compute subscripts

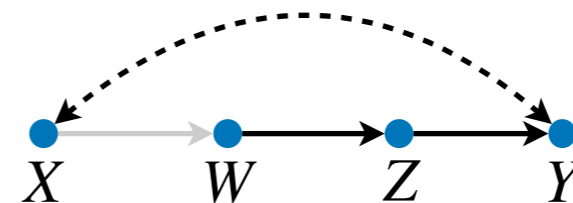
# Counterfactual Ancestors (cont)

**Definition – Ancestors (of a counterfactual).** Let  $Y_{\mathbf{x}}$  be such that  $Y \in \mathbf{V}$ ,  $\mathbf{X} \subseteq \mathbf{V}$ . Then, the set of (counterfactual) ancestors of  $Y_{\mathbf{x}}$ , denoted  $An(Y_{\mathbf{x}})$ , consist of each  $W_{\mathbf{z}}$ , such that  $W \in An_{\mathcal{G}_{\bar{\mathbf{x}}}}(Y) \setminus \mathbf{X}$  (which includes  $Y$  itself), and  $\mathbf{z} = \mathbf{x} \cap An_{\mathcal{G}_{\bar{\mathbf{x}}}}(W)$ .

Example (3). Consider the causal diagram:



- $An(Y_w) = \{Y_w, Z_w\}$



$$An_{\mathcal{G}_{\bar{w}}}(Y) \setminus \{W\} = \{Y, Z\}$$

*Get ancestors*

$$\{w\} = \{w\} \cap An_{\mathcal{G}_{\bar{w}}}(Y)$$

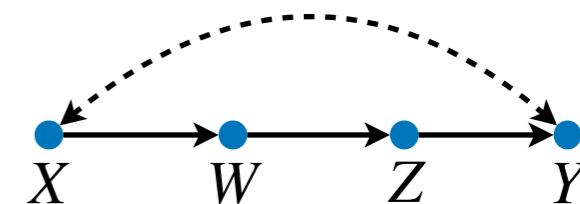
$$\{w\} = \{w\} \cap An_{\mathcal{G}_{\bar{w}}}(Z)$$

*Compute subscripts*

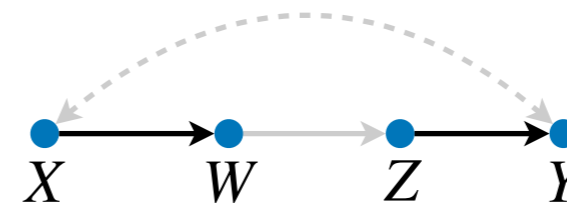
# Counterfactual Ancestors (cont)

**Definition – Ancestors (of a counterfactual).** Let  $Y_{\mathbf{x}}$  be such that  $Y \in \mathbf{V}$ ,  $\mathbf{X} \subseteq \mathbf{V}$ . Then, the set of (counterfactual) ancestors of  $Y_{\mathbf{x}}$ , denoted  $An(Y_{\mathbf{x}})$ , consist of each  $W_{\mathbf{z}}$ , such that  $W \in An_{\mathcal{G}_{\bar{\mathbf{X}}}}(Y) \setminus \mathbf{X}$  (which includes  $Y$  itself), and  $\mathbf{z} = \mathbf{x} \cap An_{\mathcal{G}_{\bar{\mathbf{X}}}}(W)$ .

Example (4). Consider the causal diagram:



- $An(W_{xz}) = \{W_x\}$



$$An_{\mathcal{G}_{\bar{XZ}}}(W) \setminus \{X, Z\} = \{W\}$$

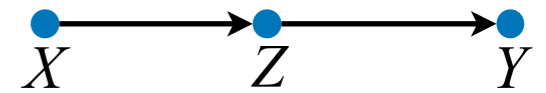
*Get ancestors*

$$\{x\} = \{xz\} \cap An_{\mathcal{G}_{\bar{XZ}}}(W)$$

*Compute subscripts*

# Exclusion Restrictions

- Semantically, one can consider counterfactuals  $Y_{\mathbf{t}}$  for arbitrary  $Y \in \mathbf{V}$  and  $\mathbf{T} \subseteq \mathbf{V}$ .
- However, the variations allowed for counterfactual variables usually depend on the topology and the sparsity of the causal system.
- Consider for instance  $Y_z$  and  $Y_{zx}$  in the chain graph above. Note that for every unit  $\mathbf{U} = \mathbf{u}$ , the variables  $Y_z$  and  $Y_{zx}$  get always the same value. Formally:



$$\mathcal{F}_z = \begin{cases} X_z \leftarrow f_X(U_x) \\ Z_z \leftarrow z \\ Y_z \leftarrow f_Y(z, U_y) \end{cases}$$

$$\mathcal{F}_{zx} = \begin{cases} X_{zx} \leftarrow x \\ Z_{zx} \leftarrow z \\ Y_{zx} \leftarrow f_Y(z, U_y), \end{cases}$$

$$Y_{zx}(\mathbf{u}) = f_Y(z, u_y) = Y_z(\mathbf{u}).$$

# Exclusion Operator

**Lemma – Exclusion Operator.** Consider a causal diagram  $G$  and a counterfactual variable  $Y_{\mathbf{x}}$ . Let the subscript restriction operator be labeled with  $\eta$ , defined by

$$\eta(Y_{\mathbf{x}}) := Y_{\mathbf{z}},$$

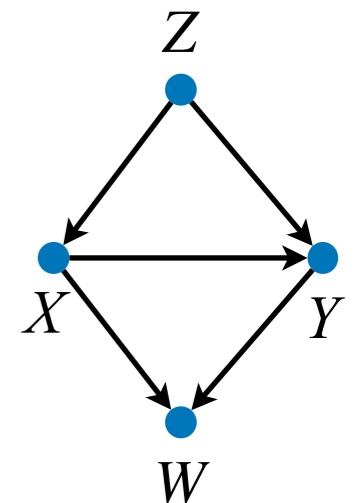
where  $\mathbf{Z} = \mathbf{X} \cap \text{An}_{G_{\bar{\mathbf{X}}}}(Y)$  and  $\mathbf{z} = \mathbf{x} \cap \mathbf{Z}$ . Then,  $\eta(Y_{\mathbf{x}})$  and  $Y_{\mathbf{x}}$  are the same counterfactual variable in any model compatible with  $G$ .

For a set  $\mathbf{Y}_*$ , define

$$\eta(\mathbf{Y}_*) = \bigcup_{Y_{\mathbf{t}} \in \mathbf{Y}_*} \eta(Y_{\mathbf{t}}).$$

# Exclusion Operator Examples

Consider the causal diagram to the right. The following are some counterfactual variables that can be minimized using the exclusion operator:



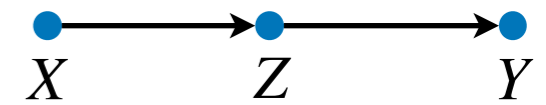
$$\begin{array}{ll}
 \eta(X_{zyw}) = X_z & \{Z, Y, W\} \cap An_{\mathcal{G}_{ZYW}}(X) = \{Z\} \\
 \eta(Y_{xz}) = Y_{xz} & \{X, Z\} \cap An_{\mathcal{G}_{XZ}}(Y) = \{X, Z\} \\
 \eta(W_{zxy}) = W_{xy} & \{Z, X, Y\} \cap An_{\mathcal{G}_{ZXY}}(W) = \{X, Y\} \\
 \eta(W_{zy}) = W_{zy} & \{Z, Y\} \cap An_{\mathcal{G}_{ZY}}(W) = \{Z, Y\}
 \end{array}$$

# Causal Diagrams for Counterfactuals: Ancestral Multi-World Networks (AMWN)

**Input:** graph  $G$ , ctf. variables  $\mathbf{Y}_*$

**Output:** Ctf. ancestral graph  $G_A(\mathbf{Y}_*)$  over ctf. variables  $\mathbf{Y}_*$

**Example:**



- compute all  $An(\mathbf{Y}_*)$  & add them as nodes to initial  $G'$
- add directed arrows witnessing ancestry
- **for each** edge  $U \rightarrow V$  in  $G$ 
  - **if** more than one instance of  $V$  in  $G'$

$$Y_z \perp Z_x \mid Y \quad \text{X}$$

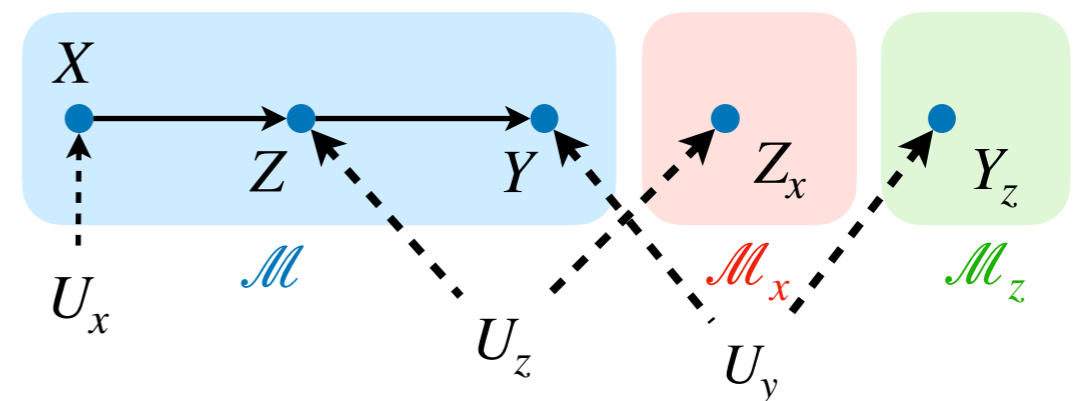
$$An(Y_z) = \{Y_z\}$$

$$An(Z_x) = \{Z_x\}$$

$$An(Y) = \{Y, Z, X\}$$

add node  $U$  with edges to each copy  $V_x$  of  $V$

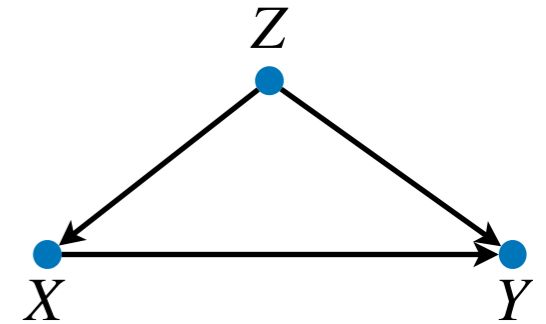
- **return**  $G'$



# ETT Derivation

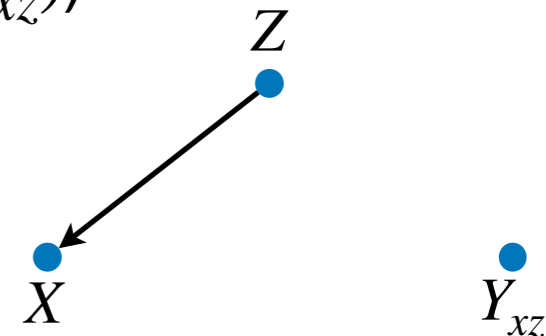
$$P(y_x | x') = \sum_z P(y_x | z, x')P(z | x') \text{ (Conditioning on } Z)$$

$$= \sum_z P(y_x | z_x, x')P(z | x') \text{ (exclusion: } \{X\} \cap An(Z) = \emptyset)$$



$$= \sum_z P(y_{xz} | z_x, x')P(z | x') \text{ (consistency } (Z_x = z \Rightarrow Y_x = Y_{xz}))$$

$$= \sum_z P(y_{xz} | z, x')P(z | x') \text{ (exclusion: } \{X\} \cap An(Z) = \emptyset)$$



$$= \sum_z P(y_{xz} | z, x)P(z | x') \text{ (indep. } Y_{xz} \perp X | Z \text{ in } \mathcal{G}_A)$$

$\mathcal{G}_A(Y_{xz}, X, Z)$

$$= \sum_z P(y | z, x)P(z | x') \text{ (consistency } (Z = z, X = x \Rightarrow Y_{xz} = Y))$$

can be computed from observational data!

# NDE Derivation

$$P(y_{x'w_x}) = \sum_w P(y_{x'w}, w_x) \quad (\text{CUT})$$

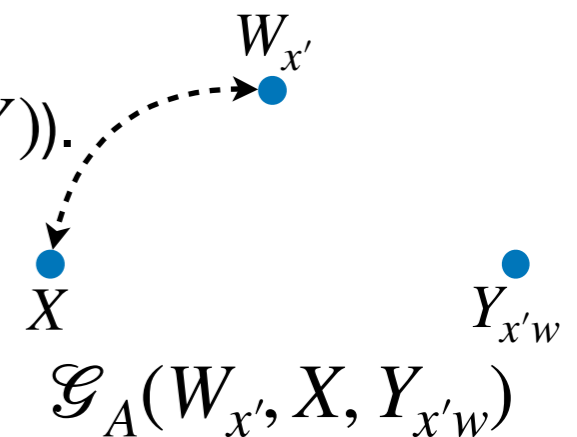
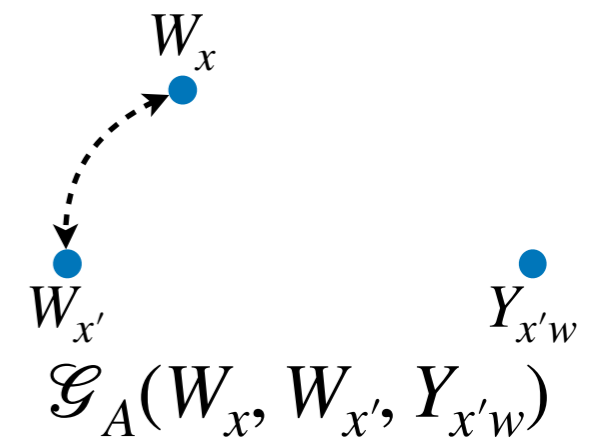
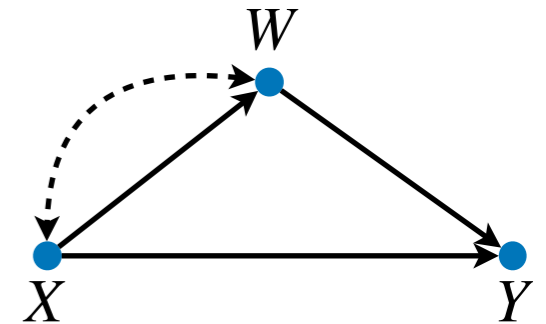
$$= \sum_w P(y_{x'w} | w_x) P(w_x) \quad (\text{Chain rule})$$

$$= \sum_w P(y_{x'w} | w_{x'}) P(w_x) \quad (\text{indep.: } (Y_{x'w} \perp W_x, W_{x'}) \text{ in } \mathcal{G}_A)$$

$$= \sum_w P(y_{x'w} | w_{x'}, x') P(w_x) \quad (\text{indep.: } (Y_{x'w} \perp X | W_{x'}) \text{ in } \mathcal{G}_A)$$

$$= \sum_w P(y_{x'w} | w, x') P(w_x) \quad (\text{consist.: } (X = x' \Rightarrow W_{x'} = W))$$

$$= \sum_w P(y | w, x') P(w_x) \quad (\text{consist.: } (W = w, X = x' \Rightarrow Y_{x'w} = Y))$$



# Counterfactual Calculus

**Theorem.** Let  $\mathcal{G}$  be a causal diagram, then for  $Y, X, Z, W, T \subseteq V$ , the following rules hold for the probability distributions generated by any model compatible with  $\mathcal{G}$ :

**Rule 1 (Consistency rule – Observation/intervention exchange)**

$$P(\mathbf{y}_{T_x}, \mathbf{x}_T, \mathbf{w}_*) = P(\mathbf{y}_T, \mathbf{x}_T, \mathbf{w}_*)$$

**Rule 2 (Independence Rule – Adding/removing counterfactual observations)**

$$P(\mathbf{y}_r \mid \mathbf{x}_t, \mathbf{w}_*) = P(\mathbf{y}_r \mid \mathbf{w}_*) \quad \text{if } (\mathbf{Y}_r \perp \mathbf{X}_t \mid \mathbf{W}_*) \text{ in } \mathcal{G}_A,$$

**Rule 3 (Exclusion Rule – Adding/removing interventions)**

$$P(\mathbf{y}_{xz}, \mathbf{w}_*) = P(\mathbf{y}_z, \mathbf{w}_*) \quad \text{if } \mathbf{X} \cap An(\mathbf{Y}) = \emptyset \text{ in } \mathcal{G}_{\bar{z}},$$

where  $\mathcal{G}_A$  is the counterfactual ancestral graph  $\mathcal{G}_A(\mathbf{Y}_r, \mathbf{X}_t, \mathbf{W}_*)$ .