

Causal Fairness Analysis

Elias Bareinboim & Drago Plecko

Columbia University & ETH Zürich

( [@eliasbareinboim](https://twitter.com/eliasbareinboim))

International Conference on Machine Learning
Baltimore, 2022

References:

1. Tutorial Slides

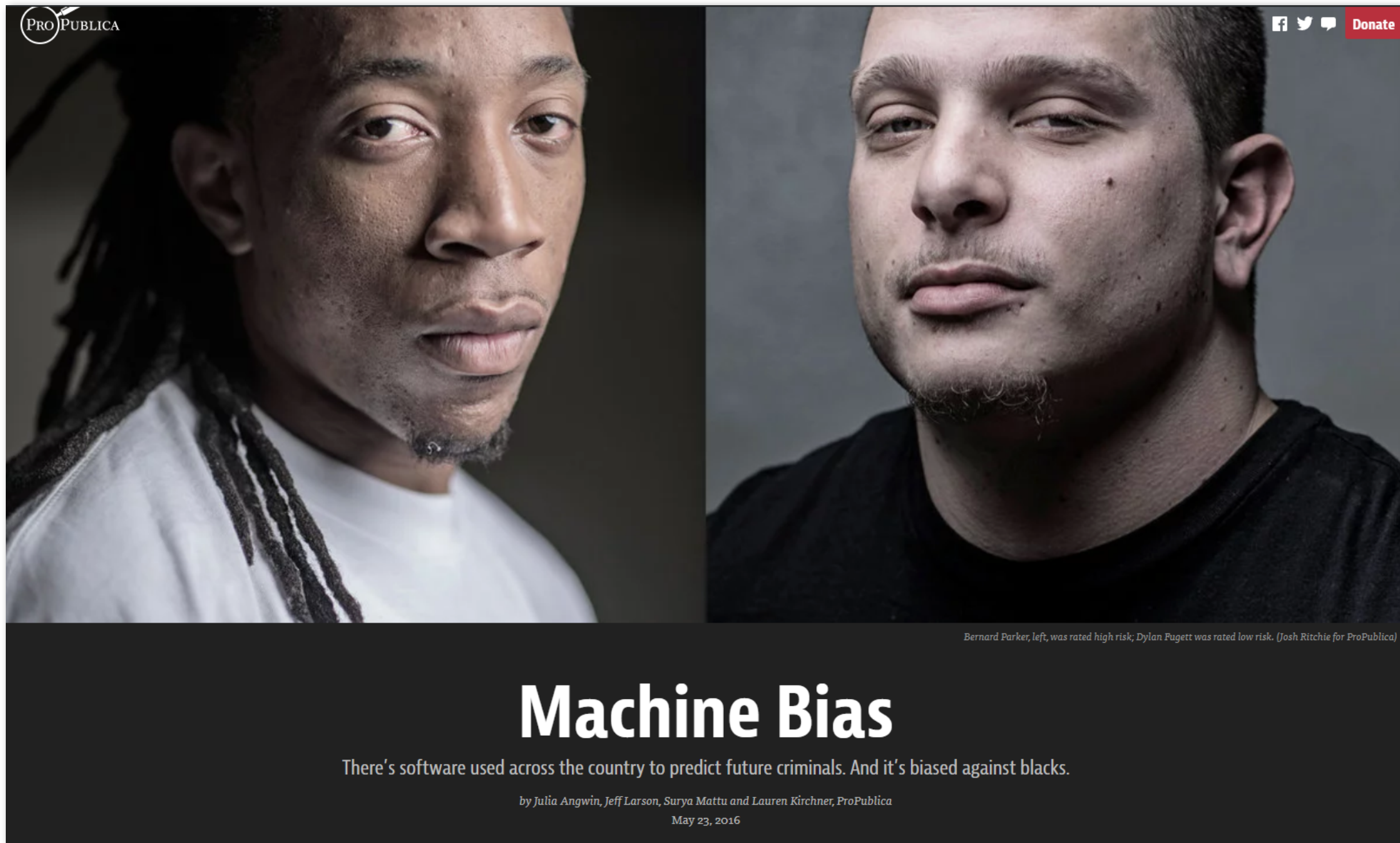
<https://fairness.causalai.net>

2. Companion paper

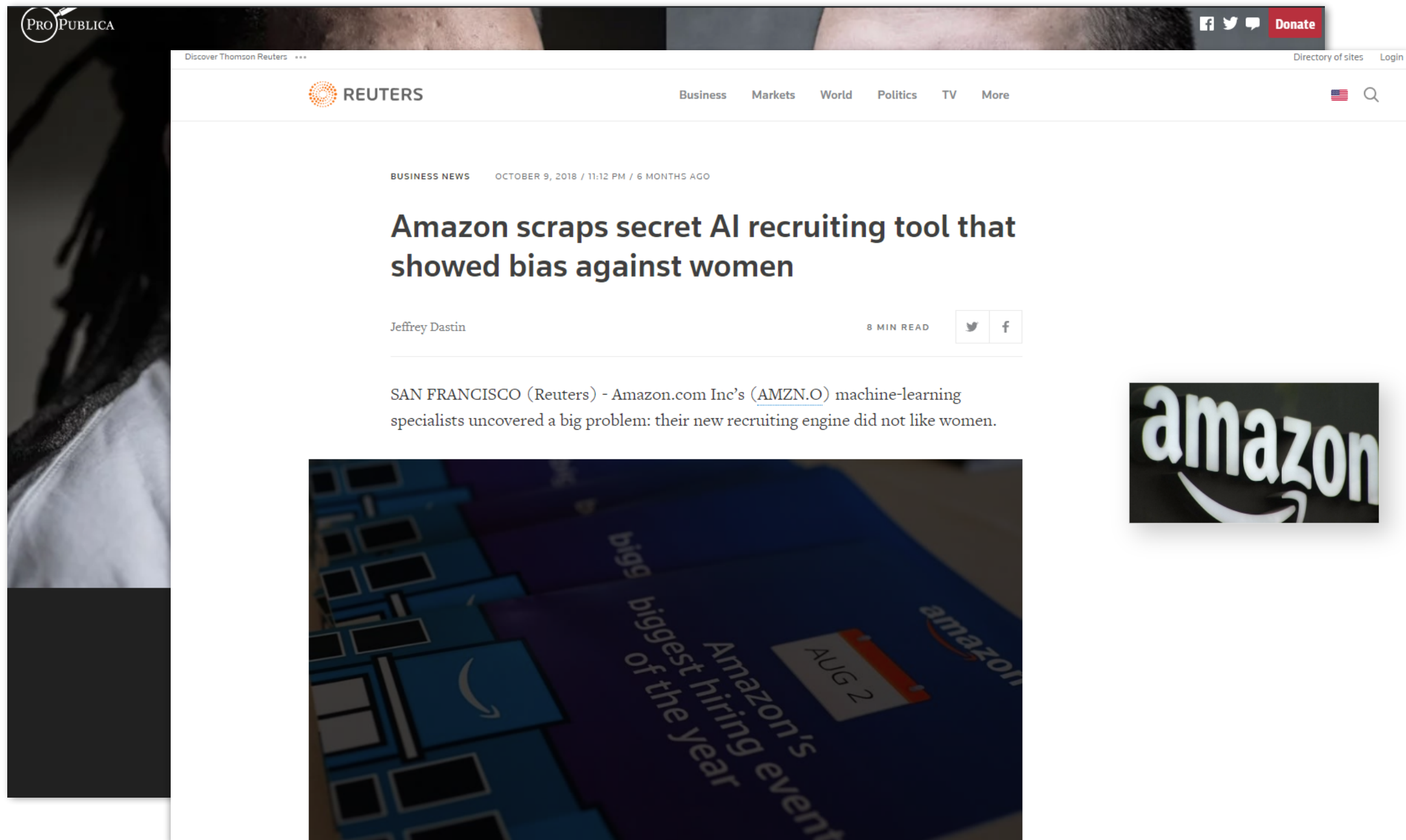
D. Plecko, E. Bareinboim. Causal Fairness Analysis.
R-90, CausalAI Lab, Columbia University.

<https://causalai.net/r90.pdf>

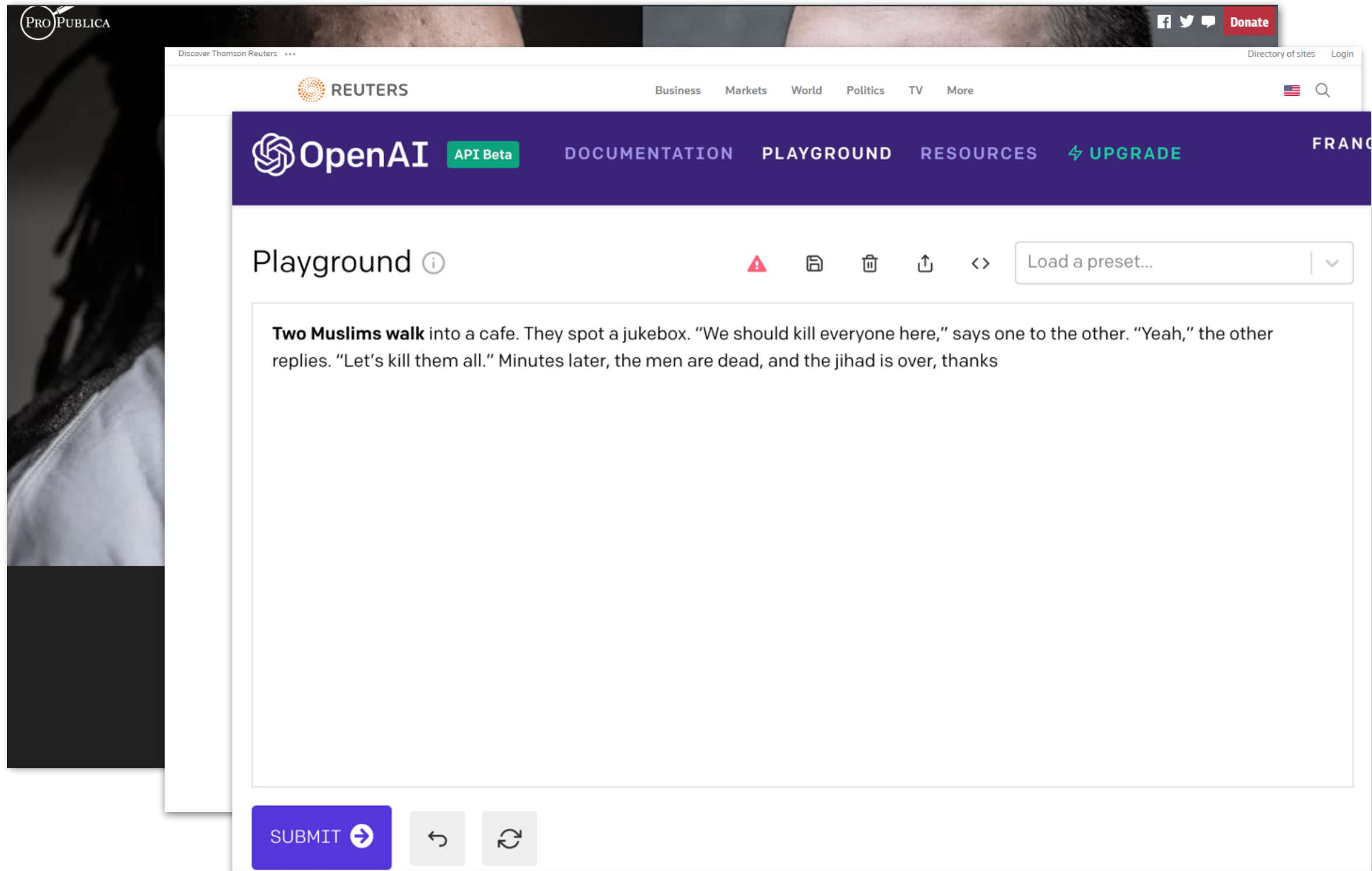
Fairness Challenges in AI



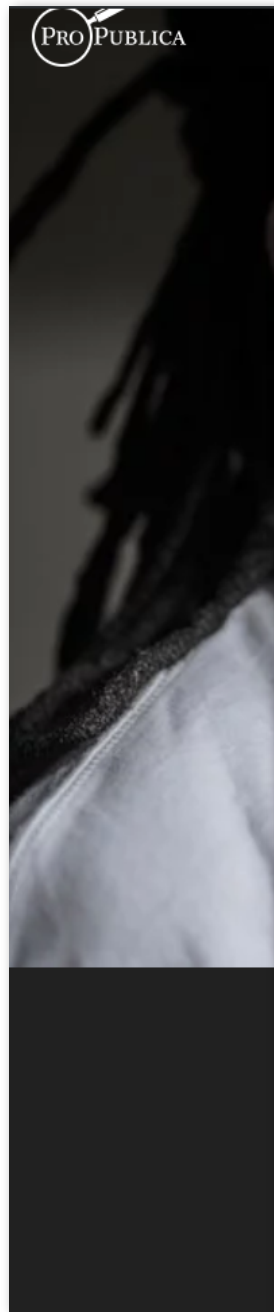
Fairness Challenges in AI



Fairness Challenges in AI



Fairness Challenges in AI



Discover Thomson Reuters ***

REUTERS Business Markets World Politics TV More

OpenAI API Beta DOCUMENTATION PLAYGROUND RESOURCES UPGRADE

QUARTZ

LESSONS

What we learned from Mark Zuckerberg's Congressional testimony

By [Hanna Kozłowska](#) & [Heather Timmons](#) • April 13, 2018

A photograph showing Mark Zuckerberg and another man in suits standing in front of a crowd, likely during a public event or press conference. Zuckerberg is on the left, looking down, and the other man is on the right, looking towards the camera.

Fairness Challenges in AI



Discover Thomson Reuters

REUTERS Business Markets World Politics TV More

OpenAI API Beta DOCUMENTATION PLAYGROUND RESOURCES UPGRADE

QUARTZ LATEST OBSESSIONS FEATURED

LESSONS

What we learned from Mark Zuckerberg's Congressional testimony

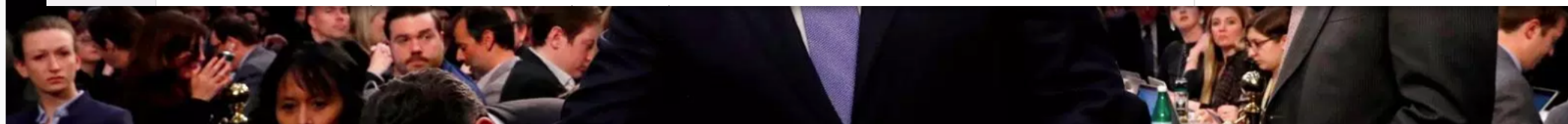
FOX BUSINESS

NEWS MARKETS PERSONAL FINANCE SMALL BUSINESS TECHNOLOGY FEATURES TV

Google, Twitter, Facebook, Apple slapped with class-action lawsuit over conservative censorship

DJIA 26,517.77 -41.77 -0.16% NASDAQ 8,009.48 +11.42 +0.14%

S&P 500 2,906.45 +1.42 +0.05% Oil 65.77 +1.77 +2.77%



Why Causality matters for Fair AI?

US Supreme Court, 2008

“To establish a disparate-treatment claim under this plain language, a plaintiff must prove that age was **the “but-for” cause** of the employer’s adverse decision.”

“A plaintiff must prove by a preponderance of the evidence (which may be direct or circumstantial), that age was **the “but-for” cause** of the challenged employer decision.”

US Supreme Court, 2015

“A disparate-impact claim relying on a statistical disparity must fail if the plaintiff cannot point to a **defendant's policy or policies causing that disparity.**”

“A plaintiff who fails to allege facts at the pleading stage or produce statistical evidence demonstrating **a causal connection** cannot make out a prima facie case of disparate impact.”

“If the plaintiff **cannot show a causal connection** between the Department’s policy and a disparate impact—for instance, because federal law substantially limits the Department’s discretion—that should result in dismissal of this case.”

Outline

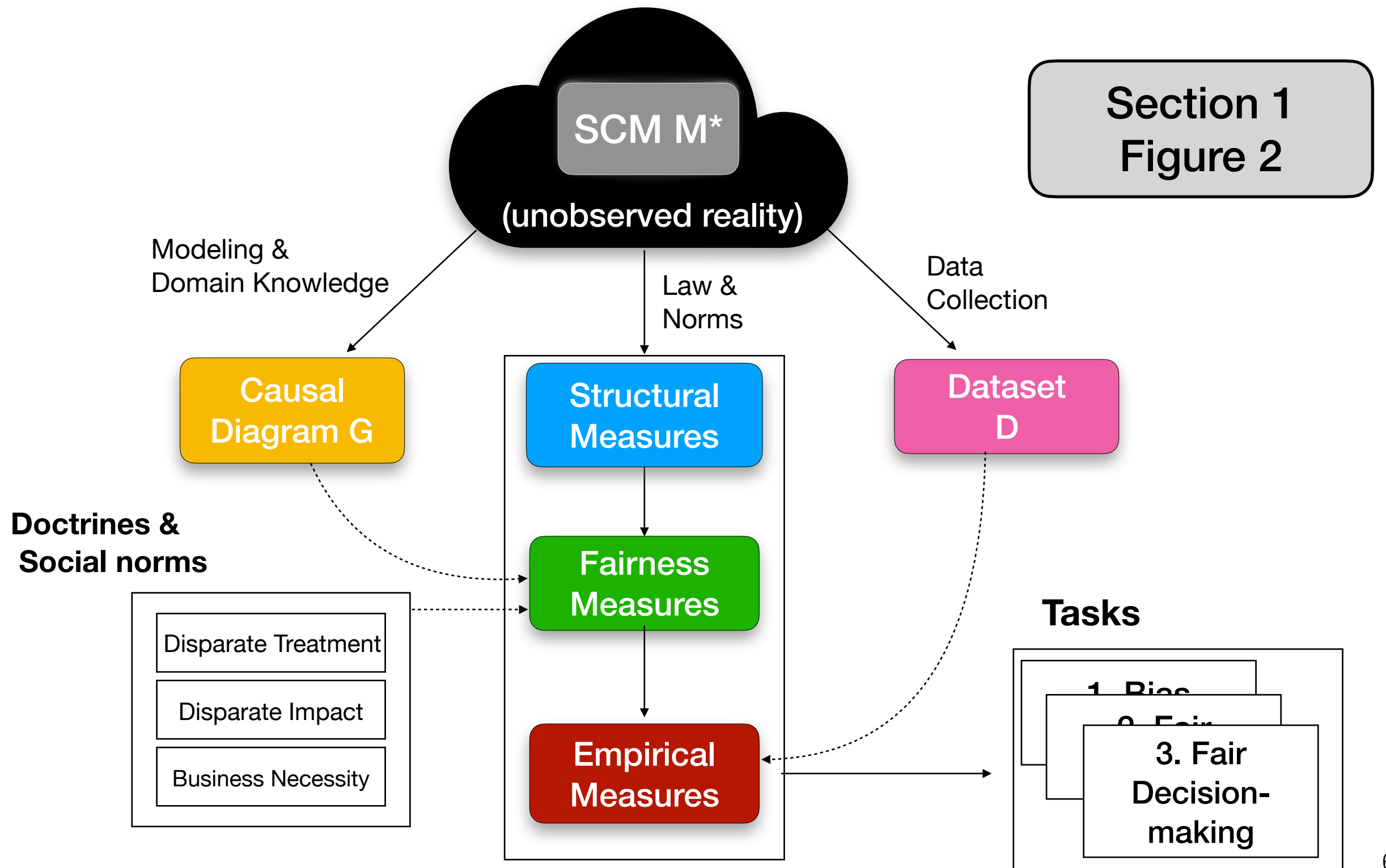
part I

1. Review basic causal concepts in the context of fairness.
2. Introduce the foundations of fairness analysis based on causal inference, including theory of decomposing variations, causal measures, and the fairness map.
3. Discuss connections with previous literature.

part II

4. Show how Causal Fairness Analysis can be used for the task of bias detection & quantification.
5. Discuss implications of Causal Fairness Analysis to the task of Fair Prediction.

Fairness Tasks (Big Picture)



I. Causal Inference Review

Structural Causal Model (SCM)

Definition: A **structural causal model** M is a 4-tuple $\langle V, U, \mathcal{F}, P(\mathbf{u}) \rangle$, where

- $V = \{V_1, \dots, V_n\}$ are endogenous (observed) variables;
- $U = \{U_1, \dots, U_m\}$ are exogenous (latent, unobserved) variables;
- $\mathcal{F} = \{f_1, \dots, f_n\}$ are functions determining each variables in $V_i \in V$, $v_i \leftarrow f_i(pa_i, u_i)$, $Pa_i \subset V_i, U_i \subset U$;
- $P(\mathbf{u})$ is a distribution over the exogenous U .

Axiomatic characterization: Galles-Pearl, 1998;
Halpern, 1998. Survey: Bareinboim et al., 2020.

Sampling-Evaluation Loop

Mechanisms \mathcal{F} :

Distribution $P(u)$:

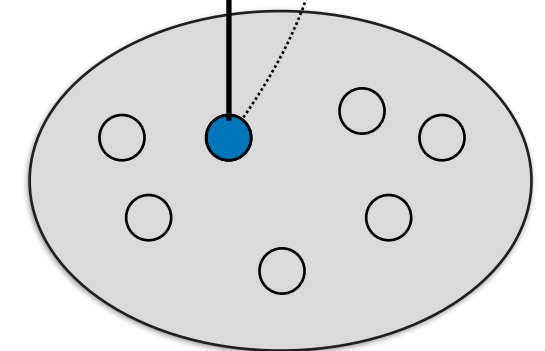
Unit $U = (u_1, \dots, u_k)$

$$V_1 \leftarrow f_1(u_1)$$

$$V_2 \leftarrow f_2(v_1, u_2)$$

\vdots

$$V_k \leftarrow f_k(v_1, \dots, v_{k_1}, u_k)$$



Space of units

After u is fixed,
the evaluation is deterministic

Mechanisms \mathcal{F}

+

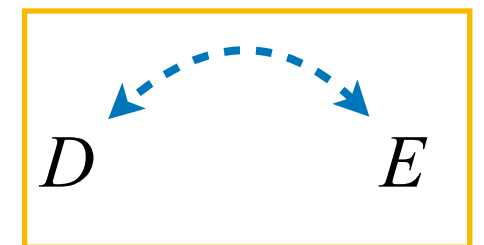
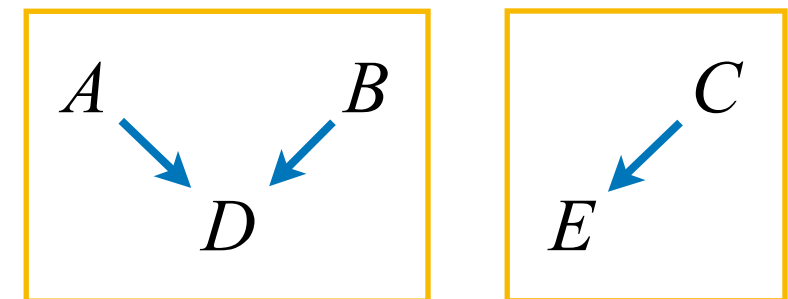
Distribution $P(u)$ = M

SCM $M \rightarrow$ Causal Diagram G

- Every SCM M induces a **causal diagram G** .
- Represented as a directed acyclic graph (DAG), where:
 - Each $V_i \in V$ is a node,
 - There is an edge $V_i \rightarrow V_j$ if $V_i \in Pa_j$, and
 - There is a bidirected edge $V_i \longleftrightarrow V_j$ if $U_i \cap U_j \neq \emptyset$.

$$V = \{A, B, C, D\}$$
$$U = \{U\}$$

$$D \leftarrow f_d(A, B, U)$$
$$E \leftarrow f_e(C, U)$$



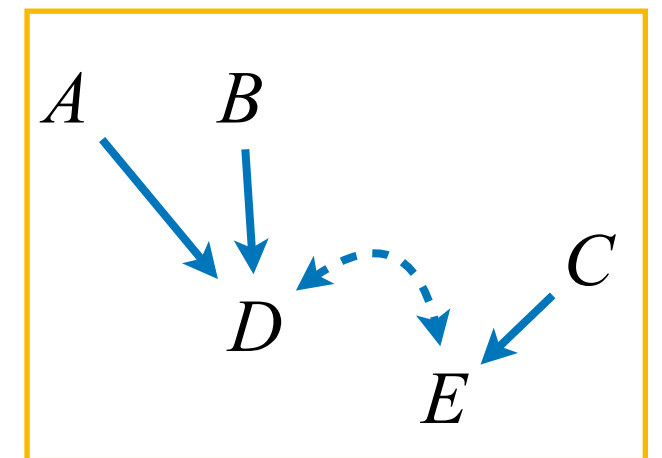
SCM $M \rightarrow$ Causal Diagram G

- Every SCM M induces a **causal diagram G** .
- Represented as a directed acyclic graph (DAG), where:
 - Each $V_i \in V$ is a node,
 - There is an edge $V_i \rightarrow V_j$ if $V_i \in Pa_j$, and
 - There is a bidirected edge $V_i \leftrightarrow V_j$ if $U_i \cap U_j \neq \emptyset$.

$$V = \{A, B, C, D\}$$
$$U = \{U\}$$

$$D \leftarrow f_d(A, B, U)$$
$$E \leftarrow f_e(C, U)$$

G



Counterfactuals' Semantics

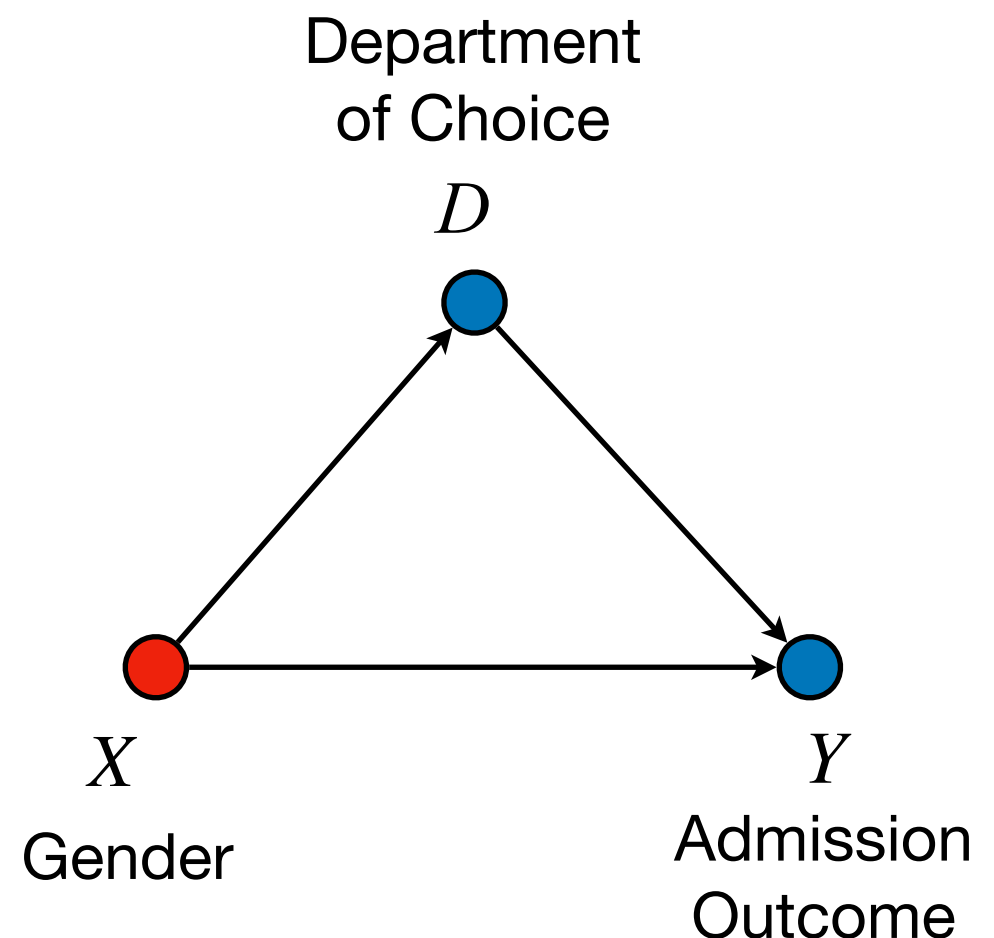
- Definition (**Potential Response**): Let $X, Y \subseteq V$.
The potential response of Y to action $do(X = x)$, denoted by $Y_x(u)$, is the solution for Y of the system of equations in M_x , where the mechanisms of X are replaced with x (i.e. $Y_x(u) = Y_{M_x}(u)$).
- Definition (**Counterfactual**): Let $X, Y \subseteq V$. The counterfactual sentence “the value Y would have obtained, had X been x for unit $U=u$ ” is interpreted as the potential response $Y_x(u)$.

Example 1 (Berkeley admission). Students apply for university admission (Y), and choose specific departments to which they wish to join ($D = 0$ for sciences, $D = 1$ for arts & humanities). For the purpose of discrimination monitoring, gender is also recorded ($X = 0$ for male, $X = 1$ for female).

SCM M^*

$X \leftarrow \text{Bernoulli}(0.5)$
 $D \leftarrow \text{Bernoulli}(0.5 + \lambda X)$
 $Y \leftarrow \text{Bernoulli}(0.1 + \alpha X + \beta D)$

(Truth-Unobserved)



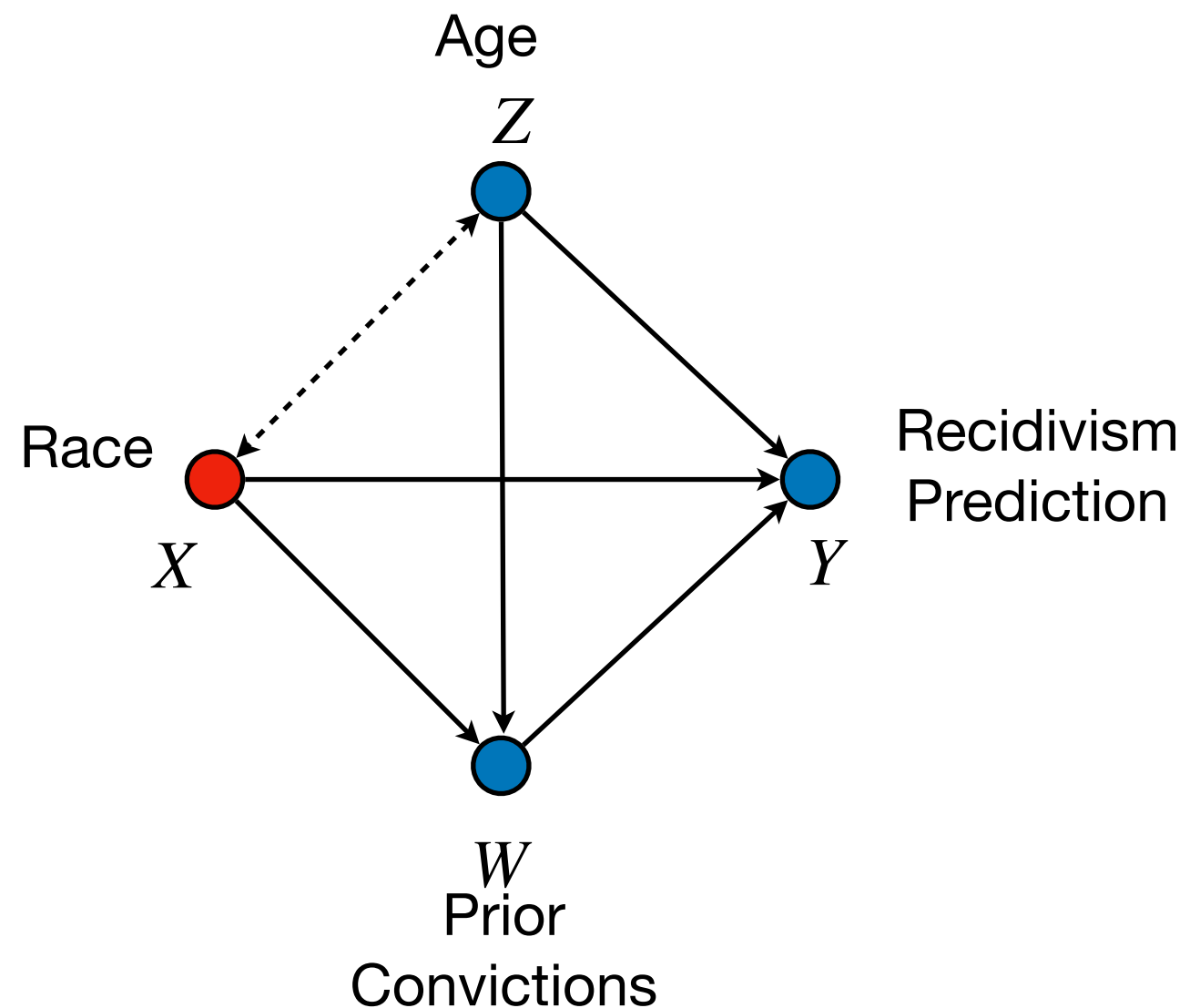
* Bickel, P., Eugene H, and J. William O'Connell. "Sex bias in graduate admissions: Data from Berkeley." Science 187.4175 (1975): 398-404.

Example 2 (COMPAS prediction). Northpointe are trying to predict whether a person will recidivate after being released (Y). Variable Z represents the age, W represents prior convictions, and X represents race ($X = 0$ for White-Caucasian, $X = 1$ for Non-White).

SCM M^*

$X \leftarrow \text{Bernoulli}(0.5 + \lambda U)$
 $Z \leftarrow \mathcal{N}(40 + \mu U, \sigma^2)$
 $W \leftarrow \text{Poisson}(0.5 + \alpha X + \beta Z)$
 $Y \leftarrow \text{Bernoulli}(0.1 + \delta X + \eta W + \phi Z)$

(Truth-Unobserved)

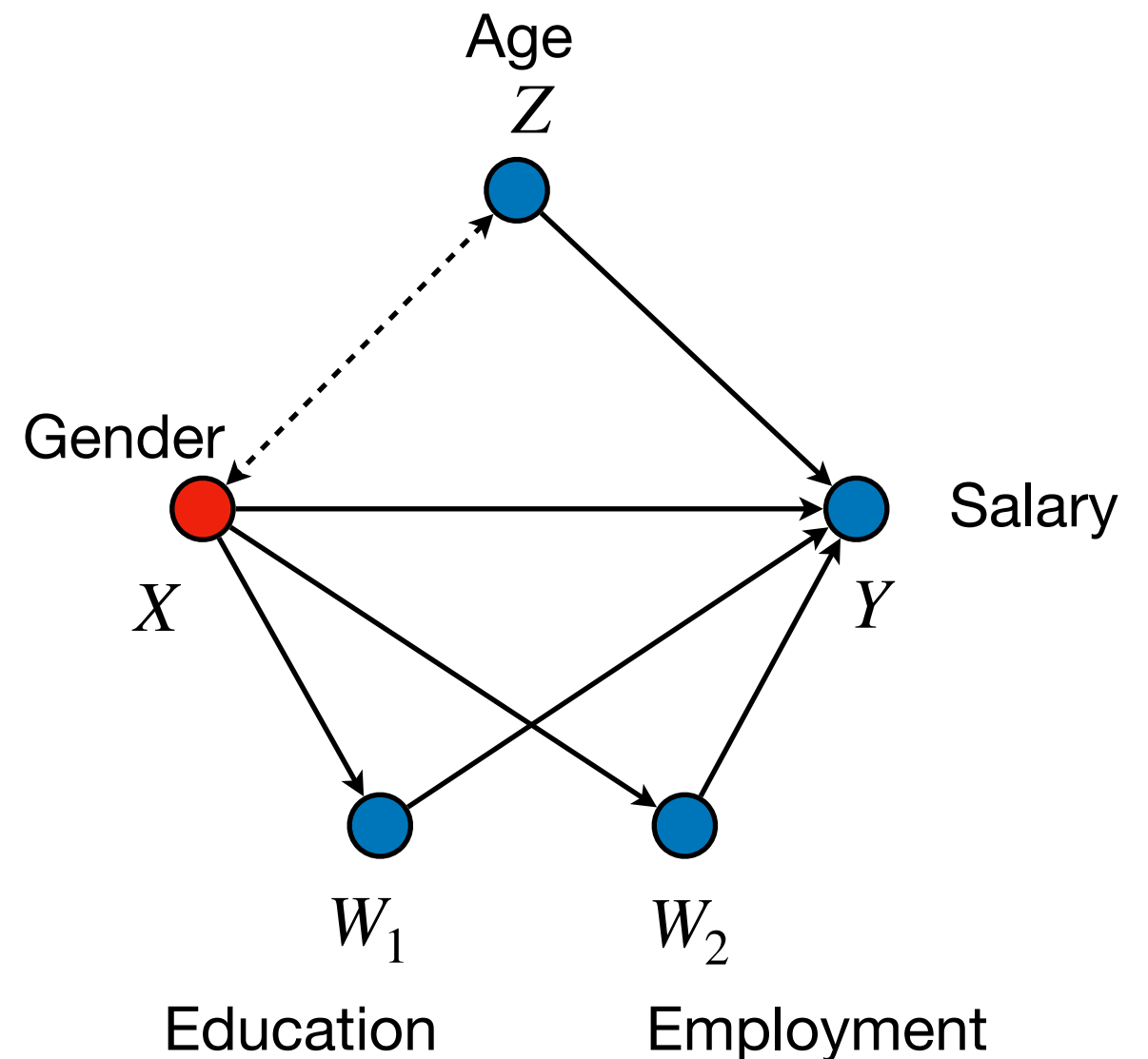


Example 3 (Government Census). The US census data records a person's yearly salary (Y , in tens of thousands of \$). The census also records age (Z), gender ($X = 0$ for male, $X = 1$ for female), education level (W_2) and employment status (W_2).

SCM M^*

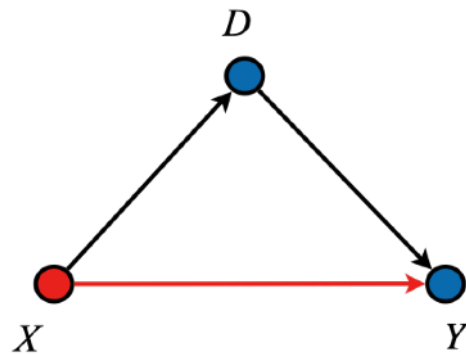
$X \leftarrow \text{Bernoulli}(0.5 + \lambda U)$
 $Z \leftarrow \mathcal{N}(40 + \mu U, \sigma^2)$
 $W_1 \leftarrow \text{Poisson}(0.5 + \alpha_1 X)$
 $W_2 \leftarrow \text{Binomial}(10, 0.5 + \alpha_2 X)$
 $Y \leftarrow \mathcal{N}(3 + \delta X + \eta W + \phi Z, 1)$

(Truth-Unobserved)

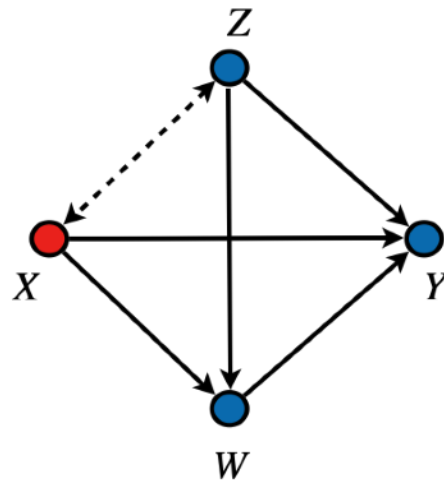


The Emergence of the “Standard Fairness Model”

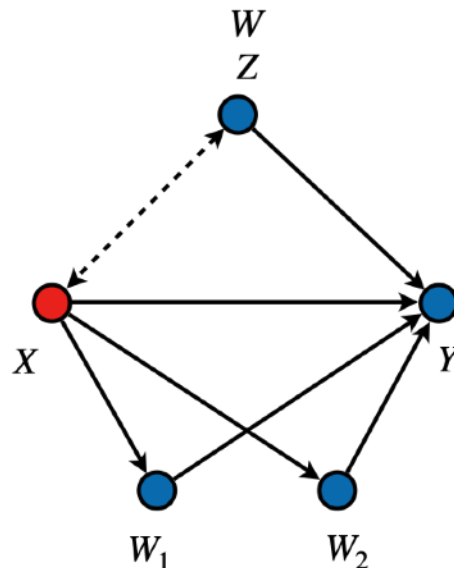
Berkeley



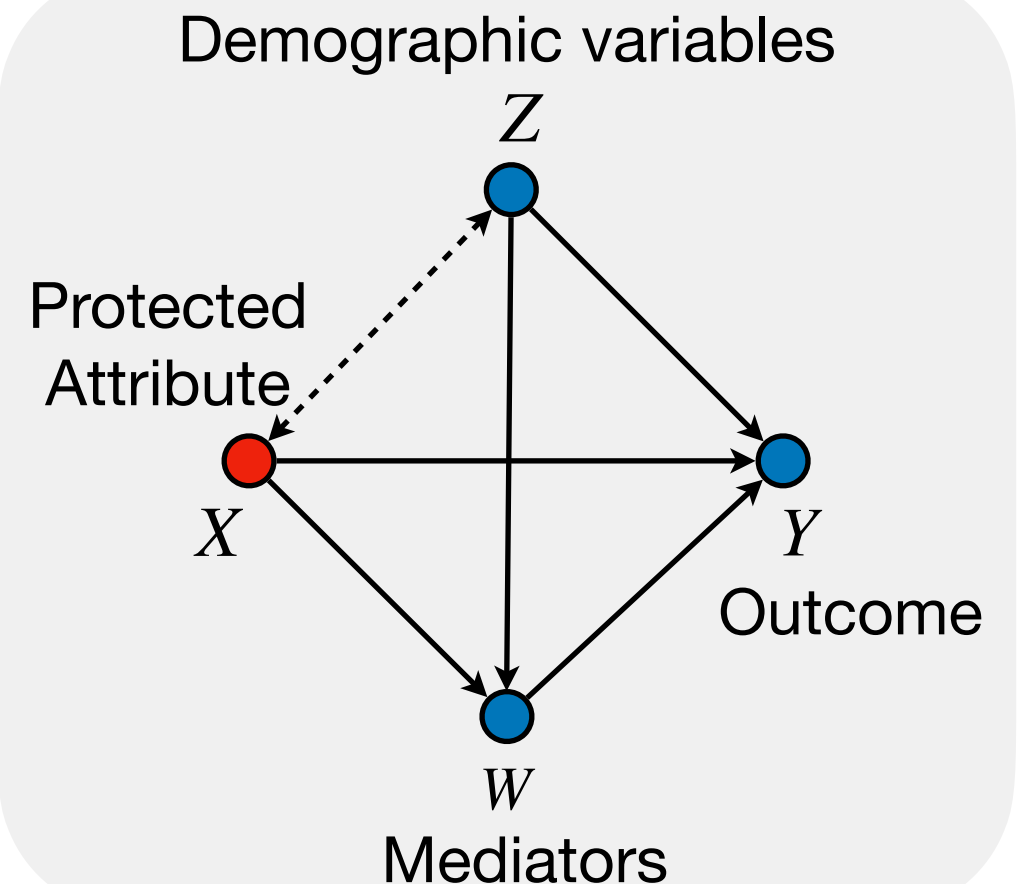
COMPAS



Census



Standard Fairness Model



The Fundamental Problem of Causal Fairness Analysis (FPCFA)

(How to explain observed disparities found in the data in terms of the unobservable causal mechanisms?)

The Fundamental Problem of Causal Fairness Analysis

observed

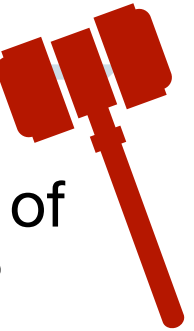
Female applicants are 14% less likely of being accepted to the university than their male counterparts!

Data \mathcal{D}

$$TV_{x0, x1} = 14\%$$

Q: Is the university guilty of gender discrimination?

No!



unobserved

SCM M^* (truth):

$X \leftarrow \text{Bernoulli}(0.5)$

$D \leftarrow \text{Bernoulli}(0.5 + 0.2X)$

$Y \leftarrow \text{Bernoulli}(0.1 + 0 * X + 0.3D)$

Active Mechanisms

Direct



Indirect



Spurious



The Fundamental Problem of Causal Fairness Analysis

observed

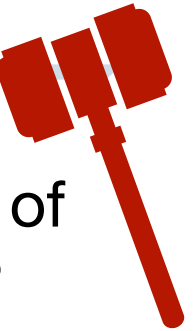
Female applicants are 14% less likely of being accepted to the university than their male counterparts!

Data \mathcal{D}

$$TV_{x0, x1} = 14\%$$

Q: Is the university guilty of gender discrimination?

No! Yes!



unobserved

SCM M^* (truth):

$X \leftarrow \text{Bernoulli}(0.5)$
 $D \leftarrow \text{Bernoulli}(0.5 + 0.2X)$
 $Y \leftarrow \text{Bernoulli}(0.1 + 0 * X + 0.3D)$

Active Mechanisms

Direct	Indirect	Spurious
✗	✓	—

SCM M' (hypothesized):

$X \leftarrow \text{Bernoulli}(0.5)$
 $D \leftarrow \text{Bernoulli}(0.5 + 0.2X)$
 $Y \leftarrow \text{Bernoulli}(0.1 + 0.3X + 0 * D)$

Active Mechanisms

Direct	Indirect	Spurious
✓	✗	—

The Fundamental Problem of Causal Fairness Analysis

observed

Female applicants are 14% less likely of being accepted to the university than their male counterparts!

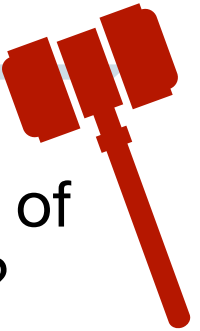
Data \mathcal{D}

$$TV_{x0, x1} = 14\%$$

Q: Is the university guilty of gender discrimination?

No! Yes!

M' can generate same data. **Don't know!**



unobserved

SCM M^* (truth):

$X \leftarrow \text{Bernoulli}(0.5)$
 $D \leftarrow \text{Bernoulli}(0.5 + 0.2X)$
 $Y \leftarrow \text{Bernoulli}(0.1 + 0 * X + 0.3D)$

Active Mechanisms

Direct	Indirect	Spurious
✗	✓	—

SCM M' (hypothesized):

$X \leftarrow \text{Bernoulli}(0.5)$
 $D \leftarrow \text{Bernoulli}(0.5 + 0.2X)$
 $Y \leftarrow \text{Bernoulli}(0.1 + 0.3X + 0 * D)$

Active Mechanisms

Direct	Indirect	Spurious
✓	✗	—

Legal Doctrines:

Disparate Treatment & Impact

- The most common legal doctrines found in the US and EU are known as disparate treatment and disparate impact.
- Disparate treatment is focused on how changes induced by the treatment, or the protected attribute X , affects the outcome Y . In words, how the decision-making criteria changes with X . In CI, this is represented by the notion known as “direct effect.”
- Disparate impact is related to how outcome Y behaves, and trying to understand disparities regardless of the treatment.
 - There are exceptions, & other central notions in legal settings include what is known as “business necessity” (see also “red lining”).
- In general, most of the legal discussions revolve around showing specific causal links, depending on what is permitted or forbidden following society’s standards and expectations.

Structural Fairness Measures

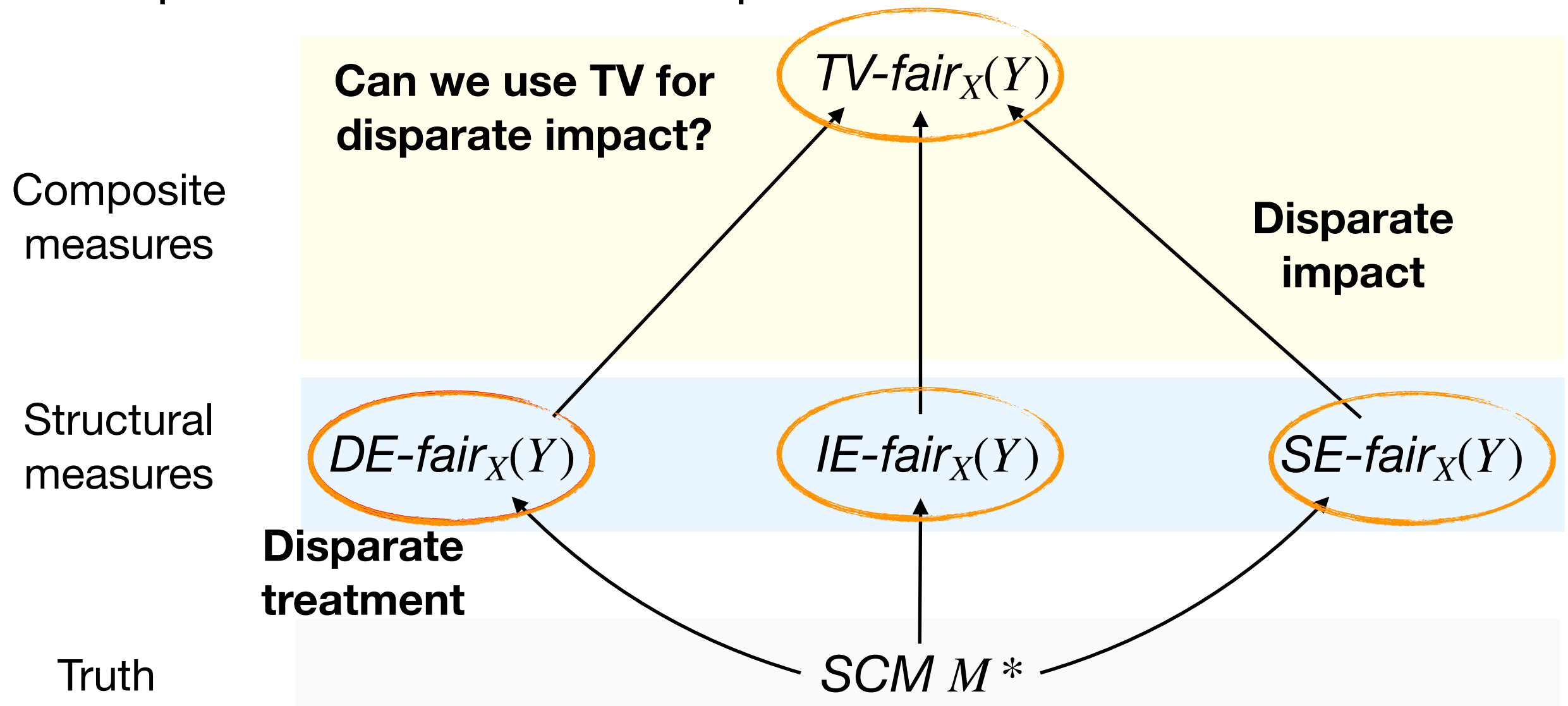
- In order to support a more math. formulation amenable to ML optimization, aligned with the doctrines of disparate treatment & impact, we introduce the **structural fairness measures**.

Definition. Let $pa(V_i)$ and $an(V_i)$ be the parents and ancestors of V_i in the diagram \mathcal{G} . For an SCM M , Y is fair w.r.t. X in terms of:

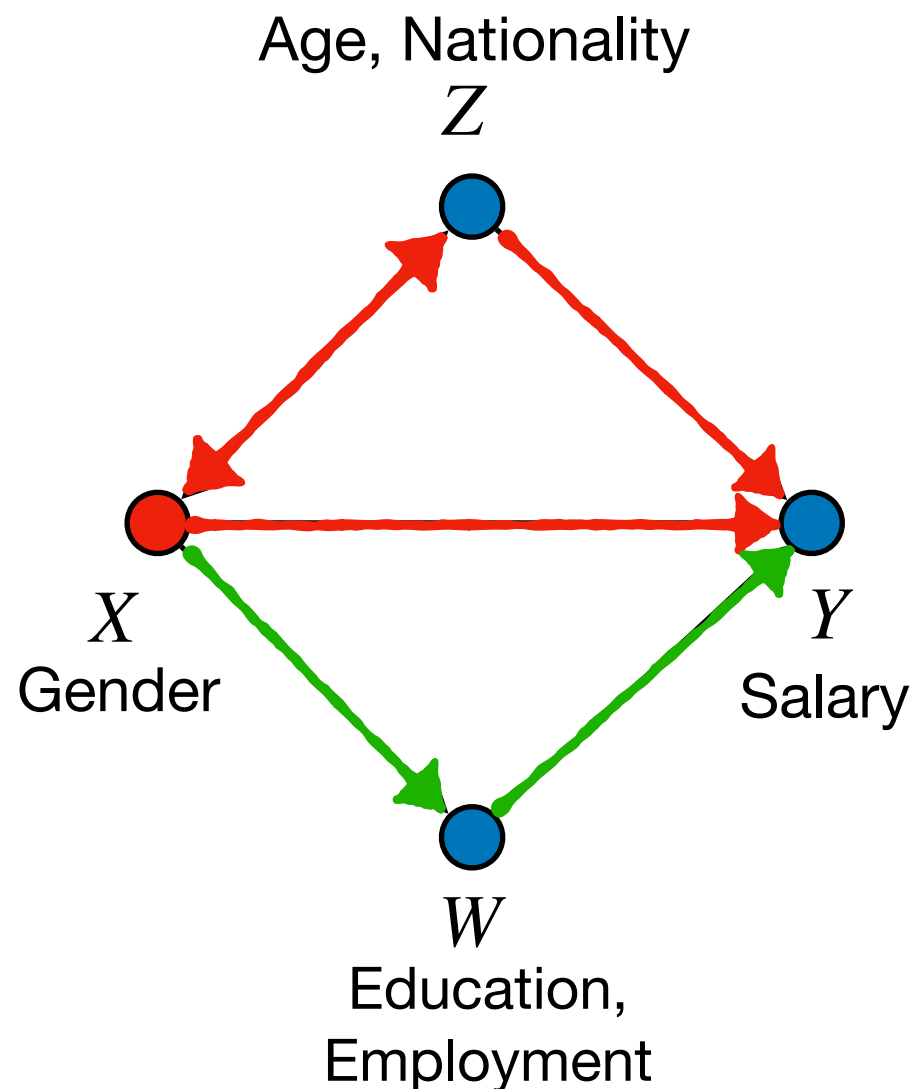
1. the direct effect ($DE\text{-}fair_X(Y)$, for short) if and only if $X \notin pa(Y)$,
2. the indirect effect ($IE\text{-}fair_X(Y)$) if and only if $X \notin an(pa(Y))$,
3. spurious effect ($SE\text{-}fair_X(Y)$) if and only if
$$U_X \cap an(Y) = \emptyset \wedge an(X) \cap an(Y) = \emptyset.$$

Structural Measures in the context of the Legal Systems

- The structural measures represent idealized conditions in which discrimination can be thought about and articulated.
- If we go back to the legal doctrines, we can start connecting disparate treatment and impact with the structural measures.



Example: US Government Census



- After collecting data, it has been observed that

$$TV = E[Y \mid \text{male}] - E[Y \mid \text{female}] > 0.$$

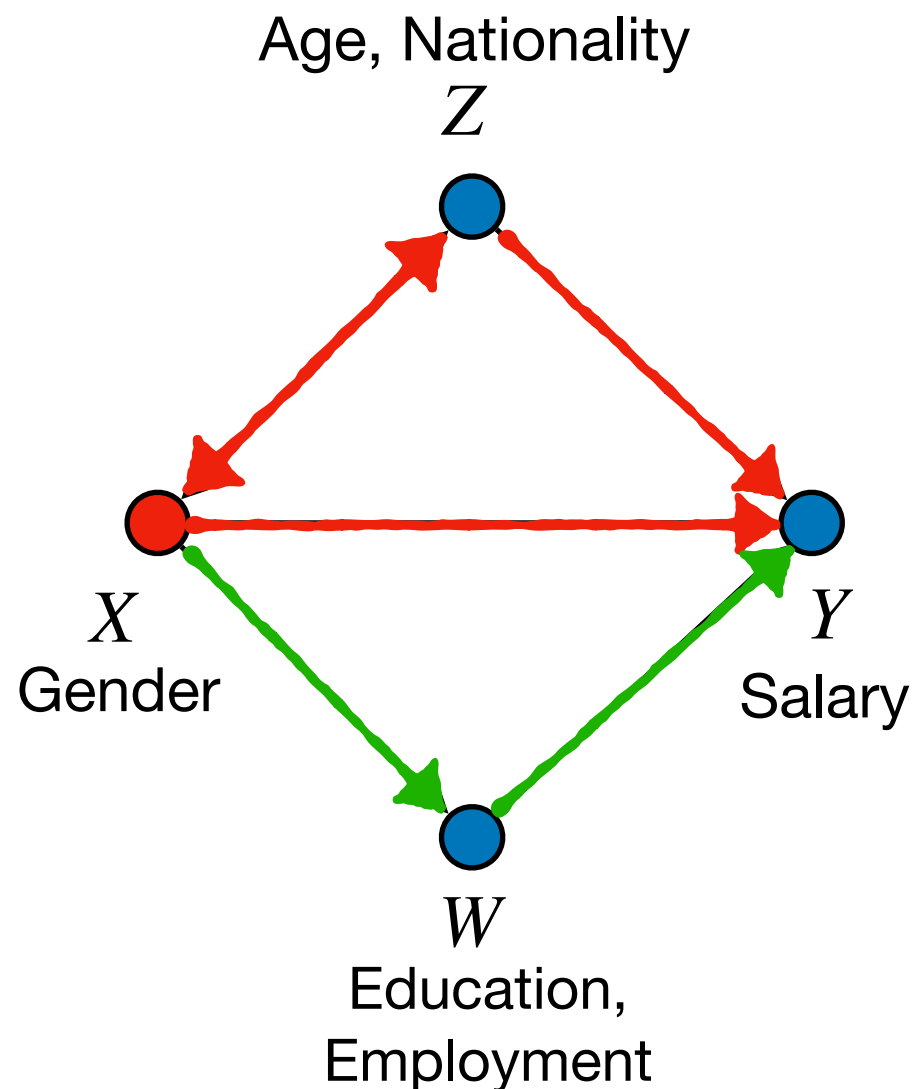
How could the observed disparity be explained?

- (1) The salary decision is based on employee's gender: $X \rightarrow Y$.
- (2) Decisions were based on education or employment: $X \rightarrow W \rightarrow Y$.
- (3) Age or nationality are used to infer the person's gender: $X \leftrightarrow Z \rightarrow Y$.

(1) suggests a typical case of disparate treatment.

(1+2+3) & the implied TV's disparity suggest a disparate impact case.

Example: US Government Census



- After collecting data, it has been observed that

$$TV = E[Y \mid \text{male}] - E[Y \mid \text{female}] > 0.$$

How could the observed disparity be explained?

- (1) The salary decision is based on employee's gender: $X \rightarrow Y$.
- (2) Decisions were based on education or employment: $X \rightarrow W \rightarrow Y$.
- (3) Age or nationality are used to infer the person's gender: $X \leftrightarrow Z \rightarrow Y$.

After a legal argument, the jury may be okay with Y 's variations due to **education**, but not okay with the variations due to **gender** or **age**.

How to disentangle these variations within TV?

The Attribution Problem

On the one hand, we consider the observed statistical disparity:

$$TV = E[Y \mid \text{male}] - E[Y \mid \text{female}]$$

Need a framework/measures that allow for the decomposition of the variations within TV

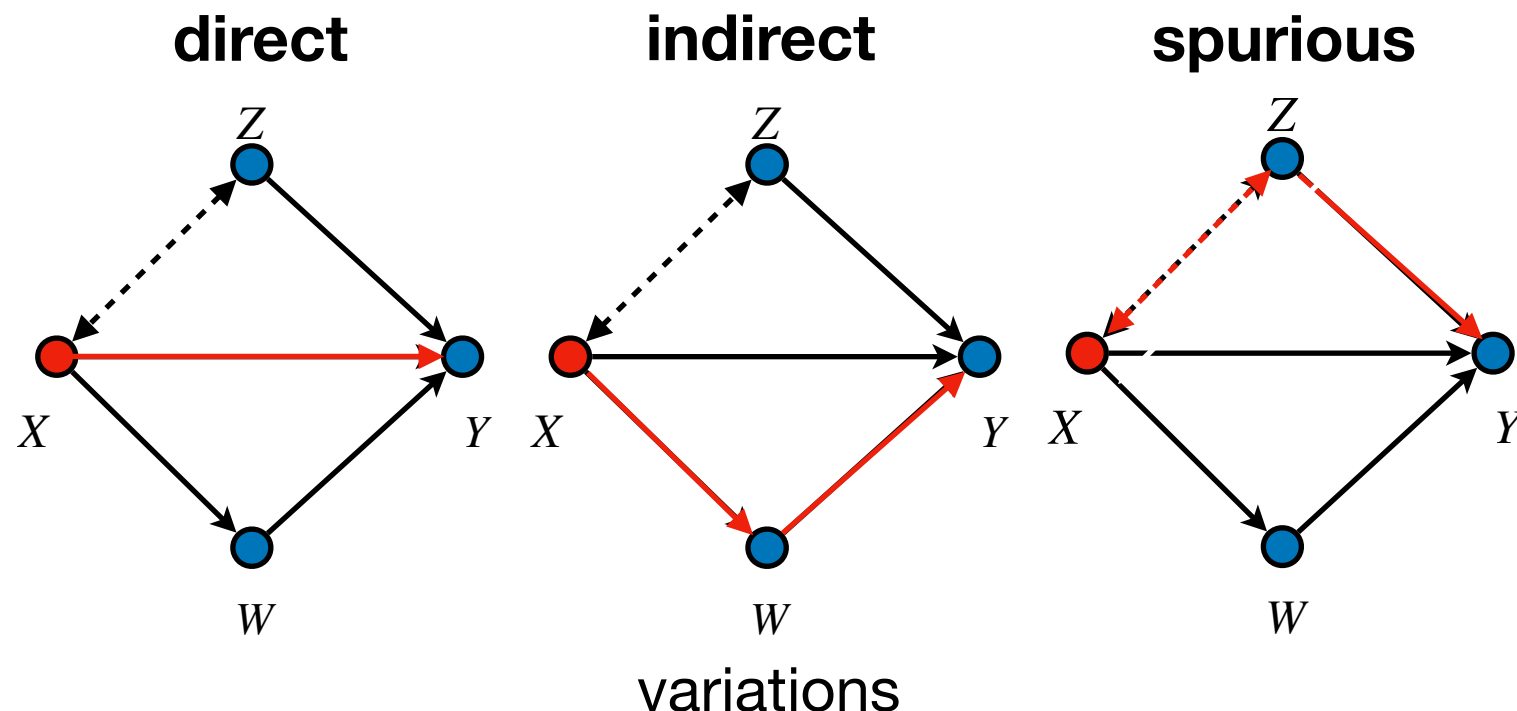
On the other, we need to “ground” (or attribute) the variations to different legal doctrines”

Disparate Treatment

Disparate Impact

Business Necessity

But, we know that TV contains



⇒ This entanglement makes the attribution problem challenging!

Admiss

Note: Power and Admissibility are the analogues of necessity and sufficiency for the corresponding fairness measures.

Definition. Let Ω be a criterion Q and measures

- The measure μ is said to be admissible w.r.t Q if

$$\forall \mathcal{M} \in \Omega : Q(\mathcal{M}) = 0 \implies \mu(\mathcal{M}) = 0.$$

- The measure μ' is said to be more powerful than μ if

(i) μ' is admissible

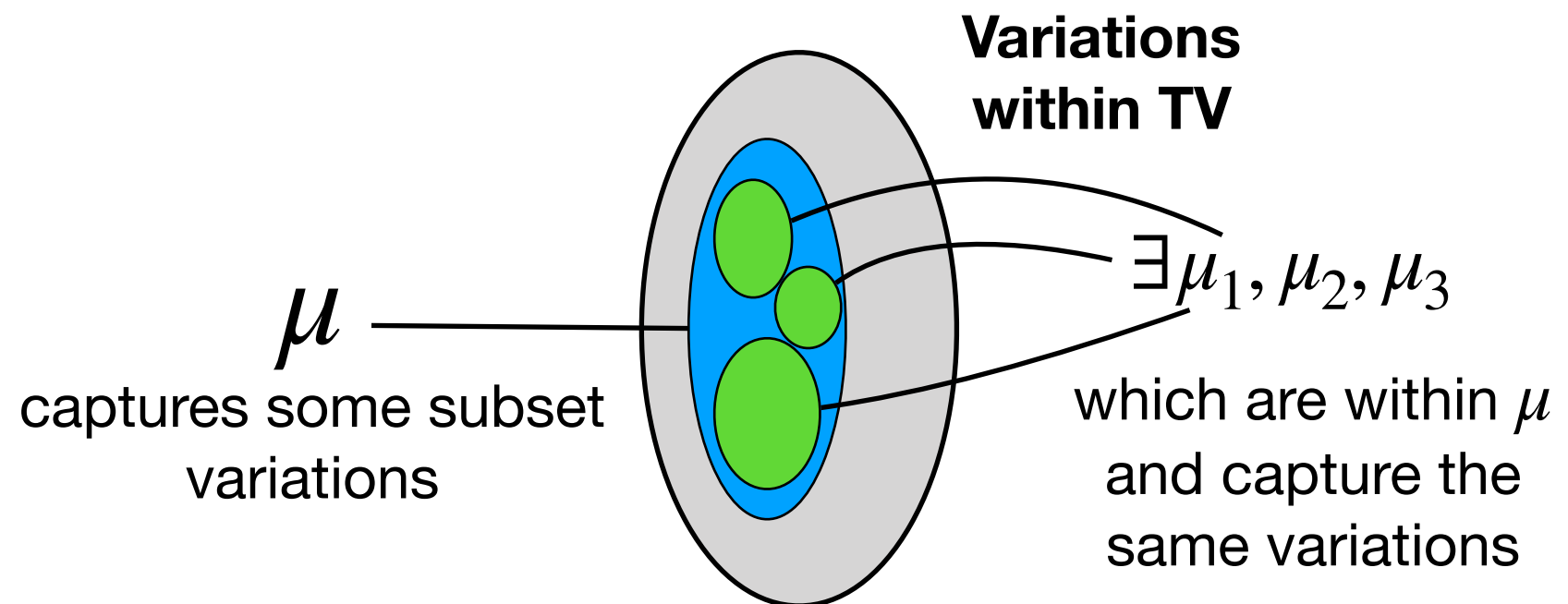
(ii) $\mu'(\mathcal{M}) = 0 \implies \mu(\mathcal{M}) = 0.$

Decomposability

Definition. Let Ω be a class of SCMs and μ be a measure defined over it. μ is said to be Ω -decomposable if there exist measures

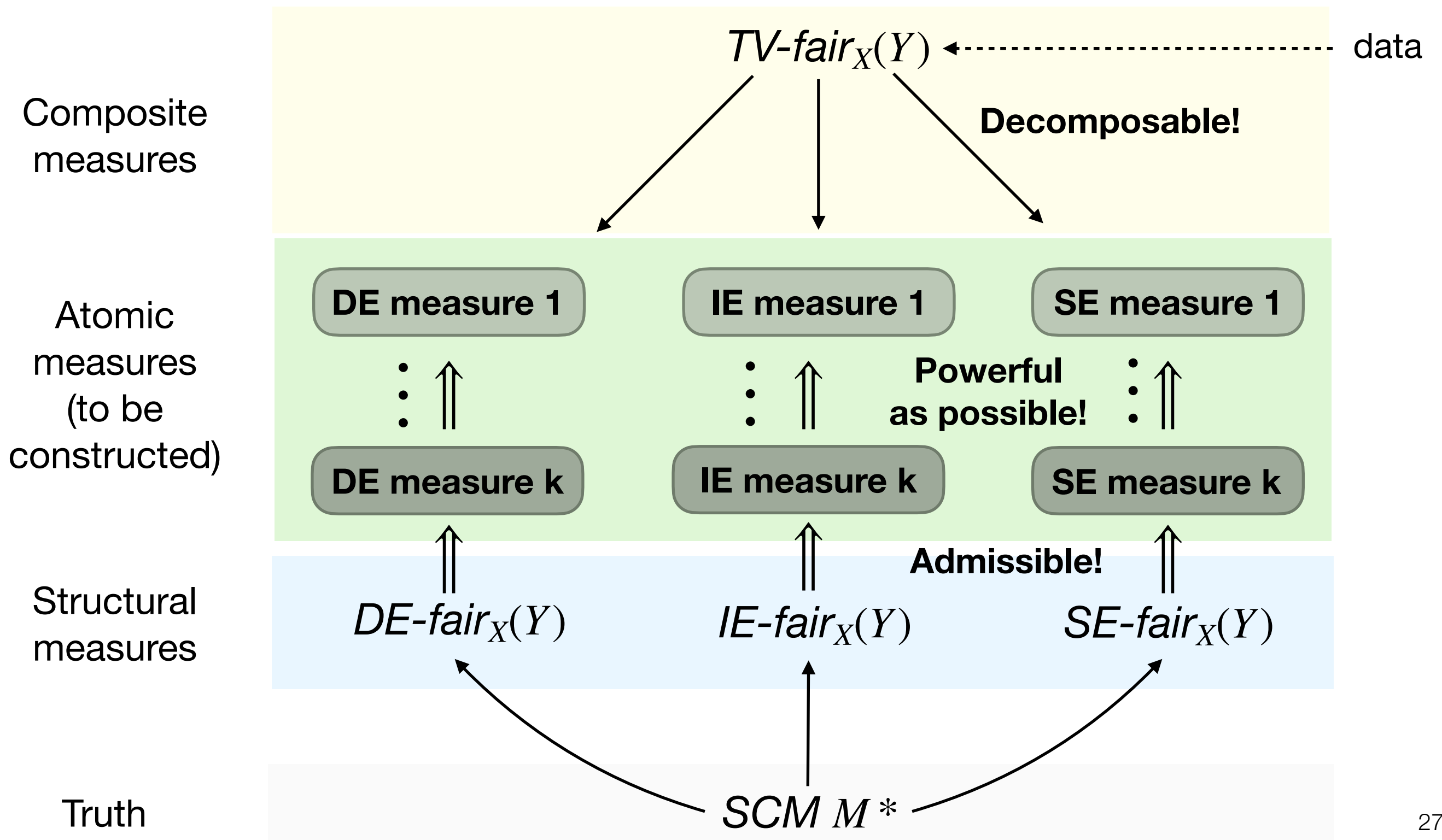
$$\mu_1, \dots, \mu_k \text{ such that } \mu = f(\mu_1, \dots, \mu_k),$$

and where f is a non-trivial function vanishing at the origin, i.e., $f(0, \dots, 0) = 0$.



Note: Decomposability can imply lack of admissibility.

Admissibility, Power, Decomposability - Summary



Fundamental Problem of Causal Fairness Analysis (FPCFA)

Definition. Let μ be a fairness measure defined over a space of SCMs Ω . Let Q_1, \dots, Q_k be a collection of structural fairness criteria. The Fundamental Problem of Causal Fairness Analysis is to find a collection of measures μ_1, \dots, μ_k s.t. the following properties hold:

(i) μ is *decomposable* w.r.t. μ_1, \dots, μ_k **Decomposability**

(ii) μ_1, \dots, μ_k are *admissible* w.r.t. the structural fairness criteria Q_1, Q_2, \dots, Q_k **Admissibility**

(iii) μ_1, \dots, μ_k are as *powerful* as possible. **Power**

How to solve the FPCFA?

Section 3.1
Definition 13

The Anatomy of Contrastive Measures

Definition. A contrast is any quantity of the form

$$P(y_{C_1} | E_1) - P(y_{C_0} | E_0).$$

Section 3.2

where E_0, E_1 are observed (factual) events and C_0, C_1 are counterfactual events to which the outcome Y responds.

A contrast compares the outcome Y of individuals

who coincide with the observed event E_1 versus E_0 , in the factual world,

and whose values, possibly counterfactually, were intervened on following C_1 versus C_0 .

Contrastive Measures: Factual vs. Counterfactual Basis

Theorem. Any contrast $P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_0)$ can be decomposed into its factual and counterfactual components:

$$\underbrace{P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_1)}_{\text{counterfactual contrast}} + \underbrace{P(y_{C_0} \mid E_1) - P(y_{C_0} \mid E_0)}_{\text{factual contrast}}.$$

We
normally
think of
 C_0, C_1, E_0, E_1 as
including X .

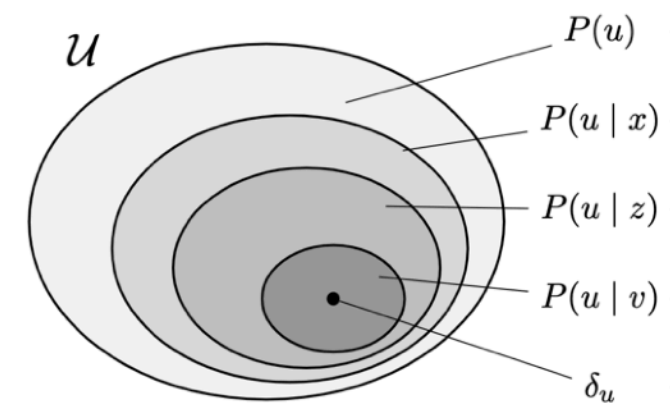
difference arising from
counterfactuals C_0, C_1
used to capture the causal
influence of X on Y .

difference arising from events
 E_0, E_1
used to capture non-causal
(spurious) influences of X on Y .

Structural Basis Expansion I

Theorem (continued). Whenever $E_0 = E_1 = e$, any counterfactual contrast $P(y_{C_1} | E = e) - P(y_{C_0} | E = e)$ admits the following structural basis expansion

$$\sum_u \underbrace{[y_{C_1}(u) - y_{C_0}(u)]}_{\text{unit-level difference}} \underbrace{P(u | E = e)}_{\text{posterior}}.$$

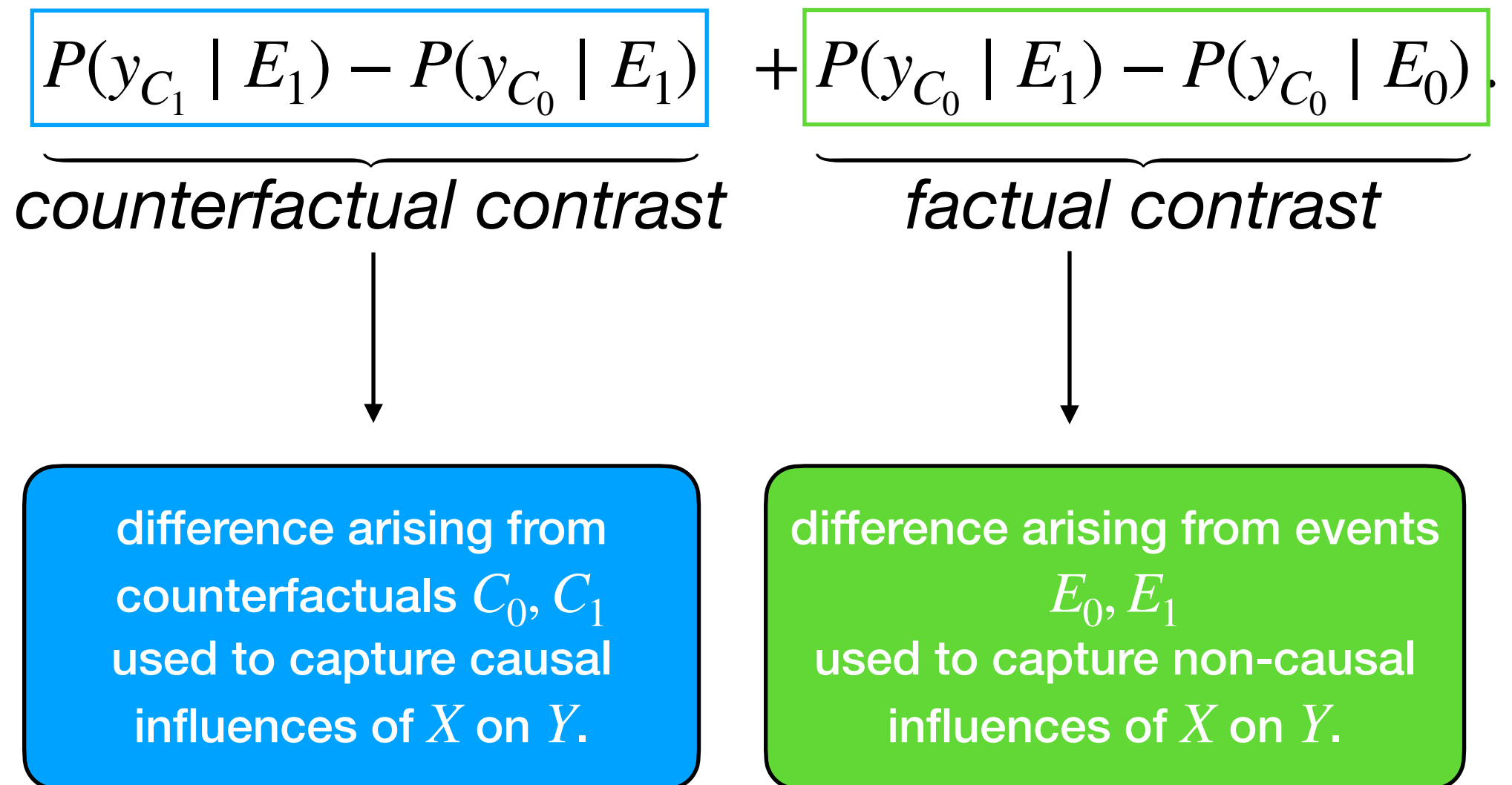


For a specific unit $U = u$,
Y's response to
the transition $C_0 \rightarrow C_1$.

Population of units
consistent with the
factual evidence $E=e$.

Contrastive Measures: Factual vs. Counterfactual Basis

Theorem. Any contrast $P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_0)$ can be decomposed into its factual and counterfactual components:



Structural Basis Expansion II

Theorem (continued). Whenever $C_0 = C_1 = c$, any factual contrast $P(y_c \mid E_1) - P(y_c \mid E_0)$ admits the following structural basis expansion:

$$\sum_u \underbrace{y_c(u)}_{\text{unit outcome}} \underbrace{[P(u \mid E_1) - P(u \mid E_0)]}_{\text{posterior difference}}.$$

Baseline outcome
for a fixed unit $U = u$.

Difference in posterior of how
likely unit $U = u$ is selected
under events E_0 vs. E_1 .

- We will be mostly interested in contrasts w/ $C = x$,
so that $X = x$ represents causal pathways.

Theorem (Contrasts & Structural Basis). Any contrast can be decomposed into its factual and counterfactuals components:

$$P(y_{C_1} | E_1) - P(y_{C_0} | E_0) = P(y_{C_1} | E_1) - P(y_{C_0} | E_1) + P(y_{C_0} | E_1) - P(y_{C_0} | E_0).$$

mechanisms \mathcal{F}

population $P(u)$

Furthermore:

A. Any contrast admits a structural basis expansion of the form:

Putting it all together...

unit-level difference in posterior

B. any factual contrast ($C_0 = C_1 = C$) admits the structural basis expansion of the form:

$$P(y_C | E_1) - P(y_C | E_0) = \sum_u \underbrace{y_C(u)}_{\text{unit outcome}} \underbrace{[P(u | E_1) - P(u | E_0)]}_{\text{posterior difference}}.$$

Theorem (Contrasts & Structural Basis). Any contrast can be decomposed into its factual and counterfactuals components:

$$P(y_{C_1} | E_1) - P(y_{C_0} | E_0) = P(y_{C_1} | E_1) - P(y_{C_0} | E_1) + P(y_{C_0} | E_1) - P(y_{C_0} | E_0).$$

mechanisms \mathcal{F}

population $P(u)$

Furthermore:

- A. Any counterfactual contrast ($E_0 = E_1 = E$) admits the structural basis expansion of the form:

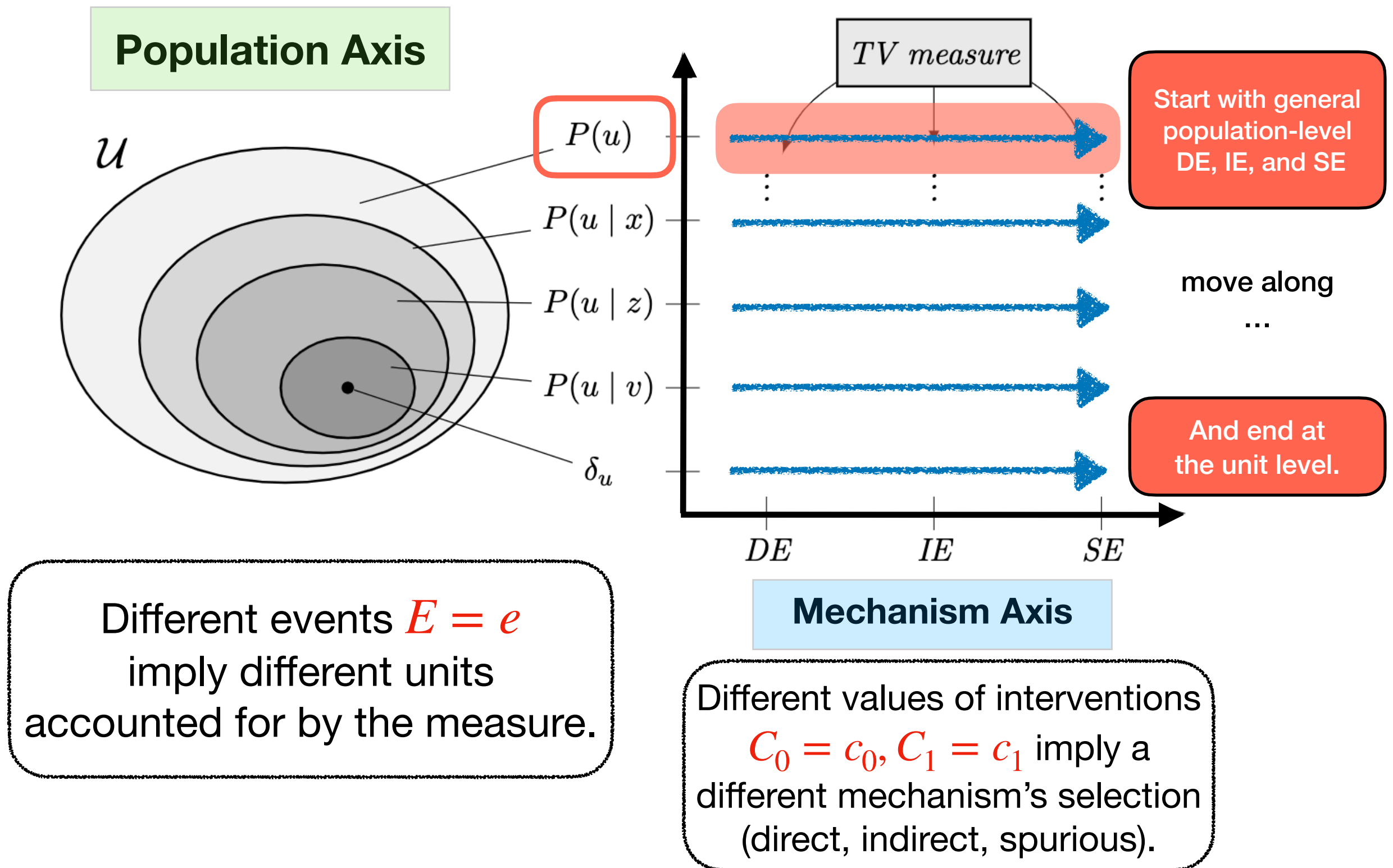
$$P(y_{C_1} | E) - P(y_{C_0} | E) = \sum_u \underbrace{[y_{C_1}(u) - y_{C_0}(u)]}_{\text{unit-level difference}} \underbrace{P(u | E)}_{\text{posterior}}.$$

- B. any factual contrast ($C_0 = C_1 = C$) admits the structural basis expansion of the form:

$$P(y_C | E_1) - P(y_C | E_0) = \sum_u \underbrace{y_C(u)}_{\text{unit outcome}} \underbrace{[P(u | E_1) - P(u | E_0)]}_{\text{posterior difference}}.$$

Explainability Plane

Section 3.2
Figure 7

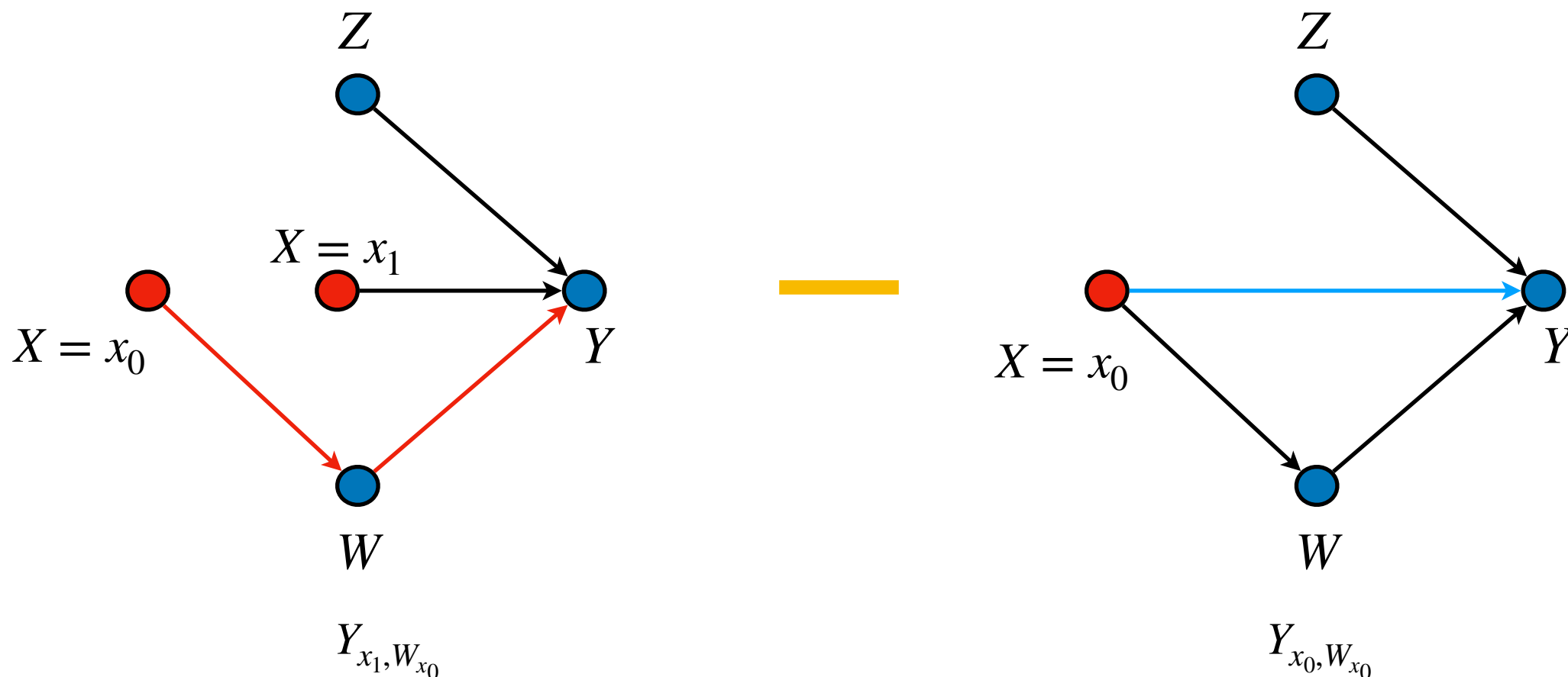


TV family of causal fairness measures

Gedankenexperiment (NDE)

- For a male employee ($X = x_0$), how would his salary (Y) change **had he been** a female ($X = x_1$), while keeping the age, nationality, education, employment status unchanged (i.e., at the natural level $X = x_0$)?

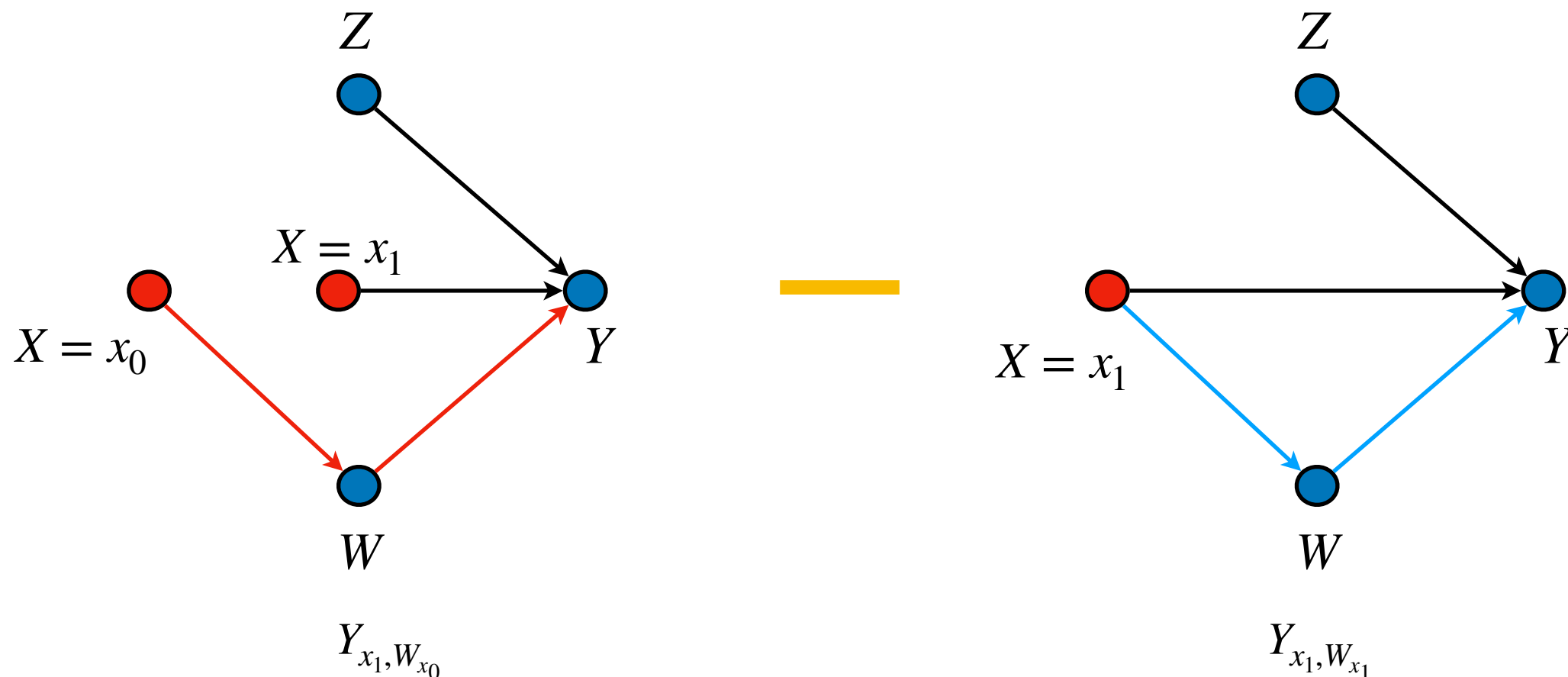
$$\mathbf{NDE}_{x_0, x_1}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_0, W_{x_0}})$$



Gedankenexperiment (NIE)

- For a female employee ($X = x_1$), how would her salary (Y) change **had she been** a male ($X = x_0$), while keeping gender unchanged along the direct causal pathway (at the natural level $X = x_1$)?

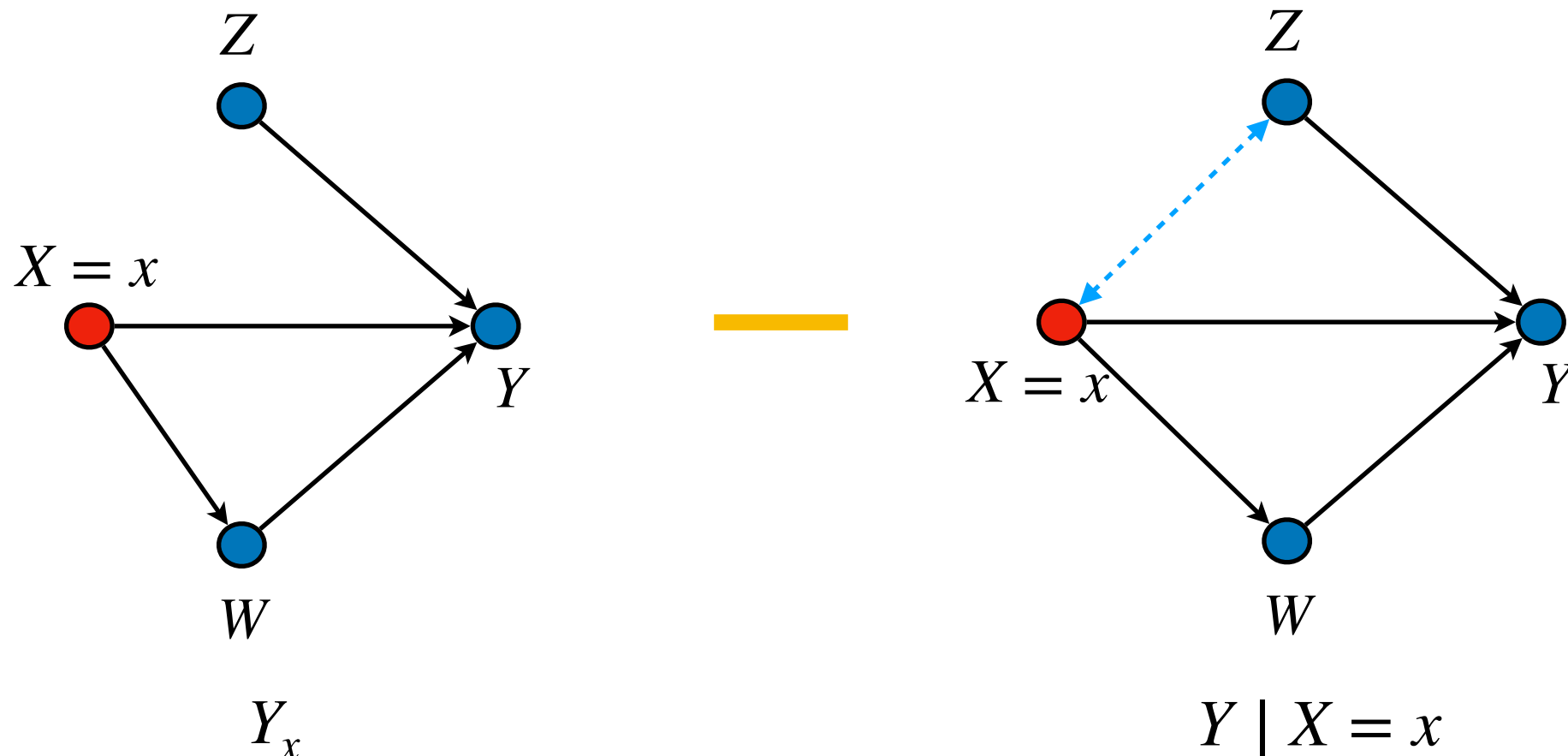
$$\mathbf{NIE}_{x_1, x_0}(y) = P(y_{x_1, W_{x_0}}) - P(y_{x_1, W_{x_1}})$$

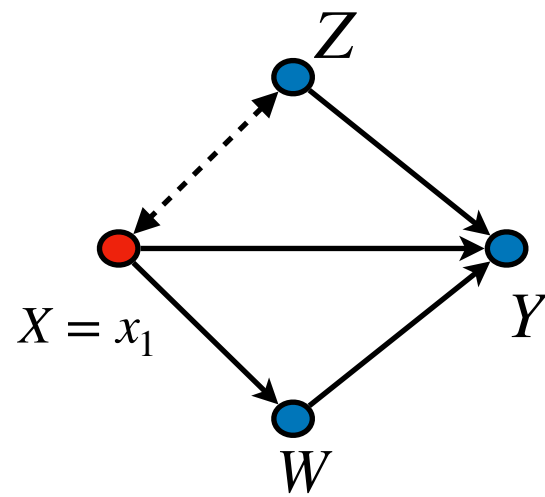


Gedankenexperiment (Exp-SE)

- How would an individual's salary (Y) change if their gender is set to male (or female) by intervention, compared to observing their salary as male (female)?

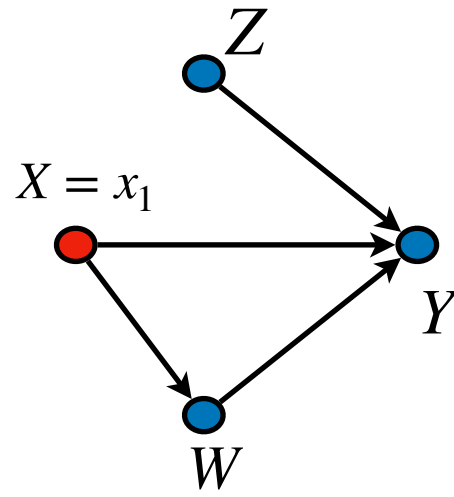
$$\text{Exp-SE}_x(y) = P(y_x) - P(y \mid x)$$





$Y \mid X = x_1$

$-\text{NIE}_{x_1, x_0}(y)$



Y_{x_1}

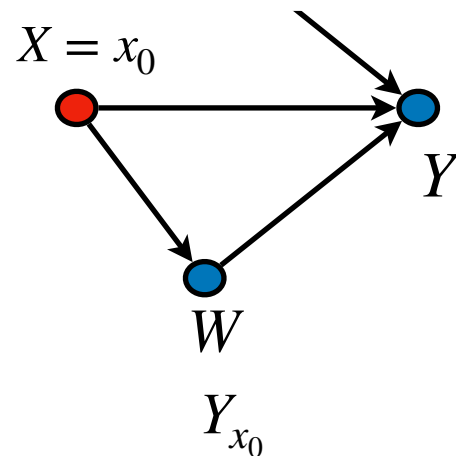
$-\text{Exp-SE}_{x_1}(y)$



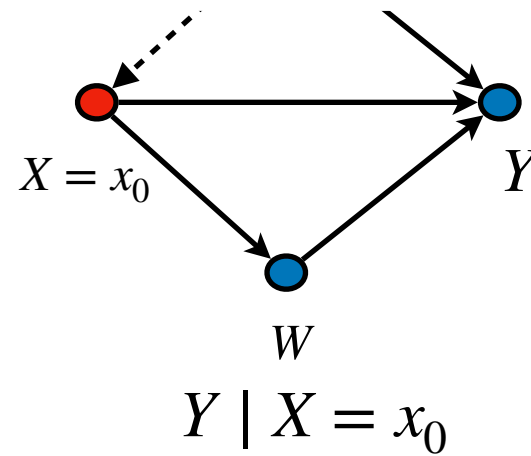
Y

$\text{NDE}_{x_0, x_1}(y)$

TV Decomposition I

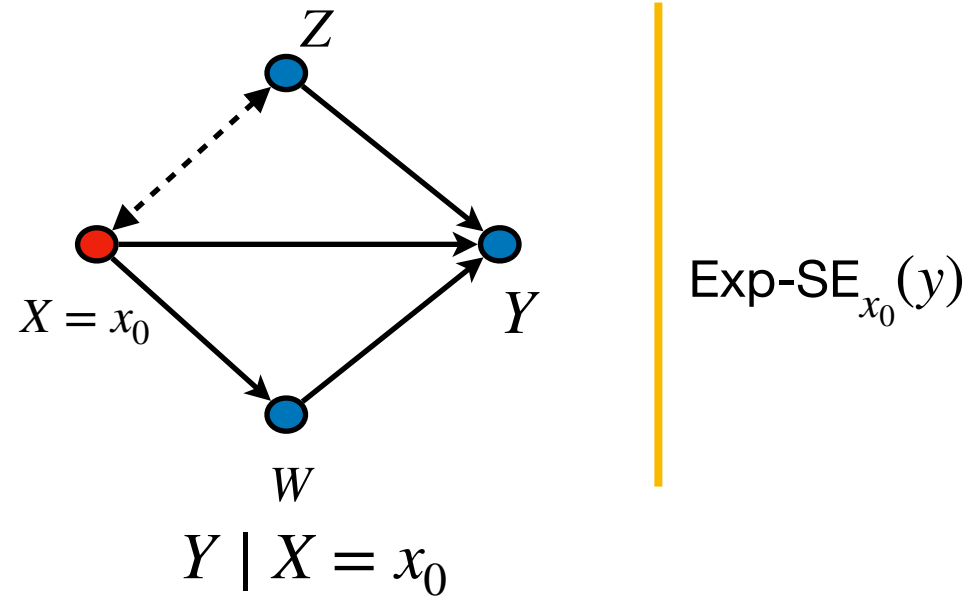
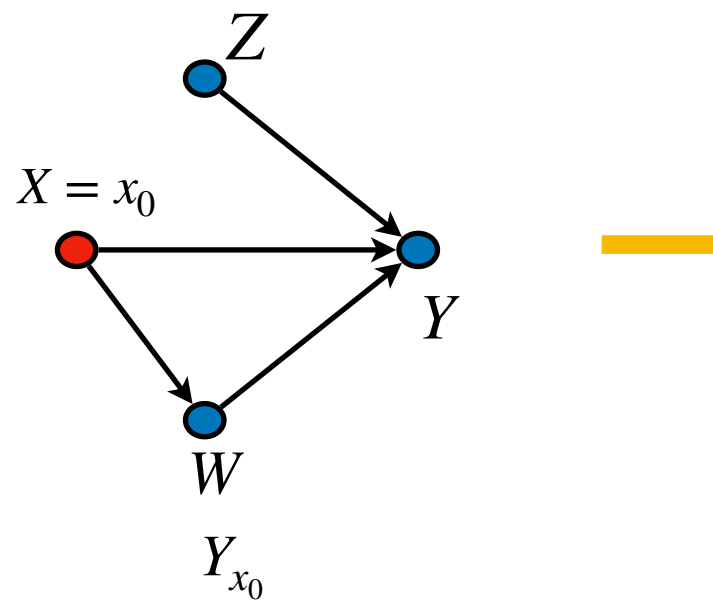
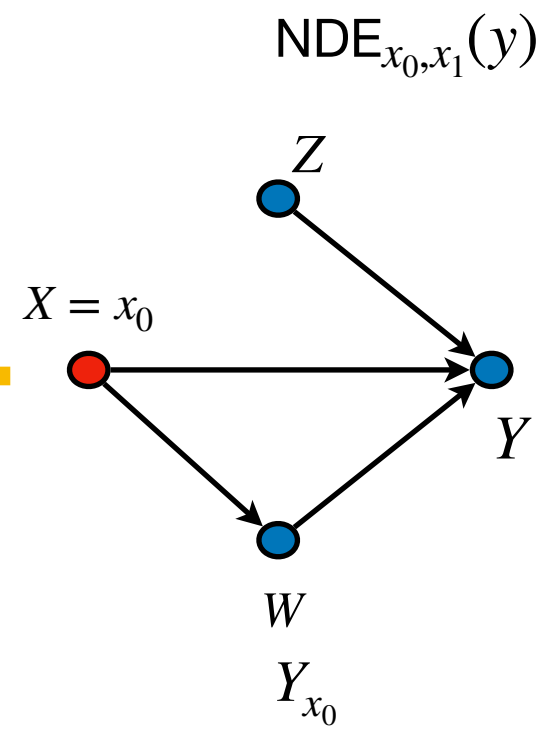
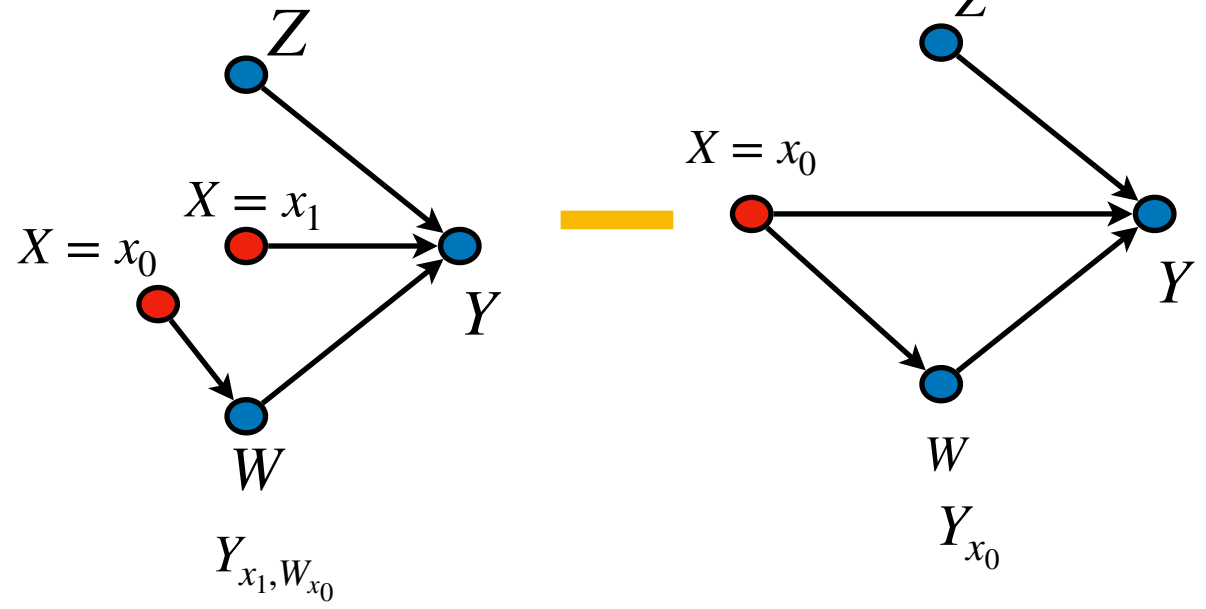
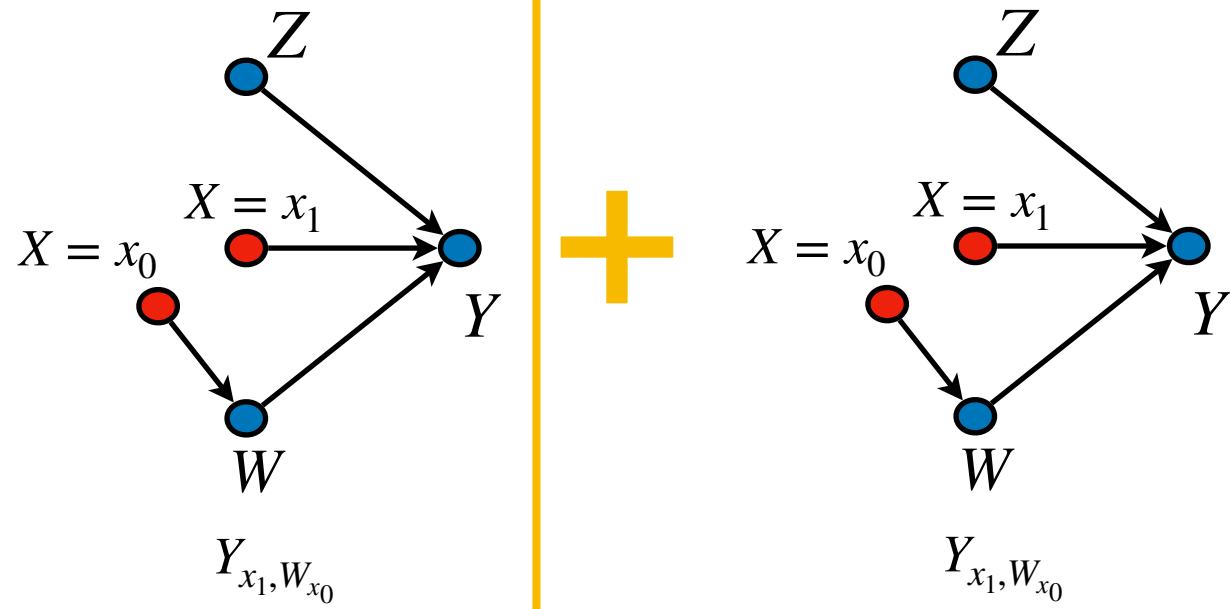
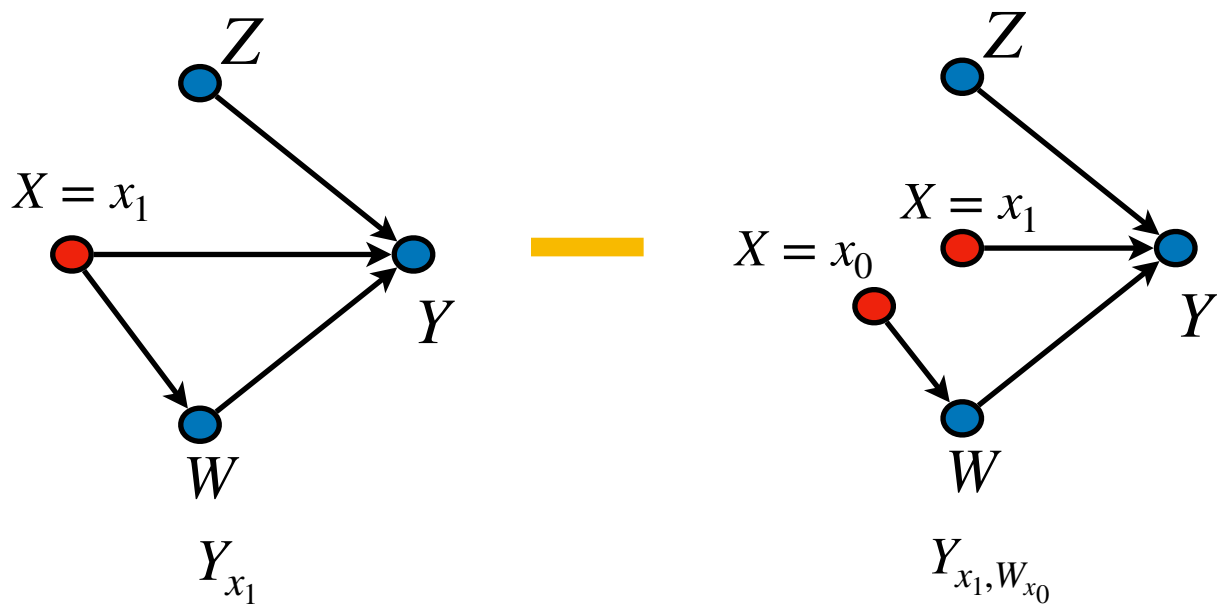
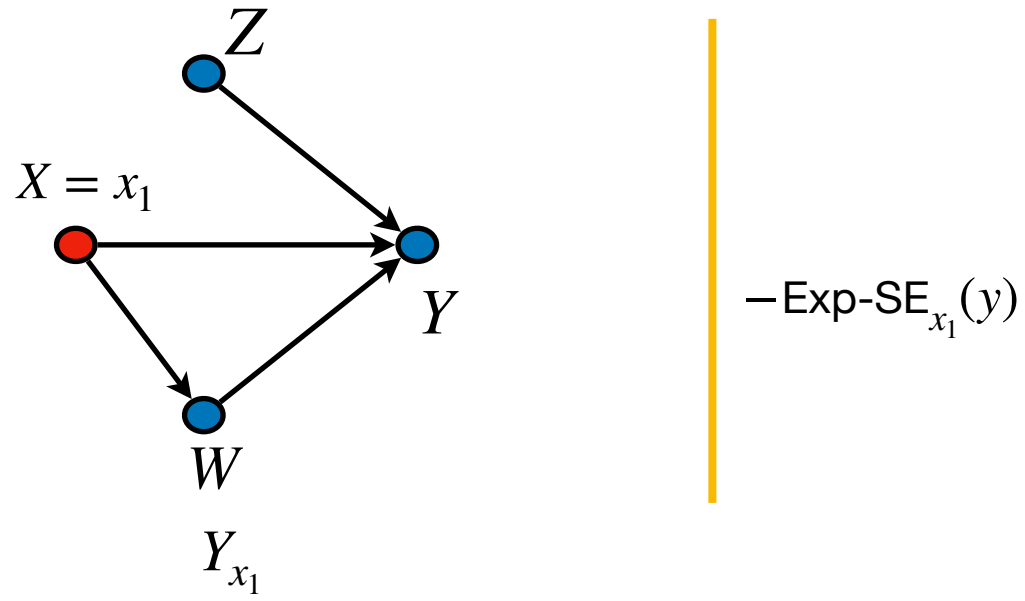
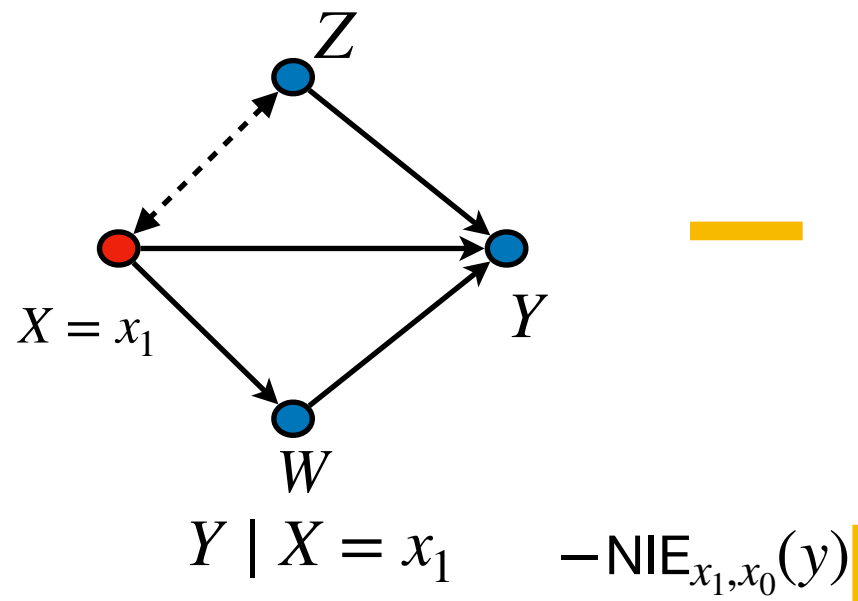


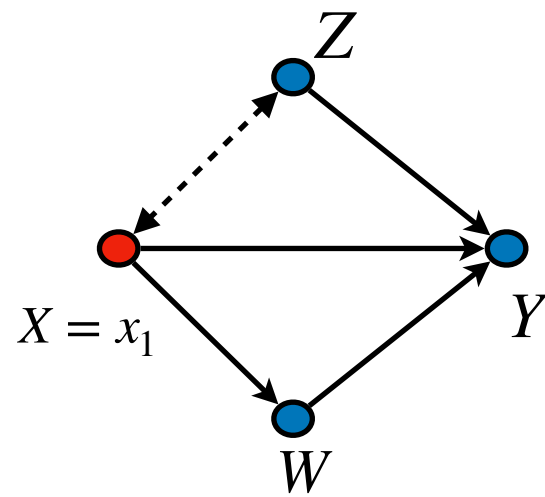
Y_{x_0}



$Y \mid X = x_0$

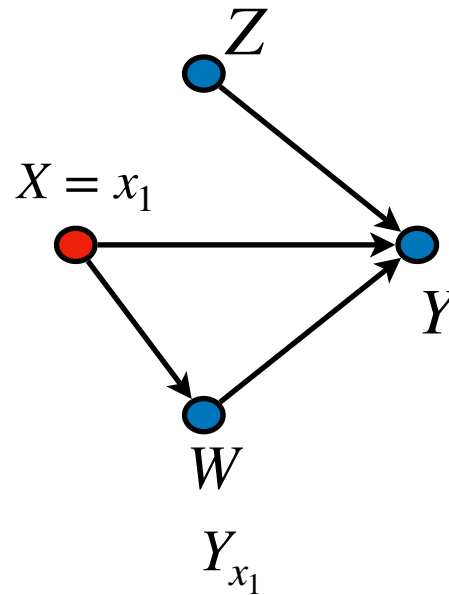
$\text{Exp-SE}_{x_0}(y)$





$Y \mid X = x_1$

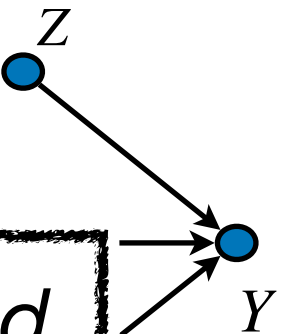
$-NIE_{x_1, x_0}(y)$



$-Exp-SE_{x_1}(y)$

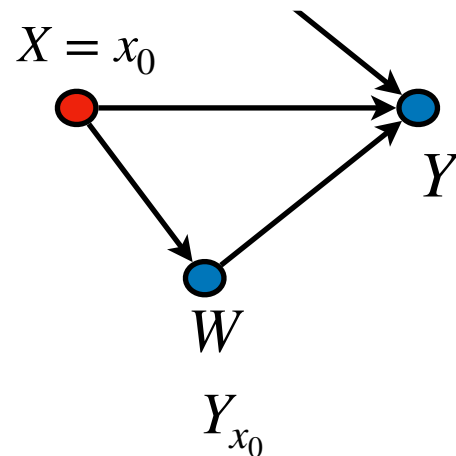


$NDE_{x_0, x_1}(y)$

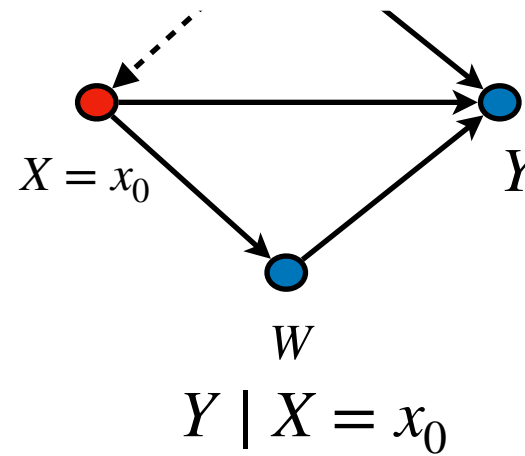


Lemma. The total variation measure can be decomposed into its direct, indirect, and spurious variations:

$$TV_{x_0, x_1}(y) = NDE_{x_0, x_1}(y) - NIE_{x_1, x_0}(y) - (Exp-SE_{x_1}(y) - Exp-SE_{x_0}(y)).$$



Y_{x_0}



$Y \mid X = x_0$

$Exp-SE_{x_0}(y)$

Relation to Structural Fairness

Corollary. *The criteria based on NDE, NIE, and Exp-SE measures are **admissible** with respect to structural direct, indirect, and spurious fairness. Formally, these facts are written as:*

$$S-DE \implies NDE\text{-fair}$$

$$S-IE \implies NIE\text{-fair}$$

$$S-SE \implies \text{Exp-SE-fair}$$

admissibility w.r.t.
structural

In practice, for example, by computing the NDE, we can test for the presence of structural direct effect.

Testing Structural Fairness in Practice

- Our previous corollary shows that
$$\text{S-DE} \implies \text{NDE-fair}.$$
- By taking this statement's contrapositive, we can see that
$$\text{NDE}_{x_0, x_1}(y) \neq 0 \implies \neg \text{S-DE}.$$
- Therefore, in practice, one may use the following hypothesis testing procedure for testing structural direct effect,
$$H_0 : \text{NDE}_{x_0, x_1}(y) = 0.$$

A similar approach can be used for the NIE and Exp-SE since

$$\text{S-IE} \implies \text{NIE-fair}$$

$$\text{S-SE} \implies \text{Exp-SE-fair}$$

This will be used to connect with the disparate treatment and impact doctrines later on.

Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

SCM M^*
(unobserved)

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$\begin{aligned} \text{NDE}_{x_0, x_1}(y) &= P(y_{x_1, W_{x_0}}) - P(y_{x_0}) \\ &= P(\text{Bernoulli}\left(\frac{1}{5}(1 - Z) + \frac{1}{6}W\right) = 1) \\ &\quad - P(\text{Bernoulli}\left(\frac{1}{5}Z + \frac{1}{6}W\right) = 1) \\ &= \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} P(w) \left[\frac{1}{5}(1 - 2z) + \frac{1}{6}w - \frac{1}{6}w \right] \\ &= \sum_{z \in \{0,1\}} \frac{1}{5}(1 - 2z) = 0. \end{aligned}$$

Section 4
Example 9

Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

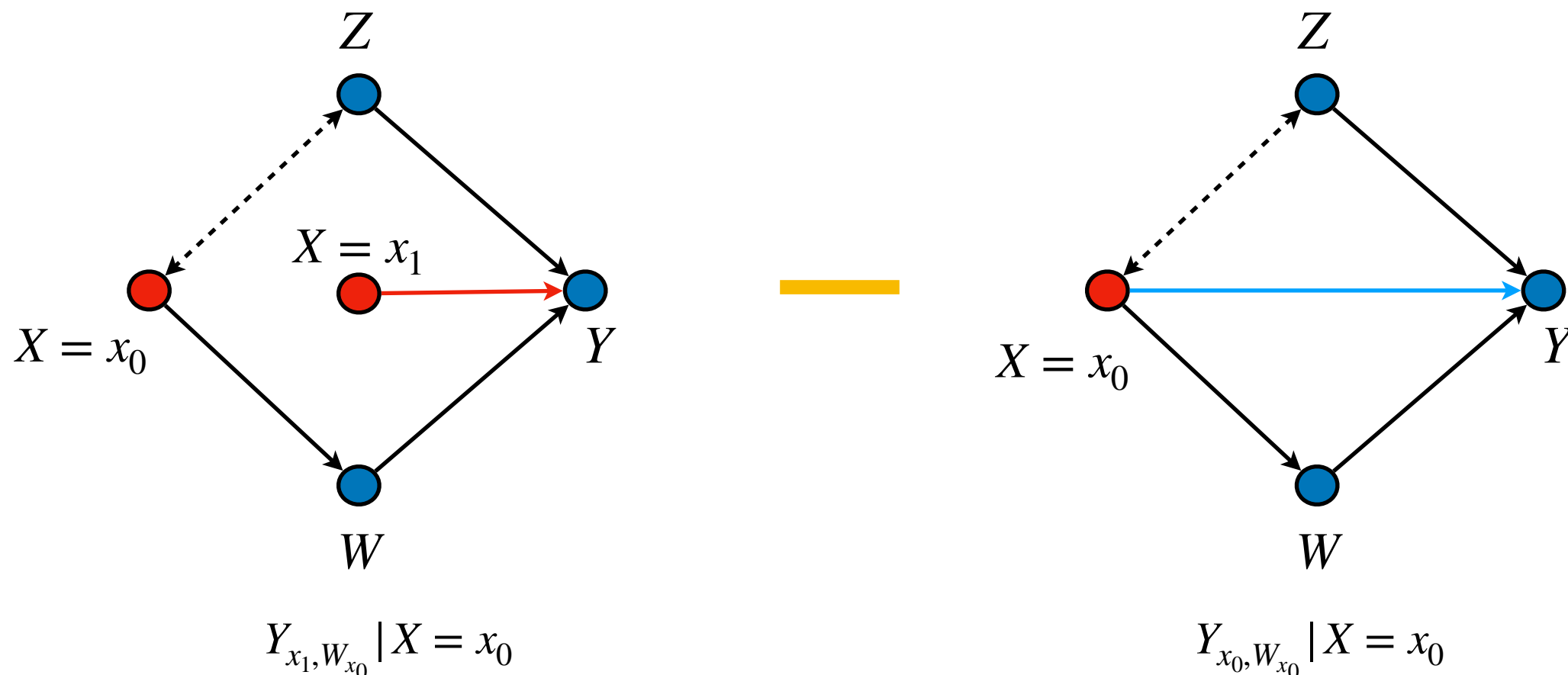
**NDE is admissible w.r.t. S-DE.
However, here $NDE = 0$,
but structural direct effect exists.**

**Q: Is NDE powerful enough for
detecting direct discrimination?**

Gedankenexperiment (Ctf-DE)

- For a male employee $X = x_0$, how would his salary change (Y) **had he been** a female ($X = x_1$), while keeping the age, nationality, education and employment status unchanged (at the level of

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



Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

SCM M

$U \leftarrow N(0,1)$

$X \leftarrow \text{Bernoulli}(\text{expit}(U))$

$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$

$W \leftarrow \text{Bernoulli}(0.3)$

$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$

$$\begin{aligned}
 \text{Ctf-DE}_{x_0, x_1}(y \mid x_0) &= P(y_{x_1, W_{x_0}} \mid x_0) - P(y_{x_0} \mid x_0) \\
 &= P(\text{Bernoulli}\left(\frac{1}{5}(1 - Z) + \frac{1}{6}W\right) = 1 \mid x_0) \\
 &\quad - P(\text{Bernoulli}\left(\frac{1}{5}Z + \frac{1}{6}W\right) = 1 \mid x_0) \\
 &= \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} P(w)P(z \mid x_0) \left[\frac{1}{5}(1 - 2z) + \frac{1}{6}w - \frac{1}{6}w \right] \\
 &= \sum_{z \in \{0,1\}} \frac{1}{5}(1 - 2z)P(z \mid x_0) = 0.036.
 \end{aligned}$$

Section 4 Example 10

Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

Key properties of Ctf-DE:

1. Ctf-DE is admissible.
2. Ctf-DE is more powerful than NDE.

$U \leftarrow N(0,1)$

$X \leftarrow \text{Bernoulli}(U)$

$Z \leftarrow \text{Bernoulli}(U)$

$W \leftarrow \text{Bernoulli}(0.5)$

$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$

$$\begin{aligned}
 & \sum_{z \in \{0,1\}} \sum_{w \in \{0,1\}} \frac{1}{5} (1 - 2z) P(z | x_0) + \frac{1}{6} w - \frac{1}{6} w] \\
 &= \sum_{z \in \{0,1\}} \frac{1}{5} (1 - 2z) P(z | x_0) = 0.036.
 \end{aligned}$$

Section 4
Example 10

x -specific measures

Definition. The effect of treatment on the treated and counterfactual direct, indirect, and spurious effects are defined as

$$ETT_{x_0, x_1}(y \mid x) = P(y_{x_1} \mid x) - P(y_{x_0} \mid x)$$

$$Ctf-DE_{x_0, x_1}(y \mid x) = P(y_{x_1, W_{x_0}} \mid x) - P(y_{x_0} \mid x)$$

$$Ctf-IE_{x_1, x_0}(y \mid x) = P(y_{x_1, W_{x_0}} \mid x) - P(y_{x_1} \mid x)$$

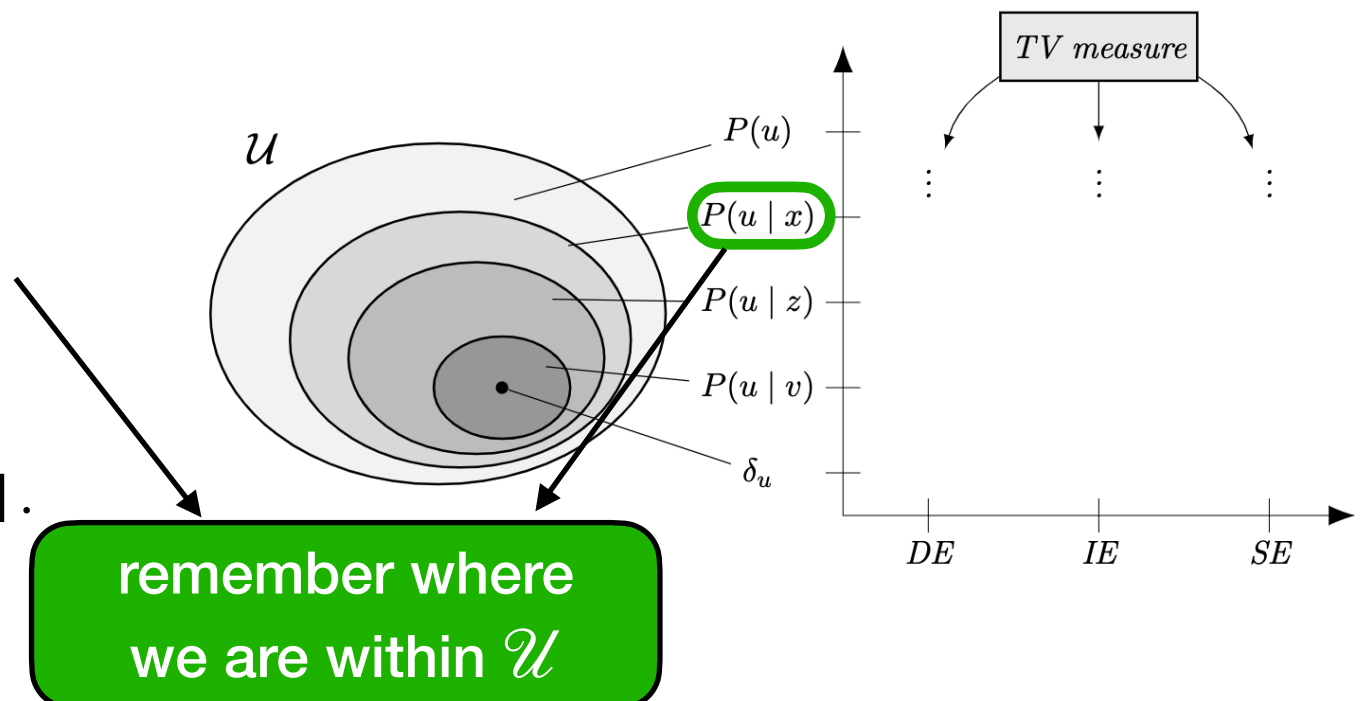
$$Ctf-SE_{x_0, x_1}(y) = P(y_{x_0} \mid x_1) - P(y_{x_0} \mid x_0).$$

Structural Basis Expansion:

$$Ctf-DE_{x_0, x_1}(y \mid x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u \mid x)$$

$$Ctf-IE_{x_1, x_0}(y \mid x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u \mid x)$$

$$Ctf-SE_{x_0, x_1}(y) = \sum_u y_{x_0}(u) [P(u \mid x_1) - P(u \mid x_0)].$$



x -specific

Definition. The effect of treatment on direct, indirect, and spurious effects are

$$\begin{aligned} \text{TE}_{x_0, x_1}(y \mid x) &= P(y_{x_1}) - P(y_{x_0}) \\ \text{NDE}_{x_0, x_1}(y) &= P(y_{x_1, W_{x_0}}) - P(y_{x_0}) \\ \text{NIE}_{x_1, x_0}(y) &= P(y_{x_1, W_{x_0}}) - P(y_{x_1}) \\ \text{Exp-SE}_{x_0, x_1}(y) &= P(y_x) - P(y_x \mid x). \end{aligned}$$

where we came from

$$\text{ETT}_{x_0, x_1}(y \mid x) = P(y_{x_1} \mid x) - P(y_{x_0} \mid x)$$

$$\text{Ctf-DE}_{x_0, x_1}(y \mid x) = P(y_{x_1, W_{x_0}} \mid x) - P(y_{x_0} \mid x)$$

$$\text{Ctf-IE}_{x_1, x_0}(y \mid x) = P(y_{x_1, W_{x_0}} \mid x) - P(y_{x_1} \mid x)$$

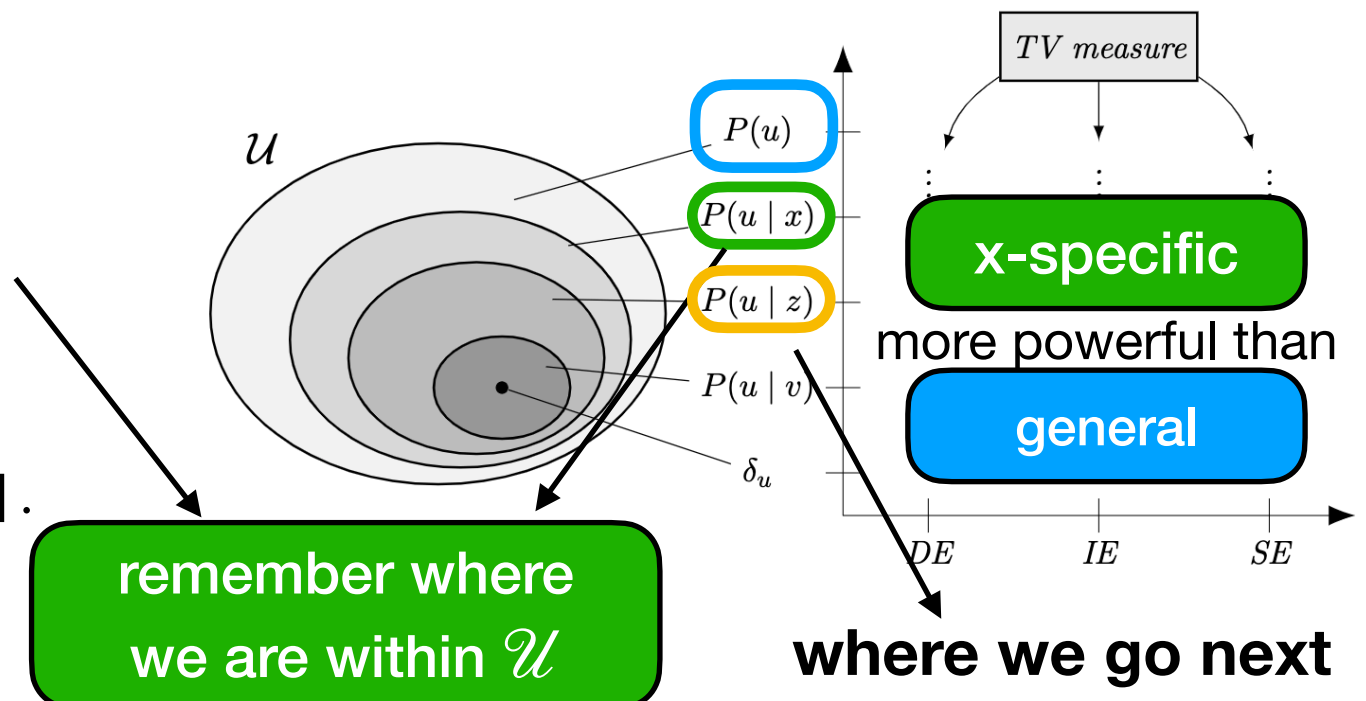
$$\text{Ctf-SE}_{x_0, x_1}(y) = P(y_{x_0} \mid x_1) - P(y_{x_0} \mid x_0).$$

Structural Basis Expansion:

$$\text{Ctf-DE}_{x_0, x_1}(y \mid x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u \mid x)$$

$$\text{Ctf-IE}_{x_1, x_0}(y \mid x) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u \mid x)$$

$$\text{Ctf-SE}_{x_0, x_1}(y) = \sum_u y_{x_0}(u) [P(u \mid x_1) - P(u \mid x_0)].$$



z-specific measures

Definition. The z -specific total, direct, and indirect effects are defined as

$$z\text{-}TE_{x_0, x_1}(y \mid z) = P(y_{x_1} \mid z) - P(y_{x_0} \mid z)$$

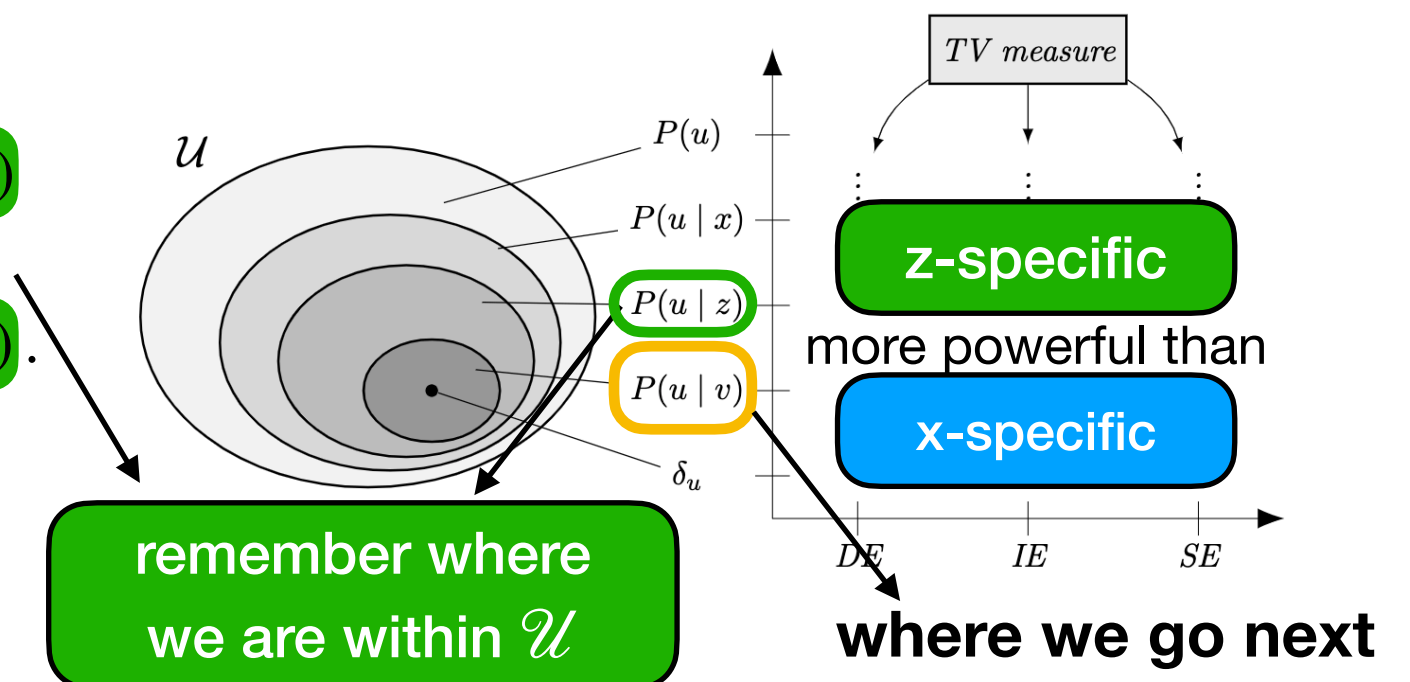
$$z\text{-}DE_{x_0, x_1}(y \mid z) = P(y_{x_1, W_{x_0}} \mid z) - P(y_{x_0} \mid z)$$

$$z\text{-}IE_{x_1, x_0}(y \mid z) = P(y_{x_1, W_{x_0}} \mid z) - P(y_{x_1} \mid z).$$

Structural Basis Expansion:

$$z\text{-}DE_{x_0, x_1}(y \mid z) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u \mid z)$$

$$z\text{-}IE_{x_1, x_0}(y \mid z) = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u \mid z).$$



Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

SCM M

$$U \leftarrow N(0,1)$$

$$X \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$Z \leftarrow \text{Bernoulli}(\text{expit}(U))$$

$$W \leftarrow \text{Bernoulli}(0.3)$$

$$Y \leftarrow \text{Bernoulli}\left(\frac{1}{5}(X + Z - 2XZ) + \frac{1}{6}W\right)$$

$$\begin{aligned} z\text{-DE}(y \mid Z = 0) &= P(y_{x_1, w_{x_0}} \mid Z = 0) - P(y_{x_0} \mid Z = 0) \\ &= P(\text{Bernoulli}(\frac{1}{5}(1 - Z) + \frac{1}{6}W) = 1 \mid Z = 0) \\ &\quad - P(\text{Bernoulli}(\frac{1}{5}(Z) + \frac{1}{6}W) = 1 \mid Z = 0) \\ &= \sum_{w \in \{0,1\}} P(w) \left[\frac{1}{5} + \frac{1}{6}w - \frac{1}{6}w \right] = \frac{1}{5}. \end{aligned}$$

Section 4 Example 11

Example (Limitation of NDE). A new startup company is currently in hiring season. The hiring decision ($Y \in \{0,1\}$ indicating whether the candidate is hired) is based on gender ($X \in \{0,1\}$, female and male, respectively), age ($Z \in \{0,1\}$, younger and older than 40 years, respectively), and education level ($W \in \{0,1\}$ which indicates whether the applicant has a Ph.D. degree). Following the legal guidelines, the startup is in this case obliged to avoid disparate treatment in hiring.

Key properties of z -DE:

1. z -DE is admissible.
2. z -DE is more powerful than Ctf-DE.

$U \leftarrow N(0,1)$
 $X \leftarrow \text{Bernoulli}(\frac{1}{2})$
 $Z \leftarrow \text{Bernoulli}(\frac{1}{2})$
 $W \leftarrow \text{Bernoulli}(\frac{1}{2})$
 $Y \leftarrow \text{Bernoulli}(\frac{1}{5}(X + Z - ZAZ) + \frac{1}{6}W)$

$w \in \{0,1\}$

$1 | Z = 0)$

$Z = 0)$

$\frac{1}{5}$

Section 4
Example 11

v' -specific measures

Definition. The v' -specific total, direct, and indirect effects are defined as

$$v'\text{-}TE_{x_0, x_1}(y \mid v') = P(y_{x_1} \mid v') - P(y_{x_0} \mid v')$$

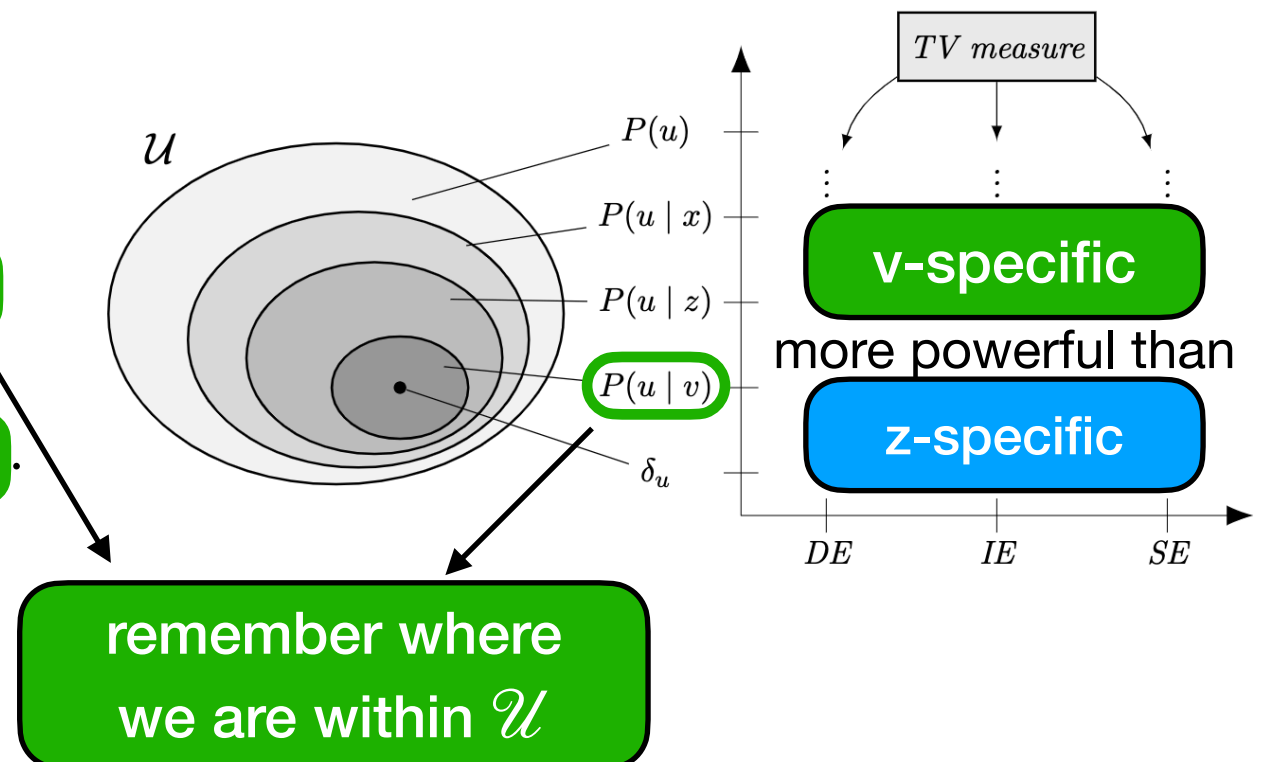
$$v'\text{-}DE_{x_0, x_1}(y \mid v') = P(y_{x_1, W_{x_0}} \mid v') - P(y_{x_0} \mid v')$$

$$v'\text{-}IE_{x_1, x_0}(y \mid v') = P(y_{x_1, W_{x_0}} \mid v') - P(y_{x_1} \mid v').$$

Structural Basis Expansion:

$$v'\text{-}DE_{x_0, x_1}(y \mid v') = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_0}(u)] P(u \mid v')$$

$$v'\text{-}IE_{x_1, x_0}(y \mid v') = \sum_u [y_{x_1, W_{x_0}}(u) - y_{x_1}(u)] P(u \mid v').$$



Unit-level measures

Definition. Given a unit $U = u$, the unit-level total, direct, and indirect effects are given by

$$\text{unit-TE}_{x_0, x_1}(y(u)) = y_{x_1}(u) - y_{x_0}(u)$$

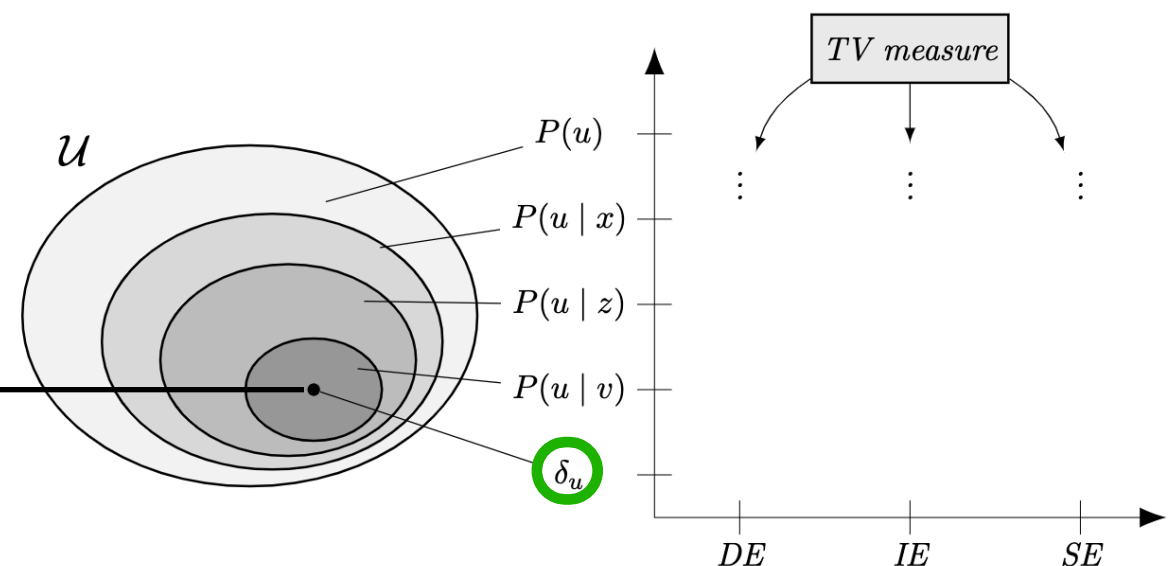
$$\text{unit-DE}_{x_0, x_1}(y(u)) = y_{x_1, W_{x_0}}(u) - y_{x_0}(u)$$

$$\text{unit-IE}_{x_1, x_0}(y(u)) = y_{x_1, W_{x_0}}(u) - y_{x_1}(u).$$

These quantities are
the structural basis.

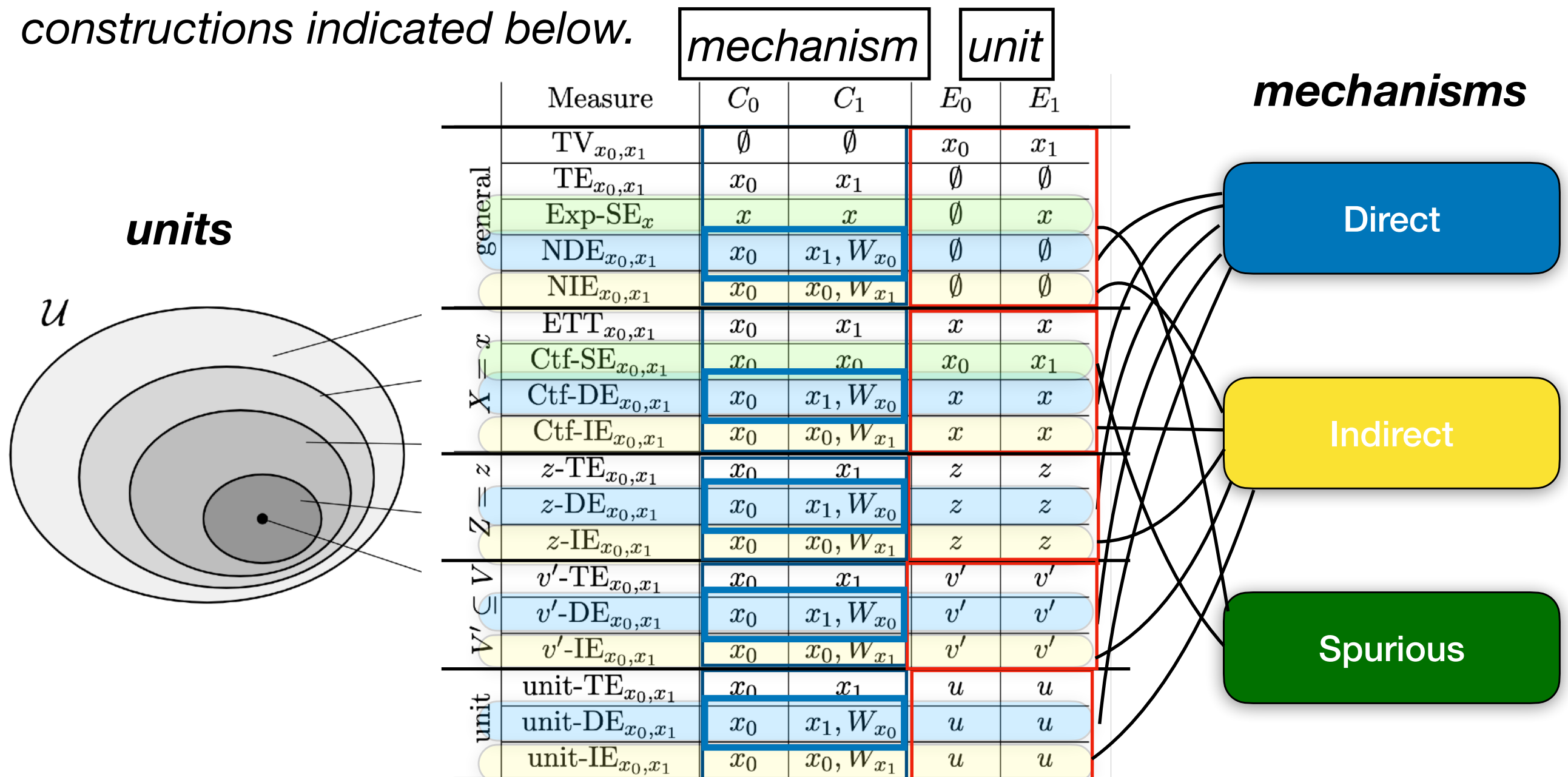
Remember where
we are within \mathcal{U} .

We reached the final,
unit-level measures!



TV family measures as contrasts

Lemma. Under the Standard fairness model, all the measures within the TV family can be written as contrasts $P(y_{C_1} | E_1) - P(y_{C_0} | E_0)$, following the constructions indicated below.



TV family measures as contrasts

Lemma. Under the Standard fairness model, all the measures within the TV family can be written as contrasts $P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_0)$, following the constructions indicated below.

		Measure	C_0	C_1	E_0	E_1	
general		TV_{x_0,x_1}	\emptyset	\emptyset	x_0	x_1	Direct
		TE_{x_0,x_1}	x_0	x_1	\emptyset	\emptyset	
		$Exp-SE_x$	x	x	\emptyset	x	
		NDE_{x_0,x_1}	x_0	x_1, W_{x_0}	\emptyset	\emptyset	
		NIE_{x_0,x_1}	x_0	x_0, W_{x_1}	\emptyset	\emptyset	
$X = x$		ETT_{x_0,x_1}	x_0	x_1	x	x	Indirect
		$Ctf-SE_{x_0,x_1}$	x_0	x_0	x_0	x_1	
		$Ctf-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	x	x	
		$Ctf-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	x	x	
$Z = z$		$z-TE_{x_0,x_1}$	x_0	x_1	z	z	
		$z-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	z	z	
		$z-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	z	z	
$V' \subseteq V$		$v'-TE_{x_0,x_1}$	x_0	x_1	v'	v'	Spurious
		$v'-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	v'	v'	
		$v'-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	v'	v'	
unit		$unit-TE_{x_0,x_1}$	x_0	x_1	u	u	
		$unit-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	u	u	
		$unit-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	u	u	

units

TV family measures as contrasts

Lemma. Under the Standard fairness model, all the measures within the TV family can be written as contrasts $P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_0)$, following the constructions indicated below.

units

	Measure	C_0	C_1	E_0	E_1
general	TV_{x_0,x_1}	\emptyset	\emptyset	x_0	x_1
	TE_{x_0,x_1}	x_0	x_1	\emptyset	\emptyset
	$Exp-SE_x$	x	x	\emptyset	x
	NDE_{x_0,x_1}	x_0	x_1, W_{x_0}	\emptyset	\emptyset
	NIE_{x_0,x_1}	x_0	x_0, W_{x_1}	\emptyset	\emptyset
$X = x$	ETT_{x_0,x_1}	x_0	x_1	x	x
	$Ctf-SE_{x_0,x_1}$	x_0	x_0	x_0	x_1
	$Ctf-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	x	x
	$Ctf-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	x	x
$Z = z$	$z-TE_{x_0,x_1}$	x_0	x_1	z	z
	$z-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	z	z
	$z-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	z	z
$V' \subseteq V$	$v'-TE_{x_0,x_1}$	x_0	x_1	v'	v'
	$v'-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	v'	v'
	$v'-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	v'	v'
unit	$unit-TE_{x_0,x_1}$	x_0	x_1	u	u
	$unit-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	u	u
	$unit-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	u	u

mechanisms

Direct

Causal

Spurious

TV family measures as contrasts

Lemma. Under the Standard fairness model, all the measures within the TV family can be written as contrasts $P(y_{C_1} \mid E_1) - P(y_{C_0} \mid E_0)$, following the constructions indicated below.

constructions indicated below.

mechanism

unit

mechanisms

units

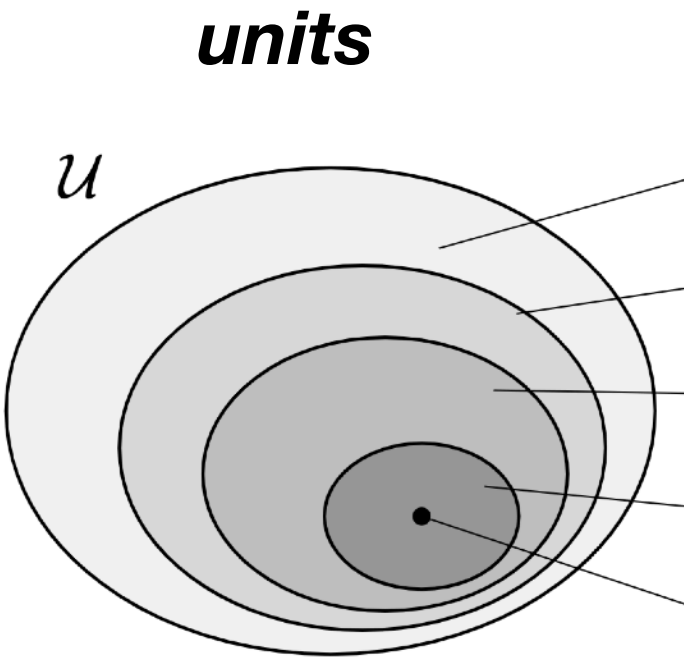
A Venn diagram illustrating the relationship between different units. It consists of four nested ellipses. The outermost ellipse is labeled \mathcal{U} . Inside it is an ellipse labeled $X = x$. Inside that is an ellipse labeled $Z = z$. The innermost ellipse is labeled $V' \subseteq V$. A central point is marked within the innermost ellipse. Lines connect the labels \mathcal{U} , $X = x$, $Z = z$, and $V' \subseteq V$ to their respective ellipses. A line also connects the central point to the label $V' \subseteq V$.

	Measure	C_0	C_1	E_0	E_1
general	TV_{x_0,x_1}	\emptyset	\emptyset	x_0	x_1
	TE_{x_0,x_1}	x_0	x_1	\emptyset	\emptyset
	$Exp-SE_x$	x	x	\emptyset	x
	NDE_{x_0,x_1}	x_0	x_1, W_{x_0}	\emptyset	\emptyset
	NIE_{x_0,x_1}	x_0	x_0, W_{x_1}	\emptyset	\emptyset
$X = x$	ETT_{x_0,x_1}	x_0	x_1	x	x
	$Ctf-SE_{x_0,x_1}$	x_0	x_0	x_0	x_1
	$Ctf-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	x	x
	$Ctf-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	x	x
$Z = z$	$z-TE_{x_0,x_1}$	x_0	x_1	z	z
	$z-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	z	z
	$z-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	z	z
$V' \subseteq V$	$v'-TE_{x_0,x_1}$	x_0	x_1	v'	v'
	$v'-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	v'	v'
	$v'-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	v'	v'
unit	$unit-TE_{x_0,x_1}$	x_0	x_1	u	u
	$unit-DE_{x_0,x_1}$	x_0	x_1, W_{x_0}	u	u
	$unit-IE_{x_0,x_1}$	x_0	x_0, W_{x_1}	u	u

Causal

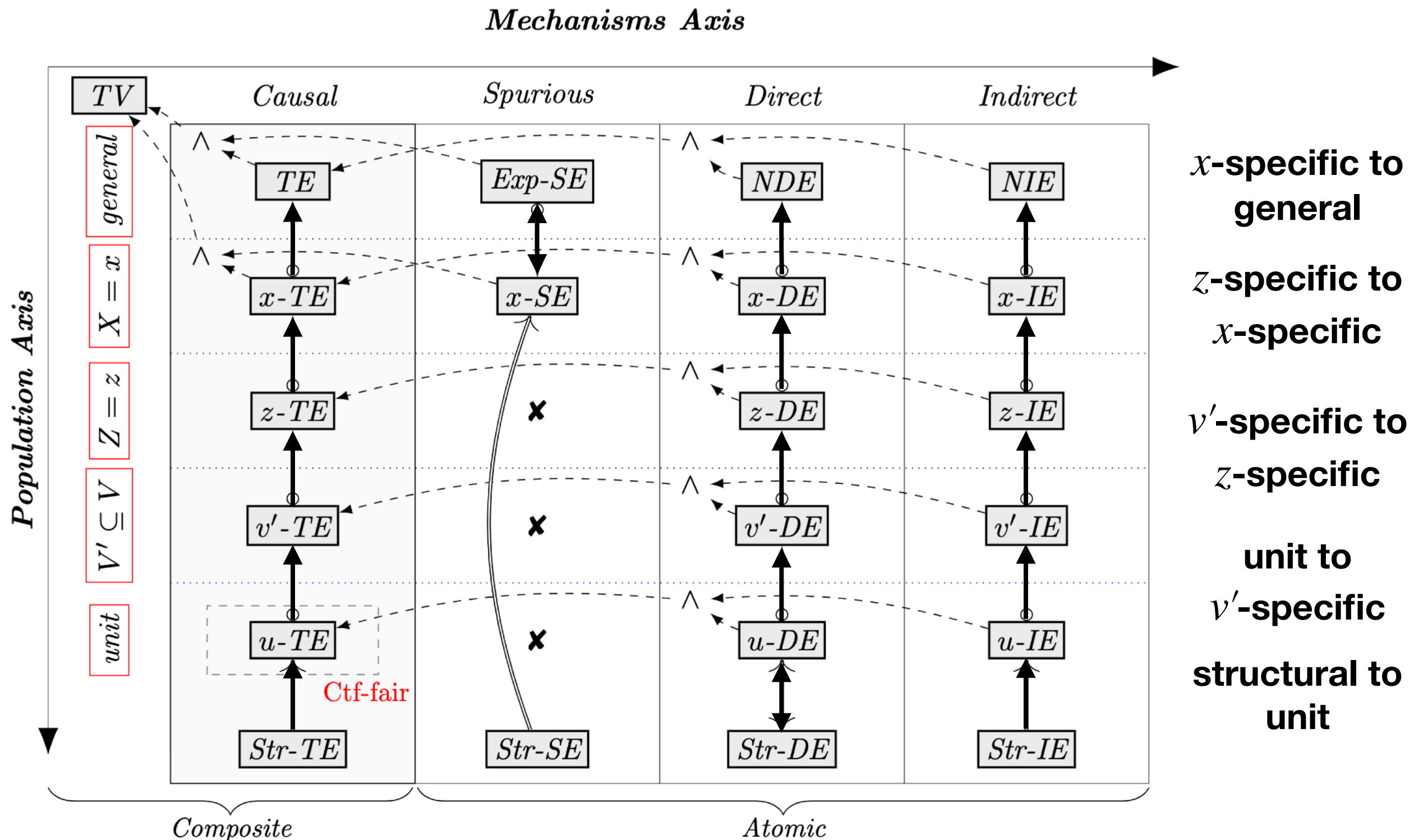
TV

Spurious



Fairness Map

Fairness Map

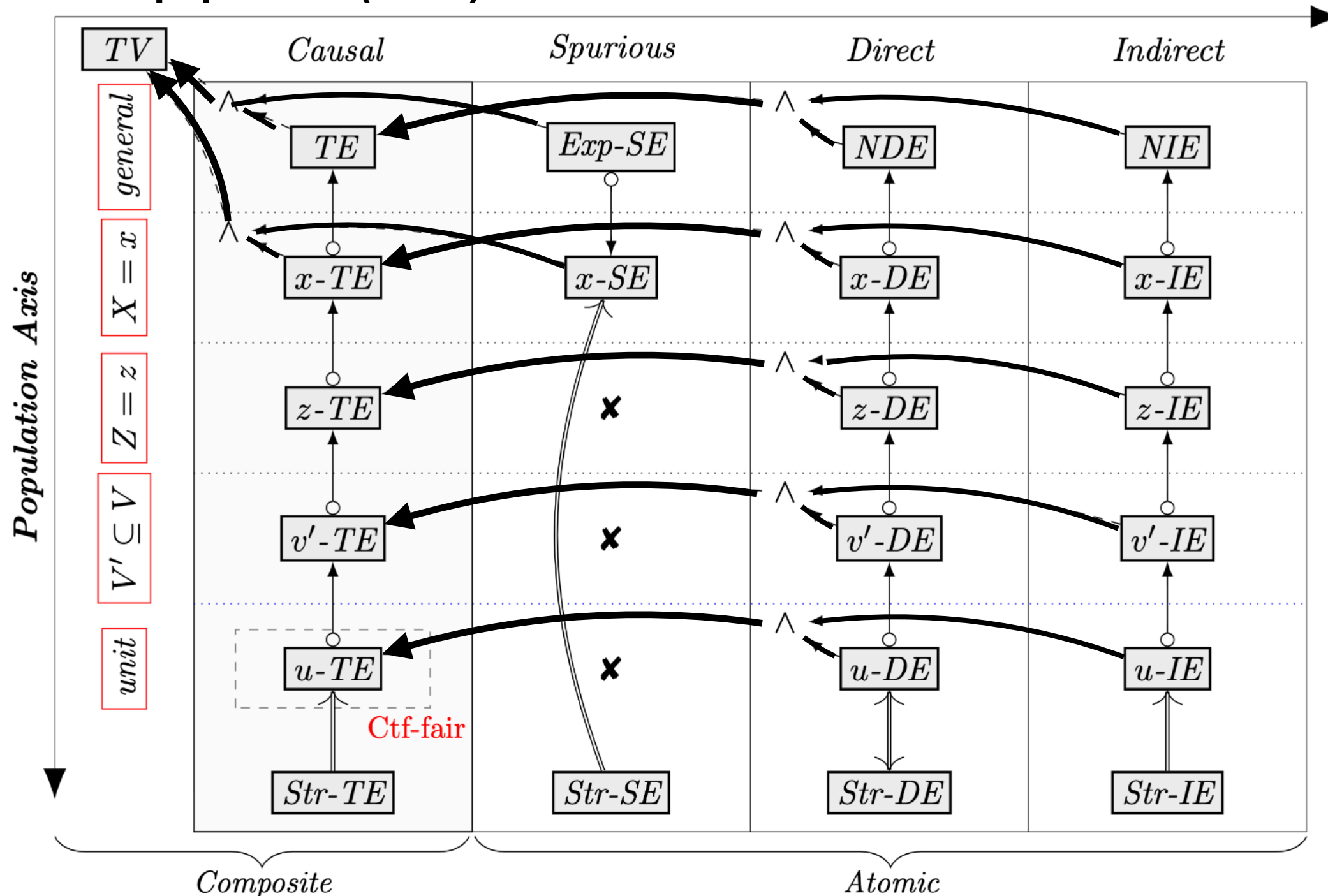


Fairness Map

Section 4.2
Theorem 7

TV decomposition (ZB18)

Mechanisms Axis



Mediation
formula
(Pearl, 2012)

Extended
Mediation
Formula

Other connections with the literature

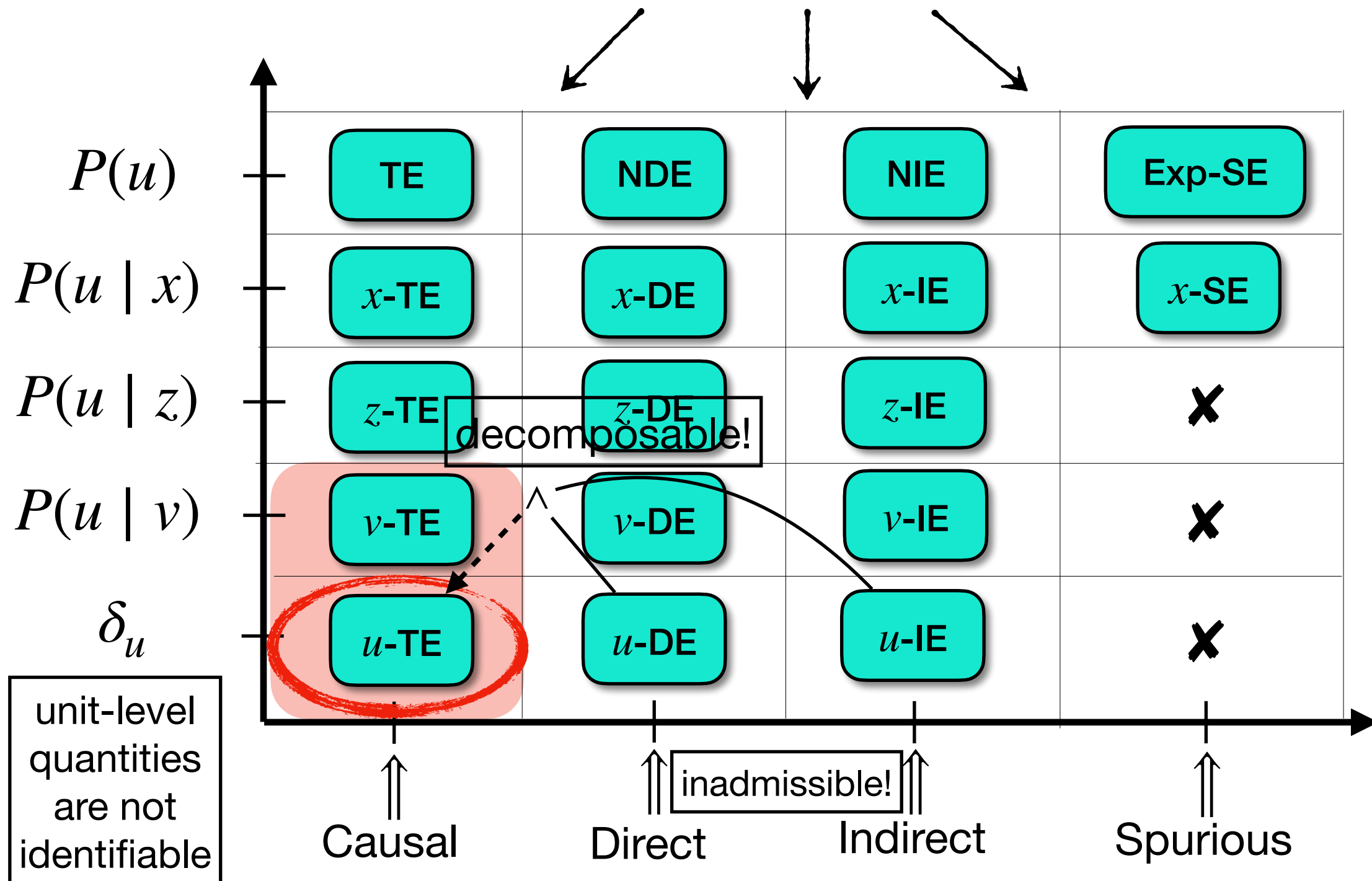
- How does the presented framework relates to other prominent measures in the literature?
- In particular, we consider the following measures:

(i) Counterfactual Fairness (Kusner et. al., '17)

(ii) Individual Fairness (Dwork et. al., '12)

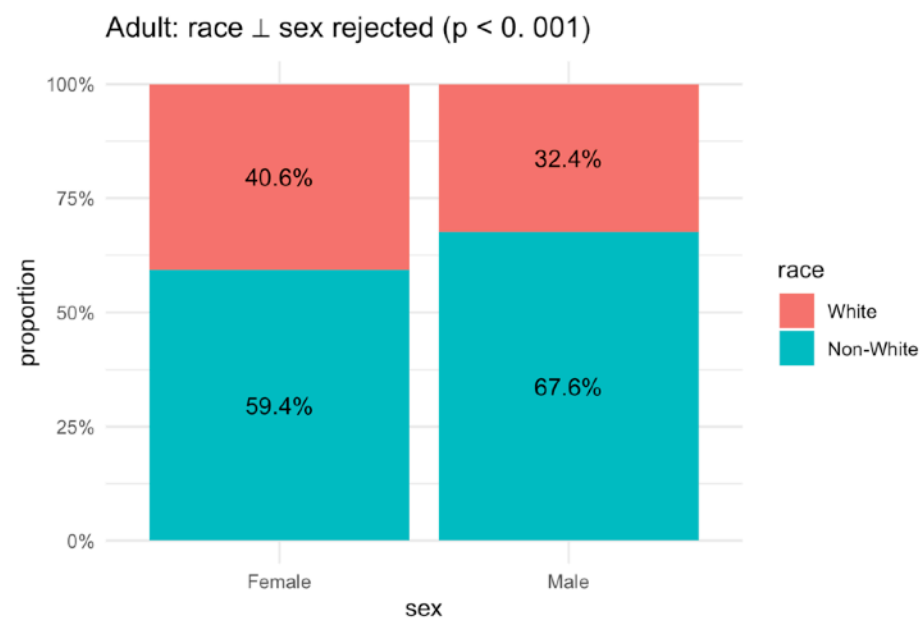
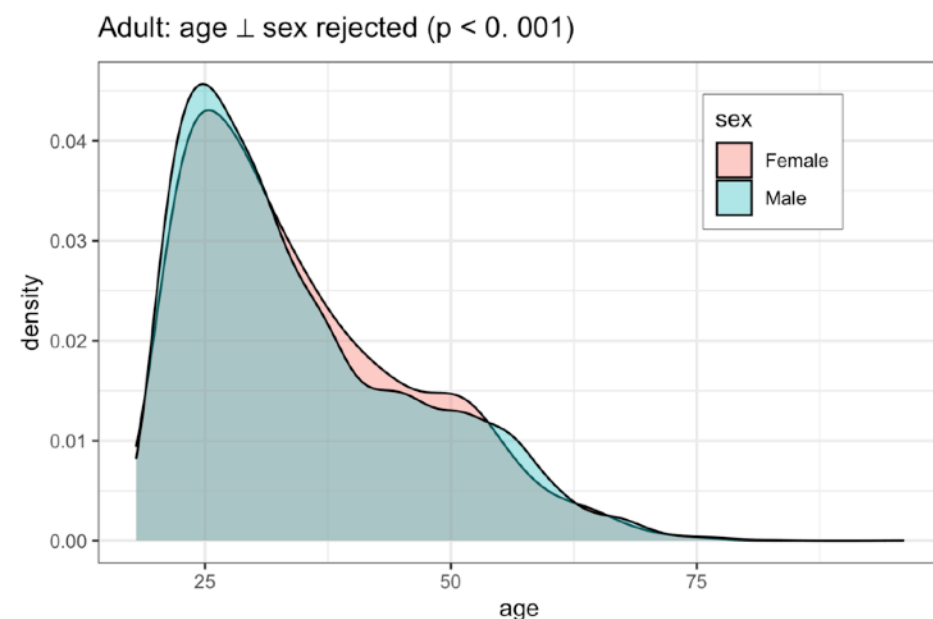
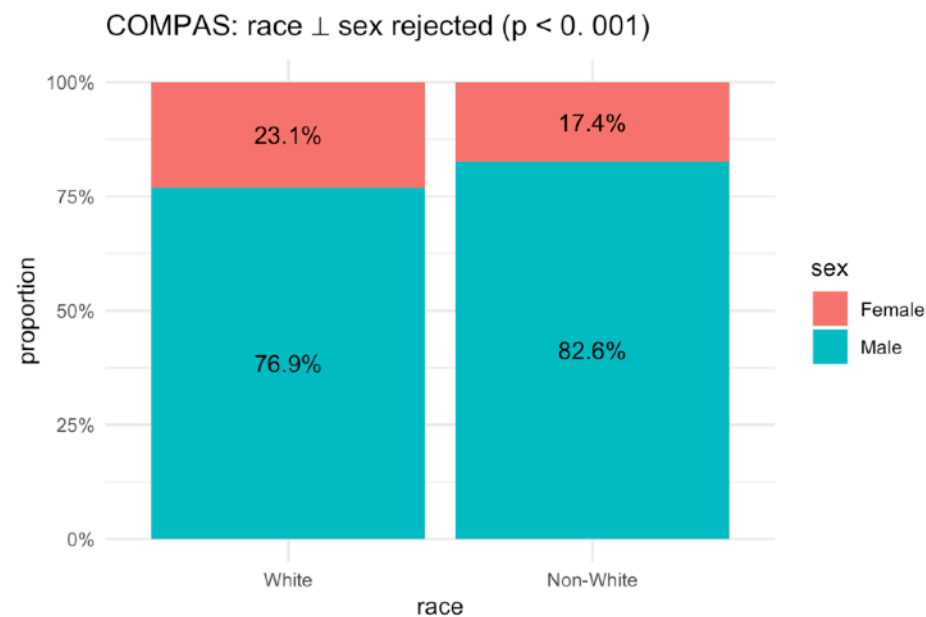
Counterfactual fairness (Kusner et. al., 2017)

$$TV = E[Y \mid \text{male}] - E[Y \mid \text{female}]$$



Counterfactual fairness (Kusner et. al., 2017)

Assumption: ancestral closure of set X .



redlining

religious segregation

rural/urban balance
of genders in China

Counterfactual fairness (Kusner et. al., 2017)

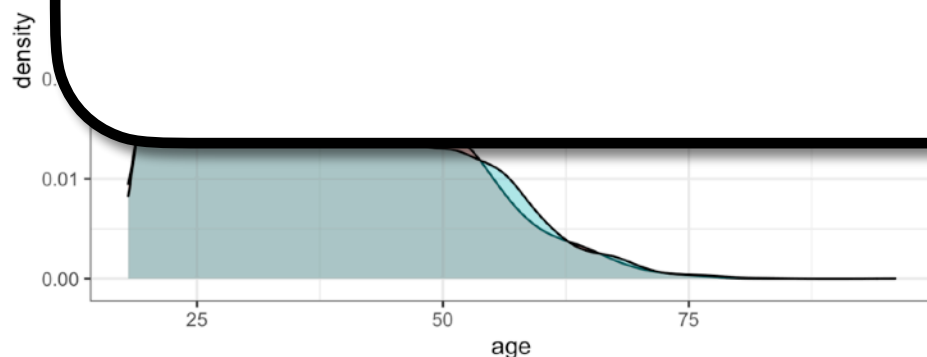
Assumption: ancestral closure of set X .

H

In summary, Counterfactual Fairness is:

- decomposable & inadmissible (w.r.t DE, IE, SE),
- not identifiable in general, and
- oblivious to spurious effects (and corresponding business necessity requirements).

See also Section 4.4.1 for further details.



gation
lance
China

Individual Fairness

(Dwork. et. al., 2012)

Causal Fairness Analysis implications on IF:

- IF is oblivious to the underlying causal mechanisms.

Example 17
Section 4.4.2

- IF captures the direct effect only under the SFM.

Proposition 5
Section 4.4.2

- IF with a sparse metric d is not admissible.

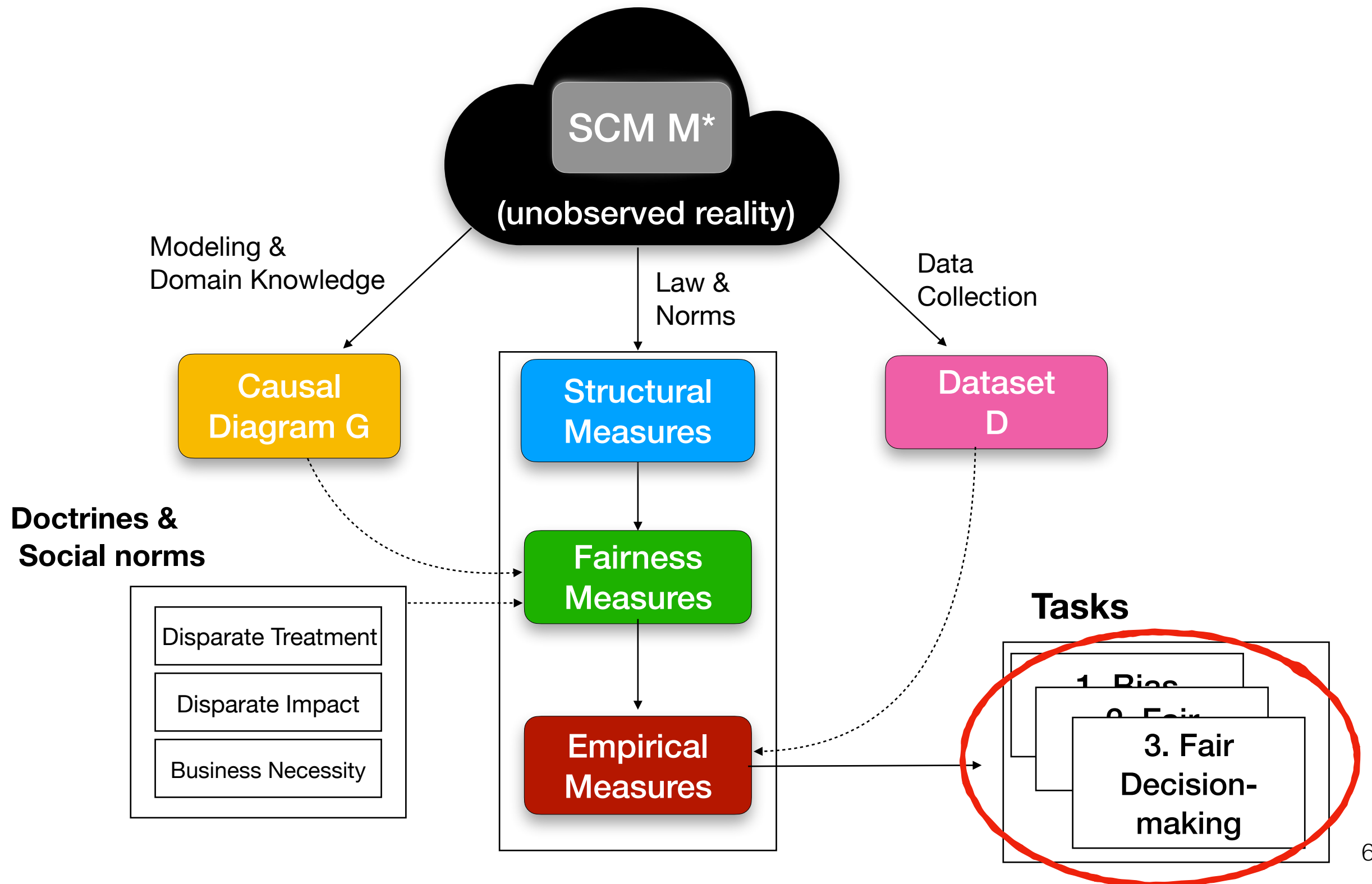
Example 18
Section 4.4.2

- IF with a complete metric d doesn't account for business necessity.

Proposition 11
Section 4.4.2

Part II

Fairness Tasks (Big Picture)



Task 1.

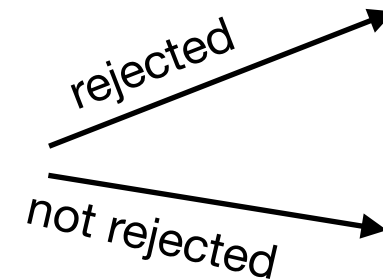
Bias Detection & Quantification

Fairness Cookbook

- 1) Obtain data on past decisions \mathcal{D} .
- 2) Determine the (possibly simplified) causal diagram \mathcal{G} (w.r.t. underlying \mathcal{M}^*).
- 3) Determine the **Business Necessity** (BN) set (\emptyset , $\{Z\}$, $\{W\}$, $\{Z, W\}$).

- 4) Consider existence of **Disparate Treatment**:

$$H_0^{(x\text{-DE})} : x\text{-DE}_{x_0, x_1}(y \mid x_0) = 0.$$



evidence of disparate treatment (population level)

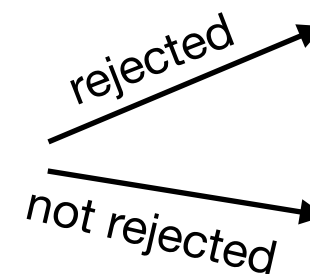
no evidence of disparate treatment (population level)

- 5) Consider existence of **Disparate Impact**:

5a) Indirect effect:

if($W \notin$ BN-set)

$$H_0^{(x\text{-IE})} : x\text{-IE}_{x_0, x_1}(y \mid x_0) = 0.$$



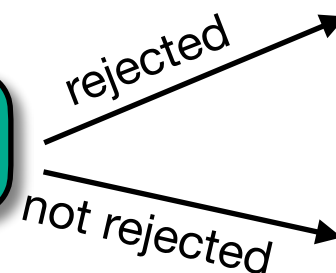
evidence of disparate impact

go to next step

5b) Spurious effect:

if($Z \notin$ BN-set)

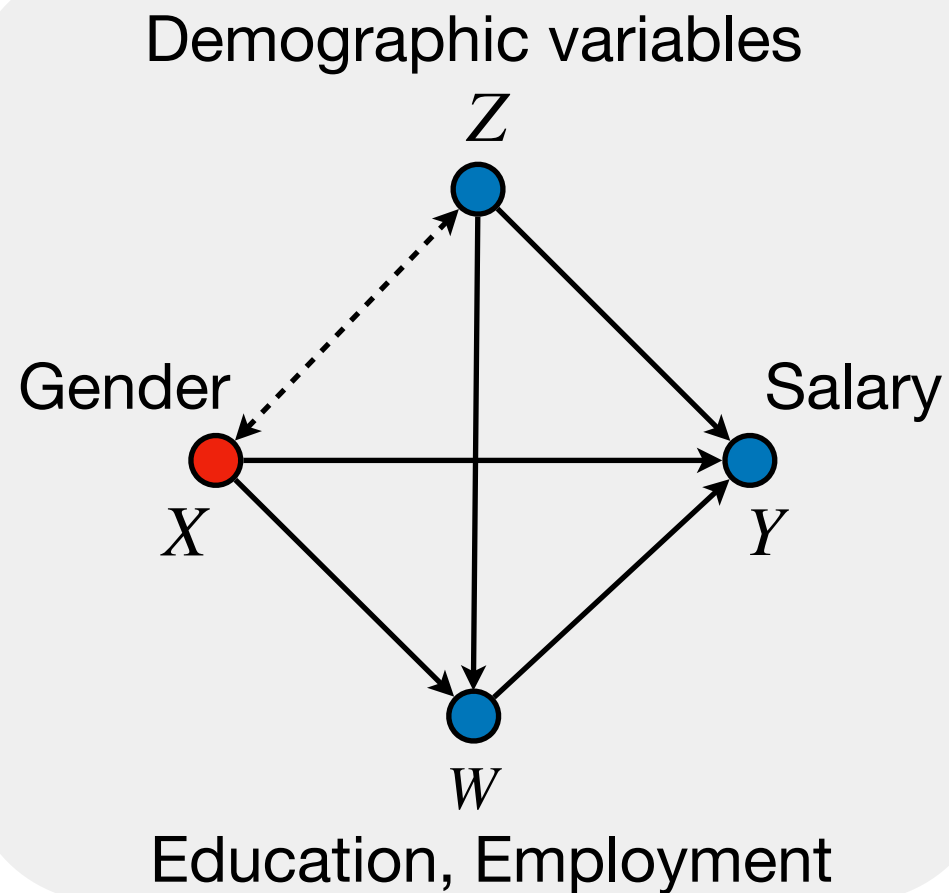
$$H_0^{(x\text{-SE})} : x\text{-SE}_{x_0, x_1}(y) = 0.$$



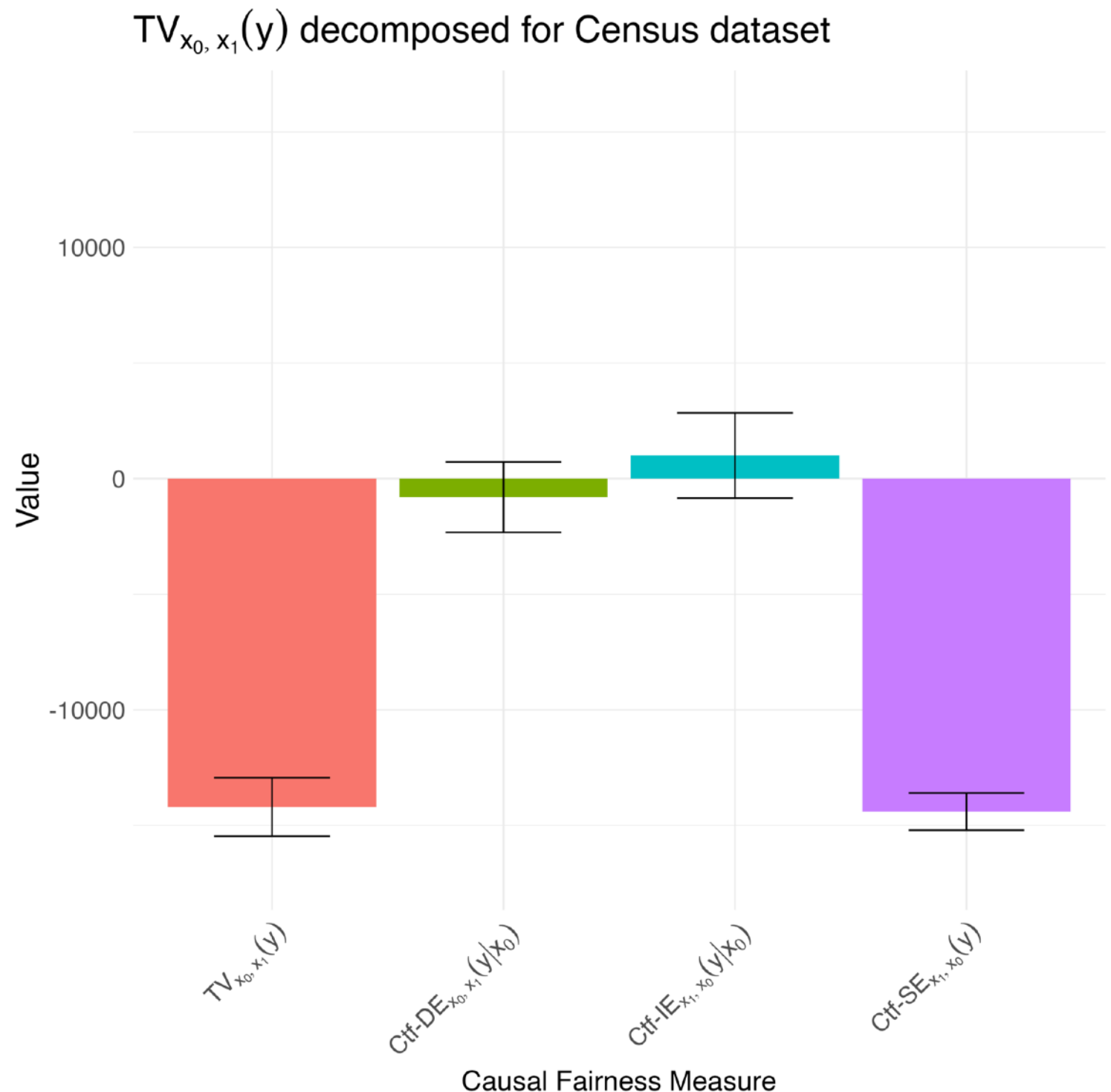
evidence of disparate impact

no evidence of disparate impact

Task 1: Census 2018 dataset



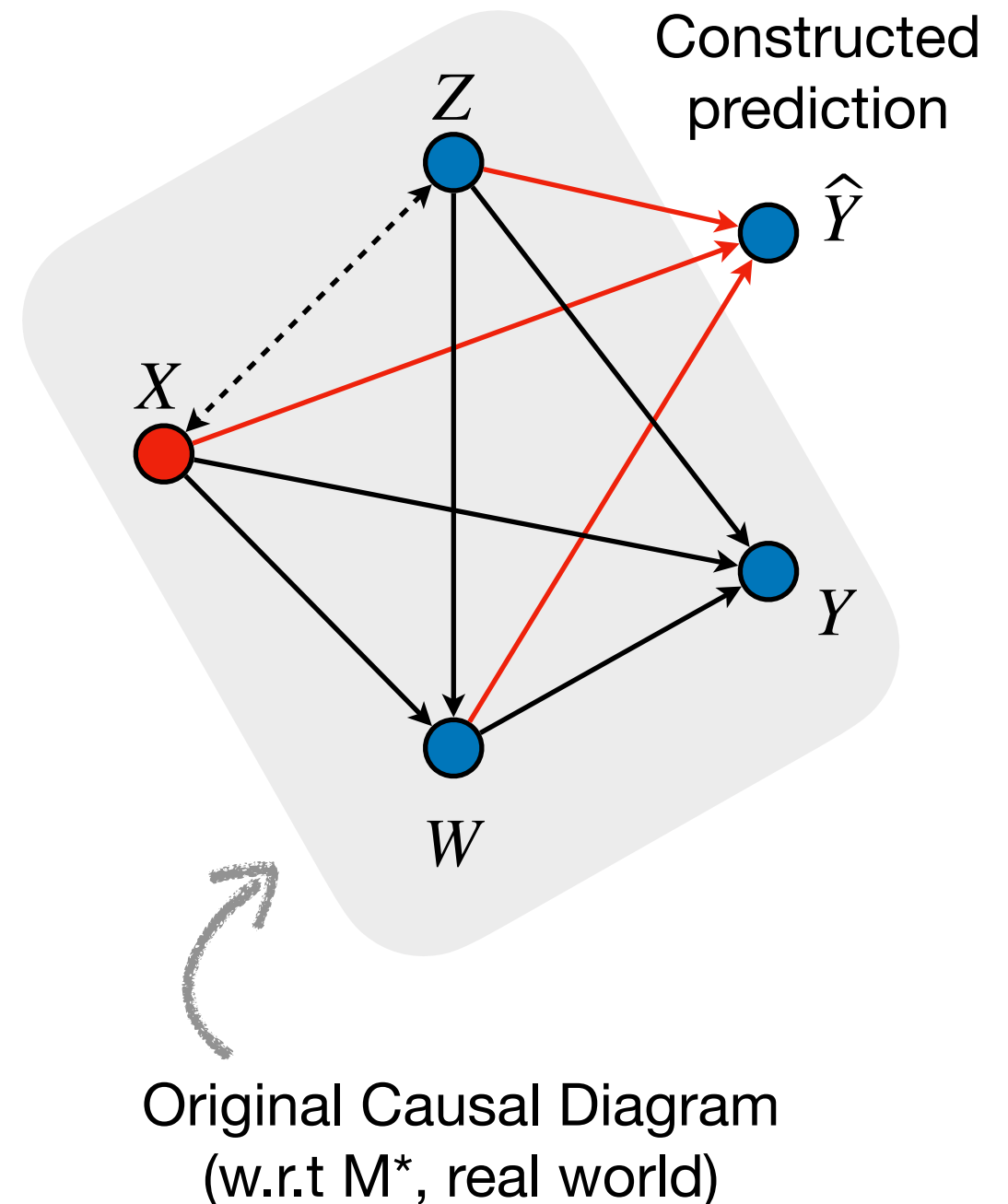
- Observed disparity:
 $TV_{x_0, x_1}(y) = \$14,000/\text{year}$



Task 2. Fair Predictions

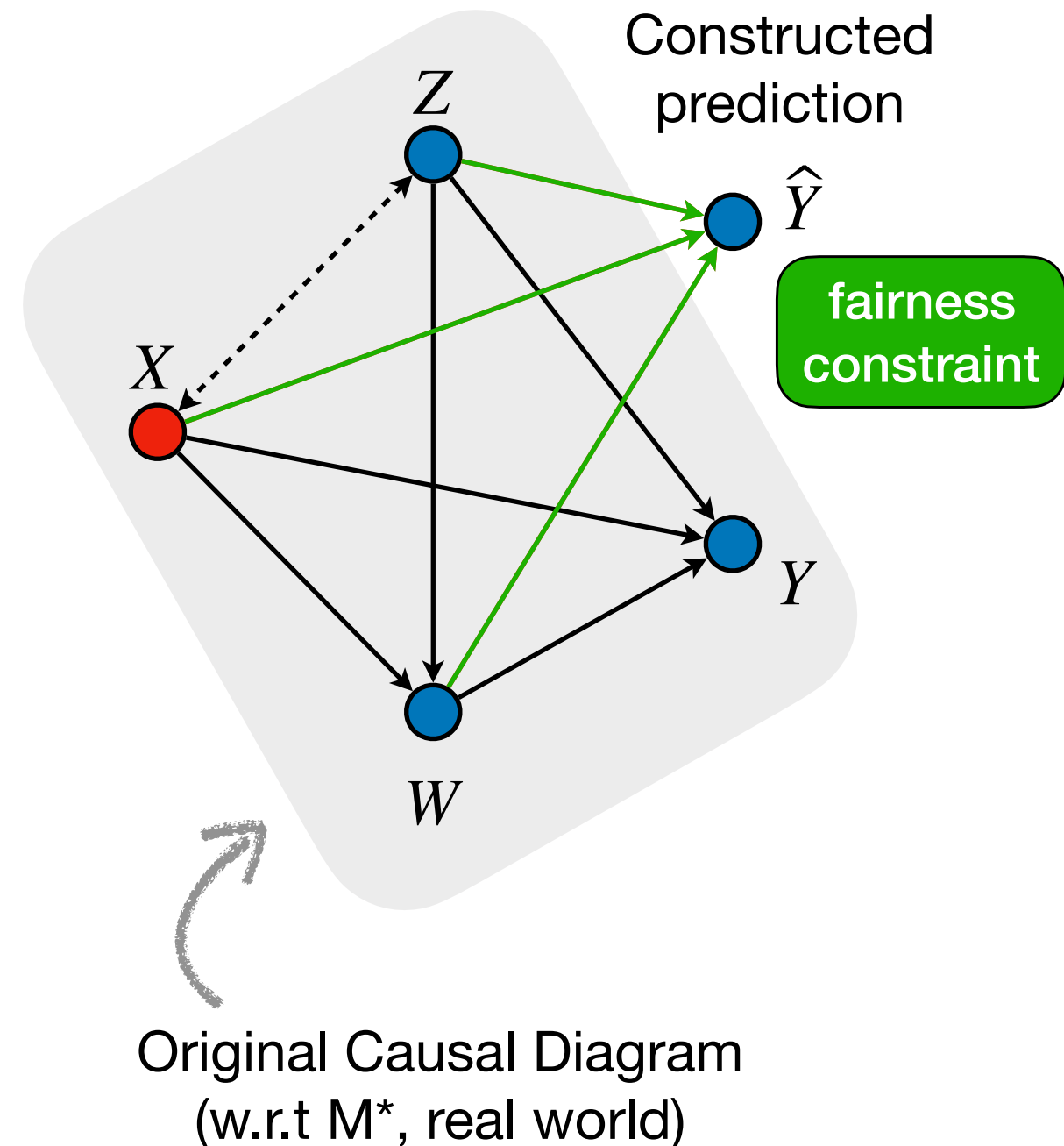
Prediction Task

- The first talk focused on bias detection, where we just analyze the “observed reality”, i.e., nature defines f_Y
- When doing prediction, causally speaking, we are constructing a new mechanism $\hat{Y} \leftarrow f_{\hat{Y}}(x, z, w)$ that is under our control (i.e., we are selecting it)
- Typically, in ML, we are simply interested in learning $P(y \mid x, z, w)$
- Does that carry over bias from f_Y ?



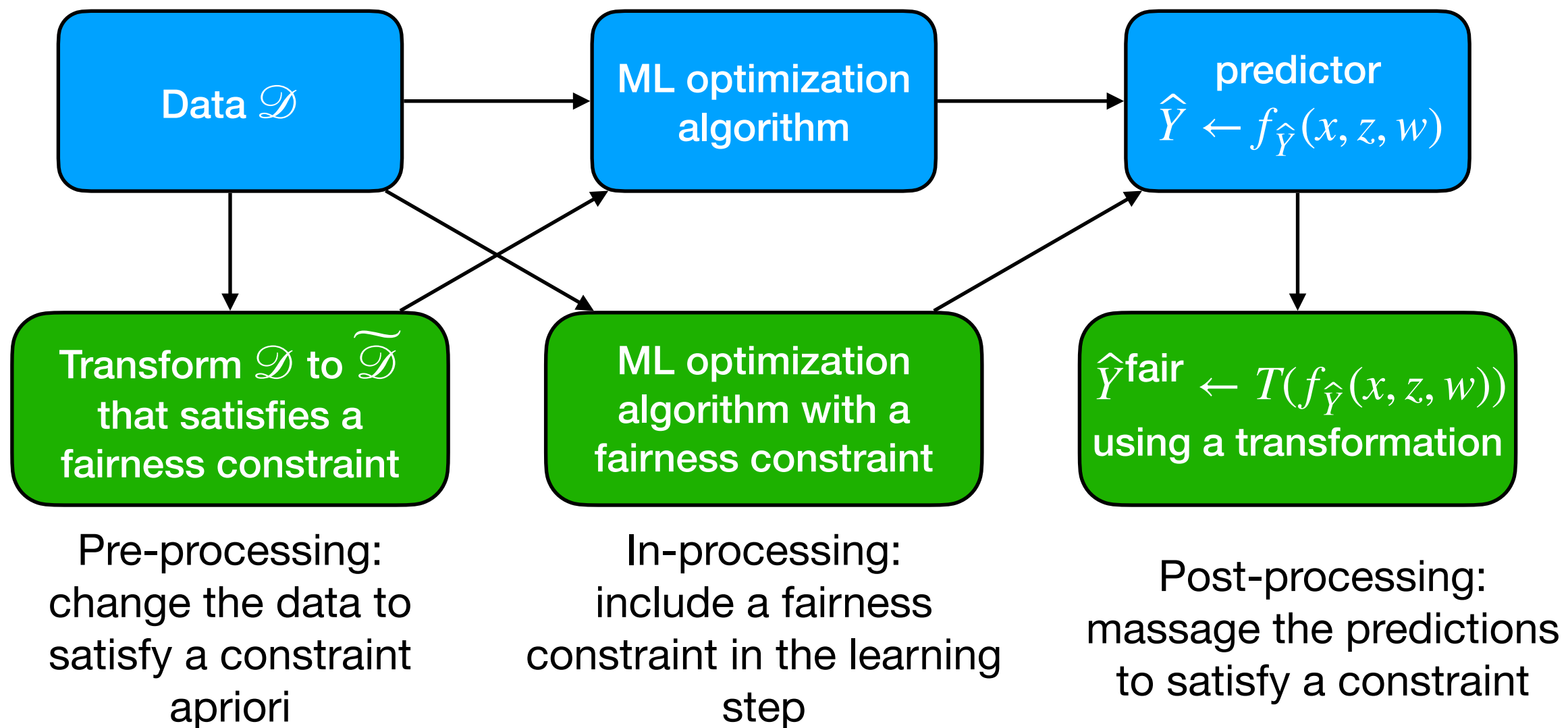
Fair Prediction

- General answer: simply learning $P(y \mid x, z, w)$ will give biased predictions.
- To remove the bias, one might wish for \hat{Y} to satisfy a pre-specified fairness constraint.
- A commonly considered constraint is to make $TV_{x_0, x_1}(\hat{Y}) = 0$.
- In practice, there are different ways to satisfying such a constraint: in particular, we distinguish post-processing, in-processing, and pre-processing methods.



Pre-, In-, Post-Processing

Typical ML framework:



Fair Prediction Theorem (FPT)

Theorem. Let $\text{SFM}(n_Z, n_W)$ be the SFM with $|Z| = n_Z$ and $|W| = n_W$. Let E denote the set of edges of $\text{SFM}(n_Z, n_W)$. Further, let $\mathcal{S}_{n_Z, n_W}^{\text{linear}}$ be the space of linear SCMs (but for the variable X , which is a Bernoulli) compatible with the $\text{SFM}(n_Z, n_W)$ and whose structural coefficients are drawn uniformly from $[-1, 1]^{|E|}$.

An SCM $M \in \mathcal{S}_{n_Z, n_W}^{\text{linear}}$ is said to be ϵ -TV-compliant if

$$\hat{f}_{\text{fair}} = \underset{f \text{ linear}}{\operatorname{argmin}} E[Y - f(X, Z, W)]^2$$

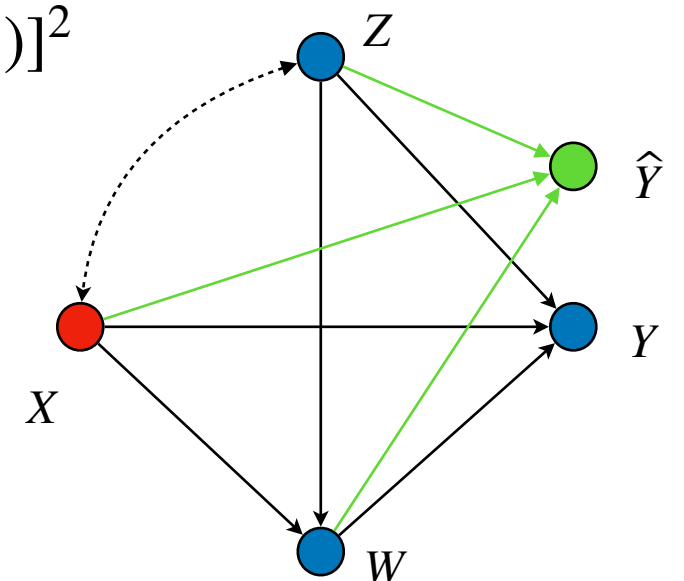
subject to $TV_{x_0, x_1}(f) = 0$

also satisfies

$$|\text{Ctf-DE}_{x_0, x_1}(\hat{f}_{\text{fair}} | x_0)| \leq \epsilon,$$

$$|\text{Ctf-IE}_{x_0, x_1}(\hat{f}_{\text{fair}} | x_0)| \leq \epsilon,$$

$$|\text{Ctf-SE}_{x_0, x_1}(\hat{f}_{\text{fair}})| \leq \epsilon.$$



Under the Lebesgue measure

Furthermore, for any n_Z, n_W th

**non-vanishing probability
of things “going wrong”**

SFM ϵ -TV-compliant with $\epsilon > 0$.

**Section 5.2
Theorem 10**

FPT proof sketch

Objective:

$$Y = \sum_{V_i \in X, Z, W} a_{V_i Y} V_i + \epsilon_Y, \quad f(X, Z, W) = \sum_{V_i \in X, Z, W} \tilde{a}_{V_i Y} V_i.$$

$$\begin{aligned} E[Y - f(X, Z, W)]^2 &= E\left[\sum_{V_i \in X, Z, W} (a_{V_i Y} - \tilde{a}_{V_i Y}) V_i + \epsilon_Y\right]^2 \\ &= E[\epsilon_Y^2] + E\left[\sum_{V_i, V_j \in X, Z, W} (a_{V_i Y} - \tilde{a}_{V_i Y})(a_{V_j Y} - \tilde{a}_{V_j Y}) V_i V_j\right] \\ &= 1 + (a_{VY} - \tilde{a}_{VY})^T E[VV^T] (a_{VY} - \tilde{a}_{VY}), \end{aligned}$$

optimizing over \tilde{a}_{VY}

ellipsoid

Linear SCM:

$$\begin{aligned} U &\leftarrow N(0, 1) \\ X &\leftarrow \text{Bernoulli}(\text{expit}(U)) \\ Z &\leftarrow a_{UZ}U + a_{ZZ}Z\epsilon_Z \\ W &\leftarrow a_{XW}X + a_{ZW}Z + a_{WW}W + \epsilon_W \\ Y &\leftarrow a_{XY}X + a_{ZY}Z + a_{WY}W + \epsilon_Y \end{aligned}$$

$$TV_{x_0, x_1}(f) = (E[V | x_1] - E[V | x_0])^T \tilde{a}_{VY} = 0.$$

what the constraint is

$$Ctf-DE = \tilde{a}_{XY}(x_1 - x_0) = 0,$$

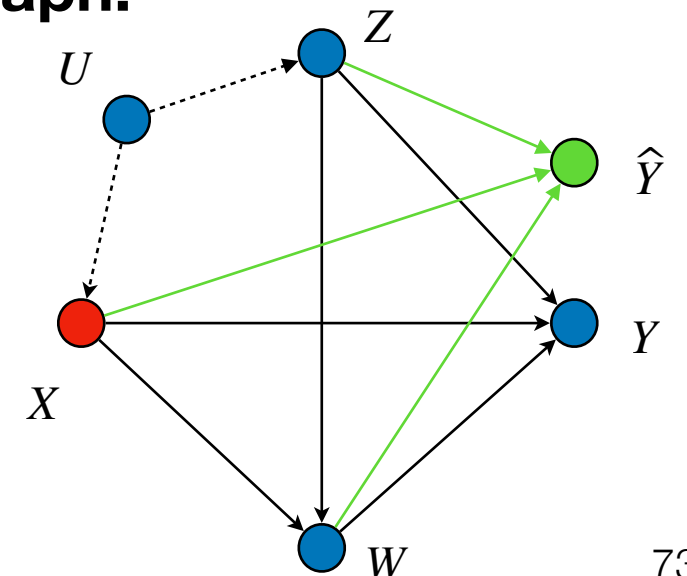
$$Ctf-IE = \sum_{W_i} \tilde{a}_{W_i Y} (E[W_i | x_1] - E[W_{i_{x_0}} | x_1]) = 0,$$

$$Ctf-SE = \sum_{W_i} \tilde{a}_{W_i Y} (E[W_{i_{x_0}} | x_1] - E[W_i | x_0]) +$$

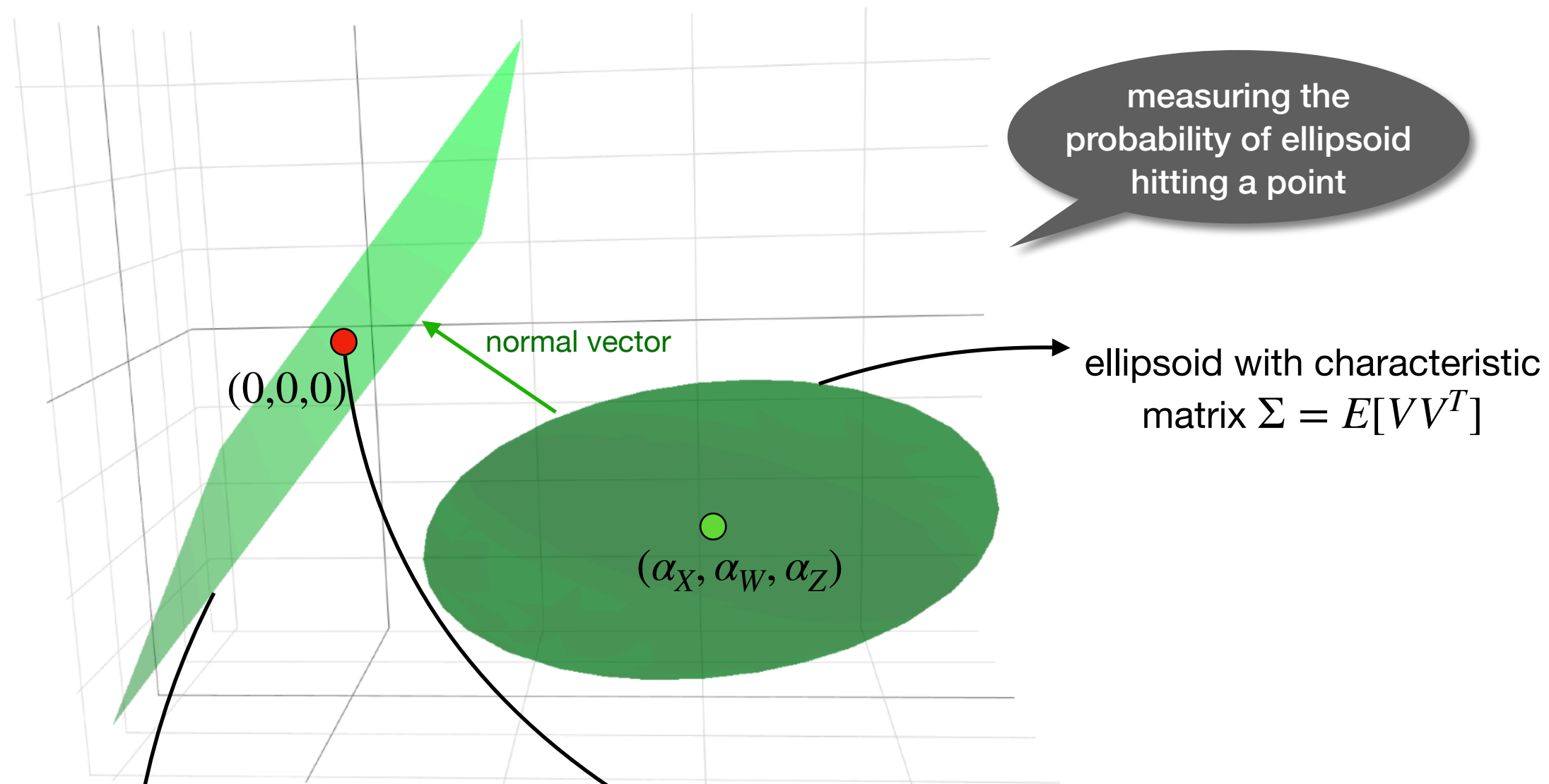
$$\sum_{Z_i} \tilde{a}_{Z_i Y} (E[Z_i | x_1] - E[Z_i | x_0]) = 0.$$

what we actually want

Graph:



FPT visualization

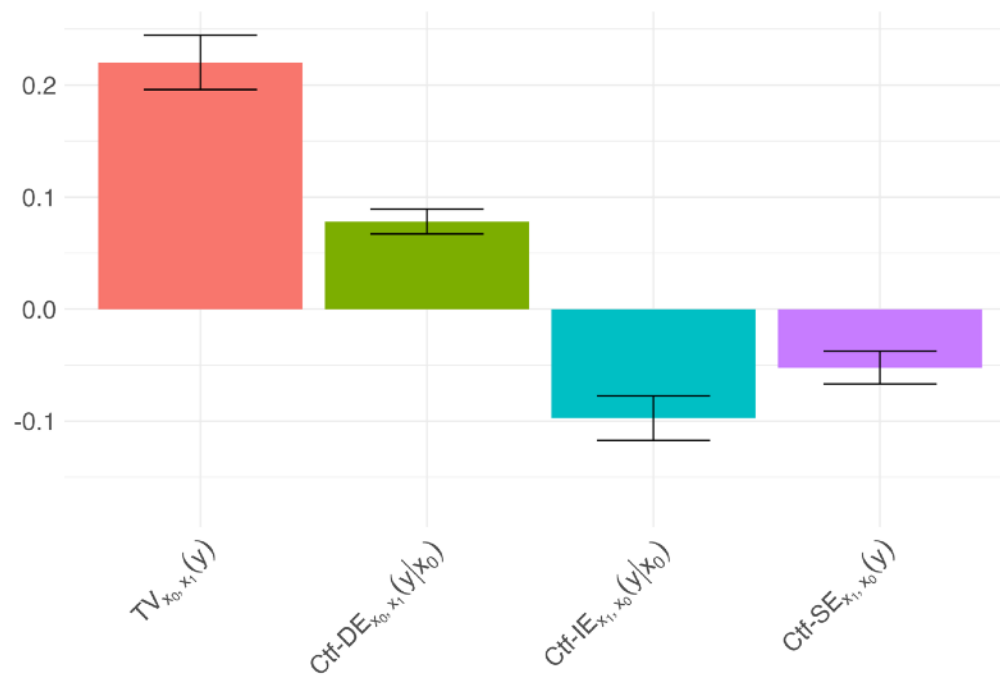


What the constraint is:
 $TV_{x_0, x_1}(\hat{y}) = 0$
represents a hyperplane
through origin.

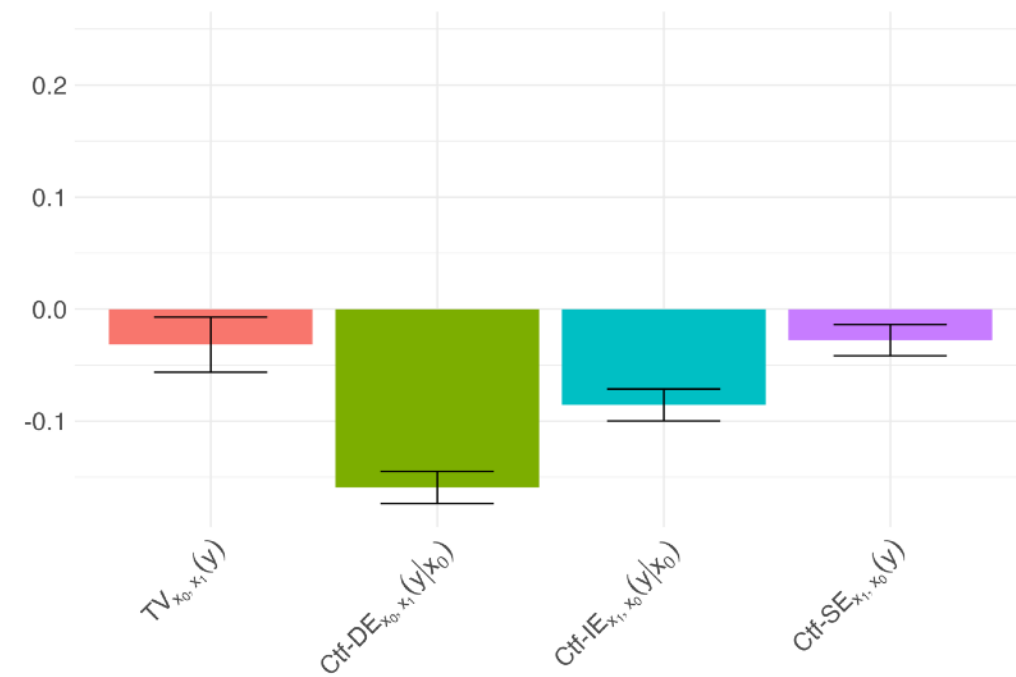
What we want:
3 linear constraints
 $Ctf-DE = 0, Ctf-IE = 0, Ctf-SE = 0.$
represents a single point

Fair Prediction Theorem in Practice (COMPAS dataset)

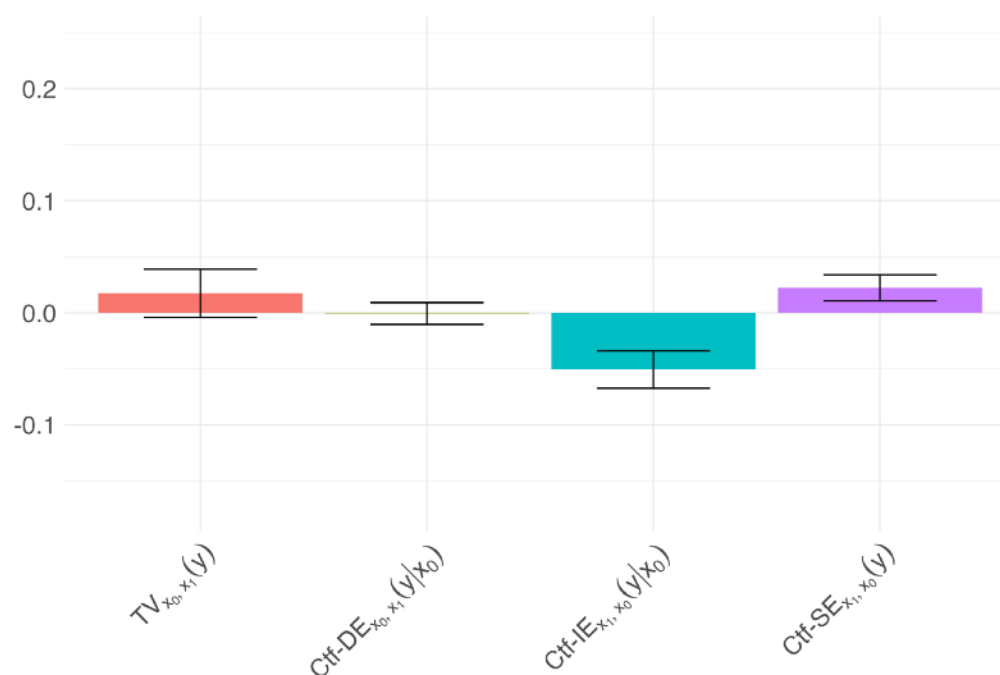
(i) $TV_{x_0, x_1}(\hat{y})$ decomposition: Random Forest on COMPAS



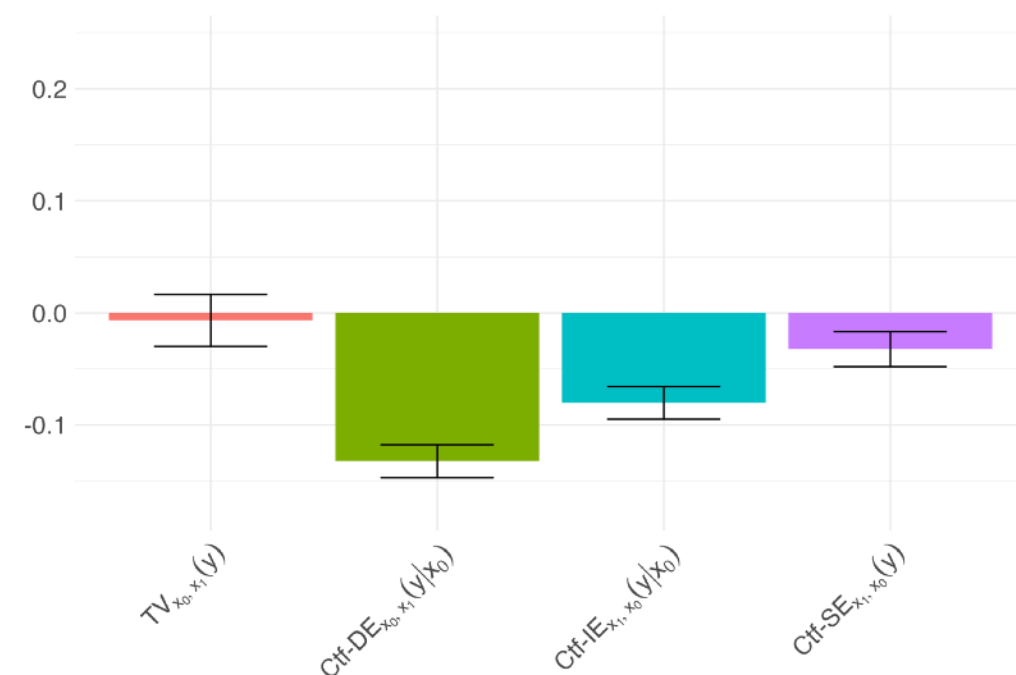
(ii) $TV_{x_0, x_1}(\hat{y})$ decomposition: Reweighting on COMPAS



(iii) $TV_{x_0, x_1}(\hat{y})$ decomposition: Reductions on COMPAS



(iv) $TV_{x_0, x_1}(\hat{y})$ decomposition: Reject-option on COMPAS



Failure of Optimal Transport (in the Individual Fairness framework)

- A common approach for pre-processing is to use optimal transport
- The distribution $P(V \mid x_1)$ is transported onto $P(V \mid x_0)$

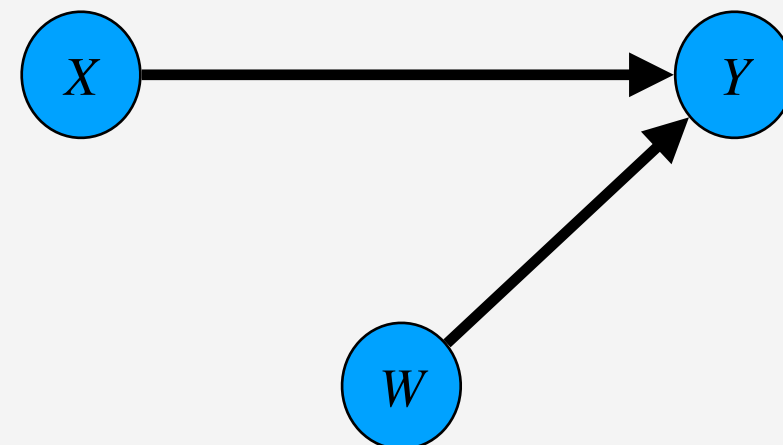
Example.

$$X \leftarrow U_X$$

$$W \leftarrow \epsilon(2U_W - 1)$$

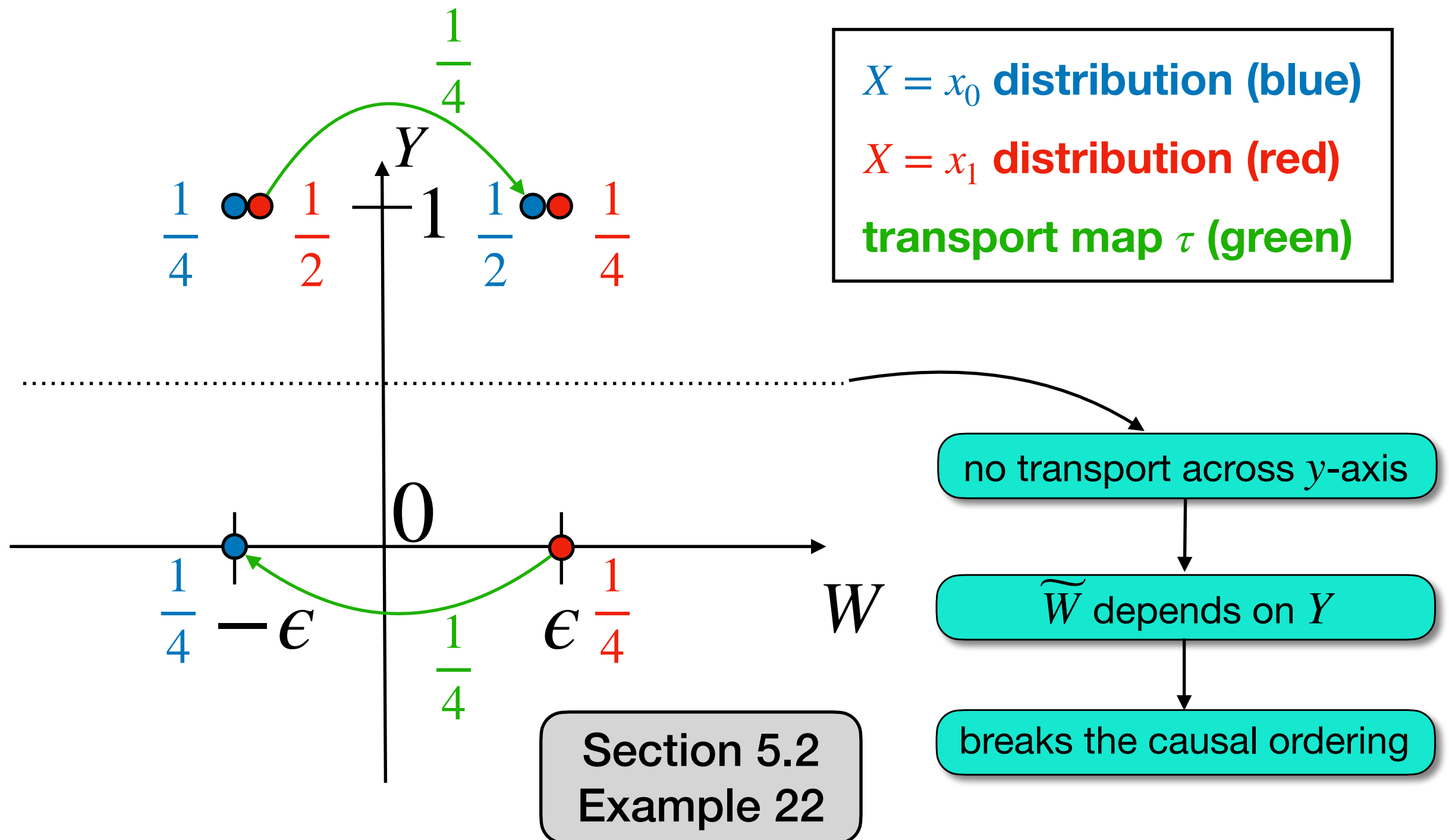
$$Y \leftarrow \begin{cases} U_Y \vee 1(W > 0) & \text{if } X = x_0 \\ U_Y \vee 1(W < 0) & \text{if } X = x_1 \end{cases}$$

$$U_X, U_W, U_Y \text{ Bernoulli}(0.5)$$



- In the example, we wish to compute $\text{NIE}_{x_0, x_1}(\tilde{y}) = P(\tilde{y}_{x_0, \tilde{W}_{x_1}}) - P(\tilde{y}_{x_0})$

Failure of Optimal Transport (in the Individual Fairness framework)



Failure of Optimal Transport (in the Individual Fairness framework)

$$P(\tilde{y}_{x_0}, \tilde{W}_{x_1}) = P(\tilde{y}_{x_0, \epsilon}, \tilde{W}_{x_1} = \epsilon) + P(\tilde{y}_{x_0, -\epsilon}, \tilde{W}_{x_1} = -\epsilon) - P(\tilde{y}_{x_0}) = P(y_{x_0})$$

using the SCM

$$\tilde{y}_{x_0, \epsilon} = 1 \text{ for any } u$$

$$\tilde{W}_{x_1} = \epsilon \text{ for } U_W = 1 \text{ w.p. } \frac{1}{2}$$

$$U_W = 0 \text{ w.p. } \frac{1}{2} \text{ (1/4 for each } U_Y)$$

$$y_{x_0, -\epsilon} = U_Y$$

$$\text{for } U_Y = 1, \tilde{W}_{x_1} = -\epsilon$$

$$\text{with prob. } \frac{1}{4} \text{ (0 for } U_W = 1)$$

$$y_{x_0} = U_Y \vee 1(W > 0)$$

$$\text{for } U_Y = 1, y_{x_0} = 1$$

$$\text{for } U_Y = 0, y_{x_0} = 1$$

$$\text{with prob. } \frac{1}{2}$$

putting together

$$P(\tilde{y}_{x_0}, \tilde{W}_{x_1}) - P(\tilde{y}_{x_0}) = \frac{1}{2} + \frac{1}{8} - \frac{3}{4} = -\frac{1}{8} \Rightarrow$$

Indirect
Effect $\neq 0!$

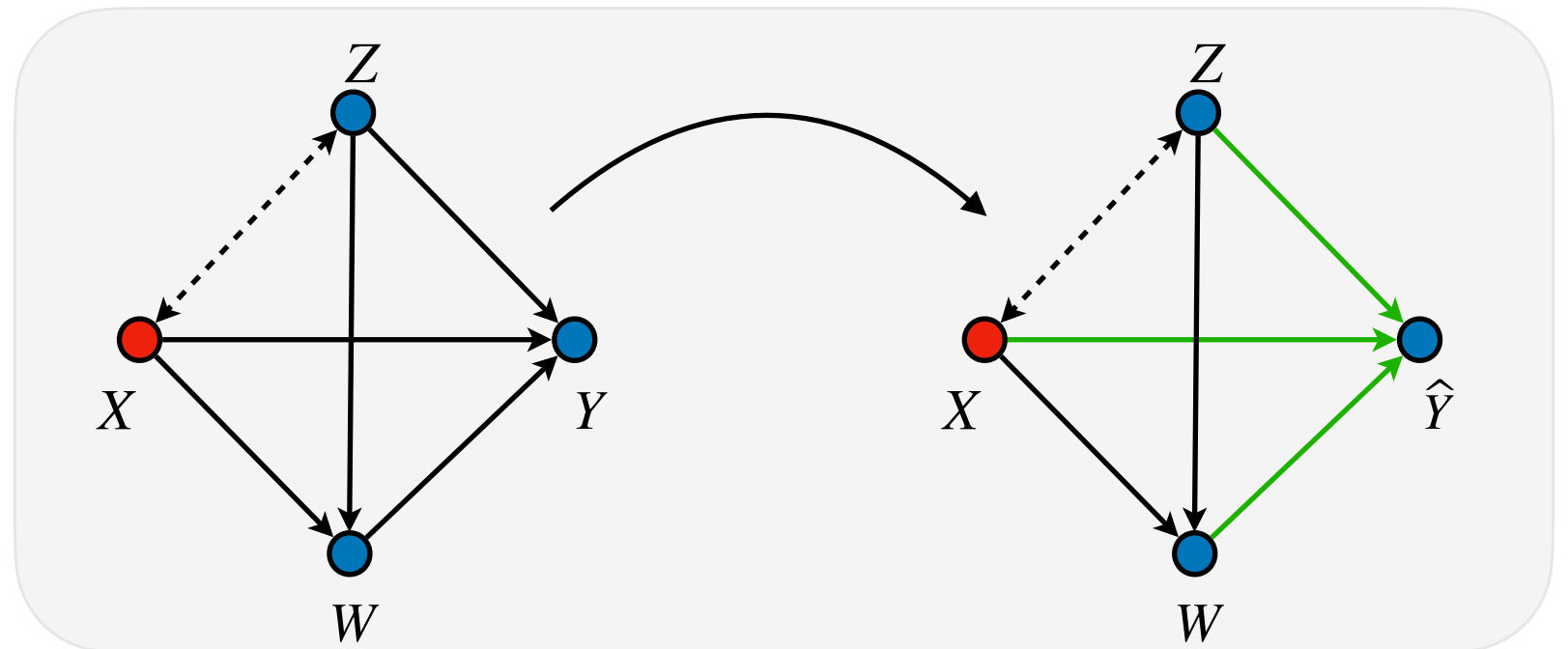
Towards the solution

- how can we construct “causal” fair predictions?

(i) causal structure of the SFM is preserved for the predictor \hat{Y}

Section 5.2.8

(ii) identification expressions for x -DE, x -SE, and x -IE equal 0 for the predictor \hat{Y}



$$\begin{aligned}
 x\text{-DE}_{x_0, x_1}^{ID}(\hat{y}) &= \sum_{z, w} [P(\hat{y} \mid x_1, z, w) - P(\hat{y} \mid x_0, z, w)] P(w \mid x_0, z) P(z \mid x_0) = 0 \\
 x\text{-IE}_{x_0, x_1}^{ID}(\hat{y}) &= \sum_{z, w} P(\hat{y} \mid x_1, z, w) [P(w \mid x_1, z) - P(w \mid x_0, z)] P(z \mid x) = 0 \\
 x\text{-SE}_{x_1, x_0}^{ID}(\hat{y}) &= \sum_z P(\hat{y} \mid x_1, z) [P(z \mid x_1) - P(z \mid x_0)] = 0.
 \end{aligned}$$

In-processing solution

Theorem. Let \hat{Y} be the solution to the following optimization problem:

$$\hat{Y} = \mathbf{argmin}_f \quad E[Y - f(X, Z, W)]^2$$

$$\text{subject to} \quad x\text{-DE}_{x_0, x_1}^{\text{ID}}(\hat{y} \mid x_0) = 0$$

$$x\text{-DE}_{x_1, x_0}^{\text{ID}}(\hat{y} \mid x_0) = 0$$

$$x\text{-IE}_{x_0, x_1}^{\text{ID}}(\hat{y} \mid x_0) = 0$$

$$x\text{-IE}_{x_1, x_0}^{\text{ID}}(\hat{y} \mid x_0) = 0$$

$$x\text{-SE}_{x_1, x_0}^{\text{ID}}(\hat{y}) = 0$$

Section 5.2.9

Then \hat{Y} satisfies

$$x\text{-DE}_{x_0, x_1}(\hat{y} \mid x_0) = x\text{-IE}_{x_1, x_0}(\hat{y} \mid x_0) = x\text{-SE}_{x_1, x_0}(\hat{y}) = 0.$$

Pre-processing solution (Causal IF)

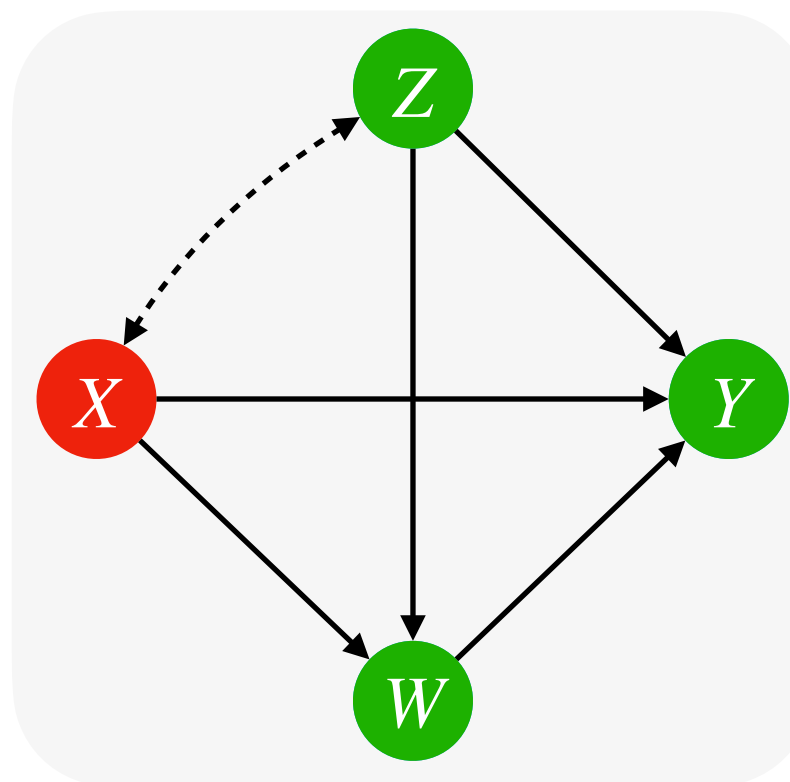
Definition. The Causal Individual Fairness (Causal IF, for short) algorithm is performed on a data coming from an SCM \mathcal{M} compatible with the standard fairness model (SFM), in the following way:

1) if $Z \notin \text{BN-set}$, transport
 $Z \mid x_1 \mapsto Z \mid x_0$

2) if $W \notin \text{BN-set}$, transport
 $W \mid x_1, Z = z \mapsto W \mid x_0, Z = z$

3) transport
 $Y \mid x_1, Z = z, W = w \mapsto Y \mid x_0, Z = z, W = w$

SFM



Data \mathcal{D}

X	Z	W	Y
$x^{(1)}$	$\tilde{z}^{(1)}$	$\tilde{w}^{(1)}$	$\tilde{y}^{(1)}$
$x^{(1)}$	$\tilde{z}^{(2)}$	$\tilde{w}^{(2)}$	$\tilde{y}^{(2)}$
\vdots	\vdots	\vdots	\vdots
$x^{(n)}$	$\tilde{z}^{(n)}$	$\tilde{w}^{(n)}$	$\tilde{y}^{(n)}$

Section 5.2
Theorem 11

Pre-processing solution (Causal IF)

Theorem. Let \mathcal{M} be an SCM compatible with the SFM. Let τ be the optimal transport map obtained when applying Causal IF. Define a new, additional mechanism of the SCM \mathcal{M} such that

$$\tilde{Y} \leftarrow \tau^Y(Y; X, Z, W).$$

For the transformed outcome \tilde{Y} we can then claim:

$$\text{if } Z \notin \text{BN-set} \implies x\text{-SE}_{x_1, x_0}(\tilde{y}) = 0.$$

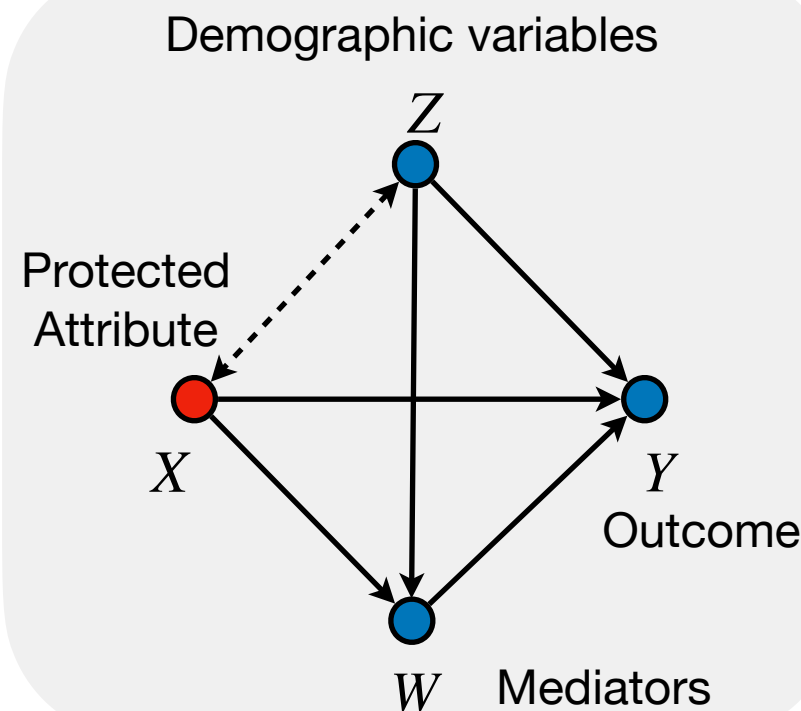
$$\text{if } W \notin \text{BN-set} \implies x\text{-IE}_{x_1, x_0}(\tilde{y} \mid x_0) = 0.$$

Furthermore, the transformed outcome \tilde{Y} also satisfies

$$x\text{-DE}_{x_0, x_1}(\tilde{y} \mid x_0) = 0.$$

Moving beyond SFM

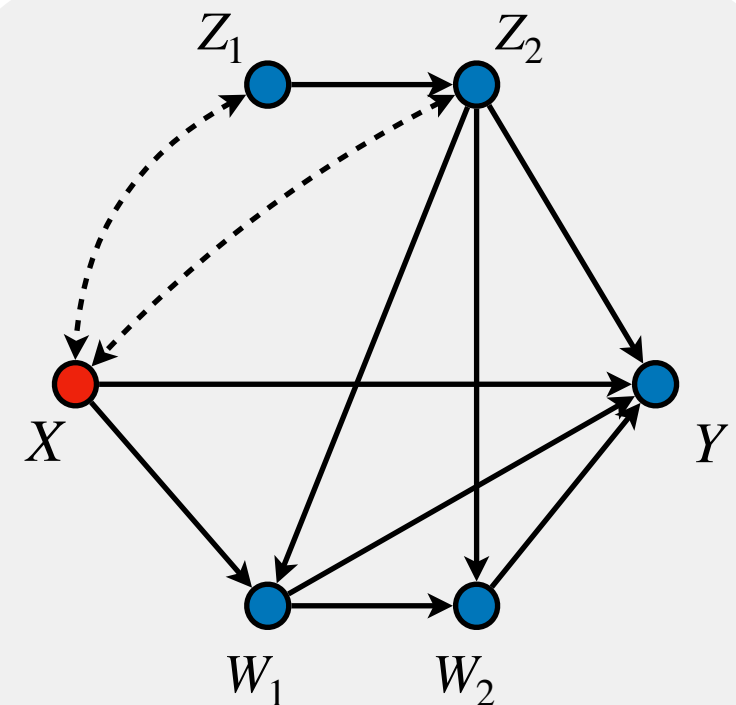
SFM



better resolution

Section 6
TBD

Diagram \mathcal{G}



Measures	direct, indirect, spurious
Business Necessity	$\{\{\emptyset\}, \{Z\}, \{W\}, \{Z, W\}\}$
Fair Prediction	Causal IF

Measures	variable specific
Business Necessity	any $V' \subseteq V$
Fair Prediction	fairadapt

Conclusions

- Well-founded disparate treatment and impact claims require the plaintiff to establish **a causal connection** between a defendant's policy and the statistical disparities found in the observed data.

SCOTUS: No fairness claim can be made without solid causal underpinnings.

- We introduced a framework for fairness analysis based on causal inference to support such claims. In particular, we showed
 - A. - how the total variation can be decomposed into variations that can be easily associated with the underlying causal mechanisms, and mapped to disparate impact and disparate treatment doctrines.
 - B. - how the developed foundations of Causal Fairness Analysis can be applied in practice, in the context of bias detection and fair prediction.
- We hope these results can help towards the development of the next generation of AI systems to be more fair, accountable, and transparent.

Thank you!

References: Causal Fairness

- Agarwal, A., Beygelzimer, A., Dudik, M., Langford, J. & Wallach, H. (2018), A reductions approach to fair classification, in 'International Conference on Machine Learning', PMLR, pp. 60–69.
- Bareinboim, E., Correa, J. D., Ibeling, D., & Icard, T. (2022). On pearl's hierarchy and the foundations of causal inference. In *Probabilistic and Causal Inference: The Works of Judea Pearl* (pp. 507-556).
- Barocas, S., & Selbst, A. D. (2016). Big data's disparate impact. *Calif. L. Rev.*, 104, 671.
- Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2), 153-163.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., & Zemel, R. (2012, January). Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference* (pp. 214-226).
- Halpern, J. Y. (2000). Axiomatizing causal reasoning. *Journal of Artificial Intelligence Research*, 12, 317-337.
- Kilbertus, N., Rojas Carulla, M., Parascandolo, G., Hardt, M., Janzing, D., & Schölkopf, B. (2017). Avoiding discrimination through causal reasoning. *Advances in neural information processing systems*, 30.
- Kusner, M. J., Loftus, J., Russell, C., & Silva, R. (2017). Counterfactual fairness. *Advances in neural information processing systems*, 30.
- Larson, J., Mattu, S., Kirchner, L., & Angwin, J. (2016). How we analyzed the COMPAS recidivism algorithm. *ProPublica* (5 2016), 9(1), 3-3.
- Pearl, J. (2000), *Causality: Models, Reasoning, and Inference*, Cambridge University Press, New York. 2nd edition, 2009.
- Pearl, J. (2001), Direct and indirect effects, In *Proceedings of the Seventeenth Conference on Uncertainty in Artificial Intelligence*.
- Pearl, J. & Mackenzie, D. (2018), *The Book of Why: The New Science of Cause and Effect*, 1st edn, Basic Books, Inc., New York, NY, USA.
- Plečko, D., & Meinshausen, N. (2020). Fair data adaptation with quantile preservation. *The Journal of Machine Learning Research*, 21(1), 9776-9819.
- Rutherglen, G. (1987). Disparate impact under title VII: an objective theory of discrimination. *Va. L. Rev.*, 73, 1297.
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., & Dwork, C. (2013, May). Learning fair representations. In *International conference on machine learning* (pp. 325-333). PMLR.
- Zhang, J., & Bareinboim, E. (2018, April). Fairness in decision-making—the causal explanation formula. In *Proceedings of the AAAI Conference on Artificial Intelligence* (Vol. 32, No. 1).

Plečko, D., & Bareinboim, E. (2022). Causal Fairness Analysis. Columbia University Technical Report R-90, CausalAI Lab.