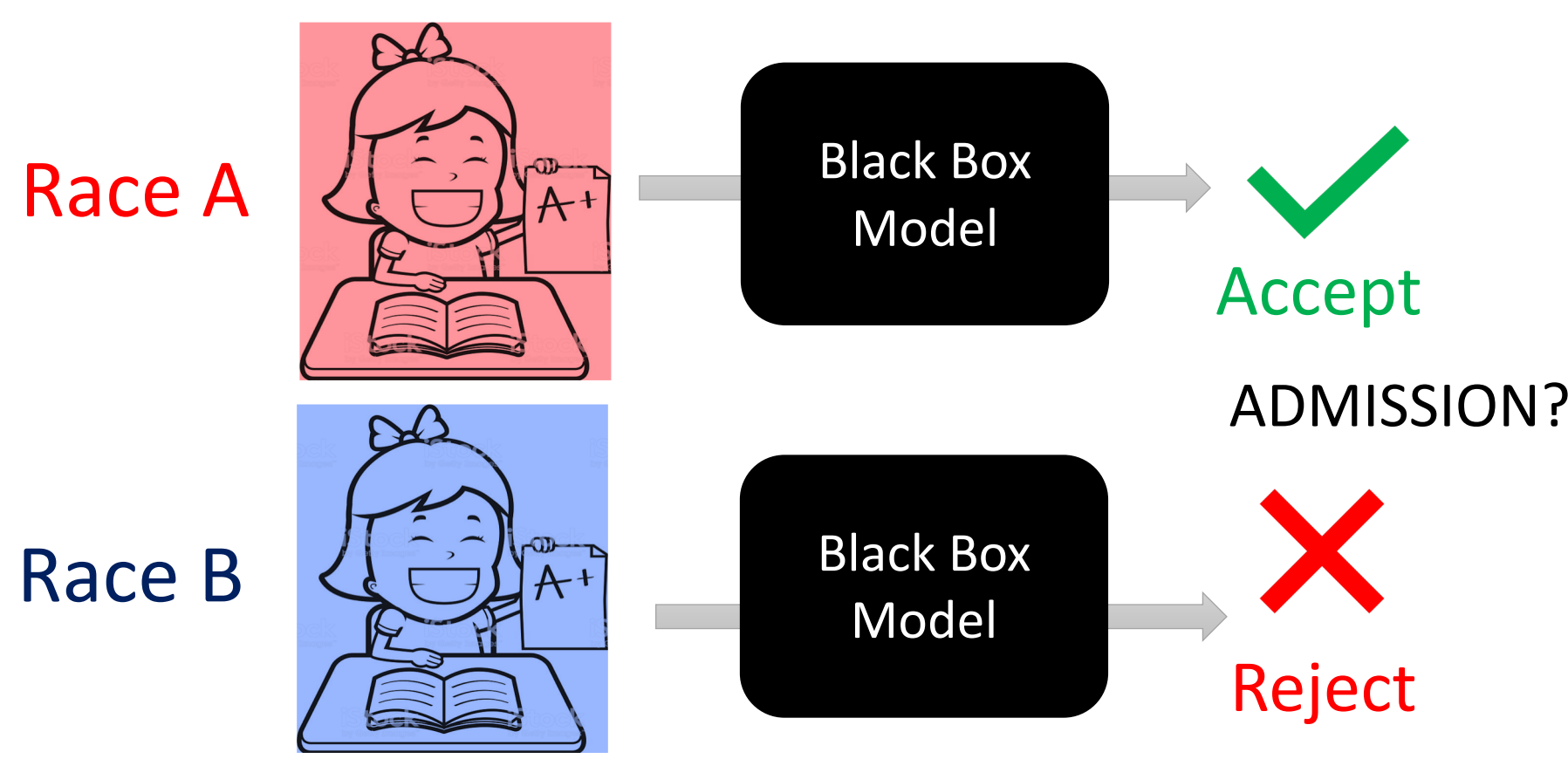


An Information-Theoretic Quantification of Discrimination with Exempt Features

Sanghamitra Dutta, Praveen Venkatesh, Piotr Mardziel, Anupam Datta, Pulkit Grover

Motivation



Fairness without Domain Knowledge?

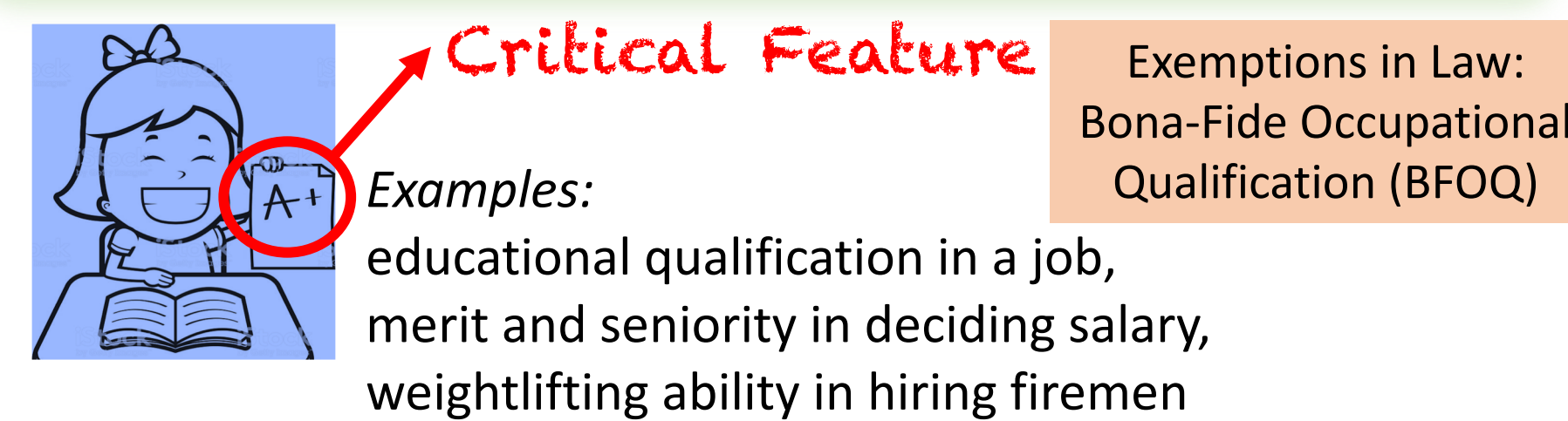
Statistical Parity: Accepts unqualified members of protected group [Zemel et. al. '13][Hardt et. al.'16]

Equalized Odds: Agreement with true labels, may be affected by label bias [Hinnefeld '18][Barocas & Selbst '16]



Decision may not place weight on critical features.

Goal: Quantify "non-exempt" discrimination while allowing for exemptions due to critical features



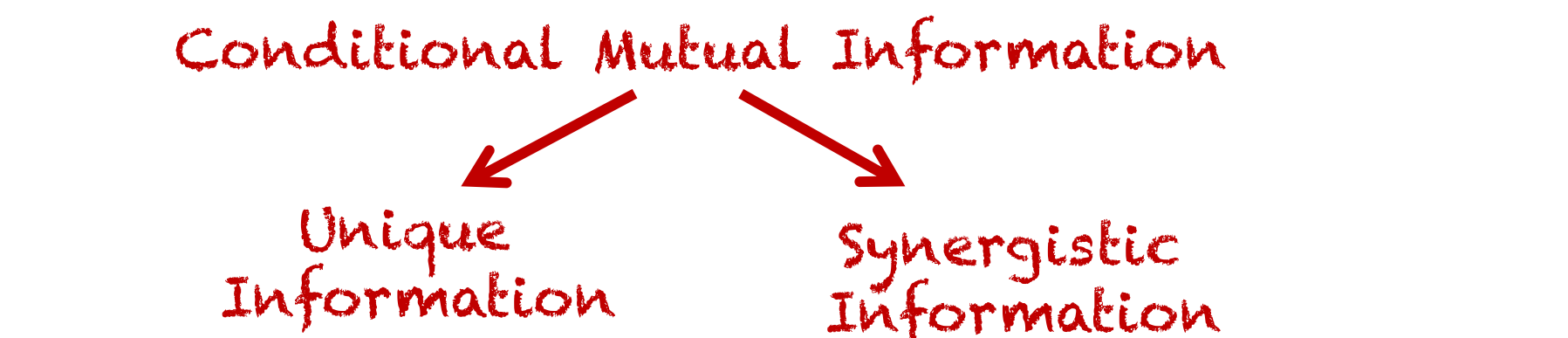
Contributions

Quantification of "non-exempt" discrimination while allowing for exemptions due to critical features

Axiomatic approach, Counterexamples to existing works

E.g., unfair by Conditional Statistical Parity, but fair by Counterfactual Fairness

Lead us to examine Partial Information Decomposition (PID)

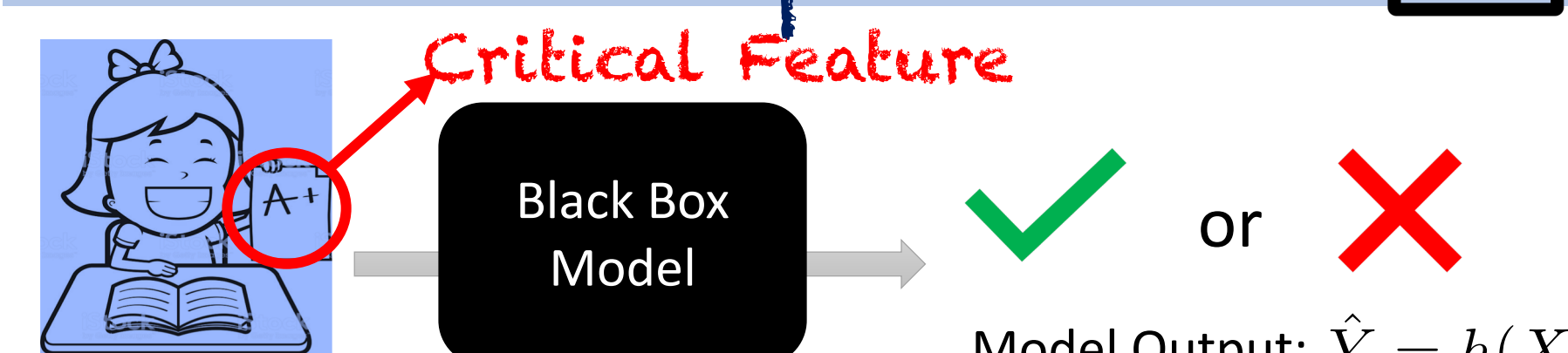


Propose a novel "counterfactual" measure of non-exempt discrimination

Difficult to realize in practice?

Observational measures of "non-exempt" discrimination (impossibility, utility and limitations)

Problem Setup



Total Features: $X = (X_c, X_g)$

Critical

Non-Critical/General

Protected Attribute: Z

Structural Causal Model



Counterfactual Causal Influence (CCI) [Kusner et. al.'17][Datta et. al.'17][Russell et. al.'17]

$\hat{Y} = f(Z, U_X)$
 $Z \perp U_X$

Information-Theoretic Equivalent of CCI

Total Discrimination: $I(Z; W)$

$W = [f(Z, u_a), f(Z, u_b), \dots] \quad \forall u \quad \Pr(U_X = u) > 0$

Thought Experiments

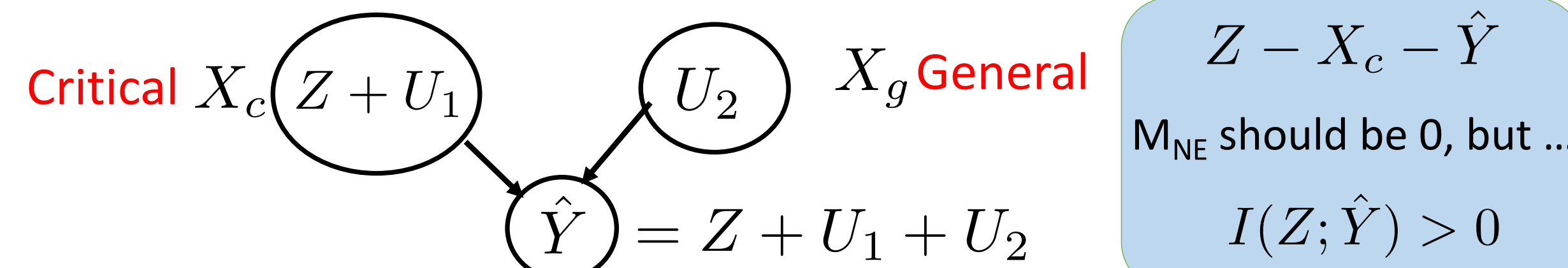
Total Discrimination: $I(Z; W)$



Candidate Measures of Non-Exempt Discrimination

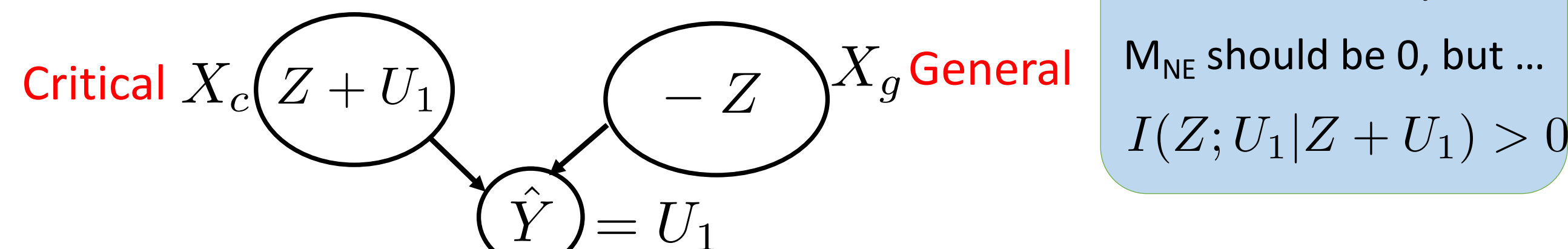
Candidate 1: $I(Z; \hat{Y})$

Counterexample: Hiring actors for a male role (BFOQ Defense)



Candidate 2: $I(Z; \hat{Y}|X_c)$

Counterexample: College admissions (U_1 is the true ability of a candidate)



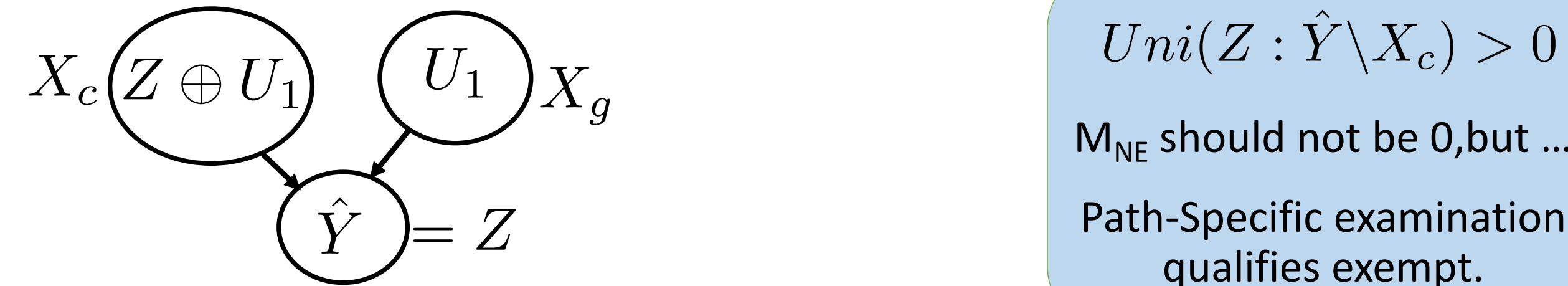
Candidate 3: $Uni(Z; \hat{Y} \setminus X_c)$

Counterexample: Expensive Housing ad shown to high-income Race A and low-income Race B (largely irrelevant to latter)

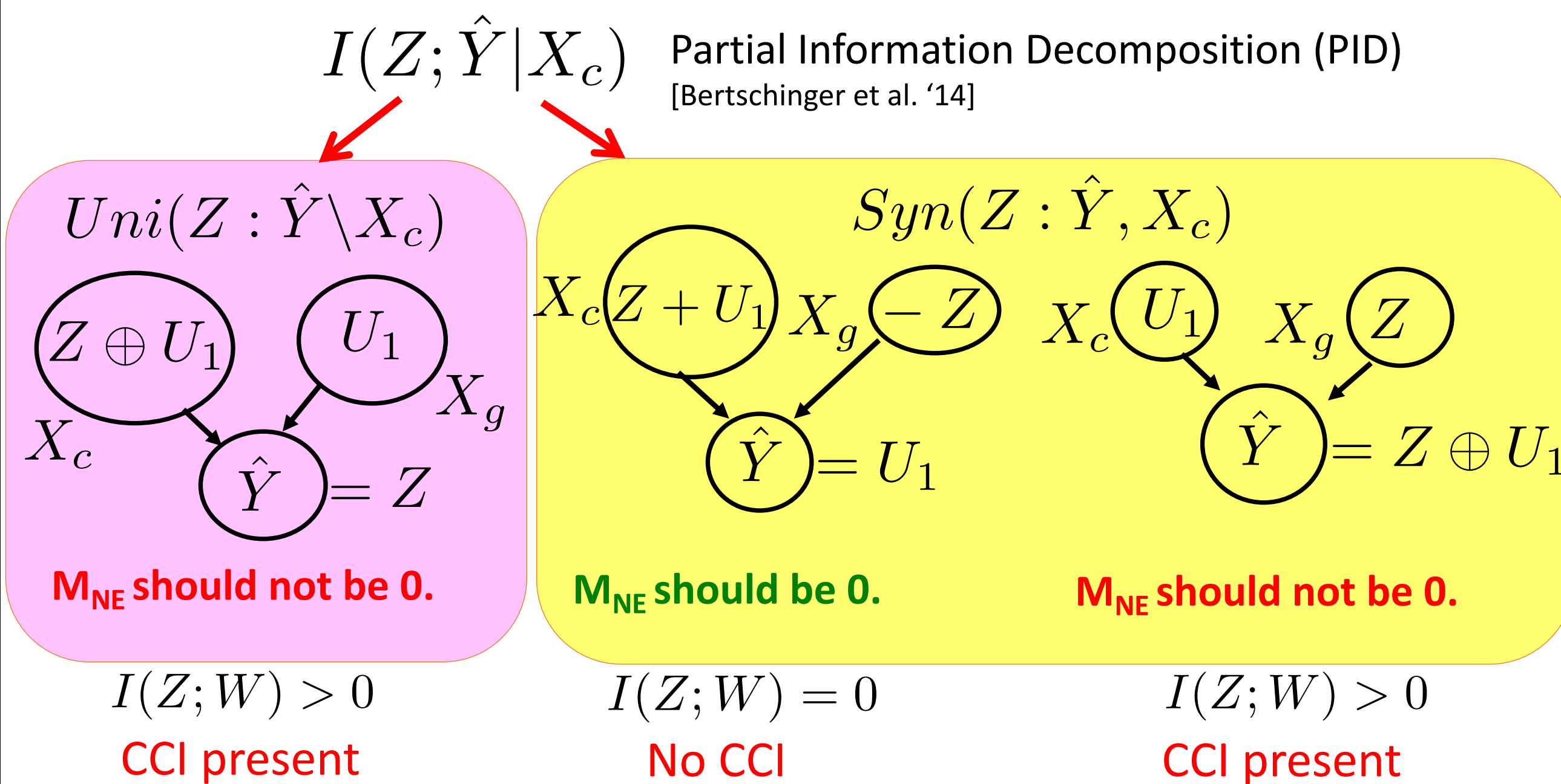


Candidate 4: Path-Specific Causal Influence

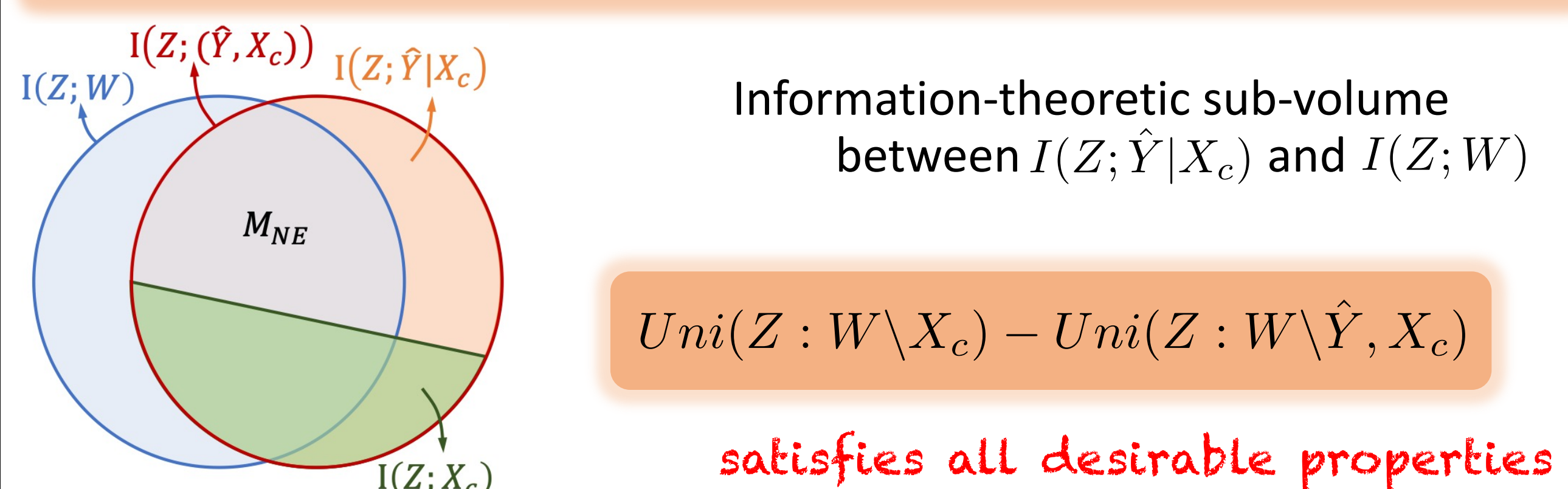
Counterexample: Synergy between critical and general



Understanding Scenarios where $I(Z; \hat{Y}|X_c) > 0$



Proposed Counterfactual Measure of Non-Exempt Discrimination



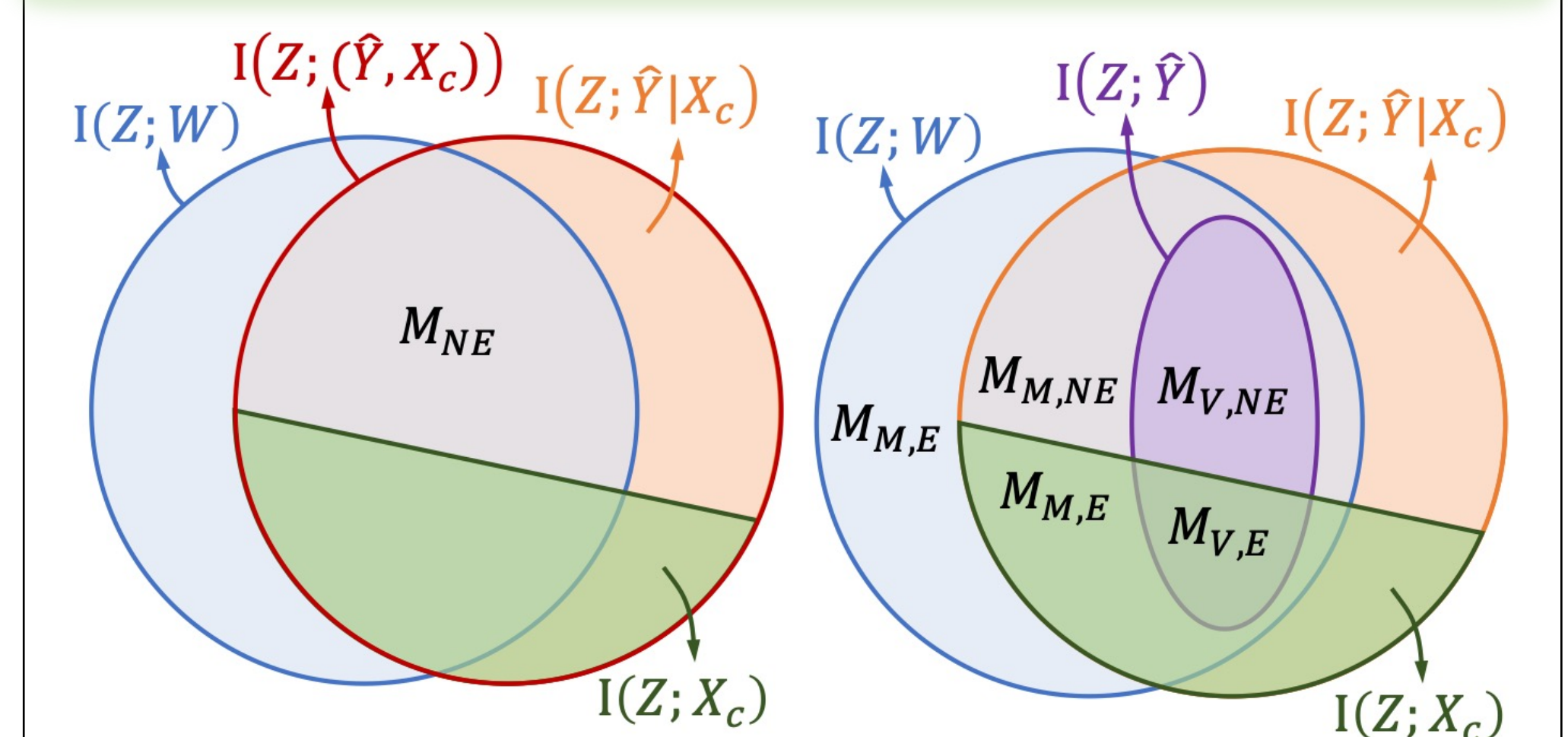
Main Results

Desirable Properties

- Property 1: M_{NE} should be 0 if $X_c = X$.
- Property 2: M_{NE} should be greater than 0 if $Uni(Z; \hat{Y} \setminus X_c) > 0$.
- Property 3: M_{NE} should be greater than 0 in the canonical example of masked discrimination.
- Property 4: M_{NE} should be 0 if $I(Z; W) = 0$.

Theorem 1: Our proposed measure satisfies all four properties.

Theorem 2: Total discrimination $I(Z; W)$ decomposes into four non-negative components: visible exempt, visible non-exempt, masked exempt, masked non-exempt.



Observational Measures: Impossibility, Utility, Limitation

Theorem 3: No observational measure can satisfy Properties 3 and 4 together.

- $Uni(Z; \hat{Y} \setminus X_c)$: Satisfies all properties except the property of non-exempt masked discrimination (Prop. 3).
- $I(Z; \hat{Y}|X_c)$: Captures masked discrimination but gives false positives under cancellation (Prop. 4).
- $I(Z; \hat{Y}|X_c, X')$: Satisfies only Prop. 1, others are satisfied partially with some counterexamples.

A Simple Case Study

Goal: Decide whether to show ads for an editorial job
Protected attribute Z: native English speaker or not

- X_1 : a score based on online writing samples **Critical**
- X_2 : a score based on browsing history, e.g., interest in English websites as compared to websites of other languages
- X_3 : a preference score based on geographical proximity.

$X_1 = Z + U_1$ $X_2 = Z + U_2$ $X_3 = U_3$

Ground Truth: $Y = \mathcal{I}(S \geq 1)$ where $S = X_1 + X_2 + X_3$

$Z \sim \text{Bern}(1/2)$ $U_1, U_2, U_3 \sim \mathcal{N}(0, 1)$

Equal Weight for all Features

Histograms of predicted scores for all candidates (Red: Z=0, Blue: Z=1)

Histogram of predicted scores for those satisfying critical necessity (Red: Z=0, Blue: Z=1)

