

An Information-Theoretic Quantification of Discrimination with Exempt Features

Sanghamitra Dutta, Praveen Venkatesh, Piotr Mardziel,
Anupam Datta, Pulkit Grover
Carnegie Mellon University

February 4, 2020

Abstract

The needs of a business (e.g., hiring) may require the use of certain features that are critical in a way that any discrimination arising due to them should be exempted. In this work, we propose a novel information-theoretic decomposition of the total discrimination (in a counterfactual sense) into a non-exempt component, which quantifies the part of the discrimination that cannot be accounted for by the critical features, and an exempt component, which quantifies the remaining discrimination. Our decomposition enables selective removal of the non-exempt component if desired. We arrive at this decomposition through examples and counterexamples that enable us to first obtain a set of desirable properties that any measure of non-exempt discrimination should satisfy. We then demonstrate that our proposed quantification of non-exempt discrimination satisfies all of them. This decomposition leverages a body of work from information theory called Partial Information Decomposition (PID). We also obtain an impossibility result showing that no observational measure of non-exempt discrimination can satisfy all of the desired properties, which leads us to relax our goals and examine alternative observational measures that satisfy only some of these properties. We then perform a case study using one observational measure to show how one might train a model allowing for exemption of discrimination due to critical features.

1 Introduction

As artificial intelligence becomes ubiquitous, it is important to understand whether a machine-learned model is perpetuating existing biases, and if so, how we can engineer fairness into such a model. The field of fair machine learning [1–13] provides several measures of fairness, and uses them to reduce discrimination based on *protected attributes*, e.g., as a regularizer during training [2, 5].

In particular applications, there are some features that are *critical* in a way that they are required to be weighed strongly in the decision *even if* they perpetuate bias, e.g., educational qualification for a job, merit and seniority in deciding salary etc. [14]. Hence, the discrimination arising due to these features can be exempted. In this work, our goal is to formalize and quantify the *non-exempt discrimination*, i.e., the part of the discrimination that cannot be accounted for by critical features, and selectively remove it if desired.

While such categorization of features is application-dependent and might require domain knowledge and ethical evaluation, such exemptions do exist. E.g., the US Equal Pay Act ([15]) exempts for any difference in salary based

The authors are with the Department of Electrical and Computer Engineering, Carnegie Mellon University. Author Contacts: S. Dutta (sanghamd@andrew.cmu.edu), P. Venkatesh (vpraveen@cmu.edu), P. Mardziel (piotrm@cmu.edu), A. Datta (danupam@cmu.edu), P. Grover (pulkit@cmu.edu)

Part of this work is accepted at AAAI 2020

Table 1: Observational Measures (M_{NE}) of Non-Exempt Discrimination (Utility and Limitations)

Desirable Properties	$\text{Uni}(Z : \hat{Y} X_c)$	$\text{I}(Z; \hat{Y} X_c)$	$\text{I}(Z; \hat{Y} X_c, X')$
1. Complete exemption if $X_c = X$.	Yes	Yes	Yes
2. Detects unique information about Z in \hat{Y} not in X_c .	Yes	Yes	Not Always
3. Detects Non-Exempt Masked Discrimination.	No	Masked by $g(X_c)$	Masked by $g(X_c, X')$
4. No causal influence from Z to $\hat{Y} \Rightarrow M_{NE} = 0$.	Yes	Not Always	Not Always

on gender that can be explained by merit and seniority. Similarly, the US employment discrimination law [16] contains a Bona Fide Occupational Qualification (BFOQ) defense where discrimination based on protected attributes may be exempted if the discrimination is due to a BFOQ reasonably necessary to the normal operation of that particular business, or other reasonable differentials. E.g., fire departments may require firemen to be able to lift a given weight to demonstrate that they will be able to carry fire victims out of a burning building. This feature is therefore allowed to be weighed strongly in hiring even if it is correlated with protected attributes. Similarly, UK employment discrimination law also allows exemptions based on the privacy and decency of the people the employer would be dealing with, e.g., staff in a care home [17].

In this work, we assume that the critical features are known (similar to [14], [18], [19]). Let X_c and X_g denote the critical and the non-critical or general features, respectively. We denote the protected attribute(s) as Z and the model output as \hat{Y} . Note that \hat{Y} is a function of the entire feature vector $X = (X_c, X_g)$.

Why should a model use the “general” features at all for prediction if they are not critical? The general features can improve accuracy, or reduce the candidate pool, e.g., if 60% of applicants clear a test but resources are available to interview only 10%. Not using the general features at all may reduce accuracy or produce a very large candidate pool. Our goal is to use both critical and general features in a way that maximizes accuracy (to the extent possible) while preventing only the non-exempt discrimination.¹

In this work, our contributions are as follows:

1. Quantification of Non-Exempt Discrimination: As a first step towards this quantification, we propose an information-theoretic quantification of the total discrimination (exempt and non-exempt) that is 0 if and only if the “counterfactual causal influence” [21] is 0, i.e., the model is *counterfactually fair*. Intuitively speaking, we extend the idea of “proxy-use” [22] from white-box models to black-box models, where we regard a model as being discriminatory if a virtual component (P) is formed inside the model that has high mutual information about Z (i.e., P is a virtual proxy of Z) and that also causally influences the final output \hat{Y} . Interestingly, note that this discrimination may not exhibit itself entirely in $\text{I}(Z; \hat{Y})$, which is the “statistically visible” information about Z in \hat{Y} because of “statistical masking effects,” e.g., $\hat{Y} = P + G$ where $G \perp Z$.

Next, we quantify the *non-exempt* part of this discrimination. Our quantification leverages a body of work in information theory called Partial Information Decomposition (PID). We consider examples and thought experiments to arrive at a set of desirable properties that any measure of non-exempt discrimination should satisfy, and then provide a measure that satisfies them (see Theorem 1). First, we require the measure to be 0 if all the features are in the exempt set X_c . Next, it is desirable that the measure be non-zero if \hat{Y} has any “unique” information about Z that is not present in X_c because then that information content is also attributed to X_g . However, because of statistical masking effects, even if this unique information is 0, there may still be non-exempt masked discrimination. Lastly, the measure should not capture false positives, e.g., it should be 0 if such virtual proxies cancel each other such that

¹Example (inspired by [20]): To choose a “good” employee, an employer could evaluate standardized test scores and reference letters (human-graded performance reviews). Both features are “job-related” in that they have statistical correlation with the prediction goal and can help improve accuracy. However, test-scores, a critical feature, should be weighed strongly in the decision *even if* biased whereas reference letters may be used only to the extent that they do not discriminate.

the final model output has no counterfactual causal influence of Z .

2. Decomposition of Total Discrimination: Next, we propose the decomposition of total discrimination into four non-negative components, namely, exempt and non-exempt visible discrimination and exempt and non-exempt masked discrimination (see Theorem 2).

3. An Impossibility Result: We show that no purely observational measure of non-exempt discrimination can satisfy all our desirable properties (see Theorem 3).

4. Observational Relaxations: Relaxing our requirements, we obtain purely observational measures that satisfy some of the desirable properties (summarized in Table 1) and then use one of them, namely, conditional mutual information, to demonstrate how to selectively reduce non-exempt discrimination in practice through a case study.

Related Work: We are aware that the idea of using conditional mutual information as a metric for non-exempt discrimination has surfaced in another work [23], where the focus is on conditional debiasing of neural networks using novel estimators. Other observational measures of non-exempt discrimination have also been discussed in [14], [24], [19], [25]. In this work, our focus is on an axiomatic examination of such measures and their relationship with the concept of *counterfactual fairness*² which has not received detailed attention. We also examine and acknowledge the utility and limitations of our observational measures (e.g., see an impossibility result in Theorem 3).

Causal approaches for fairness have been explored in [21], [18], [26], [27], [22], including impossibility results on purely observational measures [18, 22]. The main novelty arises from our adoption of a proxy-use viewpoint for black-box models *that allows for feature exemptions*. The decomposition of total discrimination into exempt and non-exempt components is tricky: one might be tempted to examine specific causal paths from Z to \hat{Y} that pass (or do not pass) through X_c , and deem those influences as the two measures. However, as the PID literature notes, discrimination can also arise from synergistic information [28–30] about Z in both X_c and X_g , that cannot be attributed to any one of them alone, *i.e.*, $I(Z; X_c)$ and $I(Z; X_g)$ may both be 0 but $I(Z; X_c, X_g)$ may not (see Counterexample 3). Purely causal measures (that do not rely on the PID framework) can attribute such discrimination entirely to X_c . We contend that such synergistic information, if influencing the decision, must be included in the *non-exempt* component of discrimination because, operationally, both X_c and X_g are contributors. We also note that identifying synergy is important: synergy arises frequently in machine-learning [31].

In a sense, this work treads a middle ground between two schools of thought, namely, *demographic parity* [2, 10], which enforces the criterion $Z \perp \hat{Y}$, and *equalized odds* [3, 10], which enforces $Z \perp \hat{Y} | Y$ (directly or through practical relaxations) where Y denotes the true labels of the historic dataset. Our selective quantification of non-exempt discrimination helps address one of the major criticisms against demographic parity, namely, that it can deliberately choose unqualified members from the protected group [32], e.g., by disregarding the critical features if they are correlated with Z . Another strength of our approach is that it does not use the true labels for fairness (unlike equalized odds). The use of true labels has been criticized in [20] because “often the best labels for different classifications will be open to debate,” e.g., if the labels themselves are biased. This work also shares intellectual connections with *individual fairness* [1] in the sense that it enables individuals with similar X_c to be treated similarly, if desired.

Background on Partial Information Decomposition (PID): Here, we provide a brief background on the PID framework [29, 30] to help follow this paper. Appendix A provides more details and specific properties used in the proofs.

The PID framework decomposes the mutual information $I(Z; (A, B))$ about a random variable Z contained in

²Our measure of total (exempt and non-exempt) discrimination is zero if and only if the “counterfactual causal influence” of Z on \hat{Y} is zero (see Lemma 1).

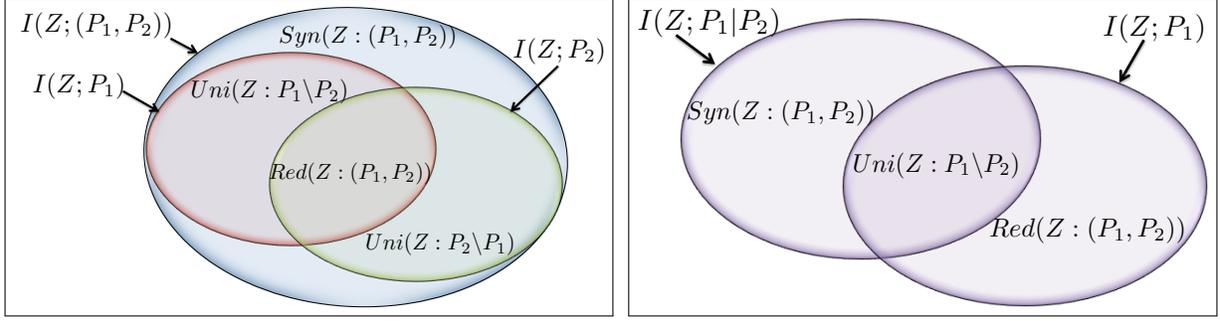


Figure 1: (Left) Mutual information $I(Z; (P_1, P_2))$ is decomposed into 4 non-negative terms, namely, $Uni(Z : P_1 \setminus P_2)$, $Uni(Z : P_2 \setminus P_1)$, $Red(Z : (P_1, P_2))$ and $Syn(Z : (P_1, P_2))$. (Right) Relation between $I(Z; P_1|P_2)$ and $I(Z; P_1)$ showing that $Uni(Z : P_1 \setminus P_2)$ is the common component between the two.

the tuple (A, B) into four *non-negative* terms as follows:

$$I(Z; (A, B)) = Uni(Z : A \setminus B) + Uni(Z : B \setminus A) + Red(Z : (A, B)) + Syn(Z : (A, B)). \quad (1)$$

Here, $Uni(Z : A \setminus B)$ denotes the unique information about Z that is present only in A and not in B . Likewise, $Uni(Z : B \setminus A)$ is the unique information about Z that is present only in B and not in A . $Red(Z : (A, B))$ denotes the redundant information about Z , present in both A and B , and $Syn(Z : (A, B))$ denotes the synergistic information not present in either of A or B individually, but present jointly in (A, B) (see Fig. 1 for illustrations).

Example 1 (Partial Information Decomposition). Let $Z = (Z_1, Z_2, Z_3)$, $Z_i \sim i.i.d. \text{Bern}(1/2)$. Let $A = (Z_1, Z_2, Z_3 \oplus N)$ where \oplus denotes XOR, $B = (Z_2, N)$, and $N \sim \text{Bern}(1/2)$ is independent of Z . Here, $I(Z; (A, B)) = 3$ bits.

Observe that, the unique information about Z that is contained only in A and not in B is effectively contained in Z_1 and is given by $Uni(Z : A \setminus B) = I(Z; Z_1) = 1$ bit. The redundant information about Z that is contained in both A and B is effectively contained in Z_2 and is given by $Red(Z : (A, B)) = I(Z; Z_2) = 1$ bit. Lastly, the synergistic information about Z that is not contained in either A or B alone, but is contained in both of them together is effectively contained in the tuple $(Z_3 \oplus N, N)$, and is given by $Syn(Z : (A, B)) = I(Z; (Z_3 \oplus N, N)) = 1$ bit. This accounts for the 3 bits in $I(Z; (A, B))$. Here, B does not have any unique information about Z that is not contained in A .

Existing literature suggests different definitions for the individual PID terms [29, 30]. However, irrespective of the exact definition of these individual terms, the following identities always hold (for all the definitions):

$$I(Z; A) = Uni(Z : A \setminus B) + Red(Z : (A, B)). \quad (2)$$

$$I(Z; A | B) = Uni(Z : A \setminus B) + Syn(Z : (A, B)). \quad (3)$$

Given the three independent equations (1), (2) and (3) in four unknowns (the four PID terms), defining one of the terms (e.g., $Uni(Z : A \setminus B)$) is sufficient to obtain the other three. For completeness, we include the definition of unique information from [29] (that also allows for estimation via convex optimization [33]). The author is referred to [29] for more details and insights on this particular definition. To follow our paper, only an intuitive understanding of the concept of unique information is sufficient.

Definition 1 (Unique Information). [29] Let Δ be the set of all joint distributions on (Z, A, B) and Δ_p be the set of joint distributions with the same marginals on (Z, A) and (Z, B) as their true distribution, i.e., $\Delta_p = \{Q \in \Delta : q(z, a) = \Pr(Z=z, A=a) \text{ and } q(z, b) = \Pr(Z=z, B=b)\}$ Then, $Uni(Z : A \setminus B) = \min_{Q \in \Delta_p} I_Q(Z; A|B)$.

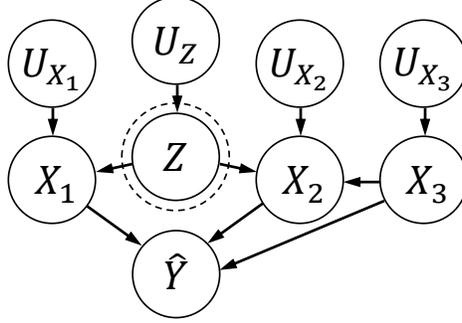


Figure 2: An SCM with protected attribute Z , features $X = \{X_1, X_2, X_3\}$, and output \hat{Y} . Z does not have any parents and \hat{Y} is completely determined by $\{X_1, X_2, X_3\}$.

Definition 1 uniquely defines $\text{Red}(Z : (P_1, P_2))$ and $\text{Syn}(Z : (P_1, P_2))$ using (2) and (3). The key intuition behind this definition is that unique and synergistic information should only depend on the marginal distribution of the pairs (Z, P_1) and (Z, P_2) . This is motivated from an operational perspective that if P_1 has unique information about Z (with respect to P_2), then there must be a situation where P_1 can use this information to perform better at predicting Z than P_2 (see also [29]).

2 System Model and Assumptions

Definition 2 (Structural Causal Model: $\text{SCM}(U, V, \mathcal{F})$). A structural causal model (U, V, \mathcal{F}) consists of a set of latent (unobserved) and mutually independent variables U which are not caused by any variable in the set of observable variables V , and a collection of deterministic functions (structural assignments) $\mathcal{F} = \{f_1, f_2, \dots\}$, one for each $V_i \in V$, such that: $V_i = f_i(V_{pa_i}, U_i)$. Here $V_{pa_i} \subseteq V \setminus V_i$ are the parents of V_i , and $U_i \subseteq U$. The structural assignment graph (SAG) of $\text{SCM}(U, V, \mathcal{F})$ has one vertex for each V_i , and directed edges to V_i from each parent in V_{pa_i} , and is always a directed acyclic graph.

For our problem, the latent variables U represent possibly unknown social factors. The observables V consist of the protected attributes Z , the features $X = \{X_c, X_g\}$ and the output \hat{Y} (see Fig. 2). For simplicity, we assume ancestral closure of the protected attributes, *i.e.*, the parents of any $V_i \in Z$ also lie in Z and hence Z is not caused by any of the features in X ($V_i \in Z$ are source nodes in the SAG). Therefore, $Z = f_z(U_Z)$ for $U_Z \subseteq U$. Any feature X_j in X is a function of its corresponding latent variable and its parents, which are again functions of their own latent variables and parents. Note that, X can also be written as $f(Z, U_X)$ for some deterministic $f(\cdot)$, where $f(\cdot)$ may be constant in some of its arguments, and $Z \perp\!\!\!\perp U_X$ (see [34, Proposition 6.3]). This holds because the underlying graph is acyclic. A model takes $X = \{X_c, X_g\}$ as its input and produces an output \hat{Y} which depends only on X . Therefore, $\hat{Y} = h(Z, U_X)$ for some function $h(\cdot)$.

For completeness, we define Counterfactual Causal Influence (CCI) inspired from [21], [26], [35], [36], [37], [38], [39].

Definition 3 (Counterfactual Causal Influence: $\text{CCI}(Z \rightarrow \hat{Y})$). If $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$ where U_X are latent variables that do not cause Z in the true SCM, and Z', Z are *i.i.d.*, then

$$\text{CCI}(Z \rightarrow \hat{Y}) = \mathbb{E}_{Z, Z', U_X} [|h(Z, U_X) - h(Z', U_X)|].$$

Remark 1. Statistical independence does not imply absence of causal effects. *E.g.*, $\hat{Y} = Z \oplus U_X$ where $Z, U_X \sim \text{i.i.d. Bern}(1/2)$. Here, $\hat{Y} \perp\!\!\!\perp Z$, but Z still has a causal effect on \hat{Y} . If we vary Z while fixing all other sources of

randomness in \hat{Y} as constants (i.e., fixing $U_X = u_x$), then \hat{Y} also varies. This is in fact an example of masked discrimination, where $I(Z; \hat{Y}) = 0$, but Z causally influences \hat{Y} .

Next, we define a variable W as follows:

Definition 4 (Variable W). We define a variable $W = [h(Z, u_x^{(1)}), \dots, h(Z, u_x^{(k)})]$, where $\{u_x^{(1)}, \dots, u_x^{(k)}\}$ is the set of all values with $\Pr(U_X = u_x) > 0$.

Here, W is a deterministic function of Z alone, consisting of all the functional forms that $\hat{Y} = h(Z, U_X)$ takes for all values u_x attainable by U_X .

Lemma 1 (Information-Theoretic Equivalent of CCI). Let $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$. Then $\text{CCI}(Z \rightarrow \hat{Y}) \neq 0$ if and only if $I(Z; W) > 0$.

The proof is provided in Appendix B.1.

Remark 2. We also show that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ (or, $I(Z; W) = 0$) is equivalent to the counterfactual fairness criterion of [21] (proved in the Appendix B.2). Therefore, in this work, we will regard $I(Z; W)$ as an **information-theoretic quantification of the total discrimination (exempt and non-exempt)**.

3 Main Results

We formally state the desirable properties, intuitively stated in Section 1, and then introduce our proposed measure that satisfies all of them (Theorem 1 in Section 3.1). While the proof is presented in the Appendix C.1, in Section 3.2 we present the main intuition behind our proposed measure through several examples, counterexamples and thought experiments, that also help us arrive at the desirable properties. Our proposed measure leads to a non-negative decomposition of the total discrimination $I(Z; W)$ into four components, i.e., statistically visible and masked portions, each with exempt and non-exempt components (see Section 3.3). Lastly, in Section 3.4, we demonstrate how to modify our measure to account for other kinds of masked discrimination under different sociological contexts.

3.1 Desirable Properties and Proposed Measure

We introduce a set of desirable properties for any measure of non-exempt discrimination (M_{NE}). Firstly, we require the measure to be 0 if all the features are in the exempt set X_c :

Property 1 (Complete Exemption). M_{NE} should be 0 if all features are categorized into X_c , i.e., $X_c = X$ and $X_g = \phi$.

Next, it is desirable that the measure be non-zero if \hat{Y} has any unique information about Z that is not present in X_c because then that information is also attributed to X_g .

Property 2 (Non-Exempt Visible Discrimination). M_{NE} should be strictly greater than 0 if $\text{Uni}(Z : \hat{Y} \setminus X_c) > 0$.

However, as we discussed in Section 2, statistical masking can sometimes prevent the entire non-exempt discrimination component from exhibiting itself in $\text{Uni}(Z : \hat{Y} \setminus X_c)$. As an extreme example, consider the following scenario.

Example 2. Let $\hat{Y} = Z \oplus f(U_X)$ for some function $f(\cdot)$ on $X_c = U_X$ with Z and $f(U_X)$ being i.i.d. $\text{Bern}(1/2)$. E.g., an ad for expensive housing is presented to white people ($Z = 1$) with income above a threshold ($f(U_X) = 1$), and also to black people ($Z = 0$) with income below a threshold ($f(U_X) = 0$) (while being largely irrelevant to the latter).

Not all forms of masked effects are undesirable. An example is if the only available features are $X_g = (Z, U_X)$, where Z is the race and U_X is $\text{Bern}(1/2)$, a random coin flip. Then, performing $\hat{Y} = Z \oplus U_X$ randomizes the race, and can be a preventive measure against discrimination even if $\text{CCI}(Z \rightarrow \hat{Y}) > 0$. In the following property, we will assume that the discrimination (masked/unmasked) is exempt if the Markov chain $Z - X_c - \hat{Y}$ holds. This property only accounts for masking that is entirely due to X_c , e.g., $\hat{Y} = Z + f(X_c)$ for some function $f(\cdot)$ where $\text{CCI}(Z \rightarrow f(X_c)) = 0$ and exempts other forms of masking (revisited in Remark 3).

Property 3 (Non-Exempt Masking). *A measure M_{NE} should be non-zero in the canonical example of masked discrimination, i.e., Example 2 even if $I(Z; \hat{Y}) = 0$. However, M_{NE} should be 0 if $Z - X_c - \hat{Y}$ form a Markov chain.*

Remark 3. *In general, one might also choose to consider a subset of latent factors $\tilde{U} \subseteq U_X$ such that any statistical masking arising due to these latent variables is also undesirable. Then, the Markov chain in Property 3 may be modified to $Z - X_c - (\hat{Y}, \tilde{U})$, and the proposed measure can be modified accordingly, as also elaborated further in Section 3.4.*

Lastly, the measure should also not capture false positives, e.g., it should be 0 if such virtual proxies cancel each other causing the final model output to have no counterfactual causal influence of Z , leading to the following property.

Property 4 (Cancellation of Influence). *M_{NE} should be 0 if $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ (or equivalently, $I(Z; W) = 0$).*

Now, we introduce our proposed measure and then show that it satisfies all these desirable properties (see Theorem 1).

Definition 5 (Non-Exempt Discrimination). *Our proposed measure of non-exempt discrimination is given by:*

$$M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c). \quad (4)$$

Remark 4. *The proposed measure is essentially the volume of the overlap between $I(Z; W)$ and $I(Z; \hat{Y} | X_c)$, that becomes 0 when either of them is 0 (see Fig. 3).*

Theorem 1 (Properties). *Properties 1, 2, 3 and 4 are satisfied by $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c)$.*

The proof is provided in Appendix C.1. The next result provides an equivalent definition of non-exempt discrimination (based on the definition of unique information proposed in [29]).

Lemma 2 (Non-Exempt Discrimination Equivalence). *The proposed measure $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c)$ is equal to $I(W; \hat{Y} | X_c)$.*

The proof is provided in Appendix C.1.

3.2 Main Intuition behind the Proposed Measure

We examine some candidate measures (M_{NE}) of non-exempt discrimination through examples and counterexamples, leading to our proposed measure.

Candidate Measure 1. $M_{NE} = I(Z; \hat{Y})$.

Counterexample 1. Let $X_c = Z + U_{X_1}$ where $Z \sim \text{Bern}(1/2)$, $U_{X_1} \sim \mathcal{N}(0, \sigma_1^2)$ and $X_g = U_{X_2}$ where $U_{X_2} \sim \mathcal{N}(0, \sigma_2^2)$ and is independent of U_{X_1} . The decision of the model is $\hat{Y} = \text{sgn}(X_c + X_g - 0.5) = \text{sgn}(Z + U_{X_1} + U_{X_2} - 0.5)$.

Here, $I(Z; \hat{Y})$ is non-zero and so is $I(Z; X_c)$. However, $I(Z; \hat{Y} | X_c) = 0$ (see the Markov chain $Z - X_c - \hat{Y}$). The information that \hat{Y} contains about Z is redundant information also contained in X_c . Therefore, the discrimination here should be exempted because it arises entirely from X_c .

Candidate Measure 2. $M_{NE} = I(Z; \hat{Y} | X_c)$.

This measure resolves Counterexample 1. It also has some provision for selectively capturing the non-exempt component: it is 0 in Counterexample 1, consistent with the intuition that there is no non-exempt discrimination. However, the following example exposes some of its limitations.

Counterexample 2 (Cancellation of Influence). Let $X_c = Z + U_X$ and $X_g = Z$ where Z denotes the gender and U_X denotes the student's knowledge. The model's decision on a student's ability is $\hat{Y} = X_c - X_g = U_X$.

The influences of Z along two different causal paths cancel each other in the final output, so that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ (and, $I(Z; W) = 0$). Thus, there is no discrimination in the outcome \hat{Y} (this is true even if the features in X_c were not exempt; see Remark 2). However, the measure $M = I(Z; \hat{Y} | X_c)$ is positive for this example, leading to a false positive in detecting discrimination. These two examples serve as our motivation behind Properties 1 and 4. The next candidate resolves both these examples.

Candidate Measure 3. $M_{NE} = \text{Uni}(Z : \hat{Y} \setminus X_c)$.

This measure resolves Counterexample 1: \hat{Y} and X_c have redundant information about Z , but there is no unique information about Z in \hat{Y} that is not in X_c . Thus $\text{Uni}(Z : \hat{Y} \setminus X_c) = 0$, consistent with the conclusion that the discrimination in Counterexample 1 should be exempt. $\text{Uni}(Z : \hat{Y} \setminus X_c)$ is also 0 in Counterexample 2. In fact, $\text{Uni}(Z : \hat{Y} \setminus X_c)$ captures the non-exempt discrimination that is statistically visible in $I(Z; \hat{Y})$, leading to Property 2.

Counterexample 3 (Masked Discrimination). Refer to Example 2 in Section 3.1 where $\hat{Y} = Z \oplus f(X_c)$.

Here $Z \perp \hat{Y}$, i.e., $I(Z; \hat{Y}) = 0$, making the model “appear to have no discrimination.” However, when examined more deeply, the model racially discriminates against half of the population (high-income black people) for whom the ad is relevant. This is also demonstrated by the fact that $\text{CCI}(Z \rightarrow \hat{Y}) \neq 0$ and the Markov chain $Z - X_c - \hat{Y}$ does not hold. $\text{Uni}(Z : \hat{Y} \setminus X_c)$ fails to capture such *non-exempt masked discrimination*. In fact, this example motivates Property 3. $\text{Uni}(Z : \hat{Y} \setminus X_c)$ does not satisfy this property as it has to be zero whenever $I(Z; \hat{Y}) = 0$.

Inspired from $\text{CCI}(Z \rightarrow \hat{Y})$, another possible candidate for quantifying non-exempt discrimination is a causal, path-specific examination (see also [27], [21], [18]) by varying Z only along the direct paths through X_g and comparing if it causes any difference in the decision.

Candidate Measure 4. Let $\hat{Y} = h(Z, U_X)$ in the true causal model. Assume a new causal graph with a new source node Z' having an independent and identical distribution as Z where we replace all direct edges from Z to X_g with an edge from Z' to X_g . Let $\tilde{h}(Z, Z', U_X)$ be the model output in the new causal graph. A candidate measure is $M_{NE} = \mathbb{E}_{Z, Z', U_X} [|h(Z, U_X) - \tilde{h}(Z, Z', U_X)|]$.

Counterexample 4 (Non-zero Unique Information). Suppose that $X_c = Z \oplus U_{X_1}$ and $X_g = U_{X_1}$ where Z and U_{X_1} are i.i.d. $\text{Bern}(1/2)$. Let $\hat{Y} = X_c \oplus X_g = Z$.

In this example, \hat{Y} has unique information about Z that is not contained in X_c , implying non-exempt visible discrimination. However, a path-specific examination would conclude that the causal influence of Z is only propagating through X_c , and hence should be exempt. Following the PID literature, here \hat{Y} receives synergistic information about Z from both X_c and X_g , that cannot be attributed to X_c alone ($I(Z; X_c) = 0$). From an operational perspective, \hat{Y} and X_c together lead to a better estimate of Z than X_c alone which means X_g is definitely a contributor to the discrimination, and thus $M_{NE} > 0$. We therefore seek a measure under which such discrimination qualifies as non-exempt. Motivated by this example, we now consider another candidate measure that is derived from $I(Z; W)$.

Candidate Measure 5. $M_{NE} = \text{Uni}(Z : W \setminus X_c)$.

While this measure resolves all the examples so far, it may not always satisfy Property 1.

Counterexample 5. Suppose that $X = X_c = Z \oplus U_X$, and $\hat{Y} = X_c = Z \oplus U_X$.

In this scenario, this measure is not 0 even though the discrimination is completely exempt. This motivates our proposed measure $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c)$, which accounts for, and effectively removes, such exempt components from $\text{Uni}(Z : W \setminus X_c)$, and finally satisfies all the desirable properties.

M_{NE} being non-zero actually implies that both $I(Z; W) > 0$ and $I(Z; \hat{Y}|X_c) > 0$ (overlapping volume). However, this is only a one-way implication. $I(Z; W)$ and $I(Z; \hat{Y}|X_c)$ both being non-zero does not necessarily capture non-exempt discrimination.

Example 3. Let $Z = (Z_1, Z_2)$, $X_c = (Z_1 \oplus U_{X_1}, Z_2)$, $X_g = (Z_1, U_{X_2})$ and $\hat{Y} = (U_{X_1}, Z_2 \oplus U_{X_2})$ where $Z_1, Z_2, U_{X_1}, U_{X_2}$ are i.i.d. $\text{Bern}(1/2)$.

This example should be exempt because Z_2 already appears in X_c , and is hence exempt. Our proposed measure also suggests the same conclusion. However, both $I(Z; W)$ and $I(Z; \hat{Y}|X_c)$ are non-zero here.

3.3 Understanding the Overall Decomposition

This work enables an information-theoretic decomposition of the total discrimination $I(Z; W)$ into non-exempt and exempt components, namely, M_{NE} and $I(Z; W) - M_{NE}$ respectively. Alongside, $I(Z; W)$ can also be decomposed into statistically visible and masked components, namely, $I(Z; \hat{Y})$ and $I(Z; W) - I(Z; \hat{Y})$ respectively. Combining these two decompositions leads to an overall four-way decomposition of $I(Z; W)$ as shown in Theorem 2 (see Fig. 3).

Theorem 2 (Non-negative Decomposition of Total Discrimination). *The total discrimination can be decomposed into four non-negative components as follows:*

$$I(Z; W) = M_{V,NE} + M_{V,E} + M_{M,NE} + M_{M,E}. \quad (5)$$

Here $M_{V,NE} = \text{Uni}(Z : \hat{Y} \setminus X_c)$ is the visible, non-exempt component and $M_{V,E} = \text{Red}(Z : (\hat{Y}, X_c))$ is the visible, exempt component. These two terms add to form $I(Z; \hat{Y})$ which is the total statistically visible discrimination. Likewise, $M_{M,NE} = M_{NE} - M_{V,NE}$ is the masked, non-exempt component, and $M_{M,E} = I(Z; W) - I(Z; \hat{Y}) - M_{M,NE}$ is the masked, exempt component.

The proof is in Appendix C.2.

Lemma 3 (Masked Discrimination). *The total masked discrimination $I(Z; W) - I(Z; \hat{Y})$ is equal to $\text{Uni}(Z : W \setminus \hat{Y})$.*

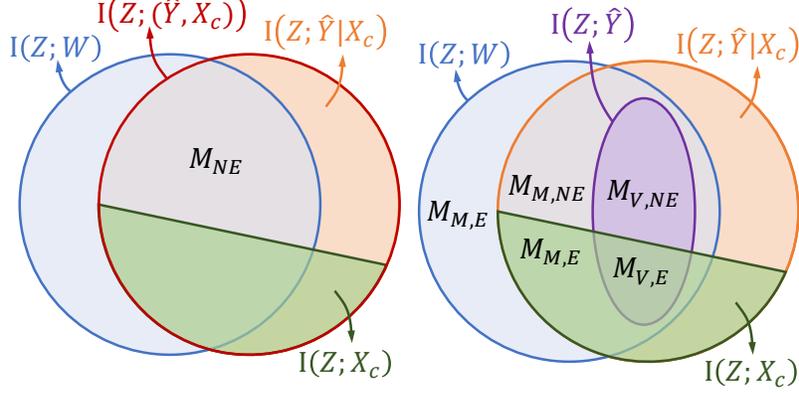


Figure 3: Information-theoretic decomposition of total discrimination, $I(Z; W)$: (Left) The red full-circle denotes $I(Z; (X_c, \hat{Y}))$ which is equal to $I(Z; X_c) + I(Z; \hat{Y} | X_c)$. Both $I(Z; X_c)$ and $I(Z; \hat{Y} | X_c)$ are denoted by sub-volumes within the red full-circle. The volume of overlap between $I(Z; W)$ and $I(Z; \hat{Y} | X_c)$ is our proposed measure of non-exempt discrimination M_{NE} . (Right) Note that, $I(Z; \hat{Y})$ (total statistically visible discrimination) is the purple circle that is entirely contained inside $I(Z; W)$ and $I(Z; \hat{Y} | X_c)$. This leads to a four-way decomposition of $I(Z; W)$: the visible non-exempt component $M_{V,NE} = \text{Uni}(Z : \hat{Y} \setminus X_c)$, the visible exempt component $M_{V,E} = \text{Red}(Z : (\hat{Y}, X_c))$, the masked non-exempt component $M_{M,NE} = M_{NE} - M_{V,NE}$, and the masked exempt component $M_{M,E} = I(Z; W) - I(Z; \hat{Y}) - M_{M,NE}$. Also note that $I(Z; \hat{Y})$ has an intersection with $I(Z; \hat{Y} | X_c)$, but both $I(Z; \hat{Y})$ and $I(Z; \hat{Y} | X_c)$ also have components (volumes) outside the intersection which allows either of them to be greater or less than the other in our Venn diagram.

The proof is provided in Appendix C.2.

Lemma 4 (Masked Discrimination Implications). *The following two statements are equivalent:*

- $I(Z; \hat{Y} | U_X) - I(Z; \hat{Y}) > 0$.
 - \exists a random variable G of the form $G = g(U_X)$ such that $I(Z; \hat{Y} | G) > I(Z; \hat{Y})$.
- Either of these statements imply $I(Z; W) - I(Z; \hat{Y}) > 0$.

The proof is provided in Appendix C.2.

3.4 Modifying the Proposed Measure to Account for More Masked Effects

Different forms of statistical masking can have different implications under different sociological contexts, e.g., $\hat{Y} = Z \oplus U_X$ may be undesirable if U_X is the income (recall Example 2) but not necessarily unfair if U_X is the random flip of a coin. In our proposed measure, we only accounted for statistical masking effects caused by the critical features X_c . However, there may be scenarios where we might want to capture masking effects by other variables also, e.g., X_g . Let us understand this using the following example.

Example 4. Let $X_c = (U_{X_1}, U_{X_2})$ and $X_g = (Z, U_{X_3})$, where all the latent random variables are i.i.d. $\text{Bern}(1/2)$. Now the output \hat{Y} can take different forms, such as $Z \oplus f_1(X_c) = Z \oplus U_{X_1}$, or $Z \oplus f_1(X_c) \oplus f_2(X_g) = Z \oplus U_{X_1} \oplus U_{X_3}$ or $Z \oplus f_2(X_g) = Z \oplus U_{X_3}$.

By our proposed measure, only $\hat{Y} = Z \oplus U_{X_1} \oplus U_{X_2}$ and $\hat{Y} = Z \oplus U_{X_1}$ are considered non-exempt. Masking by X_g (e.g., $\hat{Y} = Z \oplus U_{X_3}$) or masking by a combination of X_c and X_g (e.g., $\hat{Y} = Z \oplus U_{X_1} \oplus U_{X_3}$) is exempted ($Z - X_c - \hat{Y}$ is a Markov chain). Statistical masking of Z by $f_2(X_g)$ is viewed more like randomization, e.g., using

a coin flip to prevent discrimination, whereas masking by $f_1(X_c)$ is like discriminating against high-income black people (Example 2).

In general, which masking effects should be accounted for depends on the problem design. In some scenarios, one may be interested in not exempting masking effects due to some latent variables. Let $\tilde{U}_X \subseteq U_X$ be the set of latent random variables such that any statistical masking effect derived from them should be accounted for. Then, we may redefine Property 3 as follows: the measure $M'_{NE} = 0$ if $Z - X_c - (\hat{Y}, \tilde{U}_X)$ is a Markov chain. This leads to the following modified measure of non-exempt discrimination.

Definition 6 (Modified Non-Exempt Discrimination). $M'_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c, \tilde{U}_X)$.

This measure is the volume of overlap between $I(Z; W)$ and $I(Z; \hat{Y}, \tilde{U}_X | X_c)$. Using this measure in Example 4 leads to the conclusion that all the cases are non-exempt if \tilde{U}_X is chosen as $(U_{X_1}, U_{X_2}, U_{X_3})$. This unravels the statistical masking by $U_{X_1}, U_{X_2}, U_{X_3}$ and exposes the discriminatory component Z lying underneath. Again, in some examples, accounting for only some latent factors makes sense:

Example 5. Let $X_c = Z + U_{X_1}$ and $X_g = (U_{X_1}, U_{X_2})$ where all the latent variables are independent with $U_{X_1} \sim \mathcal{N}(0, 1000)$ and all others distributed as $\mathcal{N}(0, 1)$. The output \hat{Y} can take different forms, such as, $\hat{Y} = Z + U_{X_1}$, or $Z + U_{X_1} + U_{X_2}$, or $Z + U_{X_2}$.

When $\hat{Y} = Z + U_{X_1}$, the output is entirely derived from X_c and hence should be exempt. Here, $Z - X_c - \hat{Y}$ is a Markov chain but $Z - X_c - (\hat{Y}, U_{X_1})$ is not. For this example, it does not make sense to try to unravel masked effects of U_{X_1} over Z , or include it in \tilde{U}_X . When $\hat{Y} = Z + U_{X_1} + U_{X_2}$, it should also be exempt for the same reason. However, $\hat{Y} = Z + U_{X_2}$ is not necessarily exempt because it contains unique information about Z not present in X_c (X_g helps unmask and expose $Z + U_{X_2}$). Here, $Z - X_c - \hat{Y}$ is not a Markov chain. To unravel the masked effect caused by U_{X_2} and expose Z entirely, one may include it in \tilde{U}_X .

4 Observational Relaxations for Practical Application in Training

Theorem 3 (Impossibility of Observational Measures). *No observational measure of non-exempt discrimination can distinguish between Example 6, a case of no discrimination and Example 7, a case of non-exempt discrimination.*

Example 6. Let us assume that there exists a scenario where $X_c = Z \oplus U_{X_1}$, $X_g = Z$ and $\hat{Y} = X_c \oplus X_g = U_{X_1}$ where Z and U_{X_1} are both independent and identically distributed as $\text{Bern}(1/2)$.

Example 7. Let us assume that there exists another scenario where $X_c = U_{X_1}$, $X_g = Z$ and $\hat{Y} = X_c \oplus X_g = Z \oplus U_{X_1}$ where Z and U_{X_1} are both independent and identically distributed as $\text{Bern}(1/2)$.

In Example 6, the influences of Z cancel each other and there is no discrimination (Property 4). However, in Example 7, there is non-exempt masked discrimination (Property 3). But, for both these examples, the joint distribution of the observables (Z, X_c, X_g, \hat{Y}) is the same which means that no observational measure can distinguish between these two cases. This also completes the proof.

Nevertheless, because counterfactual measures are difficult to realize in practice, we examine the following observational measures of non-exempt discrimination that satisfy only a few of Properties 1-4.

1. $\text{Uni}(Z : \hat{Y} \setminus X_c)$: This measure satisfies Properties 1, 2 and 4 (proved in Appendix D.1). However, it does not quantify any masked discrimination.
2. $I(Z; \hat{Y} | X_c)$: This measure satisfies Properties 1, 2, and 3 (proved in Appendix D.1). However, it can lead to false positives for Property 4 (absence of $\text{CCI}(Z \rightarrow \hat{Y})$), e.g., in Counterexample 2.
3. $I(Z; \hat{Y} | X_c, X')$: X' consists of features of X_g suspected of masking Z . This is somewhat of a heuristic relaxation that only satisfies Property 1 but partly satisfies all the rest with some exceptions, *i.e.*, it exempts synergistic

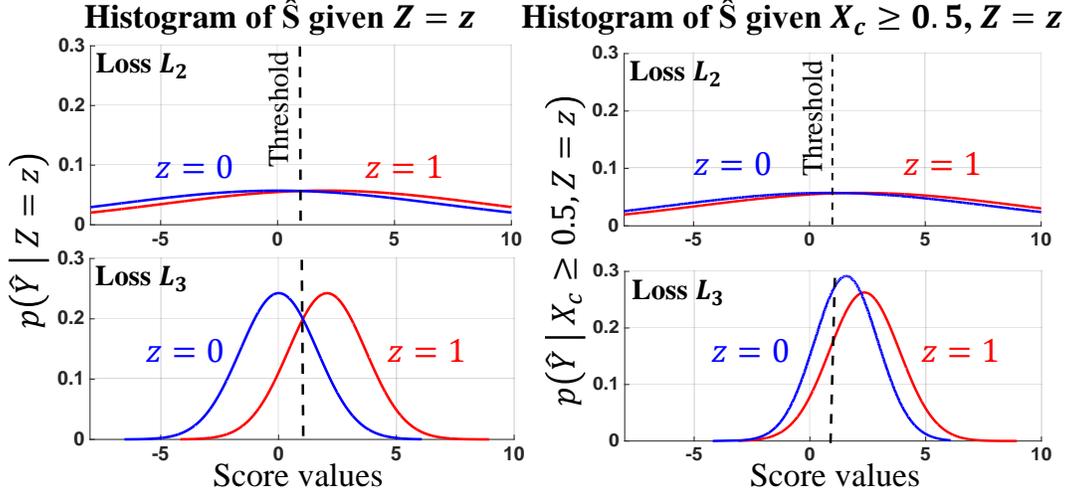


Figure 4: Histogram of Predicted Scores ($\hat{S} = -w^T X/b$): (Left) $p(\hat{S}|Z=i)$ for $i=0, 1$; (Right) $p(\hat{S}|X_c \geq 0.5, Z=i)$ for $i = 0, 1$. Regularizing with $I(Z; \hat{Y})$ (L_2) brings $p(\hat{S}|Z)$ closer for $Z=0$ and 1 by placing higher weight on a less important feature (proximity score). This increases the variance and reduces accuracy (see Table 2). Regularizing with $I(Z; \hat{Y}|X_c)$ (L_3) makes $p(\hat{S}|X_c \geq 0.5, Z)$ approach each other for $Z=0$ and 1 , aiming to give similar prediction scores to individuals with similar X_c ($\lambda=10$ for these plots).

Table 2: Observations after training a classifier ($w_1 X_1 + w_2 X_2 + w_3 X_3 + b \geq 0$) using three loss functions with different fairness criteria (100 simulations of 7000 iterations each with batch size 200).

Loss (λ)	$-\frac{w_1}{b}$	$-\frac{w_2}{b}$	$-\frac{w_3}{b}$	Acc%
$L_1 (-)$	1.08	1.08	1.08	98.5
$L_2 (4)$	1.07	1.07	3.76	81.1
$L_2 (10)$	1.01	1.03	13.9	70.2
$L_3 (4)$	1.46	0.73	1.91	89.6
$L_3 (10)$	2.05	0.02	2.57	80.8

information about Z in (X_c, X') that can show up in \hat{Y} , and cause non-zero $\text{Uni}(Z : \hat{Y} \setminus X_c)$. It is able to detect more masked discrimination than $I(Z; \hat{Y}|X_c)$, i.e., when the mask is of the form $G = g(X_c, X')$. However, it can lead to false positives for Property 4 (absence of $\text{CCI}(Z \rightarrow \hat{Y})$).

Case Study: The goal is to decide whether to show ads for an editorial job requiring English proficiency, based on whether a score generated from internet activity is above a threshold. $Z \sim \text{Bern}(1/2)$ is a protected attribute denoting whether a person is a native English speaker or not. Now, consider three features $X = (X_1, X_2, X_3)$, such that: (i) X_1 : a score based on online writing samples; (ii) X_2 : a score based on browsing history, e.g., interest in English websites as compared to websites of other languages; and (iii) X_3 : a preference score based on geographical proximity. Let $X_c = X_1$ and $X_g = (X_2, X_3)$.

Suppose the true SCM is as follows: $X_1 = Z + U_1$, $X_2 = Z + U_2$, $X_3 = U_3$ and the historic scores of selected candidates are $S = X_1 + X_2 + X_3$ where $U_1, U_2, U_3 \sim i.i.d. \mathcal{N}(0, 1)$. Let the historic true labels be $Y = \mathbb{1}(S \geq 1)$ indicating whether $S \geq 1$ or not. We train a classifier of the form $\hat{Y} = 1/(1 + e^{-(w^T X + b)})$ (logistic regression). The classifier decides to show the ads if $\hat{Y} \geq 0.5$, i.e., if $w^T X + b \geq 0$ (equivalent to $\hat{S} = -\frac{w^T X}{b} \geq 1$). We train using the following loss functions:

Loss L_1 : $\min_{w,b} L_{\text{Cross Entropy}}(Y, \hat{Y})$.

Loss L_2 : $\min_{w,b} L_{\text{Cross Entropy}}(Y, \hat{Y}) + \lambda \tilde{I}(Z; \hat{Y})$, where λ is a regularizer and $\tilde{I}(Z; \hat{Y}) = -\frac{1}{2} \log(1 - \rho_{Z, \hat{Y}}^2)$ is an approximate expression of mutual information where $\rho_{Z, \hat{Y}}$ is the correlation between Z and \hat{Y} . This approximation is exact if Z and \hat{Y} are jointly Gaussian.

Loss L_3 : $\min_{w,b} L_{\text{Cross Entropy}}(Y, \hat{Y}) + \lambda \tilde{I}(Z; \hat{Y} | X_c)$, where the range of X_c is first divided into multiple discrete bins, and $\tilde{I}(Z; \hat{Y} | X_c)$ is $\sum_i \Pr(X_c \in \text{Bin } i) \tilde{I}(Z; \hat{Y} | X_c \in \text{Bin } i) = -\frac{1}{2} \sum_i \Pr(X_c \in \text{Bin } i) \log(1 - \rho_{Z, \hat{Y}, i}^2)$ and $\rho_{Z, \hat{Y}, i}$ is the conditional correlation of \hat{Y} and Z given that X_c is in the i -th discrete bin.

Observations: For L_1 , the separation boundary is very close to that based on the historic scores. But, because the past scores are correlated with browsing history (X_2), there is a danger that even when a non-native speaker has good writing score, they may not be shown an ad due to their browsing history. Regularizing with $I(Z; \hat{Y})$ (Loss L_2) does not work well because the model begins to weigh both X_1 and X_2 less, and many proficient candidates are dropped in favour of a less-important feature, namely, proximity (X_3), also reducing the accuracy (see Table 2). However, regularizing with $I(Z; \hat{Y} | X_c)$ (Loss L_3) is able to reduce the importance (weight) of browsing history relative to online scores, leading to an intermediate accuracy between L_1 and L_2 for same λ (see also Fig. 4). In a sense, our measure enables individuals with similar X_c to be treated similarly.

5 Discussion

This work provides a novel information-theoretic quantification of fairness under exemptions by adopting an axiomatic approach. We note that our properties, as stated, do not lead to a unique measure of non-exempt discrimination. They provide a qualitative separation of exempt and non-exempt discrimination, but, in line with much of the literature on fairness, do not quantify its “scaling.” However, it is not obvious what properties one can use to constrain this scaling, and remains an open question to pursue as future work. In fact, we believe that there is value in the fact that the properties do not yield a unique measure: this allows for tuning the measure for the needs of an application. E.g., Shannon established uniqueness on Shannon’s entropy with respect to some properties in [40] but the needs of the application can still drive the use of alternate measures, e.g. Renyi entropy [41] that weighs outliers differently than Shannon entropy.

While our properties do not quantify the scaling, the measure we propose does capture important aspects of the problem, e.g., it captures both masked and statistically visible components when they are present together, that existing measures such as $I(Z; \hat{Y})$ or $\text{Uni}(Z : \hat{Y} \setminus X_c)$ do not. E.g., let $X_c = U \sim \mathcal{N}(0, 1)$, $X_g = Z \sim \text{Bern}(1/2)$, and $\hat{Y} = Z + U$, i.e., Z is partially masked by U even though the visible discrimination is nonzero (a modification of Counterexample 3). Here, our measure is equal to the Shannon entropy $H(Z)$, whereas $I(Z; \hat{Y})$ or $\text{Uni}(Z : \hat{Y} \setminus X_c)$ are smaller than $H(Z)$ because they do not account for the masked component.

We also acknowledge that given the probability distribution on the data, an SCM is not always unique [34] making it difficult to use counterfactual measures in practice (as also noted for other results in the field, e.g., [18], [21]). To address this, we also propose observational relaxations of our measure and analyze what they capture and what they miss (see Table 1). In practice, this can inform which measure can be used when, e.g., $I(Z; \hat{Y} | X_c)$ can be used when cancellation of influences (Counterexample 2) does not occur (i.e., if the SCM satisfies certain faithfulness assumptions). Similarly, $\text{Uni}(Z : \hat{Y} \setminus X_c)$ may be used when accounting for masked discrimination is not required. Since the assumptions in relaxing the measure to observational ones are explicitly identified, corrections can be made if it is found that these assumptions are not satisfied. Finally, in scenarios where the SCM is known or can be evaluated from the data (see Chapters 4 and 7 in [34]), the proposed measure exactly captures the non-exempt discrimination.

References

- [1] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 214–226, 2012.
- [2] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A reductions approach to fair classification. *arXiv preprint arXiv:1803.02453*, 2018.
- [3] Moritz Hardt, Eric Price, Nati Srebro, et al. Equality of opportunity in supervised learning. In *Advances in Neural Information Processing Systems*, pages 3315–3323, 2016.
- [4] Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Optimized pre-processing for discrimination prevention. In *Advances in Neural Information Processing Systems*, pages 3992–4001, 2017.
- [5] Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier with prejudice remover regularizer. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 35–50. Springer, 2012.
- [6] Hao Wang, Berk Ustun, and Flavio P Calmon. Repairing without retraining: Avoiding disparate impact with counterfactual distributions. *arXiv preprint arXiv:1901.10501*, 2019.
- [7] Aditya Krishna Menon and Robert C Williamson. The cost of fairness in binary classification. In *Conference on Fairness, Accountability and Transparency*, pages 107–118, 2018.
- [8] Junpei Komiyama and Hajime Shima. Two-stage algorithm for fairness-aware machine learning. *arXiv preprint arXiv:1710.04924*, 2017.
- [9] Michele Donini, Luca Oneto, Shai Ben-David, John S Shawe-Taylor, and Massimiliano Pontil. Empirical risk minimization under fairness constraints. In *Advances in Neural Information Processing Systems*, pages 2791–2801, 2018.
- [10] AmirEmad Ghassami, Sajad Khodadadian, and Negar Kiyavash. Fairness in supervised learning: An information theoretic approach. In *IEEE International Symposium on Information Theory*, pages 176–180, 2018.
- [11] Hoda Heidari, Claudio Ferrari, Krishna Gummadi, and Andreas Krause. Fairness behind a veil of ignorance: A welfare analysis for automated decision making. In *Advances in Neural Information Processing Systems*, pages 1265–1276, 2018.
- [12] Jiachun Liao, Chong Huang, Peter Kairouz, and Lalitha Sankar. Learning generative adversarial representations (gap) under fairness and censoring constraints. *arXiv preprint arXiv:1910.00411*, 2019.
- [13] Kush R Varshney. Trustworthy machine learning and artificial intelligence. *XRDS: Crossroads, The ACM Magazine for Students*, 25(3):26–29, 2019.
- [14] Faisal Kamiran, Indrė Žliobaitė, and Toon Calders. Quantifying explainable discrimination and removing illegal discrimination in automated decision making. *Knowledge and information systems*, 35(3):613–644, 2013.

- [15] US Equal Employment Opportunity Commission. The Equal Pay Act of 1963. <https://www.eeoc.gov/laws/statutes/epa.cfm>, 1963.
- [16] Katie Manley. The bfoq defense: Title vii’s concession to gender discrimination. *Duke J. Gender L. & Pol’y*, 16:169, 2009.
- [17] Dept. of Trade and Industry. Genuine Occupational Qualifications, A Good Practice Guide for Employers, 2008.
- [18] Niki Kilbertus, Mateo Rojas Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding discrimination through causal reasoning. In *Advances in Neural Information Processing Systems*, pages 656–666, 2017.
- [19] Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. Interventional Fairness: Causal Database Repair for Algorithmic Fairness. In *Proceedings of the 2019 International Conference on Management of Data, SIGMOD ’19*, pages 793–810. ACM, 2019.
- [20] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016.
- [21] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems*, pages 4066–4076, 2017.
- [22] Anupam Datta, Matthew Fredrikson, Gihyuk Ko, Piotr Mardziel, and Shayak Sen. Use privacy in data-driven systems: Theory and experiments with machine learnt programs. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1193–1210, 2017.
- [23] Anonymous. Conditional Debiasing for Neural Networks, 2019.
- [24] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. Fairness constraints: Mechanisms for fair classification. *arXiv preprint arXiv:1507.05259*, 2015.
- [25] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic decision making and the cost of fairness. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’17*, pages 797–806. ACM, 2017.
- [26] Chris Russell, Matt J Kusner, Joshua Loftus, and Ricardo Silva. When worlds collide: integrating different counterfactual assumptions in fairness. In *Advances in Neural Information Processing Systems*, pages 6414–6423, 2017.
- [27] Silvia Chiappa and Thomas PS Gillam. Path-specific counterfactual fairness. *arXiv preprint arXiv:1802.08139*, 2018.
- [28] Praveen Venkatesh, Sanghamitra Dutta, and Pulkit Grover. Information flow in computational systems. *arXiv preprint arXiv:1902.02292*, 2019.
- [29] Nils Bertschinger, Johannes Rauh, Eckehard Olbrich, Jürgen Jost, and Nihat Ay. Quantifying unique information. *Entropy*, 16(4):2161–2183, 2014.
- [30] Paul L Williams and Randall D Beer. Nonnegative decomposition of multivariate information. *arXiv preprint arXiv:1004.2515*, 2010.

- [31] Tycho Tax, Pedro Mediano, and Murray Shanahan. The partial information decomposition of generative neural network models. *Entropy*, 19(9):474, 2017.
- [32] Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pages 325–333. JMLR.org, 2013.
- [33] Pradeep Kr Banerjee, Johannes Rauh, and Guido Montúfar. Computing the unique information. In *IEEE International Symposium on Information Theory*, pages 141–145, 2018.
- [34] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: foundations and learning algorithms*. MIT press, 2017.
- [35] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [36] Anupam Datta, Shayak Sen, and Yair Zick. Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems. In *IEEE Symposium on Security and Privacy*, pages 598–617, 2016.
- [37] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pages 1885–1894. JMLR.org, 2017.
- [38] Philip Adler, Casey Falk, Sorelle A Friedler, Tionney Nix, Gabriel Rybeck, Carlos Scheidegger, Brandon Smith, and Suresh Venkatasubramanian. Auditing black-box models for indirect influence. *Knowledge and Information Systems*, 54(1):95–122, 2018.
- [39] Andreas Henelius, Kai Puolamäki, Henrik Boström, Lars Asker, and Panagiotis Papapetrou. A peek into the black box: exploring classifiers by randomization. *Data mining and knowledge discovery*, 28(5-6):1503–1529, 2014.
- [40] Claude Elwood Shannon. A mathematical theory of communication. *Bell system technical journal*, 27(3):379–423, 1948.
- [41] Alfréd Rényi. On measures of entropy and information. In *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability*, volume 1. The Regents of the University of California, 1961.
- [42] Johannes Rauh, Pradeep Kr Banerjee, Eckehard Olbrich, and Jürgen Jost. Unique information and secret key decompositions. *arXiv preprint arXiv:1901.08007*, 2019.

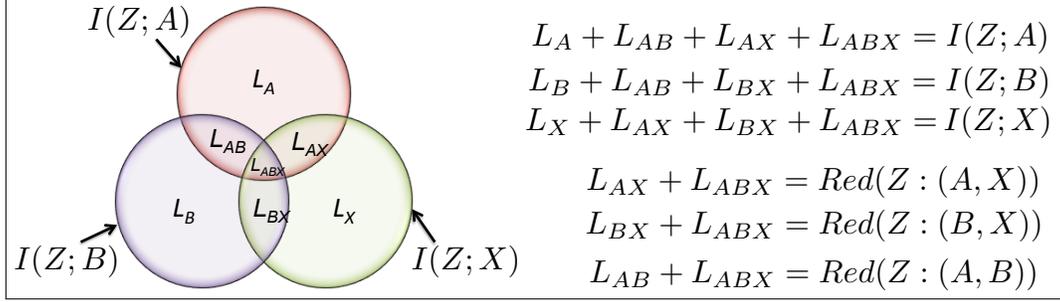


Figure 5: Partial Information Decomposition explained using Venn diagrams

A Partial Information Decomposition: Relevant Properties

Here, we state and prove some important lemmas related to PID that are useful for the proofs in the rest of the paper.

Lemma 5 (Triangle inequality of unique information). *For any four random variables Z , A , B and X , if $\text{Uni}(Z : A \setminus X) > \text{Uni}(Z : B \setminus X)$, then it implies that, $\text{Uni}(Z : A \setminus B) > 0$.*

Proof of Lemma 5:. In [42, Proposition 2], the authors show that for any (Z, A, B, X) ,

$$\text{Uni}(Z : A \setminus X) \leq \text{Uni}(Z : B \setminus X) + \text{Uni}(Z : A \setminus B). \quad (6)$$

Thus,

$$\text{Uni}(Z : A \setminus B) \geq \text{Uni}(Z : A \setminus X) - \text{Uni}(Z : B \setminus X) > 0. \quad (7)$$

□

Lemma 6 (Monotonicity of unique information with more excluded variables). [29, Lemma 25] *Unique information about Z in A that is not present in B is non-increasing as more variables are added to the set B , i.e.,*

$$\text{Uni}(Z : A \setminus B_1) \geq \text{Uni}(Z : A \setminus (B_1 \cup B_2)). \quad (8)$$

Proof of Lemma 6. This result is stated and proved in [29, Lemma 25].

□

Lemma 7 (Zero-synergy property of deterministic functions). *Let $f(Z)$ be any deterministic function of Z , and let X be any random variable. Then,*

$$\text{Uni}(Z : f(Z) \setminus X) = \text{I}(Z; f(Z)|X) \text{ and } \text{Syn}(Z : (f(Z), X)) = 0. \quad (9)$$

Proof of Lemma 7:. Recall from the definition of unique information that Δ denotes the set of all joint distributions of (X, Y, Z) and Δ_p is the set of all such joint distributions that have the same marginals for (Z, Y) and (Z, X) as the true distribution, i.e.,

$$\Delta_p = \{Q \in \Delta : q(z, y) = \Pr(Z = z, Y = y) \text{ and } q(z, x) = \Pr(Z = z, X = x)\}. \quad (10)$$

We first show that if $Y = f(Z)$, then $q(y|z)$ becomes a point measure, and hence Δ_p is only a singleton set which only consists of the true distribution.

Observe that, for any $Q \in \Delta_p$,

$$\begin{aligned}
q(x, y, z) &= q(z)q(y|z)q(x|y, z) && \text{[chain rule of probability]} \\
&= \Pr(Z = z) \Pr(Y = y|Z = z)q(x|y, z) && \text{[}q(z, y) = \Pr(Z = z, Y = y)\text{]} \\
&= \begin{cases} \Pr(Z = z)q(x|y, z), & \text{if } y = f(z) \\ 0, & \text{otherwise} \end{cases} && \text{[}\Pr(Y = y|Z = z) = 1 \text{ only if } y = f(z)\text{]} \\
&= \begin{cases} \Pr(Z = z)q(x|z), & \text{if } y = f(z) \\ 0, & \text{otherwise} \end{cases} \\
&= \begin{cases} \Pr(Z = z) \Pr(X = x|Z = z), & \text{if } y = f(z) \\ 0, & \text{otherwise} \end{cases} && \text{[}q(x|z) = \Pr(X = x|Z = z)\text{]} \\
&= \Pr(X = x, Y = y, Z = z). && (11)
\end{aligned}$$

Thus,

$$\text{Uni}(Z : f(Z) \setminus X) = \min_{Q \in \Delta_p} I_Q(Z; Y|X) = I(Z; Y|X) = I(Z; f(Z)|X). \quad (12)$$

□

Lemma 8 (Zero-redundancy property of non-descendants). *Let $\text{CCI}(Z \rightarrow B) = 0$. Then, $\text{Uni}(Z : A \setminus B) = I(Z; A)$.*

Proof of Lemma 8. Because $\text{CCI}(Z \rightarrow B) = 0$, we have $I(Z; B) = 0$. Now, observe that,

$$\begin{aligned}
\text{Uni}(Z : B \setminus A) + \text{Red}(Z : (B, A)) &= 0 && \text{[by PID since } I(Z; B) = 0\text{]} \\
\implies \text{Red}(Z : (B, A)) &= 0 && \text{[non-negativity of PID terms]} \\
\implies \text{Uni}(Z : A \setminus B) + \text{Red}(Z : (B, A)) &= \text{Uni}(Z : A \setminus B) \\
\implies I(Z; A) &= \text{Uni}(Z : A \setminus B) && \text{[by PID].}
\end{aligned} \quad (13)$$

□

Lemma 9 (Synergy equivalence for non-descendants). *Let $\text{CCI}(Z \rightarrow B) = 0$. Then,*

$$I(Z; A|B) > I(Z; A) \iff \text{Syn}(Z : (A, B)) > 0. \quad (14)$$

Proof of Lemma 9. Observe that,

$$I(Z; A|B) > I(Z; A) \quad (15)$$

$$\iff \text{Uni}(Z : A \setminus B) + \text{Syn}(Z : (A, B)) > I(Z; A) \quad \text{[by PID]} \quad (16)$$

$$\iff \text{Uni}(Z : A \setminus B) + \text{Syn}(Z : (A, B)) > \text{Uni}(Z : A \setminus B) \quad \text{[using Lemma 8]} \quad (17)$$

$$\iff \text{Syn}(Z : (A, B)) > 0. \quad (18)$$

□

Lemma 10 (Maximal conditional mutual information). *Let $A = f(Z, U_X)$ where $Z \perp U_X$ and $B = g(U_X)$ for some deterministic functions $f(\cdot)$ and $g(\cdot)$ respectively. Then,*

$$I(Z; A|U_X) \geq I(Z; A|B). \quad (19)$$

Proof of Lemma 10. Observe that,

$$\begin{aligned}
I(Z; U_X|A, B) &\geq 0 && \text{[non-negativity property]} \\
\implies H(Z|A, B) - H(Z|A, B, U_X) &\geq 0 && \text{[by definition]} \\
\implies H(Z|A, B) - H(Z|A, U_X) &\geq 0 && [B = g(U_X)] \\
\implies H(Z) - H(Z|A, U_X) &\geq H(Z) - H(Z|A, B) \\
\implies H(Z|U_X) - H(Z|A, U_X) &\geq H(Z|B) - H(Z|A, B) && [Z \perp U_X \text{ and } Z \perp B] \\
\implies I(Z; A|U_X) &\geq I(Z; A|B). && (20)
\end{aligned}$$

□

Lemma 11 (Property of Variable W). *Let $X_c = f_c(Z, U_X)$ where $Z \perp U_X$ and $f_c(\cdot)$ is a deterministic function. Also let $\tilde{Y} = h(Z, \tilde{U}_X)$ where \tilde{U}_X is i.i.d. as U_X , with $Z \perp \tilde{U}_X$ and $h(\cdot)$ is a deterministic function. Now, $W = [h(Z, u_x^{(1)}), \dots, h(Z, u_x^{(k)})]$, where $\{u_x^{(1)}, \dots, u_x^{(k)}\}$ is the set of all values with $\Pr(U_X = u_x) > 0$. Then,*

$$I(Z; \tilde{Y}|X_c) = I(W; \tilde{Y}|X_c).$$

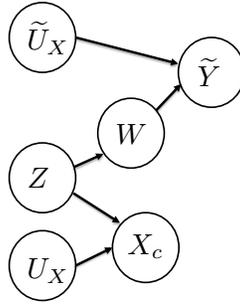


Figure 6: Structural Causal Model corresponding to Lemma 11: In this SCM, $I(Z; \tilde{Y}|X_c, W) = 0$.

Proof of Lemma 11. First observe that, \tilde{Y} is a deterministic function of \tilde{U}_X and W , i.e., it is the i -th element of W when $\tilde{U}_X = u_x^{(i)}$. And, \tilde{U}_X is independent of both X_c and Z . Thus, $\Pr(\tilde{Y} = y|W = w, X_c = x_c)$ and $\Pr(\tilde{Y} = y|W = w, X_c = x_c, Z = z)$ are both equal to

$$\Pr(\tilde{Y} = y|W = w) = \begin{cases} \Pr(\tilde{U}_X = u_x^{(i)}), & y = w[i] \\ 0, & \text{otherwise.} \end{cases} \quad (21)$$

This implies that (see also Fig. 6),

$$I(Z; \tilde{Y}|X_c, W) = 0. \quad (22)$$

Now,

$$I(Z; \tilde{Y}|X_c) = I(Z, W; \tilde{Y}|X_c) - I(W; \tilde{Y}|X_c, Z) \quad (23)$$

$$= I(W; \tilde{Y}|X_c) + I(Z; \tilde{Y}|X_c, W) - I(W; \tilde{Y}|X_c, Z) \quad (24)$$

$$= I(W; \tilde{Y}|X_c) + 0 - I(W; \tilde{Y}|X_c, Z) \quad [\text{see (22)}] \quad (25)$$

$$= I(W; \tilde{Y}|X_c) + 0 - 0 \quad [W \text{ is a deterministic function of } Z] \quad (26)$$

$$= I(W; \tilde{Y}|X_c). \quad (27)$$

□

B Counterfactual Causal Influence (CCI) and its connections to Counterfactual Fairness

B.1 Proofs of Lemmas in Section 2

Here, we first provide a proof of Lemma 1 and then show the connections of CCI to counterfactual fairness [21]. For ease of reading, we repeat the statements of the lemma here again.

Lemma 1 (Information-Theoretic Equivalent of CCI). *Let $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$. Then $\text{CCI}(Z \rightarrow \hat{Y}) \neq 0$ if and only if $I(Z; W) > 0$.*

Proof of Lemma 1. Observe that,

$$\begin{aligned} \text{CCI}(Z \rightarrow P) &= \mathbb{E}_{Z, Z', U_X} [|h(Z, U_X) - h(Z', U_X)|] \\ &= \sum_{z_1, z_2, u_x} \Pr(Z = z_1, Z' = z_2, U_X = u_x) |h(z_1, u_x) - h(z_2, u_x)| \\ &= \sum_{z_1, z_2, u_x} \Pr(Z = z_1) \Pr(Z' = z_2) \Pr(U_X = u_x) |h(z_1, u_x) - h(z_2, u_x)| \quad [\text{from independence}]. \end{aligned} \quad (28)$$

The summation consist of non-negative terms. Therefore, $\text{CCI}(Z \rightarrow P) = 0$, if and only if all the terms in the summation are zero, *i.e.*, for all z_1, z_2 and u_x with $\Pr(Z = z_1), \Pr(Z = z_2), \Pr(U_X = u_x) > 0$, $|h(z_1, u_x) - h(z_2, u_x)| = 0$. This is also equivalent to $h(Z, u_x)$ being constant over all possible values of random variable $Z = z$.

Now, observe that $I(Z; W) = H(W) - H(W|Z) = H(W)$. This can be 0 if and only if $h(Z, u_x)$ is constant over all possible $Z = z$ with $\Pr(Z = z) > 0$, and for all $U_X = u_x$ with $\Pr(U_X = u_x) > 0$. Thus, the first and second statements have both way implication. □

B.2 Connections to Counterfactual Fairness

Now, let $\hat{Y}_{Z \leftarrow z_1}(U)$ denote the value of \hat{Y} when the value of Z is set as z_1 by an intervention.

Definition 7 (Counterfactual Fairness [21]). *A predictor \hat{Y} is counterfactually fair if under any context $X = x$ and $Z = z_1$, we have*

$$\begin{aligned} &\Pr(\hat{Y}_{Z \leftarrow z_1}(U) = y | \text{evidence } X = x \text{ when } Z = z_1) \\ &= \Pr(\hat{Y}_{Z \leftarrow z_2}(U) = y | \text{evidence } X = x \text{ when } Z = z_1), \end{aligned} \quad (29)$$

for all y and for any value z_2 attainable by Z .

Next, we show that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ is equivalent to the counterfactual fairness criterion of [21].

Lemma 12. $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ is equivalent to counterfactual fairness.

Proof of Lemma 12. Let $X = f(Z, U_X)$ and $\hat{Y} = r(X) = r \circ f(Z, U_X) = h(Z, U_X)$. Recall from Lemma 1, that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ implies that for all z_1, z_2, u_x with $\Pr(Z = z_1), \Pr(Z = z_2), \Pr(U_X = u_x)$, we have $r \circ f(z_1, u_x) = r \circ f(z_2, u_x)$. Now,

$$\begin{aligned}
& \Pr(\hat{Y}_{Z \leftarrow z_1}(U) = y \mid \text{evidence } X = x \text{ when } Z = z_1) \\
&= \Pr(r \circ f(z_1, U_X) = y \mid \text{evidence } X = x \text{ when } Z = z_1) \\
&= \Pr(r \circ f(z_1, U_X) = y \mid U_X \in \mathcal{S}_1) && \text{[where } \mathcal{S}_1 = \{u_x : x = f(u_x, z_1)\}] \\
&= \Pr(r \circ f(z_2, U_X) = y \mid U_X \in \mathcal{S}_1) && \text{[since } r \circ f(z_1, u_x) = r \circ f(z_2, u_x)] \\
&= \Pr(\hat{Y}_{Z \leftarrow z_2}(U) = y \mid \text{evidence } X = x \text{ when } Z = z_1). \tag{30}
\end{aligned}$$

Thus, we show that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ implies counterfactual fairness.

Now, we prove the implication in the other direction. Suppose that the counterfactual fairness criterion holds. Therefore, we have

$$\begin{aligned}
& \Pr(r \circ f(z_1, U_X) = y \mid U_X \in \mathcal{S}_1) && \text{[where } \mathcal{S}_1 = \{u_x : x = f(u_x, z_1)\}] \\
&= \Pr(r \circ f(z_2, U_X) = y \mid U_X \in \mathcal{S}_1). \tag{31}
\end{aligned}$$

Or,

$$\Pr(r \circ f(z_1, U_X) = y, U_X \in \mathcal{S}_1) = \Pr(r \circ f(z_2, U_X) = y, U_X \in \mathcal{S}_1). \tag{32}$$

Or,

$$\sum_{u_x \in \mathcal{S}_1} p(u_x) \mathbb{1}(r \circ f(z_1, u_x) = y) = \sum_{u_x \in \mathcal{S}_1} p(u_x) \mathbb{1}(r \circ f(z_2, u_x) = y). \tag{33}$$

For a particular y , observe that $\mathbb{1}(r \circ f(z_1, u_x) = y)$ is the same for all $u_x \in \mathcal{S}_1$ because $f(z_1, u_x) = x$ for all $u_x \in \mathcal{S}_1$. Thus, for (33) to hold, all $u_x \in \mathcal{S}_1$ should also satisfy $\mathbb{1}(r \circ f(z_2, u_x) = y) = \mathbb{1}(r \circ f(z_1, u_x) = y)$, implying $r \circ f(z_2, u_x) = r \circ f(z_1, u_x)$. This is equivalent to $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ using Lemma 1. \square

C Appendix to Section 3

Here, we provide the proofs of the results stated in Section 3. For the ease of reading, we again repeat the statements of the results.

C.1 Proofs of results in Section 3.1 (Theorem 1 and Lemma 2)

Theorem 1 (Properties). *Properties 1, 2, 3 and 4 are satisfied by $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c)$.*

Proof of Theorem 1. Here, we formally show that our proposed measure satisfies all the four desirable properties mentioned in Section 3. We restate each of the properties again and then show that it is satisfied.

Property 1 (Complete Exemption). *M_{NE} should be 0 if all features are categorized into X_c , i.e., $X_c = X$ and $X_g = \phi$.*

Observe that, when $X = X_c$,

$$\begin{aligned}
M_{NE} &= \text{Uni}(Z : W \setminus X) - \text{Uni}(Z : W \setminus X, \hat{Y}) \\
&= I(Z; W | X) - I(Z; W | X, \hat{Y}) && \text{[Using Lemma 7]} \\
&= H(Z | X) - H(Z|W, X) - H(Z|X, \hat{Y}) + H(Z|W, X, \hat{Y}) && \text{[By Definition]} \\
&= I(Z; \hat{Y} | X) - I(Z; \hat{Y} | X, W) \\
&\leq I(Z; \hat{Y} | X) && \text{[Non-negativity of Mutual Information]} \\
&= 0 && \text{[}\hat{Y} \text{ is a deterministic function of } X\text{].}
\end{aligned} \tag{34}$$

Property 2 (Non-Exempt Visible Discrimination). M_{NE} should be strictly greater than 0 if $\text{Uni}(Z : \hat{Y} \setminus X_c) > 0$.

Let $\tilde{Y} = h(Z, \tilde{U}_X)$ where \tilde{U}_X and U_X are i.i.d., and $\tilde{U}_X \perp Z$. Also let us represent $X_c = f_c(Z, U_X)$. This holds because X_c is fully determined by the latent factors in the SCM. Now, observe that,

$$\begin{aligned}
\text{Uni}(Z : \hat{Y} \setminus X_c) &= \min_{Q \in \Delta_p} I_Q(Z; \hat{Y} | X_c) \\
&\leq I(Z; \tilde{Y} | X_c) \text{ [since } (Z, \tilde{Y}) \text{ and } (Z, \hat{Y}) \text{ have same joint distribution]} \\
&= I(W; \tilde{Y} | X_c) \text{ [Using Lemma 11]} \\
&= M_{NE}. \text{ [Using Lemma 2]}
\end{aligned} \tag{35}$$

Thus, $M_{NE} \geq \text{Uni}(Z : \hat{Y} \setminus X_c)$ and is thus strictly greater than 0 if $\text{Uni}(Z : \hat{Y} \setminus X_c) > 0$.

Property 3 (Non-Exempt Masking). A measure M_{NE} should be non-zero in the canonical example of masked discrimination, i.e., Example 2 even if $I(Z; \hat{Y}) = 0$. However, M_{NE} should be 0 if $Z - X_c - \hat{Y}$ form a Markov chain.

Observe that,

$$\begin{aligned}
M_{NE} &= \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus X_c, \hat{Y}) \\
&= I(Z; W | X_c) - I(Z; W | X_c, \hat{Y}) && \text{[Using Lemma 7]} \\
&= H(Z | X_c) - H(Z|W, X_c) - H(Z|X_c, \hat{Y}) + H(Z|W, X_c, \hat{Y}) && \text{[By Definition]} \\
&= I(Z; \hat{Y} | X_c) - I(Z; \hat{Y} | X_c, W) \\
&\leq I(Z; \hat{Y} | X_c) && \text{[Non-negativity of Mutual Information].}
\end{aligned} \tag{36}$$

Now, suppose that the Markov chain $Z - X_c - \hat{Y}$ hold. Then $I(Z; \hat{Y} | X_c) = 0$ implying that $M_{NE} = 0$. Therefore, this property is satisfied.

Property 4 (Cancellation of Influence). M_{NE} should be 0 if $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ (or equivalently, $I(Z; W) = 0$).

In Lemma 1, we demonstrated that $\text{CCI}(Z \rightarrow \hat{Y}) = 0$ is equivalent to the condition that $I(Z; W) = 0$. Now observe that,

$$\text{Uni}(Z : W \setminus X_c) = I(Z; W) - \text{Red}(Z : (W, X_c)) \leq I(Z; W) = 0. \tag{37}$$

Similarly,

$$\text{Uni}(Z : W \setminus X_c, \hat{Y}) = \text{I}(Z; W) - \text{Red}(Z : (W, (X_c, \hat{Y}))) \leq \text{I}(Z; W) = 0. \quad (38)$$

Thus, $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus X_c, \hat{Y}) = 0$ whenever $\text{I}(Z; W) = 0$.
Therefore, this property is satisfied. \square

Lemma 2 (Non-Exempt Discrimination Equivalence). *The proposed measure $M_{NE} = \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus \hat{Y}, X_c)$ is equal to $\text{I}(W; \hat{Y} | X_c)$.*

Proof of Lemma 2.

$$\begin{aligned} M_{NE} &= \text{Uni}(Z : W \setminus X_c) - \text{Uni}(Z : W \setminus X_c, \hat{Y}) \\ &= \text{I}(Z; W | X_c) - \text{I}(Z; W | X_c, \hat{Y}) && \text{[Using Lemma 7]} \\ &= H(W | X_c) - H(W | Z, X_c) - H(W | X_c, \hat{Y}) + H(W | Z, X_c, \hat{Y}) && \text{[By Definition]} \\ &= H(W | X_c) - H(W | X_c, \hat{Y}) && \text{[} W \text{ is a deterministic function of } Z \text{]} \\ &= \text{I}(W; \hat{Y} | X_c). \end{aligned} \quad (39)$$

\square

C.2 Proofs of results in Section 3.3 (Theorem 2, Lemma 3 and Lemma 4)

Theorem 2 (Non-negative Decomposition of Total Discrimination). *The total discrimination can be decomposed into four non-negative components as follows:*

$$\text{I}(Z; W) = M_{V,NE} + M_{V,E} + M_{M,NE} + M_{M,E}. \quad (5)$$

Here $M_{V,NE} = \text{Uni}(Z : \hat{Y} \setminus X_c)$ is the visible, non-exempt component and $M_{V,E} = \text{Red}(Z : (\hat{Y}, X_c))$ is the visible, exempt component. These two terms add to form $\text{I}(Z; \hat{Y})$ which is the total statistically visible discrimination. Likewise, $M_{M,NE} = M_{NE} - M_{V,NE}$ is the masked, non-exempt component, and $M_{M,E} = \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - M_{M,NE}$ is the masked, exempt component.

Proof of Theorem 2. First consider $M_{V,NE} = \text{Uni}(Z : \hat{Y} \setminus X_c)$ and $M_{V,E} = \text{Red}(Z : (\hat{Y}, X_c))$. Because all PID terms are non-negative by definition, both $M_{V,NE}$ and $M_{V,E}$ are non-negative.

Now, consider $M_{M,E}$. Observe that,

$$\begin{aligned} M_{M,E} &= \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - M_{M,NE} \\ &= \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - M_{NE} + M_{V,NE} \\ &= \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - \text{Uni}(Z : W \setminus X_c) + \text{Uni}(Z : W \setminus X_c, \hat{Y}) + \text{Uni}(Z : \hat{Y} \setminus X_c) && \text{[By Definition]} \\ &= \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - \text{Uni}(Z : W \setminus X_c) + \text{Uni}(Z : \hat{Y} \setminus X_c) + \text{Uni}(Z : W \setminus X_c, \hat{Y}) && \text{[Rearrangement]} \\ &\geq \text{I}(Z; W) - \text{I}(Z; \hat{Y}) - \text{Uni}(Z : W \setminus \hat{Y}) + \text{Uni}(Z : W \setminus X_c, \hat{Y}) && \text{[Triangle inequality (Lemma 5)]} \\ &\geq 0 + \text{Uni}(Z : W \setminus X_c, \hat{Y}) && \text{[Markov Chain } Z - W - \hat{Y} \text{]} \\ &\geq 0 && \text{[Non-negativity property].} \end{aligned} \quad (40)$$

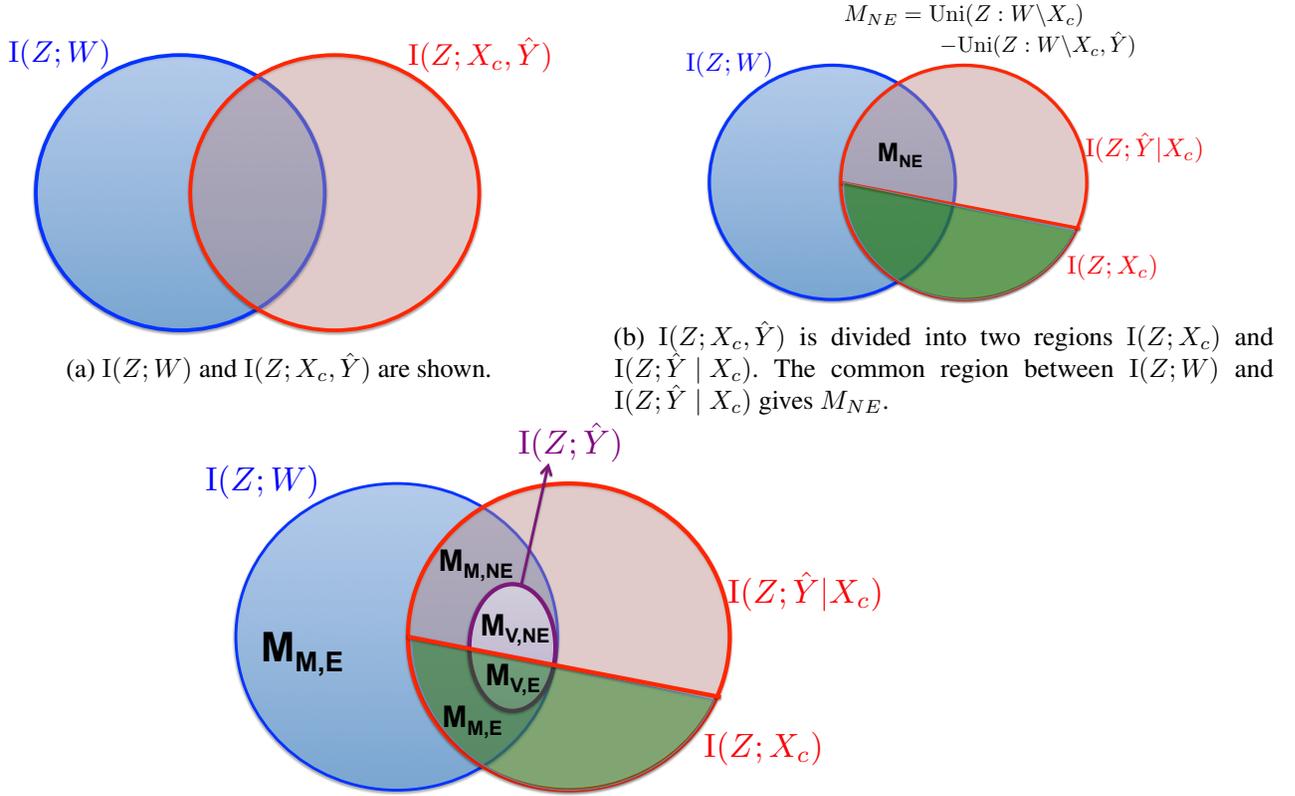


Figure 7: The figure shows the overall decomposition of $I(Z; W)$ into four components, namely $M_{M,NE}$ (masked non-exempt), $M_{V,NE}$ (visible non-exempt), $M_{M,E}$ (masked exempt) and $M_{V,E}$ (visible exempt) respectively.

Lastly, we consider $M_{M,NE} = M_{NE} - \text{Uni}(Z : \hat{Y} \setminus X_c)$. Let $\tilde{Y} = h(Z, \tilde{U}_X)$ where \tilde{U}_X and U_X are i.i.d., and $\tilde{U}_X \perp\!\!\!\perp Z$. Also let us represent $X_c = f_c(Z, U_X)$. This holds because X_c is fully determined by the latent factors in the SCM. Now, observe that,

$$\begin{aligned}
\text{Uni}(Z : \hat{Y} \setminus X_c) &= \min_{Q \in \Delta_p} I_Q(Z; \hat{Y} | X_c) \\
&\leq I(Z; \tilde{Y} | X_c) \quad [\text{since } (Z, \tilde{Y}) \text{ and } (Z, \hat{Y}) \text{ have same joint distribution}] \\
&= I(W; \tilde{Y} | X_c) \quad [\text{Using Lemma 11}] \\
&= M_{NE}. \quad [\text{Using Lemma 2}]
\end{aligned} \tag{41}$$

Therefore, $M_{M,NE}$ is non-negative. □

Lemma 3 (Masked Discrimination). *The total masked discrimination $I(Z; W) - I(Z; \hat{Y})$ is equal to $\text{Uni}(Z : W \setminus \hat{Y})$.*

Proof of Lemma 3. From the Markov chain $Z - W - \hat{Y}$, we have $I(Z; \hat{Y} | W) = 0$, implying $\text{Uni}(Z : \hat{Y} \setminus W) = 0$. Thus,

$$\begin{aligned}
I(Z; W) - I(Z; \hat{Y}) \\
&= \text{Uni}(Z : W \setminus \hat{Y}) - \text{Uni}(Z : \hat{Y} \setminus W) = \text{Uni}(Z : W \setminus \hat{Y}) \geq 0.
\end{aligned} \tag{42}$$

□

Lemma 4 (Masked Discrimination Implications). *The following two statements are equivalent:*

- $I(Z; \hat{Y} | U_X) - I(Z; \hat{Y}) > 0$.
 - \exists a random variable G of the form $G = g(U_X)$ such that $I(Z; \hat{Y} | G) > I(Z; \hat{Y})$.
- Either of these statements imply $I(Z; W) - I(Z; \hat{Y}) > 0$.*

Proof of Lemma 4. First, we show that the second statement implies the first statement. Suppose there exists a $G = g(U_X)$ such that $I(Z; \hat{Y} | G) > I(Z; \hat{Y})$.

Observe that,

$$\begin{aligned}
&\sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x)) \\
&= \sum_{u_x} \Pr(U_X = u_x) (\mathbb{H}(h(Z, u_x)) - \mathbb{H}(h(Z, u_x) | Z)) \quad [\text{by definition}] \\
&= \sum_{u_x} \Pr(U_X = u_x) \mathbb{H}(h(Z, u_x)) \quad [h(Z, u_x) \text{ is a function of } Z] \\
&= \sum_{u_x} \Pr(U_X = u_x) \mathbb{H}(\hat{Y} | U_X = u_x) \quad [\text{since } Z \perp\!\!\!\perp U_X] \\
&= \mathbb{H}(\hat{Y} | U_X) \quad [\text{by definition}] \\
&= \mathbb{H}(\hat{Y} | U_X) - \mathbb{H}(\hat{Y} | Z, U_X) \quad [\mathbb{H}(\hat{Y} | Z, U_X) = 0 \text{ as } \hat{Y} = h(Z, U_X)] \\
&= I(Z; \hat{Y} | U_X) \quad [\text{by definition}] \\
&\geq I(Z; \hat{Y} | G) \quad [\text{using Lemma 10 for any } G = g(U_X)].
\end{aligned} \tag{43}$$

This together with the second statement of this lemma, *i.e.*, $I(Z; \hat{Y}|G) > I(Z; \hat{Y})$ implies that $\sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x)) > I(Z; \hat{Y})$ which is essentially the first statement. Thus, we prove that the second statement implies the first statement.

Next, we show that the first statement also implies the second statement. This is equivalent to showing that if $I(Z; \hat{Y}|G) \leq I(Z; \hat{Y})$ for all G of the form $G = g(U_X)$, then $\sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x))$ is not greater than $I(Z; \hat{Y})$.

Suppose $I(Z; \hat{Y}|G) \leq I(Z; \hat{Y})$ for all $G = g(U_X)$. Choosing $G = U_X$, we therefore have

$$I(Z; \hat{Y}|U_X) \leq I(Z; \hat{Y}). \quad (44)$$

On the other hand, from Lemma 10, $I(Z; \hat{Y}|U_X) \geq I(Z; \hat{Y}|g(U_X))$ for any deterministic function $g(\cdot)$ which could even be a constant. Choosing $g(\cdot)$ to be a constant function, we obtain

$$I(Z; \hat{Y}|U_X) \geq I(Z; \hat{Y}). \quad (45)$$

From (44) and (45), we get

$$I(Z; \hat{Y}|U_X) = I(Z; \hat{Y}). \quad (46)$$

Next, observe that,

$$\begin{aligned} I(Z; U_X|\hat{Y}) &= I(Z; \hat{Y}|U_X) - I(Z; \hat{Y}) + I(Z; U_X) && \text{[by definition and regrouping]} \\ &= I(Z; \hat{Y}|U_X) - I(Z; \hat{Y}) && [Z \perp U_X] \\ &= 0 && \text{[using (46)].} \end{aligned} \quad (47)$$

From (47), we get $Z \perp U_X|\hat{Y}$. Thus,

$$\begin{aligned} \Pr(Z = z, U_X = u|\hat{Y} = y) &= \Pr(Z = z|\hat{Y} = y) \Pr(U_X = u|\hat{Y} = y). \\ \implies \Pr(Z = z|U_X = u, \hat{Y} = y) &= \Pr(Z = z|\hat{Y} = y) && \text{[Chain rule]} \end{aligned} \quad (48)$$

$$\implies \Pr(Z = z, \hat{Y} = y|U_X = u_x) = \Pr(Z = z, \hat{Y} = y) \quad \text{[Chain rule]} \quad (49)$$

$$\implies \Pr(\hat{Y} = y|U_X = u_x) = \Pr(\hat{Y} = y) \quad \text{[Marginal]} \quad (50)$$

Therefore, for all u_x with $\Pr(U_X = u_x) > 0$,

$$\begin{aligned} I(Z; h(Z, u_x)) &= I(Z; \hat{Y}|U_X = u_x) \\ &= \sum_{z,y} \Pr(Z = z, \hat{Y} = y|U_X = u_x) \log \left(\frac{\Pr(Z = z, \hat{Y} = y|U_X = u_x)}{\Pr(Z = z|U_X = u_x) \Pr(\hat{Y} = y|U_X = u_x)} \right) \\ &= \sum_{z,y} \Pr(Z = z, \hat{Y} = y) \log \left(\frac{\Pr(Z = z, \hat{Y} = y)}{\Pr(Z = z) \Pr(\hat{Y} = y)} \right) \quad \text{[using (49), (50) and } Z \perp U_X] \\ &= I(Z; \hat{Y}) \quad \text{[by definition].} \end{aligned} \quad (51)$$

Thus, $\sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x)) = I(Z; \hat{Y})$ and is not greater. Therefore, the first statement also implies the second statement.

Thus, we prove that the first and second statements are equivalent.

Now, from the Markov chain $Z - W - h(Z, u_x)$, we have $I(Z; W) \geq I(Z; h(Z, u_x))$ for each u_x , leading to $I(Z; W) \geq \sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x))$. Consequently $I(Z; W)$ is strictly greater than $I(Z; \hat{Y})$ whenever $\sum_{u_x} \Pr(U_X = u_x) I(Z; h(Z, u_x)) > I(Z; \hat{Y})$. \square

Remark 5 (Conditioning to capture masked discrimination). *In general, conditioning on a random variable G leading to $I(Z; \hat{Y}|G) > I(Z; \hat{Y})$ can sometimes detect masked discrimination, if conditioning exposes more discrimination than what was already visible. For example, $I(Z; \hat{Y}|X_c)$ can detect masked discrimination if the mask is of the form $g(X_c)$. However, conditioning on any random variable G leading to $I(Z; \hat{Y}|G) > I(Z; \hat{Y})$ cannot always be interpreted as a case of masked discrimination because this can sometimes lead to false positives in detecting discrimination, e.g., if $\hat{Y} = U_1^{(X)}$ and G is chosen as $U_1^{(X)} \oplus Z$, then $I(Z; \hat{Y}|G) > I(Z; \hat{Y}) = 0$ even though there is no discrimination ($\text{CCI}(Z \rightarrow \hat{Y}) = 0$). In Lemma 4, we therefore had to include an additional criterion that $\text{CCI}(Z \rightarrow G) = 0$, or equivalently $G = g(U_X)$ for an equivalence with our proposed definition. Such a $G = g(U_X)$ may be difficult to determine in practice from observational data alone, because observational data is a function of both Z and U_X .*

Remark 6 (Conditioning to capture non-exempt masked discrimination). *Extending a similar equivalence for non-exempt masked discrimination is not straightforward. For instance, a criterion such as $I(Z; \hat{Y}|G, X_c) > I(Z; \hat{Y}|X_c)$ can lead to false positives even if $\text{CCI}(Z \rightarrow G) = 0$. E.g., suppose $\hat{Y} = U_{X_1} \oplus U_{X_2}$, $X_c = Z \oplus U_{X_2}$ and $X_g = U_{X_1}$ where Z, U_{X_1} and U_{X_2} are all independent Bern(1/2). Here, $\text{CCI}(Z \rightarrow \hat{Y}) = 0$, which suggests that there is no discrimination. However, defining $G = X_g = U_{X_1}$, we can get a false positive if we check for $I(Z; \hat{Y}|G, X_c) > I(Z; \hat{Y}|X_c)$.*

D Appendix to Section 4

D.1 Additional Results: Fairness Properties of the Observational Relaxations

Lemma 13 (Fairness Properties of $\text{Uni}(Z : \hat{Y} \setminus X_c)$). *The measure $\text{Uni}(Z : \hat{Y} \setminus X_c)$ satisfies three desirable properties, namely, 1, 2 and 4.*

Proof of Lemma 13. Property 1 is satisfied because \hat{Y} is a deterministic function of the entire X , and hence the Markov chain $Z - X - \hat{Y}$ holds. Thus $I(Z; \hat{Y} | X_c) = 0$, also implying $\text{Uni}(Z : \hat{Y} \setminus X_c) = 0$

Property 2 is trivially satisfied because the property itself requires that $\text{Uni}(Z : \hat{Y} \setminus X_c) > 0$.

For Property 4, observe that,

$$\begin{aligned}
& \text{CCI}(Z \rightarrow \hat{Y}) = 0 \\
& \implies I(Z; \hat{Y}) = 0 \\
& \implies \text{Uni}(Z : \hat{Y} \setminus X_c) + \text{Red}(Z : (\hat{Y}, X_c)) = 0 \quad \text{[by PID]} \\
& \implies \text{Uni}(Z : \hat{Y} \setminus X_c) = 0 \quad \text{[non-negativity of PID terms]}. \tag{52}
\end{aligned}$$

□

Lemma 14 (Fairness Properties of $I(Z; \hat{Y}|X_c)$). *The measure $I(Z; \hat{Y}|X_c)$ satisfies Properties 1, 2 and 3.*

Proof of Lemma 14. Property 1 is satisfied because \hat{Y} is a deterministic function of the entire X , and hence the Markov chain $Z - X - \hat{Y}$ holds.

For Property 2, observe that

$$\begin{aligned}
& \text{Uni}(Z : \hat{Y} \setminus X_c) > 0 \\
& \implies I(Z; \hat{Y}|X_c) > 0 \quad \text{[since } I(Z; \hat{Y}|X_c) \geq \text{Uni}(Z : \hat{Y} \setminus X_c)\text{]}. \tag{53}
\end{aligned}$$

Lastly, suppose there is a random variable G is of the form $G = g(X_c)$ such that $I(Z; \hat{Y}|G, X_c) > \text{Uni}(Z : \hat{Y} \setminus X_c)$. Observe that,

$$\begin{aligned} I(Z; \hat{Y}|X_c) &= I(Z; \hat{Y}|G, X_c) && \text{[because } G = g(X_c)\text{]} \\ \implies I(Z; \hat{Y}|X_c) &> \text{Uni}(Z : \hat{Y} \setminus X_c) \geq 0. \end{aligned} \tag{54}$$

Thus, the claim holds.

Lastly, Property 3 is satisfied because if the Markov chain $Z - X_c - \hat{Y}$ holds, then $I(Z; \hat{Y}|X_c) = 0$. \square

D.2 Complete results of the simulation, with standard deviations

The full tabulation of results from the simulation can be found in Table 3.

Table 3: Observations after training a classifier ($w_1X_1 + w_2X_2 + w_3X_3 + b \geq 0$) using three loss functions with different fairness criteria (100 simulations of 7000 iterations each with batch size 200)

Setup (λ)	$-\frac{w_1}{b}$ (SD.)	$-\frac{w_2}{b}$ (SD.)	$-\frac{w_3}{b}$ (SD.)	Accuracy (SD.)
Loss L_1 (N.A.)	1.083 (0.002)	1.083 (0.003)	1.075 (0.003)	98.46 (0.003)
Loss L_2 ($\lambda = 4$)	1.072 (0.017)	1.074 (0.018)	3.758 (0.15)	81.10 (0.013)
Loss L_2 ($\lambda = 10$)	1.007 (0.145)	1.026 (0.152)	13.86 (0.947)	70.17 (0.014)
Loss L_3 ($\lambda = 4$)	1.455 (0.020)	0.734 (0.012)	1.908 (0.03)	89.58 (0.009)
Loss L_3 ($\lambda = 10$)	2.048 (0.029)	0.018 (0.022)	2.57 (0.031)	80.76 (0.013)