# Fairness Under Feature Exemptions: Counterfactual and Observational Measures

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

**Abstract**

With the growing use of machine learning algorithms in highly consequential domains, the quantification and removal of bias in decision making with respect to protected attributes, such as gender, race, etc., is becoming increasingly important. While quantifying bias is essential, sometimes the needs of a business (e.g., hiring) may require the use of certain features that are critical in a way that any bias that can be explained by them might need to be exempted. For instance, in hiring a software engineer, a standardized coding-test score may be a critical feature that is weighed strongly in the decision even if it introduces bias, whereas other features, such as name, zip code, or reference letters may be used to improve decision-making, but only to the extent that they do not add bias. In this work, we propose a novel information-theoretic decomposition of the total bias (a quantification inspired from counterfactual fairness) into two components: a non-exempt component which quantifies the part of the bias that cannot be accounted for by the critical features, and an exempt component which quantifies the remaining bias. This decomposition is important: it allows one to check if the bias arose purely due to the critical features (inspired from the business necessity defense of disparate impact law) and also enables selective removal of the non-exempt component of bias if desired. We arrive at this decomposition through canonical examples that lead to a set of desirable properties (axioms) that any measure of non-exempt bias should satisfy. We then demonstrate that our proposed counterfactual measure of non-exempt bias satisfies all of them. Our quantification bridges ideas of causality, Simpson's paradox, and 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 bias 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 perform case studies to show how one can audit existing models as well as train new models while reducing non-exempt bias.

## I. INTRODUCTION

As artificial intelligence becomes ubiquitous, it is important to understand whether the output of a machine-learnt model is unfairly biased with respect to *protected attributes* such as gender, race, etc., and if so, how we can engineer fairness into such a model. The field of fair machine learning provides several measures for fairness [2]–[28], and uses them to reduce bias, e.g., as a regularizer during training [6], [10]. In several 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 give rise to bias. Examples of such critical features might be weightlifting ability for a firefighter's job, educational qualification for an academic job, coding skills for a software engineering job, merit and seniority in deciding salary, etc. In an attempt to preserve the importance of the critical features in the decision making, one might choose to exempt the bias created by them. On the other hand, racial bias in mortgage lending decisions arising due to zip code (a non-critical feature) [29], or bias in promotion/transfer decisions arising from aptitude tests[1] are examples of non-exempt bias. In this work, our goal is to formalize and quantify the *non-exempt bias*, i.e., the part of the bias that cannot be accounted for by the critical features. This quantification is important for two reasons: (i) it allows one to check if the bias arose purely due to the critical features (inspired from the "business necessity defense" in the disparate impact law, i.e., Title VII of the Civil Rights Act of 1964 [31]); and (ii) it enables selective removal of the non-exempt component if desired.

In this work, we assume that the critical features or business necessities are known (similar to [4], [17]; this discussion is revisited in Section VIII). We let $X_c$ and $X_g$ denote the critical and the non-critical (or general) features, and $X$ denote the entire set of features. We also denote the protected attribute(s) by $Z$, the true label by $Y$, and the model output by $\hat{Y}$ which is a function of the entire feature vector $X$. While we acknowledge that such categorization of features is application-dependent and might require domain knowledge and ethical evaluation, such exemptions do exist in law. E.g., the US Equal Pay Act [32] exempts for difference in salary based on gender that can be explained by merit and seniority. Similarly, the US employment discrimination law contains a business necessity defense [30] where bias about protected attributes may be exempted if the bias can be justified as "necessary to the normal operation of that particular business." For example, a standardized coding-test score may be a critical feature in hiring software engineers. Similarly, weightlifting ability might be a critical feature in hiring firefighters so that they are able to carry fire victims out of a burning building. The critical feature is therefore required to be weighed strongly in hiring even if it is correlated with some protected attributes.

[1] In the landmark employment discrimination court-case of Griggs v. Duke Power [30], the US Supreme Court deemed certain aptitude tests as not job-related and hence not business necessities, ruling against the employer.

TABLE I: Observational Measures ($M_{NE}$) of Non-Exempt Bias (Utility and Limitations)

| | Desirable Properties | $\text{Uni}(Z:\hat{Y} \mid X_c)$ | $I(Z;\hat{Y} \mid X_c)$ | $I(Z;\hat{Y} \mid X_c, X')$ |
|---|---|---|---|---|
| 1. | No counterfactual causal influence from $Z$ to $\hat{Y} \Rightarrow M_{NE} = 0$. | Yes | Not Always | Not Always |
| 2. | $M_{NE}$ detects unique information about $Z$ in $\hat{Y}$ not in $X_c$. | Yes | Yes | Not Always |
| 3. | $M_{NE}$ detects non-exempt masked bias. | No | Masked by $g(X_c)$ | Masked by $g(X_c, X')$ |
| 4. | $M_{NE}$ is equal to total bias (in a counterfactual sense) if $X_c = \phi$ and $X_g = X$. | No | No | No |
| 5. | $M_{NE}$ is non-increasing as more features are added to $X_c$ from $X_g$. | Yes | No | No |
| 6. | $M_{NE}$ is 0 (complete exemption) if $X_c = X$ and $X_g = \phi$. | Yes | Yes | Yes |

*Why should we use the "general" features at all for prediction if they are not critical?* General features can improve performance metrics such as accuracy of the model, or even help reduce the candidate pool, e.g., if 60% applicants clear a test, but resources are available to interview only 10%. Not using the general features at all can reduce accuracy, or produce a very large candidate pool. In this work, our proposition is to use both critical and general features in a way that maximizes accuracy (to the extent possible) while preventing non-exempt bias. For instance (inspired from [31]), to choose a "good" employee, an employer could evaluate standardized test scores and also reference letters (human-graded performance reviews). All these features are "job-related" in that they have statistical correlation with the prediction goal, and can help improve the accuracy. However, test scores, a critical feature, may need to be weighed strongly in the decision, *even if* biased, whereas, reference letters may be used only to the extent that they do not discriminate.

This work treads a middle ground between two popular measures of fairness that do not use domain knowledge, namely, *statistical parity* [3], [6], [12], [27], which enforces the criterion $Z \perp\!\!\!\perp \hat{Y}$, and *equalized odds* [7], [12], [27], which enforces $Z \perp\!\!\!\perp \hat{Y}|Y$ (directly or through practical relaxations). Our selective quantification of non-exempt bias (using domain knowledge to identify critical features) helps address one of the major criticisms against statistical parity. The criticism is that it can lead to the selection of unqualified members from the protected group [7], [22], e.g., by disregarding the critical features if they are correlated with the protected attribute $Z$. In fact, in our case study in Section VII, we observe that the weight of the critical feature is significantly reduced in the decision making when one uses statistical parity as a regularizer with the loss function because the critical feature is correlated with $Z$ (also see Canonical Example 1 in Section III-C). On the other hand, equalized odds suffers from label bias [26], [29], [33], [34] because it is based on agreement with the true labels. In fact, we demonstrate (Canonical Example 2 in Section III-C) that if the historic labels themselves reinforce bias from the non-critical features, then even if we obtain a perfect classifier after training on the historic data, which satisfies equalized odds, it can reinforce undesirable non-exempt bias[2].

## A. Contributions

Our main contribution in this work is the quantification of non-exempt bias based on a rigorous axiomatic approach. As a first step towards this quantification, we propose an information-theoretic quantification (see Definition 4 in Section II-B) of the total bias (exempt and non-exempt) that is 0 if and only if the model is *counterfactually fair* [16]. Counterfactual fairness [16], [18] is a causal notion of fairness where the features $X$, the protected attribute $Z$ and the model output $\hat{Y}$ are assumed to be observables in a Structural Causal Model (SCM) (defined formally in Section II; see Definition 2). The model is deemed *counterfactually fair* if $Z$ has no *counterfactual causal influence* on $\hat{Y}$, i.e., $\hat{Y}$ does not change if we are able to vary $Z$ in the SCM in a manner that other independent latent factors remain constant (defined formally in Section II; see Definition 3).

Interestingly, note that the total bias (in a counterfactual sense) may not exhibit itself entirely in the mutual information $I(Z;\hat{Y})$, which is the *statistically visible information*[3] about $Z$ in $\hat{Y}$, because of "statistical masking effects" (also relates to Simpson's paradox [35]). Consider an example inspired from [16], [20], [26] where a software engineering job advertisement is shown only to a) men with coding skills above a threshold, and b) women with coding skills below a threshold. That is, the decision $\hat{Y} = Z \oplus G$ where $\oplus$ denotes XOR, $G$ is the binary variable denoting whether coding skills are above a threshold (that does not have a causal influence of $Z$ in this example), and $G, Z$ are i.i.d. Bern($\frac{1}{2}$). This decision is biased against the high-skilled women for whom the ad is relevant, but $I(Z;\hat{Y}) = 0$ here, thus failing to capture this bias. Intuitively, our quantification of total bias also extends the idea of *proxy-use* [20] from *white-box models*[4] to black-box models. Proxy-use [20] examines "white-box" models, i.e., models with clearly defined constituents (e.g., decision trees) and regards a model as having bias if (i) there is a constituent that has high mutual information about $Z$ (a proxy of $Z$); and (ii) this constituent also causally influences the output $\hat{Y}$ (i.e., varying the constituent while keeping other constituents constant does not change the output). In this work, the total bias captures the intuitive notion of a virtual constituent or proxy of $Z$ that causally influences the final

---

[2]Our quantification does not use the true labels for fairness (unlike equalized odds), addressing the criticism in [31] which says that " [...] often the best labels for different classifications will be open to debate."

[3]This is a quantification of bias inspired from statistical parity which deems a model fair if and only if $\hat{Y} \perp\!\!\!\perp Z$. Note that, $I(Z;\hat{Y}) = 0$ if and only if $\hat{Y} \perp\!\!\!\perp Z$.

[4]White-box models [20] are the type of models where one can clearly explain how they behave, how they produce predictions and what the influencing variables or sub-components of the model are, e.g., decision trees, linear regression, etc.

output $\hat{Y}$ (this intuition is revisited to understand Scenario 2 in Section II-B). For instance, a virtual constituent $Z$ is formed in the example of masked bias in ads that causally influences $\hat{Y}$ even though $I(Z; \hat{Y}) = 0$.

Next, we quantify the *non-exempt* part of this total bias, i.e., the part that cannot be explained by the critical features $(X_c)$. Building on the extension of proxy-use [20] for black-box models as discussed above, we aim to quantify the influence of a discriminatory virtual constituent or proxy of $Z$, if formed inside the black-box model, on the model output $\hat{Y}$, and that cannot be attributed entirely to the critical features (this idea is revisited for an intuitive understanding of the canonical examples in Section II-B.). To quantify this *non-exempt bias*, we consider toy examples and thought experiments to first arrive at a set of desirable properties (axioms) that any measure of non-exempt bias should satisfy, and then provide a measure that satisfies them (see Theorem 1). These desirable properties can be intuitively described as follows. If the model is counterfactually fair, e.g., if the virtual constituents or proxies of $Z$ cancel each other leading to a final model output that has no counterfactual causal influence of $Z$, then it is desirable that the non-exempt bias is also 0. Next, it is desirable that the measure be non-zero if $\hat{Y}$ has any "unique" statistically visible 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 bias* that needs to be captured, e.g., in the aforementioned example of software-engineering-job ads (also revisited in Canonical Example 4 in Section III-B where we discuss our rationale for the properties). The next three properties are more intuitive. If all the features are in the non-critical set, then the measure should be equal to the total bias since no bias is exempt. For a fixed set of features $X$ and a fixed model, as more features become categorized as critical, the measure of non-exempt bias should not increase, i.e., it either decreases or stays the same. Ultimately, if all the features are in the critical set $X_c$, then we require the measure of non-exempt bias to be 0 since then the total bias is exempt.

Our proposed measure of non-exempt bias, that satisfies all these desirable properties, is *counterfactual* in nature, i.e., it depends on the true SCM, and hence, is not *observational*[5] in general. We also show the theoretical impossibility of any observational measure in satisfying all the desirable properties together (see Theorem 3). We note that in some applications, counterfactual measures can be realized or approximated with assumptions on the causal model. However, for more general use in practical applications, we also propose several observational relaxations of our measure that satisfy only some of these properties. Nevertheless, we believe that a counterfactual measure and its properties are crucial in understanding the utility and the limitations of different observational measures and informing which measure to choose in practice (summarized in Table I; detailed discussion in Section VI).

To summarize, our contributions in this work are as follows:

**1. Quantification of Non-Exempt Bias**: We propose a novel counterfactual measure of non-exempt bias that captures the bias that cannot be explained by the critical features. Our quantification attempts to capture the intuitive notion of whether a discriminatory virtual constituent or proxy [20] of $Z$ is formed inside the black-box model that influences the output $\hat{Y}$ and that cannot be attributed entirely to the critical features $(X_c)$. We adopt a rigorous axiomatic approach where we first arrive at a set of desirable properties that any measure of non-exempt bias should satisfy by analyzing several canonical examples (thought experiments). Next, we show that the proposed measure satisfies these properties (see Theorem 1). Our quantification leverages a body of work in information theory called Partial Information Decomposition (PID), as well as, causality.

**2. Overall Decomposition of Total Bias into Statistically Visible and Masked components**: Our quantification finally leads us to an overall decomposition of the total bias into four non-negative components, namely, exempt and non-exempt *statistically visible* bias and exempt and non-exempt *masked* bias (see Theorem 2). The exempt and non-exempt *statistically visible* biases add up to give $I(Z; \hat{Y})$ which is the total statistically visible bias.

**3. An Impossibility Result**: We show that no purely observational measure of non-exempt bias 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 I) and then use them in case studies to demonstrate how to (i) audit existing models; and also (ii) train new models that selectively reduce non-exempt bias.

**Our contribution in the context of related works:** Causal approaches for fairness have been explored in [16]–[20], [36], [37], including impossibility results on purely observational measures [17], [20]. Our main novelty lies in using a rigorous axiomatic approach based on realistic examples and thought experiments for quantifying non-exempt and exempt bias separately, thereby allowing for exemptions due to critical features. The decomposition of total bias into exempt and non-exempt components is tricky. For instance, following the ideas of path-specific counterfactual fairness [19], 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 (exempt and non-exempt) measures. However, we provide a counterexample (see Canonical Example 6 in Section III-B) to show that bias can also arise from synergistic information 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 be. Purely causal measures (that do not rely on the PID framework) can attribute such bias entirely to $X_c$. We contend that such synergistic information, if influencing the decision,

---

[5]Observational measures are those that can be estimated from the probability distribution of the data without knowledge of the underlying SCM.

must be included in the *non-exempt* component of bias because both $X_c$ and $X_g$ are contributors. We note that identifying such synergy is important: synergy arises frequently in machine-learning and other related applications [35], [38], [39].

Some observational measures for quantifying non-exempt bias have been introduced previously in [2], [4] where the authors propose a decomposition of statistically visible discrimination (statistical parity) into explainable and non-explainable components (see also subsequent works [5], [40]–[42] that build on this idea). They examine the difference in the expected model output ($\hat{Y}$) for candidates of different races/genders ($Z$) after conditioning on specific subsets of features[6] (this relates to dependence between $Z$ and $\hat{Y}$ after conditioning on specific features; also referred to as conditional statistical parity [40]). In this context, in this work, we provide simple yet relevant counterexamples showing that conditioning may not always faithfully capture non-exempt bias. E.g., Canonical Example 3 in Section III-B) is deemed *unfair* by conditional mutual information (or conditional statistical parity), but is *fair* by counterfactual fairness [16], [18]. We use these examples as motivation to decompose conditional mutual information into unique and synergistic information using PID, separating two kinds of "statistical dependence" which conditioning alone fails to do (see Section II-A). We refer to Section III-C for more detailed discussion on existing measures that have some provision for exemption, namely, conditional statistical parity [40], [42], justifiable fairness [41], as well as a related causal measure of path-specific counterfactual fairness [19]. Our problem also differs from *sub-group fairness* [26] where the sub-populations in consideration are based on the protected attributes alone, e.g., $Z = (Z_1, Z_2)$ with $Z_1$ being gender, and $Z_2$ being race, and does not consider exemptions with respect to the other (non-protected) attributes.

We also note that the idea of using correlation-based observational approximations of bias (e.g., correlation between $Z$ and $\hat{Y}$ to represent statistical parity) as a regularizer during training has been proposed earlier [10]. In this context, our main contribution here is on first arriving at a measure of non-exempt bias (that happens to be non-observational), and then proposing 3 observational measures for applications in *both* auditing existing models and training new models with reduced non-exempt bias. For auditing, we use alternate non-correlation-based estimators for unique information, mutual information, and conditional mutual information from the `dit` package [43]. For training, we rely on simplistic correlation-based approximations for mutual information and conditional mutual information along the lines of [10] for ease of computation. For unique information, we introduce novel correlation-based regularizers for training in Section VII, leveraging a Gaussian approximation for PID [44].

### B. Paper Outline

The rest of the paper is organized as follows. Section II introduces the background, system model and assumptions underlying our problem formulation, i.e., how to quantify the non-exempt bias. Section III-A first states all the desirable properties that a measure of non-exempt bias should satisfy, and then introduces our proposed counterfactual measure that satisfies all of them (Theorem 1 in Section III-A). This is followed by a rationale behind the desirable properties through canonical examples and thought experiments in Section III-B. We also discuss the utility and limitations of some existing measures, namely, path-specific counterfactual fairness [19], conditional statistical parity [40], and justifiable fairness [41] in Section III-C. Next, Section IV provides insights on the overall decomposition of the total bias (in a counterfactual sense) into exempt and non-exempt components, with each of them being further decomposed into *statistically visible* and *masked* components (Theorem 2 in Section IV). Section V provides an impossibility result on observational measures, stating that no observational measure can satisfy all of the desirable properties. Nonetheless, since counterfactual measures are often difficult to realize in practice, we propose several observational relaxations of our proposed counterfactual measure in Section VI (that only satisfy some of the desirable properties), and discuss their utility and limitations. Next, in Section VII, we use our observational measures to conduct case studies on both artificial and real datasets to demonstrate practical application in training. Finally, we conclude with a discussion in Section VIII.
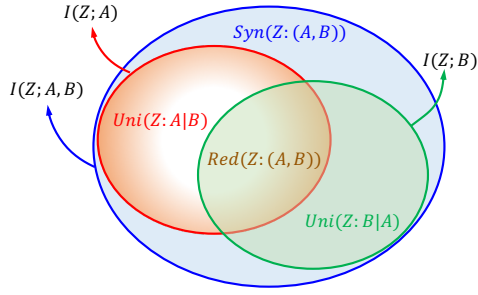
## II. PRELIMINARIES

Here, we first provide a brief background on Partial Information Decomposition (PID) in Section II-A to help follow the paper. Appendix B provides more details on the specific properties used in the proofs. Next, we introduce our system model and assumptions in Section II-B. We use the following notations: (i) $X = (X_1, X_2, \ldots, X_n)$ denotes a tuple [45], i.e., an ordered set of elements $X_1, X_2, \ldots, X_n$; (ii) $\phi$ denotes the empty tuple (no elements); (iii) For tuple with a single element, the bracket is omitted for brevity, i.e., $(X_1) = X_1$; (iv) $(X, A)$ is equivalent to the new tuple $(X_1, X_2, \ldots, X_n, A)$ formed by appending the element $A$ at the end of tuple $X$; (v) $X_1 \in X$ means $X_1$ is an element of tuple $X$; (vi) $S \subseteq X$ means the set of elements in tuple $S$ form a subset of the set of elements in tuple $X$; and (vii) $X \backslash X_2$ denotes a new tuple formed by removing element $X_2$ from $X$ without changing the order of other elements, i.e., $(X_1, X_3, X_4, \ldots, X_n)$.

### A. Background on Partial Information Decomposition (PID)

The PID framework [46]–[48] decomposes the mutual information $\mathrm{I}(Z; (A, B))$ about a random variable $Z$ contained in the tuple $(A, B)$ into four *non-negative* terms as follows (also see Fig. 1):

$$\mathrm{I}(Z; (A, B)) = \mathrm{Uni}(Z : A|B) + \mathrm{Uni}(Z : B|A) + \mathrm{Red}(Z : (A, B)) + \mathrm{Syn}(Z : (A, B)). \tag{1}$$

---

[6]Conditional mutual information (conditioned on the critical feature(s)) as a measure of non-exempt bias has surfaced in [42] with a focus on novel estimators.

(a) Venn diagram showing PID of $I(Z; (A, B))$

| | $I(Z; A\|B)$ | $I(Z; B)$ | $I(Z; A, B)$ |
|---|---|---|---|
| $I(Z; A)$ | $Uni(Z : A\|B)$ | $Red(Z : (A, B))$ | |
| $I(Z; B\|A)$ | $Syn(Z : (A, B))$ | $Uni(Z : B\|A)$ | |
| $I(Z; A, B)$ | | | |

(b) Tabular Representation of PID of $I(Z; (A, B))$

Fig. 1: Mutual information $I(Z; (A, B))$ is decomposed into 4 non-negative terms, namely, $Uni(Z : A|B)$, $Uni(Z : B|A)$, $Red(Z : (A, B))$ and $Syn(Z : (A, B))$. Also note that, $I(Z; (A, B)) = I(Z; B) + I(Z; A \mid B)$, each of which is in turn a sum of two PID terms. $Red(Z : (A, B))$ is the sub-volume between $I(Z; A)$ and $I(Z; B)$, and $Uni(Z : A|B)$ is the sub-volume between $I(Z; A \mid B)$ and $I(Z; A)$.

Here, $Uni(Z : A|B)$ denotes the unique information about $Z$ that is present only in $A$ and not in $B$. Likewise, $Uni(Z : B|A)$ is the unique information about $Z$ that is present only in $B$ and not in $A$. The term $Red(Z : (A, B))$ denotes the redundant information about $Z$ that is 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 $(A, B)$. *All four of these terms are non-negative. Also notice that,* $Red(Z : (A, B))$ *and* $Syn(Z : (A, B))$ *are symmetric in $A$ and $B$.* Before defining these PID terms formally, let us understand them through an intuitive scenario.

**Scenario 1** (Understanding Partial Information Decomposition). *Let $Z = (Z_1, Z_2, Z_3)$ with $Z_1, Z_2, Z_3 \sim$ i.i.d. Bern(½). Let $A = (Z_1, Z_2, Z_3 \oplus N)$, $B = (Z_2, N)$, $N \sim$ Bern(½) is independent of $Z$. Here, $I(Z; (A, B)) = 3$ bits.*

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 $Uni(Z : A|B) = I(Z; Z_1) = 1$ 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 $A$, i.e., $Uni(Z : B|A) = 0$.

Irrespective of the formal definition of these individual terms, the following identities also hold (see Fig. 1b):

$$I(Z; A) = Uni(Z : A|B) + Red(Z : (A, B)). \tag{2}$$

$$I(Z; A \mid B) = Uni(Z : A|B) + Syn(Z : (A, B)). \tag{3}$$

**Remark 1** (An Interpretation of PID as Information-Theoretic Sub-Volumes). *Equations (1), (2) and (3) have been represented in a tabular fashion in Fig. 1b. Notice that, $Uni(Z : A|B)$ can be viewed as the information-theoretic sub-volume of the intersection between $I(Z; A)$ and $I(Z; A \mid B)$. Similarly, $Red(Z : (A, B))$ is the sub-volume between $I(Z; A)$ and $I(Z; B)$.*

These equations also demonstrate that $Uni(Z : A|B)$ and $Red(Z : (A, B))$ are the information contents that exhibit themselves in $I(Z; A)$ which is the statistically visible information content about $Z$ present in $A$. Because both these PID terms are non-negative, if any one of them is non-zero, we will have $I(Z; A) > 0$. Similarly, $Uni(Z : B|A)$ and $Red(Z : (A, B))$ also exhibit themselves in $I(Z; B)$. On the other hand, $Syn(Z : (A, B))$ is the information content that does not exhibit itself in $I(Z; A)$ or $I(Z; B)$ individually, i.e., these terms can still be 0 even if $Syn(Z : (A, B)) > 0$. But, $Syn(Z : (A, B))$ exhibits itself in $I(Z; (A, B))$. Notice that,

$$I(Z; (A, B)) = \underbrace{Uni(Z : A|B) + Red(Z : (A, B))}_{I(Z; A)} + \underbrace{Uni(Z : B|A) + Syn(Z : (A, B))}_{I(Z; B|A)}$$
$$= \underbrace{Uni(Z : B|A) + Red(Z : (A, B))}_{I(Z; B)} + \underbrace{Uni(Z : A|B) + Syn(Z : (A, B))}_{I(Z; A|B)}.$$

Given three independent equations (1), (2) and (3) in four unknowns (the four PID terms), defining any one of the terms (e.g., $Uni(Z : A|B)$) is sufficient to obtain the other three. For completeness, we include the definition of unique information from [46] (that also allows for estimation via convex optimization [49]) with the specific properties used in the proofs in Appendix B. To follow the paper, only an intuitive understanding is sufficient.

**Definition 1** (Unique Information [46]). *Let $\Delta$ be the set of all joint distributions on $(Z, A, B)$ and $\Delta_p$ 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,* $\text{Uni}(Z : A|B) = \min_{Q \in \Delta_p} \text{I}_Q(Z; A \mid B)$, *where* $\text{I}_Q(Z; A \mid B)$ *is the conditional mutual information when* $(Z, A, B)$ *have joint distribution $Q$.*

The key intuition behind this definition is that the unique information should only depend on the marginal distribution of the pairs $(Z, A)$ and $(Z, B)$. This is motivated from an **operational** perspective that if $A$ has unique information about $Z$ (with respect to $B$), then there must be a situation where one can predict $Z$ better using $A$ than $B$ (more details in [46, Section 2]). Therefore, all the joint distributions in the set $\Delta_p$ with the same marginals essentially have the same unique information, and the distribution $Q^*$ that minimizes $\text{I}_Q(Z; A \mid B)$ is the joint distribution that has no synergistic information leading to $\text{I}_{Q^*}(Z; A \mid B) = \text{Uni}(Z : A|B)$. Definition 1 also defines $\text{Red}(Z : (A, B))$ and $\text{Syn}(Z : (A, B))$ using (2) and (3).

*B. System Model and Assumptions*

Here, we will introduce our system model and assumptions. We start with an introduction to Structural Causal Model (SCM).

**Definition 2** (Structural Causal Model: $\text{SCM}(U, V, \mathcal{F})$ [35]). *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) $\mathcal{F} = (F_1, F_2, \ldots)$, 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 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.*
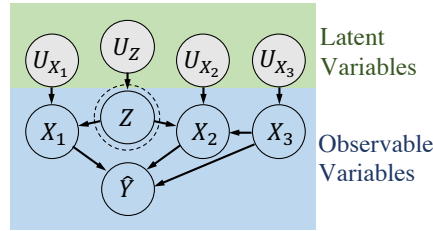


Fig. 2: An SCM with protected attribute $Z$, features $X = (X_1, X_2, X_3)$, and output $\hat{Y}$. Here $X$ and $\hat{Y}$ are the observables, and $U_Z$ and $U_X = (U_{X_1}, U_{X_2}, U_{X_3})$ are the latent social factors. $Z$ does not have any parents in the SCM and $\hat{Y}$ is completely determined by $X = (X_1, X_2, X_3)$.

**Our System Model:** For our problem, consistent with several other works on fairness [16], [17], [19], the latent variables $U$ represent possibly unknown social factors. The observables $V$ consist of the protected attributes $Z$, the features $X$ 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 graph). 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 ($U_{X_j}$) and its parents, which are again functions of their own latent variables and parents. Therefore, each $X_j$ can also be written as $f_j(Z, U_X)$ for some deterministic $f_j(\cdot)$, where $U_X = U \setminus U_Z$ denotes the latent factors in $U$ that do not cause $Z$ (see a formal proof in [35, Proposition 6.3]). Here, $f_j(\cdot)$ may be constant in some of its arguments. This claim holds because the underlying graph is acyclic, and hence the structural assignments of the ancestors of $X_j$ can be substituted recursively into one another until all observables except $Z$ are substituted by latent variables. Also note that, $Z \perp\!\!\!\perp U_X$. A model takes $X$ (which consists of critical features $X_c$ and general features $X_g$) as its input and produces an output $\hat{Y}$ which is a deterministic function of $X$, *i.e.*, $\hat{Y} = r(X)$ where $X$ is itself a deterministic function of $(Z, U_X)$. Therefore, $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$.

Next, we introduce the concept of Counterfactual Causal Influence (CCI) ( [16], [18], [50]–[54]), which will help us understand the well-known causal definition of fairness called *counterfactual fairness* [16].

**Definition 3** (Counterfactual Causal Influence: $\text{CCI}(Z \to \hat{Y})$). *Consider the aforementioned system model. Let $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$ where $U_X$ are latent variables that do not cause $Z$ in the true SCM. Then,*

$$\text{CCI}(Z \to \hat{Y}) = \mathbb{E}_{Z, Z', U_X} \left[|h(Z, U_X) - h(Z', U_X)|\right] \text{ where } Z', Z \text{ are i.i.d.}$$

Counterfactual causal influence quantifies the change in $\hat{Y} = h(Z, U_X)$ if we only vary $Z$ while keeping the other latent factors ($U_X$) unchanged. A model is said to satisfy *counterfactual fairness* [16], [18] if and only if the output $\hat{Y}$ has no counterfactual causal influence of $Z$ (we formally derive that $\text{CCI}(Z \to \hat{Y}) = 0$ is equivalent to counterfactual fairness [16] in Lemma 6 in Appendix A-B). What this means is that a model is *counterfactually fair* if and only if the output $\hat{Y} = h(Z, U_X)$ does

(a) Model is not counterfactually fair as $\hat{Y}$ has counterfactual causal influence of $Z$.

(b) Model is counterfactually fair after cancelling out the influence of $Z$ from $X_1$.

(c) Model is counterfactually fair even though it uses an entirely unrelated feature.
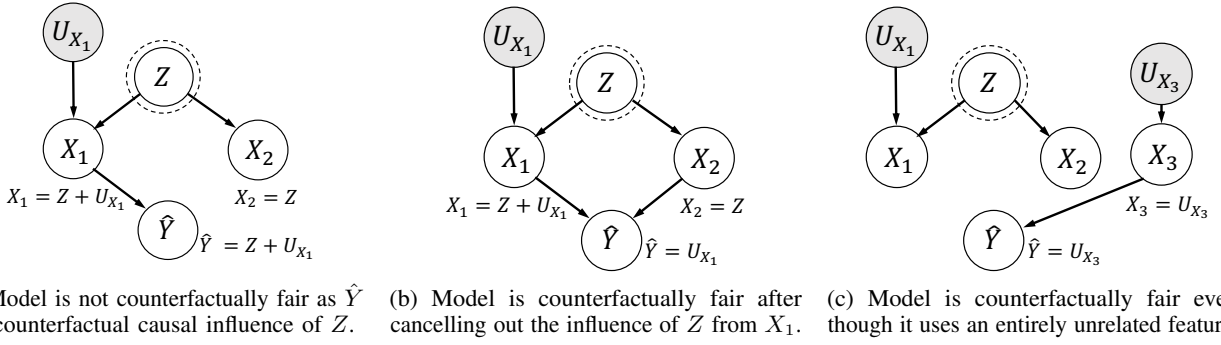
Fig. 3: Illustration of Scenario 2 for understanding the concept of counterfactual fairness: Different models are used to make hiring decisions on data corresponding to the same SCM with $Z$ denoting the protected attribute, $U_{X_1}$ denoting inner ability, $X_1 = Z + U_{X_1}$ denoting interview score, and $X_3$ denoting an alternate feature, e.g., location.

not change with $Z$ while keeping the other latent factors ($U_X$) unchanged. It captures the intuitive notion that no virtual constituent or proxy of $Z$ influences the output (inspired from the work on proxy-use [20]). In other words, $\hat{Y} \perp\!\!\!\perp Z|U_X$ (proved in Lemma 1), i.e.,

$$\Pr(\hat{Y} = y|Z = z, U_X = u_x) = \Pr(\hat{Y} = y|Z = z', U_X = u_x) \ \forall z, z', y, u_x.$$

This notion of fairness also leads us to propose an information-theoretic quantification of total bias (exempt and non-exempt) that is 0 if and only if the counterfactual causal influence of $Z$ on $\hat{Y}$ is 0 (equivalence is demonstrated in Lemma 1 with the proof in Appendix A-A).

**Definition 4** (Total Bias). *The total bias in a model is defined as* $\mathrm{I}(Z; (\hat{Y}, U_X))$.

Notice that, $\mathrm{I}(Z; (\hat{Y}, U_X)) = \mathrm{I}(Z; \hat{Y}|U_X) + \underbrace{\mathrm{I}(Z; U_X)}_{=0 \text{ since } Z \perp\!\!\!\perp U_X} = \mathrm{I}(Z; \hat{Y}|U_X)$.

**Lemma 1** (Equivalences of CCI). *Consider the aforementioned system model. Let* $\hat{Y} = h(Z, U_X)$ *for some deterministic function* $h(\cdot)$ *and* $Z \perp\!\!\!\perp U_X$. *Then,* $\mathrm{CCI}(Z \rightarrow \hat{Y}) = 0$ *if and only if* $\mathrm{I}(Z; (\hat{Y}, U_X)) = 0$.

**Remark 2** (Advantage of our Information-Theoretic Quantification). *One might wonder why such an information-theoretic quantification of counterfactual causal influence (or, total bias) is necessary. The information-theoretic quantification of total bias enables analytical decomposition into exempt and non-exempt components that better satisfy our intuitive understanding. Our non-exempt bias intuitively attempts to capture whether discriminatory proxies are formed inside the black-box model that cannot be entirely attributed to the critical features* $X_c$. *The decomposition of counterfactual causal influence (Definition 3) into exempt and non-exempt components is not straightforward. For instance, following the ideas of path-specific counterfactual fairness [19], 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, bias can also arise from synergistic information about* $Z$ *in both* $X_c$ *and* $X_g$, *that cannot be attributed to any one of them alone, i.e.,* $\mathrm{I}(Z; X_c)$ *and* $\mathrm{I}(Z; X_g)$ *may both be 0 but* $\mathrm{I}(Z; X_c, X_g)$ *may not be (see Canonical Example 6). Purely causal measures can attribute such bias entirely to* $X_c$. *We contend that such synergistic information, if influencing the decision, must be included in the non-exempt component of bias because both* $X_c$ *and* $X_g$ *are contributors to the proxy. Information-theoretic equivalences of other existing notions of fairness, e.g., statistical parity, equalized odds, etc. have also been used in the broader literature on fairness [8], [10], [12], [27], [55].*

For a better understanding of counterfactual fairness, we now consider an intuitive scenario (inspired from [16]).

**Scenario 2** (Understanding Counterfactual Fairness). *Suppose a company makes its decisions about hiring based on a feature* $X_1$ *which denotes an interview score. In the SCM, this feature* $X_1 = Z + U_{X_1}$ *where* $Z$ *denotes the protected attribute and* $U_{X_1}$ *denotes the inner ability which is independent of* $Z$. *An output* $\hat{Y} = X_1$ *is not counterfactually fair because it has counterfactual causal influence of the protected attribute* $Z$ *(Fig. 3a). The total bias* $\mathrm{I}(Z; (\hat{Y}, U_X))$ *is also non-zero, capturing the intuitive notion that a proxy of* $Z$ *influences the output. On the other hand, suppose the model now uses another feature* $X_2 = Z$ *and produces the output* $\hat{Y} = X_1 - X_2 = U_{X_1}$. *This model is now deemed counterfactually fair (Fig. 3b), and its total bias* $\mathrm{I}(Z; (\hat{Y}, U_X))$ *is zero. No proxy of* $Z$ *influences the output any longer.*

**Remark 3** (Accuracy vs Counterfactual Fairness). *The goals of fairness and accuracy on a given dataset are not always aligned [9], [56]. For instance, suppose the model in Scenario 2 takes decisions only based on a new feature* $X_3 = U_{X_3}$ *that is derived entirely from some latent factor that is unrelated with the ability to perform the job (see Fig. 3c). Or, even worse,*
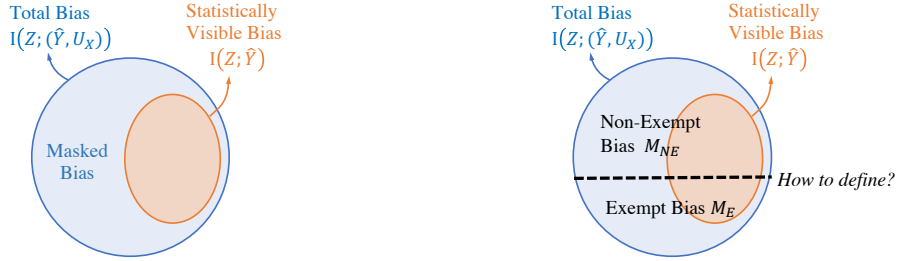
Fig. 4: Decomposition of Total Bias: (Left) Total bias (information-theoretic quantification of counterfactual causal influence) is shown in blue. The statistically visible bias and masked bias are two sub-components of the total bias. (Right) Our goal is to decompose the total bias into exempt and non-exempt components.

*suppose a model is hiring based on a random coin flip. Such a model may be highly inaccurate and absurd but it is still counterfactually fair because it has no counterfactual causal influence of $Z$ and does not cause any disparate impact with respect to $Z$. In this work, we will assume that a model has absolutely no bias (exempt or non-exempt) if and only if there is no counterfactual causal influence of $Z$ on $\hat{Y}$. We will also run into some toy examples that might have lower accuracy, but from a fairness-point-of-view, it will be desirable that they are deemed fair if there is no counterfactual causal influence of $Z$.*

Next, we propose two definitions, namely, statistically visible bias and masked bias. Statistically visible bias is an information-theoretic quantification inspired from a well-known observational definition of fairness called *statistical parity* [3], [6].

**Definition 5** (Statistically Visible Bias). *The statistically visible bias in a model is defined as* $\mathrm{I}(Z; \hat{Y})$.

Statistical parity deems a model fair if and only if $Z \perp\!\!\!\perp \hat{Y}$, i.e.,

$$\Pr(\hat{Y} = y | Z = z) = \Pr(\hat{Y} = y | Z = z') \quad \forall y, z, z'.$$

Thus, a model is said to be fair by statistical parity if and only if its statistically visible bias $\mathrm{I}(Z; \hat{Y}) = 0$.

**Remark 4** (Statistical Parity vs Counterfactual Fairness). *Statistical parity (or independence) does not imply absence of causal effects. E.g., consider $\hat{Y} = Z \oplus U_X$ where $Z, U_X \sim$ i.i.d. Bern(½). Here, $\hat{Y} \perp\!\!\!\perp Z$, but $Z$ still has a causal effect on $\hat{Y}$. If we vary $Z$ keeping all other sources of randomness in $\hat{Y}$ constant (i.e., fixing $U_X = u_x$), then $\hat{Y}$ also varies. This is, in fact, an example of masked bias, where $\mathrm{I}(Z; \hat{Y}) = 0$, but $Z$ has counterfactual causal influence on $\hat{Y}$.*

**Definition 6** (Masked Bias). *The masked bias in a model is defined as* $\mathrm{I}(Z; (\hat{Y}, U_X)) - \mathrm{I}(Z; \hat{Y})$.

The masked bias is the difference between the total bias and the statistically visible bias. Notice that, $\mathrm{I}(Z; \hat{Y}, U_X) - \mathrm{I}(Z; \hat{Y}) = \mathrm{I}(Z; U_X \mid \hat{Y})$, implying that masked bias is non-negative. We will revisit masked bias in Section IV.

**Goal:** In this work, $\mathrm{I}(Z; (\hat{Y}, U_X))$ will serve as our *information-theoretic quantification of the total bias (exempt and non-exempt)* as we discussed in Definition 4 (also recall Lemma 1 and Remark 2). Our *goal* is to appropriately decompose the total bias $\mathrm{I}(Z; (\hat{Y}, U_X))$ into an exempt component ($M_E$) and a non-exempt component ($M_{NE}$), which can and cannot be explained by the critical features $X_c$ (also see Fig. 4). Intuitively, the total bias captures the idea of a virtual constituent or proxy of $Z$ that has a causal influence on the output $\hat{Y}$. We would like the exempt and non-exempt components of total bias to be able to capture and mathematically quantify our intuitive notion of what part of the virtual constituent or proxy can and cannot be attributed to the critical features $X_c$ alone.

Before proceeding further, we also clarify our terminology here. We say that there is *no bias* when $\mathrm{I}(Z; \hat{Y}, U_X) = 0$. Alternately, we call *the bias to be exempt* if only the non-exempt component is 0, though $\mathrm{I}(Z; \hat{Y}, U_X)$ may be zero or non-zero. Table II summarizes all the important notations to help follow the rest of the paper.

TABLE II: Summary of Notations

| Symbol | Description | Observable or Not |
|---|---|---|
| $X_c$ | Tuple of Critical features | Observable |
| $X_g$ | Tuple of Non-critical or general features | Observable |
| $X$ | Tuple of all input features (critical and general) | Observable |
| $Z$ | Protected attribute (s) | Observable |
| $U_X$ (Note that, $Z \perp\!\!\!\perp U_X$) | Tuple of latent social factors that do not cause $Z$ | Not observable in general |
| $\hat{Y} = r(X) = h(Z, U_X)$ | Model output | Observable |

(a) Canonical Example 1: Hiring with Biased Critical Feature (Desirable: $M_{NE} = 0$)



(b) Canonical Example 2: Hiring with Biased General Feature (Desirable: $M_{NE} > 0$)



(c) Canonical Example 3: Counterfactually Fair Hiring (Desirable: $M_{NE} = 0$)



(d) Canonical Example 4: Non-Exempt Masked Bias in Hiring Ads I (Desirable: $M_{NE} > 0$)



(e) Canonical Example 5: Non-Exempt Masked Bias in Hiring Ads II (Desirable: $M_{NE} > 0$)



(f) Canonical Example 6: Bias Amplification by Unmasking (Desirable: $M_{NE} > 0$)
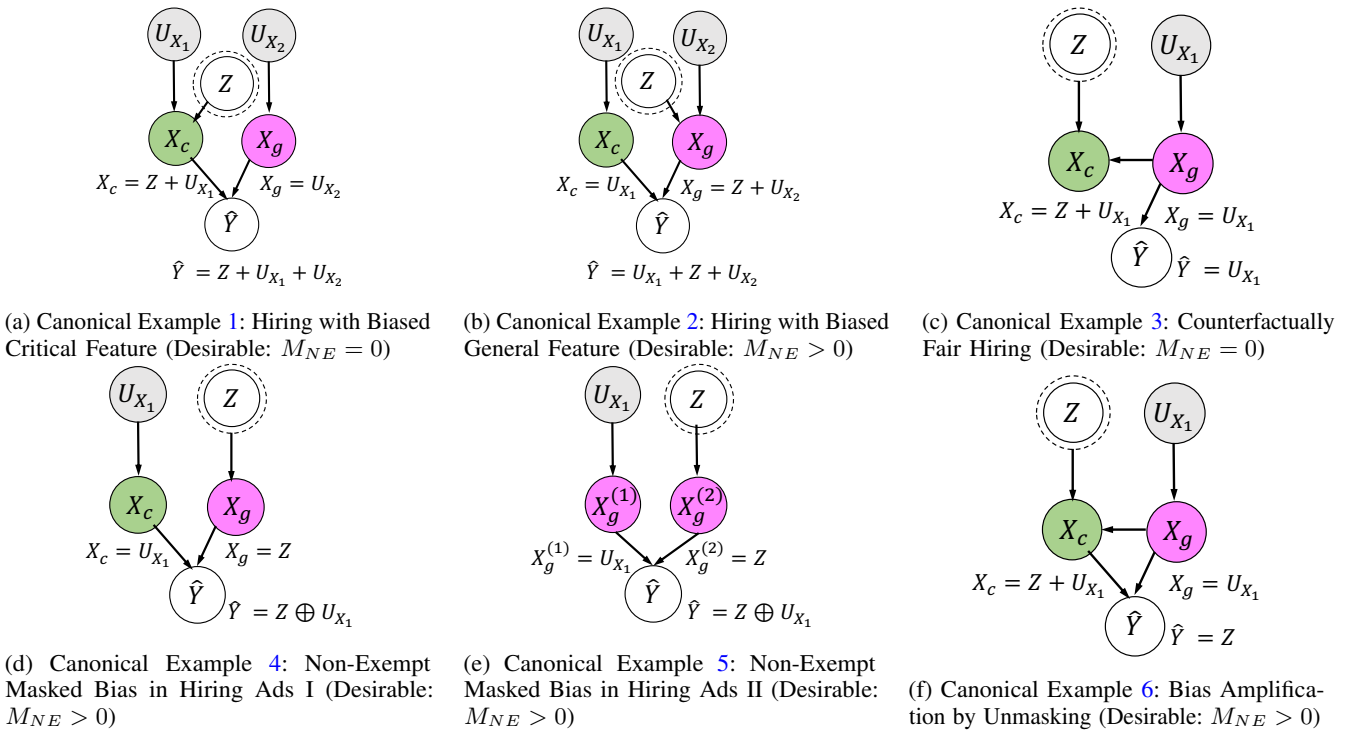
Fig. 5: Thought experiments to motivate desirable properties of non-exempt bias: In all the figures, $Z$ denotes the protected attribute, e.g., gender, race, etc., and $U_{X_1}, U_{X_2}$ denotes other latent social factors independent of $Z$. The critical feature is denoted by $X_c$, the non-critical/general feature is $X_g$, and the model output (hiring decision) is $\hat{Y}$.

## III. MAIN RESULTS

In Section III-A, we first formally state the desirable properties that a measure of non-exempt bias $(M_{NE})$ should satisfy. These properties were only intuitively stated in Section I. Next, we introduce our proposed measure that satisfies all these properties (Theorem 1 in Section III-A). In Section III-B, we discuss in detail on how we arrive at these desirable properties through several canonical examples (summarized in Table III and Fig. 5), that helps us quantify our intuitive notion of non-exempt bias. In Section III-C, we examine measures in existing literature that have some provision for exemptions, namely, path-specific counterfactual fairness [19], conditional statistical parity [40], and justifiable fairness [41], and understand their limitations.

### A. Desirable Properties Leading to Our Proposed Measure of Non-Exempt Bias

It is desirable that our measure of *non-exempt* bias $(M_{NE})$ is able to capture the intuition of a virtual constituent or proxy of $Z$ being formed inside a given black-box model that a) causally influences the output $\hat{Y}$; and b) cannot be attributed to the critical features $X_c$ alone. To arrive at a set of desirable properties for a measure of *non-exempt bias* $(M_{NE})$, we examine candidate measures and examine their utility and limitations through canonical examples (see Fig. 5). While we discuss the rationale for each of these properties in more detail in Section III-B, here we state the properties and provide a brief intuition for each of them. For simplicity, assume that the protected attribute $Z$ as well as all the other independent latent variables $U_{X_1}, U_{X_2}, \ldots$ are i.i.d. Bern(½) in our canonical examples.

The first candidate measure of non-exempt bias we consider is based on conditional mutual information (CMI), and is: $M_{NE} = \mathrm{I}(Z; \hat{Y} \mid X_c)$ (Candidate Measure 1 in Section III-B). Inspired from the concept of conditional statistical parity [40], this measure assumes that there is no non-exempt bias if and only if the hiring decision $\hat{Y}$ and the protected attribute $Z$ (e.g., gender) are independent, conditioned on the critical feature $X_c$ (e.g., coding-test score for a software engineering job). While this measure is intuitively appealing (further demonstrated through Canonical Examples 1 and 2 in Section III-B; see Fig. 5a and Fig. 5b), it has a limitation: it can sometimes *falsely detect non-exempt bias when there is none. E.g.*, consider a scenario where the model is counterfactually fair, and hence there is no bias (exempt or non-exempt; Canonical Example 3 in Section III-B, Fig. 5c): the critical feature, namely, the coding-test score for a software engineering job is biased, i.e., $X_c = Z + U_{X_1}$ with $U_{X_1}$ being the latent inner ability of a candidate. However, the model is able to distill out the latent inner ability $U_{X_1}$ using all the features and take hiring decisions entirely based on them, i.e., $\hat{Y} = U_{X_1}$. Here, $M_{NE} = \mathrm{I}(Z; \hat{Y} \mid X_c) > 0$ when it is desirable that $M_{NE}$ be 0. This canonical example motivates the following property:

**Property 1** (Zero Influence). *$M_{NE}$ should be 0 if $\mathrm{CCI}(Z \to \hat{Y}) = 0$ (or equivalently, $\mathrm{I}(Z; \hat{Y}, U_X) = 0$).*

This limitation of $\mathrm{I}(Z; \hat{Y} \mid X_c)$ leads us to examine PID, decomposing $\mathrm{I}(Z; \hat{Y} \mid X_c)$ into two components: unique information $\mathrm{Uni}(Z : \hat{Y} \mid X_c)$ and synergistic information $\mathrm{Syn}(Z : (\hat{Y}, X_c))$. The sub-component $\mathrm{Uni}(Z : \hat{Y} \mid X_c)$ always satisfies Property 1 (proof in Lemma 13 in Appendix B), even though $\mathrm{I}(Z; \hat{Y} \mid X_c)$ sometimes may not do so because of the synergistic component (which caused false detection of non-exempt bias in the previous scenario). This leads us to examine another candidate measure of non-exempt bias, namely, $M_{NE} = \mathrm{Uni}(Z : \hat{Y} \mid X_c)$ (Candidate Measure 2 in Section III-B). For example, consider hiring for a software-engineering job using coding-test score (critical feature) and aptitude-test score (non-critical/general feature). It is desirable that $M_{NE}$ be non-zero if $\hat{Y}$ has any unique information about $Z$ that is not present in $X_c$ (coding test) because then that information content is also attributed to $X_g$ (also see Section III-B4 to further motivate this property).

**Property 2** (Non-Exempt Statistically Visible Bias). *$M_{NE}$ should be strictly greater than $0$ if $\hat{Y}$ has any unique information about $Z$ not present in $X_c$. Thus,* $\mathrm{Uni}(Z : \hat{Y}|X_c) > 0$ *should imply that $M_{NE} > 0$.*

However, this property alone does not capture all scenarios where $M_{NE}$ is desired to be non-zero. Statistical masking can sometimes prevent the entire non-exempt bias from exhibiting itself in $\mathrm{Uni}(Z : \hat{Y}|X_c)$ as demonstrated in the following scenario. Suppose an ad for a job is shown selectively to: a) men with high coding-test scores and b) women with low coding-test scores (Canonical Examples 4 and 5 in Section III-B; see Fig. 5d and 5e). Such a model might seem "facially neutral", i.e., with no statistically visible dependence between $Z$ and $\hat{Y}$ ($\mathrm{I}(Z; \hat{Y}) = 0$), but is clearly unfair to high-scoring women candidates. Since $\mathrm{Uni}(Z : \hat{Y}|X_c) \leq \mathrm{I}(Z; \hat{Y})$ (recall (2) in Section II-A and non-negativity of all PID terms), we have $\mathrm{Uni}(Z : \hat{Y}|X_c) = 0$ for this canonical example, showing that it fails to capture such "non-exempt masked bias." In essence, $\mathrm{Uni}(Z : \hat{Y}|X_c)$ is therefore a lower bound for non-exempt bias $M_{NE}$, i.e., $\mathrm{Uni}(Z : \hat{Y}|X_c) > 0 \implies M_{NE} > 0$ but not necessarily the other way round (making this candidate measure a "lower bound" for $M_{NE}$). The next property attempts to find an upper bound for $M_{NE}$.

Notice that, in the previous Canonical Examples 4 and 5, $\hat{Y}$ has a virtual constituent $Z$ influencing it, that is not due to the critical features $X_c$. However, the influence of $Z$ does not exhibit itself in the statistically visible bias $\mathrm{I}(Z; \hat{Y})$. To resolve this issue, we now consider a non-observational, causal candidate measure inspired from path-specific counterfactual fairness [19] that specifically examines causal paths from $Z$ to $\hat{Y}$ in the SCM (Candidate Measure 3 in Section III-B). This measure implies there is no non-exempt bias if all paths from $Z$ to $\hat{Y}$ in the SCM pass through $X_c$. However, we identify scenarios where this approach can also fail to quantify non-exempt bias, e.g., in Canonical Example 6 in Section III-B (Fig. 5f). Here the critical feature, coding-test score is $X_c = Z + U_{X_1}$, and the non-critical feature, aptitude-test score is $X_g = U_{X_1}$. The model amplifies the bias in the hiring decision by cancelling $U_{X_1}$, i.e., $\hat{Y} = Z$. For this example, even though we have the causal path from $Z$ to $\hat{Y}$ passing through $X_c$, we contend that here both $X_c$ and $X_g$ jointly have information about $Z$ that cannot be attributed to $X_c$ alone. Therefore, it is desirable that we have a measure of non-exempt bias $M_{NE}$ which is non-zero for this example ($\mathrm{Uni}(Z : \hat{Y}|X_c)$ and $\mathrm{I}(Z; \hat{Y}|X_c)$ are also non-zero for this example).

From a causal point of view, here $U_{X_1}$ is a "confounder" for both $X_c$ and $\hat{Y}$ (separately influences both $X_c$ and $\hat{Y}$ along different paths). Intuitively, a scenario when there is no non-exempt bias would be: (i) All causal paths from $Z$ to $\hat{Y}$ in the SCM pass through $X_c$; and also (ii) No $U_{X_i}$ acts as a confounder for both $X_c$ and $\hat{Y}$. This leads to the intuition that to be able to say there is no non-exempt bias, one might be able to split $U_X$ into two subsets $U_a$ and $U_b$ (further functional generalizations discussed in Section VIII), such that: (i) $U_a$ consists of the latent factors that do not influence $\hat{Y}$ at all, or influence it only through $X_c$ without acting as confounder; (ii) $U_b$ consists of the remaining latent factors, that only influence $\hat{Y}$ and not $X_c$; and (iii) The Markov chain $(Z, U_a) - X_c - (\hat{Y}, U_b)$ holds[7]. To understand this, one may also refer to Canonical Example 1 in Section III-B (Fig. 5a). This leads to the following property (see Section III-B5 to further motivate this property).

**Property 3** (Non-Exempt Masked Bias). *$M_{NE}$ should be non-zero in the canonical example of non-exempt masked bias: $X_1 = Z$, $X_2 = U_X$, and $\hat{Y} = Z \oplus U_X$ with $Z, U_X \sim$ i.i.d. Bern(½) and $X_1 \in X_g$. However, $M_{NE}$ should be $0$ if $(Z, U_a) - X_c - (\hat{Y}, U_b)$ form a Markov chain for some subsets $U_a, U_b \subseteq U_X$ such that $U_a = U_X \backslash U_b$.*

Properties 2 and 3 provide lower and upper bounds (respectively) on our measure of non-exempt bias, i.e., it is desirable that:

$$\mathrm{Uni}(Z : \hat{Y}|X_c) \leq M_{NE} \leq \min_{U_a, U_b \text{ s.t. } U_a = U_X \backslash U_b} \mathrm{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c).$$

This observation is important in itself: the unique information measure, being a lower bound, never falsely detects non-exempt bias when there is none, and thus can serve as a conservative estimate of non-exempt bias.

The next three properties are more intuitive. Consider the scenario where no feature is deemed critical (i.e., $X_c = \phi$) and all features are non-critical, e.g., hiring for a manager's role using aptitude-test and coding-test scores. Here, one would like $M_{NE}$ to be equal to the total bias $\mathrm{I}(Z; (\hat{Y}, U_X))$, i.e., no bias is exempt because no feature is deemed critical.

**Property 4** (Absence of Exemptions). *If no feature is deemed critical ($X_c = \phi$), then a measure $M_{NE}$ should be equal to the total bias, i.e., $\mathrm{I}(Z; (\hat{Y}, U_X))$.*

---

[7]Notice that, this condition implies $Z - X_c - \hat{Y}$ but not the other way round.

Next, suppose that the same model is being used for a software-engineering role where coding-test score is deemed as a critical feature but aptitude-test score is not. For a fixed set of features and a fixed model $\hat{Y} = h(Z, U_X)$, it is desirable that $M_{NE}$ either decreases or stays the same as more features are removed from the set $X_g$ and added to $X_c$.

**Property 5** (Non-Increasing with More Exemptions). *For a fixed set of features $X$ and a fixed model $\hat{Y} = h(Z, U_X)$, a measure $M_{NE}$ should be non-increasing if a feature is removed from $X_g$ and added to $X_c$.*

Lastly, suppose that the model is used for an even more specific role where both coding test and aptitude test are deemed as critical features. If all the features are in the exempt set $X_c$, we require the measure $M_{NE}$ to be 0.

**Property 6** (Complete Exemption). *$M_{NE}$ should be 0 if all features are exempt, i.e., $X_c = X$ and $X_g = \phi$.*

These six properties lead to a novel measure of non-exempt bias that satisfies all of them (proved in Theorem 1).

**Definition 7** (Non-Exempt Bias). *Our proposed measure of non-exempt bias is given by:*

$$M_{NE}^* = \min_{U_a, U_b} \mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \text{ such that } U_a = U_X \backslash U_b.$$

Note that, for the rest of the paper, we use the notation $M_{NE}$ to denote any candidate measure of non-exempt bias, and $M_{NE}^*$ to specifically denote our proposed measure in Definition 7.

**Theorem 1** (Properties). *Properties 1-6 are satisfied by our proposed measure*

$$M_{NE}^* = \min_{U_a, U_b} \mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \text{ such that } U_a = U_X \backslash U_b.$$

*Proof Sketch:* A detailed proof is provided in Appendix C-A. Here, we provide a brief proof sketch. For Property 1, note that

$$M_{NE}^* \leq \mathrm{Uni}(Z : \hat{Y}, U_X|X_c) \leq \mathrm{I}(Z; (\hat{Y}, U_X)), \tag{4}$$

where the last step holds as unique information is also a component of mutual information (see (2) in Section II-A). For Property 2, we show that $M_{NE}^* \geq \mathrm{Uni}(Z : \hat{Y}|X_c)$ using a monotonicity property of unique information [57, Lemma 31]. Lastly, for Property 3, we have $\mathrm{I}(Z, U_a; \hat{Y}, U_b|X_c) = 0$ for some $U_a, U_b$, implying that $\mathrm{Uni}(Z, U_a : \hat{Y}, U_b|X_c)$ is also 0 for those $U_a, U_b$ because unique information is a component of conditional mutual information (see (3) in Section II-A). For Property 4, we show that when $X_c = \phi$, we have $M_{NE}^* = \min_{U_a, U_b \text{ s.t. } U_a = U_X \backslash U_b} \mathrm{I}(Z, U_a; \hat{Y}, U_b) = \mathrm{I}(Z; (\hat{Y}, U_X))$. Property 5 is derived using another monotonicity property of unique information [57, Lemma 32]. For Property 6,

$$M_{NE}^* \leq \mathrm{Uni}(Z, U_X : \hat{Y}|X) \overset{(a)}{\leq} \mathrm{I}(Z, U_X; \hat{Y}|X) \overset{(b)}{=} 0, \tag{5}$$

where (a) holds because unique information is a component of conditional mutual information (see (3) in Section II-A) and (b) holds as $\hat{Y}$ is a deterministic function of $X$.

**Remark 5** (On Exhaustive Set of Properties leading to a Unique Measure). *We note that our properties do not quantify how exactly the non-exempt bias should "scale" when the measure is nonzero since they are only conditions on when this bias is nonzero, or on the monotonicity of this bias. Hence, these properties do not lead to a unique measure. Also, note that this is an issue with all measures of fairness in that they go to zero based on an intuitive notion of fairness but their exact scaling when they are non-zero is not unique. Neither do we claim that the proposed list of desirable properties (axioms) are exhaustive. In general, it is difficult to prove that a proposed set of properties (or, axioms) is exhaustive for a problem. E.g., Shannon established uniqueness on entropy with respect to **some** properties in [58] but the needs of the application can still drive the use of alternate measures. E.g. Renyi measures [21], [55], [59]–[61] have been found to be useful in security and privacy applications because they weigh outliers differently. Therefore, we believe, that there may be value in the measure not being unique so that it can be tuned to the needs of the application, as well as, motivate future work in this direction. Nonetheless, our properties do capture important aspects of the problem, e.g., non-exempt masked and non-exempt statistically visible biases, as discussed in Section IV.*

**Remark 6.** *We note that the proposed measure is counterfactual (non-observational) in nature, i.e., it requires knowledge of the true SCM. While we are able to compute the measure in our case study on artificial datasets (known SCM) in Section VII, we acknowledge that even after knowledge of the true SCM, there may be computational challenges if the number of latent variables is large. However, one must note that it is important to arrive at measures that satisfy all desirable properties, however hard they might be to compute: (i) It makes the shortcomings of other measures more explicit, informing which computable/estimable definition to choose in a given situation; (ii) It opens the avenue of obtaining relaxations that may be easier to estimate; (iii) One can begin exploring research directions to reduce the difficulty/complexity (statistical and/or computational) of estimating these measures.*

**Remark 7** (On Simplicity of Examples). *We note that, at a first glance, our examples might seem simple, and real world models will only be more complex due to a mix of causal and statistical relationships. These simple examples help us isolate many of these individual causal and statistical relationships, and examine them carefully. E.g., scenarios where only one of non-exempt masked, non-exempt visible, exempt masked or exempt visible bias is present or none of them is present (see Fig. 6). When both non-exempt masked and non-exempt statistically visible biases are present together, we are able to quantify both of them appropriately (discussed further in Section IV). Thus, developing an axiomatic understanding of such simple examples is an essential first step in understanding the complex interplay of various relationships in a real dataset. Indeed, examining toy examples (thought experiments) is a common practice in several works in existing fairness literature [16], [17], [26], [34], [41], some of which have also inspired our examples in this work. Furthermore, our quantification of non-exempt bias is not limited to black-box models alone, but also applies to "white-box" models [20], e.g., decision trees, linear classifiers, etc., and also to non-AI-based decisions as long as the decision is as a deterministic function of the input features, i.e., $\hat{Y} = h(X)$.*
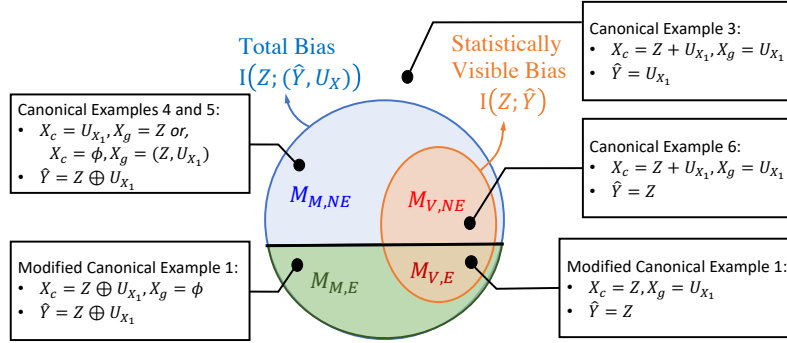


Fig. 6: Our examples isolate different kinds of scenarios, namely, masked non-exempt ($M_{M,NE}$), masked exempt ($M_{M,E}$), visible non-exempt ($M_{V,NE}$), and visible exempt ($M_{V,E}$), as well as scenarios where there is no total bias (more in Section IV).

### B. Detailed Rationale Behind the Desirable Properties Leading to A Measure of Non-Exempt Bias

Here we provide detailed rationale behind all our desirable properties using canonical examples (summarized in Table III). We start by examining two canonical examples that help us motivate the basic intuition behind *non-exempt bias*. These examples also help us understand the limitations of *statistical parity* [3], [6] and *equalized odds* [7] which are two popular measures of fairness that do not have provision for critical feature exemptions.

*1) Limitations of Statistical Parity:* As discussed in Section II, a model is deemed fair by statistical parity if $Z \perp\!\!\!\perp \hat{Y}$, i.e., $I(Z; \hat{Y}) = 0$. However, the following example exposes some of its limitations.

**Canonical Example 1** (Hiring with Biased Critical Feature). *Let $X_c = Z + U_{X_1}$ be a coding-test score[8] and $X_g = U_{X_2}$ be an aptitude-test score. Here the protected attribute $Z \sim \text{Bern}(\frac{1}{2})$ denotes gender, $U_{X_1} \sim \text{Bern}(\frac{1}{2})$ denotes inner ability to code and $U_{X_2} \sim \text{Bern}(\frac{1}{2})$ denotes knowledge. An algorithm is deciding whether to hire software engineers based on a score $\hat{Y} = Z + U_{X_1} + U_{X_2}$. This is shown in Fig. 5a. Here + denotes addition (not to be confused with the binary OR).*

First notice that this model will be deemed *unfair* by both statistical parity and counterfactual fairness. Statistical parity is violated because $Z$ and $\hat{Y}$ are not independent, i.e., the statistically visible bias $I(Z; \hat{Y}) > 0$. Consequently, the total bias $I(Z; (\hat{Y}, U_X))$ is also non-zero since $I(Z; (\hat{Y}, U_X)) \geq I(Z; \hat{Y}) > 0$, violating counterfactual fairness. However, for this example, the coding-test score is a critical feature (bonafide requirement) for the job. Therefore, one may feel that any bias in $\hat{Y}$ that is explainable by the coding-test score may be exempted. An attempt to ensure statistical parity for such an example, e.g., by reducing the importance (weight) of the critical feature in the decision making, violates the bonafide requirement of the job. Intuitively, even though the virtual constituent or proxy of $Z$, namely, $Z + U_{X_1}$, influences the output $\hat{Y}$, it is entirely explainable by $X_c$. Thus, for such an example, it is desirable that a measure of discrimination (non-exempt bias $M_{NE}$) be 0.

*2) Limitations of Equalized Odds:* Equalized odds [7], [12] is another popular measure of fairness that attempts to address this limitation of statistical parity by using the true labels (or true final-decision scores) to represent the job requirements. Equalized odds states that a model is fair if

$$\Pr(\hat{Y} = y | Z = z, Y = \tilde{y}) = \Pr(\hat{Y} = y | Z = z', Y = \tilde{y}) \ \forall z, z', y, \tilde{y}.$$

This criterion is also equivalent to $\hat{Y} \perp\!\!\!\perp Z | Y$, or, $I(Z; \hat{Y} | Y) = 0$. Indeed, in the previous example (Canonical Example 1), if the true final-decision scores already incorporate this critical requirement in them, e.g., $Y = Z + U_{X_1} + U_{X_2}$, then $I(Z; \hat{Y} | Y) = 0$,

---

[8]The influence of $Z$ on score in the SCM can arise due to various factors, e.g., historical lack of opportunities or sampling bias due to candidates of one protected group not applying enough etc. For instance, there may be a hidden node representing opportunity such that $Z$ influences the score only though that hidden node, and the score becomes independent of $Z$ given opportunity. We adopt a simplistic representation here for ease of understanding (also see [62]).

TABLE III: Summary of Canonical Examples and Candidate Measures of Non-Exempt Bias

| Canonical Examples | Candidate Measure 1: $I(Z;\hat{Y}\|X_c)$ | Candidate Measure 2: $Uni(Z:\hat{Y}\|X_c)$ | Candidate Measure 3: Path-Specific Causality | Proposed Measure: $\min_{U_a,U_b} Uni((Z,U_a):(\hat{Y},U_b)\|X_c)$ such that $U_a = U_X \backslash U_b$. |
|---|---|---|---|---|
| 1. Hiring with Biased Critical Feature<br>• $X_c = Z + U_{X_1}$ and $X_g = U_{X_2}$.<br>• $\hat{Y} = Z + U_{X_1} + U_{X_2}$.<br>**Desirable:** $M_{NE} = 0$ | ✓ | ✓ | ✓ | ✓ |
| 2. Hiring with Biased General Feature<br>• $X_c = U_{X_1}$ and $X_g = Z + U_{X_2}$.<br>• $\hat{Y} = Z + U_{X_1} + U_{X_2}$.<br>**Desirable:** $M_{NE} > 0$ | ✓ | ✓ | ✓ | ✓ |
| 3. Counterfactually Fair Hiring<br>• $X_c = Z + U_{X_1}$ and $X_g = U_{X_1}$.<br>• $\hat{Y} = U_{X_1}$.<br>**Desirable:** $M_{NE} = 0$ | ✗ | ✓ | ✓ | ✓ |
| 4. Non-Exempt Masked Bias in Hiring Ads I<br>• $X_c = U_{X_1}$ and $X_g = Z$.<br>• $\hat{Y} = Z \oplus U_{X_1}$.<br>**Desirable:** $M_{NE} > 0$ | ✓ | ✗ | ✓ | ✓ |
| 5. Non-Exempt Masked Bias in Hiring Ads II<br>• $X_c = \phi$ and $X_g = (Z, U_{X_1})$.<br>• $\hat{Y} = Z \oplus U_{X_1}$.<br>**Desirable:** $M_{NE} > 0$ | ✗ | ✗ | ✓ | ✓ |
| 6. Bias Amplification by Unmasking<br>• $X_c = Z + U_{X_1}$ and $X_g = U_{X_1}$.<br>• $\hat{Y} = Z$.<br>**Desirable:** $M_{NE} > 0$ | ✓ | ✓ | ✗ | ✓ |

and the model is deemed *fair* by equalized odds. While equalized odds is a reasonable quantification in scenarios where the true label (or true final-decision score) is indeed a justified representation of the job requirements, the measure $I(Z;\hat{Y} \mid Y)$ has often been criticized to be affected by label bias, as we demonstrate through this example.

**Canonical Example 2** (Hiring with Biased General Feature). *Let $X_c = U_{X_1}$ denote the coding-test score and $X_g = \begin{cases} U_{X_2} + 1, & Z = 0 \\ U_{X_2}, & Z = 1 \end{cases}$ denote the aptitude-test score (biased). This can be rewritten as $X_g = Z(U_{X_2}+1)+(1-Z)U_{X_2} = Z + U_{X_2}$, where $Z \sim Bern(\frac{1}{2})$ denotes gender, $U_{X_1} \sim Bern(\frac{1}{2})$ denotes the inner ability to code and $U_{X_2} \sim Bern(\frac{1}{2})$ denotes knowledge. Now suppose, the historic dataset has true decision scores given by $Y = U_{X_1} + Z + U_{X_2}$. This is shown in Fig. 5b.*

In this scenario, suppose we choose a perfect predictor, i.e., $\hat{Y} = Y = U_{X_1} + Z + U_{X_2}$. The perfect predictor always satisfies equalized odds because $I(Z;\hat{Y} \mid Y) = 0$ if $\hat{Y} = Y$. However, if examined deeply, this model is propagating bias from aptitude-test score, a non-critical/general feature, which is discriminatory and non-exempt. Intuitively, a virtual constituent or proxy of $Z$, i.e., $Z + U_{X_2}$, is being formed from $X_g$ that is influencing the output $\hat{Y}$. For such an example[9], it is desirable that a measure of discrimination (non-exempt bias $M_{NE}$) is not zero.

*3) Motivation for Conditional Mutual Information and its Limitations:* Next, we start out with the aim of finding a suitable measure of non-exempt bias ($M_{NE}$) that resolves both these canonical examples. Notice that, both these examples can be resolved by a notion of *conditional statistical parity* [40], which deems a model as fair if and only if $Z \perp\!\!\!\perp \hat{Y}|X_c$, i.e.,

$$\Pr(\hat{Y} = y|X_c = x_c, Z = z) = \Pr(\hat{Y} = y|X_c = x_c, Z = z') \quad \forall y, x_c, z, z'.$$

This idea also connects with Simpson's paradox [35] which refers to a statistical trend that appears in several different groups of data but disappears or reverses when these groups are combined. In Canonical Example 1, $Z$ and $\hat{Y}$ are not independent but they become so when conditioned on $X_c$, i.e., $I(Z;\hat{Y}) > I(Z;\hat{Y} \mid X_c)$. In Canonical Example 2, $I(Z;\hat{Y}) < I(Z;\hat{Y} \mid X_c)$. This notion of *conditional statistical parity* leads us to propose the following quantification of non-exempt bias ($M_{NE}$).

**Candidate Measure of Non-Exempt Bias 1.** $M_{NE} = I(Z;\hat{Y} \mid X_c)$.

---

[9]The example can be made more realistic if $U_{X_1}, U_{X_2}$ are i.i.d. $\mathcal{N}(0,1)$. Now suppose, the historic dataset has true labels given by $Y = \text{sgn}\left(Z + U_{X_1} + U_{X_2} - 0.5\right)$ which is binary. A perfect classifier $\hat{Y} = Y$, that satisfies equalized odds, is still discriminatory because it is influenced by $Z$ in its decision, that is arising from a non-critical feature.

This measure resolves both Canonical Examples 1 and 2. However, the following example exposes some of its limitations.

**Canonical Example 3** (Counterfactually Fair Hiring). *Let $Z \sim Bern(\frac{1}{2})$ be gender, $U_{X_1} \sim Bern(\frac{1}{2})$ be the inner ability of a candidate, and $X_c = \begin{cases} U_{X_1}, & Z = 0 \\ U_{X_1} + 1, & Z = 1 \end{cases}$ be the coding-test score (critical feature). This can be rewritten as $X_c = Z(U_{X_1} + 1) + (1 - Z)U_{X_1} = Z + U_{X_1}$. However, instead of only using the biased test score, suppose the company chooses to conduct thorough evaluation of their online code samples, leading to another score that distills out their inner ability, i.e., $X_g = U_{X_1}$. Suppose the model for hiring that maximizes accuracy turns out to be $\hat{Y} = X_g = U_{X_1}$. This is shown in Fig. 5c.*

Notice that, this model is deemed *fair* by counterfactual fairness because the total bias $I(Z; (\hat{Y}, U_X)) = 0$. This means that the output $\hat{Y}$ has no counterfactual causal influence of $Z$. Even though the bias from $X_c$ is legally exempt, the trained black-box model happens to base its decisions on another available non-critical/general feature that has no counterfactual causal influence of $Z$. Thus, there is no bias in the outcome $\hat{Y}$ (this is true even if the features in $X_c$ were not exempt). Therefore, it is desirable that the non-exempt bias $M_{NE}$ is also 0. This is also consistent with the intuition that here no virtual constituent or proxy of $Z$ influences the output. However, the candidate measure $I(Z; \hat{Y} \mid X_c) = I(Z; U_{X_1} \mid Z + U_{X_1})$ is non-zero here, leading to a false positive conclusion in detecting non-exempt bias.

**Remark 8** (Cancellation of Paths). *A similar situation arises if $X_c = Z + U_{X_1}$, $X_g = Z$ and $\hat{Y} = X_c - X_g = U_{X_1}$. Even though the bias from $X_c$ may be exempt, the trained model ends up removing the counterfactual causal influence of $Z$ from the decisions to make them counterfactually fair in a manner similar to the example of interviews (recall Scenario 2 in Section II; also shown in Fig. 3b). The influences of $Z$ along two different causal paths cancel each other in the final output, so that $CCI(Z \rightarrow \hat{Y}) = 0$ (and, $I(Z; (\hat{Y}, U_X)) = 0$). Since the total bias $I(Z; (\hat{Y}, U_X)) = 0$, the question of non-exempt or exempt bias does not arise. However, the candidate measure $I(Z; \hat{Y} \mid X_c)$ is non-zero here.*

This example also serves as a rationale for the property of zero influence, i.e., Property 1 which states that $M_{NE}$ should be 0 if the total bias is 0. We aim to find a measure that resolves all of these examples (summarized in Fig. 5).

*4) Motivation for Unique Information and its Limitations:* We notice that conditioning on the critical feature $X_c$ can increase or decrease mutual information. For instance, in Canonical Example 1, we have $I(Z; \hat{Y}) > 0$ but $I(Z; \hat{Y} \mid X_c) = 0$. In Canonical Example 3, $I(Z; \hat{Y} \mid X_c) > 0$ but $I(Z; \hat{Y}) = 0$. For both these examples, it is desirable that $M_{NE} = 0$. This motivates us to consider another candidate measure of non-exempt bias that is equal to the information-theoretic sub-volume of intersection between $I(Z; \hat{Y})$ and $I(Z; \hat{Y} \mid X_c)$ (recall Fig. 1b), that goes to 0 when any one of them is 0. This is a quantity that is derived from the PID literature, and is called the *unique information* of $Z$ in $\hat{Y}$ that is not present in $X_c$.

**Candidate Measure of Non-Exempt Bias 2.** $M_{NE} = \text{Uni}(Z : \hat{Y} | X_c)$.

We now show that this measure resolves the examples discussed so far, namely, Canonical Example 1 (Fig. 5a), Canonical Example 2 (Fig. 5b), Canonical Example 3 (Fig. 5c) and a (similar) example in Remark 8. We start with Canonical Example 1 (hiring with biased critical feature), where $\hat{Y} = Z + U_{X_1} + U_{X_2}$ and $X_c = Z + U_{X_1}$. Recall that the total mutual information (statistically visible bias) can be decomposed as follows:

$$I(Z; \hat{Y}) = \text{Uni}(Z : \hat{Y} | X_c) + \text{Red}(Z : (\hat{Y}, X_c)) \text{ (from (2) in Section II-A).}$$

For this example, we notice that even though $I(Z; \hat{Y}) > 0$, we have $\text{Uni}(Z : \hat{Y} | X_c) = 0$. This is because,

$$I(Z; \hat{Y} \mid X_c) = \text{Uni}(Z : \hat{Y} | X_c) + \text{Syn}(Z : (\hat{Y}, X_c)) \text{ (from (3) in Section II-A),}$$

and $I(Z; \hat{Y} \mid X_c) = 0$ for Canonical Example 1. In Canonical Example 1, the entire statistically visible bias $I(Z; \hat{Y})$ is essentially redundant information between $\hat{Y}$ and $X_c$ which is exempted.

Next, we revisit Canonical Example 2 ($\hat{Y} = U_{X_1} + Z + U_{X_2}$ and $X_c = U_{X_1}$) where it is intuitive that the measure of non-exempt bias should be non-zero. $\text{Uni}(Z : \hat{Y} | X_c)$ is non-zero here (see Supporting Derivation 1 in Appendix C-B), consistent with our intuition. As a proof sketch, recall the tabular representation in Fig. 1b. $\text{Red}(Z : (\hat{Y}, X_c))$ is the sub-volume of intersection between $I(Z; X_c)$ and $I(Z; \hat{Y})$, and hence goes to zero because $I(Z; X_c) = 0$. This leads to $\text{Uni}(Z : \hat{Y} | X_c) = I(Z; \hat{Y})$ which is non-zero here.

Lastly, $\text{Uni}(Z : \hat{Y} | X_c)$ is also 0 in Canonical Example 3 (counterfactually fair hiring) and the (similar) example of cancellation of paths in Remark 8. More importantly, we note that, while conditional mutual information $I(Z; \hat{Y} \mid X_c)$ may be non-zero even if the the total bias or counterfactual causal influence is 0 (as in Canonical Example 3), unique information is not. *In Lemma 13 in Appendix B, we show that $\text{Uni}(Z : \hat{Y} | X_c)$ is always zero if the total bias or counterfactual causal influence is 0, i.e., $I(Z; (\hat{Y}, U_X)) = 0$.* In fact, $\text{Uni}(Z : \hat{Y} | X_c)$ is a sub-volume or component of the previous candidate measure $I(Z; \hat{Y} \mid X_c)$, that is guaranteed to be 0 if the total bias is zero.

These examples serve as our rationale for the property of non-exempt statistically visible bias, i.e., Property 2 which states that $M_{NE}$ should be 0 if $\text{Uni}(Z : \hat{Y} | X_c) > 0$. $\text{Uni}(Z : \hat{Y} | X_c)$, however, is not sufficient as a candidate measure as it fails to capture *non-exempt masked bias*, as we will demonstrate in Canonical Example 4. Thus, Property 2 is only a lower bound,

i.e., sometimes $M_{NE}$ may still need to be non-zero even when $\text{Uni}(Z : \hat{Y}|X_c) = 0$. Property 2 only captures the non-exempt *statistically visible bias* that cannot be accounted for by $X_c$ alone.

**Canonical Example 4** (Non-Exempt Masked Bias in Hiring Ads I). *An ad for a software-engineering job is only presented to men $(Z = 1)$ with a coding-test score above a threshold $(U_{X_1} = 1)$, and to women $(Z = 0)$ with a coding-test score below a threshold $(U_{X_1} = 0)$ with $Z$ and $U_{X_1}$ being i.i.d. Bern(½). Here, $X_c = U_{X_1}$ and $X_g = Z$. The model output is given by $\hat{Y} = Z \oplus U_{X_1}$. This example is shown in Fig. 5d.*

This model discriminates against half of the population (high-scoring women) for whom the ad may be relevant. This is also supported by the fact that that the total bias $I(Z; (\hat{Y}, U_X)) > 0$. Intuitively, here a virtual constituent or proxy $(Z)$ is formed inside the black-box model that influences the output and that is derived entirely from $X_g$. For such an example, it is desirable that the non-exempt bias $M_{NE}$ should not be 0. In fact, this example demonstrates that there may be non-exempt bias even when the statistically visible bias $I(Z; \hat{Y}) = 0$. Here, $\text{Uni}(Z : \hat{Y}|X_c)$ fails to capture the masked bias because it has to be zero whenever $I(Z; \hat{Y}) = 0$ (using (2) in Section II-A).

Let us revisit the candidate measure $I(Z; \hat{Y} \mid X_c)$. This measure resolves all the examples discussed so far (1-4) except giving a false positive conclusion in Canonical Example 3. Notice that, $I(Z; \hat{Y} \mid X_c)$ is zero if and only if $Z - X_c - \hat{Y}$ form a Markov chain. While the Markov chain $Z - X_c - \hat{Y}$ may not always hold even when it is desirable for $M_{NE}$ to be zero as in Canonical Example 3, we have seen that in all the examples so far (1-4) where the Markov chain $Z - X_c - \hat{Y}$ holds, it has been desirable that $M_{NE}$ be zero (possible one-way implication). Assuming that the Markov chain $Z - X_c - \hat{Y}$ is a sufficient condition for $M_{NE}$ to be zero, we proposed the following property of non-exempt masked bias in our prior work [1].
$M_{NE}$ *should be non-zero in the example of non-exempt masked bias, i.e., Canonical Example 4 even if $I(Z; \hat{Y}) = 0$. But, $M_{NE}$ should be 0 if the Markov chain $Z - X_c - \hat{Y}$ holds.*

**Remark 9** (Relation to our prior work [1]). *In our prior work [1], this property, in conjunction with Properties 1, 2 and 6, leads to a measure that quantifies only a sub-volume of $I(Z; \hat{Y} \mid X_c)$ that no longer gives false positive conclusion in Canonical Example 3 while still resolving all the other examples discussed so far. The measure proposed in [1] is essentially the information-theoretic sub-volume of the intersection between $I(Z; \hat{Y} \mid X_c)$ and total bias $I(Z; (\hat{Y}, U_X))$, which goes to 0 whenever either of them is 0 (details are provided in Appendix C-C)[10].*

The property of non-exempt masked bias stated in [1] is built on the rationale that in the example of non-exempt masked bias in hiring ads (Canonical Example 4 where $\hat{Y} = Z \oplus U_{X_1}$), instead of $U_{X_1}$ being the coding-test score, if $U_{X_1}$ is a random coin flip used to randomize the race, then this scenario may not necessarily be regarded as non-exempt. Then, we would have $X_c = \phi$ and $X_g = (Z, U_{X_1})$, and the Markov chain $Z - X_c - \hat{Y}$ would hold, deeming this example as *exempt*. In [1], the goal was to only account for non-exempt masked bias in $M_{NE}$ when the "mask" is either a critical feature or arises exclusively from the critical features, e.g., Canonical Example 4 while any mask from the non-critical/general features were viewed more like these random coin flips. But what if the user wishes to also account for masked bias if the mask is arising from $X_g$ as well, as demonstrated in the following modified version of the example?

**Canonical Example 5** (Non-Exempt Masked Bias in Hiring Ads II). *An ad for a job is only presented to men $(Z = 1)$ with a coding-test score above a threshold $(U_{X_1} = 1)$, and to women $(Z = 0)$ with a coding-test score below a threshold $(U_{X_1} = 0)$ with $Z$ and $U_{X_1}$ being i.i.d. Bern(½). The model output is given by $\hat{Y} = Z \oplus U_{X_1}$. Here, $Z \in X_g$ but $U_{X_1}$ is not be a critical feature for the job.*

Canonical Example 5 with $X_c = \phi$ and $X_g = (Z, U_{X_1})$ will be deemed *exempt* by [1] because the Markov chain $Z - X_c - \hat{Y}$ holds. However, here the virtual constituent or proxy $Z$ is arising from $X_g$ and is being masked by another feature of $X_g$, i.e., $U_{X_1}$. If $U_{X_1}$ denotes coding-test score and $\hat{Y}$ denotes the decision of showing hiring ads, then the model is again unfair to high-scoring women. This argument is also supported by the fact that the total bias is non-zero (not counterfactually fair). Since $X_c = \phi$, no bias is exempt, and a measure of non-exempt bias should ideally capture the total bias in this model.

In this work, we would like to arrive at an alternate criterion (modification of the property of non-exempt masked bias in [1]) that can capture non-exempt masked bias irrespective of whether the "mask" arises from the critical or general features. What this means is that any scenario deemed exempt by the property of non-exempt masked bias in [1] will also be deemed exempt by our modified property[11] but it is desirable that our modified property also accounts for scenarios, such as Canonical Example 5, that is sometimes deemed exempt by the former property even though intuitively, it may not be reasonable to do so.

*5) Leveraging Latent Variables to Understand Non-Exempt Masked Bias:* One commonality that we notice in the examples so far (1-5) is that whenever it is desirable that $M_{NE}$ be zero, either there is no counterfactual causal influence of $Z$ on $\hat{Y}$ (i.e., $\text{CCI}(Z \to \hat{Y}) = 0$) or the influence of $Z$ on $\hat{Y}$ has propagated *only* along paths that pass through $X_c$. In scenarios where $\text{CCI}(Z \to \hat{Y}) \neq 0$, one may choose to define another candidate measure of non-exempt bias that is inspired from the notion of

---

[10]One might also wonder why a measure of the form of a product, i.e., $M_{NE} = I(Z; \hat{Y} \mid X_c) \times I(Z; (\hat{Y}, U_X))$ does not work instead. We discuss a counterexample for such a product measure in [1] that we also include in Appendix C-C here for completeness.

[11]We show in Lemma 2 that the Markov chain in our modified property, i.e., $(Z, U_a) - X_c - (\hat{Y}, U_b)$ also implies $Z - X_c - \hat{Y}$, but the opposite implication is not true.
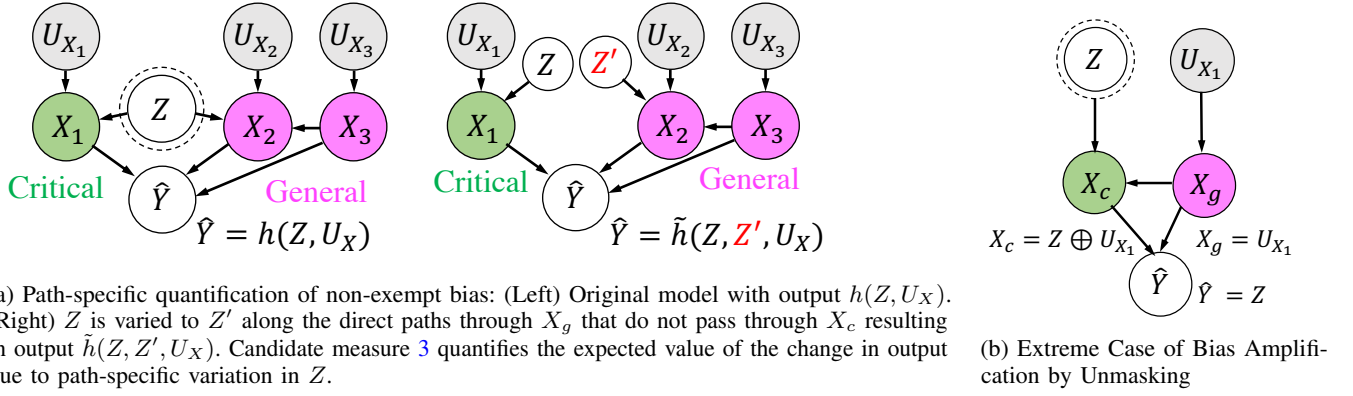
(a) Path-specific quantification of non-exempt bias: (Left) Original model with output $h(Z, U_X)$. (Right) $Z$ is varied to $Z'$ along the direct paths through $X_g$ that do not pass through $X_c$ resulting in output $\tilde{h}(Z, Z', U_X)$. Candidate measure 3 quantifies the expected value of the change in output due to path-specific variation in $Z$.

(b) Extreme Case of Bias Amplification by Unmasking

Fig. 7: Path-specific quantification of non-exempt bias (Candidate Measure 3) and its limitation

*path-specific counterfactual fairness* [19] (also see [16], [17]). This candidate measure for quantifying non-exempt bias is a causal, path-specific quantification by varying $Z$ only along the paths through $X_g$ that do not pass through $X_c$ and comparing if it causes any change in the model output (also see Fig. 7a).

**Candidate Measure of Non-Exempt Bias 3.** *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 relevant direct edges from $Z$ to $X_g$ with an edge from $Z'$ to $X_g$. Let $\hat{Y} = \tilde{h}(Z, Z', U_X)$ in the new causal graph. A candidate measure is $M_{NE} = \mathbb{E}_{Z,Z',U_X} \left[ |h(Z, U_X) - \tilde{h}(Z, Z', U_X)| \right]$.*

This measure, when used in conjunction with $\text{CCI}(Z \to \hat{Y}) = 0$, resolves the examples so far (1-5). For Canonical Example 1, it is zero and for Canonical Example 2, it is non-zero, as desired. For Canonical Example 3, $\text{CCI}(Z \to \hat{Y}) = 0$, and hence there is no need for a path-specific examination. For the example of non-exempt masked bias (Canonical Examples 4 and 5), *this measure is $0$ in spite of the statistically visible bias $\text{I}(Z; \hat{Y})$ being* $0$. However, the following example exposes some of its limitations.

**Canonical Example 6** (Bias Amplification by Unmasking). *Let $U_{X_1}$ be the inner ability of a candidate, and suppose that $X_c = Z + U_{X_1}$ denote the coding test score. Also let $X_g = U_{X_1}$ be the aptitude-test score where $Z$ and $U_{X_1}$ are i.i.d. Bern(½). Let the hiring decision be based on $\hat{Y} = X_c - X_g = Z$. This is shown in Fig. 5f with a more extreme modification in Fig. 7b.*

The bias in this example will be deemed *exempt* by a causal path-specific examination. However, this model has statistically visible bias ($\text{I}(Z; \hat{Y}) > 0$) that cannot be attributed to $X_c$ alone. Following the PID literature, here $X_c$ and $X_g$ have synergistic information about $Z$ that ultimately appears in $\hat{Y}$ which in itself is the virtual constituent or proxy of $Z$ being formed in this model. This synergistic information cannot be attributed to $X_c$ alone because $\text{I}(Z; X_c)$ is much smaller that $\text{I}(Z; \hat{Y})$. This is further supported by the argument that $X_g$ 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 bias. Thus, $M_{NE}$ should be greater than 0. Also, note that, here $\text{Uni}(Z : \hat{Y}|X_c) > 0$ (Supporting Derivation 2 in Appendix C-B) because it is this "joint" information about $Z$ in $(X_c, X_g)$ that ultimately appears in $\hat{Y}$ that cannot be attributed to $X_c$ alone.

Ideally, we would like a property and a measure that captures the intuition in this example. From a causal perspective, here $U_{X_1}$ is a confounder [35] to both $X_c$ and $\hat{Y}$, i.e., an extraneous variable that influences both of them along separate paths. A scenario when there is no non-exempt bias would be: (i) All causal paths from $Z$ to $\hat{Y}$ in the SCM pass through $X_c$; and also (ii) No $U_{X_i}$ acts as a confounder for both $X_c$ and $\hat{Y}$. This leads to the intuition that to be able to say $M_{NE} = 0$, one might be able to divide $U_X$ into two subsets $U_a$ and $U_b$ (further functional generalizations discussed in Section VIII), such that: (i) $U_a$ consists of the latent factors that do not influence $\hat{Y}$ at all, or influence it only through $X_c$ without acting as confounder; (ii) On the other hand, $U_b$ consists of the remaining latent factors, that only influence $\hat{Y}$ and not $X_c$; and (iii) The Markov chain $(Z, U_a) - X_c - (\hat{Y}, U_b)$ holds.

To understand this better, we again revisit Canonical Example 1 (visualization in Fig. 5a). Intuitively, the total bias in this example is exempt because $Z$ was already masked by $U_{X_1}$ in $X_c$, and the mask remained untampered in the final output $\hat{Y}$ with only additional independent masks added inside the black-box model. Here, neither $Z - X_c - (\hat{Y}, U_X)$ nor $(Z, U_X) - X_c - \hat{Y}$ hold, but $(Z, U_{X_1}) - X_c - (\hat{Y}, U_{X_2})$ does. A Markov chain of the form $(Z, U_a) - X_c - (\hat{Y}, U_b)$ also implies both the criterion $(Z, U_a) - X_c - \hat{Y}$ and $Z - X_c - (\hat{Y}, U_b)$ (see Lemma 2 with proof in Appendix C-A). One can interpret $U_a$ as the latent variables that either do not influence $\hat{Y}$ at all or already mask $Z$ in $X_c$ and remain untampered in the final output $\hat{Y}$. On the other hand, $U_b$ consists of the remaining latent variables that contribute to "additional masking inside the black-box model."

This leads us to propose the following criterion for $M_{NE}$ that also serves as our main rationale for Property 3: *$M_{NE}$ should be $0$ if $(Z, U_a) - X_c - (\hat{Y}, U_b)$ form a Markov chain for some subsets $U_a, U_b \subseteq U_X$ such that $U_a = U_X \backslash U_b$.*

**Lemma 2.** *The Markov chain* $(Z, U_a) - X_c - (\hat{Y}, U_b)$ *implies that the following Markov chains also hold: (i)* $Z - X_c - \hat{Y}$; *(ii)* $(Z, U_a) - X_c - \hat{Y}$; *and (ii)* $Z - X_c - (\hat{Y}, U_b)$.

The Markov chain $(Z, U_a) - X_c - (\hat{Y}, U_b)$ holding implies $M_{NE}=0$, but the Markov chain not holding for all $U_a, U_b$ such that $U_a = U_X \backslash U_b$ does not necessarily imply that $M_{NE} \neq 0$. This criterion $(Z, U_a) - X_c - (\hat{Y}, U_b)$ implying $M_{NE} = 0$ only attempts to provide an upper bound on $M_{NE}$, i.e., it is desirable that $M_{NE} \leq \min_{U_a, U_b \text{ s.t. } U_a=U_X \backslash U_b} \mathrm{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c)$ such that $U_a = U_X \backslash U_b$. The measure $\min_{U_a, U_b \text{ s.t. } U_a=U_X \backslash U_b} \mathrm{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c)$ does not suffice in itself as a measure of non-exempt bias because it again does not satisfy Property 1. To see this, notice that $\min_{U_a, U_b \text{ s.t. } U_a=U_X \backslash U_b} \mathrm{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c) \geq \mathrm{I}(Z; \hat{Y} \mid X_c)$ (see proof of Lemma 2), and thus, it also gives a false positive conclusion about non-exempt bias in Canonical Example 3 (counterfactually fair hiring). Instead, $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b) \mid X_c)$ is a sub-component of $\mathrm{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c)$ that satisfies Property 1. Our desirable properties ultimately leads us to our proposed measure of non-exempt bias, given by:

$$M_{NE}^* = \min_{U_a, U_b} \mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \text{ such that } U_a = U_X \backslash U_b.$$

*6) Our Proposed Measure Resolves all the Canonical Examples:* To develop intuition on what our proposed measure captures, we will now discuss how this measure resolves all of the examples in this work. We group "similar" examples together.

- **Scenarios where total bias** $\mathrm{I}(Z; (\hat{Y}, U_X))$ **is zero:** This applies to Canonical Example 3 and the related example in Remark 8. Because $\min_{U_a, U_b \text{ s.t. } U_a=U_X \backslash U_b} \mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \leq \mathrm{Uni}(Z : (\hat{Y}, U_X)|X_c) \leq \mathrm{I}(Z; (\hat{Y}, U_X))$ (see proof of Theorem 1 in Appendix C-A), it satisfies Property 1 and goes to 0 whenever total bias is 0.

- **Scenarios where** $Z$ **is already masked in** $X_c$ **and remains so in the output (with or without additional independent masks):** This applies to Canonical Example 1. We will examine the value of $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c)$ for different choices of $U_a \subseteq U_X$ to find the minimum. First notice that, if $U_a = \phi$ (and $U_b = U_X$), we have

$$\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = \mathrm{Uni}(Z : (\hat{Y}, U_X)|X_c) \overset{(a)}{\geq} \mathrm{Uni}(Z : Z|X_c) > 0$$

(see Supporting Derivation 3 in Appendix C-B; (a) holds from a monotonicity property of unique information because $Z$ can be obtained from deterministic local operations on $(\hat{Y}, U_X)$). This is in agreement with the intuition that $U_{X_1}$ should not belong to the set of candidate masks ($U_b$) that need to be accounted for. Next, if $U_a = U_{X_1}$ (and $U_b = U_{X_2}$), we have $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = 0$ (implied from the Markov chain $(Z, U_{X_1}) - X_c - (\hat{Y}, U_{X_2})$). Since unique information is non-negative, we therefore have $\min_{U_a, U_b \text{ s.t. } U_a=U_X \backslash U_b} \mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = 0$. In essence, the pair $(U_a^*, U_b^*)$ that minimizes $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c)$ is such that $U_a^* = U_{X_1}$, and the candidate masks that need to be accounted for, i.e., $U_b^* = U_{X_2}$.
  Now, what happens to the value of $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c)$ if the accountable mask $U_{X_2}$ is instead in $U_a$? We have

$$\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \overset{(a)}{\geq} \mathrm{Uni}(U_{X_2} : \hat{Y}|X_c) \overset{(b)}{=} \mathrm{I}(U_{X_2}; \hat{Y}),$$

  which is strictly greater than 0. This agrees with the intuition that $U_{X_2}$ should belong to the candidate set of masks that one should account for ($U_b$). Here (a) holds using two monotonicity properties of unique information (see Properties 10 and 9 in Appendix B) and (b) holds because $\mathrm{I}(U_{X_2}; X_c) = 0$, leading to $\mathrm{Red}(U_{X_2} : (\hat{Y}, X_c)) = 0$.

- **Scenarios where non-exempt statistically visible bias is present, i.e.,** $\mathrm{Uni}(Z : \hat{Y}|X_c) > 0$**:** This applies to Canonical Example 2 and Canonical Example 6. Because $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \geq \mathrm{Uni}(Z : \hat{Y}|X_c)$ (see proof of Theorem 1 in Appendix C-A), our proposed $M_{NE}^*$ satisfies Property 2, and is thus non-zero whenever $\mathrm{Uni}(Z : \hat{Y}|X_c) > 0$.

- **Scenarios where non-exempt masked bias is present:** This applies to Canonical Example 4 and Canonical Example 5. In the proof of Theorem 1 in Appendix C-A, we show that the proposed measure satisfies Property 3 (non-exempt masked bias), and is thus non-zero for these canonical examples of non-exempt masked bias.
  We note that Canonical Example 2 is an interesting case where both non-exempt statistically visible bias and non-exempt masked bias are present. Here, $M_{NE}^*$ is strictly greater than the non-exempt statistically visible bias ($\mathrm{Uni}(Z : \hat{Y}|X_c)$), and this difference can be interpreted as a quantification of the non-exempt masked bias. First notice that,

$$\mathrm{Uni}(Z : \hat{Y}|X_c) \overset{(a)}{=} \mathrm{I}(Z; \hat{Y}) = \mathrm{H}(Z) - \mathrm{H}(Z|\hat{Y}) = \mathrm{H}(Z) - \mathrm{H}(Z|U_{X_1} + Z + U_{X_2}) = 1 - \frac{3}{4}h_b(1/3) \text{ bits.} \quad (6)$$

The full derivation is in Supporting Derivation 4 in Appendix C-B. Here $h_b(\cdot)$ is the binary entropy function [63] given by $h_b(p) = -p \log_2(p) - (1 - p) \log_2(1 - p)$ and (a) holds because $\mathrm{I}(Z; U_{X_1}) = 0$, implying $\mathrm{Red}(Z : (\hat{Y}, U_{X_1})) = 0$ as well. Now, we will examine the value of $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c)$ for different choices of $U_a$ to find the minimum. The full derivation for all of these cases is in Supporting Derivation 4 in Appendix C-B. Here, we only mention the key step. Let $U_a = \phi$ (and $U_b = U_X$). Then,

$$\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = \mathrm{Uni}(Z : (\hat{Y}, U_{X_1}, U_{X_2})|U_{X_1}) \overset{(a)}{=} \mathrm{I}(Z; U_{X_1} + Z + U_{X_2}, U_{X_1}, U_{X_2}) = 1 \text{ bit.} \quad (7)$$

Here (a) holds again because $I(Z; U_{X_1}) = 0$, implying the redundant information is 0 as well (using (2) in Section II-A). Next, for $U_a = U_{X_2}$ (and $U_b = U_{X_1}$), we have,

$$\text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) = \text{Uni}((Z, U_{X_2}) : (\hat{Y}, U_{X_1}) | U_{X_1}) \stackrel{(a)}{=} I((Z, U_{X_2}); (\hat{Y}, U_{X_1})) = 3/2 \text{ bit.} \tag{8}$$

Here (a) holds again because $I((Z, U_{X_2}); U_{X_1}) = 0$, implying the redundant information is 0 as well. Next, for $U_a = U_{X_1}$ (and $U_b = U_{X_2}$), we have,

$$\text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) = \text{Uni}((Z, U_{X_1}) : (\hat{Y}, U_{X_2}) | U_{X_1}) \stackrel{(b)}{=} I((Z, U_{X_1}); (\hat{Y}, U_{X_2}) \mid U_{X_1}) = 1 \text{ bit.} \tag{9}$$

Here (b) holds because $\text{Syn}((Z, U_{X_1}) : (A, B)) = 0$ if one of the terms $A$ or $B$ is a deterministic function of $(Z, U_{X_1})$ (using Lemma 14 in Appendix B) and hence unique information becomes equal to the conditional mutual information (see (3) in Section II-A). Lastly, for $U_a = U_X$ (and $U_b = \phi$), we have,

$$\text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) = \text{Uni}((Z, U_{X_1}, U_{X_2}) : \hat{Y} | U_{X_1}) \stackrel{(b)}{=} I((Z, U_{X_1}, U_{X_2}); \hat{Y} \mid U_{X_1}) = 3/2 \text{ bit.} \tag{10}$$

Here (b) holds again using Lemma 14 in Appendix B. Thus, we obtain that,

$$M_{NE}^* = \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) = 1 \text{ bit,}$$

which is strictly greater than $\text{Uni}(Z : \hat{Y} | X_c) = 1 - \frac{3}{4} h_b(1/3)$ bits, accounting for both non-exempt statistically visible and non-exempt masked biases.

As noted in Remark 5, our properties are insufficient to arrive at a unique functional form for the measure of non-exempt bias. It is easiest to understand this issue by contrasting it with Shannon's discussion on entropy as a measure for uncertainty. First, we do not have a counterpart of "additivity" of entropy (see Property 3 in Section 6 of [58]) which allows Shannon to arrive at the logarithmic scaling in entropy. Second, we also do not provide an operational meaning for this measure (such as that provided by the lossless source coding theorem for entropy [63]), which further supports the logarithmic scaling. This is a direction of meaningful future work (further functional generalizations discussed in Section VIII). We note that this is the case with almost all existing measures of fairness (with the notable exceptions of [21], [55], [60]). Exploring more deeply the desirable attributes of the influence of a virtual constituent or proxy of $Z$ that influences the model output and that cannot be attributed to the critical features $X_c$ alone (inspired from the work on proxy-use [20]) could be a starting point towards deriving an exact operational meaning for our proposed measure. Nonetheless, our measure does satisfy all six desirable properties, and also captures important nuances of the problem, e.g., both non-exempt masked bias and non-exempt statistically visible bias when they are present together (revisited in Section IV). Our examples also help us understand the utility and limitations of some existing measures that have some provision for exemptions, as we discuss next.

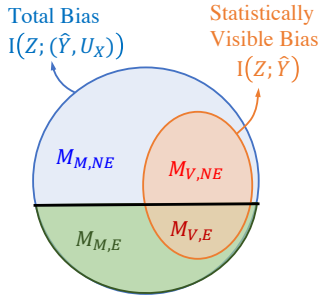## C. Understanding Existing Measures of Fairness with Provision for Exemptions

**Conditional Statistical Parity:** This definition [40], [42] is equivalent to $I(Z; \hat{Y} \mid X_c) = 0$. Therefore, it has similar utility and limitations as Candidate Measure 1 ($I(Z; \hat{Y} \mid X_c)$). It resolves some limitations of both statistical parity and equalized odds. However, it gives a false positive conclusion in detecting non-exempt bias in Canonical Example 3 (the example of counterfactually fair hiring), where there is no causal influence of $Z$ on $\hat{Y}$ but $I(Z; \hat{Y} \mid X_c) > 0$. Because this is an observational measure, it is not able to distinguish between scenarios where there is causal influence of $Z$ on $\hat{Y}$ (non-exempt masked bias in hiring ads; Canonical Example 4) and where there is not (Canonical Example 3), even if $I(Z; \hat{Y} \mid X_c) > 0$ in both (elaborated further in relation to our impossibility result in Remark 12 Section V). It also fails to capture non-exempt masked bias when the mask arises from the general features as in Canonical Example 5.

**Justifiable Fairness:** A model is said to be justifiably fair [41] if $I(Z; \hat{Y} \mid X_s) = 0$ for all sets $X_s \subseteq X$ such that $X_c \subseteq X_s$. This measure addresses several concerns of the previously stated measures, including capturing several forms of non-exempt masked bias. However, it also gives false positive conclusion in Canonical Example 3 (counterfactually fair college admissions), which shows no causal influence of $Z$ on $\hat{Y}$ but $I(Z; \hat{Y} \mid X_c) > 0$. Because this is an observational measure, it is not able to distinguish between scenarios where there is causal influence of $Z$ on $\hat{Y}$ and where there is not, even if $I(Z; \hat{Y} \mid X_c) > 0$ in both (elaborated further in relation to our impossibility result in Remark 12 Section V).

Another limitation of such an individual feature-based conditioning arises when the causal effects of both $Z$ and an independent latent factor are present in the same feature, e.g., different digits of a zip-code, and it is not known in advance whether to condition on the entire zip-code or its sub-portions like the individual digits.

**Scenario 3** (Special Case of Canonical Example 5). *Let $X_g = [Z, U_{X_1}]$ be a single multivariate feature, e.g., two bits of a number and $X_c = \phi$, and the output be $\hat{Y} = Z \oplus U_{X_1}$ where $Z$ and $U_{X_1}$ are i.i.d. Bern(1/2).*

In this example, as long as one treats $X_g$ as a single feature, the model will be deemed *justifiably fair* because $I(Z; \hat{Y} \mid X_g) = 0$ and $I(Z; \hat{Y}) = 0$. But, this is a case of non-exempt masked bias. It is necessary to have an advance suspicion of this possible

(a) Venn diagram representation of overall decomposition

|  | $I(Z; \hat{Y})$ | $I(Z; U_X | \hat{Y})$ | $I(Z; (\hat{Y}, U_X))$ |
|---|---|---|---|
| $M_{NE}$ | $M_{V,NE} = Uni(Z:\hat{Y}|X_c)$ | $M_{M,NE} = M_{NE} - M_{V,NE}$ | |
| $M_E$ | $M_{V,E} = Red(Z:(\hat{Y}, X_c))$ | $M_{M,E} = M_E - M_{V,E}$ | |
|  | | | $I(Z; (\hat{Y}, U_X))$ |

(b) Tabular representation of overall decomposition

Fig. 8: Overall decomposition of total bias $\mathrm{I}(Z; (\hat{Y}, X_c))$ into four non-negative components, namely, non-exempt visible bias $M_{V,NE}$, exempt visible bias $M_{V,E}$, non-exempt masked bias $M_{M,NE}$ and exempt masked bias $M_{M,E}$.

nature of the true SCM to be able to condition on the two bits of $X_g$ separately. This definition captures the non-exempt masked bias in this example if the sub-portions of any single feature are defined in advance.

**Path-Specific Counterfactual Fairness:** Path-specific counterfactual fairness [19] is a purely causal notion of fairness which exempts the causal influence of $Z$ along selected paths. Based on this idea, we proposed Candidate Measure 3 in Section III-B. However, Canonical Example 6 (the example of discrimination by unmasking) captures some of its limitations, when there is synergistic or joint information about $Z$ present in $X_c$ and $X_g$ that appears in $\hat{Y}$ that cannot be attributed to any one of them alone. Furthermore, sometimes the influence of $Z$ can cancel along two paths so that the final output has no influence of $Z$, e.g., the example in Remark 8. For such scenarios, this measure alone can lead to false positive conclusions about non-exempt bias, and might need to be used in conjunction with a measure of total bias (e.g., $\mathrm{CCI}(Z \to \hat{Y})$).

## IV. UNDERSTANDING THE OVERALL DECOMPOSITION

In this section, we demonstrate how our proposed quantification enables a *non-negative* information-theoretic decomposition of the total bias $\mathrm{I}(Z; (\hat{Y}, U_X))$ into four components, that can be interpreted as: statistically visible non-exempt bias, statistically visible exempt bias, masked non-exempt bias and masked exempt bias (also see Fig. 8).

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

$$\mathrm{I}(Z; (\hat{Y}, U_X)) = M_{V,NE} + M_{V,E} + M_{M,NE} + M_{M,E}. \qquad (11)$$

*Here $M_{V,NE} = \mathrm{Uni}(Z : \hat{Y}|X_c)$ and $M_{V,E} = \mathrm{Red}(Z : (\hat{Y}, X_c))$. These two terms add to form $\mathrm{I}(Z; \hat{Y})$ which is the total statistically visible bias. Next, $M_{M,NE} = M^*_{NE} - M_{V,NE}$ where $M^*_{NE}$ is our proposed measure of non-exempt bias (Definition 7), and $M_{M,E} = \mathrm{I}(Z; \hat{Y}, U_X) - \mathrm{I}(Z; \hat{Y}) - M_{M,NE}$. All of these components are non-negative.*

The decomposition of total bias into a summation of these four terms is trivial. What remains to be shown is that these four terms are non-negative (details provided in Appendix D-A).

**Interpretation of the four components:** Here $M_{V,NE} = \mathrm{Uni}(Z : \hat{Y}|X_c)$ can be interpreted as the non-exempt statistically visible bias (as also motivated in Section III-B). The remaining part of the statistically visible bias (recall Definition 5), i.e., $\mathrm{I}(Z; \hat{Y}) - \mathrm{Uni}(Z : \hat{Y}|X_c) = \mathrm{Red}(Z : (\hat{Y}, X_c))$ then becomes the exempt statistically visible bias ($M_{V,E}$). This also agrees with the intuition that redundant information about $Z$ visible in both $\hat{Y}$ and $Z$ represents the exempt statistically visible bias.

Now that we have a measure of non-exempt bias ($M^*_{NE}$) and a measure of non-exempt statistically visible bias ($M_{V,NE}$), we can interpret their difference as the non-exempt masked bias, i.e., $M_{M,NE} = M^*_{NE} - M_{V,NE} = M^*_{NE} - \mathrm{Uni}(Z : \hat{Y}|X_c)$. It also agrees with the intuition that non-exempt masked bias is the part of non-exempt bias that $\mathrm{Uni}(Z : \hat{Y}|X_c)$ alone fails to capture. For instance, recall Canonical Example 4 where $\hat{Y} = Z \oplus U_{X_1}$ and $X_c = U_{X_1}$. Here, $\mathrm{I}(Z; \hat{Y}) = 0$, implying $M_{V,NE} = \mathrm{Uni}(Z : \hat{Y}|X_c) = 0$. But, $M^*_{NE} = 1$ bit (supporting derivation in Appendix C-A; see the proof of Theorem 1 under Property 3). Therefore, the non-exempt masked bias $M_{M,NE} = M^*_{NE} - M_{V,NE} = 1$ bit here, which is in agreement with our intuition of non-exempt masked bias. Lastly, the remaining component $M_{M,E} = \mathrm{I}(Z; \hat{Y}, U_X) - \mathrm{I}(Z; \hat{Y}) - M_{M,NE}$ is interpreted as the exempt masked bias. For instance, recall Canonical Example 1 where $\hat{Y} = X_c = Z + U_{X_1} + U_{X_2}$ with $Z, U_{X_1}, U_{X_2} \sim i.i.d.$ Bern(½). Here, the total bias $\mathrm{I}(Z; \hat{Y}, U_X) = 1$ bit, but the statistically visible bias $\mathrm{I}(Z; \hat{Y}) = 0.5$ bits which means that there is masked bias present. Our intuition is that this masked bias should be entirely exempt because there is no non-exempt bias in this example. This is in agreement with the value that we obtain, i.e., $M_{M,E} = \mathrm{I}(Z; \hat{Y}, U_X) - \mathrm{I}(Z; \hat{Y}) - M_{M,NE} = 0.5$ bits. This is because $M_{M,NE}$ and $M_{V,NE}$ are both non-negative sub-components of $M^*_{NE}$, and $M^*_{NE} = 0$ (from the Markov chain $(Z, U_{X_1}, U_{X_2}) - X_c - \hat{Y}$).

**Remark 10** (On conditioning to capture masked bias). *Conditioning on a random variable $G$ leading to $\mathrm{I}(Z; \hat{Y} | G) > \mathrm{I}(Z; \hat{Y})$ can sometimes detect masked bias, if conditioning exposes more bias than what was already visible. For example, $\mathrm{I}(Z; \hat{Y} | X_c)$*

*can detect masked bias if the mask is of the form* $g(X_c)$, *e.g., in Canonical Example 4 (a special case of the canonical example of masking with* $X_c = U_{X_1}$ *and* $\hat{Y} = Z \oplus U_{X_1}$*). However, conditioning on any random variable G leading to* $\mathrm{I}(Z; \hat{Y} \mid G) > \mathrm{I}(Z; \hat{Y})$ *cannot always be interpreted as a case of masked bias because this can sometimes lead to a false positive conclusion in detecting masked bias, e.g., in Canonical Example 3 where* $\hat{Y} = U_{X_1}$ *and* $X_c = Z + U_{X_1}$. *If G is chosen as* $X_c$, *then* $\mathrm{I}(Z; \hat{Y} \mid X_c) > \mathrm{I}(Z; \hat{Y})$ *even though there is no bias here at all (recall* $\mathrm{CCI}(Z \to \hat{Y}) = 0$*). For completeness, we therefore include another result here (Lemma 3) that clarifies when conditioning can correctly capture masked bias.*

**Lemma 3** (Conditioning to Capture Masked Bias)**.** *The following two statements are equivalent:*
- *Masked bias* $\mathrm{I}(Z; (\hat{Y}, U_X)) - \mathrm{I}(Z; \hat{Y}) > 0$.
- $\exists$ *a random variable G of the form* $G = g(U_X)$ *such that* $\mathrm{I}(Z; \hat{Y} \mid G) - \mathrm{I}(Z; \hat{Y}) > 0$.

Without knowledge of the true causal model, such a $G = g(U_X)$ may be difficult to determine from observational data alone, because the observational data can be a function of both $Z$ and $U_X$. This serves as the motivation behind our impossibility result on observational measures, that we state next.

## V. IMPOSSIBILITY RESULT

**Theorem 3** (Impossibility of Observational Measures)**.** *No observational measure of non-exempt bias simultaneously satisfies all six desirable properties.*

*Proof of Theorem 3.* Observe the two examples here:

**Example 1** (A Case of No Bias)**.** *Let* $X_c = Z \oplus U_{X_1}$, $X_g = Z$ *and* $\hat{Y} = X_c \oplus X_g = U_{X_1}$ *where* $Z$ *and* $U_{X_1}$ *are both independent and identically distributed as Bern(½).*

**Example 2** (A Case of Non-Exempt Bias)**.** *Let* $X_c = U_{X_1}$, $X_g = Z$ *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 Bern(½).*

In Example 1, the influences of $Z$ cancel each other and there is no total bias. So, the non-exempt bias should be zero by Property 1 (Zero Influence). However, Example 2 is the canonical example of non-exempt masked bias where there is non-exempt bias present, and hence the non-exempt bias should be non-zero by Property 3 (Non-Exempt Masked Bias). But, for both of 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 proves the result. $\square$

**Remark 11** (Alternative Examples)**.** *In fact, we can show that no observational measure can satisfy Property 3. Consider a scenario of no bias given by:* $X_c = \phi$, $X_g = (Z \oplus U_{X_1}, Z)$ *and* $\hat{Y} = U_{X_1}$. *For this example, the Markov chain* $Z - X_c - (\hat{Y}, U_{X_1})$ *holds implying that* $M_{NE} = 0$ *by Property 3. Alternatively, consider a scenario of non-exempt bias given by:* $X_c = \phi$, $X_g = (U_{X_1}, Z)$ *and* $\hat{Y} = Z \oplus U_{X_1}$ *which is again a variant of the canonical example of non-exempt masked discrimination. Let* $Z$ *and* $U_{X_1}$ *be independent and identically distributed as Bern(½). Then, no purely observational measure can distinguish between these two scenarios because* $(Z, X_c, X_g, \hat{Y})$ *have the same joint distribution.*

**Remark 12** (Revisiting Conditional Statistical Parity and Justifiable Fairness)**.** *For both Examples 1 and 2, we observe that conditional mutual information* $\mathrm{I}(Z; \hat{Y} \mid X_c) > 0$. *Because* $\mathrm{I}(Z; \hat{Y} \mid X_c)$ *is an observational measure, it fails to distinguish between whether there is causal influence of Z or not in* $\hat{Y}$. *Existing observational definitions of fairness, e.g., conditional statistical parity and justifiable fairness would also not be able to distinguish between these two examples. One needs counterfactual measures to be able to distinguish between them, such as the counterfactual measure proposed in this work.*

Nevertheless, because counterfactual measures are difficult to realize in practice, we examine the following observational measures of non-exempt bias that satisfy only a few of Properties 1-6.

## VI. OBSERVATIONAL RELAXATIONS OF OUR PROPOSED COUNTERFACTUAL MEASURE: UTILITY AND LIMITATIONS

In this section, we propose three observational measures of non-exempt bias and discuss their utility and limitations.

**Observational Measure 1.** $M_{NE} = \mathrm{Uni}(Z : \hat{Y} | X_c)$.

**Utility:** This measure satisfies several desirable properties as stated here:

**Lemma 4.** *[Fairness Properties of* $\mathrm{Uni}(Z : \hat{Y} | X_c)$*] The measure* $\mathrm{Uni}(Z : \hat{Y} | X_c)$ *satisfies Properties 5, 6, 1, and 2.*

The proof is in Appendix E. Importantly, note that, $\mathrm{Uni}(Z : \hat{Y} | X_c)$ satisfies Property 1 which $\mathrm{I}(Z; \hat{Y} \mid X_c)$ does not (recall Canonical Example 3). Thus, $\mathrm{Uni}(Z : \hat{Y} | X_c)$ does not give false positive conclusions in detecting non-exempt bias if a model is counterfactually fair.

This measure may be preferred over our other observational measures when one wants to prioritize avoiding false positive quantification of non-exempt bias when a model is counterfactually fair. Recall that, $\mathrm{Uni}(Z : \hat{Y} | X_c)$ is a measure of non-exempt, statistically visible bias. *It correctly captures the entire non-exempt bias when non-exempt masked bias is absent.*

**Limitations:** It does not quantify any non-exempt masked bias (Property 3). This is because $\text{Uni}(Z:\hat{Y}|X_c)$ is a sub-component of the statistically visible bias $\text{I}(Z;\hat{Y})$, and hence always goes to 0 whenever the statistically visible bias $\text{I}(Z;\hat{Y}) = 0$ (recall Canonical Examples 4 and 5). It also does not satisfy Property 4 because when $X_c = \phi$, we have $\text{Uni}(Z:\hat{Y}|X_c) = \text{I}(Z;\hat{Y})$, which is only the statistically visible bias but not the total bias in a counterfactual sense (i.e., $\text{I}(Z;\hat{Y}, U_X)$).

**Observational Measure 2.** $M_{NE} = \text{I}(Z;\hat{Y} \mid X_c)$.

**Utility:** This measure also satisfies several desirable properties, as stated here:

**Lemma 5.** *[Fairness Properties of* $\text{I}(Z;\hat{Y} \mid X_c)$*] The measure* $\text{I}(Z;\hat{Y} \mid X_c)$ *satisfies Properties 6 and 2.*

The proof is in Appendix E. We note that, while it does not satisfy Property 3 in its entirely, it does capture some scenarios of non-exempt masked bias. E.g., it can detect the non-exempt masked bias in Canonical Example 4 which $\text{Uni}(Z:\hat{Y}|X_c)$ is not able to, even though they both fail to detect the non-exempt masked bias in Canonical Example 5. In general, $\text{I}(Z;\hat{Y} \mid X_c)$ can detect non-exempt masked bias when the "mask" is entirely derived from the critical features, i.e., $G = g(X_c)$.

**Limitations:** It can sometimes lead to false positive conclusion about non-exempt bias, e.g., in Canonical Example 3 (does not satisfy Property 1). It also does not satisfy Property 5 because clearly $\text{I}(Z;\hat{Y} \mid X_c)$ may be greater or less that $\text{I}(Z;\hat{Y})$ (recall Canonical Example 4). It also does not satisfy Property 4 because when $X_c = \phi$, we have $\text{I}(Z;\hat{Y} \mid X_c) = \text{I}(Z;\hat{Y})$, which is only the statistically visible bias but not the total bias in a counterfactual sense (i.e., $\text{I}(Z;\hat{Y}, U_X)$).

**Observational Measure 3.** $M_{NE} = \text{I}(Z;\hat{Y} \mid X_c, X')$ where $X'$ consists of certain features in $X_g$.

**Utility and Limitations:** This is somewhat of a heuristic relaxation that only satisfies Property 6. However, while it does not satisfy any of the other properties in their entirety, it can still lead to the desirable quantification in several examples where the previous two measures may not be successful if $X'$ is chosen appropriately. For example, recall Canonical Example 5 where $\hat{Y} = Z \oplus U_{X_1}$ with $X_g = (Z, U_{X_1})$. With some partial knowledge or assumption about the SCM, if we choose $X' = U_{X_1}$, then $\text{I}(Z;\hat{Y} \mid X_c, X') > 0$ for this example even though $\text{I}(Z;\hat{Y} \mid X_c) = 0$. Thus, this measure is able to detect some more scenarios of non-exempt masked bias that $\text{I}(Z;\hat{Y} \mid X_c)$ cannot, i.e., when the mask is of the form $G = g(X_c, X')$. It can also sometimes avoid false positive quantification of non-exempt bias if $X'$ is chosen appropriately, e.g., in Canonical Example 3 if $X' = U_{X_1}$. Thus, under partial knowledge or assumption about the true SCM, this measure can correctly capture the non-exempt bias in many scenarios where the previous two measures may not be successful.

Lastly, one may also consider using various combinations of these measures, e.g., $\text{Uni}(Z:\hat{Y}|X_c) + \text{I}(Z;\hat{Y} \mid X')$, or $\text{I}(Z;\hat{Y} \mid X_c) + \text{I}(Z;\hat{Y} \mid X')$, or $\text{Uni}(Z:\hat{Y}|X_c) + \text{Syn}(Z:(\hat{Y},X'))$, that can also approximate our proposed measure in several scenarios if $X'$ is chosen appropriately based on partial knowledge or assumptions about the true SCM.

## VII. CASE STUDIES DEMONSTRATING PRACTICAL APPLICATION IN AUDITING AND TRAINING

Here, we discuss some case studies to demonstrate application of our proposed techniques on both simulated and real data.

### A. Case Study on Simulated Data

We present our case study on simulated data first. The benefit of using simulated data is that the true causal model (SCM) is known. The knowledge of the SCM enables the following: (i) we can exactly compute our proposed causal measure of non-exempt bias ($M^*_{NE}$), as well as, demonstrate the decomposition of total bias into four components during auditing a pre-trained model; (ii) we can also compare the performance of different observational measure of non-exempt bias when used as a regularizer during training. Assuming the SCM is not available during training (but available during auditing), we examine the tradeoff between accuracy and the actual causal non-exempt bias ($M^*_{NE}$) when each of these observational measures are used as a regularizer, under various experimental scenarios.

In this case study, an algorithm has to decide whether to show ads for a job using a score generated from internet activity. We will consider four different experimental scenarios, each with a known SCM. To demonstrate application in **auditing**, we first train a Deep-Neural-Network (DNN) model with no fairness regularizer for each of the four scenarios, and then use our techniques for computing the total bias ($\text{I}(Z;(\hat{Y}, U_X))$), as well as, decompose the total bias into four components, namely, visible and masked, exempt and non-exempt biases. We use the dit [43] package to compute all of these quantities from the empirical distribution of the test data after the model has been trained, and after appropriately discretizing continuous random variables as required. Note that, to compute unique information, the package solves an optimization problem [43].

To demonstrate application in **training**, we train a DNN model $\hat{Y} = h(X)$ for classification with different **observational regularizers** and examine the tradeoff between accuracy and the *actual* non-exempt bias (as measured by our causal measure of non-exempt bias $M_{NE}*$), when each of these observational regularizers are used. For simplicity and ease of computation during training, we rely on simple correlation-based estimates (inspired from [10]) of mutual information and conditional mutual information. Further, we introduce a novel regularizer for approximating unique information, leveraging a Gaussian approximation for PID in [44]. We train using the following loss functions:

- Loss $L_1$ (**Statistical Parity** using Mutual Information regularizer $\text{I}(Z;\hat{Y})$ (denoted as MI)):

$$\min_{w,b} L_{\text{Cross Entropy}}(Y, \hat{Y}) + \lambda \widetilde{\text{I}}(Z;\hat{Y}),$$

where (i) $\lambda$ is the regularization constant; and (ii) $\widetilde{I}(Z;\hat{Y}) = -\frac{1}{2}\log{(1-\rho_{Z,\hat{Y}}^2)}$ 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 [63].

- Loss $L_2$ (Proposed Unique Information-based (observational) regularizer $\text{Uni}(Z:\hat{Y}|X_c)$ (denoted as Uniq)):

$$\min_{w,b} L_{\text{Cross Entropy}}(Y,\hat{Y})+\lambda\widetilde{\text{Uni}}(Z:\hat{Y}|X_c),$$

where $\widetilde{\text{Uni}}(Z:\hat{Y}|X_c)$ is given by:

$$\widetilde{\text{Uni}}(Z:\hat{Y}|X_c) = \widetilde{I}(Z;\hat{Y}) - \min\{\widetilde{I}(Z;\hat{Y}),\widetilde{I}(Z;X_c)\} = -\frac{1}{2}\log{(1-\rho_{Z,\hat{Y}}^2)} - \min\{-\frac{1}{2}\log{(1-\rho_{Z,\hat{Y}}^2)}, -\frac{1}{2}\log{(1-\rho_{Z,X_c}^2)}\}.$$

We note that, in general, $\text{Uni}(Z:\hat{Y}|X_c \geq I(Z;\hat{Y}) - \min\{I(Z;\hat{Y}),I(Z;X_c)\}$, where the lower bound is tight if all of the random variables are jointly Gaussian [44]. Similarly, the correlation-based approximations are also exact under Gaussian assumptions [63].

- Loss $L_3$ (Proposed Conditional Mutual Information regularizer $I(Z;\hat{Y}|X_c)$ (denoted as CMI)):

$$\min_{w,b} L_{\text{Cross Entropy}}(Y,\hat{Y})+\lambda\widetilde{I}(Z;\hat{Y}\mid X_c),$$

where again (i) $\lambda$ is the regularization constant; and (ii) $\widetilde{I}(Z;\hat{Y}\mid X_c)$ is given by:

$$\widetilde{I}(Z;\hat{Y}\mid X_c) = \sum_{i=1}^{n}\text{Pr}(X_c \in \text{Bin } i)\widetilde{I}(Z;\hat{Y}\mid X_c \in \text{Bin } i) = -\frac{1}{2}\sum_{i=1}^{n}\text{Pr}(X_c \in \text{Bin } i)\log{(1-\rho_{Z,\hat{Y},i}^2)},$$

where the range of $X_c$ is divided into $n$ discrete bins, and $\rho_{Z,\hat{Y},i}$ is the conditional correlation of $\hat{Y}$ and $Z$ given $X_c$ is in the $i$-th discrete bin.

- Loss $L_4$ (Another Proposed Heuristic regularizer $I(Z;\hat{Y}|X_c,X')$ (denoted as CMI')):

$$\min_{w,b} L_{\text{Cross Entropy}}(Y,\hat{Y})+\lambda\widetilde{I}(Z;\hat{Y}\mid X_c,X'),$$

where again (i) $\lambda$ is the regularization constant; and (ii) $\widetilde{I}(Z;\hat{Y}\mid X_c,X')$ is given by:

$$\widetilde{I}(Z;\hat{Y}\mid X_c,X') = \sum_{i=1}^{n}\text{Pr}(X_c,X' \in \text{Bin } i)\widetilde{I}(Z;\hat{Y}\mid X_c,X' \in \text{Bin } i) = -\frac{1}{2}\sum_{i=1}^{n}\text{Pr}(X_c,X' \in \text{Bin } i)\log{(1-\rho_{Z,\hat{Y},i}^2)},$$

where the range of the joint random variables $(X_c,X')$ is divided into $n$ discrete bins, and $\rho_{Z,\hat{Y},i}$ is the conditional correlation of $\hat{Y}$ and $Z$ given $(X_c,X')$ is in the $i$-th discrete bin.

- Loss $L_5$ (**Equalized Odds** using regularizer $I(Z;\hat{Y}|Y)$ (denoted as EO)):

$$\min_{w,b} L_{\text{Cross Entropy}}(Y,\hat{Y})+\lambda\widetilde{I}(Z;\hat{Y}\mid Y),$$

where again (i) $\lambda$ is the regularization constant; and (ii) $\widetilde{I}(Z;\hat{Y}\mid Y)$ is given by:

$$\widetilde{I}(Z;\hat{Y}\mid Y) = \sum_{i=1}^{n}\text{Pr}(Y \in \text{Bin } i)\widetilde{I}(Z;\hat{Y}\mid Y \in \text{Bin } i) = -\frac{1}{2}\sum_{i=1}^{n}\text{Pr}(Y \in \text{Bin } i)\log{(1-\rho_{Z,\hat{Y},i}^2)}.$$

The range of $Y$ is divided into $n$ discrete bins, and $\rho_{Z,\hat{Y},i}$ is the correlation of $\hat{Y}$ and $Z$ given $Y$ is in the $i$-th bin.

Now, we discuss the four scenarios (SCMs) and the corresponding results.

**Experimental Scenario 1 (All four biases present):** The decision of showing ads for a reporter's job requiring English proficiency, is based on three features $X = (X_1,X_2,X_3)$: (i) $X_1$: a score based on online writing samples (critical feature $X_c = X_1$); (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. $Z$ is a protected attribute denoting whether a person is a native English speaker or not, distributed as Bern(½). Suppose that the true SCM is as follows: $X_1 = Z + U_{X_1}$, $X_2 = Z + U_{X_2}$, and $X_3 = U_{X_3}$, where $U_{X_1},U_{X_2},U_{X_3} \sim i.i.d. \mathcal{N}(0,\sigma^2)$ denote latent writing ability, interests, and geographical proximity, respectively. The true labels, based on previous candidates, are given by $Y = \mathbb{1}(X_1 + X_2 + X_3 \geq 1)$. Here, the critical feature $X_c = X_1$ and the general features are $X_g = (X_2,X_3)$. The results are provided in Fig. 9a and Fig. 10a.

**Experimental Scenario 2 (Masking by critical feature):** The decision of showing ads for an editor's job in a newspaper company is based on four features: (i) $X_1$: a relevant score based on online writing samples (critical feature $X_c = X_1$); (ii) $X_2$: a score based on browsing history, e.g., awareness of current events; (iii) $X_3$: a score based on proofreading and reviewing experience; and (iv) $X_4$: a preference score based on activity in social media, e.g., political and ideological alignment with the newspaper company. Let the protected attribute $Z$ be political inclination, distributed as Bern(½). Suppose the true

(a) Experimental Scenario 1 (All four biases present)



(b) Experimental Scenario 2 (Masking by critical feature)



(c) Experimental Scenario 3 (Masking by general feature)
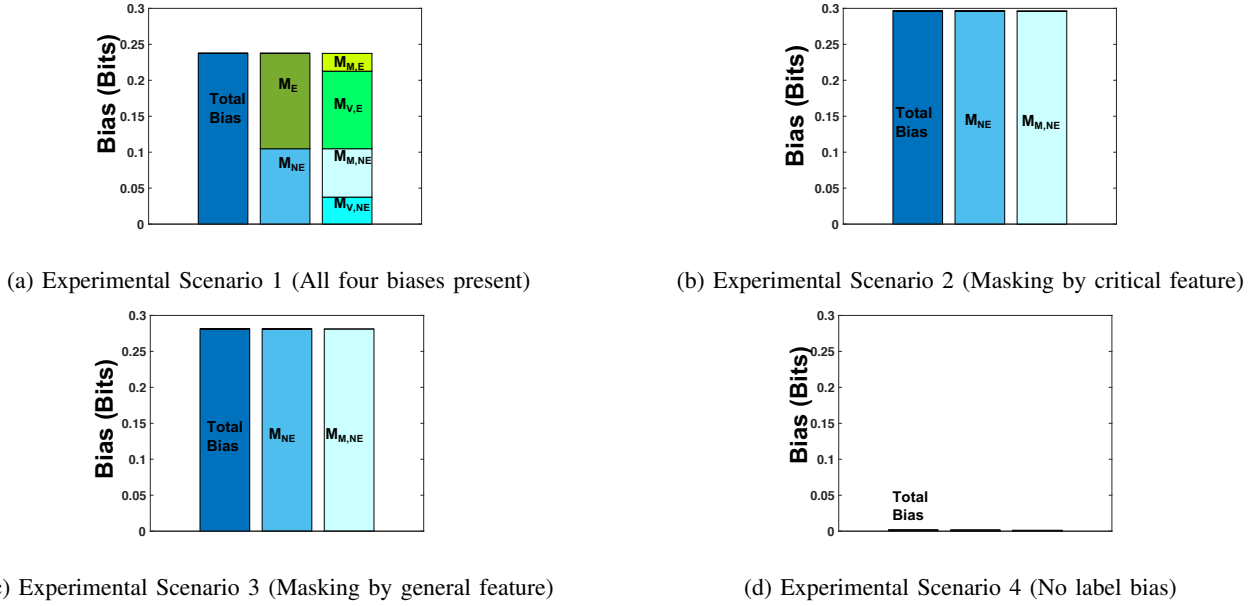


(d) Experimental Scenario 4 (No label bias)

Fig. 9: **Observations from Auditing.** The different types of biases after training a model with no fairness regularizer for all the experimental scenarios: $M_E$ and $M_{NE}(= M_{NE}*)$ denote the exempt and non-exempt biases, respectively. $M_{V,E}$, $M_{M,E}$, $M_{V,NE}$, and $M_{M,NE}$ denote the visible and masked exempt bias and visible and masked non-exempt bias, respectively. Because the SCM is known, all of these quantities can be computed. For each of the four experimental scenarios, the test accuracy is close to $99\%$ (model output is very similar to the true label). We observe that the bias decomposition for the model output $\hat{Y}$ is also quite similar to what one might intuitively expect for the true label $Y$. In Experimental Scenario 1, biased critical and general features are used in the true label. We also observe all four biases $M_{V,E}$, $M_{M,E}$, $M_{V,NE}$, and $M_{M,NE}$ are present in output $\hat{Y}$. In Experimental Scenarios 2 and 3, the bias in $\hat{Y}$ is dominated by non-exempt, masked bias $M_{M,NE}$, and the other components are negligible. In Experimental Scenario 4, the total bias is significantly less in comparison to the other three scenarios (intuitively agrees with the fact that the true labels that have no bias at all).

SCM is as follows: $X_1 = U_{X_1} + U_{X_3}$, $X_2 = U_{X_2}$, $X_3 = U_{X_3}$, and $X_4 = U_{X_2} - Z$, where $U_{X_1} \sim \text{Bern}(\frac{1}{2})$ denotes if the writing ability is above a threshold, and $U_{X_2}, U_{X_3} \sim i.i.d.\ \mathcal{N}(0,\sigma^2)$ denote interests and proofreading skill-level, respectively. Suppose that the historic true labels are given by $Y = \mathbb{1}((X_1 + X_4)^2 \geq 0.5)$, i.e., primarily high online-writing scores and high social-media-based-preference scores, but to appear "facially neutral" with respect to political inclination, the ad is also shown to candidates with low social-media-based-preference scores and low writing scores for whom the ad may be irrelevant. Here, the critical feature $X_c = X_1$ and the general features are $X_g = (X_2, X_3, X_4)$. The results are provided in Fig. 9b and Fig. 10b.

**Experimental Scenario 3 (Masking by general feature):** Consider another example similar to the previous one. The decision of showing ads for a website-manager's job in a newspaper company is based on three features, none of them critical: (i) $X_1$: a score based on online writing samples; (ii) $X_2$: a score based on browsing history, e.g., awareness of current events; and (iii) $X_3$: a preference score based on activity in social media, e.g., political alignment with the newspaper. The protected attribute $Z$ is political inclination, distributed Bern($\frac{1}{2}$). Suppose the true SCM is as follows: $X_1 = U_{X_1} + U'_{X_1}$, $X_2 = U_{X_2}$, and $X_3 = U_{X_2} - Z$, where $U_{X_1} \sim \text{Bern}(\frac{1}{2})$ denotes if writing ability is above a threshold, and $U'_{X_1}, U_{X_2} \sim i.i.d.\ \mathcal{N}(0,\sigma^2)$ denote proofreading skill and interests. Suppose that the true labels are given by $Y = \mathbb{1}((X_1 + X_3)^2 \geq 0.5)$, i.e., primarily high online-writing scores and high social-media-based-preference scores, but to appear "facially neutral" with respect to political inclination, the ad is also shown to candidates with low social-media-based-preference scores and low writing scores. Here, all the features are non-critical: $X_g = (X_1, X_2, X_3)$. The results are provided in Fig. 9c and Fig. 10c.

**Experimental Scenario 4 (No label bias):** The decision of showing ads for an editor's job is based on four features: (i) $X_1$: a score based on online writing samples (critical feature $X_c = X_1$); (ii) $X_2$: a score based on browsing history, e.g., awareness of current events; (iii) $X_3$: a preference score based on geographical proximity; and (iv) $X_4$: a score based on browsing history, e.g., interest in English websites as compared to websites of other languages. Let $Z \sim \text{Bern}(\frac{1}{2})$ be the protected attribute denoting whether the candidate is a native English speaker. Suppose the true SCM is as follows: $X_1 = Z + U_{X_1} + U'_{X_1}$, $X_2 = U_{X_2}$, $X_3 = U_{X_3}$, and $X_4 = Z + U_{X_2}$, where $U_{X_1} \sim \text{Bern}(\frac{1}{2})$ denotes whether writing skill is above a threshold, and $U'_{X_1}, U_{X_2}, U_{X_3} \sim i.i.d.\ \mathcal{N}(0,\sigma^2)$ denote proofreading skill, interests, and proximity. Suppose that the true labels do not have label bias and are given by $Y = \mathbb{1}(U_{X_1} + U_{X_2} \geq 0.5)$. Here, the critical feature is $X_c = X_1$ and the general features are $X_g = (X_2, X_3, X_4)$. The results are provided in Fig. 9d and Fig. 10d.
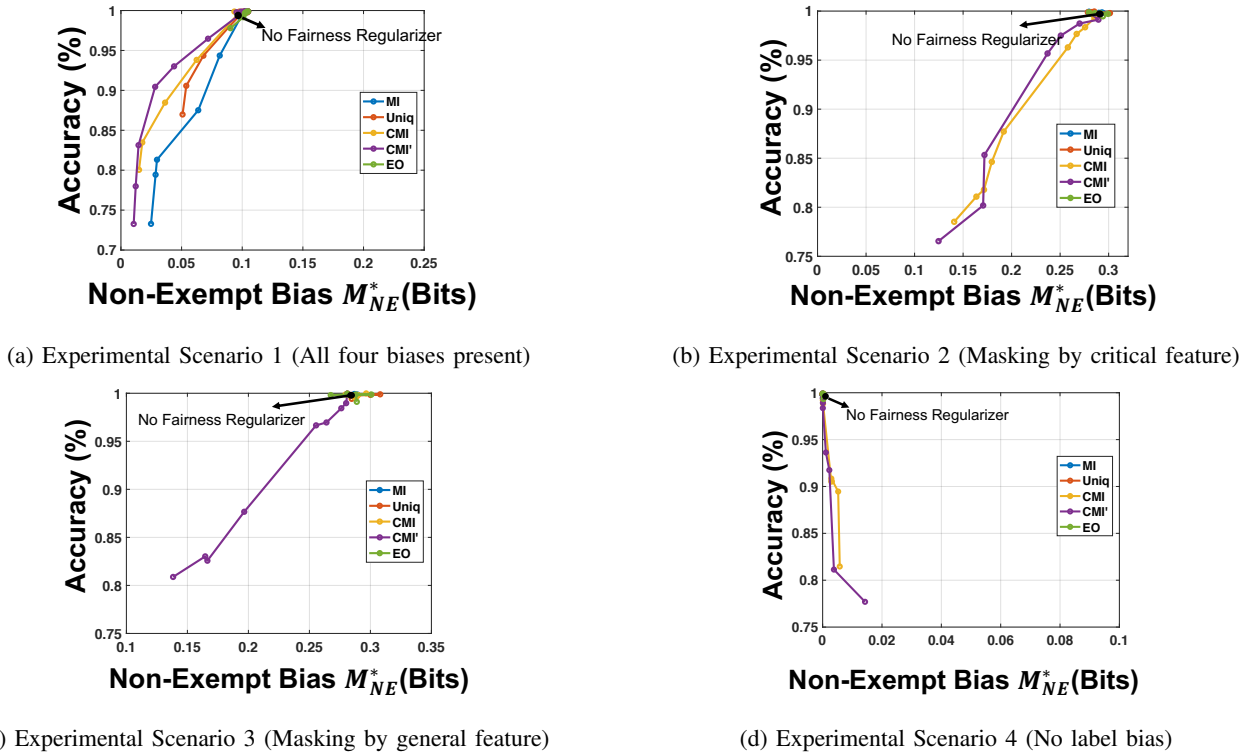
(a) Experimental Scenario 1 (All four biases present)

(b) Experimental Scenario 2 (Masking by critical feature)

(c) Experimental Scenario 3 (Masking by general feature)

(d) Experimental Scenario 4 (No label bias)

Fig. 10: **Observations from training:** For each experimental scenario, we train a model using each of the five observational regularizers: MI (Statistical Parity), Uniq, CMI, CMI', and EO (Equalized Odds) for different values of regularization constant $\lambda$. The tradeoff between test accuracy and the actual non-exempt bias ($M_{NE}*$) computed using the `dit` package is shown. In Experimental Scenario 1, the model output (no fairness) has all four types of biases, $M_{V,E}$, $M_{M,E}$, $M_{V,NE}$, and $M_{M,NE}$. We observe that, all three of Uniq, CMI, and CMI' attain better tradeoff between accuracy and non-exempt bias as compared to EO (Equalized Odds) and MI (Statistical Parity). Equalized Odds does not affect the accuracy or the non-exempt bias much, even for high values of the regularization constant. Statistical Parity attempts to reduce both exempt and non-exempt biases, and ends up reducing accuracy a lot for same values of non-exempt bias as compared to Uniq, CMI, and CMI'. CMI and CMI' are slightly better than Uniq because they also partially quantify non-exempt, masked bias. For CMI'($= I(Z; \hat{Y}|X_c, X')$), we choose $X' = X_3$ (location, a general feature that has no causal influence of $Z$, but is suspected to "mask" $Z$ in the final output) which leads to a better trade-off than CMI. In Experimental Scenarios 2 and 3, the bias in the model output (no fairness) is dominated by non-exempt, masked bias. This bias is missed by MI (Statistical Parity), EO (Equalized Odds), and Uniq. Consequently, they do not affect the accuracy or the non-exempt bias much, even for high values of regularization constant. For Experimental Scenario 2, only CMI and CMI' (with $X' = X_3$) are able to detect the non-exempt, masked bias, and lead to alternate models with reduced accuracy and also reduced non-exempt bias. For Experimental Scenario 3, only CMI' (with $X' = X_1$, the general feature that masks $Z$ in the final output) detects the non-exempt bias, and reduces it. In Experimental Scenario 4, the model output (no fairness) has almost negligible non-exempt bias because the true labels do not have any bias at all. We observe that, MI, EO, and Uniq also do not affect the accuracy much even for high values of regularization constant (which is desirable). However, CMI, and CMI' (with $X' = X_2$) falsely detect bias here, when there is no non-exempt bias actually present. In an attempt to reduce the falsely detected bias, they lead to alternate models with significantly reduced accuracy, and slightly increased non-exempt bias.

**Summary of Results:** We present results for auditing and training in Fig. 9 and Fig. 10 with detailed explanations. Our proposed regularizers, namely, Uniq, CMI and CMI' attain better trade-off between accuracy and non-exempt bias than MI (Statistical Parity) and EO (Equalized Odds) in Experimental Scenario 1. CMI and CMI' are also able to detect certain scenarios of non-exempt, masked bias that Uniq, MI and EO fail to detect, e.g., in Experimental Scenario 2 where the masking is by the critical feature $X_c$. Experimental Scenario 3 demonstrates additional scenarios of non-exempt, masked bias, e.g., masking by $X_g$, where even CMI is unable to detect this bias, and only CMI' succeeds (by choosing $X'$ based on certain knowledge/suspicion of the causal model). However, Experimental Scenario 4 denotes a scenario of false detection of bias by CMI and CMI'. In essence, Uniq is a somewhat conservative measure of non-exempt bias which can miss non-exempt, masked bias, but never does false detection of bias. On the other hand, CMI and CMI' can sometimes detect certain scenarios of non-exempt, masked bias, but can also sometimes falsely detect bias. This is expected: these are observational measures attempting to approximate a causal measure, a fundamentally impossible task. However, these examples illustrate how knowledge of aspects of the SCM (e.g., whether the bias is predominantly masked bias) can be used to inform the choice of the observational measure.

*B. Case Study on Real Data:* `Adult` *Dataset*

The Adult dataset [64], also known as the Census income dataset, consists of 14 features (e.g. age, educational qualification), and the true labels denote whether the income is greater than $50k. This dataset is widely used in existing fairness literature (e.g., [22]), because it is representative of data used in highly consequential applications, such as, lending, showing expensive ads, etc. Here, we choose gender as the protected attribute ($Z$) for analyzing the Adult dataset. Our set of input features ($X$) consists of all the other features except gender, and our critical feature ($X_c$) is working-hours per-week.

We train a deep neural network (multi-layer perceptron) on this dataset, with all features, except gender, as input (with one hot encoding of all categorical variables). The input layer is followed by three hidden layers, each having 32 neurons with ReLu activation and dropout probability 0.2. Finally, the output layer consists of a single neuron with sigmoidal activation that produces an output value between 0 and 1 (likelihood of income being $> 50k$ possibly leading to a loan decision).

Since the true causal model is not known, we cannot compute the exact value of the total bias or non-exempt bias ($M_{NE}^*$) as in the previous case study. However, our observational measures can still provide valuable insights as we demonstrate here (see Fig. 11). We consider five setups for auditing: (i) No fairness: model trained with no fairness regularizer; (ii) Statistical Parity: model trained with $I(Z; \hat{Y})$ as regularizer; (iii) CMI Regularizer: model trained with $I(Z; \hat{Y}|X_c)$ regularizer; (iv) Uniq Regularizer: model trained with $\mathrm{Uni}(Z : \hat{Y}|X_c)$ as regularizer; and (iv) Equalized Odds: model trained with $I(Z; \hat{Y}|Y)$ regularizer. For each of these setups, we choose the same value of the regularization constant $\lambda = 4$, and similar correlation-based estimates for the regularizers as in the previous case study.

After training these models, we audit/evaluate the trained models by computing the following observational quantities on the empirical distribution of the test data using the `dit` [43] package: MI (statistically visible bias:$I(Z; \hat{Y})$), CMI (conditional mutual information $I(Z; \hat{Y}|X_c)$), as well as, the decomposition of CMI into Unique Information (Uniq) given by $\mathrm{Uni}(Z : \hat{Y}|X_c)$) and Synergistic Information (Syn) given by $\mathrm{Syn}(Z : (\hat{Y}, X_c))$). Recall that Uniq is the non-exempt statistically visible bias, while Syn can correspond to either non-exempt masked bias or false detection of bias (recall our impossibility result; one might need some knowledge of the causal model to be certain). As discussed in the caption of Fig. 11, the correlation-based estimates serve as relatively good approximations and reduce the respective statistical dependences as one would intuitively expect to see.
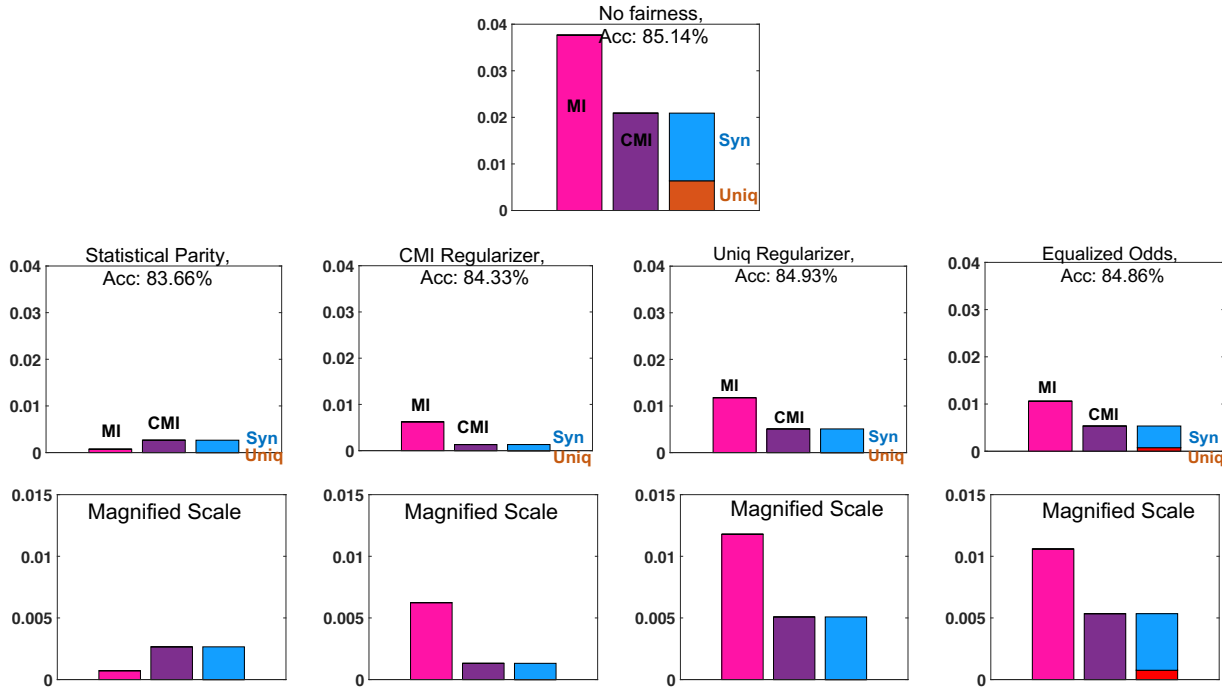


Fig. 11: For the model with no fairness, we see a high value of MI as well as CMI (consisting of both Uniq and Syn). When the model is trained for statistical parity, the MI reduces as expected, but interestingly CMI is now higher than MI. Next, when CMI is used as a regularizer, we notice that CMI (and its sub-components Uniq and Syn) reduce as expected, but MI is higher than CMI. For Uniq as a regularizer, we notice that MI or CMI are not reduced that much, but only Uniq is minimized selectively. Lastly, for equalized odds, we observe that the trained model still has some Uniq (non-exempt, visible bias). These experiments also demonstrate that the correlation-based estimates for the regularizers are relatively good approximations for this real dataset and actually reduce the respective statistical dependences as one would intuitively expect.

### C. Case Study on Real Data: `German Credit` Dataset

We also perform a similar case study on the German Credit Dataset [64]. This dataset consists of 20 features (e.g., status of a checking account, credit amount, present employment, etc.), and the true labels denote whether a customer is good or bad. Our critical feature ($X_c$) is the number of existing credits at this bank, and the protected attribute ($Z$) is gender. Our set of all features ($X$) consist of all features except gender and marital status.

We train a deep neural network (multi-layer perceptron) on this dataset, with all features, except gender and marital status as input (with one hot encoding of categorical variables). The input layer is followed by two hidden layers, each having 124 neurons with ReLu activation and dropout probability 0.5. Finally, the output layer consists of a single neuron with sigmoidal activation that produces an output value between 0 and 1 (likelihood of being a good customer).

The causal model is again not known, similar to the previous case. However, similar to the case study on the Adult dataset, we train the model using different observational regularizers, and audit/evaluate the trained models. As discussed in the caption of Fig. 12, the correlation-based estimates reduce the respective statistical dependences as one would intuitively expect to see.
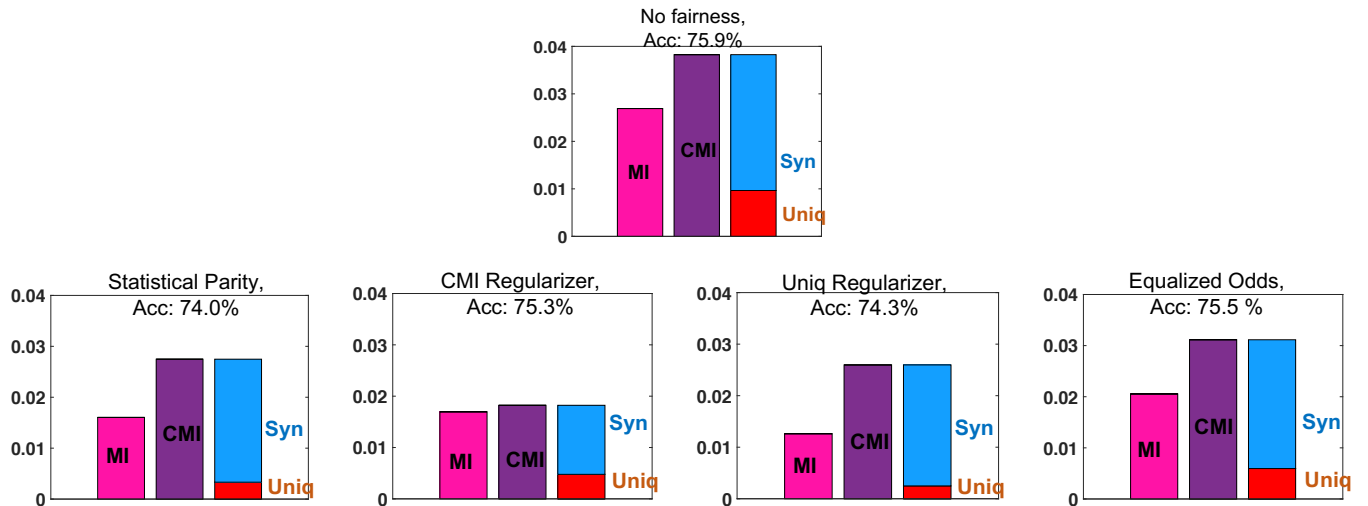


Fig. 12: The experimental results demonstrate that the correlation-based estimates for the regularizers behave as expected. When the model is trained for statistical parity, MI reduces as expected without significantly reducing CMI. Next, when CMI is used as a regularizer, we notice that CMI (and its sub-components Uniq and Syn) reduce as expected. For Uniq as a regularizer, we notice that MI or CMI are not reduced that much, but only Uniq is minimized selectively. Lastly, for equalized odds, we observe that the trained model still retains quite a bit of MI, CMI, Uniq and Syn, as compared to the model with no fairness.

## VIII. DISCUSSION AND CONCLUSION

**On Choice of Critical Features and Connections with Explainability:** In this work, as also in some existing works on fairness [4], [17], we assume that the critical features are known. We adopt a viewpoint stated in [65] which suggests that "We can't just rely on the math; we still need a human person applying human judgements." Since most of these exemptions are embedded in law and social science [30]–[32], we believe that fairness researchers need to collaborate with social scientists and lawyers in order to determine which set of features can be designated as critical for a particular application.

This work also shares close connections with the field of *explainability* in machine learning [15], [51], [66], and motivates several related research problems, e.g., how to check or explain if certain features contributed to the bias in a model, or how to incorporate exemptions in applications, such as, image processing, where certain neurons in an intermediate hidden layer might need to be exempted instead of the input layer because they often have more interpretability [66].

**On Better Understanding of Observational Measures:** Our proposed counterfactual measure and the desirable properties help in evaluation of observational measures in practice, and understand their utility and limitation, i.e., what they capture and miss. Finally, in applications where when the true SCM is known or can be evaluated from the data [35, Chapters 4,7], the proposed measure exactly captures the non-exempt bias.

**On Uniqueness, Operational Meaning and Further Generalizations:** We acknowledge that we do not prove uniqueness of our measure with respect to the desirable properties, and neither do we show that the properties are exhaustive (recall Remark 5 in Section III-B). This is an interesting direction of future work. However, there may also be value in the fact that the properties do not yield a unique measure: this allows for tuning the measure based on the application. E.g., Shannon established uniqueness on entropy with respect to **some** properties in [58] but subsequent applications have still led to the use of modified measures, e.g. Renyi entropy [21], [55], [59], [60].

Deriving the exact operational meaning of our proposed counterfactual measure is also an interesting direction of future work. Nonetheless, the proposed measure does satisfy our stated desirable properties and capture important aspects of the problem, e.g., statistically visible and masked biases. Furthermore, our measure can also be modified to account for further functional generalizations. First notice, that our proposed Property 3 is a special case of the following statement:

*If $(Z, f_a(U_X)) - X_c - (\hat{Y}, f_b(U_X))$ form a Markov chain for any deterministic functions $f_a(\cdot)$ and $f_b(\cdot)$ such that $f_a(U_X) \perp\!\!\!\perp f_b(U_X)$ and $\mathrm{H}(U_X) = \mathrm{H}(f_a(U_X)) + \mathrm{H}(f_b(U_X))$, then $M_{NE} = 0$.*

To account for this more general property, our proposed measure might be modified as follows:

$$\min_{f_a(U_X), f_b(U_X)} \mathrm{Uni}((Z, f_a(U_X)) : (\hat{Y}, f_b(U_X)) | X_c), \tag{12}$$

such that $f_a(U_X) \perp\!\!\!\perp f_b(U_X)$ and $\mathrm{H}(U_X) = \mathrm{H}(f_a(U_X)) + \mathrm{H}(f_b(U_X))$. This measure also satisfies all the other desirable properties. In this work, we restrict ourselves to $f_a(U_X)$ and $f_b(U_X)$ being disjoint subsets of $U_X$ for simplicity, computability and ease of understanding. Future work will explore how different assumptions on the SCM restrict the class of $f_a$ and $f_b$.

**On Understanding Other Forms of Masked Bias:** Let us revisit the discussion from Section III-B that not all forms of masked discrimination are necessarily undesirable. E.g., if $U_{X_1}$ is a random coin flip in Canonical Example 5, then performing $\hat{Y} = Z \oplus U_{X_1}$ randomizes the race, and can even be regarded as a preventive measure against discrimination. However, keeping the mathematics of the example same, if $U_{X_1}$ instead denotes whether one's income is above a threshold, then the model is unfair. It is an interesting future direction to examine how to quantify non-exempt discrimination while allowing the user with more flexibility on what latent factors are allowed to mask $Z$.

**On Estimation of Mutual Information, Conditional Mutual Information and Unique Information:** In general, it is difficult to directly incorporate these information-theoretic measures as a regularizer with the loss function (see [67], [68] and the references therein). Examining alternate methods of incorporating our proposed measures as regularizer (using or building upon techniques proposed in [21], [27], [42], [55], [68]) is an interesting direction of future work.

## ACKNOWLEDGMENT

## REFERENCES

[1] S. Dutta, P. Venkatesh, P. Mardziel, A. Datta, and P. Grover, "An information-theoretic quantification of discrimination with exempt features," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.

[2] I. Žliobaite, F. Kamiran, and T. Calders, "Handling conditional discrimination," in *2011 IEEE 11th International Conference on Data Mining*. IEEE, 2011, pp. 992–1001.

[3] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, "Fairness through awareness," in *Proceedings of the 3rd innovations in theoretical computer science conference*. ACM, 2012, pp. 214–226.

[4] F. Kamiran, I. Žliobaitė, and T. Calders, "Quantifying explainable discrimination and removing illegal discrimination in automated decision making," *Knowledge and information systems*, vol. 35, no. 3, pp. 613–644, 2013.

[5] T. Calders, A. Karim, F. Kamiran, W. Ali, and X. Zhang, "Controlling attribute effect in linear regression," in *2013 IEEE 13th international conference on data mining*. IEEE, 2013, pp. 71–80.

[6] A. Agarwal, A. Beygelzimer, M. Dudík, J. Langford, and H. Wallach, "A reductions approach to fair classification," *arXiv preprint arXiv:1803.02453*, 2018.

[7] M. Hardt, E. Price, N. Srebro *et al.*, "Equality of opportunity in supervised learning," in *Advances in neural information processing systems*, 2016, pp. 3315–3323.

[8] F. Calmon, D. Wei, B. Vinzamuri, K. N. Ramamurthy, and K. R. Varshney, "Optimized pre-processing for discrimination prevention," in *Advances in Neural Information Processing Systems*, 2017, pp. 3992–4001.

[9] A. K. Menon and R. C. Williamson, "The cost of fairness in binary classification," in *Conference on Fairness, Accountability and Transparency*, 2018, pp. 107–118.

[10] T. Kamishima, S. Akaho, H. Asoh, and J. Sakuma, "Fairness-aware classifier with prejudice remover regularizer," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2012, pp. 35–50.

[11] M. Donini, L. Oneto, S. Ben-David, J. S. Shawe-Taylor, and M. Pontil, "Empirical risk minimization under fairness constraints," in *Advances in Neural Information Processing Systems*, 2018, pp. 2791–2801.

[12] A. Ghassami, S. Khodadadian, and N. Kiyavash, "Fairness in supervised learning: An information theoretic approach," in *2018 IEEE International Symposium on Information Theory (ISIT)*, 2018, pp. 176–180.

[13] M. B. Zafar, I. Valera, M. G. Rodriguez, and K. P. Gummadi, "Fairness constraints: Mechanisms for fair classification," *arXiv preprint arXiv:1507.05259*, 2015.

[14] Z. Lipton, J. McAuley, and A. Chouldechova, "Does mitigating ML's impact disparity require treatment disparity?" in *Advances in Neural Information Processing Systems*, 2018, pp. 8125–8135.

[15] K. R. Varshney, "Trustworthy machine learning and artificial intelligence," *XRDS: Crossroads, The ACM Magazine for Students*, vol. 25, no. 3, pp. 26–29, 2019.

[16] M. J. Kusner, J. Loftus, C. Russell, and R. Silva, "Counterfactual fairness," in *Advances in Neural Information Processing Systems*, 2017, pp. 4066–4076.

[17] N. Kilbertus, M. R. Carulla, G. Parascandolo, M. Hardt, D. Janzing, and B. Schölkopf, "Avoiding discrimination through causal reasoning," in *Advances in Neural Information Processing Systems*, 2017, pp. 656–666.

[18] C. Russell, M. J. Kusner, J. Loftus, and R. Silva, "When worlds collide: Integrating different counterfactual assumptions in fairness," in *Advances in Neural Information Processing Systems*, 2017, pp. 6414–6423.

[19] S. Chiappa, "Path-specific counterfactual fairness," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, 2019, pp. 7801–7808.

[20] A. Datta, M. Fredrikson, G. Ko, P. Mardziel, and S. 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*. ACM, 2017, pp. 1193–1210.

[21] J. Liao, C. Huang, P. Kairouz, and L. Sankar, "Learning generative adversarial representations (gap) under fairness and censoring constraints," *arXiv preprint arXiv:1910.00411*, 2019.

[22] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, "Learning fair representations," in *International Conference on Machine Learning*, 2013, pp. 325–333.

[23] S. Yeom, A. Datta, and M. Fredrikson, "Hunting for discriminatory proxies in linear regression models," in *Advances in Neural Information Processing Systems*, 2018, pp. 4568–4578.

[24] T. Speicher, H. Heidari, N. Grgic-Hlaca, K. P. Gummadi, A. Singla, A. Weller, and M. B. Zafar, "A unified approach to quantifying algorithmic unfairness: Measuring individual & group unfairness via inequality indices," in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 2018, pp. 2239–2248.

[25] H. Wang, H. Hsu, M. Diaz, and F. P. Calmon, "To split or not to split: The impact of disparate treatment in classification," *arXiv preprint arXiv:2002.04788*, 2020.

[26] M. Kearns, S. Neel, A. Roth, and Z. S. Wu, "Preventing fairness gerrymandering: Auditing and learning for subgroup fairness," in *International Conference on Machine Learning*, 2018, pp. 2564–2572.

[27] J. Cho, G. Hwang, and C. Suh, "A fair classifier using mutual information," 2020.

[28] M. J. Kusner, C. Russell, J. R. Loftus, and R. Silva, "Causal Interventions for Fairness," *arXiv preprint arXiv:1806.02380*, 2018.

[29] J. H. Hinnefeld, P. Cooman, N. Mammo, and R. Deese, "Evaluating fairness metrics in the presence of dataset bias," *arXiv preprint arXiv:1809.09245*, 2018.

[30] S. S. Grover, "The business necessity defense in disparate impact discrimination cases," *Ga. L. Rev.*, vol. 30, p. 387, 1995.

[31] S. Barocas and A. D. Selbst, "Big data's disparate impact," *Calif. L. Rev.*, vol. 104, p. 671, 2016.

[32] EEOC Website, "US Equal Pay Act," https://www.eeoc.gov/laws/statutes/epa.cfm.

[33] J. H. Hinnefeld, "Measuring model fairness," https://www.youtube.com/watch?v=V3tmxMf2UH8.

[34] S. Yeom and M. C. Tschantz, "Discriminative but not discriminatory: A comparison of fairness definitions under different worldviews," *arXiv preprint arXiv:1808.08619*, 2018.

[35] J. Peters, D. Janzing, and B. Schölkopf, *Elements of Causal Inference: Foundations and Learning Algorithms*. MIT press, 2017.

[36] J. Zhang and E. Bareinboim, "Fairness in decision-making—the causal explanation formula," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[37] R. Nabi and I. Shpitser, "Fair inference on outcomes," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.

[38] T. Tax, P. Mediano, and M. Shanahan, "The partial information decomposition of generative neural network models," *Entropy*, vol. 19, no. 9, p. 474, 2017.

[39] P. Venkatesh, S. Dutta, and P. Grover, "Information flow in computational systems," *IEEE Transactions on Information Theory (To Appear)*, 2020.

[40] S. Corbett-Davies, E. Pierson, A. Feller, S. Goel, and A. Huq, "Algorithmic decision making and the cost of fairness," in *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '17. ACM, 2017, pp. 797–806. [Online]. Available: http://doi.acm.org/10.1145/3097983.3098095

[41] B. Salimi, L. Rodriguez, B. Howe, and D. Suciu, "Interventional Fairness: Causal Database Repair for Algorithmic Fairness," in *Proceedings of the 2019 International Conference on Management of Data*, ser. SIGMOD '19. ACM, 2019, pp. 793–810.

[42] Anonymous, "Conditional debiasing for neural networks."

[43] R. G. James, C. J. Ellison, and J. P. Crutchfield, "dit: a Python package for discrete information theory," *The Journal of Open Source Software*, vol. 3, no. 25, p. 738, 2018.

[44] A. B. Barrett, "Exploration of synergistic and redundant information sharing in static and dynamical Gaussian systems," *Physical Review E*, vol. 91, no. 5, p. 052802, 2015.

[45] E. W. Weisstein, "n-tuple," https://mathworld.wolfram.com/n-Tuple.html, from MathWorld–A Wolfram Web Resource.

[46] N. Bertschinger, J. Rauh, E. Olbrich, J. Jost, and N. Ay, "Quantifying unique information," *Entropy*, vol. 16, no. 4, pp. 2161–2183, 2014.

[47] P. L. Williams and R. D. Beer, "Nonnegative decomposition of multivariate information," *arXiv preprint arXiv:1004.2515*, 2010.

[48] V. Griffith and C. Koch, "Quantifying synergistic mutual information," in *Guided Self-Organization: Inception*. Springer, 2014, pp. 159–190.

[49] P. K. Banerjee, J. Rauh, and G. Montúfar, "Computing the unique information," in *2018 IEEE International Symposium on Information Theory (ISIT)*, 2018, pp. 141–145.

[50] L. Breiman, "Random forests," *Machine learning*, vol. 45, no. 1, pp. 5–32, 2001.

[51] A. Datta, S. Sen, and Y. Zick, "Algorithmic transparency via quantitative input influence: Theory and experiments with learning systems," in *2016 IEEE Symposium on Security and Privacy (SP)*, 2016, pp. 598–617.

[52] P. W. Koh and P. Liang, "Understanding black-box predictions via influence functions," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 2017, pp. 1885–1894.

[53] P. Adler, C. Falk, S. A. Friedler, T. Nix, G. Rybeck, C. Scheidegger, B. Smith, and S. Venkatasubramanian, "Auditing black-box models for indirect influence," *Knowledge and Information Systems*, vol. 54, no. 1, pp. 95–122, 2018.

[54] A. Henelius, K. Puolamäki, H. Boström, L. Asker, and P. Papapetrou, "A peek into the black box: exploring classifiers by randomization," *Data mining and knowledge discovery*, vol. 28, no. 5-6, pp. 1503–1529, 2014.

[55] J. Liao, L. Sankar, O. Kosut, and F. P. Calmon, "Robustness of maximal $\alpha$-leakage to side information," *arXiv preprint arXiv:1901.07105*, 2019.

[56] S. Dutta, D. Wei, H. Yueksel, P.-Y. Chen, S. Liu, and K. R. Varshney, "An information-theoretic perspective on the relationship between fairness and accuracy," *arXiv preprint arXiv:1910.07870*, 2019.

[57] P. K. Banerjee, E. Olbrich, J. Jost, and J. Rauh, "Unique informations and deficiencies," in *2018 56th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, 2018, pp. 32–38.

[58] C. E. Shannon, "A mathematical theory of communication," *Bell system technical journal*, vol. 27, no. 3, pp. 379–423, 1948.

[59] A. Rényi, "On measures of entropy and information," in *Proceedings of the 4th Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1. The Regents of the University of California, 1961.

[60] I. Issa, A. B. Wagner, and S. Kamath, "An operational approach to information leakage," *IEEE Transactions on Information Theory*, 2019.

[61] I. Mironov, "Rényi differential privacy," in *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. IEEE, 2017, pp. 263–275.

[62] N. Kilbertus, P. Ball, M. Kusner, A. Weller, and R. Silva, "The sensitivity of counterfactual fairness to unmeasured confounding," in *35th Conference on Uncertainty in Artificial Intelligence, UAI 2019*. Association for Uncertainty in Artificial Intelligence, 2019.

[63] T. M. Cover and J. A. Thomas, *Elements of Information Theory*. John Wiley & Sons, 2012.

[64] D. Dua and C. Graff, "UCI machine learning repository," 2017. [Online]. Available: http://archive.ics.uci.edu/ml

[65] P. Pandey, "Is your machine learning model biased?" https://towardsdatascience.com/is-your-machine-learning-model-biased-94f9ee176b67.

[66] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas *et al.*, "Interpretability Beyond Feature Attribution: Quantitative Testing with Concept Activation Vectors (TCAV)," in *International Conference on Machine Learning*, 2018, pp. 2668–2677.

[67] D. Pál, B. Póczos, and C. Szepesvári, "Estimation of Rényi entropy and mutual information based on generalized nearest-neighbor graphs," in *Advances in Neural Information Processing Systems*, 2010, pp. 1849–1857.

[68] S. Mukherjee, H. Asnani, and S. Kannan, "CCMI: Classifier based Conditional Mutual Information estimation," *arXiv preprint arXiv:1906.01824*, 2019.

[69] J. Rauh, P. Kr. Banerjee, E. Olbrich, and J. Jost, "Unique Information and Secret Key Decompositions," in *2019 IEEE International Symposium on Information Theory (ISIT)*, 2019, pp. 3042–3046.

## APPENDIX A
## COUNTERFACTUAL CAUSAL INFLUENCE (CCI) AND ITS CONNECTION TO COUNTERFACTUAL FAIRNESS

### A. Proof of Lemma 1

Here, we first provide a proof of Lemma 1 which shows that our proposed quantification of total bias is zero if and only if $\text{CCI}(Z \to \hat{Y}) = 0$. For ease of reading, we repeat the statement of the lemma here again.

**Lemma 1** (Equivalences of CCI). *Consider the aforementioned system model. Let $\hat{Y} = h(Z, U_X)$ for some deterministic function $h(\cdot)$ and $Z \perp\!\!\!\perp U_X$. Then, $\text{CCI}(Z \to \hat{Y}) = 0$ if and only if $\text{I}(Z; (\hat{Y}, U_X)) = 0$.*

*Proof of Lemma 1.* From the definition of CCI (Definition 3 in Section II-B),

$$\text{CCI}(Z \to \hat{Y}) = \mathbb{E}_{Z,Z',U_X} \left[ |h(Z, U_X) - h(Z', U_X)| \right]$$
$$= \sum_{z_1, z_2, u_x} \Pr(Z = z_1, Z' = z_2, U_X = u_x) |h(z_1, u_x) - h(z_2, u_x)|$$
$$= \sum_{z_1, z_2, u_x} \Pr(Z = z_1) \Pr(Z' = z_2) \Pr(U_X = u_x) |h(z_1, u_x) - h(z_2, u_x)| \quad \text{[from independence].} \quad (13)$$

The summation consist of non-negative terms. Therefore, $\text{CCI}(Z \to \hat{Y}) = 0$, *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 equivalent to $h(z, u_x)$ being constant over all possible values of $z$ with $\Pr(Z = z) > 0$ given a fixed value of $u_x$, and this should happen over all values of $u_x$ with $\Pr(U_X = u_x)$.

Now, observe that,

$$\text{I}(Z; (\hat{Y}, U_X)) = \text{I}(Z; \hat{Y} \mid U_X) + \text{I}(Z; U_X) \quad (14)$$
$$= \text{I}(Z; \hat{Y} \mid U_X) \quad [Z \perp\!\!\!\perp U_X] \quad (15)$$
$$= \text{H}(\hat{Y} \mid U_X) - \text{H}(\hat{Y} \mid U_X, Z) \quad [\text{By Definition}] \quad (16)$$
$$= \text{H}(\hat{Y} \mid U_X). \quad [\hat{Y} \text{ is completely determined by } Z \text{ and } U_X] \quad (17)$$

$\text{H}(\hat{Y} \mid U_X)$ can be 0 *if and only if* $h(z, u_x)$ is constant over all possible values of $z$ with $\Pr(Z = z) > 0$ given a fixed value of $u_x$, and this should happen over all $u_x$ with $\Pr(U_X = u_x) > 0$. Thus, $\text{CCI}(Z \to \hat{Y}) = 0$ if and only if $\text{I}(Z; (\hat{Y}, U_X)) = 0$. □

### B. Connections to Counterfactual Fairness

We note that the concept of counterfactual causal influence (often referred to as only "influence") is derived from a separate body of work [50]–[54]) outside the fairness literature. The original definition of counterfactual fairness in [16] was stated differently (without using CCI), although the connection with CCI has been hinted at in [18]. Here, for the sake of completeness, we will formally show in Lemma 6 that $\text{CCI}(Z \to \hat{Y}) = 0$ is equivalent to the counterfactual fairness criterion proposed in [16]. What this means is that, our proposed quantification of total bias is also 0 if and only if a model is counterfactually fair.

First, we clarify the differences in notation between our work and [16]. In our work, $X = f(Z, U_X)$ and $\hat{Y} = r(X) = r \circ f(Z, U_X) = h(Z, U_X)$ where $h = r \circ f$. In [16], $\hat{Y}_{Z \leftarrow z_1}(U)$ denotes the random variable $\hat{Y}$ when the value of $Z$ is fixed as $z_1$ by an intervention, i.e., $\hat{Y}_{Z \leftarrow z_1}(U) = h(z_1, U_X)$. Alongside, we also clarify that the event that $X$ takes the value $x$ when $Z$ is fixed as $z_1$ refers to the event that $U_X$ takes a value from the set $\mathcal{S}(x, z_1) = \{u_x : x = f(z_1, u_x), \Pr(U_X = u_x) > 0\}$ because $X = f(Z, U_X)$.

**Definition 8** (Counterfactual Fairness given $X = x$ and $Z = z_1$ [16]). *A predictor $\hat{Y}$ is counterfactually fair given the protected attribute $Z = z_1$ and the observed variable $X = x$, if we have,*

$$\Pr(\hat{Y}_{Z \leftarrow z_1}(U) = y \mid X \text{ takes the value } x \text{ when } Z \text{ is fixed as } z_1)$$
$$= \Pr(\hat{Y}_{Z \leftarrow z_2}(U) = y \mid X \text{ takes the value } x \text{ when } Z \text{ is fixed as } z_1), \quad (18)$$

*for all attainable $y$ and $z_2$. In our notations, this definition is equivalent to the following: Given the sensitive attribute $Z = z_1$ and the observed variable $X = x$,*

$$\Pr(h(z_1, U_X) = y \mid U_X \in \mathcal{S}(x, z_1)) = \Pr(h(z_2, U_X) = y \mid U_X \in \mathcal{S}(x, z_1)), \quad (19)$$

*for all attainable $y$ and $z_2$, where $\mathcal{S}(x, z_1) = \{u_x : \ x = f(z_1, u_x), \ \Pr(U_X = u_x) > 0\}$.*

Next, we show that $\text{CCI}(Z \to \hat{Y}) = 0$ is equivalent to the counterfactual fairness criterion of [16].

**Lemma 6.** $\text{CCI}(Z \to \hat{Y}) = 0$ *is equivalent to counterfactual fairness (Definition 8) for all $X = x$ and $Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$.*

*Proof of Lemma 6.* Suppose that, $\text{CCI}(Z \to \hat{Y}) = 0$. Recall from Lemma 1, that $\text{CCI}(Z \to \hat{Y}) = 0$ is equivalent to the criterion that $h(z_1, u_x) = h(z_2, u_x)$ for all attainable $z_1$, $z_2$ given a particular value of $u_x$, and this should hold for all $u_x$ with $\Pr(U_X = u_x) > 0$. Therefore, for any particular $X = x$ and $Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$,

$$\Pr(h(z_1, U_X) = y \mid \ U_X \in \mathcal{S}(x, z_1)) = \Pr(h(z_2, U_X)) = y \mid \ U_X \in \mathcal{S}(x, z_1)), \tag{20}$$

because $h(z_1, u_x) = h(z_2, u_x)$ for all $u_x \in \mathcal{S}(x, z_1)$. Thus, we show that $\text{CCI}(Z \to \hat{Y}) = 0$ implies counterfactual fairness.

Now, we prove the implication in the other direction. Suppose that the counterfactual fairness criterion (19) holds for all $X = x$ and $Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$.

First consider any particular $X = x$ and $Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$. Since $\Pr(X = x, Z = z_1) > 0$, there exists at least one $u_x$ with $\Pr(U_X = u_x) > 0$ such that $x = f(z_1, u_x)$. So, the set $\mathcal{S}(x, z_1)$ is non-empty. Equation (19) implies that,

$$\Pr(h(z_1, U_X) = y \mid \ U_X \in \mathcal{S}(x, z_1)) = \Pr(h(z_2, U_X)) = y \mid \ U_X \in \mathcal{S}(x, z_1)) \ \forall \text{ attainable } y, z_2. \tag{21}$$

This leads to,

$$\Pr(h(z_1, U_X) = y, \ U_X \in \mathcal{S}(x, z_1)) = \Pr(h(z_2, U_X) = y, \ U_X \in \mathcal{S}(x, z_1)) \ \forall \text{ attainable } y, z_2. \tag{22}$$

Or,

$$\sum_{u_x \in \mathcal{S}(x, z_1)} \Pr(U_X = u_x) \mathbb{1}(h(z_1, u_x) = y) = \sum_{u_x \in \mathcal{S}(x, z_1)} \Pr(U_X = u_x) \mathbb{1}(h(z_2, u_x) = y) \ \forall \text{ attainable } y, z_2. \tag{23}$$

Now, observe that, $f(z_1, u_x) = x$ for all $u_x \in \mathcal{S}(x, z_1)$, and thus $h(z_1, u_x) = r \circ f(z_1, u_x)$ takes the same value for all $u_x \in \mathcal{S}(x, z_1)$. Let $h(z_1, u_x) = \tilde{y}$ for all $u_x \in \mathcal{S}(x, z_1)$. Then, for (23) to hold, we need,

$$\sum_{u_x \in \mathcal{S}(x, z_1)} \Pr(U_X = u_x)(1 - \mathbb{1}(h(z_2, u_x) = \tilde{y})) = 0 \ \forall \text{ attainable } z_2.$$

This holds if and only if $\mathbb{1}(h(z_2, u_x) = \tilde{y}) = 1$ for all $u_x \in \mathcal{S}(x, z_1)$ and for all attainable $z_2$. Thus, the counterfactual fairness criterion (19) *for a particular* $X = x, Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$ implies that for all $u_x \in \mathcal{S}(x, z_1)$,

$$h(z_2, u_x) = h(z_1, u_x) \ \ \forall \text{ attainable } z_2. \tag{24}$$

Because the counterfactual criterion (19) holds for all $X = x, Z = z_1$ with $\Pr(X = x, Z = z_1) > 0$, we therefore have (24) hold for all

$$u_x \in \cup_{\{x, z_1 : \Pr(X = x, Z = z_1) > 0\}} \mathcal{S}(x, z_1).$$

Now, because $U_X$ is independent of $Z$, for any $u_x^*$ with $\Pr(U_X = u_x^*) > 0$, there always exists some $x^*$ such that $x^* = f(z_1, u_x^*)$, and $\Pr(X = x^*, Z = z_1) \geq \Pr(U_X = u_x^*, Z = z_1) > 0$. Thus, $u_x^* \in S(x^*, z_1)$ for some $(x^*, z_1)$ with $\Pr(X = x^*, Z = z_1) > 0$. Thus,

$$\{u_x : \Pr(U_X = u_x) > 0\} \subseteq \cup_{\{x, z_1 : \Pr(X = x, Z = z_1) > 0\}} \mathcal{S}(x, z_1),$$

implying that $h(z_2, u_x) = h(z_1, u_x)$ for all attainable $z_1, z_2$ given a particular value of $u_x$, and this holds for all $u_x$ with $\Pr(U_X = u_x) > 0$. This is equivalent to $\text{CCI}(Z \to \hat{Y}) = 0$ (recall Lemma 1). $\qquad\square$

## APPENDIX B
### RELEVANT INFORMATION-THEORETIC PROPERTIES

**Lemma 7** (Conditional DPI). *For all $(A, A', B, X_c)$ such that $(B, X_c) - A - A'$ form a Markov chain, we have the following conditional form of the Data Processing Inequality (DPI):* $\text{I}(A; B \mid X_c) \geq \text{I}(A'; B \mid X_c)$.

*Proof of Lemma 7.* From the Markov chain, we have $\text{I}(A'; (B, X_c) \mid A) = 0$. Because, $\text{I}(A'; (B, X_c) \mid A) = \text{I}(A'; X_c \mid A) + \text{I}(A'; B \mid A, X_c)$ by chain rule and mutual information is non-negative, we also have $\text{I}(A'; B \mid A, X_c) = 0$. Now, similar to the proof of DPI, we have:

$$\text{I}(A'; B \mid X_c) + \text{I}(A; B \mid A', X_c) = \text{I}(A; B \mid X_c) + \text{I}(A'; B \mid A, X_c) = \text{I}(A; B \mid X_c), \tag{25}$$

because $I(A'; B \mid A, X_c) = 0$. This leads to $I(A; B \mid X_c) \geq I(A'; B \mid X_c)$. $\qquad \square$

**Lemma 8** (Triangle Inequality of Unique Information). *For all $(Z, B, A, X_c)$, we have:*

$$\mathrm{Uni}(Z : A|X_c) \leq \mathrm{Uni}(Z : A|B) + \mathrm{Uni}(Z : B|X_c).$$

This result is derived in [69, Proposition 2].

**Lemma 9** (Monotonicity under local operations on $Z$). *Let $Z' = f(Z)$ where $f(\cdot)$ is a deterministic function. Then, we have:*

$$\mathrm{Uni}(Z : B|X_c) \geq \mathrm{Uni}(Z' : B|X_c).$$

This result is derived in [57, Lemma 31]. We include a proof for completeness.

*Proof of Lemma 9.* Let $P'$ be the true joint distribution of $(Z', B, X_c)$ and $P$ be the true joint distribution of $(Z, B, X_c)$. Also let $Q^* = \arg\min_{Q \in \Delta_P} I_Q(Z; B \mid X_c)$ where $\Delta_P$ is the set of all joint distributions of $(Z, B, X_c)$ with the same marginals between $(Z, B)$ and $(Z, X_c)$ as the true joint distribution $P$. Let us also define

$$Q'^*(z', b, x_c) = \sum_z \Pr(z' \mid z) Q^*(z, b, x_c),$$

where $\Pr(z' \mid z)$ is the true conditional distribution of $Z' = f(Z)$ given $Z$.

Now, observe that,

$$
\begin{aligned}
\mathrm{Uni}(Z : B|X_c) &= \min_{Q \in \Delta_P} I_Q(Z; B \mid X_c) && \text{[By Definition of Unique Information]} \\
&= I_{Q^*}(Z; B \mid X_c) && \text{[By Definition of $Q^*$]} \\
&\overset{(a)}{\geq} I_{Q'^*}(Z'; B \mid X_c) \\
&\overset{(b)}{\geq} \min_{Q' \in \Delta_{P'}} I_{Q'}(Z'; B \mid X_c) \\
&= \mathrm{Uni}(Z' : B|X_c) && \text{[By Definition of Unique Information].} && (26)
\end{aligned}
$$

Here (a) holds using the conditional form of the Data Processing inequality (Lemma 7) as follows. Consider the random variables $(Z, B, X_c)$ following distribution $Q^*$ and $Z' = f(Z)$. Then, $(B, X_c) - Z - Z'$ form a Markov chain. Also note that (b) holds because $Q'^*$ belongs to $\Delta_{P'}$ which is the set of all joint distributions of $(Z', B, X_c)$ with the same marginals between $(Z', B)$ and $(Z', X_c)$ as the true joint distribution $P'$. $\qquad \square$

**Lemma 10** (Monotonicity under local operations on $B$). *Let $B' = f(B)$ where $f(\cdot)$ is a deterministic function. Then, we have:*

$$\mathrm{Uni}(Z : B|X_c) \geq \mathrm{Uni}(Z : B'|X_c).$$

This result is derived in [57, Lemma 31]. We include a proof for completeness.

*Proof of Lemma 10.* Let $P'$ be the true joint distribution of $(Z, B', X_c)$ and $P$ be the true joint distribution of $(Z, B, X_c)$. Also let $Q^* = \arg\min_{Q \in \Delta_P} I_Q(Z; B \mid X_c)$ where $\Delta_P$ is the set of all joint distributions of $(Z, B, X_c)$ with the same marginals between $(Z, B)$ and $(Z, X_c)$ as the true joint distribution $P$. Let us also define

$$Q'^*(z, b', x_c) = \sum_b \Pr(b' \mid b) Q^*(z, b, x_c),$$

where $\Pr(b' \mid b)$ is the true conditional distribution of $B' = f(B)$ given $B$.

Now, observe that,

$$
\begin{aligned}
\mathrm{Uni}(Z : B|X_c) &= \min_{Q \in \Delta_P} I_Q(Z; B \mid X_c) && \text{[By Definition of Unique Information]} \\
&= I_{Q^*}(Z; B \mid X_c) && \text{[By Definition of $Q^*$]} \\
&\overset{(a)}{\geq} I_{Q'^*}(Z; B' \mid X_c) \\
&\overset{(b)}{\geq} \min_{Q' \in \Delta_{P'}} I_{Q'}(Z; B' \mid X_c) \\
&= \mathrm{Uni}(Z : B'|X_c) && \text{[By Definition of Unique Information].} && (27)
\end{aligned}
$$

Here (a) holds using the conditional form of the Data Processing inequality (Lemma 7) as follows. Consider the random variables $(Z, B, X_c)$ following distribution $Q^*$ and $B' = f(B)$. Then, $(Z, X_c) - B - B'$ form a Markov chain. Also note that (b) holds because $Q'^*$ belongs to $\Delta_{P'}$ which is the set of all joint distributions of $(Z, B', X_c)$ with the same marginals between $(Z, B')$ and $(Z, X_c)$ as the true joint distribution $P'$. □

**Lemma 11** (Monotonicity under adversarial side information). *For all $(A, B, X_c, X_c')$, we have:*

$$\text{Uni}(A : B|(X_c, X_c')) \leq \text{Uni}(A : B|X_c).$$

This result is derived in [57, Lemma 32].

**Lemma 12** (Maximal conditional mutual information). *Let $A = f(Z, U_X)$ where $Z \perp\!\!\!\perp U_X$ and $B = g(U_X)$ for some deterministic functions $f(\cdot)$ and $g(\cdot)$ respectively. Then,*

$$I(Z; A \mid U_X) \geq I(Z; A \mid B). \tag{28}$$

*Proof of Lemma 12.* Observe that,

$$\begin{aligned}
I(Z; U_X \mid A, B)) &\geq 0 && \text{[non-negativity property]} \\
\implies H(Z \mid A, B) - H(Z \mid A, B, U_X) &\geq 0 && \text{[by definition]} \\
\implies H(Z \mid A, B) - H(Z \mid A, U_X) &\geq 0 && [B = g(U_X)] \\
\implies H(Z) - H(Z \mid A, U_X) &\geq H(Z) - H(Z \mid A, B) \\
\implies H(Z \mid U_X) - H(Z \mid A, U_X) &\geq H(Z \mid B) - H(Z \mid A, B) && [Z \perp\!\!\!\perp U_X \text{ and } Z \perp\!\!\!\perp B] \\
\implies I(Z; A \mid U_X) &\geq I(Z; A \mid B). \tag{29}
\end{aligned}$$

□

**Lemma 13** (Absence of counterfactual causal influence). *Let $\hat{Y} = h(Z, U_X)$ where $Z \perp\!\!\!\perp U_X$ and $X_c = g(Z, U_X)$ for some deterministic functions $h(\cdot)$ and $g(\cdot)$ respectively. Then $\text{CCI}(Z \to \hat{Y}) = 0$ implies $\text{Uni}(Z : (\hat{Y}, U_X)|X_c) = 0$ and also $\text{Uni}(Z : \hat{Y}|X_c) = 0$.*

*Proof of Lemma 13.* $\text{CCI}(Z \to \hat{Y}) = 0$ is equivalent to $I(Z; (\hat{Y}, U_X)) = 0$ (using Lemma 1). Now,

$$\text{Uni}(Z : (\hat{Y}, U_X)|X_c) \overset{(a)}{\leq} I(Z; (\hat{Y}, U_X)) = 0,$$

where (a) holds from (2) in Section II-A and non-negativity of PID. Also,

$$\text{Uni}(Z : \hat{Y}|X_c) \overset{(a)}{\leq} I(Z; \hat{Y}) \overset{(b)}{\leq} I(Z; (\hat{Y}, U_X)) = 0,$$

where (a) holds from (2) in Section II-A and non-negativity of PID terms, and (b) holds from the chain rule and non-negativity of mutual information.

□

**Lemma 14** (Zero-synergy property of deterministic functions). *Let $f(Z)$ be any deterministic function of $Z$, and let $X_c$ be any random variable. Then,*

$$\text{Syn}(Z : (f(Z), X_c)) = \text{Syn}(Z : (X_c, f(Z))) = 0. \tag{30}$$

*This leads to $\text{Uni}(Z : f(Z)|X_c) = I(Z; f(Z)|X_c)$ and $\text{Uni}(Z : X_c|f(Z)) = I(Z; X_c|f(Z))$.*

*Proof of Lemma 14:.* Recall from the definition of $\text{Uni}(Z : B|X_c)$ that $\Delta$ denotes the set of all joint distributions of $(Z, B, X_c)$ and $\Delta_p$ is the set of all such joint distributions that have the same marginals for $(Z, B)$ and $(Z, X_c)$ as the true distribution, *i.e.*,

$$\Delta_p = \{Q \in \Delta : q(z, b) = \Pr(Z = z, B = b) \text{ and } q(z, x_c) = \Pr(Z = z, X_c = x_c)\}. \tag{31}$$

We first show that if $B = f(Z)$, then $\Delta_p$ is only a singleton set which only consists of the true distribution. Observe that, for any $Q \in \Delta_p$,

$$
\begin{aligned}
q(z, b, x_c) &= q(z)q(b|z)q(x_c|b, z) && \text{[chain rule of probability]} \\
&= \Pr(Z = z) \Pr(B = b|Z = z)q(x_c|b, z) && [q(z, b) = \Pr(Z = z, B = b)] \\
&= \begin{cases} \Pr(Z = z)q(x_c|b, z), & \text{if } b = f(z) \\ 0, & \text{otherwise} \end{cases} && [\Pr(B = b|Z = z) = 1 \text{ only if } b = f(z)] \\
&= \begin{cases} \Pr(Z = z)q(x_c|z), & \text{if } b = f(z) \\ 0, & \text{otherwise} \end{cases} && [b \text{ is entirely determined by } z] \\
&= \begin{cases} \Pr(Z = z)\Pr(X_c = x_c|Z = z), & \text{if } y = f(z) \\ 0, & \text{otherwise} \end{cases} && [q(x_c|z) = \Pr(X_c = x_c|Z = z)] \\
&= \Pr(Z = z, B = b, X_c = x_c).
\end{aligned}
\tag{32}
$$

Thus, for $B = f(Z)$,

$$
\text{Uni}(Z : B|X_c) = \min_{Q \in \Delta_p} \text{I}_Q(Z; B|X_c) = \text{I}(Z; B|X_c).
\tag{33}
$$

This leads to $\text{Syn}(Z : (f(Z), X_c)) = \text{I}(Z; f(Z)|X_c) - \text{Uni}(Z : f(Z)|X_c) = 0$ (using (3) in Section II-A). Note that, $\text{Syn}(Z : (f(Z), X_c))$ is symmetric between $f(Z)$ and $X_c$. □

## APPENDIX C
## APPENDIX TO SECTION III

Here, we provide the proofs of the results as well as additional discussion to supplement Section III. For convenience, we repeat the statements of the results.

### A. Proof of Theorem 1 and Lemma 2

**Theorem 1** (Properties). *Properties 1-6 are satisfied by our proposed measure*

$$
M_{NE}^* = \min_{U_a, U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \text{ such that } U_a = U_X \backslash U_b.
$$

*Proof of Theorem 1.* Here, we formally show that our proposed measure satisfies all the four desirable properties. We restate each of the properties again and then show that they are is satisfied.

**Property 4** (Absence of Exemptions). *If no feature is deemed critical ($X_c = \phi$), then a measure $M_{NE}$ should be equal to the total bias, i.e., $\text{I}(Z; (\hat{Y}, U_X))$.*

When $X_c = \phi$, we have $\text{Uni}(Z, U_a : \hat{Y}, U_b|X_c) = \text{I}(Z, U_a; \hat{Y}, U_b)$. We are required to show that $\min_{U_a, U_b \text{ s.t. } U_a = U_X \backslash U_b} \text{I}(Z, U_a; \hat{Y}, U_b)$ is equal to $\text{I}(Z; (\hat{Y}, U_X))$. Note that,

$$
\begin{aligned}
\text{I}(Z, U_a; \hat{Y}, U_b) &= \text{H}(\hat{Y}, U_b) - \text{H}(\hat{Y}, U_b \mid Z, U_a) && \text{[By Definition]} \\
&= \text{H}(\hat{Y} \mid U_b) + \text{H}(U_b) - \text{H}(U_b \mid Z, U_a) - \text{H}(\hat{Y} \mid U_b, Z, U_a) && \text{[Chain Rule]} \\
&= \text{H}(\hat{Y} \mid U_b) + \text{H}(U_b) - \text{H}(U_b \mid Z, U_a) && [\hat{Y} \text{ is entirely determined by } Z, U_a, U_b] \\
&= \text{H}(\hat{Y} \mid U_b) && [Z, U_a, U_b \text{ are mutually independent]} \\
&\geq \text{H}(\hat{Y} \mid U_X) && \text{[conditioning reduces entropy]} \\
&= \text{H}(\hat{Y} \mid U_X) - \text{H}(\hat{Y} \mid Z, U_X) + \text{I}(Z; U_X) && [\hat{Y} \text{ entirely determined by } Z, U_X, \text{ and } Z \perp\!\!\!\perp U_X] \\
&= \text{I}(Z; \hat{Y} \mid U_X) + \text{I}(Z; U_X) && \text{[By Definition]} \\
&= \text{I}(Z; (\hat{Y}, U_X)). && \text{[By Chain Rule]}
\end{aligned}
\tag{34}
$$

Thus, $\text{I}(Z, U_a; \hat{Y}, U_b) \geq \text{I}(Z; (\hat{Y}, U_X))$ with equality when $U_b = U_X, U_a = \phi$.

**Property 5** (Non-Increasing with More Exemptions). *For a fixed set of features $X$ and a fixed model $\hat{Y} = h(Z, U_X)$, a measure $M_{NE}$ should be non-increasing if a feature is removed from $X_g$ and added to $X_c$.*

Let $X_c'$ denote the additional feature that is to be removed from $X_g$ and is to be added to $X_c$. From Lemma 11, we have,

$$
\text{Uni}((Z, U_a) : (\hat{Y}, U_b)|(X_c, X_c')) \leq \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c),
\tag{35}
$$

for any $U_a, U_b$. Thus,

$$\min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | (X_c, X'_c)) \leq \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c). \tag{36}$$

**Property 6** (Complete Exemption). $M_{NE}$ *should be* $0$ *if all features are exempt, i.e.,* $X_c = X$ *and* $X_g = \phi$.

Observe that, when $X = X_c$,

$$\begin{aligned} M^*_{NE} &= \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X) \\ &\leq \text{Uni}(Z, U_X : \hat{Y} | X) \\ &\leq \text{I}(Z, U_X; \hat{Y} \mid X) && \text{[(3) in Section II-A and non-negativity of PID terms]} \\ &= \text{H}(\hat{Y} \mid X) - \text{H}(\hat{Y} \mid Z, U_X, X) && \text{[By Definition]} \\ &= 0. && \text{[$\hat{Y}$ is a deterministic function of $X$]} \end{aligned} \tag{37}$$

**Property 1** (Zero Influence). $M_{NE}$ *should be* $0$ *if* $\text{CCI}(Z \to \hat{Y}) = 0$ *(or equivalently,* $\text{I}(Z; \hat{Y}, U_X) = 0$*).*

$$\begin{aligned} M^*_{NE} &= \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) \\ &\leq \text{Uni}(Z : (\hat{Y}, U_X) | X_c) \\ &\leq \text{I}(Z; (\hat{Y}, U_X)). \quad \text{[(2) in Section II-A and non-negativity of PID terms]} \end{aligned} \tag{38}$$

Thus, $\text{I}(Z; (\hat{Y}, U_X)) = 0$ implies $M_{NE} = 0$.

**Property 2** (Non-Exempt Statistically Visible Bias). $M_{NE}$ *should be strictly greater than* $0$ *if* $\hat{Y}$ *has any unique information about* $Z$ *not present in* $X_c$*. Thus,* $\text{Uni}(Z : \hat{Y} | X_c) > 0$ *should imply that* $M_{NE} > 0$.

$$\begin{aligned} M^*_{NE} &= \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) \\ &= \text{Uni}((Z, U^*_a) : (\hat{Y}, U^*_b) | X_c) && \text{[for some $(U^*_a, U^*_b)$]} \\ &\geq \text{Uni}(Z : (\hat{Y}, U^*_b) | X_c) && \text{[Using Lemma 9]} \\ &\geq \text{Uni}(Z : \hat{Y} | X_c). && \text{[Using Lemma 10]} \end{aligned} \tag{39}$$

Thus, $\text{Uni}(Z : \hat{Y} | X_c) > 0$ implies that $M_{NE} > 0$.

**Property 3** (Non-Exempt Masked Bias). $M_{NE}$ *should be non-zero in the canonical example of non-exempt masked bias:* $X_1 = Z$, $X_2 = U_X$, *and* $\hat{Y} = Z \oplus U_X$ *with* $Z, U_X \sim$ *i.i.d. Bern(½) and* $X_1 \in X_g$*. However,* $M_{NE}$ *should be* $0$ *if* $(Z, U_a) - X_c - (\hat{Y}, U_b)$ *form a Markov chain for some subsets* $U_a, U_b \subseteq U_X$ *such that* $U_a = U_X \setminus U_b$.

First we will show that $M^*_{NE} > 0$ for the canonical example of non-exempt bias where $\hat{Y} = Z \oplus U_{X_1}$ where $Z$ lies in the non-critical/general features and $U_{X_1}$ can be either critical or non-critical.

**Case 1:** $X_c = U_{X_1}$, $X_g = Z$ and $\hat{Y} = Z \oplus U_{X_1}$ with $Z, U_{X_1} \sim$ i.i.d. Bern(½).

We will check the value of $\text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c)$ for different choices of $U_a$ to find the minimum.

For $U_a = \phi$ and $U_b = U_{X_1}$, we have

$$\begin{aligned} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) &= \text{Uni}(Z : (\hat{Y}, U_{X_1}) | X_c) && \text{[Substituting the variables]} \\ &= \text{I}(Z; (\hat{Y}, U_{X_1})) - \text{Red}(Z : ((\hat{Y}, U_{X_1}), X_c)) && \text{[Using (2) in Section II-A]} \\ &\overset{(a)}{=} \text{I}(Z; (\hat{Y}, U_{X_1})) \\ &= 1 \text{ bit.} \end{aligned} \tag{40}$$

Here (a) holds because $\text{Red}(Z : ((\hat{Y}, U_{X_1}), X_c)) \leq \text{I}(Z; X_c)$ (using (2) in Section II-A and non-negativity of PID terms), and here $\text{I}(Z; X_c) = 0$.

For $U_a = U_{X_1}$ and $U_b = \phi$, we have

$$\begin{aligned} \text{Uni}((Z, U_a) : (\hat{Y}, U_b) | X_c) &= \text{Uni}((Z, U_{X_1}) : \hat{Y} | X_c) && \text{[Substituting the variables]} \\ &= \text{I}((Z, U_{X_1}); \hat{Y} \mid X_c) && \text{[Lemma 14 as $\hat{Y}$ is deterministic function of $f(Z, U_{X_1})$]} \\ &= 1 \text{ bit.} \end{aligned} \tag{41}$$

Thus,

$$M^*_{NE} = \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = 1 \text{ bit,}$$

which is strictly greater than 0.

**Case 2:** $X_c = \phi$, $X_g = (Z, U_{X_1})$ and $\hat{Y} = Z \oplus U_{X_1}$ with $Z, U_{X_1} \sim i.i.d.$ Bern(½).

Since $X_c = \phi$, we can use Property 4 (proved above) to compute

$$M^*_{NE} = \text{I}(Z; (\hat{Y}, U_X)) = 1 \text{ bit,}$$

which is strictly greater than 0. Thus, our proposed measure is non-zero in the canonical example of non-exempt masked bias. Now, we move on to the proof of the next part of this property.

Suppose that $(Z, U_a) - X_c - (\hat{Y}, U_b)$ form a Markov chain for some subsets $U_a, U_b \subseteq U_X$ such that $U_a = U_X \setminus U_b$. Then, $\text{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c) = 0$, implying that $\text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = 0$ for those subsets $U_a, U_b \subseteq U_X$ because unique information is a sub-component of conditional mutual information. Therefore,

$$M^*_{NE} = \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \leq 0.$$

Again, using the fact that unique information is non-negative, we have,

$$M^*_{NE} = \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) \geq 0.$$

Thus, $M^*_{NE} = 0$.

$\square$

**Lemma 2.** *The Markov chain $(Z, U_a) - X_c - (\hat{Y}, U_b)$ implies that the following Markov chains also hold: (i) $Z - X_c - \hat{Y}$; (ii) $(Z, U_a) - X_c - \hat{Y}$; and (ii) $Z - X_c - (\hat{Y}, U_b)$.*

*Proof of Lemma 2.* We note that the terms $\text{I}(Z; \hat{Y} \mid X_c)$, $\text{I}(Z; (\hat{Y}, U_b) \mid X_c)$ and $\text{I}((Z, U_a); \hat{Y} \mid X_c)$ are all less than or equal to $\text{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c)$ using the chain rule and non-negativity of conditional mutual information.

Thus, if $\text{I}((Z, U_a); (\hat{Y}, U_b) \mid X_c) = 0$, then all those three terms are also 0.

$\square$

### B. Supporting Derivations

Here, we include the supporting derivations for some of our statements in Section III-A and Section III-B.

**Supporting Derivation 1:** $\text{Uni}(Z : \hat{Y}|X_c) > 0$ **for Canonical Example 2 (discrimination in admissions).**

*Proof.* Recall that for this example, $X_c = U_{X_1}$, $X_g = Z \oplus U_{X_2}$, and $\hat{Y} = U_{X_1} + Z + U_{X_2}$ with $Z, U_{X_1}, U_{X_2} \sim i.i.d.$ Bern(½). The claim can be verified as follows:

$$\begin{aligned}
\text{Uni}(Z : \hat{Y}|X_c) &= \text{I}(Z; \hat{Y}) - \text{Red}(Z : (\hat{Y}, X_c)) &&\text{[using (2) in Section II-A]} \\
&\overset{(a)}{\geq} \text{I}(Z; \hat{Y}) - \text{I}(Z; X_c) \\
&\overset{(b)}{=} \text{I}(Z; \hat{Y}) \\
&\overset{(c)}{>} 0,
\end{aligned}$$

where (a) holds because $\text{Red}(Z : (\hat{Y}, X_c)) \leq \text{I}(Z; X_c)$ (using (2) in Section II-A and non-negativity of all PID terms) and (b) holds because $\text{I}(Z; X_c) = 0$. Lastly, (c) holds because $\hat{Y}$ and $Z$ are not independent of each other for this specific example. $\square$

**Supporting Derivation 2:** $\text{Uni}(Z : \hat{Y}|X_c) > 0$ **for Canonical Example 6 (discrimination by unmasking).**

*Proof.* Recall that for this example, $X_c = Z \oplus U_{X_1}$, $X_g = U_{X_1}$ and $\hat{Y} = Z$ with $Z, U_{X_1} \sim i.i.d.$ Bern(½).

The claim can be verified as follows:

$$\begin{aligned}
\text{Uni}(Z : \hat{Y}|X_c) &= \text{I}(Z; \hat{Y}) - \text{Red}(Z : (\hat{Y}, X_c)) &&\text{[using (2) in Section II-A]} \\
&\overset{(a)}{\geq} \text{I}(Z; \hat{Y}) - \text{I}(Z; X_c) \\
&\overset{(b)}{=} 1 \text{ bit,}
\end{aligned}$$

where (a) holds because $\text{Red}(Z : (\hat{Y}, X_c)) \leq \text{I}(Z; X_c)$ (using (2) in Section II-A and non-negativity of all PID terms) and (b) holds because $\text{I}(Z; X_c) = 0$. $\square$

**Supporting Derivation 3:** $\text{Uni}(Z : (\hat{Y}, U_X)|X_c) > 0$ **in Canonical Example 1.**

*Proof.* Consider Canonical Example 1.

$$\mathrm{Uni}(Z : (\hat{Y}, U_X)|X_c) = \mathrm{Uni}(Z : (Z + U_{X_1} + U_{X_2}, U_X)|Z + U_{X_1}) \qquad \text{[Substituting the variables]}$$

$$\overset{(a)}{\geq} \mathrm{Uni}(Z : Z|Z + U_{X_1})$$

$$\overset{(b)}{=} \mathrm{I}(Z; Z \mid Z + U_{X_1})$$

$$\overset{(c)}{>} 0.$$

Here, (a) holds because $Z$ is a deterministic function of $(Z + U_{X_1} + U_{X_2}, U_X)$ and unique information is non-increasing under local operations of $B$ (see Lemma 10 in Appendix B). Next, (b) holds because if we consider $\Delta_p$, the set of joint distributions of $(Z, Z, Z + U_{X_1})$, such that the marginals $(Z, Z)$ and $(Z, Z + U_{X_1})$ are the same as the marginals of the true joint distribution, we find that there is only one distribution in this set, which is exactly the true distribution. Thus, $\mathrm{Uni}(Z : Z|Z + U_{X_1}) = \min_{Q \in \Delta_p} \mathrm{I}_Q(Z; Z \mid Z + U_{X_1}) = \mathrm{I}(Z; Z \mid Z + U_{X_1})$. Lastly (c) holds because,

$$\mathrm{I}(Z; Z \mid Z + U_{X_1}) = \mathrm{H}(Z|Z + U_{X_1}) - \mathrm{H}(Z|Z, Z + U_{X_1}) = \mathrm{H}(Z|Z + U_{X_1}) = \sum_{t=0,1,2} \mathrm{H}(Z|Z + U_{X_1} = t)\Pr(Z + U_{X_1}=t).$$

Using the fact that $Z, U_{X_1} \sim i.i.d.$ Bern($\frac{1}{2}$), we can compute $\mathrm{H}(Z|Z + U_{X_1} = 0) = 0$, $\mathrm{H}(Z|Z + U_{X_1} = 1) = h_b(\frac{1}{2}) = 1$, and $\mathrm{H}(Z|Z + U_{X_1} = 2) = 0$. Here, $h_b(\cdot)$ is the binary entropy function [63] given by $h_b(p) = -p\log_2(p) - (1-p)\log_2(1-p)$. Also note that, $\Pr(Z + U_{X_1} = 1) = \frac{1}{2}$. So, $\mathrm{I}(Z; Z \mid Z + U_{X_1}) = 0.5$ bits. $\square$

**Supporting Derivation 4: Exact computation of $\mathrm{Uni}(Z : \hat{Y}|X_c)$ and $M_{NE}^*$ for Canonical Example 2.**

$$\mathrm{Uni}(Z : \hat{Y}|X_c) \overset{(a)}{=} \mathrm{I}(Z; \hat{Y})$$

$$= \mathrm{H}(Z) - \mathrm{H}(Z|\hat{Y})$$

$$= \mathrm{H}(Z) - \mathrm{H}(Z|U_{X_1} + Z + U_{X_2})$$

$$= \mathrm{H}(Z) - \sum_{t=0,1,2,3} \mathrm{H}(Z|U_{X_1} + Z + U_{X_2} = t)\Pr(U_{X_1} + Z + U_{X_2} = t)$$

$$\overset{(b)}{=} 1 - \frac{3}{4}h_b(\frac{1}{3}) \text{ bits.} \tag{42}$$

Here (a) holds because $\mathrm{I}(Z; U_{X_1}) = 0$, implying $\mathrm{Red}(Z : (\hat{Y}, U_{X_1})) = 0$ as well (using (2) in Section II-A and non-negativity of PID terms). Lastly, (b) holds because $Z, U_{X_1}, U_{X_2} \sim i.i.d.$ Bern($\frac{1}{2}$). So, we can exactly compute $\mathrm{H}(Z|U_{X_1} + Z + U_{X_2} = 0) = 0$, $\mathrm{H}(Z|U_{X_1} + Z + U_{X_2} = 1) = h_b(\frac{1}{3})$, $\mathrm{H}(Z|U_{X_1} + Z + U_{X_2} = 2) = h_b(\frac{1}{3})$, and $\mathrm{H}(Z|U_{X_1} + Z + U_{X_2} = 3) = 0$. Here, $h_b(\cdot)$ is the binary entropy function [63] given by $h_b(p) = -p\log_2(p) - (1-p)\log_2(1-p)$. Also note that, $\Pr(U_{X_1} + Z + U_{X_2} = 1) = \Pr(U_{X_1} + Z + U_{X_2} = 2) = \frac{3}{8}$.

Now, we will examine the value of $\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c)$ for different choices of $U_a$ to find the minimum.

Let $U_a = \phi$ (and $U_b = U_X$). Then,

$$\mathrm{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = \mathrm{Uni}(Z : (\hat{Y}, U_{X_1}, U_{X_2})|U_{X_1})$$

$$\overset{(a)}{=} \mathrm{I}(Z; U_{X_1} + Z + U_{X_2}, U_{X_1}, U_{X_2})$$

$$= \mathrm{I}(Z; U_{X_1}, U_{X_2}) + \mathrm{I}(Z; U_{X_1} + Z + U_{X_2} \mid U_{X_1}, U_{X_2}) \qquad \text{[Chain Rule]}$$

$$= \mathrm{I}(Z; U_{X_1} + Z + U_{X_2} \mid U_{X_1}, U_{X_2}) \qquad \text{[$Z$ is independent of $U_{X_1}, U_{X_2}$]}$$

$$= \mathrm{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}, U_{X_2}) - \mathrm{H}(U_{X_1} + Z + U_{X_2} \mid Z, U_{X_1}, U_{X_2}) \qquad \text{[By Definition]}$$

$$= \mathrm{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}, U_{X_2}) \qquad \text{[Deterministic Function]}$$

$$= \sum_{u_1, u_2 \in \{0,1\}} \mathrm{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1} = u_1, U_{X_2} = u_2)\Pr(U_{X_1} = u_1, U_{X_2} = u_2)$$

$$= \sum_{u_1, u_2 \in \{0,1\}} h_b(\frac{1}{2})\Pr(U_{X_1} = u_1, U_{X_2} = u_2)$$

$$= 1 \text{ bit.} \tag{43}$$

Here (a) holds again because $\mathrm{I}(Z; U_{X_1}) = 0$, implying the redundant information is 0 as well (using (2) in Section II-A).

Next, for $U_a = U_{X_2}$ (and $U_b = U_{X_1}$), we have,

$$
\begin{aligned}
\text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) &= \text{Uni}((Z, U_{X_2}) : (\hat{Y}, U_{X_1})|U_{X_1}) \\
&\overset{(a)}{=} \text{I}((Z, U_{X_2}); (\hat{Y}, U_{X_1})) \\
&= \text{I}((Z, U_{X_2}); U_{X_1}) + \text{I}((Z, U_{X_2}); \hat{Y} \mid U_{X_1}) && \text{[Chain Rule]} \\
&= \text{I}((Z, U_{X_2}); \hat{Y} \mid U_{X_1}) && [Z, U_{X_2} \text{ is independent of } U_{X_1}] \\
&= \text{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}) - \text{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}, (Z, U_{X_2})) && \text{[By Definition]} \\
&= \text{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}) && \text{[Deterministic Function]} \\
&= \sum_{u_1 = 0,1} \text{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1} = u_1) \Pr(U_{X_1} = u_1) \\
&= {}^1\!/_4 \log_2 4 + {}^1\!/_2 \log_2 2 + {}^1\!/_4 \log_2 4 \\
&= {}^3\!/_2 \text{ bit.} && (44)
\end{aligned}
$$

Here (a) holds again because $\text{I}((Z, U_{X_2}); U_{X_1}) = 0$, implying the redundant information is 0 as well (using (2) in Section II-A).

Next, for $U_a = U_{X_1}$ (and $U_b = U_{X_2}$), we have,

$$
\begin{aligned}
\text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) &= \text{Uni}((Z, U_{X_1}) : (\hat{Y}, U_{X_2})|U_{X_1}) \\
&\overset{(b)}{=} \text{I}((Z, U_{X_1}); (\hat{Y}, U_{X_2}) \mid U_{X_1}) \\
&= \text{I}((Z, U_{X_1}); U_{X_2} \mid U_{X_1}) + \text{I}((Z, U_{X_1}); \hat{Y} \mid U_{X_1}, U_{X_2}) && \text{[Chain Rule]} \\
&= \text{I}((Z, U_{X_1}); \hat{Y} \mid U_{X_1}, U_{X_2}) && \text{[Mutual Independence]} \\
&= \text{H}(\hat{Y} \mid U_{X_1}, U_{X_2}) - \text{H}(\hat{Y} \mid (Z, U_{X_1}), U_{X_1}, U_{X_2}) && \text{[By Definition]} \\
&= \text{H}(\hat{Y} \mid U_{X_1}, U_{X_2}) && \text{[Deterministic Function]} \\
&= \text{H}(U_{X_1} + Z + U_{X_2} \mid U_{X_1}, U_{X_2}) \\
&= 1 \text{ bit.} && (45)
\end{aligned}
$$

Here (b) holds because $\text{Syn}((Z, U_{X_1}) : (A, B)) = 0$ if one of the terms $A$ or $B$ is a deterministic function of $(Z, U_{X_1})$ (using Lemma 14 in Appendix B) and hence unique information becomes equal to the conditional mutual information (see (3) in Section II-A).

Lastly, for $U_a = U_X$ (and $U_b = \phi$), we have,

$$
\begin{aligned}
\text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) &= \text{Uni}((Z, U_{X_1}, U_{X_2}) : \hat{Y}|U_{X_1}) \\
&\overset{(b)}{=} \text{I}((Z, U_{X_1}, U_{X_2}); \hat{Y} \mid U_{X_1}) \\
&= \text{H}(\hat{Y} \mid U_{X_1}) - \text{H}(\hat{Y} \mid (Z, U_{X_1}, U_{X_2}), U_{X_1}) && \text{[By Definition]} \\
&= \text{H}(\hat{Y} \mid U_{X_1}) && \text{[Deterministic Function]} \\
&= {}^1\!/_4 \log_2 4 + {}^1\!/_2 \log_2 2 + {}^1\!/_4 \log_2 4 \\
&= {}^3\!/_2 \text{ bit.} && (46)
\end{aligned}
$$

Here (b) holds again using Lemma 14 in Appendix B.

Thus, we obtain that,

$$
M_{NE}^* = \min_{U_a, U_b \text{ s.t. } U_a = U_X \setminus U_b} \text{Uni}((Z, U_a) : (\hat{Y}, U_b)|X_c) = 1 \text{ bit.}
$$

This is strictly greater than $\text{Uni}(Z : \hat{Y}|X_c) = 1 - \frac{3}{4} h_b({}^1\!/_3)$ bits, accounting for both non-exempt statistically visible and non-exempt masked biases.

## C. Discussion on Other Candidate Measures

**Why the product of the two measures $\text{I}(Z; \hat{Y} \mid X_c)$ and $\text{I}(Z; (\hat{Y}, U_X))$ does not work?**

One might recall that the measure $\text{I}(Z; \hat{Y} \mid X_c)$ resolved most of the examples except in Canonical Example 3 where the output $\hat{Y}$ had no counterfactual causal influence of $Z$ and yet this measure gave a false positive conclusion about non-exempt bias. This leads us to examine another candidate measure, i.e., product of $\text{I}(Z; \hat{Y} \mid X_c)$ and $\text{I}(Z; (\hat{Y}, U_X))$ where the latter is always 0 whenever there is no counterfactual causal influence of $Z$ on $\hat{Y}$.

**Candidate Measure of Non-Exempt Bias 4.** $M_{NE} = \text{I}(Z; \hat{Y} \mid X_c) \times \text{I}(Z; (\hat{Y}, U_X))$.

**Canonical Example 7.** *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. Bern(½).*

This example should be exempt because $Z_2$ already appears in $X_c$, and is hence exempt. However, both $\mathrm{I}(Z; (\hat{Y}, U_X))$ and $\mathrm{I}(Z; \hat{Y} \mid X_c)$ are non-zero for this example. This leads us to examine another candidate measure, which is essentially the common information-theoretic volume between $\mathrm{I}(Z; (\hat{Y}, U_X))$ and $\mathrm{I}(Z; \hat{Y} \mid X_c)$, i.e., a measure of the common reason that can make both $\mathrm{I}(Z; (\hat{Y}, U_X)) > 0$ and $\mathrm{I}(Z; \hat{Y} \mid X_c) > 0$ (overlapping volume).

**Measure proposed in [1]: Information-theoretic sub-volume of the intersection between** $\mathrm{I}(Z; \hat{Y} \mid X_c)$ **and** $\mathrm{I}(Z; (\hat{Y}, U_X))$**:**
The previous Canonical Example demonstrates that both these measures $\mathrm{I}(Z; \hat{Y} \mid X_c)$ and $\mathrm{I}(Z; (\hat{Y}, U_X))$ can be non-zero for different reasons leading to a false positive conclusion using Candidate Measure 4. Intuitively, we need to identify the common reason that makes them non-zero, if any. This motivates us to examine another candidate (Candidate Measure 5) which is the information-theoretic sub-volume of the intersection between these two measures, as shown in Fig. 13.

**Candidate Measure of Non-Exempt Bias 5.** $M_{NE} = \mathrm{Uni}(Z : (\hat{Y}, U_X)|X_c) - \mathrm{Uni}(Z : (\hat{Y}, U_X)|(X_c, \hat{Y}))$.
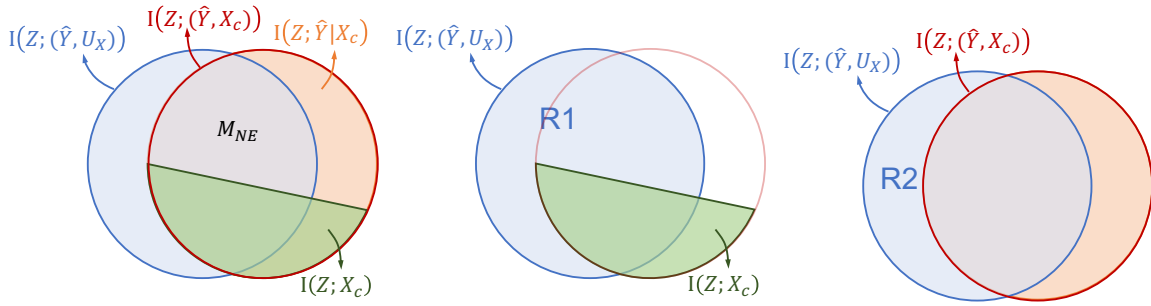


Fig. 13: (Left) Notice that the blue full-circle denotes $\mathrm{I}(Z; (\hat{Y}, U_X))$ and the red full-circle denotes $\mathrm{I}(Z; (\hat{Y}, X_c))$. The term $\mathrm{I}(Z; (\hat{Y}, X_c))$ is equal to the sum of $\mathrm{I}(Z; X_c)$ (green half-circle) and $\mathrm{I}(Z; \hat{Y} \mid X_c)$ (orange half-circle). The candidate measure $(M_{NE})$ is the intersecting volume between $\mathrm{I}(Z; (\hat{Y}, U_X))$ and $\mathrm{I}(Z; \hat{Y} \mid X_c)$. Next, we show pictorially that this intersecting volume is given by $R1 - R2$ where $R1$ is shown in the middle figure and $R2$ is shown in the rightmost figure. (Middle) Notice that $R1 = \mathrm{Uni}(Z : (\hat{Y}, U_X)|X_c)$. (Right) Notice that $R2 = \mathrm{Uni}(Z : (\hat{Y}, U_X)|(\hat{Y}, X_c))$.

**Limitations of Candidate Measure 5:** This measure does resolve many of the examples and satisfies several desirable properties (discussed more in [1]). However, it fails to capture certain types of non-exempt masked bias when the mask arises from $X_g$, e.g., scenarios like Canonical Example 5 in Section III-B, where non-exempt masked bias is present even though $Z - X_c - \hat{Y}$ form a Markov chain.

APPENDIX D
APPENDIX TO SECTION IV

*A. Proof of Theorem 2 and Lemma 3*

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

$$\mathrm{I}(Z; (\hat{Y}, U_X)) = M_{V,NE} + M_{V,E} + M_{M,NE} + M_{M,E}. \tag{11}$$

*Here* $M_{V,NE} = \mathrm{Uni}(Z : \hat{Y}|X_c)$ *and* $M_{V,E} = \mathrm{Red}(Z : (\hat{Y}, X_c))$. *These two terms add to form* $\mathrm{I}(Z; \hat{Y})$ *which is the total statistically visible bias. Next,* $M_{M,NE} = M_{NE}^* - M_{V,NE}$ *where* $M_{NE}^*$ *is our proposed measure of non-exempt bias (Definition 7), and* $M_{M,E} = \mathrm{I}(Z; \hat{Y}, U_X) - \mathrm{I}(Z; \hat{Y}) - M_{M,NE}$. *All of these components are non-negative.*

*Proof of Theorem 2.* First consider $M_{V,NE} = \mathrm{Uni}(Z : \hat{Y}|X_c)$ and $M_{V,E} = \mathrm{Red}(Z : (\hat{Y}, X_c))$. 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,

$$
\begin{aligned}
M_{M,E} &= \mathrm{I}(Z;(\hat{Y},U_X)) - \mathrm{I}(Z;\hat{Y}) - M_{M,NE} && \text{[By Definition]}\\
&= \mathrm{I}(Z;\hat{Y}) + \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{I}(Z;\hat{Y}) - M_{M,NE} && \text{[Chain Rule for mutual information]}\\
&= \mathrm{I}(Z;U_X \mid \hat{Y}) - M_{M,NE}\\
&= \mathrm{I}(Z;U_X \mid \hat{Y}) - M^*_{NE} + M_{V,NE} && \text{[By Definition]}\\
&= \mathrm{I}(Z;U_X \mid \hat{Y}) - \min_{U_a,U_b \text{ s.t. } U_a = U_X \setminus U_b} \mathrm{Uni}((Z,U_a):(\hat{Y},U_b)|X_c) + \mathrm{Uni}(Z:\hat{Y}|X_c) && \text{[By Definition]}\\
&\geq \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{Uni}(Z:(\hat{Y},U_X)|X_c) + \mathrm{Uni}(Z:\hat{Y}|X_c)\\
&\geq \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{Uni}(Z:(\hat{Y},U_X)|\hat{Y}) && \text{[Triangle Inequality (Lemma 8)]}\\
&\geq \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{I}(Z;(\hat{Y},U_X) \mid \hat{Y}) && \text{[(3) in Section II-A]}\\
&= \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{I}(Z;U_X \mid \hat{Y}) - \mathrm{I}(Z;\hat{Y} \mid U_X,\hat{Y}) && \text{[Chain Rule for mutual information]}\\
&= 0. && (47)
\end{aligned}
$$

Lastly, we consider $M_{M,NE} = \min_{U_a,U_b \text{ s.t. } U_a = U_X \setminus U_b} \mathrm{Uni}((Z,U_a):(\hat{Y},U_b)|X_c) - \mathrm{Uni}(Z:\hat{Y}|X_c)$.

$$
\begin{aligned}
M_{NE} &= \min_{U_a,U_b \text{ s.t. } U_a = U_X \setminus U_b} \mathrm{Uni}((Z,U_a):(\hat{Y},U_b)|X_c) - \mathrm{Uni}(Z:\hat{Y}|X_c)\\
&= \mathrm{Uni}((Z,U_a^*):(\hat{Y},U_b^*)|X_c) - \mathrm{Uni}(Z:\hat{Y}|X_c) && \text{[for some } (U_a^*,U_b^*)]\\
&\geq \mathrm{Uni}(Z:(\hat{Y},U_b^*)|X_c) - \mathrm{Uni}(Z:\hat{Y}|X_c) && \text{[Using Lemma 9]}\\
&\geq \mathrm{Uni}(Z:\hat{Y}|X_c) - \mathrm{Uni}(Z:\hat{Y}|X_c) && \text{[Using Lemma 10]}\\
&= 0. && (48)
\end{aligned}
$$

$\square$

**Lemma 3** (Conditioning to Capture Masked Bias). *The following two statements are equivalent:*
- *Masked bias* $\mathrm{I}(Z;(\hat{Y},U_X)) - \mathrm{I}(Z;\hat{Y}) > 0$.
- $\exists$ *a random variable $G$ of the form $G = g(U_X)$ such that* $\mathrm{I}(Z;\hat{Y} \mid G) - \mathrm{I}(Z;\hat{Y}) > 0$.

*Proof of Lemma 3.* Before proceeding, note that, $\mathrm{I}(Z;\hat{Y},U_X) = \mathrm{I}(Z;U_X) + \mathrm{I}(Z;\hat{Y} \mid U_X) = \mathrm{I}(Z;\hat{Y} \mid U_X)$ because $Z$ is independent of $U_X$. This also leads to the masked bias being equal to $\mathrm{I}(Z;\hat{Y} \mid U_X) - \mathrm{I}(Z;\hat{Y})$.

First, we show that the first statement implies the second statement. Suppose that, masked bias $\mathrm{I}(Z;\hat{Y} \mid U_X) - \mathrm{I}(Z;\hat{Y}) > 0$. Then, we can choose the function $G = U_X$ such that $\mathrm{I}(Z;\hat{Y} \mid G) - \mathrm{I}(Z;\hat{Y}) > 0$. Thus, the implication holds.

We will now show that the second statement also implies the first statement. First note that, using Lemma 12, for any deterministic $g(\cdot)$, we always have $\mathrm{I}(Z;\hat{Y} \mid U_X) \geq \mathrm{I}(Z;\hat{Y} \mid g(U_X))$. Now, suppose there exists a $G = g(U_X)$ such that $\mathrm{I}(Z;\hat{Y} \mid G) > \mathrm{I}(Z;\hat{Y})$. Then, $\mathrm{I}(Z;\hat{Y} \mid U_X) \geq \mathrm{I}(Z;\hat{Y} \mid g(U_X)) > \mathrm{I}(Z;\hat{Y})$, implying masked bias is present.

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

$\square$

## APPENDIX E
## APPENDIX TO SECTION VI

**Lemma 4.** *[Fairness Properties of $\mathrm{Uni}(Z:\hat{Y}|X_c)$] The measure $\mathrm{Uni}(Z:\hat{Y}|X_c)$ satisfies Properties 5, 6, 1, and 2.*

*Proof of Lemma 4.* Property 5 is satisfied using Lemma 11 in Appendix B (originally derived in [57, Lemma 32]).

Property 6 is satisfied because $\hat{Y}$ is a deterministic function of the entire $X$, and hence the Markov chain $Z - X - \hat{Y}$ holds. Thus $\mathrm{I}(Z;\hat{Y} \mid X_c) = 0$, also implying $\mathrm{Uni}(Z:\hat{Y}|X_c) = 0$.

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

For Property 1, observe that,

$$
\begin{aligned}
\mathrm{CCI}(Z \to \hat{Y}) &= 0\\
\implies \mathrm{I}(Z;\hat{Y}) &= 0\\
\implies \mathrm{Uni}(Z:\hat{Y}|X_c) + \mathrm{Red}(Z:(\hat{Y},X_c)) &= 0 && \text{[Using (2) in Section II-A]}\\
\implies \mathrm{Uni}(Z:\hat{Y}|X_c) &= 0 && \text{[Non-negativity of PID terms].} \quad (49)
\end{aligned}
$$

$\square$

**Lemma 5.** *[Fairness Properties of $\mathrm{I}(Z;\hat{Y} \mid X_c)$] The measure $\mathrm{I}(Z;\hat{Y} \mid X_c)$ satisfies Properties 6 and 2.*

*Proof of Lemma 5.* Property 6 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

$$\text{Uni}(Z : \hat{Y}|X_c) > 0$$
$$\implies \text{I}(Z; \hat{Y} \mid X_c) > 0 \qquad \text{[Using (3) in Section II-A and non-negativity of PID terms].} \qquad (50)$$

$\square$