

Discrimination Exposed?

On the Reliability of Explanations for Discrimination Detection

JULIAN SKIRZYŃSKI, University of California, San Diego, USA

DAVID DANKS, University of California, San Diego, USA

BERK USTUN, University of California, San Diego, USA

Explanations are often cast as tools to uncover algorithmic discrimination. Given a model, we can explain its predictions to identify the rationale behind the model’s predictions. We can present these explanations to decision subjects to let them contest potentially discriminatory outcomes. We can also present them to auditors to flag biased models. These beliefs – which have motivated rules and regulations surrounding explanation – are founded on inherently unverifiable assumptions. These include assumptions about the causal relationship between the inputs of a model and protected membership, the reliability of explanation to reveal salient information, and the ability of consumers or auditors to use information to make accurate claims about discrimination. In this work, we evaluate the viability of these beliefs under best-case assumptions. We consider a simple task where we can associate each prediction with a ground truth label. We design a user study where we can train participants to detect discrimination using explanations and evaluate the accuracy of claims surrounding explanations. We evaluate detection performance as we control the saliency of proxies of protected attributes, human knowledge about protected class, and their knowledge of causal mechanisms. Our results show that explanations fail to reliably flag unfair predictions and underscore the need for alternative safeguards to detect discrimination.

1 Introduction

Machine learning models are routinely used to automate decisions that affect people – be it to approve a loan [79], an insurance claim [37], or a public service [77]. Over the past decade, it has become clear that deploying models can lead to discrimination, as their predictions or performance can change across *protected attributes* such as sex, age, or race [10, 71]. In applications like lending and hiring, such effects arise inadvertently due to indirect discrimination [72]. This type of discrimination occurs when models exclude protected attributes (e.g., sex) but assign predictions through proxies (e.g., `credit_history`).

Many rules and regulations to protect consumers from discrimination in these sensitive domains revolve around explainability. In effect, multiple jurisdictions reference “discrimination” as a core reason for a “right to an explanation” in “high-risk” applications (e.g., EU [73, 74], Brazil [12], Korea [38] and proposed legislation in the United States [1, 2]). Our reliance on explainability stems from a widely-held belief that explanations can reveal that “*an algorithmic decision is affected by a (legally) protected attribute.*”[78]. In the event that this belief were true, post-hoc explanation methods provide a substantial benefit. Namely, they could safeguard against discrimination in ways that are easy to operationalize [6, 8, 23, 27, 48, 54, 81] – e.g., to audit black-box models without interfering in model development, or to provide decision subjects with information to contest adverse decisions.

Despite explanations being central to enforcing anti-discrimination laws, there is little evidence they can fulfill this function effectively. Simply put, we currently do not know the answers to questions such as “If we provide consumers with an explanation, can they effectively detect proxies?” or “If we ask auditors to check for proxies using explanations, could they effectively detect proxies?” or “How sensitive is this to causal assumptions or access to data?” This is surprising since the right to an explanation in a major consumer application was enacted over fifty years ago [see e.g., the adverse action provision in ECOA 70]. In this case, evidence is lacking because evaluating explanations requires technical validation and usability testing. The algorithms must produce faithful, relevant explanations. Users must be able to understand and utilize them effectively. In discrimination detection tasks, we face yet another barrier as any

claim is subject to assumptions related to chance and causality (e.g., which variable is a proxy, whether it affected a given decision, etc.).

In this paper, we aim to test if explanations can assist humans in detecting discrimination, and characterize the conditions under which this assistance is meaningful. Our goal is to produce evidence to inform policy or compliance – either that we need to consider an alternative mechanism or that we need to impose additional conditions on explanations. Our approach seeks to distill the most basic assumptions behind non-direct discrimination and create a minimal setup that enacts them. We also aim to identify and control for confounding factors and explanation *failure modes* to attribute detection performance directly to the explanations. Our main contributions include:

1. We present a formal model for discrimination detection with explanations. Our model highlights the assumptions we require to evaluate the reliability of claims. We use it to highlight the assumptions and failure modes of relying on explanations to support claims of indirect discrimination.
2. We design a user study to evaluate the reliability of discrimination detection with explanations. Our design provides a sandbox environment for key failure modes related to human interaction and provides full control over our task – a machine-learning model, causal assumptions, and explanations.
3. We conduct controlled human-subject experiments. Our results show that participants fail to perform reliably irrespective of which explanations they see and how much knowledge about the problem they have. By showing that explanations fail to deliver on a simple task, these results stress the need for alternative solutions.

Related Work We study the value of explanations as a safeguard for algorithmic discrimination in domains such as lending and hiring [5, 31, 52]. In these domains, fair treatment requires models to output similar predictions across protected groups (i.e., treatment parity). In practice, models may violate this principle as a result of indirect discrimination via proxy variables [see e.g., 72, for a review]. These issues have motivated an extensive stream of work to detect and mitigate discrimination – e.g., methods to train models that do not discriminate [see e.g., 85], to identify proxies in a third-party audit [see e.g., 4], and to enable reporting group or individual discrimination [21]. Our work formalizes discrimination by adopting a causal notion of fairness [see e.g., 43, 60] - e.g., “would my prediction change if I belonged to a different protected group.” [39]

Our work is related to a stream of research on how humans interact with explanations [see e.g., 9, 14–18, 44, 45, 80, 83]. Many works study if and how explanations impact decision-making [9, 14–19, 34, 44, 45, 47, 76, 86, 87]. Studies on counterfactual explanations [see e.g., 22, 25, 30, 41, 42, 67–69, 82] show marginal improvements in decision-making [22, 49, 50, 75, 82] and debugging model behavior [3, 55, 64]. There is less work on using explanations to assess discrimination. As we discuss, one of the key challenges of studying this question is a mismatch in *scope*. In particular, assessing discrimination involves questions about causality at a population level. In contrast, explanations provide answers about model behavior at the instance level. The few studies on using explanations to detect discrimination at the individual level focus on cases where models use protected characteristics [see e.g., 26, 59]. In these settings, explanations were found to help people spot which predictions are discriminatory. However, studies that concern realistic scenarios where the model might at most use *proxy variables*, reveal a different picture. First, Goyal et al. [35] demonstrate that explanations can perpetuate discrimination as users cannot reliably determine discrimination on the basis of reliance on proxy variables. Second, multiple studies show that explanation properties significantly influence whether people perceive a model as fair [7, 11, 51, 56, 63, 65, 66, 84]. These perceptions vary depending on the prediction task [7], explanation type [11, 51, 63, 84], and information content [65]. Our work bridges these research areas by

examining whether explanations work in the tasks envisioned by regulators, where users need to detect discrimination of individual predictions based on proxies.

2 Framework

We consider a task where (un)fairness involves whether a model’s predictions change based on a *protected attribute* A (e.g., gender). Specifically, we examine whether altering the protected attribute while keeping other features constant would result in different model outputs for individual predictions. We formalize this task through causal relationships between features and outcomes in a directed acyclic graph shown in Fig. 1. The model h is a deterministic function $h : B \times X \rightarrow \hat{Y}$ that predicts an outcome Y (e.g., repayment). B denotes the proxy variable, and X denotes inputs that are independent of the protected attribute (e.g., $X = \text{income}$). The model satisfies two common assumptions:

1. *Indirect Discrimination*. The model does not use the protected attribute as input, but its predictions may change as a result of a variable B (e.g., $B = \text{credit_history}$) that is a *proxy* for the protected attribute. [4, 72]
2. *Business Necessity*. The proxy B can improve predictive accuracy, else the model owner could simply remove it from the list of features [32]

These assumptions are met by the vast majority of models in applications where we would be uncertain about discrimination. In effect, models that used the protected attribute violate treatment disparity [10] as they would assign different predictions to different groups. In cases where the proxy did not improve accuracy, then the model owner could avoid scrutiny by training a model without it.

Characterizing Discrimination We determine the fairness of each feature vector based on a notion of *counterfactual fairness* [43]. Counterfactual fairness measures the likelihood that the prediction for (a, b) would change if we were to swap the protected attribute.

Definition 1. A prediction $\hat{Y} = h(x, b)$ is δ -counterfactually fair if the chances of obtaining this prediction under the current value of the protected attribute $A = a$ and under another value $A = a'$ are at most δ away from each other:

$$\left| \underbrace{\Pr(\hat{Y}_{A \leftarrow a} = h(x, b) \mid X = x, B = b, A = a)}_{\text{Current Prediction where } A = a} - \underbrace{\Pr(\hat{Y}_{A \leftarrow a'} = h(x, b) \mid X = x, B = b, A = a)}_{\text{Counterfactual Prediction when } A = a'} \right| \leq \delta$$

Here, $\hat{Y}_{A \leftarrow a}$ is the current output of the classifier, $\hat{Y}_{A \leftarrow a'}$ is the *potential output* in a *counterfactual* world where we set the protected attribute of the individual to $A = a'$, and $\delta \in [0, 1]$ is a *fairness threshold* that represents the maximum degree to which a prediction could change as a result of this intervention.

We can compute $\Pr(\hat{Y}_{A \leftarrow a} = h(x, b) \mid X = x, B = b, A = a) = \Pr(\hat{Y} = h(x, b) \mid X = x, B = b, A = a) = 1$ since the there is no intervention required, the model is deterministic, and we assume no random effects. We can compute $\Pr(\hat{Y}_{A \leftarrow a'} = h(x, b) \mid X = x, B = b, A = a)$ by setting the protected attribute A to a' and propagating its effect on the

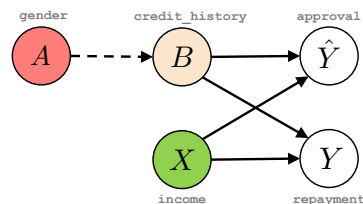


Fig. 1. Causal diagram for discrimination detection. Model $h : B \times X \rightarrow \hat{Y}$ returns prediction \hat{Y} of an outcome variable Y given input proxy B and features X . We seek to determine if model predictions change with respect protected attribute A through its proxy B , which is assumed to be related to the outcome Y . For example, in loan approval predictions (\hat{Y}), the model uses an individual’s income (X) and credit history (B) as inputs. Gender (A) could affect credit history due to differences in credit scores or the intensity of credit usage found between men and women [see e.g. 53].

proxy B . Given the causal structure in Fig. 1, we can express this term as:

$$\Pr(\hat{Y}_{A \leftarrow a'} = h(x, b) \mid X = x, B = b, A = a) = \sum_{b' \in B} \underbrace{\Pr(\hat{Y} = h(x, b) \mid X = x, B = b', A = a')}_{\text{Prediction for } b'} \cdot \underbrace{\Pr(B = b' \mid A = a')}_{\text{Proxy Strength}} \quad (1)$$

As a result, we can write the following corollary:

Corollary 1. A prediction $\hat{Y} = h(x, b)$ is δ -counterfactually fair if the chance it flips, $\phi_{x,b,a}$, as we intervene on the protected attribute $A = a$ and change it to $A = a'$ is at most δ :

$$\phi_{x,b,a} = \underbrace{\left| 1 - \sum_{b' \in B} \Pr(\hat{Y} = h(x, b) \mid X = x, B = b', A = a') \cdot \Pr(B = b' \mid A = a') \right|}_{\text{Chance the prediction flips if we intervene on } A} \leq \delta$$

We can refer to quantity in Definition 1 and Corollary 1 as the *level of discrimination*. The maximum value of discrimination we tolerate is defined by the fairness threshold δ . This threshold can be set on a task-by-task basis. For example, if we are using a model to screen resumes in a job application, then we could set $\delta = 0.2$ to reflect the “4/5ths rule” in U.S. employment discrimination law [28].

Discrimination Detection with Explanations Many rules and regulations, such as the adverse action requirement [61], mandate explanations as an anti-discrimination measure, based on the assumption that they help users identify and contest unfair predictions. Given a particular prediction for individual i , say *Alice*, the user would evaluate if $h(x_{Alice}, b_{Alice})$ flips when a_{Alice} is changed to some a' – and whether that occurrence is likely enough (δ) to state the prediction is discriminatory.

We evaluate such claims by formalizing our problem as a detection task. We associate each instance i with feature vector (x_i, b_i) with two binary labels: (a) a “ground-truth” label that reflects actual discrimination in the prediction; and (b) a “prediction” label that denotes the claim a prediction is discriminatory derived from analyzing it alongside the explanation. Given a model h , we set the ground truth label for prediction $h(x_i, b_i)$ as an indicator the prediction is not δ -counterfactually fair:

$$g_i | h, \delta := \mathbb{I}[\phi_{x_i, b_i, a_i} > \delta] \quad (2)$$

We let the prediction label $\hat{g}_i | h, \mathcal{E}_i$ denote user’s claim about discrimination upon seeing explanation \mathcal{E}_i . In what follows, we write $g_i := g_i | h, \delta$, $\hat{g}_i := \hat{g}_i | h, \mathcal{E}_i$, and $\phi_i := \phi_{x_i, b_i, a_i}$ when their dependencies are clear from context.

The chance a given prediction flips is a number. This chance is fixed for every individual with features (x, b, a) . Whether the prediction actually flips if we intervene on the protected attribute is random. The results may be different for Alice, for Bob, for Charlie, etc., even if $(x_{Alice}, b_{Alice}, a_{Alice}) = (x_{Bob}, b_{Bob}, a_{Bob}) = \dots$. In a hypothetical process where we draw random samples that indicate if the prediction flips or not – $G_i \sim \text{Bernoulli}(\phi_{x, b, a})$ – the expected proportion of samples that indicate a flip would converge to our ground truth discrimination label g_i . Therefore, we can interpret g_i in terms of hypothetical proportions. Given a set of N individuals with features x and proxies \hat{b}_i , where each \hat{b}_i is drawn based on $A = a'_i$, a model would that is δ -counterfactually fair would yield a different prediction for δN individuals.

We expect users to approximate the reasoning about hypothetical proportions above. However, since they only see a particular prediction for instance i , we interpret $\hat{g}_i | h, \mathcal{E}_i$ as their willingness to bet that this prediction changes under an intervention on A . If $\hat{g}_i | h, \mathcal{E}_i = 1$ then their personal probability that the prediction flips exceeded their internal

threshold $\delta^{internal}$, and they are willing to make a bet [see e.g., 24, for more details about this interpretation].¹ We write this as:

$$\hat{g}_i|h, \mathcal{E}_i \approx \mathbb{I}[\phi_{x,b,a} > \delta^{internal}]$$

Measures We evaluate detection performance through standard performance measures for binary classification that vary as a function of the fairness threshold δ_{\min} , namely: PPV(δ_{\min}), which indicates the internal reliability of discrimination claims; TPR(δ_{\min}), which measures how often participants correctly identify discriminatory cases; and FPR(δ_{\min}), which tracks false alarms on fair cases. These metrics are computed by comparing participant’s claims $\hat{g}_i|h, \mathcal{E}_i$ to ground truth labels $g_i|h, \delta_{\min}$ over all values of $\delta_{\min} \in [0, 1]$. We expect the following:

- Instance-level Detection: Explanations are a perfect mechanism to support individual claims when the claims are aligned with ground-truth labels. In this case, we should have that $\hat{g}_i|h, \mathcal{E}_i = g_i|h, \delta_{\min}$ for *some* fairness threshold δ_{\min} , likely individualized per user and equal to $\delta^{internal}$, and any explanation \mathcal{E}_i . Under such a threshold, the claims should neither be selective and miss discrimination nor be overly sensitive and raise false alarms (see Eq. (No Missed Signals) and Eq. (No False Alarms)). Given a model h , and a set of n individuals $S = \{(x_i, b_i)\}_{i=0}^n$, we can evaluate the reliability of discrimination claims by reporting the empirical PPV, TPR and FPR. It is desired they reach perfect scores of PPV = 100%, TPR = 100%, and FPR = 0%. In reality, we would say explanations help detect fairness if they are high enough, e.g., 90% for PPV and TPR and 0% for FPR.
- Model-level Detection: If explanations work at the instance level, they can also support discrimination determinations at the model level. Perfect individual-level detection means auditors can accurately calculate the proportion of discriminatory predictions, and check if it exceeds a model-level threshold τ_h (preferably matching δ_{\min} for consistency). Notably, people may perform this comparison even if individual detection is imperfect. It is sufficient to estimate if the model discriminates more often than $\tau_h\%$ of the time or not. A model that clearly discriminates can tolerate many false alarms while still being correctly identified as discriminatory. Conversely, a clearly fair model can withstand some missed discriminatory cases. The closer the true discrimination rate is to τ_h , the more precise individual detection needs to be. For this use case it suffices users can distinguish between fair models and discriminatory models.

$$\mathbb{E}[\hat{g}_i = 1 \mid g_i = 1] = 100\% \quad (\text{No Missed Signals})$$

$$\mathbb{E}[\hat{g}_i = 1 \mid g_i = 0] = 0\% \quad (\text{No False Alarms})$$

Failure Modes Users may fail to detect discrimination with explanations. These failures can be the result of multiple false beliefs they hold or issues with explanations themselves. Here, we list these failure modes explicitly so that we could later control for them in our study design. Given model h and an explanation method, the user may claim $\hat{g}_i \neq g_i$ because:

REMARK 1 (RECOVERY). *There exist many explanations $\mathcal{E}_i, \mathcal{E}'_i$ such that $\mathcal{E}_i \neq \mathcal{E}'_i$ for the same prediction $h(x_i, b_i)$ [13, 40]. Specifically, explanation \mathcal{E}_i may fail to reveal $h(x_i, b_i) \neq h(x_i, b'_i)$ for some $b'_i \neq b_i$. Explanation \mathcal{E}'_i could reveal this dependence. If the user is shown explanation \mathcal{E}_i , they might conclude that $\forall b'_i \neq b_i, h(x_i, b_i) = h(x_i, b'_i)$ leading to the erroneous determination that $\Pr(\hat{Y}_{A \leftarrow a'_i} = h(x_i, b_i) \mid x_i, b_i, a_i) = 1$ and hence $\phi_{x,b,a} = 0$ and $\hat{g}_i = 0$ for any $\delta > 0$. In reality, we could have $\Pr(\hat{Y}_{A \leftarrow a'_i} = h(x_i, b_i) \mid x_i, b_i, a_i) < 1$ and $g_i = 1$ for some low enough δ .*

¹If one prefers a different interpretation of probability statements, then $\hat{g}_i|h, \mathcal{E}_i$ can be reinterpreted; for example, $\hat{g}_i = 1$ could be understood as indicating a sufficiently large change in subjective strength of belief.

REMARK 2 (MISINTERPRETATION). Given \mathcal{E}_i that clearly conveys $h(x_i, b_i) \neq h(x_i, b'_i)$, users may still not know they should compute $\phi_{x,b,a}$ to obtain \hat{g}_i . Their claims may effectively become random with respect to g_i .

REMARK 3 (MISSPECIFIED BELIEFS ABOUT CAUSAL MECHANISM). User may operate using a different probability function $Pr_{user}()$ such that they incorrectly estimate the causal relationship in $\phi_{x,b,a}$, because they misrepresent the proxy strength $Pr_{user}(B | A) \neq Pr(B | A)$. Using $Pr_{user}()$ to compute $\phi_{x,b,a}$ leads to inaccurate result due to Eq. (1). For some δ this may lead to an incorrect determination that $\phi_{x,b,a} > \delta$ and $\hat{g}_i = 1$, when in reality $\phi_{x,b,a} \leq \delta$ and $g_i = 0$. It may also lead to a determination that $\phi_{x,b,a} \leq \delta$ and $\hat{g}_i = 0$, when in reality $\phi_{x,b,a} > \delta$ and $g_i = 1$.

REMARK 4 (KNOWLEDGE OF PROTECTED CLASS). Let a_i be the true protected attribute value for instance i . The user may incorrectly assume the protected attribute value is $a'_i \neq a_i$. This may lead to an error in computing $\phi_{x,b,a}$. This is because $Pr(\hat{Y}_{A \leftarrow a'_i} = h(x_i, b_i) | x_i, b_i, a_i) \neq Pr(\hat{Y}_{A \leftarrow a_i} = h(x_i, b_i) | x_i, b_i, a'_i)$. The computed $\phi_{x,b,a}$ may be too high or too low. If it's too high and we fix δ to be high enough, this may in turn lead to an incorrect determination that $\phi_{x,b,a} > \delta$ and $\hat{g}_i = 1$, when in reality $\phi_{x,b,a} \leq \delta$ and $g_i = 0$. If the computed $\phi_{x,b,a}$ is too low and we fix δ to be low enough, this may also lead to a determination that $\phi_{x,b,a} \leq \delta$ and $\hat{g}_i = 0$, when in reality $\phi_{x,b,a} > \delta$ and $g_i = 1$.

REMARK 5 (MISSPECIFIED BELIEFS ABOUT CAUSAL STRUCTURE). Let $\mathcal{G} = (V, E)$ be the true causal DAG from Fig. 1 where $V = \{A, B, X, Y\}$ and E includes edges (A, B) and (B, Y) but not (A, X) . Individuals may believe in an incorrect causal DAG $\mathcal{G}_{user} = (V, E_{user})$. If they believe $(A, B) \notin E_{user}$, this may lead them to conclude that $Pr(\hat{Y}_{A \leftarrow a'_i} = h(x_i, b_i) | x_i, b_i, a_i)$ only depends on $Pr(B)$ based on Eq. (1) and hence $Pr(\hat{Y}_{A \leftarrow a'_i} = h(x_i, b_i) | x_i, b_i, a_i) = Pr(\hat{Y}_{A \leftarrow a_i} = h(x_i, b_i) | x_i, b_i, a_i)$ for any $a_i \neq a'_i$. If so, $\phi_{x,b,a} = c$ for some $c \in [0, 1]$. In reality, $\phi_{x,b,a}$ differs based on the value of $a_i \in A$ since $Pr(B = b_i | A = a_i) \neq Pr(B = b_i \text{ mid } A = a'_i)$ for $a_i \neq a'_i$. Then, users may conclude $c > \delta$ for some δ and $\hat{g}_i = 1$ when in reality $\phi_{x,b,a} < \delta < c$ and $g_i = 0$. These inequalities could also be inverted. They may also believe $(A, X) \in E_{user}$ and incorrectly derive the formula for $Pr(\hat{Y}_{A \leftarrow a_i} = h(x_i, b_i) | x_i, b_i, a_i)$. As a result, they may misattribute discrimination, e.g., when E_i shows $h(x_i, b_i) \neq h(x'_i, b_i)$ for some $x'_i \neq x_i$ and they assess $\phi_{x,b,a} > \delta$ when in reality $\phi_{x,b,a} = 0$.

These failure modes are the reason it is challenging to make falsifiable claims about explanations. Each time we may try to conclude explanations fail, it is possible to attribute the failure to one of the listed causes. Therefore, our experimental design focused on minimizing these confounds. Our task design and technical solutions address REMARK 1 and REMARK 2. The remaining three failure modes stem from human reasoning and prior beliefs. They will change across individuals and auditors. To understand their impact, we thus ran an empirical study over multiple participants.

3 Experimental Design

We describe an experimental task to evaluate the reliability of explanations as a tool for aiding discrimination detection. We consider a simple task where: (1) we can endow participants with the skills that we expect from auditors and verify their understanding through comprehension checks; (2) we can manipulate and elicit participant's beliefs in the causal model from Fig. 1; (3) we can collect data to evaluate fairness under different assumptions and use cases (e.g., for all $\delta_{\min} \in [0, 1]$, with or without access to protected attributes, fitting the proxy strength $Pr(B | A)$ to match participant's beliefs $Pr_{user}(B | A)$).

3.1 Robot Classification Task

We consider a task where participants are asked to audit a model that predicts the reliability of fictional robots for NASA. The model was created to inform NASA's purchasing decisions by identifying which robots are reliable versus defective.

While robot reliability is determined by their body parts, the two manufacturers, Company X and Company S, design their robots with slightly different components. This difference could lead to discrimination in the model’s predictions with respect to the manufacturing company. Since NASA is legally prohibited from making decisions based on the company, participants must determine if the model’s predictions are inadvertently discriminatory or not.

We cast the identity of the company as our protected attribute A . We assume that the model predicts that a robot is reliable using a set of four salient characteristics shown in Fig. 2, namely: Antenna, HeadShape, BodyShape, BaseType. We represent the input variables as:

$$B := \mathbb{I}[\text{Antenna} = \text{Yes}]$$

$$X_1 := \mathbb{I}[\text{HeadShape} = \text{Round}]$$

$$X_2 := \mathbb{I}[\text{BodyShape} = \text{Round}]$$

$$X_3 := \mathbb{I}[\text{BaseType} = \text{Wheels}]$$

In this setup, we have $2^4 = 16$ distinct combinations of input variables (B, X) , and 32 distinct robots (A, B, X) . We control all quantities that affect the degree of discrimination by specifying the model’s predictions at each input and the prevalence of each robot. We include the table with this data in Table 2.

We can arbitrarily increase the number of distinct robots to show participants by introducing spurious features, such as $\text{Paint} \in (\text{Red}, \text{Blue})$. In this way, we can ensure that participants are shown new kinds of robots. This is crucial for three reasons: it prevents learning effects from seeing the same robot multiple times, ensures decisions are based on feature relationships rather than memorized patterns, and better simulates real-world auditing where each case presents unique characteristics.

We predict the reliability of each robot using a linear classification model that outputs “Reliable” whenever the robot has an Antenna and one of the following conditions: a Round HeadShape, a Round BodyShape, or Wheels:

$$h(B, X) = \text{sign}(6B + 4X_1 + 4X_2 + 3X_3 - 8)$$

$$\hat{Y} = \mathbb{I}[B \text{ AND } (X_1 \text{ OR } X_2 \text{ OR } X_3)].$$

The reliability of each robot Y is generated by the random process, under which the model reaches an accuracy of 88%:

$$A, X_1, X_2, X_3 \sim \text{Bernoulli}(0.5)$$

$$B | A \sim \text{Bernoulli}(p_{B|A}) \quad \text{where } p_{B|A} \text{ is a set in Table 1}$$

$$Y \sim \text{Logistic}(B + X_1 + X_2 + X_3).$$

3.2 Discrimination

Under the causal model and features we defined in our task, predictions have at most three levels of discrimination $\phi_{x,b,a}$. Computing $\phi_{x,b,a}$ (see Corollary 1) requires knowing the value of $\Pr(\hat{Y} = h(x, b) | x, \hat{b}, a')$ under all $\hat{b} \in B$. Since $b \in \{0, 1\}$, \hat{b} can be written as either b or $1 - b$ without the loss of generality. If $1 - b$ does not flip the prediction, both distributions in ?? are 1, and the prediction is perfectly fair for any δ . For our variables, this occurs for instances (x_i, b) where $x_i \in \{(0, 0, 0), (1, 1, 0), (1, 1, 1)\}$. Otherwise, $\phi_{x,b,a} = 1 - \Pr(B = b | A = a')$ which shows that the level of

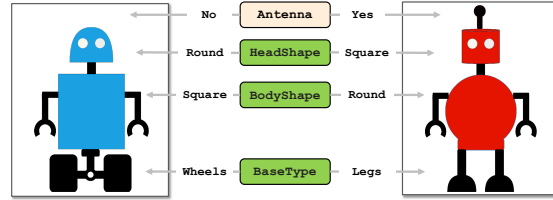


Fig. 2. Overview of robot characteristics. We show two robots to cover all possible values of each characteristic. Our model predicts that each robot is reliable or defective using dummy variables $B = \mathbb{I}[\text{Antenna} = \text{Yes}]$, $X_1 = \mathbb{I}[\text{HeadShape} = \text{Round}]$, $X_2 = \mathbb{I}[\text{BodyShape} = \text{Round}]$ and $X_3 = \mathbb{I}[\text{BaseType} = \text{Wheels}]$.

discrimination depends solely on $\Pr(B | A)$. We vary the strength of this relationship across three regimes (see Table 1) to evaluate how proxy strength affects discrimination detection and claims $\hat{g}_i|_{h,\delta}$. This variation is crucial because real-world proxies range from weak correlations (e.g., zip codes as proxies for race) to almost perfect proxies (e.g., height as a proxy for gender). By testing different proxy strengths, we can assess whether participants’ performance varies with proxy obviousness. In what follows, we also remain agnostic about the value of δ and evaluate the potential to detect discrimination over all possible thresholds $\delta_{\min} \in [0, 1]$.

3.3 Explanations

To detect discrimination of classifier h on instance (x_i, b_i) and provide \hat{g}_i , the user needs to compute the discrimination level ϕ_{x_i, b_i, a_i} and compare it to their fairness threshold $\delta^{internal}$.

Regime	Proxy Strength		Discrimination Level			
	$A = 0$	$A = 1$	$A = 0, h(x, b) \neq h(x, 1 - b)$	$A = 1, h(x, b) \neq h(x, 1 - b)$	$A = a, h(x, b) = h(x, 1 - b)$	
Weak	5%	10%		10%	5%	0%
Medium	5%	55%		55%	45%	0%
Strong	5%	95%		95%	90%	0%

Table 1. Overview of parameters determining discrimination claims under each proxy regime. Here, Proxy strength denotes $\Pr(B = 1 | A)$, whereas Discrimination level lists possible values of $\phi_{x, b, a}$.

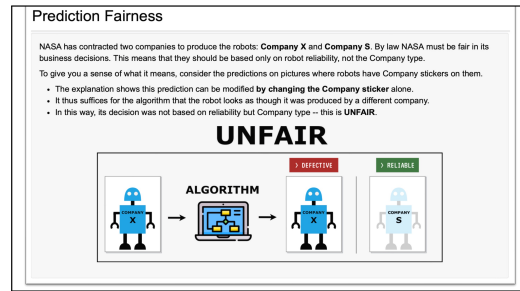
Under the correct assumptions about the probabilities and the causal structure of the problem, evaluating ϕ_{x_i, b_i, a_i} entails knowing if $h(x_i, b_i) = h(x_i, 1 - b_i)$ (because of the the formula Corollary 1). Explanations could potentially reveal this information. We test if they do by manipulating the type of explanations we show to participants, i.e., either \mathcal{E}_i that uses b and could provide insight into whether $h(x_i, b_i) = h(x_i, 1 - b_i)$ or \mathcal{E}'_i that does not use b , and provides no insight about it. In this way, we address REMARK 1 and test if Recovery could affect detection. Our design also addresses REMARK 2, that assumes that users might not know they should assess ϕ_{x_i, b_i, a_i} based on the provided \mathcal{E}_i . We do so by training participants on how to use the explanation method to detect discriminatory examples in cases when the model uses protected attributes. Specifically, participants are showed examples of predictions and explanations that are discriminatory and told that this is because the predictions rely on the information about the protected attribute (the company). Then, they take a comprehension check where they label predictions and explanations that contain or not contain the protected attribute as fair or not.

3.4 Procedure

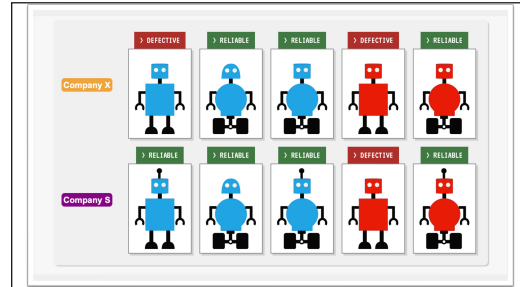
We implemented our task into an online user study where we can evaluate how well participants can detect discrimination using explanations. Our study consists of four phases shown in Fig. 3. This implementation is in principle explanation-agnostic and may be adapted to any explanation method by changing the instructions and the visual materials.

Our overall design offers key advantages that address the remaining belief-related failure modes. First, it allows post-hoc evaluation across different fairness thresholds δ_{\min} or different probabilities in the causal structure of the problem from Fig. 1 without additional experiments. We can compare participant claims \hat{g}_i to $g_i|_{h, \delta_{\min}}$ for any δ_{\min} and any underlying probabilities. This property means we only need to elicit participants’ beliefs about the prediction problem to recompute these probabilities. We accomplish this by directly measuring participants’ beliefs about proxy strength (REMARK 3) and protected attributes (REMARK 4). We also assess the impact of false causal beliefs (REMARK 5) by comparing claims \hat{g}_i to g_i in the most beneficial scenario, where we assume participants have both the correct knowledge about protected attributes and the causal mechanism.

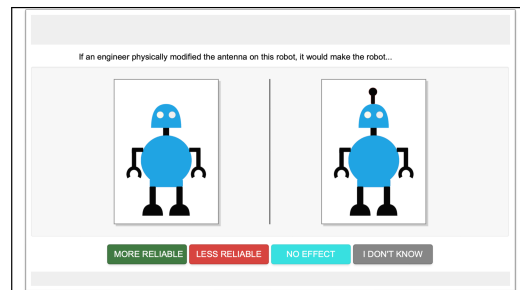
1. **Training:** Participants were introduced to four key elements of the study: robots, their components, a reliability prediction model, and the concept of discrimination. We used counterfactual explanations as the explanation method and presented them visually by highlighting modifiable robot parts. To explain discrimination, we used examples of robots with company stickers, establishing that predictions based on manufacturer identity were illegal. Participants completed a screening test where predictions were either discriminatory because they could be changed with company stickers or fair because they depended on robot parts. Participants then had three attempts to pass a comprehension quiz or were otherwise dropped from the study.



2. **Anchoring:** We presented participants with a set of robots to anchor their beliefs on the strength of the proxy and its impact on reliability. Each participant saw 10 robots from each company. We arranged the robots so that robots from each company shared the same features. We then assigned reliability labels and antenna to robots to anchor their beliefs on the impact of the proxy. The set contained two defective robots from Company X and one from Company S. All robots from Company X had no antenna. Company S had 1/3/5 robots with antennas depending on the regime. Participants were explicitly told which feature distinguished the sets and were informed that proxy-based predictions *could* be discriminatory since the antenna can behave like the company sticker.



3. **Elicitation:** Participants were elicited for their beliefs on the protected class c_i and the effect of the proxy u_i on each possible robot. Participants saw a total of 16 robots for (X, B) . We elicited beliefs regarding protected class by asking them to predict its manufacturer or state they don't know. We coded these as $\{0, 1, ?\}$. We elicited beliefs regarding the impact of the proxy by asking them how adding (or removing) an antenna from the robot would change reliability, allowing them to answer (more, less, no effect, or unknown), coded as $\{-1, 0, 1, ?\}$. Given the participant's beliefs in the protected attribute, we could recompute performance with participant-assumed attributes; we could also estimate $\Pr_{\text{user}}(B | A)$ to match their belief in the causal mechanism of the proxy. This allowed us to address REMARK 4 and REMARK 3. By storing reliability beliefs, we could analyze if these beliefs affect discrimination claims.



4. **Auditing:** Participants judged if predictions were discriminatory. They were shown an image of a robot, its prediction (always Defective), and one or more of the closest counterfactuals. The participant aimed to select whether the prediction was fair or unfair. This phase consisted of 16 rounds with all seven unique defective robots shown in different colors (2 robots appeared twice). We collected discrimination claims $\hat{g}_i \in \{0, 1\}$.

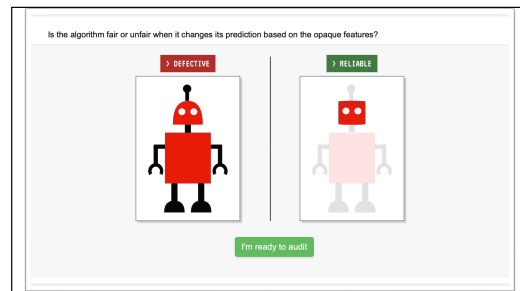


Fig. 3. All four phases of our experiment with their description.

4 Experiment

Our experiment sought to characterize the viability and effectiveness of explanations in detecting algorithmic discrimination. In particular, we sought to determine if individuals could use explanations to make reliable discrimination claims across use cases in consumer protection. Our specific research questions include:

- RQ1** Can participants use explanations to make reliable claims for discrimination at an instance level? If so, this would suggest that explanations are an effective mechanism to exercise individual rights (e.g., to contest predictions that are unfair).
- RQ2** Can participants who are shown explanations make reliable claims for discrimination at a model level? If so, this would suggest that explanations could serve as an effective mechanism to audit models.
- RQ3** How does the reliability of claims depend on the information that is available to participants? In particular, explanations may be a viable mechanism only in use cases where participants have perfect information on the protected attributes of each instance (e.g., in a third-party audit).
- RQ4** How does the reliability of claims depend on the correctness of causal assumptions (e.g., does the strength of the proxy match their beliefs)? In particular, explanations may be a viable mechanism only in settings where participants have correct beliefs about the strength of the proxy variable .
- RQ5** How does the reliability of detection change if we could provide participants with multiple explanations for each prediction? If so, this would speak to the importance of diverse explanations [see, e.g., 58]
- RQ6** Do participants behave in ways that are consistent and predictable? For example, will participants in each experiment make identical claims? In this case, inconsistency would highlight a need for standardization.

4.1 Setup

We used a study design with $2 \times 3 = 6$ conditions in which we varied the strength of the proxy variable $\in \{\text{Weak Proxy, Medium Proxy, Strong Proxy}\}$ and the format of counterfactual explanations $\in \{\text{Single, Multiple}\}$.

1. **Single:** Participants were shown a single explanation for each prediction. This mimics real-world scenarios where participants might be given “the best explanation” or just *some* explanation and need to decide about discrimination. In this setup, an explanation might show no dependence on the proxy, but the prediction could still heavily rely on it, making it potentially discriminatory.
2. **Multiple:** Participants were presented with two competing explanations for each prediction, with one explanation always containing the proxy variable when it existed. This setup represents a scenario with maximum insight into the model’s decision-making process. In this setup, the participants know exactly which predictions depend on the proxy and are potentially discriminatory.

Participants in each condition were shown a different set of robots to anchor their beliefs on proxy strength. The sets differed by the number of robots in `Company S` that had antennas: 1 robot for the Weak Proxy conditions, 3 robots for the Medium Proxy conditions and all five robots for the Strong Proxy conditions. Our evaluation also considered different levels of knowledge in the task:

1. **Auditor Baseline:** Participants have no information about the true protected attributes and estimate the distribution of the proxy based on the anchoring robot set. This is a realistic assumption where the protected attributes are not readily available, and auditors have internal estimates of the true distributions.
2. **Known Protected Attribute:** Participants have perfect information about the protected attributes according to their elicited beliefs. This maps to an information regime where the auditor has access to the protected attributes (e.g., filing claims from consumers, or a third-party audit where the protected attributes are stored according to the law, such as audits (in New York) of employment decisions [36]).

3. **Known Causal Mechanism:** Participants have perfect information about the causal mechanism, i.e., the conditional distribution of the proxy matches their elicited beliefs and we set $\Pr(B | A) = \Pr_{\text{user}}(B | A)$. This is an idealized assumption and allows us to estimate best-case performance.

Counterfactual Explanations Participants in our task audit discrimination using *counterfactual explanations*. A *counterfactual explanation* (CE) returns a set of changes to the input features that result in a different prediction. For example, when a loan application is denied, a CE might state “If your income were \$5,000 higher and credit_history was 2 years longer, the loan would be approved.” Given a classifier $f : \mathcal{X} \rightarrow \{0, 1\}$ that assigns a prediction $f(x) = 0$ for $x = (x_1, \dots, x_n)$, a CE specifies a set of feature changes $\mathcal{E}(x) = \{(i, x_i^{cf}) : x_i^{cf} \neq x_i\}$ that define a new instance $x^{cf} = (x_1^{cf}, \dots, x_n^{cf})$ where: $x_j^{cf} = \begin{cases} x_i^{cf} & \text{if } (j, x_j^{cf}) \in \mathcal{E}(x) \\ x_j & \text{otherwise} \end{cases}$ such that $f(x^{cf}) = 1$. When the set is minimal, we say

that $\mathcal{E}(x)$ is a *closest counterfactual*. Given our task, the possible closest CEs for an arbitrary prediction $h(x, b)$ span the following cases: $\mathcal{E}(x, b) = \{1 - b\}$, $\mathcal{E}(x, b) = \{1 - x_i\}$, $\mathcal{E}(x, b) = \{1 - x_i, 1 - x_j\}$, and $\mathcal{E}(x, b) = \{1 - b, 1 - x_i\}$, where $x_i, x_j \in x, i \neq j$. We write \mathcal{E}_k to denote the CE for a particular instance (x_k, b_k) .

Our interest in counterfactual explanations stems from three main benefits. First, they directly relate to participant claims \hat{g}_i , and the fact they involve evaluating $\phi_{x,b,a}$ because they list the exact changes needed to flip the prediction. Second, CEs are easy to grasp as we can present them visually (e.g., by highlighting feature changes). Third, we can confirm that participants understand their guarantees and limitations (e.g., via a comprehension quiz). These tasks are far more difficult to achieve when, for example, we explain predictions using a feature attribution method in which guarantees are ambiguous and prone to misinterpretation [29].

Procedure We recruited 126 participants through Prolific (20-23 per condition). All participants were fluent English speakers from the United States, comprising 74 females and 52 males, age 19-74 (mean = 35). Each experiment lasted 32 minutes on average. We assigned each participant to one of the 6 conditions. Participant beliefs about the proxy strength were established through exposure to a different set of robots in the Anchoring stage. Participants also received a set of instructions on the guarantees for the explanation. Our experiment was designed to ensure that all information shown to participants (e.g., model accuracy, anchoring set) was aligned with the distributional assumptions in each regime.

We included a set of comprehension questions prior to the Auditing phase. Participants who failed this quiz three or more times were excluded from the study (10 excluded participants; exclusion rate of 8%). These quizzes ensured that participants understood how to apply each explanation and its guarantees with respect to discrimination claims.

4.2 Results

Overall, our results show that participants cannot reliably detect discrimination with explanations under any setup. The summary performance measurements of audits where participants were asked to flag discriminatory predictions based on a single explanation can be found in Fig. 6.

On the Reliability of Discrimination Detection We first consider a setting with threshold $\delta_{\min} = 0.2$ - i.e., where we wish to flag predictions that would change by over 20% given an intervention on protected group membership - given its importance in U.S. employment law [36].

As seen in Fig. 4, PPV, a measure of reliability of participant claims, indicates poor detection performance across all tested conditions. We would expect perfect, or at least very high PPV, say $\approx 90\%$, meaning that participants’ detection is generally trustworthy. To the contrary, we observe that even in the Strong Proxy condition, where the proxy was

the easiest to spot and its presence in the explanation most often indicated discrimination, PPV was as low as $48\% \pm 4\%$ (see the blue boxes in ??). It was even lower, $28\% \pm 6\%$ in the Medium Proxy condition to hit 0% in the Weak Proxy condition where all predictions were fair at $\delta_{\min} = 0.2$. This means that participants were correct in at most *half* of their discrimination claims. Further analysis revealed that this low reliability was affected by both missing most of the discriminatory predictions, and flagging fair predictions. In the Strong Proxy condition where the results were the best, TPR reached only $44\% \pm 5\%$ while maintaining substantial FPR ($33\% \pm 5\%$). This means that participants incorrectly flagged 2-3 fair predictions. They also missed at least 3 out of 5 all discriminatory predictions.

These results raise concerns about using explanations for discrimination auditing in practice. Without additional assumptions or safeguards, humans both fail to detect most of discriminatory cases, and raise multiple false alarms. This combination risks letting discriminatory practices continue and triggering unnecessary investigations that waste resources and potentially harm legitimate practices.

This poor performance is not due to the particular fairness threshold we selected. As seen in the blue line in Fig. 6, poor performance is observed systematically for all measures and almost all thresholds. This changes only at extreme values. For sufficiently high thresholds, all predictions become fair and since participants did claim discrimination, their performance drops. Conversely, at very low thresholds ($\delta_{\min} \leq 5\%$ that exemplify a “better safe than sorry” approach), most proxy-dependent predictions are discriminatory. Since participants tend to flag these predictions, they achieve high PPV ($\approx 75\%$) but still maintain poor TPR and FPR of $\approx 30\%$.

On the Sensitivity to Protected Attributes A natural question is whether the poor detection performance stems from a lack of knowledge of protected attributes. Perhaps participants reasoned about the hypothetical predictions under wrong assumptions. To answer this question, we matched participants’ attribute selections from the Elicitation phase with the corresponding predictions. This setup reflects common scenarios where an auditor is given information on the protected attribute of claims (e.g., a third-party audit processing consumer complaints).

Our results (see Fig. 5) show only marginal improvements: at $\delta_{\min} = 0.2$, PPV increased to $39\% \pm 6\%$ (Weak Proxy condition) and $37\% \pm 3\%$ (Medium Proxy condition) from the baseline of 28%, with neither change reaching significance under Mann-Whitney U test ($p > 0.1$, $U \geq 156.5$). Only the Strong Proxy condition showed significant improvement, with PPV rising to $66\% \pm 7\%$ from $48\% \pm 7\%$ ($p < 0.05$, $U = 114.5$). We found similarly slight improvements for other measures: FPR dropped by approximately 10% (equivalent to ≈ 1 prediction), and TPR decreased by 6-7%, both across all conditions. This suggests that participants sometimes chose not to flag discrimination even when their own beliefs about protected attributes would warrant it. This often occurred when participants believed changing the proxy has legitimate influence on reliability – e.g., on average, if participant believed the change in the CE affects robot reliability, they claimed the prediction is fair in 64% of the cases whereas if they thought it has no effect – in 50% of the cases.

In total, knowledge of the protected attributes played a marginal role in detection performance. Even with access to these attributes, auditors still missed many discriminatory cases and raised multiple false alarms. As seen in the

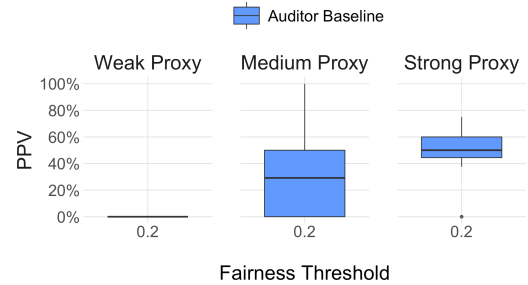


Fig. 4. Distribution of the Positive Predictive Value (PPV) at threshold $\delta_{\min} = 0.2$ used in U.S. employment law [28] across all proxy strength conditions assuming the ground truth probabilities and causal mechanism of the proxy (Auditor Baseline). As shown, PPV is low, reaching the highest value in the Strong Proxy condition where it does not exceed $\approx 50\%$. It falls below 40% in the Medium Proxy condition until it hits 0% in the Weak Proxy condition.

green lines in Fig. 6, this performance persisted across all δ_{\min} values, except for very low thresholds where most proxy-dependent predictions were discriminatory. In these cases, participants correctly focused on such predictions, leading to higher PPV (most claims were accurate), though their overall detection ability remained poor (low TPR and high FPR).

On the Sensitivity to Causal Assumptions Our experiment also allows us to evaluate how performance would improve under best-case assumptions where humans have perfect information on the causal mechanism of the proxy. In this case, we assume $\Pr(B | A)$ matches their beliefs $\Pr_{\text{user}}(B | A)$. We found that this intervention significantly improved PPV at $\delta_{\min} = 0.2$ across all conditions, as seen in green in Fig. 5. In the Strong Proxy condition, PPV went from $48\% \pm 4\%$ to $77\% \pm 7\%$ ($p < 0.001$, $U = 66.5$). In the Medium Proxy condition it went from $28\% \pm 6\%$ to $49\% \pm 8\%$ ($p \leq 0.05$, $U = 128.5$). In the Weak Proxy condition, PPV increased significantly above 0 to $61\% \pm 8\%$. This is because participants perceived a stronger proxy relationship than existed (over half of the participants assumed $\Pr(B = 0 | A = 0) = 0$), and their discrimination claims were often warranted under these beliefs.

Besides such local improvements, however, neither PPV nor TPR/FPR ever reached a value we would consider satisfactory, as seen in Fig. 6. Overall, these results point to the fact that the lack of poor performance cannot readily be remedied by domain expertise.

On the Effect of Multiple Explanations We next examined participants' performance when they were given full information about the prediction by being shown Multiple explanations. In this setup, they knew with certainty whether the prediction can be flipped with the proxy or not. Such guarantees are rarely available in reality, but we make this assumption to test if explanations *could* work in idealized circumstances.

In short, this manipulation did not lead to good performance for any $\delta_{\min} \in [0.05, 1]$ as we show in the Appendix in Fig. 8. On average, PPV was bounded by 40% across all conditions. TPR behaved irregularly but never exceeded 40%. FPR remained consistently at least 30%. The only exception occurred in the Weak Proxy condition with extreme values of $\delta_{\min} \leq 0.05$. As in this case many predictions were deemed discriminatory despite low discrimination levels $\phi_{x,b,a}$, participants detection was more reliable, with PPV reaching $77\% \pm 7\%$ and TPR $63\% \pm 9\%$ ($p < 0.01$, $U \geq 220$). However, this came at the cost of increased false positives (FPR as high as $55\% \pm 8\%$ at $\delta_{\min} = 0.2$). These results hold irrespective of the level of knowledge participants have, i.e., no knowledge (baseline), knowledge of protected attributes or knowledge about the causal mechanism of the proxy. Overall, people appear to be incapable of using explanations reliably even under idealized knowledge conditions.

On Model Audits Participants were unable to differentiate between cases when the model we tested was fair versus discriminatory. According to our design, and assuming the 20% model threshold τ_h from U.S. employment law [28], the model was fair in the Weak Proxy conditions (unless the discrimination threshold for individual predictions was very

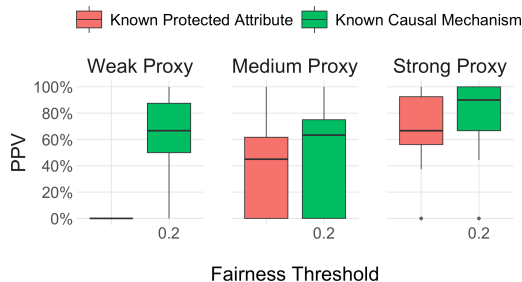


Fig. 5. Distribution of the Positive Predictive Value (PPV) at threshold $\delta_{\min} = 0.2$ used in U.S. employment law [28] across all proxy strength conditions and under different assumptions on participant knowledge: known protected attributes (red), and known causal mechanism (green). PPV is unsatisfactory even with additional knowledge. Knowing protected attributes yields a PPV of $\approx 70\%$ in the Strong Proxy condition, $\approx 50\%$ in the Medium Proxy condition, and 0% in the Weak Proxy condition. Only assuming auditors' beliefs about the causal mechanism shows PPV in the desired range of $\approx 90\%$ – but only in the Strong Proxy condition. This manipulation renders $PPV < 70\%$ in the remaining two conditions.

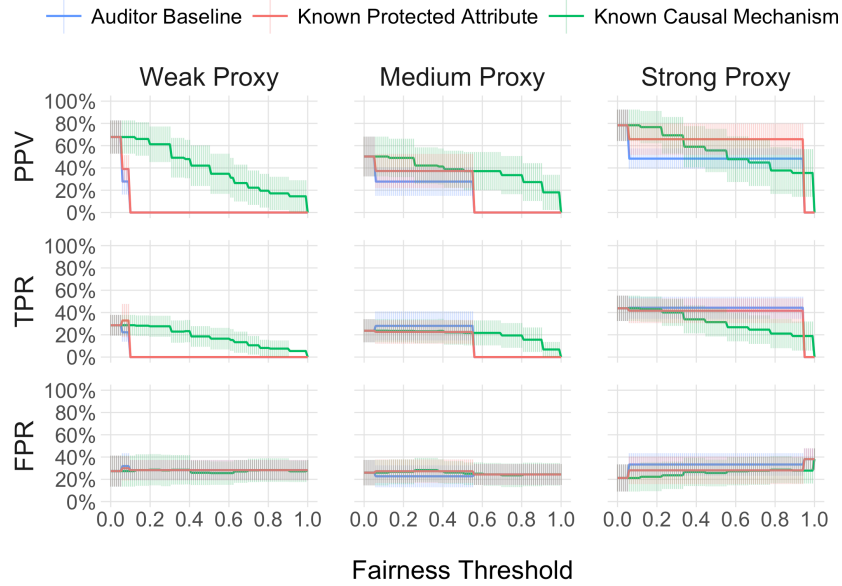


Fig. 6. Reliability of discrimination claims across all possible $\delta_{\min} \in [0, 1]$ (right). We show the confidence intervals for $PPV(\delta_{\min})$, $TPR(\delta_{\min})$ and $FPR(\delta_{\min})$ across all proxy strength conditions and under different assumptions on participant knowledge: baseline performance (blue), known protected attributes (red), and known causal mechanism (green). As shown, baseline performance is poor across all thresholds, with PPV and TPR rarely exceeding 50%. Knowledge of protected attributes yields only marginal improvements. Assuming auditors’ beliefs about the causal mechanism increases PPV to at most $\approx 70\%$ but has no effect on TPR ($\leq 40\%$). Notably, FPR remains problematic (around 30%) across all conditions and knowledge regimes.

low, e.g., $\delta_{\min} \leq 0.05$), and discriminatory otherwise. It was fair for the weak proxy because none of the predictions changed significantly between different levels of the protected attribute. It was discriminatory for stronger proxies because it required the proxy to be present in 10 out of 16 cases in order to label the robot as *Reliable*. Since these cases were in the majority and were all discriminatory for sensible δ_{\min} , this indicated the model generally discriminated. Nonetheless, participants were at most marginally affected by the proxy strength, and labeled the model discriminatory across all conditions (13/21, 10/20, and 16/21 participants across Weak Proxy, Medium Proxy, and Strong Proxy conditions, respectively). These proportions remained similar even when participants saw a comprehensive set of Multiple explanations (13/17 for Weak Proxy, 13/19 for Medium Proxy, 12/19 for Strong Proxy participants claimed the model was discriminatory). This suggests people generally equate the presence of a proxy with discrimination, regardless of its strength. If we relied on explanations to judge models globally, this would unnecessarily block deployment of multiple fair ones.

On the Consistency of Auditors and Decision Subjects Our evidence shows that participants’ claims were primarily driven by the presence of proxy variables in explanations. As expected, participants claimed discrimination 25-46% more frequently when explanations contained the proxy compared to when they did not (see Fig. 7). This effect was even more pronounced (36-60%) when participants viewed Multiple explanations. The increased exposure to explanations that contained the proxy in these conditions (14 instances versus 8 in the Single explanation conditions) led to a

30-47% increase in discrimination claims overall. These findings strongly suggest that proxy visibility directly impacts discrimination claims.

While participants were responsive to the presence of the proxy variable in the explanation, they often exercised nuance. In particular, we observed that participants consistently claimed that some predictions were “fair” even when the CE contained the proxy were judged as discriminatory. This behavior appears to be influenced by three systematic factors. First, their beliefs about robot reliability affected fairness judgments. Predictions were more likely to be labeled as fair by up to 20% when participants believed the proxy indicated higher reliability despite proxy dependence. While this pattern shows high variability ($p \approx 0.3$), it consistently appears across proxy conditions and aligns with participants’ explicit statements (e.g., *It is not unfair to say that robots with antennas work better*). The other two key factors are that participants assumed different protected attributes, which led them to state no discrimination and misrepresented the true proxy strength.

We also found that participants held false beliefs about the causal *structure* of the problem as described in REMARK 3. We

observed steady, low FPR of $\approx 30\%$ even under perfect assumptions about participant knowledge. This effect can only be attributed to labeling predictions that do not depend on the proxy as discriminatory, falsely believing other features are proxies. This sentiment can be found in participants’ answers (e.g., saying *I decided based on the body shape and the base type*). It also stems from our formulation. $\phi_{x,b,a}$ was similar across both values of the protected attribute for predictions that *could be* discriminatory and depended on the proxy (see Table 1 for predictions where $h(b, x) \neq h(b', x)$). As such, all these predictions would be discriminatory for high enough threshold values – yet we observe roughly the same FPR for $\delta_{\min} \in [0, 1]$, meaning participants also labeled predictions where $h(b, x) = h(b', x)$. In reality, we found that 36 out of 61 participants fell prey to these assumptions, including 8 participants who labeled predictions where the proxy was not present as discrimination. The remaining 28 participants saw a combination of proxy and some other features as discriminatory. This belief makes sense but shows the danger of interpreting the presence of the proxy as a single indicator of discrimination.

4.3 Discussion and Limitations

We identified two broad points that concern the lack of usability of explanations in practice:

Fundamental Detection Failure Auditing with either a single explanation or a comprehensive set of multiple explanations does not allow humans to reliably detect discrimination. Knowing the protected attribute of the audited predictions, or correctly identifying the causal mechanism of the proxy helps. However, it still does not enable detection of more than 65% of the truly discriminatory cases (TPR). It also leads to *at most 77%* correct detections (PPV), but only when one’s beliefs are treated as correct. Otherwise, reliability of detection oscillates around 50% with false alarms consistently hovering around 30% (FPR). To put that into perspective while being lenient on the participants’

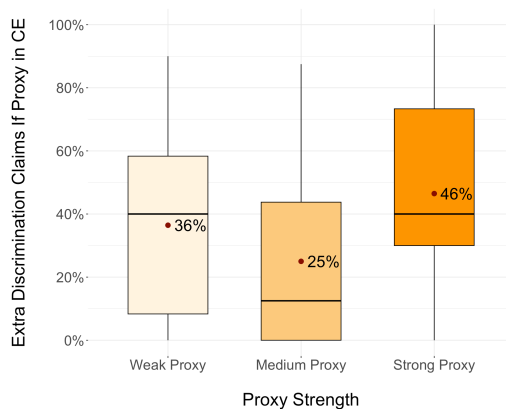


Fig. 7. Increase in discrimination claims when explanations contained the proxy versus when they did not. Mean values (red dots) show participants consistently identified the proxy as a discrimination signal across all regimes. As expected, strongest effect occurs with the strong proxy; weakest effect occurs with the medium proxy where relationships appear more irregular.

performance, this means every fourth individual that files a discrimination claim fails in court. This also means almost half of individuals whose predictions were truly discriminatory miss this.

Lack of Inter-Auditor Agreement One could try looking at the auditing performance with respect to model discrimination as more of a success. After all, the model which was discriminatory for most thresholds (when the proxy was medium and strong) would be determined as such by an average auditor. However, when it comes to individual performance, the results look much worse. First, more than half of all the participants claimed the model with the weak proxy was discriminatory when it was not (26/38 participants). Second, barely over half spot the model is discriminatory when it used a medium proxy (23/39 participants) and three quarters of the participants when the model used a strong proxy (28/40 participants). We observed a lack of overall agreement between participants who essentially operated on their own beliefs about discrimination. This led to claims that were very rarely matching (Cohen’s κ ranging from 0.05 to 0.14 across all conditions). This is also seen when we analyze predictions individually and find that every prediction was selected as discriminatory by at least 10% of the participants. Put together, if the same model or a set of predictions were analyzed by two independent auditors, it could lead to two different results. A discriminatory model could then be missed, and a fair model could be unfairly accused of discrimination.

The fundamental reason why explanations failed to aid discrimination detection is that they operate on individuals, whereas fairness must be evaluated over groups of (hypothetical) individuals. This tension is well-documented in formal definitions of fairness [62], and our experiments demonstrate how it impairs human performance. Our analysis revealed three specific challenges that emerged from this mismatch and were the direct causes of people’s failure:

Flawed Beliefs in the Causal Structure More than half of all participants (71 out of 118) fell prey to the beliefs that some features combined with the proxy are evidence of discrimination. 17 of the participants also thought that some combinations of features without the proxy can indicate discrimination. This led participants to incorrectly raise false alarms. This also led participants to not detect discrimination because they looked for “stronger proof” (e.g., one participant noted they looked for a combination of antenna and other features to claim discrimination).

Proxy Strength Misrepresentation Over half of the participants overestimated proxy strength. This is best seen by the largely improved performance (PPV and TPR) under their own beliefs in the causal mechanism when the thresholds are low. This led to many false positives in claiming discrimination. We can expect people to misrepresent the proxy strength in reality too because it is rarely observable. This misrepresentation might lead to a claim that the whole model is discriminating, while it is perfectly valid (like in the Weak Proxy conditions).

Real Outcome Interference Participants’ judgments were sometimes influenced by their beliefs about the relationship between features and desirable outcomes. This led to errors. We observed this behavior across all conditions. For instance, in the Weak Proxy condition with Multiple explanations, participants claimed predictions as fair in 52% of the cases when they thought adding a proxy makes the robot reliable, and otherwise, only in 28%. Even though the median increase was about 20%, as many as 78 out of all 118 participants made a claim like this at least once. We could also see this sentiment in participants’ responses, saying e.g., *It is not unfair to say ‘robots with antennas work better’*.

Limitations Our results are limited by two main factors that were beyond our control. First, our participants had no prior training in statistics or probability. This might have affected their judgments, making them inconsistent with respect to, e.g., proxy strength and the causal mechanism. This is especially important since fairness audits depend on probabilistic claims. Second, every study run on paid-survey platforms such as Prolific has to deal with inattentiveness or lack of motivation. Despite our best efforts, the task we introduced was abstract and gave no immediate feedback.

This could have made participants guess oftentimes and act inconsistently. They might have also had less incentive to perform thoughtfully, contrary to real auditors who may be bound by law.

5 Concluding Remarks

Our study demonstrates the fundamental limitations of using explanations for algorithmic fairness auditing. Through controlled experiments with human participants ($N = 126$), we found that explanations fail to reliably assist in discrimination detection, regardless of how much information they convey or if auditors know the protected attributes or the general causal mechanism of the proxy.

Our findings extend to real-world auditing scenarios. This is because real-world scenarios present far greater complexity, with more features, intricate relationships, and numerous plausible explanations to consider [20]. The failure modes that compromise human performance in our simple setup – flawed causal reasoning, incorrectly estimating proxy strength, and real outcome interference – are likely to persist or worsen with increased complexity. Furthermore, these individual-level failures may compound in real-world settings where multiple stakeholders must coordinate their assessments, just like the compounded in our experiment. In total, this will lead to poor discrimination detection performance in applied settings.

Our work is related to a growing body of regulations on algorithmic discrimination and transparency. In recent years, jurisdictions worldwide have adopted two main approaches. The first approach emphasizes transparency and explanation rights – see e.g., ECOA’s mandate for adverse action notices in lending [61] or provisions for a “Right to an Explanation” in for data regulation laws in the European Union [73], Brazil [12], and South Korea [38]. Mandatory fairness audits represent the second regulatory approach, e.g. in Slovenia mandates for algorithm pre-implementation [57], or in New York for third-party bias audits for automated employment decisions [36]. Similarly, the European Union’s Digital Services Act requires algorithmic audits of “very large online platforms,” including non-discrimination risk assessments [74]. Despite this momentum, there remains a lack of standardized practices for assessing algorithmic fairness as regulations provide limited guidance for how to conduct audits [46]. Our results highlight two critical insights for policy. First, there is a need for standalone regulations specifically targeting algorithmic discrimination. Current policy relying on explanations is unreliable even under controlled conditions (see also [33] for a legal discussion). Second, while the “right to explanation” serves a valuable role in accessing other rights (as exemplified in EU regulations), it should not be considered sufficient for preventing discrimination. Rather, it must be deployed alongside robust anti-discrimination measures and systematic auditing procedures that do not solely rely on human interpretation of explanations.

Acknowledgements

This work was supported by funding from the National Science Foundation under awards IIS 2040880 and IIS 2313105, and the NIH Bridge2AI Center Grant U54HG012510.

References

- [1] 116th Congress. 2019. Algorithmic Accountability Act of 2019. <https://www.congress.gov/bill/117th-congress/house-bill/6580/text>
- [2] 117th Congress. 2022. Algorithmic Accountability Act of 2022. <https://www.congress.gov/bill/117th-congress/house-bill/6580/text>
- [3] Abubakar Abid, Mert Yuksekgonul, and James Zou. 2022. Meaningfully debugging model mistakes using conceptual counterfactual explanations. In *International Conference on Machine Learning*. PMLR, 66–88.
- [4] Philip Adler, Casey Falk, Sorelle A Friedler, Tionney Nix, Gabriel Rybeck, Carlos Scheidegger, Brandon Smith, and Suresh Venkatasubramanian. 2018. Auditing black-box models for indirect influence. *Knowledge and Information Systems* 54 (2018), 95–122.
- [5] Ifeoma Ajunwa, Sorelle Friedler, Carlos E Scheidegger, and Suresh Venkatasubramanian. 2016. Hiring by algorithm: predicting and preventing disparate impact. *Available at SSRN* (2016).
- [6] Sajid Ali, Tamer Abuhmed, Shaker El-Sappagh, Khan Muhammad, Jose M Alonso-Moral, Roberto Confalonieri, Riccardo Guidotti, Javier Del Ser, Natalia Díaz-Rodríguez, and Francisco Herrera. 2023. Explainable Artificial Intelligence (XAI): What we know and what is left to attain Trustworthy Artificial Intelligence. *Information fusion* 99 (2023), 101805.
- [7] Alessa Angerschmid, Jianlong Zhou, Kevin Theuermann, Fang Chen, and Andreas Holzinger. 2022. Fairness and explanation in AI-informed decision making. *Machine Learning and Knowledge Extraction* 4, 2 (2022), 556–579.
- [8] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information fusion* 58 (2020), 82–115.
- [9] Gagan Bansal, Tongshuang Wu, Joyce Zhou, Raymond Fok, Besmira Nushi, Ece Kamar, Marco Tulio Ribeiro, and Daniel Weld. 2021. Does the whole exceed its parts? the effect of ai explanations on complementary team performance. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*. 1–16.
- [10] Solon Barocas and Andrew D Selbst. 2016. Big data’s disparate impact. *California law review* (2016), 671–732.
- [11] Reuben Binns, Max Van Kleek, Michael Veale, Ulrik Lyngs, Jun Zhao, and Nigel Shadbolt. 2018. ‘It’s Reducing a Human Being to a Percentage’ Perceptions of Justice in Algorithmic Decisions. In *Proceedings of the 2018 Chi conference on human factors in computing systems*. 1–14.
- [12] Brazil. 2020. Brazilian General Data Protection Law. https://iapp.org/media/pdf/resource_center/Brazilian_General_Data_Protection_Law.pdf
- [13] Marc-Etienne Brunet, Ashton Anderson, and Richard Zemel. 2022. Implications of Model Indeterminacy for Explanations of Automated Decisions. *Advances in Neural Information Processing Systems* 35 (2022), 7810–7823.
- [14] Zana Bućinca, Phoebe Lin, Krzysztof Z Gajos, and Elena L Glassman. 2020. Proxy tasks and subjective measures can be misleading in evaluating explainable AI systems. In *Proceedings of the 25th international conference on intelligent user interfaces*. 454–464.
- [15] Zana Bućinca, Maja Barbara Malaya, and Krzysztof Z Gajos. 2021. To trust or to think: cognitive forcing functions can reduce overreliance on AI in AI-assisted decision-making. *Proceedings of the ACM on Human-Computer Interaction* 5, CSCW1 (2021), 1–21.
- [16] Zana Bućinca, Siddharth Swaroop, Amanda E Paluch, Finale Doshi-Velez, and Krzysztof Z Gajos. 2024. Contrastive Explanations That Anticipate Human Misconceptions Can Improve Human Decision-Making Skills. *arXiv preprint arXiv:2410.04253* (2024).
- [17] Federico Maria Cau, Hanna Hauptmann, Lucio Davide Spano, and Nava Tintarev. 2023. Effects of ai and logic-style explanations on users’ decisions under different levels of uncertainty. *ACM Transactions on Interactive Intelligent Systems* 13, 4 (2023), 1–42.
- [18] Federico Maria Cau, Hanna Hauptmann, Lucio Davide Spano, and Nava Tintarev. 2023. Supporting high-uncertainty decisions through AI and logic-style explanations. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 251–263.
- [19] Valerie Chen, Q Vera Liao, Jennifer Wortman Vaughan, and Gagan Bansal. 2023. Understanding the role of human intuition on reliance in human-AI decision-making with explanations. *arXiv preprint arXiv:2301.07255* (2023).
- [20] Sasha Costanza-Chock, Inioluwa Deborah Raji, and Joy Buolamwini. 2022. Who Audits the Auditors? Recommendations from a field scan of the algorithmic auditing ecosystem. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1571–1583.
- [21] Jessica Dai, Paula Gradu, Inioluwa Deborah Raji, and Benjamin Recht. 2025. From Individual Experience to Collective Evidence: A Reporting-Based Framework for Identifying Systemic Harms. *arXiv preprint arXiv:2502.08166* (2025).
- [22] Xinyue Dai, Mark T Keane, Laurence Shaloo, Elodie Ruelle, and Ruth MJ Byrne. 2022. Counterfactual explanations for prediction and diagnosis in XAI. In *Proceedings of the 2022 AAAI/ACM Conference on AI, Ethics, and Society*. 215–226.
- [23] Arun Das and Paul Rad. 2020. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371* (2020).
- [24] Philip Dawid. 2017. On individual risk. *Synthese* 194, 9 (2017), 3445–3474.
- [25] Eoin Delaney, Arjun Pakrashi, Derek Greene, and Mark T Keane. 2023. Counterfactual explanations for misclassified images: How human and machine explanations differ. *Artificial Intelligence* 324 (2023), 103995.
- [26] Jonathan Dodge, Q Vera Liao, Yunfeng Zhang, Rachel KE Bellamy, and Casey Dugan. 2019. Explaining models: an empirical study of how explanations impact fairness judgment. In *Proceedings of the 24th international conference on intelligent user interfaces*. 275–285.
- [27] Mengnan Du, Fan Yang, Na Zou, and Xia Hu. 2020. Fairness in deep learning: A computational perspective. *IEEE Intelligent Systems* 36, 4 (2020), 25–34.
- [28] Equal Employment Opportunity Commission. 1978. Uniform Guidelines on Employee Selection Procedures). Electronic Code of Federal Regulations. <https://www.ecfr.gov/current/title-29/subtitle-B/chapter-XIV/part-1607/subject-group-ECFRdb347e844acde6> 29 CFR Part 1607.

- [29] Carlos Fernández-Loría, Foster Provost, and Xintian Han. 2022. Explaining Data-driven Decisions Made by AI Systems: The Counterfactual Approach. *MIS Quarterly* 46, 3 (2022), 1635–1660.
- [30] Maximilian Förster, Mathias Klier, Kilian Kluge, and Irina Sigler. 2020. Evaluating Explainable Artificial Intelligence - What Users Really Appreciate. In *28th European Conference on Information Systems - Liberty, Equality, and Fraternity in a Digitizing World, ECIS 2020, Marrakech, Morocco, June 15-17, 2020*, Frantz Rowe, Redouane El Amrani, Moez Limayem, Sue Newell, Nancy Pouloudi, Eric van Heck, and Ali El Quammah (Eds.). https://aisel.laisnet.org/ecis2020_rp/195
- [31] Ana Cristina Bicharra Garcia, Marcio Gomes Pinto Garcia, and Roberto Rigobon. 2024. Algorithmic discrimination in the credit domain: what do we know about it? *AI & SOCIETY* 39, 4 (2024), 2059–2098.
- [32] Talia B Gillis, Vitaly Meursault, and Berk Ustun. 2024. Operationalizing the Search for Less Discriminatory Alternatives in Fair Lending. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 377–387.
- [33] Talia B Gillis and Josh Simons. 2019. Explanation < Justification: GDPR and the Perils of Privacy. *JL & Innovation* 2 (2019), 71.
- [34] Ana Valeria González, Gagan Bansal, Angela Fan, Yashar Mehdad, Robin Jia, and Srinivasan Iyer. 2021. Do explanations help users detect errors in open-domain QA? an evaluation of spoken vs. visual explanations. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*. 1103–1116.
- [35] Navita Goyal, Connor Baumler, Tin Nguyen, and Hal Daumé III. 2024. The Impact of Explanations on Fairness in Human-AI Decision-Making: Protected vs Proxy Features. In *Proceedings of the 29th International Conference on Intelligent User Interfaces*. 155–180.
- [36] Lara Groves, Jacob Metcalf, Alayna Kennedy, Briana Vecchione, and Andrew Strait. 2024. Auditing work: Exploring the New York City algorithmic bias audit regime. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1107–1120.
- [37] Kofi Immanuel Jones and Swati Sah. 2023. The Implementation of Machine Learning In The Insurance Industry With Big Data Analytics. *International Journal of Data Informatics and Intelligent Computing* 2, 2 (2023), 21–38.
- [38] Dong Hyeon Kim and Do Hyun Park. 2024. Automated decision-making in South Korea: a critical review of the revised Personal Information Protection Act. *Humanities and Social Sciences Communications* 11, 1 (2024), 1–11.
- [39] Issa Kohler-Hausmann. 2018. Eddie Murphy and the dangers of counterfactual causal thinking about detecting racial discrimination. *Nw. UL Rev.* 113 (2018), 1163.
- [40] Satyapriya Krishna, Tessa Han, Alex Gu, Javin Pombra, Shahin Jabbari, Steven Wu, and Himabindu Lakkaraju. 2022. The Disagreement Problem in Explainable Machine Learning: A Practitioner’s Perspective. *arXiv:2202.01602 [cs]* (Feb. 2022). <http://arxiv.org/abs/2202.01602> arXiv: 2202.01602.
- [41] Ulrike Kuhl, André Artelt, and Barbara Hammer. 2022. Keep your friends close and your counterfactuals closer: Improved learning from closest rather than plausible counterfactual explanations in an abstract setting. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 2125–2137.
- [42] Ulrike Kuhl, André Artelt, and Barbara Hammer. 2023. For Better or Worse: The Impact of Counterfactual Explanations’ Directionality on User Behavior in xAI. In *World Conference on Explainable Artificial Intelligence*. Springer, 280–300.
- [43] Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. *Advances in neural information processing systems* 30 (2017).
- [44] Vivian Lai, Han Liu, and Chenhao Tan. 2020. "Why is' Chicago' decepti ve?" Towards Building Model-Driven Tutorials for Humans. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*. 1–13.
- [45] Vivian Lai and Chenhao Tan. 2019. On human predictions with explanations and predictions of machine learning models: A case study on deception detection. In *Proceedings of the conference on fairness, accountability, and transparency*. 29–38.
- [46] Khoa Lam, Benjamin Lange, Borhane Blili-Hamelin, Jovana Davidovic, Shea Brown, and Ali Hasan. 2024. A framework for assurance audits of algorithmic systems. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1078–1092.
- [47] Olesja Lammert, Birte Richter, Christian Schütze, Kirsten Thommes, and Britta Wrede. 2024. Humans in XAI: increased reliance in decision-making under uncertainty by using explanation strategies. *Frontiers in Behavioral Economics* 3 (2024), 1377075.
- [48] Markus Langer, Daniel Oster, Timo Speith, Holger Hermanns, Lena Kästner, Eva Schmidt, Andreas Sesing, and Kevin Baum. 2021. What do we want from Explainable Artificial Intelligence (XAI)?—A stakeholder perspective on XAI and a conceptual model guiding interdisciplinary XAI research. *Artificial Intelligence* 296 (2021), 103473.
- [49] Thao Le, Tim Miller, Ronal Singh, and Liz Sonenberg. 2022. Improving model understanding and trust with counterfactual explanations of model confidence. *arXiv preprint arXiv:2206.02790* (2022).
- [50] Min Hun Lee and Chong Jun Chew. 2023. Understanding the effect of counterfactual explanations on trust and reliance on ai for human-ai collaborative clinical decision making. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW2 (2023), 1–22.
- [51] Min Kyung Lee, Anuraag Jain, Hea Jin Cha, Shashank Ojha, and Daniel Kusbit. 2019. Procedural justice in algorithmic fairness: Leveraging transparency and outcome control for fair algorithmic mediation. *Proceedings of the ACM on Human-Computer Interaction* 3, CSCW (2019), 1–26.
- [52] Michelle Seng Ah Lee and Luciano Floridi. 2021. Algorithmic fairness in mortgage lending: from absolute conditions to relational trade-offs. *Minds and Machines* 31, 1 (2021), 165–191.
- [53] Geng Li. 2018. Gender-Related Differences in Credit Use and Credit Scores. *FEDS Notes* (22 June 2018). <https://doi.org/10.17016/2380-7172.2188>
- [54] Luca Longo, Mario Brcic, Federico Cabitza, Jaesik Choi, Roberto Confalonieri, Javier Del Ser, Riccardo Guidotti, Yoichi Hayashi, Francisco Herrera, Andreas Holzinger, et al. 2024. Explainable Artificial Intelligence (XAI) 2.0: A manifesto of open challenges and interdisciplinary research directions. *Information Fusion* 106 (2024), 102301.

- [55] Ana Lucic, Hinda Haned, and Maarten de Rijke. 2020. Why does my model fail? contrastive local explanations for retail forecasting. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. 90–98.
- [56] Marco Lünich and Birte Keller. 2024. Explainable Artificial Intelligence for Academic Performance Prediction. An Experimental Study on the Impact of Accuracy and Simplicity of Decision Trees on Causability and Fairness Perceptions. In *The 2024 ACM Conference on Fairness, Accountability, and Transparency*. 1031–1042.
- [57] Gianclaudio Malgieri. 2019. Automated decision-making in the EU Member States: The right to explanation and other “suitable safeguards” in the national legislations. *Computer law & security review* 35, 5 (2019), 105327.
- [58] Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. 2020. Explaining machine learning classifiers through diverse counterfactual explanations. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*. ACM, Barcelona Spain, 607–617. <https://doi.org/10.1145/3351095.3372850>
- [59] Chelsea M Myers, Evan Freed, Luis Fernando Laris Pardo, Anushay Furqan, Sebastian Risi, and Jichen Zhu. 2020. Revealing neural network bias to non-experts through interactive counterfactual examples. *arXiv preprint arXiv:2001.02271* (2020).
- [60] Hamed Nilforoshan, Johann D Gaebler, Ravi Shroff, and Sharad Goel. 2022. Causal conceptions of fairness and their consequences. In *International Conference on Machine Learning*. PMLR, 16848–16887.
- [61] Bureau of Consumer Financial Protection. 2020. Equal Credit Opportunity (Regulation B): Revocations or Unfavorable Changes to the Terms of Existing Credit Arrangements. https://files.consumerfinance.gov/f/documents/cfbp_revoking-terms-of-existing-credit-arrangement_advisory-opinion_2022-05.pdf
- [62] Drago Plečko, Elias Bareinboim, et al. 2024. Causal fairness analysis: a causal toolkit for fair machine learning. *Foundations and Trends® in Machine Learning* 17, 3 (2024), 304–589.
- [63] Emilee Rader, Kelley Cotter, and Janghee Cho. 2018. Explanations as mechanisms for supporting algorithmic transparency. In *Proceedings of the 2018 CHI conference on human factors in computing systems*. 1–13.
- [64] Max Schemmer, Joshua Holstein, Niklas Bauer, Niklas Kühl, and Gerhard Satzger. 2023. Towards meaningful anomaly detection: The effect of counterfactual explanations on the investigation of anomalies in multivariate time series. *arXiv preprint arXiv:2302.03302* (2023).
- [65] Jakob Schoeffer, Niklas Kuehl, and Yvette Machowski. 2022. “There is not enough information”: On the effects of explanations on perceptions of informational fairness and trustworthiness in automated decision-making. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1616–1628.
- [66] Avital Shulner-Tal, Tsvi Kuflik, and Doron Kliger. 2022. Fairness, explainability and in-between: understanding the impact of different explanation methods on non-expert users’ perceptions of fairness toward an algorithmic system. *Ethics and Information Technology* 24, 1 (2022), 2.
- [67] Ronal Singh, Tim Miller, Henrietta Lyons, Liz Sonenberg, Eduardo Velloso, Frank Vetere, Piers Howe, and Paul Dourish. 2023. Directive explanations for actionable explainability in machine learning applications. *ACM Transactions on Interactive Intelligent Systems* 13, 4 (2023), 1–26.
- [68] Nina Spreitzer, Hinda Haned, and Ilse van der Linden. 2022. Evaluating the Practicality of Counterfactual Explanations. In *XAI. it@ AI* IA*. 31–50.
- [69] Ilija Stepin, Jose M Alonso-Moral, Alejandro Catala, and Martin Pereira-Far ía. 2022. An empirical study on how humans appreciate automated counterfactual explanations which embrace imprecise information. *Information Sciences* 618 (2022), 379–399.
- [70] Winnie F Taylor. 1980. Meeting the Equal Credit Opportunity Act’s Specificity Requirement: Judgmental and Statistical Scoring Systems. *Buff. L. Rev.* 29 (1980), 73.
- [71] Taylor Telford. 2019. Apple Card algorithm sparks gender bias allegations against Goldman Sachs. *Washington Post* 11 (2019).
- [72] Michael Carl Tschantz. 2022. What is proxy discrimination?. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and Transparency*. 1993–2003.
- [73] European Union. 2018. General Data Protection Regulation, Art. 22. <https://gdpr-info.eu/art-22-gdpr/>
- [74] European Union. 2024. The Digital Services Act. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A32022R2065/>
- [75] Jasper van der Waa, Elisabeth Nieuwburg, Anita Cremers, and Mark Neerincx. 2021. Evaluating XAI: A comparison of rule-based and example-based explanations. *Artificial intelligence* 291 (2021), 103404.
- [76] Helena Vasconcelos, Matthew Jörke, Madeleine Grunde-McLaughlin, Tobias Gerstenberg, Michael S Bernstein, and Ranjay Krishna. 2023. Explanations can reduce overreliance on ai systems during decision-making. *Proceedings of the ACM on Human-Computer Interaction* 7, CSCW1 (2023), 1–38.
- [77] Michael Veale and Irina Brass. 2019. Administration by algorithm? Public management meets public sector machine learning. *Public management meets public sector machine learning* (2019).
- [78] Sandra Wachter, Brent Mittelstadt, and Chris Russell. 2017. Counterfactual explanations without opening the black box: Automated decisions and the GDPR. *Harv. JL & Tech.* 31 (2017), 841.
- [79] J Christina Wang and Charles B Perkins. 2019. How magic a bullet is machine learning for credit analysis? An exploration with FinTech lending data. *An Exploration with FinTech Lending Data (October 21, 2019)* (2019).
- [80] Xinru Wang and Ming Yin. 2021. Are explanations helpful? a comparative study of the effects of explanations in AI-assisted decision-making. In *26th international conference on intelligent user interfaces*. 318–328.
- [81] Richard Warner and Robert H Sloan. 2021. Making artificial intelligence transparent: Fairness and the problem of proxy variables. *Criminal Justice Ethics* 40, 1 (2021), 23–39.

- [82] Greta Warren, Ruth MJ Byrne, and Mark T Keane. 2023. Categorical and continuous features in counterfactual explanations of AI systems. In *Proceedings of the 28th International Conference on Intelligent User Interfaces*. 171–187.
- [83] Monika Westphal, Michael Vössing, Gerhard Satzger, Galit B Yom-Tov, and Anat Rafaeli. 2023. Decision control and explanations in human-AI collaboration: Improving user perceptions and compliance. *Computers in Human Behavior* 144 (2023), 107714.
- [84] Mireia Yurrita, Tim Draws, Agathe Balayn, Dave Murray-Rust, Nava Tintarev, and Alessandro Bozzon. 2023. Disentangling fairness perceptions in algorithmic decision-making: the effects of explanations, human oversight, and contestability. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems*. 1–21.
- [85] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P Gummadi. 2019. Fairness constraints: A flexible approach for fair classification. *The Journal of Machine Learning Research* 20, 1 (2019), 2737–2778.
- [86] Yunfeng Zhang, Q Vera Liao, and Rachel KE Bellamy. 2020. Effect of confidence and explanation on accuracy and trust calibration in AI-assisted decision making. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*. 295–305.
- [87] Zelun Tony Zhang, Felicitas Buchner, Yuanting Liu, and Andreas Butz. 2024. You Can Only Verify When You Know the Answer: Feature-Based Explanations Reduce Overreliance on AI for Easy Decisions, but Not for Hard Ones. In *Proceedings of Mensch und Computer 2024*. 156–170.

A Supplementary Material on Experimental Design

In this Section, we provide supplementary materials on our experimental design. This includes the exact list of robots (points the model predicted on) with their closest counterfactual explanations in Table 2, and links to our GitHub repository with the code for the experiment and the experimental data.

Features				Prevalence		Counterfactual Explanations
Antenna	HeadShape	BodyShape	BaseType	Company X	Company S	
No	Square	Square	Legs	0.0071	0.0004	{Antenna, HeadShape}, {Antenna, BaseType}, {Antenna, HeadShape}, {BodyShape, BaseType}
No	Square	Square	Wheels	0.016	0.0008	{Antenna}
No	Square	Round	Legs	0.016	0.0008	{Antenna}, {BodyShape}
No	Square	Round	Wheels	0.0297	0.0016	{Antenna}, {BodyShape}
No	Round	Square	Legs	0.016	0.0008	{Antenna}, {BaseType}
No	Round	Square	Wheels	0.0297	0.0016	{Antenna}, {BaseType}
No	Round	Round	Legs	0.0297	0.0016	{BodyShape}, {BaseType}
No	Round	Round	Wheels	0.0434	0.0023	{BodyShape}, {BaseType}
Yes	Square	Square	Legs	0.0008	0.016	{HeadShape}, {BodyShape}, {BaseType}
Yes	Square	Square	Wheels	0.016	0.0297	{Antenna}, {HeadShape}
Yes	Square	Round	Legs	0.016	0.0297	{Antenna}, {BaseType}
Yes	Square	Round	Wheels	0.0023	0.0434	{Antenna}
Yes	Round	Square	Legs	0.016	0.0297	{Antenna}, {BodyShape}
Yes	Round	Square	Wheels	0.0023	0.0434	{Antenna}
Yes	Round	Round	Legs	0.0023	0.0434	{Antenna, BodyShape}, {Antenna, BaseType}, {BodyShape, BaseType}
Yes	Round	Round	Wheels	0.0028	0.0523	{Antenna, BodyShape}, {Antenna, BaseType}

Table 2. Overview of closest counterfactual explanations over all robot types. We consider 16 robots defined by four binary attributes: Antenna, HeadShape, BodyShape, BaseType. Each combination of characteristics (row) is predicted as predicted `Reliable` if it has an Antenna and one of the following conditions: a Round HeadShape, a Round BodyShape, or Wheels. Otherwise it is predicted `Defective`. Based on this specification, we obtain closest counterfactuals that allow flipping the prediction.

A.1 Availability of data and material (data transparency)

Anonymized data from the experiments is available at https://anonymous.4open.science/r/cxai-93BB/results/results_closest_competing.

A.2 Code availability (software application or custom code)

The code for our Flask study is available at <https://anonymous.4open.science/r/cxai-93BB/>.

1. Run `pip3 install -r requirements.txt` to install the necessary requirements.
2. Then run `application.py` and open the link to the localhost to start the study.
3. Parameters listed at the top of the file can be used to run the study in different conditions.

B Supplementary Experimental Results

In this Section, we provide additional figures for our experimental results from the main text. Fig. 8 shows performance measures (PPV, TPR and FPR) across all thresholds $\delta_{\min} \in [0, 1]$ in the conditions that used Multiple explanations. We detail the results of these studies in Section 4.2. Fig. 9 shows that participants' claims depended on the presence of the proxy in the explanation also for Multiple explanations conditions. Finally, Fig. 10 shows the lack of agreement between the participants we discussed in Section 4.2, detailing how often each of the predictions used in the study was claim as discriminatory.

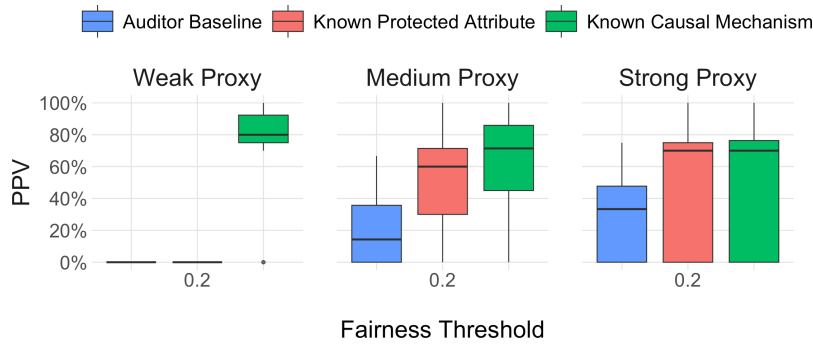


Fig. 8. Refer to Fig. 6 for the explanation of the plotted data. As shown, baseline performance (blue) is poor across all thresholds, with TPR not exceeding 50% and FPR around 30%, and sometimes exceeding this value. Knowledge of protected attributes (red) yields significant gains for PPV for medium and strong proxies but is otherwise unhelpful. Assuming auditors’ beliefs about the causal mechanism (green) provides the biggest gains for performance, especially internal reliability in terms of PPV. It still leads to largely low TPR and moderate FPR. The latter metric remains problematic (around 30%) across all conditions, indicating persistent incorrect assumptions about feature-protected attribute relationships regardless of the level of insight.

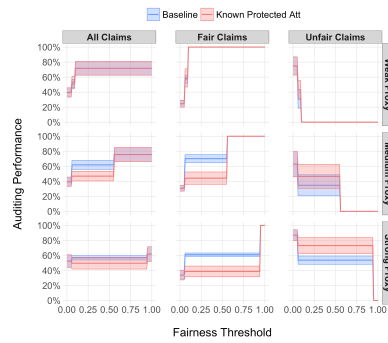


Fig. 9. Increase in discrimination claims when explanations contained the proxy versus when they did not. Mean values (red dots) show participants consistently identified the proxy as a discrimination signal across all regimes. The strongest effect occurs for the weak proxy because participants overestimated its strength.

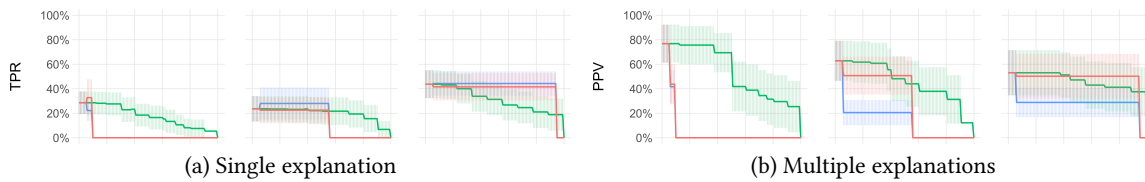


Fig. 10. Discrimination claims per individual predictions in each of the proxy regimes when participants saw a single explanation (left) and multiple explanations (right). Stars indicate when explanations contained the proxy. We can see that every prediction was judged as discriminatory by at least 10% of the participants. Participants were also not in full agreement with any of the predictions. On average, the agreement was roughly 50%.