
Learning Individually Fair Classifier with Path-Specific Causal-Effect Constraint

Yoichi Chikahara^{1,3}
¹NTT

Shinsaku Sakaue²
²The University of Tokyo

Akinori Fujino¹

Hisashi Kashima³
³Kyoto University

Abstract

Machine learning is used to make decisions for individuals in various fields, which require us to achieve good prediction accuracy while ensuring fairness with respect to sensitive features (e.g., race and gender). This problem, however, remains difficult in complex real-world scenarios. To quantify unfairness under such situations, existing methods utilize *path-specific causal effects*. However, none of them can ensure fairness for each individual without making impractical functional assumptions about the data. In this paper, we propose a far more practical framework for learning an individually fair classifier. To avoid restrictive functional assumptions, we define the *probability of individual unfairness* (PIU) and solve an optimization problem where PIU's upper bound, which can be estimated from data, is controlled to be close to zero. We elucidate why our method can guarantee fairness for each individual. Experimental results show that our method can learn an individually fair classifier at a slight cost of accuracy.

1 INTRODUCTION

Machine learning is increasingly being used to make critical decisions that severely affect people's lives (e.g., loan approvals (Khandani et al., 2010), hiring decisions (Houser, 2019), and recidivism predictions (Angwin et al., 2016)). The huge societal impact of such decisions on people's lives raises concerns about fairness because these decisions may be discriminatory with respect to *sensitive features*, including race, gender, religion, and sexual orientation.

Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

Although many researchers have studied how to make fair decisions while achieving high prediction accuracy (Dwork et al., 2012; Feldman et al., 2015; Hardt et al., 2016), it remains a challenge in complex real-world scenarios. For instance, consider hiring decisions for physically demanding jobs. Although it is discriminatory to reject applicants based on gender, since the job requires physical strength, it is sometimes **not** discriminatory to reject the applicants due to physical strength. Since physical strength is often affected by gender, rejecting applicants due to physical strength leads to gender difference in the rejection rates. Although this difference due to physical strength is **not** always unfair, it is removed when using traditional methods (e.g., Feldman et al. (2015)). Consequently, even if there is a man who has a much more physical strength than a woman, these methods might reject him to accept her, which severely reduces the prediction accuracy.

To achieve high prediction accuracy, we need to remove only an unfair difference in decision outcomes. To measure this difference, existing methods utilize a path-specific causal effect (Avin et al., 2005), which we call an *unfair effect*. Using unfair effects, the *path-specific counterfactual fairness* (PSCF) method (Chiappa and Gillam, 2019) aims to guarantee fairness for each individual; however, achieving such an individual-level fairness is possible only when the data are generated by a restricted class of functions. By contrast, *fair inference on outcome* (FIO) (Nabi and Shpitser, 2018) does not require such demanding functional assumptions; however, it cannot ensure individual-level fairness.

The goal of this paper is to propose a learning framework that guarantees individual-level fairness without making impractical functional assumptions. For this goal, we train a classifier by forcing the *probability of individual unfairness* (PIU), defined as the probability that an unfair effect is non-zero, to be close to zero. This, however, is difficult to achieve because we cannot estimate PIU from data. To overcome this difficulty, we derive its upper bound that can be estimated from data and solve a penalized optimization problem where the upper-bound value is controlled to be close to zero.

Table 1: Comparison with existing methods

Method	Individually fair	Functional assumptions
Our method	Yes	Unnecessary
PSCF	Yes	Necessary
FIO	No	Unnecessary

Our contributions are summarized as follows:

- We establish a framework that guarantees fairness for each individual without restrictive functional assumptions on the data (Table 1). To achieve this, we make the PIU value close to zero by imposing a penalty that reduces its upper bound value, which can be estimated from data.
- We elucidate why imposing such a penalty guarantees individual-level fairness in Sections 4.3.3 and 4.4. We also show how our method can be extended to address cases where there are unobserved variables called *latent confounders* in Section 4.5.
- We experimentally show that our method makes much fairer predictions for each individual than the existing methods at a slight cost of prediction accuracy.

2 PRELIMINARIES

2.1 Problem Statement

In this paper, we consider a binary classification task. We train classifier h_θ with parameter θ to predict decision outcome $Y \in \{0, 1\}$ from the features of each individual \mathbf{X} , which contains sensitive feature $A \in \{0, 1\}$.

We seek classifier parameter θ that achieves a good balance between prediction accuracy and fairness with respect to sensitive feature A . Suppose that we have loss function L_θ and penalty function G_θ , which respectively measure prediction errors and unfairness based on θ . Formally, given n training instances $\{(\mathbf{x}_i, y_i)_{i=1}^n\}$, our learning problem is formulated as follows:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n L_\theta(\mathbf{x}_i, y_i) + \lambda G_\theta(\mathbf{x}_1, \dots, \mathbf{x}_n), \quad (1)$$

where $\lambda \geq 0$ is a hyperparameter.

To achieve a high prediction accuracy, penalty function G_θ must be designed such that we can avoid imposing unnecessary penalizations. To do so, we utilize a *causal graph*, which is a directed acyclic graph (DAG) whose nodes and edges represent random variables and causal relationships, respectively (Pearl, 2009). We assume that a causal graph is provided by domain experts

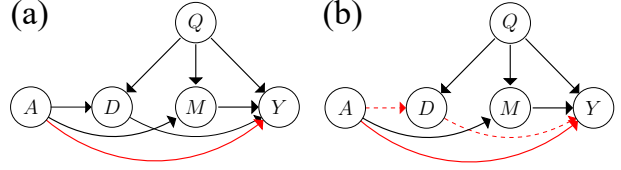


Figure 1: Causal graphs representing a scenario of hiring decisions for physically demanding jobs: Unfair pathways are (a): red solid edge $A \rightarrow Y$; (b): $A \rightarrow Y$ and red dashed pathway $A \rightarrow D \rightarrow Y$.

or can be inferred from data (Glymour et al., 2019); this assumption is common in many existing methods (Chiappa and Gillam, 2019; Kusner et al., 2017; Nabi and Shpitser, 2018; Zhang et al., 2017).

As an example of a causal graph, consider a scenario for hiring decisions for a physically demanding job. In this scenario, a causal graph might be given, as shown in Figure 1(a), where $A, Q, D, M \in \mathbf{X}$ represent gender, qualifications, the number of children, and physical strength, respectively. This graph expresses our knowledge that prediction Y is unfair only if it is based on gender A . To do so, we regard direct pathway $A \rightarrow Y$ as unfair pathway π (i.e., $\pi = \{A \rightarrow Y\}$).

In fact, in the above simple case, we can naively remove the unfairness by making a prediction without gender A ; however, this is insufficient when we consider a complex scenario with multiple unfair pathways.

For instance, as shown in Figure 1(b), we may regard not only $A \rightarrow Y$ but also the pathway through the number of children D ($A \rightarrow D \rightarrow Y$) as unfair (because it is also discriminatory to reject women because of the possibility of bearing children). In this case, a naive approach to ensure fairness is to predict without A or D . This, however, might unnecessarily decrease the prediction accuracy. For instance, consider a case where number of children D is only slightly affected by gender A (e.g., the applicants had gender-equitable opportunities to take parental leave in the past) and largely influenced by other unobserved features that are important for prediction (e.g., communication skills). Then predicting without D will seriously decrease the accuracy while contributing almost nothing to fairness.

To address such cases, given unfair pathways π , we design penalty function G_θ by quantifying the unfairness based on data. To do so, we utilize path-specific causal effects, which are described in the next section.

2.2 Path-Specific Causal Effects

A path-specific causal effect measures how largely an observed variable influences another variable via path-

ways in a causal graph (Avin et al., 2005). Although prediction Y is not observed but is given by classifier h_θ , we can utilize this measure to quantify the influence of sensitive feature A on Y via unfair pathways π .

With a path-specific causal effect, this influence is measured by the difference of the two predictions, which is obtained by modifying input features \mathbf{X} . To illustrate these predictions, consider a case where sensitive feature A is gender. Then for each woman ($A = 0$), one prediction is made by directly taking her attributes as input, and another is made with *counterfactual* attributes, which would be observed if she were male ($A = 1$); for each man, these predictions are made using the counterfactual attributes that would be if he were female (see, Appendix A.1 for details). Although such counterfactual attributes are not observed, they can be computed by a *structural equation model* (SEM).

An SEM consists of *structural equations*, each of which expresses variable $V \in \{\mathbf{X}, Y\}$ by deterministic function f_V (Pearl, 2009). Each function f_V takes as input two types of variables. One is observed variables, which are the parents of V in a causal graph, and the other is unobserved noise U_V , which expresses random variable V using deterministic function f_V .

For instance, structural equations over $D, M \in \mathbf{X}$ in the causal graph in Figure 1(b) may be formulated as

$$\begin{aligned} D &= f_D(A, Q, U_D) = A + U_D Q, \\ M &= f_M(A, Q, U_M) = 3A + 0.5Q + U_M, \end{aligned} \quad (2)$$

where U_D is multiplicative noise and U_M is additive noise. By contrast, the structural equation over prediction Y is formulated using classifier h_θ . If h_θ is deterministic, it is expressed as $Y = h_\theta(A, Q, D, M)$; otherwise, $Y = h_\theta(A, Q, D, M, U_Y)$, where U_Y is a random variable used in the classifier. See Appendix A.2 for a formal definition of SEM in our setting.

Structural equations (2) can be used to compute the (counterfactual) attributes of D and M that are observed when $A = a$ ($a \in \{0, 1\}$) as

$$D(a) = a + U_D Q, \quad M(a) = 3a + 0.5Q + U_M. \quad (3)$$

If (3) is available, we can obtain attributes $D(0)$, $D(1)$, $M(0)$, and $M(1)$ for each individual.

Using these attributes, we can compute a path-specific causal effect for each individual, which we call an unfair effect. For instance, when measuring the influence via unfair pathways $\pi = \{A \rightarrow Y, A \rightarrow D \rightarrow Y\}$ in Figure 1(b), we define an unfair effect as the difference of two predictions $Y_{A \leftarrow 1 \parallel \pi} - Y_{A \leftarrow 0}$, where $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$ are called *potential outcomes* and given as

$$\begin{aligned} Y_{A \leftarrow 0} &= h_\theta(0, Q, D(0), M(0)), \\ Y_{A \leftarrow 1 \parallel \pi} &= h_\theta(1, Q, D(1), M(0)). \end{aligned} \quad (4)$$

In (4), the inputs of $Y_{A \leftarrow 0}$ are $A = 0$, $D(0)$, and $M(0)$, all of which are given using the same value, $a = 0$. By contrast, the inputs of $Y_{A \leftarrow 1 \parallel \pi}$ are formulated based on unfair pathways π ; we use the value $a = 1$ **only** for A and D (i.e., $A = 1$ and $D(1)$), which correspond to the nodes on $\pi = \{A \rightarrow Y, A \rightarrow D \rightarrow Y\}$ (see Appendix A.3 for the formal definition).¹

In practice, however, we cannot compute an unfair effect for each individual. This is because we cannot formulate an SEM since it requires a deep understanding of true data-generating processes; consequently, for instance, we can obtain $D(a)$ and $M(a)$ in (3) **only** for either $a = 0$ or $a = 1$ but **not both**. Due to this issue, existing methods use the (conditional) expected values of unfair effects, which can be estimated from data.

3 EXISTING METHODS AND THEIR WEAKNESSES

Using unfair effects, two types of existing methods have been proposed. Unfortunately, as presented in Table 1, each has a weakness. One requires restrictive functional assumptions, and the other cannot ensure individual-level fairness. Below we describe their details.

3.1 Methods for Ensuring Individual-Level Fairness

The PSCF method (Chiappa and Gillam, 2019) aims to satisfy the following individual-level fairness criterion:

Definition 1 (Wu et al. (2019b)) *Given unfair pathways π in a causal graph, classifier h_θ achieves a (path-specific) individual-level fairness if*

$$\mathbb{E}_{Y_{A \leftarrow 0}, Y_{A \leftarrow 1 \parallel \pi}} [Y_{A \leftarrow 1 \parallel \pi} - Y_{A \leftarrow 0} | \mathbf{X} = \mathbf{x}] = 0 \quad (5)$$

holds for any value of \mathbf{x} of input features \mathbf{X} .

Condition (5) states that classifier h_θ is individually fair if the *conditional mean unfair effect* is zero, which is an average over individuals who have identical attributes for all features in \mathbf{X} . Since $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$ are expressed using classifier parameter θ as shown in (4), we need to find appropriate θ values to satisfy (5).

Unfortunately, such θ values can be found only in restricted cases. As pointed out by Wu et al. (2019b), this is because we cannot always estimate the conditional mean unfair effect in (5). For instance, when potential

¹We can also consider different potential outcomes $Y_{A \leftarrow 1}$ and $Y_{A \leftarrow 0 \parallel \pi}$, where all inputs of $Y_{A \leftarrow 1}$ are given using the value $a = 1$, and $Y_{A \leftarrow 0 \parallel \pi}$ is formulated using $a = 0$ only for the inputs that correspond to the nodes on pathways π .

outcomes are given as (4), since the conditional mean unfair effect is conditioned on A and D , estimating it requires the joint distribution of $D(0)$ and $D(1)$. This joint distribution, however, is unavailable because we cannot jointly obtain them as explained in Section 2.2.

Due to this issue, the PSCF method (and the one in (Kusner et al., 2017, Section S4)) can achieve individual-level fairness only when the data are generated from a restricted functional class of SEMs. Specifically, these existing methods assume that each variable V follows an additive noise model $V = f_V(\mathbf{pa}(V)) + U_V$, where $\mathbf{pa}(V)$ denotes the parents of V in the causal graph. Unfortunately, this model cannot express data-generating processes in many cases. For instance, it cannot express variable D in (2) due to multiplicative noise U_D . Traditionally, this assumption has been used to infer causal graphs (Hoyer et al., 2009; Shimizu et al., 2006). However, as mentioned in Glymour et al. (2019), since more recent causal graph discovery methods require much weaker assumptions (Zhang and Hyvärinen, 2009; Stegle et al., 2010), the presence of such an assumption severely restricts the scope of their applications.

3.2 Another Method for Removing Unfair Effects

To avoid the aforementioned restrictive functional assumption, the FIO method (Nabi and Shpitser, 2018) aims to remove the *mean unfair effect* over **all** individuals, which is expressed as

$$\begin{aligned} & \mathbb{E}_{Y_{A \leftarrow 0}, Y_{A \leftarrow 1} \parallel \pi} [Y_{A \leftarrow 1} \parallel \pi - Y_{A \leftarrow 0}] \\ & = \mathbb{P}(Y_{A \leftarrow 1} \parallel \pi = 1) - \mathbb{P}(Y_{A \leftarrow 0} = 1). \end{aligned} \quad (6)$$

In (6), marginal probabilities $\mathbb{P}(Y_{A \leftarrow 0} = 1)$ and $\mathbb{P}(Y_{A \leftarrow 1} \parallel \pi = 1)$ can be estimated under much weaker assumptions than the conditional mean unfair effect in (5). We detail these assumptions in Appendix B and how to derive the estimators in Appendix D.

However, removing this mean unfair effect does not imply individual-level fairness. This is because depending on input features \mathbf{X} , unfair effects might be largely positive for some individuals and largely negative for others, which is seriously discriminatory for these individuals. Note that we cannot resolve this issue simply using e.g., the mean of the absolute values of the unfair effects. This is because estimating such a quantity requires a joint distribution of $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1} \parallel \pi$; however, this joint distribution is unavailable because we cannot obtain both $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1} \parallel \pi$ for each individual without an SEM.

4 PROPOSED METHOD

4.1 Overcoming Weaknesses of Existing Methods

To resolve the weaknesses of the existing methods, we propose a framework that guarantees individual-level fairness without restrictive functional assumptions.

For this goal, we aim to train a classifier by forcing an unfair effect to be zero for **all** individuals: i.e., making potential outcomes take the same value (i.e., $Y_{A \leftarrow 0} = Y_{A \leftarrow 1} \parallel \pi = 0$ or $Y_{A \leftarrow 0} = Y_{A \leftarrow 1} \parallel \pi = 1$) with probability 1 regardless of the values of input features \mathbf{X} . This is sufficient to satisfy the individual-level fairness condition (Definition 1) because it restricts the potential outcome values more severely than the latter condition, where potential outcomes $Y_{A \leftarrow 0} = Y_{A \leftarrow 1} \parallel \pi$ can take 0 or 1 depending on \mathbf{X} 's values. Although such a fairness condition may be overly severe and might decrease the prediction accuracy, in Section 5 we experimentally show that our method can achieve comparable accuracy to the existing method for ensuring individual-level fairness (i.e., the PSCF method (Chiappa and Gillam, 2019)).

Compared with PSCF, our method has a clear advantage in that it requires much weaker assumptions. We only need to estimate the marginal potential outcome probabilities in (6), which only requires several conditional independence relations and the graphical condition on unfair pathways π (see Appendix B for our assumptions). Furthermore, we can relax these assumptions to address some cases where there are unobserved variables called latent confounders (Section 4.5).

4.2 Achieving Individual-Level Fairness with PIU

We aim to make potential outcomes take the same value for all individuals. To this end, we formulate penalty function G_θ based on the following quantity:

Definition 2 For unfair pathways π in a causal graph and potential outcomes $Y_{A \leftarrow 0}, Y_{A \leftarrow 1} \parallel \pi \in \{0, 1\}$, we define the **probability of individual unfairness (PIU)** by $\mathbb{P}(Y_{A \leftarrow 0} \neq Y_{A \leftarrow 1} \parallel \pi)$.

Intuitively, PIU is the probability that potential outcomes $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1} \parallel \pi$ take different values.

Unlike the conditional mean unfair effect in Definition 1, PIU is not conditioned on features \mathbf{X} of each individual.

Nonetheless, PIU can be used to guarantee individual-level fairness. By constraining PIU to zero, we can

guarantee that potential outcomes take the same value (i.e., $Y_{A \leftarrow 0} = Y_{A \leftarrow 1 \parallel \pi} = 0$ or $Y_{A \leftarrow 0} = Y_{A \leftarrow 1 \parallel \pi} = 1$) with probability 1 regardless of the values of \mathbf{X} , which is sufficient to ensure individual-level fairness.

Unfortunately, we cannot directly impose constraints on PIU. This is because estimating the PIU value requires the joint distribution of $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$, which is unavailable as described in Section 3.2.

To overcome this issue, instead of PIU, we utilize its upper bound that can be estimated from data. Specifically, to make the PIU value close to zero, we formulate a penalty function that forces the upper bound on PIU to be nearly zero, which is described in the next section.

4.3 Penalty By Upper Bound on PIU

4.3.1 Upper Bound Formulation

To make the PIU value small, we utilize the following upper bound on PIU:

Theorem 1 (Upper bound on PIU) *Suppose that potential outcomes $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$ are binary. Then for any joint distribution of potential outcomes $P(Y_{A \leftarrow 0}, Y_{A \leftarrow 1 \parallel \pi})$, PIU is upper bounded as follows:*

$$P(Y_{A \leftarrow 0} \neq Y_{A \leftarrow 1 \parallel \pi}) \leq 2P^I(Y_{A \leftarrow 0} \neq Y_{A \leftarrow 1 \parallel \pi}), \quad (7)$$

where P^I is an independent joint distribution, i.e., $P^I(Y_{A \leftarrow 0}, Y_{A \leftarrow 1 \parallel \pi}) = P(Y_{A \leftarrow 0})P(Y_{A \leftarrow 1 \parallel \pi})$.

The proof is detailed in Appendix C. Theorem 1 states that whatever joint distribution potential outcomes $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$ follow, the resulting PIU value is at most twice the PIU value that is approximated with independent joint distribution P^I .

Note that this upper bound can be larger than 1, and if so, the PIU value is not controlled because PIU is at most 1. However, since PIU is always smaller than its upper bound, by making the upper bound close to zero, we can guarantee that PIU is also close to zero.

4.3.2 Estimating Upper Bound

Using the observed data, we estimate the upper bound on PIU in (7), which is twice the value of the approximated PIU. Recall that this approximated PIU is the probability that potential outcomes $Y_{A \leftarrow 0}$ and $Y_{A \leftarrow 1 \parallel \pi}$ take different values when they are independent. Since potential outcomes are binary, it is expressed as the probability that potential outcome values

are $(Y_{A \leftarrow 0}, Y_{A \leftarrow 1 \parallel \pi}) = (0, 1)$ or $(1, 0)$; in other words,

$$\begin{aligned} & P^I(Y_{A \leftarrow 0} \neq Y_{A \leftarrow 1 \parallel \pi}) \\ &= P(Y_{A \leftarrow 1 \parallel \pi} = 1)(1 - P(Y_{A \leftarrow 0} = 1)) \\ & \quad + (1 - P(Y_{A \leftarrow 1 \parallel \pi} = 1))P(Y_{A \leftarrow 0} = 1). \end{aligned} \quad (8)$$

Various estimators can be used to estimate marginal probabilities $P(Y_{A \leftarrow 0} = 1)$ and $P(Y_{A \leftarrow 1 \parallel \pi} = 1)$ in (8). Among them, we utilize the computationally efficient estimator in Huber (2014), which can be computed in $O(n)$ time, where n is the number of training instances.

Let $c_\theta(\mathbf{X}) = P(Y = 1 | \mathbf{X})$ denote the conditional distribution provided by classifier h_θ ; we let $c_\theta(\mathbf{X}) = h_\theta(\mathbf{X}) \in \{0, 1\}$ if h_θ is a deterministic classifier. For instance, suppose that the causal graph is given as shown in Figure 1(b) and that the features of n individuals are provided as $\{\mathbf{x}_i\}_{i=1}^n = \{a_i, q_i, d_i, m_i\}_{i=1}^n$. Then $P(Y_{A \leftarrow 0} = 1)$ and $P(Y_{A \leftarrow 1 \parallel \pi} = 1)$ can be estimated as the following weighted averages:

$$\begin{aligned} \hat{p}_\theta^{A \leftarrow 0} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(a_i = 0) \hat{w}_i c_\theta(a_i, q_i, d_i, m_i) \text{ and} \\ \hat{p}_\theta^{A \leftarrow 1 \parallel \pi} &= \frac{1}{n} \sum_{i=1}^n \mathbf{1}(a_i = 1) \hat{w}'_i c_\theta(a_i, q_i, d_i, m_i), \end{aligned} \quad (9)$$

where $\mathbf{1}(\cdot)$ is an indicator function, and \hat{w}_i and \hat{w}'_i are the following weights for individual $i \in \{1, \dots, n\}$:

$$\begin{aligned} \hat{w}_i &= \frac{1}{\hat{P}(A = 0 | q_i)}, \\ \hat{w}'_i &= \frac{\hat{P}(A = 1 | q_i, d_i) \hat{P}(A = 0 | q_i, d_i, m_i)}{\hat{P}(A = 1 | q_i) \hat{P}(A = 0 | q_i, d_i) \hat{P}(A = 1 | q_i, d_i, m_i)}, \end{aligned}$$

where \hat{P} is a conditional distribution, which we infer in the same way as Zhang and Bareinboim (2018a), i.e., by learning a statistical model (e.g., a neural network) from the training data beforehand.² We derive the estimators (9) in Appendix D.

In (9), marginal probabilities $\hat{p}_\theta^{A \leftarrow 0}$ and $\hat{p}_\theta^{A \leftarrow 1 \parallel \pi}$ are estimated by taking a weighted average of conditional probability c_θ over individuals with $A = 0$ and $A = 1$, respectively, which is a widely used estimation technique called *inverse probability weighting* (IPW).

4.3.3 Formulating Penalty Function

To learn an individually fair classifier, we force the estimated value of the upper bound on PIU to be close

²Note that FIO infers conditional distributions not by learning statistical models beforehand but by simultaneously learning them with the predictive model of Y (Nabi and Shpitser, 2018). This is because unlike our method, it addresses not only training a classifier but also learning a generative model of joint distribution $P(\mathbf{X}, Y)$.

to zero by minimizing the objective function (1) with the following penalty function G_θ :

$$G_\theta(\mathbf{x}_1, \dots, \mathbf{x}_n) = \hat{p}_\theta^{A \leftarrow 1 \parallel \pi} (1 - \hat{p}_\theta^{A \leftarrow 0}) + (1 - \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}) \hat{p}_\theta^{A \leftarrow 0}. \quad (10)$$

For instance, if $\hat{p}_\theta^{A \leftarrow 0}$ and $\hat{p}_\theta^{A \leftarrow 1 \parallel \pi}$ are given as the weighted estimators (9), to reduce the value of G_θ , our method imposes strong penalties on the predictions for individuals whose weights \hat{w}_i and \hat{w}'_i are large.

In our experiments, we minimize the objective function using the stochastic gradient descent method (Sutskever et al., 2013). We discuss the computation time and convergence guarantees in Appendix E.

From the penalty function (10), we can see why penalizing the upper bound on PIU guarantees individual-level fairness. As the penalty parameter value goes to infinity ($\lambda \rightarrow \infty$), the marginal probabilities ($\hat{p}_\theta^{A \leftarrow 0}, \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}$) approach (0, 0) or (1, 1). This guarantees that the potential outcomes take the same value with probability 1, which is sufficient to guarantee individual-level fairness.

4.4 Comparison with Existing Fairness Constraint

To show the effectiveness of penalty function (10), we compare it with the constraint of the FIO method.

Suppose that our penalty function forces the upper bound on PIU to satisfy the following condition:

$$\hat{p}_\theta^{A \leftarrow 1 \parallel \pi} (1 - \hat{p}_\theta^{A \leftarrow 0}) + (1 - \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}) \hat{p}_\theta^{A \leftarrow 0} \leq \delta. \quad (11)$$

Here we let constant δ be $\delta \in [0, 1]$ because otherwise we cannot force the PIU value to be less than 1.

Meanwhile, the FIO constraint limits the mean unfair effect (6) to lie in

$$-\delta' \leq \hat{p}_\theta^{A \leftarrow 1 \parallel \pi} - \hat{p}_\theta^{A \leftarrow 0} \leq \delta', \quad (12)$$

where $\delta' \in [0, 1]$ is a hyperparameter. If $\delta' = 0$, it ensures $\hat{p}_\theta^{A \leftarrow 0} = \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}$.

Figure 2 shows the feasible region of our fairness condition (red) and the FIO constraint (blue), respectively, obtained by graphing the hyperbolic inequality in (11) and the linear inequality in (12). Here, to clarify their difference, we consider the case where $\delta = 2\delta'$.

When $\delta \approx 0$, our fairness condition only accepts region ($\hat{p}_\theta^{A \leftarrow 0}, \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}$) \approx (0, 0) or (1, 1), where the potential outcomes are likely to take the same value; hence, the unfair effect is likely to be zero. This demonstrates how effectively our condition removes an unfair effect for each individual. By contrast, with any δ' value, the

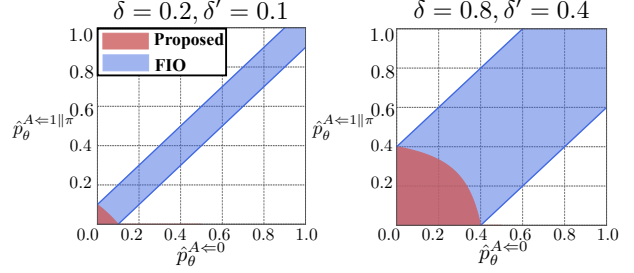


Figure 2: Feasible regions of our constraint (red) and FIO (blue) with $(\delta, \delta') = (0.2, 0.1), (0.8, 0.4)$

FIO constraint always accepts point $(\hat{p}_\theta^{A \leftarrow 0}, \hat{p}_\theta^{A \leftarrow 1 \parallel \pi}) = (0.5, 0.5)$, where it is completely uncertain whether potential outcomes take the same value as detailed in Appendix F. This implies that FIO predictions might be unfair for some individuals, indicating that FIO cannot ensure individual-level fairness.

4.5 Extension for Dealing with Latent Confounders

So far, we have assumed that the marginal probabilities of potential outcomes can be estimated from data. This assumption, however, does not hold if there is a latent confounder (Pearl, 2009), i.e., an unobserved variable that is a parent of the observed variables in the causal graph. Although this is possible in practice, inferring marginal probabilities becomes much more challenging.

However, even in the presence of latent confounders, our method can ensure individual-level fairness if the lower and upper bounds on marginal probabilities are available. For instance, the bounds of Miles et al. (2017) can be used when there is only a single mediator (i.e., a feature affected by sensitive feature A), and it takes discrete values. In such cases, we can achieve individual-level fairness by reformulating the penalty function as follows.

Suppose that marginal probabilities $P(Y_{A \leftarrow 0} = 1)$ and $P(Y_{A \leftarrow 1 \parallel \pi} = 1)$ are bounded by

$$\begin{aligned} \hat{l}_\theta^{A \leftarrow 0} &\leq P(Y_{A \leftarrow 0} = 1) \leq \hat{u}_\theta^{A \leftarrow 0}, \\ \hat{l}_\theta^{A \leftarrow 1 \parallel \pi} &\leq P(Y_{A \leftarrow 1 \parallel \pi} = 1) \leq \hat{u}_\theta^{A \leftarrow 1 \parallel \pi}, \end{aligned}$$

where $\hat{l}_\theta^{A \leftarrow 0}, \hat{u}_\theta^{A \leftarrow 0}, \hat{l}_\theta^{A \leftarrow 1 \parallel \pi},$ and $\hat{u}_\theta^{A \leftarrow 1 \parallel \pi}$ are the estimated lower and upper bounds, which we describe in Appendix G. Then for any marginal probability values, the upper bound on PIU in (7) is always smaller than twice the value of

$$G_\theta(\mathbf{x}_1, \dots, \mathbf{x}_n) = \hat{u}_\theta^{A \leftarrow 1 \parallel \pi} (1 - \hat{l}_\theta^{A \leftarrow 0}) + (1 - \hat{l}_\theta^{A \leftarrow 1 \parallel \pi}) \hat{u}_\theta^{A \leftarrow 0}. \quad (13)$$

Table 2: Test accuracy (%) on each dataset

Method	Synth	German	Adult
Proposed	80.0 ± 0.9	75.0	75.2
FIO	84.8 ± 0.6	78.0	81.2
PSCF	74.8 ± 1.6	76.0	73.4
Unconstrained	88.2 ± 0.9	81.0	83.2
Remove	76.9 ± 1.3	73.0	74.7

Therefore, by making this penalty function value nearly zero, we can achieve individual-level fairness. We experimentally confirmed that this framework makes fairer predictions than the original one in Appendix I.5.

5 EXPERIMENTS

We compared our method (**Proposed**) with the following four baselines: (1) **FIO** (Nabi and Shpitser, 2018), (2) **PSCF** (Chiappa and Gillam, 2019), which aims to achieve individual-level fairness by assuming that the data are generated by additive noise models, (3) **Unconstrained**, which does not use any constraints or penalty terms related to fairness, and (4) **Remove** (Kusner et al., 2017, Section S4), which removes unfair effects simply by making predictions without input features that correspond to the nodes on unfair pathways π . As classifiers of **Proposed**, **FIO**, **Unconstrained**, and **Remove**, we used a feed-forward neural network that contains two linear layers with 100 and 50 hidden neurons. Other settings are detailed in Appendix H.1.

Data and causal graphs: For a performance evaluation, we used a synthetic dataset and two real-world datasets: the German credit dataset and the Adult dataset (Bache and Lichman, 2013). We sampled the synthetic data from the SEM, whose formulation is described in Appendix H.2.1. To define the unfair effect, we used the causal graph in Figure 1(b). With the real-world datasets, we evaluated the performance as follows. With the German dataset, we predicted whether each loan applicant is risky (Y) from their features such as gender A and savings S . With the Adult dataset, we predicted whether an annual income exceeds \$50,000 (Y) from features such as gender A and marital status M . To measure unfairness, following (Chiappa and Gillam, 2019), we used the causal graphs in Figure 3, which we detail in Appendix H.3.1.

Accuracy and fairness: We evaluated the test accuracy and four statistics of the unfair effects: (i) the mean unfair effect (6), (ii) the standard deviation in the conditional mean unfair effects in (5), (iii) the upper bound on PIU, and (iv) the PIU.

Table 2 and Figure 4 present the test accuracy and

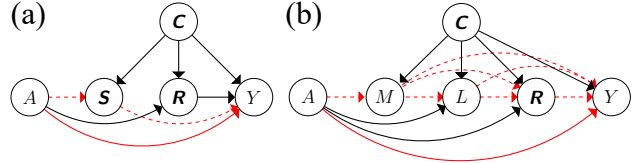


Figure 3: Causal graphs for (a) German credit dataset and (b) Adult dataset: Direct pathways with red solid edges are unfair. Red dashed pathway $A \rightarrow S \rightarrow Y$ is unfair in (a), and those that go through M (i.e., $A \rightarrow M \rightarrow \dots \rightarrow Y$) are unfair in (b).

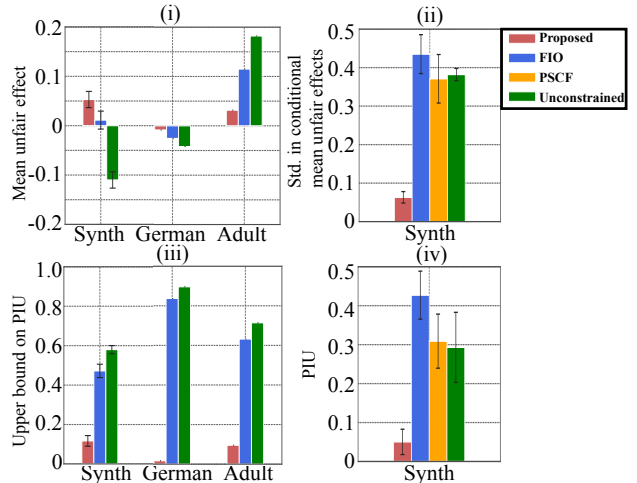


Figure 4: Four statistics of unfair effects on test data: The closer they are to zero, the fairer predictions are. Error bars of synthetic data (Synth) denote standard deviations in 10 runs with randomly generated data. With **Remove**, all statistics are zero (not shown). With **PSCF**, (i) and (iii) are not well-defined as described in Appendix H.2.3.

the four statistics of the unfair effects, respectively. In synthetic data experiments, we computed the means and the standard deviations based on 10 experiments with randomly generated training and test data. Figure 4 does not display two statistics (ii) and (iv) for the German and Adult datasets because computing them requires SEMs, which are unavailable for these real-world datasets. We do not show (i) or (iii) of **PSCF** because these are not well-defined for this method (see Appendix H.2.3 for details).

With **Proposed**, all the statistics of the unfair effects were sufficiently close to zero, demonstrating that it made fair predictions for all individuals. This is because by imposing a penalty on the upper bound on PIU, **Proposed** forced unfair effect values to be close to zero for all individuals, guaranteeing that the other statistics are close to zero.

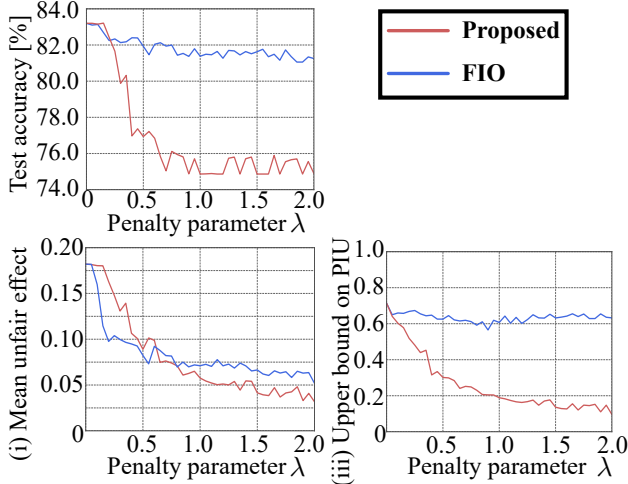


Figure 5: Test accuracy, mean unfair effect, and upper bound on PIU of (a) **Proposed** (red) and (b) **FIO** (blue) for Adult dataset when increasing penalty parameter value as $\lambda = 0, 0.05, 0.10, \dots, 2.00$.

By contrast, regarding **FIO** and **PSCF**, the unfair effect values were much larger. With **FIO**, although the mean unfair effect (i.e., (i)) was close to zero, the other statistics deviated from zero, indicating that constraining the mean unfair effect did not ensure individual-level fairness. **PSCF** failed to reduce the value of the standard deviation in the conditional mean unfair effects (i.e., (ii)). This is because the data are not generated from additive noise models (see Appendix H.2.1 for the data), violating the functional assumption of **PSCF**. Since the large values of (ii) imply that unfair effects are greatly affected by the attributes of input features \mathbf{X} , these results indicate that **FIO** and **PSCF** made unfair predictions based on these attributes. With real-world datasets, **FIO** provided large values of the upper bound on PIU (i.e., (iii)). These upper bound values cast much doubt on whether the predictions of **FIO** are fair for each individual, which is problematic in practice.

The test accuracy of **Proposed** was lower than **FIO**, higher than **Remove**, and comparable to **PSCF**. Since **FIO** imposes a much weaker fairness constraint than **Proposed**, **Remove**, and **PSCF** (i.e., the methods designed for achieving individual-level fairness), it achieved higher accuracy. By contrast, since **Remove** eliminates all informative input features that are affected by the sensitive feature to guarantee individual-level fairness, it provided the lowest accuracy. The comparison of **Proposed** and **PSCF** indicates that although our method employs a more severe fairness condition than **PSCF**, it barely sacrifices prediction accuracy, demonstrating that it strikes a better balance between individual-level fairness and accuracy.

Penalty parameter effects: Using the Adult dataset, we further compared the performance of **Proposed** with **FIO** using various penalty parameter values. Regarding unfair effects, we evaluated two statistics (i) the mean unfair effect and (iii) the upper bound on PIU since the others are unavailable with this real-world dataset, as already described above.

Figure 5 presents the results. As the penalty parameter value increased, the test accuracy of **Proposed** decreased more sharply than **FIO** since **Proposed** imposed a stronger fairness condition. However, with **Proposed**, both (i) and (iii) dropped to nearly zero, while only (i) decreased with **FIO**. This implies that while it remains uncertain whether the predictions of **FIO** are individually fair, the predictions of **Proposed** are more reliable since it can successfully reduce the upper bound on PIU (i.e., (iii)) and guarantee individual-level fairness, which is helpful for practitioners.

Additional experimental results: To further demonstrate the effectiveness of our method, we present several additional experimental results in Appendix I. Through our experiments, we show that our method also worked well using other classifier than the neural network (Appendix I.1), present the statistical significance of the test accuracy (Appendix I.2), confirm that both our method and **PSCF** achieved individual-level fairness when the data satisfy the functional assumptions (Appendix I.3), demonstrate the tightness of our upper bound on PIU (Appendix I.4), and evaluate the performance of our extended framework for addressing latent confounders (Appendix I.5).

6 RELATED WORK

Causality-based fairness: Motivated by recent developments in inferring causal graphs (Chikahara and Fujino, 2018; Glymour et al., 2019), many causality-based approaches to fairness have been proposed (Chippa and Gillam, 2019; Kilbertus et al., 2017; Kusner et al., 2017, 2019; Nabi and Shpitser, 2018; Nabi et al., 2019; Russell et al., 2017; Salimi et al., 2019; Wu et al., 2018, 2019a; Xu et al., 2019; Zhang et al., 2017; Zhang and Wu, 2017; Zhang et al., 2018; Zhang and Bareinboim, 2018a,b). However, few are designed for making individually fair predictions due to the difficulty of estimating the conditional mean unfair effect in (5). Although *path-specific counterfactual fairness* (PC-fairness) (Wu et al., 2019b) was proposed to estimate its lower and upper bounds, it only addresses measuring unfairness in data and not making fair predictions.

By contrast, we established a learning framework for making individually fair predictions under much weaker assumptions than the existing approaches.

Bounding PIU: To learn an individually fair classifier, our framework utilizes the upper bound on PIU in (7). Compared with the existing bounds described below, this upper bound has the following two advantages.

First, it can be used for binary potential outcomes. If we consider continuous potential outcomes, we can use several bounds on a functional of the joint distribution of potential outcomes (Fan et al., 2017; Firpo and Ridder, 2019) because PIU is also such a functional. However, these existing bounds cannot be used in our binary classification setting where PIU is formulated based on binary potential outcomes.

Second, it is much tighter than the result in Rubinstein and Singla (2017), which provides an upper bound on an expectation of random variables whose joint distribution is unavailable.³ Since PIU can be written as an expectation (i.e., $\mathbb{E}_{Y_{A \leftarrow 0}, Y_{A \leftarrow 1} | \pi} [\mathbf{1}(Y_{A \leftarrow 0} \neq Y_{A \leftarrow 1} | \pi)]$), we can apply this result to PIU; however, the bound becomes much looser than ours. Although the multiplicative constant in (7) is 2, this value becomes 200 with the bound of Rubinstein and Singla (2017). If we use such a loose upper bound, we need to impose an excessively severe penalty on it to ensure that PIU is close to zero. By contrast, using our upper bound, we can avoid imposing such a severe penalty, thus preventing an unnecessary decrease in prediction accuracy.

7 CONCLUSION

We proposed a learning framework for guaranteeing individual-level fairness without impractical functional assumptions. Based on the concept called path-specific causal effects, we defined PIU and derived its upper bound that can be estimated from data. By forcing this upper bound value to be nearly zero, our proposed method trains an individually fair classifier. We experimentally show that our method makes individually fairer predictions than the existing methods at a slight cost of accuracy, indicating that it strikes a better balance between fairness and accuracy.

References

- Shipra Agrawal, Yichuan Ding, Amin Saberi, and Yinyu Ye. Correlation robust stochastic optimization. In *SODA*, pages 1087–1096, 2010.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. *Machine Bias*, 2016. <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing>.
- Chen Avin, Ilya Shpitser, and Judea Pearl. Identifiability of path-specific effects. In *IJCAI*, pages 357–363, 2005.
- K. Bache and M. Lichman. UCI machine learning repository: Datasets. <http://archive.ics.uci.edu/ml/datasets>, 2013.
- Silvia Chiappa and Thomas PS Gillam. Path-specific counterfactual fairness. In *AAAI*, pages 7801–7808, 2019.
- Yoichi Chikahara and Akinori Fujino. Causal inference in time series via supervised learning. In *IJCAI*, pages 2042–2048, 2018.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *ITCS*, pages 214–226, 2012.
- Yanqin Fan, Emmanuel Guerre, and Dongming Zhu. Partial identification of functionals of the joint distribution of "potential outcomes". *Journal of Econometrics*, 197(1):42–59, 2017.
- Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *KDD*, pages 259–268, 2015.
- Sergio Firpo and Geert Ridder. Partial identification of the treatment effect distribution and its functionals. *Journal of Econometrics*, 213(1):210–234, 2019.
- Clark Glymour, Kun Zhang, and Peter Spirtes. Review of causal discovery methods based on graphical models. *Frontiers in Genetics*, 10, 2019.
- Moritz Hardt, Eric Price, Nati Srebro, et al. Equality of opportunity in supervised learning. In *NeurIPS*, pages 3315–3323, 2016.
- Kimberly A Houser. Can AI solve the diversity problem in the tech industry: Mitigating noise and bias in employment decision-making. *Stan. Tech. L. Rev.*, 22:290, 2019.
- Patrik O Hoyer, Dominik Janzing, Joris M Mooij, Jonas Peters, and Bernhard Schölkopf. Nonlinear causal discovery with additive noise models. In *NeurIPS*, pages 689–696, 2009.
- Martin Huber. Identifying causal mechanisms (primarily) based on inverse probability weighting. *Journal of Applied Econometrics*, 29(6):920–943, 2014.
- Amir E Khandani, Adlar J Kim, and Andrew W Lo. Consumer credit-risk models via machine-learning algorithms. *Journal of Banking and Finance*, 34(11): 2767–2787, 2010.
- Niki Kilbertus, Mateo Rojas Carulla, Giambattista Parascandolo, Moritz Hardt, Dominik Janzing, and Bernhard Schölkopf. Avoiding discrimination through causal reasoning. In *NeurIPS*, pages 656–666, 2017.

³Such an upper bound is known as *correlation gap* in the field of robust optimization. (Agrawal et al., 2010).

- Matt Kusner, Chris Russell, Joshua Loftus, and Ricardo Silva. Making decisions that reduce discriminatory impacts. In *ICML*, pages 3591–3600, 2019.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *NeurIPS*, pages 4066–4076, 2017.
- Caleb H Miles, Phyllis Kanki, Seema Meloni, and Eric J. Tchetgen Tchetgen. On partial identification of the natural indirect effect. *Journal of Causal Inference*, 5(2), 2017.
- Razieh Nabi and Ilya Shpitser. Fair inference on outcomes. In *AAAI*, pages 1931–1940, 2018. <https://github.com/raziehna/fair-inference-on-outcomes>.
- Razieh Nabi, Daniel Malinsky, and Ilya Shpitser. Learning optimal fair policies. In *ICML*, pages 4674–4682, 2019.
- Judea Pearl. *Causality: Models, Reasoning and Inference*. Cambridge University Press, 2009.
- Aviad Rubinfeld and Sahil Singla. Combinatorial prophet inequalities. In *SODA*, pages 1671–1687, 2017.
- Chris Russell, Matt J Kusner, Joshua Loftus, and Ricardo Silva. When worlds collide: integrating different counterfactual assumptions in fairness. In *NeurIPS*, pages 6414–6423, 2017.
- Babak Salimi, Luke Rodriguez, Bill Howe, and Dan Suciu. Interventional fairness: Causal database repair for algorithmic fairness. In *SIGMOD*, pages 793–810, 2019.
- Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, and Antti Kerminen. A linear non-gaussian acyclic model for causal discovery. *JMLR*, 7(Oct):2003–2030, 2006.
- Oliver Stegle, Dominik Janzing, Kun Zhang, Joris M Mooij, and Bernhard Schölkopf. Probabilistic latent variable models for distinguishing between cause and effect. In *NeurIPS*, pages 1687–1695, 2010.
- Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *ICML*, pages 1139–1147, 2013.
- Yongkai Wu, Lu Zhang, and Xintao Wu. On discrimination discovery and removal in ranked data using causal graph. In *KDD*, pages 2536–2544, 2018.
- Yongkai Wu, Lu Zhang, and Xintao Wu. Counterfactual fairness: Unidentification, bound and algorithm. In *IJCAI*, 2019a.
- Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. PC-fairness: A unified framework for measuring causality-based fairness. In *NeurIPS*, pages 3399–3409, 2019b.
- Depeng Xu, Yongkai Wu, Shuhan Yuan, Lu Zhang, and Xintao Wu. Achieving causal fairness through generative adversarial networks. In *IJCAI*, pages 1452–1458, 2019.
- Junzhe Zhang and Elias Bareinboim. Equality of opportunity in classification: A causal approach. In *NeurIPS*, 2018a.
- Junzhe Zhang and Elias Bareinboim. Fairness in decision-making—the causal explanation formula. In *AAAI*, 2018b.
- K Zhang and A Hyvärinen. On the identifiability of the post-nonlinear causal model. In *UAI*, pages 647–655, 2009.
- Lu Zhang and Xintao Wu. Anti-discrimination learning: a causal modeling-based framework. *International Journal of Data Science and Analytics*, 4(1):1–16, 2017.
- Lu Zhang, Yongkai Wu, and Xintao Wu. A causal framework for discovering and removing direct and indirect discrimination. In *IJCAI*, 2017.
- Lu Zhang, Yongkai Wu, and Xintao Wu. Causal modeling-based discrimination discovery and removal: criteria, bounds, and algorithms. *IEEE Transactions on Knowledge and Data Engineering*, 2018.