# Causal Fairness-Guided Dataset Reweighting using Neural Networks

Xuan Zhao
*SCHUFA Holding AG*
Germany
xuan.zhao@schufa.de

Klaus Broelemann
*SCHUFA Holding AG*
Germany
Klaus.Broelemann@schufa.de

Salvatore Ruggieri
*University of Pisa*
Italy
salvatore.ruggieri@unipi.it

Gjergji Kasneci
*Technical University of Munich*
Germany
gjergji.kasneci@tum.de

*Abstract*—The importance of achieving fairness in machine learning models cannot be overstated. Recent research has pointed out that fairness should be examined from a causal perspective, and several fairness notions based on the on Pearl's causal framework have been proposed. In this paper, we construct a reweighting scheme of datasets to address causal fairness. Our approach aims at mitigating bias by considering the causal relationships among variables and incorporating them into the reweighting process. The proposed method adopts two neural networks, whose structures are intentionally used to reflect the structures of a causal graph and of an interventional graph. The two neural networks can approximate the causal model of the data, and the causal model of interventions. Furthermore, reweighting guided by a discriminator is applied to achieve various fairness notions. Experiments on real-world datasets show that our method can achieve causal fairness on the data while remaining close to the original data for downstream tasks.

.

*Index Terms*—causal fairness, data pre-processing, adversarial reweighting.

## I. INTRODUCTION

Pre-processing data to satisfy fairness requirements is an important research question in machine learning. Models trained on biased data may learn such biases and generalize them, thus leading to discriminatory decisions against socially sensitive groups defined on the grounds of gender, race and age, or other protected grounds [1]–[5]. Many methods have been proposed to modify the training data in order to mitigate biases and to achieve specific fairness requirements [4], [6]–[11].

For reliable and effective treatment, particularly in a legal context, discrimination claims usually require demonstrating causal relationships between sensitive attributes and questionable decisions (or predictions), instead of mere associations or correlations. Compared with the fairness notions based on correlation, causality-based fairness notions and methods include additional knowledge of the causal structure of the problem. This knowledge often reveals the mechanism of data generation, which helps comprehend and interpret the influence of sensitive attributes on the output of a decision process. Causal fairness seeks to address the root causes of disparities rather than simply trying to eliminate them in a post-hoc manner.

We draw upon the ideas and concepts presented in CF-GAN [12] as the framework for our research. Instead of fair dataset generation in CFGAN, however, we propose a method which reweighs the samples to achieve fairness criteria with the help of two neural networks to reflect the causal and interventional graphs, and a discriminator to guide the reweighting. As the general requirement of modifying datasets is to preserve the data utility as much as possible for the downstream tasks. The intuition of the reweighting scheme is that in a given dataset, there are individuals who are treated 'fairer' in the causal mechanism and by assigning higher weights to these individuals, we could slightly alter the underlying causal mechanism to achieve fairness and do not influence much on the performance of downstream tasks. In this case, hopefully we could mitigate the historical bias. In addition, by analyzing the high/low weights assigned to samples, a reweighting method like ours enables for a high-level understanding the biases.

The experiments (Section IV) show that reweighed data outperform generated data in utility. In the taxonomy of pre-processing, in-processing and post-processing methods for bias mitigation [3], [13]–[16], our method falls into the category of pre-processing, as we deal with the dataset before it is given in input to the downstream learning algorithm. Thus, our approach is model-agnostic, as any pre-processing method.

We summarize our contribution as follows: (1) We formulate a novel and sample-based reweighting method for mitigating different causal bias related to sensitive groups. (2) We show that by simulating the underlying causal model that reflects the causal relations of the real data, and the causal model after the intervention, with the help of a discriminator, our reweighting approach leads to fair reweighted data. (3) We provide a thorough evaluation of the proposed technique on benchmark datasets and show the viability of our approach.

## II. PRELIMINARY

Throughout this paper, we consider a structural causal model $\mathcal{M} = \langle U, V, F \rangle$, that is learned from a dataset $\mathcal{D} = \{(s_k, x_k, y_k)\}_{k=1}^{m}$ where $s_k \in S = \{0,1\}, x_k \in X \subseteq \mathbb{R}^d, y_k \in Y = \{0,1\}$.

1) $U$ denotes exogenous variables that cannot be observed but constitute the background knowledge behind the model. $P(U)$ is a joint probability distribution of the variables in $U$.

2) $V$ denotes endogenous variables that can be observed. In our work, we set $V = \{S, X, Y\}$. $S$ represents the sensitive

attribute, $Y$ represents the outcome attribute, and $X$ represents all other attributes. Additionally, $s^+$ is used to denote $S = 1$ and $s^-$ to denote $S = 0$.

3) $F$ denotes the deterministic functions. For each $V_i \in V$, there is a corresponding function $f_{V_i}$ that maps from domains of the variables in $Pa_{V_i} \cup U_{V_i}$ to $V_i$, namely $V_i = f_{V_i}(Pa_{V_i}, U_{V_i})$. Here, $Pa_{V_i} \subseteq V \backslash V_i$ represents the parents of $V_i$, and $U_{V_i}$ also represents the parents (exogenous variables) of $V_i$, $U_{V_i} \subseteq U$.

We denote by $\mathcal{G}$ the causal graph $\mathcal{G}$ associate with $\mathcal{M}$, and assume it is a Directed Acyclic Graph (DAG).

### A. Causal Fairness Criteria

To understand causal effects in the causal model $\mathcal{M}$, we can use the do-operator [17], which represents a physical intervention that sets a variable $S \in V$ to a constant value $s$. By performing an intervention $do(S = s)$, we replace the original function $S = f_S(Pa_S, U_S)$ with $S = s$. This results in a change in the distribution of all variables that are descendants of $S$ in the causal graph. $\mathcal{M}_s$ is the interventional causal model and its corresponding graph $\mathcal{G}_s$ the interventional graph. In $\mathcal{G}_s$, edges to $S$ are deleted according to the definition of intervention and $S$ is replaced with constant $s$. The interventional distribution for $Y$ is denoted by $P(Y|do(S = s))$. Using the do-operator, we can compare the interventional distributions under different interventions to infer the causal effect of $S$ on $Y$. In this paper, we focus on the following causal causal fairness notions:

*a) Total effect:* The total effect infers the causal effect of $S$ on $Y$ through all possible causal paths from $S$ to $Y$. The total effect of the difference of $s^-$ to $s^+$ on $Y$ is given by $TE(s^+, s^-) = P(Y_{s^+}) - P(Y_{s^-})$, where $P(\cdot)$ here refers to the interventional distribution probability. Total fairness is satisfied if $|TE(s^+, s^-)| < \tau$ ($\tau$ is the fairness threshold). Note that statistical parity is similar to total effect but is fundamentally different. Statistical parity measures the conditional distributions of $Y$ change of the sensitive attribute from $s^-$ to $s^+$.

*b) Path-specific fairness:* The path-specific effect is a fine-grained assessment of causal effects, that is, it can evaluate the causal effect transmitted along certain paths. It is used to distinguish among direct discrimination, indirect discrimination, and explainable bias. It infers the causal effect of $S$ on $Y$ through a subset of causal paths from $S$ to $Y$, which is referred to as the $\pi$-specific effect denoting the subset of causal paths as $\pi$. The specific effect of a path set $\pi$ on $Y$, caused by changing the value of $S$ from $s^-$ to $s^+$ with reference to $s^-$, is given by the difference of the interventional distributions: $SE_\pi(s^+, s^-) = P(Y_{s^+|\pi, s^-|\overline{\pi}}) - P(Y_{s^-})$, where $P(Y_{s^+|\pi, s^-|\overline{\pi}})$ represents the distribution resulting from intervening $do(s^+)$ only along the paths in $\pi$ while $s-$ is used as a reference through other paths $\overline{\pi}$. If $\pi$ contains all direct edge from $S$ to $Y$, $SE_\pi(s^+, s^-)$ measures the direct discrimination. If $\pi$ contains all indirect paths from $S$ to $Y$ that pass through proxy attributes, $SE_\pi(s^+, s^-)$ evaluates

the indirect discrimination. Path-specific fairness is met if $|SE_\pi(s^+, s^-)| < \tau$.

*c) Counterfactual fairness:* The counterfactual effect of changing $S$ from $s^-$ to $s^+$ on $Y$ under certain conditions $O = o$ (where $O$ is a subset of observed attributes $O \subseteq X$) for an individual with features $o$ is given by the difference between the interventional distributions $P(Y_{s^+}|o)$ and $P(Y_{s^-}|o)$: $CE(s^+, s^-|o) = P(Y_{s^+}|o) - P(Y_{s^-}|o)$. Counterfactual fairness is met if $|CE(s^+, s^-|o)| < \tau$. Any context $O = o$ represents a certain sub-group of the population, specifically, when $O = X$, it represents specific individual(s).

### B. Causal Discovery

Methods for extracting a causal graph from given data (causal discovery) can be broadly categorized into two constraint-based and score-based methods [18], [19]. Constraint-based methods, such as [20]–[22], utilize conditional independence tests under specific assumptions to determine the Markov equivalence class of causal graphs. Score-based methods, like [23], evaluate candidate graphs using a pre-defined score function and search for the optimal graph within the space of DAGs. Such an approach is formulated as a combinatorial optimization problem:

$$\min_{\mathcal{G}} Score(\mathcal{G}; V) = \mathcal{L}(\mathcal{G}; V) + \lambda \mathcal{R}_{sparse}(\mathcal{G}),$$
$$s.t. \ \mathcal{G} \in DAG \tag{1}$$

In the realm of causal discovery, the problem can be divided into two components, which constrain the score function $Score(\mathcal{G}; V)$ and $\mathcal{G} \in DAG$. The score function is comprised of: (1) the goodness-of-fit $\mathcal{L}(\mathcal{G}; V) = \frac{1}{m}\sum_{k=1}^m l(v_k, F(v_k))$ is the loss of fitting observation of $v_k$; $F$ denotes the deterministic functions as defined earlier in Section II (2) the sparsity $\mathcal{R}_{sparse}(\mathcal{G})$ which regulates the number of edges in $\mathcal{G}$. $\lambda$ serves as a hyperparameter that controls the regularization strengths.

In this work, we assume that the given causal graph $\mathcal{G}$ is learned from a score-based causal discovery, so $\mathcal{G}$ should have goodness-of-fit and sparsity.

### C. Intervention through Controlled Neural Networks

In CausalGAN [24], a noise vector $Z$ is partitioned into $\{Z_{V_1}, Z_{V_2}, ..., Z_{V_{|V|}}\}$ to mimic the exogenous variables $U$ in the structural causal model $\mathcal{M}$ described in Section II. The generator $G(Z)$ contains $|V|$ sub-neural networks $\{G_{V_1}, G_{V_2}, ..., G_{V_{|V|}}\}$ to generate the values of each node $V_i$ in the graph. The input of $G_{V_i}$ is the output of $G_{Pa_{V_i}}$ combined with $Z_{V_i}$. Here, $G_{V_i}$ is trying to approximate the corresponding function $f_{V_i}(Pa_{V_i}, U_{V_i})$ in the causal model $\mathcal{M}$. The adversarial game is played to ensure that the generated observational distribution is not differentiable from the real observational distribution. In the work of CFGAN, two generators are used to simulate the causal model $\mathcal{M}$ and the interventional model $\mathcal{M}_s$, while two discriminators try to maintain that: (1) the generated data is close to the orginal distribution, and (2) the causal effect is mitigated. In our work,
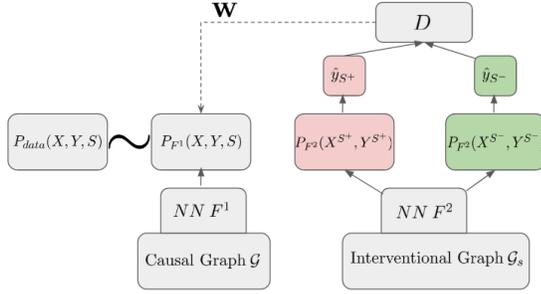
Fig. 1. The framework of reweighting: the structure of NN (neural network) $F^1$ reflects the original causal graph $\mathcal{G}$; the structure of NN $F^2$ refelects the interventional causal graph $\mathcal{G}_s$; the discriminator $D$ tells if a $\hat{y}$ estimated by $F^2$ is from the group $S^+$ or the group $S^-$. An adverserial game is played between the reweighting on the data samples and $D$ to reach a situation where $D$ is not capable of differentiating whether $y$ is from $S^+$ or $S^-$ and a specific causal fairness is reached. The weights of samples are also forwarded to $F^1$ to make sure that the reweighted empirical data distribution is close to the original data distribution from which the causal graph $\mathcal{G}$ is learned.

we also use a similar design but we do not model the noise $Z$ since our goal is not to generate fairness-aware data, but to reweigh the given data.
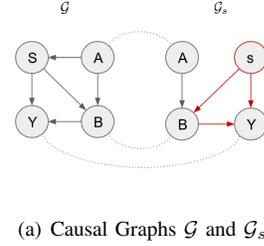
## III. A REWEIGHTING APPROACH FOR DIFFERENT CAUSAL FAIRNESS CRITIRIA
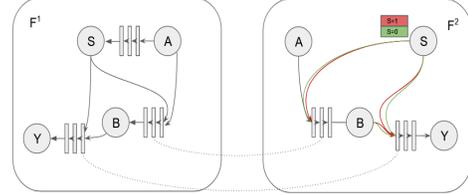
### A. Problem Formulation

As mentioned in Section II, the notation used in our work is based on the conventional approach. We are given a causal graph $\mathcal{G}$ and a dataset $\mathcal{D}$ with $m$ i.i.d. samples drawn from $P(V)$. We assume that $\mathcal{G}$ is sufficient to describe the causal relationships between the variables $V$. In this paper, we build our method on a causal graph of observational data, so we do not specifically model $U$. The problem we are facing is that from the given causal graph $\mathcal{G}$, $S$ has a causal effect on $Y$. Our method aims to achieve two objectives: (1) preserve the goodness-of-fit (mentioned in Section II-B) by maintaining the empirical reweighted data distribution close to the original data distribution for utility of the downstream tasks; and (2) ensure that $S$ cannot be used to discriminate when predicting $Y$ based on various causal criteria in the interventional model $\mathcal{M}_s$. We treat $S$ and $Y$ as binary variables in this paper. However, this can be easily extended to multi-categorical or numerical cases. Also, we focus on the causal effect of $S$ on $Y$, but the model can deal with causal effects among multiple variables. We try to reach the following causal fairness notions mentioned in Section II-A, including total fairness [25], path-specific fairness (elimination of indirect discrimination) [4], and counterfactual fairness [26].

### B. Reweighting For Causal Fairness

We propose a reweighting scheme which consists of neural networks ($F^1$, $F^2$) and one discriminator ($D$). Fig. 1 shows the framework of our method.



(a) Causal Graphs $\mathcal{G}$ and $\mathcal{G}_s$



(b) Neural Networks $F^1$ and $F^2$

Fig. 2. The Neural Networks $F^1$ and $F^2$ on total effect. $S$ is 1 or 0 for the interventional joint distributions $P_{F^2}(s^+)$ (red path) and $P_{F^2}(s^-)$ (green path), respectively. The pair of nodes connected by dashed lines indicate that they share the same function (structures and parameters of the corresponding sub-neural networks).

As shown in Section III-A, causal fairness notions measures the difference between the interventional distributions. To guarantee these notions, our method adopts two neural networks to approximate the causal relations. One neural network $F^1$ simulates the causal model $\mathcal{M}$, while the other neural network $F^2$ approximates the interventional model $\mathcal{M}_s$ according to which kind of causal effect is measured. $F^1$ aims to force the reweighted data close to the given causal graph, and $F^2$ aims to drive the interventional distributions to satisfy the specific notion defined in Section III-A. To represent the connections between the two causal models, the two neural networks share certain structures and parameters, while they differ in sub-neural networks to indicate the intervention (the edges to $S$ in the interventional graph is deleted). Then, our method adopts a discriminator $D$ trying to distinguish the two interventional distributions (reweighted) $P(Y_{s^+})$ and $P(Y_{s^-})$. Finally, the discriminator and reweighting play an adversarial game to produce weights for individuals in the dataset.

To better illustrate our design, we divide $X$ into $\{A, B\}$ and $V = \{S, A, B, Y\}$ based on the positions of the nodes in the causal graph – variables in $A$ are direct causes of $S$ and variables in $B$ are descendants of $S$ and $A$.

*1) Reweighting for Total Fairness:* The causal graph $\mathcal{G}$ is shown in Fig. 2(a). We also show the interventional graph $\mathcal{G}_s$ with the intervention $do(S = s)$ and the edge from $A$ to $S$ is deleted in $\mathcal{G}_s$, which is also altered in $F^2$. The pair of nodes connected by dashed lines indicate that they share the same function (structures and parameters of the corresponding sub-neural networks) as shown in Fig. 2(b). For parallel nodes in the two graphs, the corresponding sub-neural networks are

synchronized during the training process.

We first show our method to achieve total fairness by describing each components of our design. As mentioned in Section II-A , $|TE(s^+, s^-)| < \tau$ must hold for all possible paths from $S$ to $Y$ shown in Fig. 2(a).

*a) Neural Networks $F^1$ and $F^2$:* The feed-forward Neural Network $F^1$ is constructed to correspond with the causal graph $\mathcal{G}$. It consists of $|V| - r$ sub-neural networks ($r$ is the total number of the root nodes in $\mathcal{G}$), with each corresponding to a node in $V$ (expect for the root nodes). Similar to what is described as the design of CFGAN in Section II-C, each sub-neural network $F^1_{V_i}$ is trying to approximate the corresponding function $f_{V_i}(Pa_{V_i})$ in the causal model $\mathcal{M}$ of the given causal graph $\mathcal{G}$. When $F^1$ is properly trained, the causal model $\mathcal{M}$ is learned. Then, $F^1_{V_i}$ outputs the estimated values of $V_i$, i.e., $\hat{v}_i$. The other neural network $F^2$ is constructed to align with the interventional graph $\mathcal{G}_s$, where all the incoming edges to S are removed under the intervention $do(S = s)$. The layout of $F^2$ is analogous to $F^1$, but with the exception that the sub-neural network $F^2_S$ is designated as $F^2_S \equiv 1$ if $s = s^+$, and $F^2_S \equiv 0$ if $s = s^-$. To synchronize the two neural networks $F_1$ and $F_2$, they share the identical set of structures and parameters for every corresponding pair of sub-neural networks, i.e., $F^1_{V_i}$ and $F^2_{V_i}$ for each $V_i$ except for $S$. When $F_2$ is properly trained, the interventional model $\mathcal{M}_s$ is learned. With $\mathcal{M}$ and $\mathcal{M}_s$ learned, we could manipulate the interventional distributions to reach our goal of causal fairness.

*b) Discriminator:* $D$ is used to differentiate between the two interventional distributions $\hat{y}_{s^+} \sim P_{F^2}(Y_{s^+})$ and $\hat{y}_{s^-} \sim P_{F^2}(Y_{s^-})$. The aim of the discriminator $D$ to minimize the bias by penalizing differences between both groups.

*c) Weights:* Assuming the to-reach-causal-fairness-importance of each individual in the given dataset is known, we can assign importance to different individuals in $\mathcal{M}_s$ to improve causal fairness for any downstream task. $\mathbf{w} = (w_1, ..., w_m)$ is a sample reweighting vector with length $m$, where $w_k$ indicates the importance of the $k$-th observed sample $(s_k, x_k, y_k)$. We want to reach a balance of goodness-of-fit to the known causal graph $\mathcal{G}$ which is learned from $\mathcal{D}$ and reweighting for causal fairness.

*Recall that here we assume that the known causal graph $\mathcal{G}$ is learned from a causal discovery which means it achieves goodness-of-fit.* We do not want the reweighted data to drift too far from the original causal graph. We use hatted variables to represent the output of the neural networks of the graphs. To reach this objective, we have:

$$S_{F^1}(\mathcal{G}) = \min_{F^1} \sum_{i=1}^{m} w_i l((s_i, x_i, y_i), (s_i, \hat{x}_i, \hat{y}_i)) \qquad (2)$$

where $l((s_i, x_i, y_i), (s_i, \hat{x}_i, \hat{y}_i))$ represents the loss of fitting observation $(x_i, y_i, s_i)$. In the experiment, we use weighted MSE loss for the continuous variables and weighted cross entropy loss for the categorical variables. The problem then becomes how to learn appropriate the sample reweighting

vector $\mathbf{w}$ for the objective of causal fairness. We formulate our objective as a minmax problem to reweight with $\mathcal{M}_s$:

$$\min_{\mathbf{w}} \max_{D} \sum_{k=1}^{m} w_k (D(\hat{y}_k^{s^+}) - D(\hat{y}_k^{s^-})), \qquad (3)$$

To avoid information loss by assigning close to zero weights to some samples from the group of $S^+$, we introduce a regularization constraint to the minimization term:

$$\sum_{k=1}^{m} (w_k - 1)^2 \leqslant Tm \qquad (4)$$

Thus, by adjusting the value of $T$, we can balance between similarity and dissimilarity of the weights of samples.

Samples easily fitted with fairness constraint should contribute more to $\mathcal{G}_s$: these are the samples with less difference of discriminator outputs from $do(S = s^-)$ to $do(S = s^+)$. We therefore use downweighting on the not-hold-fairness samples, and upweighting on the hold-fairness samples. This could be achieved by assigning weights to samples based on the discriminator $D$ outputs. When the neural networks are properly trained, the discriminator should not be able to tell if the sample is from the group of $S^+$ or $S^-$ which could achieve total fairness as we describe in Section III-A.

*2) Reweighting for Path-Specific Fairness:* The notions of direct and indirect discrimination are connected to effects specific to certain paths. We concentrate on indirect discrimination, even though fulfilling the criterion for direct discrimination is comparable. As mentioned in Section II-A , $|SE_\pi(s^+, s^-)| = |P(Y_{s^+|\pi_C, s^-|\overline{\pi}_C}) - P(Y_{s^-})| < \tau$ must hold for a path set $\pi_C$ that includes paths passing through certain attributes, shown in Fig. 7(a) ( in Appendix). $F^1$ for indirect discrimination is similar to that in Section III-B1. However, the design of $F^2$ is altered because it needs to adapt to the situation where the intervention is transferred only through $\pi_C$, shown in Fig. 7(b) (in Appendix). We examine two possible states for the sub-neural network $F^2_S$: the reference state and the interventional state. Under the reference state, $F^2_S$ is constantly set to 0. On the other hand, under the interventional state, $F^2_S$ is set to 1 if $s = s^+$, and 0 if $s = s^-$. For other sub-neural networks, there are also two possible values: the reference state and the interventional state, according to the state of $F^2_S$. If a sub-neural network corresponds to a node that is not present on any path in $\pi_C$, it only accepts reference states as input and generates reference states as output. However, for any other sub-neural network $F^2_{V_j}$ that exists on at least one path in $\pi_C$, it may accept both reference and interventional states as input and generate both types of states as output.

*3) Reweighting for Counterfactual Fairness:* In the context of counterfactual fairness, interventions are made based on a subset of variables $O = o$. Both $F^1$ and $F^2$ have similar structures to those in Section III-B2. However, we only use samples in $F^2$ as interventional samples if they satisfy the condition $O = o$. This means that the interventional distribution from $F^2$ is conditioned on $O = o$ as $P_{F^2}(X_s, Y_s|o)$. The

discriminator $D$ is designed to distinguish between $\hat{y}s^+|o \sim P_{F^2}(Y_{s+}|o)$ and $\hat{y}s^-|o \sim P_{F^2}(Y_{s-}|o)$, and aims to reach $P_{F^2}(Y_{s+}|o) = P_{F^2}(Y_{s-}|o)$. During training, the value of $m$ should be adjusted based on the number of samples that are involved in the intervention.

### C. Training Algorithm

To train the network $F^1$ to minimize the loss in Equation 2, we alternately optimize the network parameters of $F^1$ and $D$ and learn the weights $\mathbf{w}$ by fixing others as known.

*a) Updating parameters of $F^1$ with fixed $\mathbf{w}$:* Fixing $\mathbf{w}$, we update $F^1$ to minimize the loss in Equation 2 for $M$ steps, using the mini-batch stochastic gradient descent algorithm.

*b) Updating $\mathbf{w}$ with fixed $F^2$ (synchronized with $F^1$):* Fixing parameters of $F^2$, we control the training data into two groups ($S^+$ and $S^-$) for intervention, and learn $\mathbf{w}$ in Equation 3. Since Equation 3 is a min-max optimization problem, we can alternately optimize the weights $\mathbf{w}$ and the parameters of $D$ of the discriminator by fixing the other one as known. Therefore, we first fix $w_i = 1$ for all $i$ and optimize $D$ to maximize the objective function in Equation 3 using the gradient penalty technique, as in WGAN with Gradient Penalty [27]. Note that when $w_i = 1$ for all $i$, Equation 2 is equivalent to the situation when there is no reweighting applied. Then, fixing the discriminator $D$, we optimize $\mathbf{w}$.

We denote $d_k = D(\hat{y}_k^{S^+}) - D(\hat{y}_k^{S^-})$ and $\mathbf{d} = (d_1, d_2, ...d_m)^T$. The optimization problem for $\mathbf{w}$ becomes a constrained least squares problem:

$$\min_{\mathbf{w}} \mathbf{d}^T\mathbf{w}, s.t. w_k \geqslant 0, \sum_{k=1}^{m} w_k = m, \sum_{k=1}^{m}(w_k - 1)^2 \leqslant Tm \quad (5)$$

## IV. EXPERIMENTAL EVALUATION

We conduct experiments on two benchmarks datasets (ADULT [28] and COMPAS [29]) to evaluate our reweighting approach and compare it with state-of-the-art methods: FairGAN [30], CFGAN [12] and Causal Inference for Social Discrimination Reasoning (CISD) [31] for total effect and indirect discrimination (please refer to Appendix (A) for more details about the datasets). CISD [31] introduces a technique for identifying causal discrimination through the use of propensity score analysis. It consists of mitigating the influence of confounding variables by reweighing samples based on the propensity scores calculated from a logistic regression. The approach, however, is purely statistical with no causal knowledge exploited. We also compare our method with CFGAN and two methods from [26] (we refer them as $CE_1$ and $CE_3$ in our paper) for counterfactual effect. $CE_1$ only uses on non-descendants of $S$ for classification. $CE_3$ is similar to $CE_1$ but presupposes an additive $U$. The reason we choose these methods is: FairGAN for statistical parity and CFGAN for causal fairness also use adversarial method to mitigate bias, similar to our design; CISD approaches causal fairness with weighting scheme. We then compare the performance of our method with the mentioned methods on total effect, indirect discrimination and counterfactual fairness with

4 different downstream classifiers: decision tree (DT) [32], logistic regression (LR) [33], support vector machine (SVM) [34] and random forest (RF) [35]. We compare the accuracy of the downstream tasks to see if the data preserves good utility, where higher accuracy indicates better utility. For the utility of the downstream task, we also compute the Wasserstein distance between the manipulated data and the original data, where a smaller Wasserstein distance indicates closer the two distributions, and better utility for the downstream tasks.

### A. The datasets and setup

Due to page limit, please refer to Appendix for the details of datasets and training.

### B. Analysis

*1) Total Effect:* In Table I, we present the total effect (TE) calculated for the original dataset and the datasets processed using various methods. The original ADULT dataset has a total effect of 0.1854 and COMPAS 0.2389, while applying FairGAN to achieve demographic parity yields almost no total effect. As mentioned in Section II-A, total effect is very similar to demographic parity. However, FairGAN is limited by its focus on statistical fairness, rather than causal fairness, and does not perform well on Wasserstein distance or downstream tasks. It is quite intuitive that if total fairness is met, total fairness should be achieved too on the condition that the causal graph is sufficient. We test it on our two datasets and the result is acceptable. CFGAN produces no total effect, but it performs worse than our method on Wasserstein distance, possibly because reweighted data could manage to stay closer to the original data distribution. Our method also outperforms CISD, which may be due to the use of a neural network instead of logistic regression to calculate weights, allowing for greater flexibility in capturing the dataset.

*A Closer Look at the Weights* After ranking the weights of samples in the Adult dataset, we observed that older individuals from Europe or Asia (e.g., Germany and India) tend to have the highest weights, while younger black individuals from Caribbean countries (e.g., Jamaica and Haiti) tend to have lower weights. This suggests that when sex is intervened from female to male, the former group is less influenced by the change, while the latter group is more influenced in terms of income. White, middle-aged individuals born in the US are assigned medium weights. To visualize it, we build a decision tree to classify top 10% individuals with highest weights and bottom 10% individuals with lowest weights using the three root nodes $\{race, native\_country, age\}$, shown in Fig. 3.

*2) Indirect Discrimination:* To address indirect discrimination (SE), we identify all possible paths except the direct one $\{S \rightarrow Y\}$ as the path $\pi_C$ and evaluated the results in Table I. Similar to total effect, FairGAN removes indirect discrimination but at the cost of significant utility loss. In contrast, both CISD and our method can effectively remove indirect discrimination while maintaining better data utility than FairGAN. Although CFGAN and CISD perform similarly using different techniques, our method outperforms both

TABLE I
THE TOTAL EFFECT (TE) AND INDIRECT DISCRIMINATION (SE) ON ADULT AND COMPAS DATASETS

| ADULT | total effect | indirect discrimination | Wasserstein distance | SVM | DT | LR | RF | COMPAS | total effect | indirect discrimination | Wasserstein distance | SVM | DT | LR | RF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| original data | 0.1854 (0.0301) | 0.1773 (0.0489) | 0 | 81.78 (1.45) | 81.77 (1.75) | 81.70 (1.63) | 81.78 (1.76) | original data | 0.2389 (0.0245) | 0.2137 (0.0985) | 0 | 65.24 (2.34) | 65.15 (1.46) | 65.10 (2.19) | 65.27 (1.09) |
| Ours (TE) | **0.0017 (0.0009)** | **0.0012 (0.0007)** | **0.71 (0.19)** | 81.12 (1.72) | 81.20 (1.86) | 81.60 (2.03) | 81.14 (1.05) | Ours (TE) | **0.0037 (0.0018)** | 0.0017 (0.0009) | **1.21 (0.32)** | 65.09 (2.75) | 65.13 (1.76) | 65.06 (2.08) | 65.11 (1.02) |
| Ours (SE) | | | **0.69 (0.23)** | 81.14 (1.58) | 80.97 (2.01) | 81.65 (1.96) | 81.17 (1.92) | Ours (SE) | | | **0.72 (0.35)** | 65.11 (1.98) | 65.14 (2.06) | 65.02 (1.12) | 65.09 (1.95) |
| FairGAN | 0.0021 (0.0007) | 0.0148 (0.0075) | 5.21 (0.78) | 79.88 (1.47) | 79.81 (1.89) | 80.36 (1.32) | 80.82 (1.65) | FairGAN | 0.0075 (0.0056) | 0.0341 (0.0075) | 3.24 (1.45) | 64.24 (1.77) | 64.15 (2.01) | 64.50 (2.75) | 64.26 (2.34) |
| CFGAN (TE) | 0.0106 (0.0008) | 0.0034 (0.0012) | 1.78 (0.65) | 80.34 (2.56) | 80.15 (1.52) | 80.07 (1.65) | 80.39 (1.32) | CFGAN (TE) | 0.0364 (0.0175) | **0.0016 (0.0025)** | 2.76 (1.65) | 64.59 (2.65) | 65.13 (2.73) | 65.02 (2.03) | 65.01 (2.45) |
| CFGAN (SE) | | | 1.89 (0.29) | 80.37 (1.56) | 80.49 (2.05) | 80.04 (1.67) | 80.24 (1.09) | CFGAN (SE) | | | 2.64 (0.91) | 64.21 (2.45) | 64.25 (1.75) | 64.80 (1.97) | 64.87 (1.54) |
| CISD (TE) | 0.0206 (0.0074) | 0.0098 (0.0045) | 2.57 (0.18) | 80.73 (1.42) | 80.74 (1.75) | 81.15 (1.82) | 81.27 (1.47) | CISD (TE) | 0.0356 (0.0246) | 0.0175 (0.0231) | 2.57 (1.61) | 65.04 (1.76) | 65.17 (1.54) | 64.80 (2.47) | 65.05 (1.75) |
| CISD (SE) | | | 2.82 (0.23) | 80.75 (1.28) | 80.72 (1.58) | 80.77 (1.96) | 81.32 (1.95) | CISD (SE) | | | 2.65 (1.56) | 64.01 (1.56) | 65.02 (1.49) | 64.09 (2.45) | 64.11 (1.32) |

TABLE II
THE COUNTERFACTUAL EFFECT (CE) ON ADULT AND COMPAS DATASETS

| ADULT | counterfactual effect $o_1$ | counterfactual effect $o_2$ | Wasserstein distance | SVM | DT | LR | RF | COMPAS | counterfactual effect $o_1$ | counterfactual effect $o_2$ | Wasserstein distance | SVM | DT | LR | RF |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| original data | 0.1254 (0.0107) | 0.1265 (0.0489) | 0 | 81.78 | 81.77 | 81.70 | 81.78 | original data | 0.2079 (0.0303) | 0.1973 (0.0484) | 0 | 65.24 (2.35) | 65.15 (1.32) | 65.10 (1.97) | 65.27 (1.35) |
| Ours | **0.0017 (0.0009)** | 0.0030 (0.0007) | **0.98 (0.10)** | 81.14 (1.32) | 81.25 (1.37) | 81.67 (1.36) | 81.29 (1.42) | Ours | **0.0027 (0.0009)** | 0.0037 (0.0109) | **1.43 (1.68)** | 65.11 (2.01) | 65.14 (1.35) | 65.07 (1.89) | 65.13 (1.68) |
| $CE_1$ | 0.0021 (0.0157) | 0.0148 (0.0025) | 4.95 (0.45) | 77.88 (1.65) | 78.81 (1.32) | 76.36 (2.06) | 77.82 (1.78) | $CE_1$ | 0.0034 (0.0007) | 0.0048 (0.0075) | 4.53 (1.56) | 63.24 (1.76) | 63.16 (1.43) | 62.19 (1.46) | 62.27 (2.21) |
| $CE_3$ | 0.1123 (0.0017) | 0.1046 (0.0057) | 3.11 (0.65) | 80.75 (1.76) | 81.59 (2.15) | 81.62 (2.01) | 81.43 (1.37) | $CE_3$ | 0.1021 (0.0981) | 0.0145 (0.0025) | 1.76 (0.89) | 64.07 (1.96) | 65.03 (1.95) | 64.01 (1.45) | 65.10 (2.89) |
| CFGAN (CE) | 0.0021 (0.0039) | **0.0027 (0.0064)** | 2.99 (0.78) | 81.03 (1.35) | 81.14 (2.19) | 81.11 (2.17) | 81.15 (1.95) | CFGAN(CE) | 0.0034 (0.0024) | **0.0031 (0.0045)** | 1.74 (1.06) | 65.04 (1.72) | 65.11 (1.74) | 64.08 (1.19) | 64.29 (1.32) |

```
|--- age <= 37.50
|   |--- class: bottom
|--- age > 37.50
|   |--- race <= 2.50
|   |   |--- race <= 1.50
|   |   |   |--- class: top
|   |   |--- race > 1.50
|   |   |   |--- age <= 57.50
|   |   |   |   |--- class: bottom
|   |   |   |--- age > 57.50
|   |   |   |   |--- class: top
|   |--- race > 2.50
|   |   |--- class: top
```

Fig. 3. The visualize of the decision tree trying to classify individuals with low or high weights. we see that age and race are the most important attributes to build the tree. The mapping of label encoder for $race$ is $\{'Amer-Indian-Eskimo':0,'Asian-Pac-Islander':1,'Black':2,'Other':3,'White':4\}$

methods in terms of Wasserstein distance, indicating the best overall utility among these approaches.

*3) Counterfactual Fairness:* To evaluate counterfactual effect (CE), we consider the conditions on two variables - race and native country (binarized) for ADULT, and sex and age (binarized) for COMPAS - resulting in four value combinations. Table II presents the results for two selections (see Appendix (B6) for more details). We find biases in the original data regarding counterfactual fairness in these two selections. $CE_1$ is counterfactually fair, but the classifier accuracy is poor because it solely employs non-descendants of the sensitive attributes for outcome attributes. $CE_3$ cannot achieve counterfactual fairness, probably due to the strong assumptions

while introducing $U$. In contrast, our method performs well on both dimensions due to its flexibility. Although CFGAN performs well in some aspects, our method outperforms it in Wasserstein distance, likely because reweighting better preserves the original distribution than generation methods.

**Summary** We find out that in general neural nets-based methods outperform due to the flexibility of neural networks to capture any function, while reweighting outperforms generation. We could see from the experiment results above, imposing strong assumptions on the $U$ and $F$ could cause unwanted problems, and we argue that is why neural nets should be explored more in causal fairness problem settings. Fairness related methods usually formalize the problem as an optimization trade-off between utility and specific fairness objectives. Nevertheless, these discussions are often based on a fixed distribution that does not align with our current situation. We think that an ideal distribution might exist where fairness and utility are in harmony. To include the reweighting scheme into the downstream tasks could be an very interesting future direction to locate this harmonious distribution.

## V. CONCLUSION, LIMITATION AND FUTURE WORK

We propose a novel approach for achieving causal fairness by dataset reweighting. Our method considers different causal fairness objectives, such as total fairness, path-specific fairness and counterfactual fairness. It consists of two feed-forward neural networks $F^1$ and $F^2$ and a discriminator $D$. The structures of $F^1$ and $F^2$ are designed based on the original causal graph $\mathcal{G}$ and interventional graph $\mathcal{G}_s$, and the discriminator $D$ is used to ensure causal fairness combined with a reweighting scheme. Our experiments on two datasets show

that the approach improves over state-of-the-art approaches for the considered causal fairness notions achieving minimal loss of utility. Moreover, by analyzing the sample weights assigned by the approach, the user can gain an understanding of the distribution of the biases in the original dataset. Future work involve analyzing the sample weights further, e.g., by using methods from the eXplainable in AI research area. As another relevant research direction, since practitioners often lack sufficient causal graphs when working with a dataset [36], an extension of our work could involve causal discovery as an integral part of the approach.

## VI. Acknowledgement

## References

[1] D. Pedreschi, S. Ruggieri, and F. Turini, "Discrimination-aware data mining," in *KDD*. ACM, 2008, pp. 560–568.

[2] I. Zliobaite, F. Kamiran, and T. Calders, "Handling conditional discrimination," in *ICDM*. IEEE Computer Society, 2011, pp. 992–1001.

[3] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in *NIPS*, 2016, pp. 3315–3323.

[4] L. Zhang, Y. Wu, and X. Wu, "A causal framework for discovering and removing direct and indirect discrimination," in *IJCAI*. ijcai.org, 2017, pp. 3929–3935.

[5] ——, "Achieving non-discrimination in prediction," in *IJCAI*. ijcai.org, 2018, pp. 3097–3103.

[6] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, "Certifying and removing disparate impact," in *KDD*. ACM, 2015, pp. 259–268.

[7] L. Zhang, Y. Wu, and X. Wu, "Causal modeling-based discrimination discovery and removal: Criteria, bounds, and algorithms," *IEEE Trans. Knowl. Data Eng.*, vol. 31, no. 11, pp. 2035–2050, 2019.

[8] H. Edwards and A. J. Storkey, "Censoring representations with an adversary," in *ICLR (Poster)*, 2016.

[9] Q. Xie, Z. Dai, Y. Du, E. H. Hovy, and G. Neubig, "Controllable invariance through adversarial feature learning," in *NIPS*, 2017, pp. 585–596.

[10] D. Madras, E. Creager, T. Pitassi, and R. S. Zemel, "Learning adversarially fair and transferable representations," in *ICML*, ser. Proceedings of Machine Learning Research, vol. 80. PMLR, 2018, pp. 3381–3390.

[11] B. H. Zhang, B. Lemoine, and M. Mitchell, "Mitigating unwanted biases with adversarial learning," in *AIES*. ACM, 2018, pp. 335–340.

[12] D. Xu, Y. Wu, S. Yuan, L. Zhang, and X. Wu, "Achieving causal fairness through generative adversarial networks," in *IJCAI*. ijcai.org, 2019, pp. 1452–1458.

[13] Y. Roh, K. Lee, S. Whang, and C. Suh, "Sample selection for fair and robust training," in *NeurIPS*, 2021, pp. 815–827.

[14] F. P. Calmon, D. Wei, B. Vinzamuri, K. N. Ramamurthy, and K. R. Varshney, "Optimized pre-processing for discrimination prevention," in *NIPS*, 2017, pp. 3992–4001.

[15] S. Aghaei, M. J. Azizi, and P. Vayanos, "Learning optimal and fair decision trees for non-discriminative decision-making," in *AAAI*. AAAI Press, 2019, pp. 1418–1426.

[16] R. Berk, H. Heidari, S. Jabbari, M. Joseph, M. J. Kearns, J. Morgenstern, S. Neel, and A. Roth, "A convex framework for fair regression," *CoRR*, vol. abs/1706.02409, 2017.

[17] J. Pearl, *Causality: Models, Reasoning and Inference*, 2nd ed. Cambridge University Press, 2009.

[18] P. Spirtes and K. Zhang, "Causal discovery and inference: Concepts and recent methodological advances," *Applied Informatics*, vol. 3, no. 1, p. 3, 2016.

[19] C. Glymour, K. Zhang, and P. Spirtes, "Review of Causal Discovery Methods Based on Graphical Models," *Frontiers in Genetics*, vol. 10, 2019.

[20] P. Spirtes, C. Meek, and T. S. Richardson, "Causal inference in the presence of latent variables and selection bias," *CoRR*, vol. abs/1302.4983, 2013.

[21] P. Spirtes and C. Glymour, "An Algorithm for Fast Recovery of Sparse Causal Graphs," *Social Science Computer Review*, vol. 9, no. 1, pp. 62–72, 1991.

[22] D. Colombo, M. H. Maathuis, M. Kalisch, and T. S. Richardson, "Learning high-dimensional directed acyclic graphs with latent and selection variables," *The Annals of Statistics*, vol. 40, no. 1, pp. 294–321, 2012.

[23] M. J. Vowels, N. C. Camgöz, and R. Bowden, "D'ya like dags? A survey on structure learning and causal discovery," *ACM Comput. Surv.*, vol. 55, no. 4, pp. 82:1–82:36, 2023.

[24] M. Kocaoglu, C. Snyder, A. G. Dimakis, and S. Vishwanath, "CausalGAN: Learning causal implicit generative models with adversarial training," in *ICLR (Poster)*. OpenReview.net, 2018.

[25] J. Zhang and E. Bareinboim, "Fairness in decision-making - the causal explanation formula," in *AAAI*. AAAI Press, 2018, pp. 2037–2045.

[26] M. J. Kusner, J. R. Loftus, C. Russell, and R. Silva, "Counterfactual fairness," in *NIPS*, 2017, pp. 4066–4076.

[27] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville, "Improved training of wasserstein GANs," in *NIPS*, 2017, pp. 5767–5777.

[28] R. Kohavi, "Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid," in *KDD*. AAAI Press, 1996, pp. 202–207.

[29] Mattu, J. Angwin, L. Kirchner, Surya, and J. Larson, "How We Analyzed the COMPAS Recidivism Algorithm," ProPublica, 2016, https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm.

[30] D. Xu, S. Yuan, L. Zhang, and X. Wu, "Fairgan: Fairness-aware generative adversarial networks," in *IEEE BigData*. IEEE, 2018, pp. 570–575.

[31] B. Qureshi, F. Kamiran, A. Karim, S. Ruggieri, and D. Pedreschi, "Causal inference for social discrimination reasoning," *J. Intell. Inf. Syst.*, vol. 54, no. 2, pp. 425–437, 2020.

[32] X. Wu, V. Kumar, J. R. Quinlan, J. Ghosh, Q. Yang, H. Motoda, G. J. McLachlan, A. Ng, B. Liu, S. Y. Philip *et al.*, "Top 10 algorithms in data mining," *Knowledge and information systems*, vol. 14, no. 1, pp. 1–37, 2008.

[33] D. R. Cox, "The regression analysis of binary sequences," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 20, no. 2, pp. 215–232, 1958.

[34] C. Cortes and V. Vapnik, "Support-vector networks," *Machine learning*, vol. 20, no. 3, pp. 273–297, 1995.

[35] T. K. Ho, "The random subspace method for constructing decision forests," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 20, no. 8, pp. 832–844, 1998.

[36] R. Binkyte-Sadauskiene, K. Makhlouf, C. Pinzón, S. Zhioua, and C. Palamidessi, "Causal discovery for fairness," *CoRR*, vol. abs/2206.06685, 2022.

[37] D. Chae, J. Kang, S. Kim, and J. Lee, "CFGAN: A generic collaborative filtering framework based on generative adversarial networks," in *CIKM*. ACM, 2018, pp. 137–146.

[38] D. Plecko, N. Bennett, and N. Meinshausen, "fairadapt: Causal reasoning for fair data pre-processing," *CoRR*, vol. abs/2110.10200, 2021.

[39] A. Paszke *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *NeurIPS*, 2019, pp. 8024–8035.

[40] S. Diamond and S. P. Boyd, "CVXPY: A python-embedded modeling language for convex optimization," *J. Mach. Learn. Res.*, vol. 17, pp. 83:1–83:5, 2016.

## Appendix

### A. Dataset and Training Details

The causal graph [37] for ADULT is shown in Fig. 4, and for COMPAS [38] in Fig. 5. Note that the causal graphs here are sourced from existing literature.

*1) Adult Dataset:* The Adult dataset was drawn from the 1994 United States Census Bureau data. It contains 65,123 samples with 11 variables. It used personal information such as education level and working hours per week to predict whether
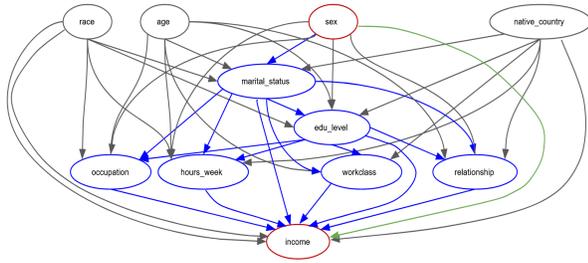
Fig. 4. The causal graph of the Adult dataset depicts the indirect path set with blue paths, while the direct path is represented by the green path.
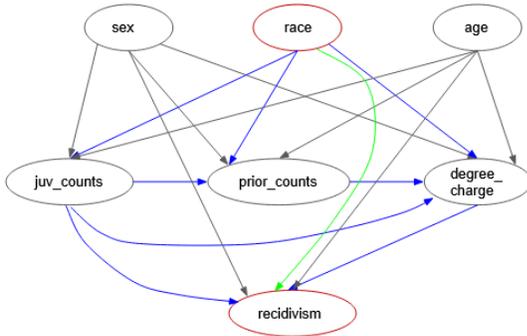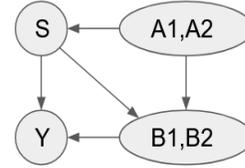


(a) Causal graph $\mathcal{G}$



(b) Neural Network $F^1$

Fig. 6. details of the connection of the neural nets of a given $\mathcal{G}$. In Fig. 6(b), each nodes are either input or output of a sub-neural nets or both. Note that we do not show the inner layers here for simplicity.



Fig. 5. The causal graph of the COMPAS dataset depicts the indirect path set with blue paths, while the direct path is represented by the green path

an individual earns more or less than \$50,000 per year. The dataset is imbalanced – the instances made less than \$50,000 constitute 25% of the dataset, and the instances made more than \$50,000 constitute 75% of the dataset. As for gender, it is also imbalanced. We use age, years of education, capital gain, capital loss, hours-per-week, etc., as continuous features, and education level, gender, etc., as categorical features. We set the batch size at 640 and train 30 epochs for convergence. We set the learning rate $\eta$ at 0.001 according to the experiment result.

*2) COMPAS Dataset:* COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) is a popular commercial algorithm used by judges and parole officers for scoring criminal defendant's likelihood of reoffending (recidivism). The COMPAS dataset includes the processed COMPAS data between 2013-2014. The data cleaning process followed the guidance in the original COMPAS repo. It Contains 6172 observations and 14 features. In our causal graph, we use 7 features. Due to the limited size of COMPAS dataset, it does not perform so well on NN based tasks.
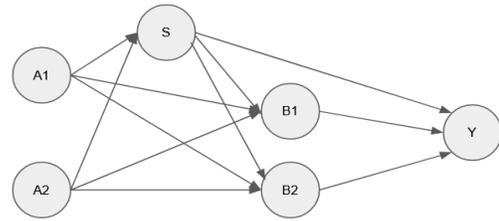
### B. Training Details

For ADULT and COMPAS datasets, some pre-processing is performed. We normalize the continuous features and use one-hot encoding to deal with the categorical features for the

input of $F^1$ and $F^2$. We use sex and race as the sensitive variable $S$ in ADULT and COMPAS respectively, income and two year recidivism as the outcome variable $Y$.

For $F^1$ and $F^2$, we apply fully connected layers. For the discriminator $D$, we use the same architecture proposed in [27]. We apply SGD algorithm with a momentum of 0.9 to update $F^1$ and $F^2$. $D$ is updated by the Adam algorithm with a learning rate 0.0001. Following [27], we adjust the learning rate $\eta$ by $\eta = \frac{0.01}{(1+10p)^{-0.75}}$, where $p$ is the training progress linearly changing from 0 to 1. We update $F^1$ and $F^2$ for 2 steps then update $D$ for 1 step. For more details of the experiment (e.g., the split of training and testing datasets, the details of architectures of the neural nets, the estimation of Wasserstein distance), please refer to the Appendix (B). We then evaluate the performance of our method of reweighting to achieve different types of causal fairness and utility.

Our test are run on an Intel(r) Core(TM) i7-8700 CPU. The networks in the experiments are built based on Pytorch [39], the optimization in Equation (5) is performed with the Python package CVXPY [40].

*1) Details of architectures of the feed-forward Neural Networks $F^1$ and $F^2$ with sub-neural networks:* To simplify our demonstration, we consider a causal graph $\mathcal{G}$ with 6 attributes $\{S, A_1, A_2, B_1, B_2, Y\}$ as shown in Fig. 6(a). And Fig. 6(b) shows the joint neural network of it.

*2) Details of WGAN-GP adaptation for our method:* In our design, we adopt the discriminator from WGAN-GP: in
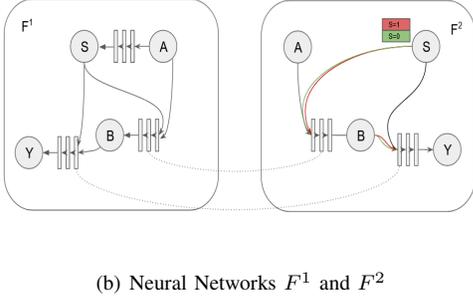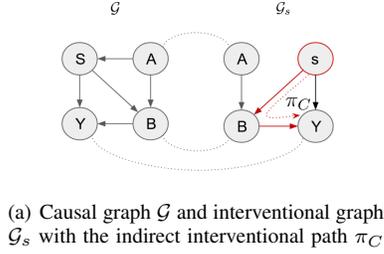
(a) Causal graph $\mathcal{G}$ and interventional graph $\mathcal{G}_s$ with the indirect interventional path $\pi_C$



(b) Neural Networks $F^1$ and $F^2$

Fig. 7. the Neural Networks $F^1$ and $F^2$ based on indirect discrimination. $S$ is 1 or 0 and the intervention is only along $\pi_C = \{S \to B \to Y\}$ for the interventional distributions $P_{F^2}(s^+)$ (red) and $P_{F^2}(s^-)$ (green) respectively. Compared with Fig. 2, we could see that the intervention is not transferred directly from $S$ to $Y$ ($\{S \to Y\}$) in Fig. 7.
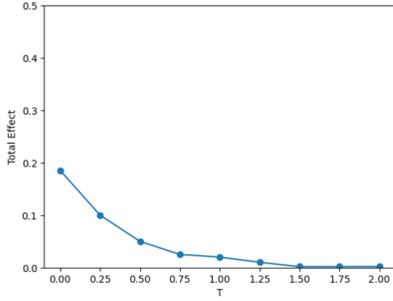


Fig. 8. Sensitivity of total effect on the change of $T$ on ADULT dataset.
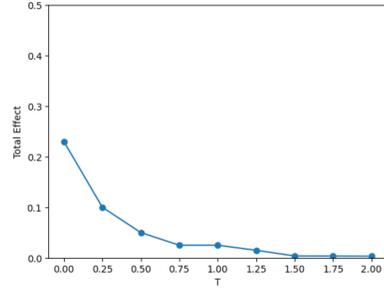


Fig. 9. Sensitivity of total effect on the change of $T$ on COMPAS dataset.

We could see from the Fig. 9 that the trend of $T$ on total effect on COMPAS dataset is similar to what is shown earlier on ADULT dataset.

*4) Details of calculating different causal effects:* We have a discriminator in our design, to calculate different causal effect, we just send the samples from different groups ($S^+$ and $S^-$), normalize the output of $D$ within groups, then get the difference as causal effect.

*5) Repetition:* We repeat experiments on each dataset five times. Before each repetition, we randomly split data into training data (80%) and test data (20%) for the computation of the standard errors of the metrics. We train 30 epochs for convergence.

*6) Choices of Counterfactual Effect:* Due to space limit, we only show two combinations of the counterfactual effect on individual features. For ADULT, we use $o_1 = \{white, us\}$, $o_2 = \{non\_white, us\}$. For COMPAS, we use $o_1 = \{male, under\_25\}$, $o_2 = \{male, above\_25\}$

*7) Details to approximate Wasserstein distance:* To approximate the Wasserstein distance, we WGAN-GP discriminator between the orginal data and data for evaluation. When the NNs are trained, we use the discriminator to approximate the Wasserstein distance between the two datasets.

the original work, the discriminator is used to differentiate between the generated and real data while we are trying to differentiate between $S^+$ and $S^-$. The difference between orginal GAN and WGAN-GP is that WGAN-GP introduces a gradient penalty term in the training objective to guarantee Wasserstein distance. Wasserstein distance itself has been used a lot in fairness realted topic to help detect or mitigate bias. Note that we choose relatively larger batch size since to approximate Wasserstein distance between two distributions requires relatively larger batch size.

*3) Sensitivity to the Choice of Hyper-Parameters:* We conduct an analysis of the sensitivity of our method to the hyper-parameters discussed in Section III, and the results are shown in Fig. 8. The figures demonstrate that our adversarial reweighting scheme's performance has low sensitivity to hyper-parameter choice when $T$ is above 1. Therefore, we set $T$ at 1.5.