

Supplementary Materials: Adjusting Machine Learning Decisions for Equal Opportunity and Counterfactual Fairness

A Proof of Theorem 1

Before proving [Thm. 1](#) about the eco-fairness and the optimality of the eco decision-maker, we first prove a lemma about the eco-fairness, which will be useful for proving [Thm. 1](#). It says that a decision-maker \hat{Y} is eco-fair if and only if there is no causal arrow from A to \hat{Y} .

Lemma 3. (*eco* \Leftrightarrow No $A \rightarrow \hat{Y}$) Assume the causal graph in [Fig. 2](#). A decision-maker \hat{Y} satisfies equal opportunities over A if and only if there is no causal arrow between A and \hat{Y} .

The intuition behind [Lemma 3](#) is that all decisions ($\hat{Y}^{\text{ml}}, \hat{Y}^{\text{eco}}, \hat{Y}^{\text{cf}}$) in the causal model [Fig. 2](#) satisfy $\hat{Y}(a, s) \perp A, S$ ([Pearl, 2009](#)). Therefore, the eco criterion ([Def. 1](#)) reduces to $\hat{Y}(a, s)$ being constant in a . In other words, \hat{Y} does not use the protected attribute to make a decision. Below we prove [Lemma 3](#).

Proof. The goal is to show that eco ([Def. 1](#)) is equivalent to

$$\hat{Y}(a, s) \stackrel{d}{=} \hat{Y}(a', s) \stackrel{d}{=} \hat{Y}(s) \quad (3)$$

for any $s \in \mathcal{S}$ and $a, a' \in \mathcal{A}$. This equation is equivalent to no causal arrow between A and \hat{Y} .

Begin with the definition of eco-fairness,

$$\hat{Y}(a', s) | A = a, S = s \stackrel{d}{=} \hat{Y}(a, s) | A = a, S = s \quad (4)$$

$$\Leftrightarrow \hat{Y}(a', s) \stackrel{d}{=} \hat{Y}(a, s) \quad \forall s, a, a' \quad (5)$$

$$\Leftrightarrow \hat{Y}(a', s) \stackrel{d}{=} \hat{Y}(s). \quad (6)$$

[Eq. 5](#) is due to the observation that $\hat{Y}(a, s) \perp S, A$ in [Fig. 2](#) ([Pearl, 2009](#)). [Eq. 6](#) is true because

$$P(\hat{Y}(s)) = \int P(\hat{Y}(A, s))P(A) dA = P(\hat{Y}(a', s)) \quad (7)$$

for any a' , where the last equality of [Eq. 7](#) relies on [Eq. 5](#). [Eq. 6](#) is equivalent to no causal arrow from A to \hat{Y} . \square

Given [Lemma 3](#), we next prove [Thm. 1](#).

Proof. The first part of [Thm. 1](#) is that \hat{Y}^{eco} is eco-fair. It is an immediate consequence of [Lemma 3](#) because $\hat{Y}^{\text{eco}}(a_{\text{new}}, s_{\text{new}}) = \hat{Y}^{\text{eco}}(s_{\text{new}})$ as is defined in [Eq. 1](#).

The second part of [Thm. 1](#) establishes the optimality of \hat{Y}^{eco} :

$$\hat{Y}^{\text{ECO}} = \arg \min_{Y^{\text{ECO}} \in \mathcal{Y}^{\text{ECO}}} \mathbb{E}_{P(A)P(S)} \left[\text{KL}(P(\hat{Y}^{\text{ml}}(A, S)) || P(Y^{\text{eco}}(A, S))) \right].$$

Begin by rewriting the goal of the proof. We will show this goal is equivalent to the definition of the eco decision-maker,

$$\begin{aligned} & \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} \mathbb{E}_{P(A)P(S)} \left[\text{KL}(P(\hat{Y}^{\text{ml}}(A, S)) || P(Y^{\text{eco}}(A, S))) \right] \\ &= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} \mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{ml}}(A, S)) \left[\log P(\hat{Y}^{\text{ml}}(A, S)) - \log P(Y^{\text{eco}}(A, S)) \right] dY dA \right] \end{aligned} \quad (8)$$

$$= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} -\mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{ml}}(A, S)) \log P(Y^{\text{eco}}(A, S)) dY dA \right] \quad (9)$$

$$= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} -\mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{ml}}(A, S)) \log P(Y^{\text{eco}}(S)) dY dA \right] \quad (10)$$

$$= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} -\mathbb{E}_{P(S)} \left[\int \left[\int P(A) P(\hat{Y}^{\text{ml}}(A, S)) dA \right] \log P(Y^{\text{eco}}(S)) dY \right] \quad (11)$$

$$= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} -\mathbb{E}_{P(S)} \left[\mathbb{E} \left[\int \left[\int P(A) P(\hat{Y}^{\text{ml}}(A, s)) dA \right] \log P(Y^{\text{eco}}(s)) dY \mid S = s \right] \right] \quad (12)$$

$$= \arg \min_{Y^{\text{eco}} \in \mathcal{Y}^{\text{eco}}} \mathbb{E}_{P(S)} \left[\mathbb{E} \left[\text{KL} \left(\int P(A) P(\hat{Y}^{\text{ml}}(A, s)) dA \parallel P(Y^{\text{eco}}(s)) \mid S = s \right) \right] \right] \quad (13)$$

$$= \int P(\hat{Y}^{\text{ml}}(A, s)) P(A) dA \quad (14)$$

$$= \hat{Y}^{\text{eco}}(a, s) \quad (15)$$

[Eq. 8](#) is due to the definition of the kl divergence. [Eq. 9](#) is due to \hat{Y}^{ml} being a given random variable. [Eq. 10](#) is due to [Lemma 3](#). [Eq. 11](#) switches the integral subject to conditions of the dominated convergence theorem. [Eq. 12](#) is due to [Fig. 2](#) and the tower property. [Eq. 13](#) is due to the definition of the kl divergence and \hat{Y}^{ml} being given. [Eq. 14](#) is because setting $P(Y^{\text{eco}}(s)) = \int P(A) P(\hat{Y}^{\text{ml}}(A, s)) dA$ has $\text{KL}(\int P(A) P(\hat{Y}^{\text{ml}}(A, s)) dA \parallel P(Y^{\text{eco}}(s))) = 0$ for all s . The expectation is also zero: $\mathbb{E}_{P(S)} \left[\text{KL}(P(\int P(A) P(\hat{Y}^{\text{ml}}(A, s)) dA) \parallel P(Y^{\text{eco}}(S))) \right] = 0$.

This calculation implies $\hat{Y}^{\text{eco}} = \int P(\hat{Y}^{\text{ml}}(A, s)) P(A) dA$ minimizes the average kl distance between the ML decision and the eco decision. In other words, the eco decision-maker maximally recovers the ML decision. Put differently, the eco decision-maker minimally modifies the ML decision to achieve eco-fairness. \square

The optimality of ftu under the population $P(A, S)$. We next present a supplementary result of [Thm. 1](#) that establishes the optimality of \hat{Y}^{FTU} under the same KL divergence objective but the population $P(A, S)$:

$$\hat{Y}^{\text{FTU}} = \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} \mathbb{E}_{P(A, S)} \left[\text{KL}(P(\hat{Y}^{\text{ml}}(A, S)) \parallel P(Y^{\text{FTU}}(A, S))) \right].$$

Begin by rewriting the goal of the proof. We will show this goal is equivalent to the definition of the ftu decision-maker,

$$\begin{aligned} & \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} \mathbb{E}_{P(A, S)} \left[\text{KL}(P(\hat{Y}^{\text{ml}}(A, S)) \parallel P(Y^{\text{FTU}}(A, S))) \right] \\ &= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} \mathbb{E}_{P(S)} \left[\int P(A | S) \int P(\hat{Y}^{\text{ml}}(A, S)) \left[\log P(\hat{Y}^{\text{ml}}(A, S)) - \log P(Y^{\text{FTU}}(A, S)) \right] dY dA \right] \quad (16) \end{aligned}$$

$$= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} -\mathbb{E}_{P(S)} \left[\int P(A | S) \int P(\hat{Y}^{\text{ml}}(A, S)) \log P(Y^{\text{FTU}}(A, S)) dY dA \right] \quad (17)$$

$$= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} -\mathbb{E}_{P(S)} \left[\int P(A | S) \int P(\hat{Y}^{\text{ml}}(A, S)) \log P(Y^{\text{FTU}}(S)) dY dA \right] \quad (18)$$

$$= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} -\mathbb{E}_{P(S)} \left[\int \left[\int P(A | S) P(\hat{Y}^{\text{ml}}(A, S)) dA \right] \log P(Y^{\text{FTU}}(S)) dY \right] \quad (19)$$

$$= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} -\mathbb{E}_{P(S)} \left[\mathbb{E} \left[\int \left[\int P(A | S) P(\hat{Y}^{\text{ml}}(A, s)) dA \right] \log P(Y^{\text{FTU}}(s)) dY \mid S = s \right] \right] \quad (20)$$

$$= \arg \min_{Y^{\text{FTU}} \in \mathcal{Y}^{\text{FTU}}} \mathbb{E}_{P(S)} \left[\mathbb{E} \left[\text{KL} \left(\int P(A | S) P(\hat{Y}^{\text{ml}}(A, s)) dA \parallel P(Y^{\text{FTU}}(s)) \mid S = s \right) \right] \right] \quad (21)$$

$$= \int P(\hat{Y}^{\text{ml}}(A, s)) P(A | S) dA \quad (22)$$

$$=\hat{Y}^{\text{FTU}}(a, s) \quad (23)$$

[Eq. 16](#) is due to the definition of the kl divergence. [Eq. 17](#) is due to \hat{Y}^{ml} being a given random variable. [Eq. 18](#) is due to [Lemma 3](#). [Eq. 19](#) switches the integral subject to conditions of the dominated convergence theorem. [Eq. 20](#) is due to [Fig. 2](#) and the tower property. [Eq. 21](#) is due to the definition of the kl divergence and \hat{Y}^{ml} being given. [Eq. 22](#) is because setting $P(Y^{\text{eco}}(s)) = \int P(A)P(\hat{Y}^{\text{ml}}(A, s))dA$ has $\text{KL}(\int P(A)P(\hat{Y}^{\text{ml}}(A, s))dA || P(Y^{\text{FTU}}(s))) = 0$ for all s . The expectation is also zero: $\mathbb{E}_{P(S)} \left[\text{KL}(P(\int P(A)P(\hat{Y}^{\text{ml}}(A, s))dA) || P(Y^{\text{FTU}}(S))) \right] = 0$.

This calculation implies $\hat{Y}^{\text{FTU}} = \int P(\hat{Y}^{\text{ml}}(A, s))P(A | S = s)dA = P(Y | S = s)$ minimizes the average kl distance between the ML decision and the eco decision. In other words, the eco decision-maker maximally recovers the ML decision. Put differently, the eco decision-maker minimally modifies the ML decision to achieve eco-fairness.

B Proof of Theorem 2

Proof. We first prove that the cf decision-maker is cf-fair.

Recall the definition of the cf decision-maker,

$$P(\hat{Y}^{\text{cf}}(a_{\text{new}}, s_{\text{new}})) = \int \int P(\hat{Y}^{\text{eco}}(S(\bar{A}))) P(S(\bar{A}) | A = a_{\text{new}}, S = s_{\text{new}}) P(\bar{A}) dS(\bar{A}) d\bar{S}. \quad (24)$$

Note the above equation is the same definition as in [Eq. 2](#) except that the notation a is changed to \bar{A} to distinguish against the notation A used later in the proof.

Following this definition, we have

$$P(\hat{Y}^{\text{cf}}(a', S(a') | A = a, S = s)) \quad (25)$$

$$= \int \int P(\hat{Y}^{\text{eco}}(S(\bar{A}))) P(S(\bar{A}) | A = a', S = S(a')) \quad (26)$$

$$\times P(\bar{A})P(S(a') | A = a, S = s) dS(a') dS(\bar{A}) d\bar{A} \quad (27)$$

$$= \int \int P(\hat{Y}^{\text{eco}}(S(\bar{A}))) P(S(\bar{A}) | A = a, S = s) P(\bar{A}) dS(\bar{A}) d\bar{A}. \quad (28)$$

[Eq. 28](#) is due to the structural equation representation of the causal graph $S \stackrel{a, s}{\leftarrow} f(A, \epsilon)$, which implies $\int P(S(\bar{A}) | A = a', S = S(a'))P(S(a') | A = a, S = s) dS(a') = P(S(\bar{A}) | A = a, S = s)$. The intuition here is that the observed A, S will provide the same information as $a', S(a')$ for the same person. They all contain the information about the background variable that affects S ,

$$P(S(\bar{A}) | S(a'), s, a) = P(S(\bar{A}) | A = a', S = S(a')) = P(S(\bar{A}) | S(a'), a', s, a).$$

Notice that the right hand side of [Eq. 28](#) does not depend on a' , which implies that

$$P(\hat{Y}^{\text{cf}}(a', S(a') | A = a, S = s)) = P(\hat{Y}^{\text{cf}}(a, S(a)) | A = a, S = s).$$

This shows that the cf decision-maker is cf-fair, hence establishes the first part of [Thm. 2](#).

We then prove demographic parity for $\hat{Y}^{\text{cf}}(A, S)$. First we notice that

$$\int P(\hat{Y}^{\text{cf}}(A, S) | A)P(S | A) dS \quad (29)$$

$$= \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A') | A, S)P(A') dS(A') dA'P(S | A) dS \quad (30)$$

$$= \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A') | A)P(A') dS(A') dA' \quad (31)$$

$$= \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A'))P(A') dS(A') dA'. \quad (32)$$

The third equality is due to $S(A') \perp A$ by applying the Pearl's twin network strategy (Pearl, 2009). We then compute the marginal distribution of $\hat{Y}^{\text{cf}}(A, S)$.

$$\int P(\hat{Y}^{\text{cf}}(A, S))P(A, S) dA dS \quad (33)$$

$$= \int \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A') | A, S)P(A') dS(A') dA' P(A, S) dS dA \quad (34)$$

$$= \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A'))P(A') dS(A') dA', \quad (35)$$

which implies $P(\hat{Y}^{\text{cf}}(A, S) | A) = P(\hat{Y}^{\text{cf}}(A, S))$, i.e., demographic parity $\hat{Y}^{\text{cf}}(A, S) \perp A$.

We finally prove the second part of Thm. 2. We show that the cf decision-maker \hat{Y}^{cf} minimizes the average kl distance between the eco decision and the cf decision $\mathbb{E}_{P(S, A)} [\text{KL}(P(\hat{Y}^{\text{eco}}(A, S)) || P(Y^{\text{cf}}(A, S)))]$. In other words, the cf decision-maker maximally recovers the eco decision. Put differently, the cf decision-maker minimally modifies the eco decision to achieve cf-fairness.

To establish this optimality of the cf decision-maker, we first notice that the cf criterion implies that any cf decision-maker must take the form of $Y^{\text{cf}}(\epsilon)$, where $S = f(A, \epsilon)$ is the structural equation of the causal model in Fig. 2. We then compute the minimizer of the kl distance.

$$\begin{aligned} & \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} \mathbb{E}_{P(A)P(S)} [\text{KL}(P(\hat{Y}^{\text{eco}}(A, S)) || P(Y^{\text{cf}}(A, S)))] \\ &= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} \mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{eco}}(A, S)) [\log P(\hat{Y}^{\text{eco}}(A, S)) - \log P(Y^{\text{cf}}(A, S))] dY dA \right] \end{aligned} \quad (36)$$

$$= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} -\mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{eco}}(A, S)) \log P(Y^{\text{cf}}(A, S)) dY dA \right] \quad (37)$$

$$= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} -\mathbb{E}_{P(S)} \left[\int P(A) \int P(\hat{Y}^{\text{eco}}(f(A, \epsilon))) \log P(Y^{\text{cf}}(\epsilon)) dY dA \right] \quad (38)$$

$$= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} -\mathbb{E}_{P(S)} \left[\int \left[\int P(A')P(\hat{Y}^{\text{eco}}(f(A', \epsilon))) dA' \right] \log P(Y^{\text{cf}}(\epsilon)) dY \right] \quad (39)$$

$$= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} -\mathbb{E}_{P(S)} \left[\mathbb{E} \left[\int \left[\int P(A')P(\hat{Y}^{\text{eco}}(f(A', \epsilon))) dA' \right] \log P(Y^{\text{cf}}(s)) dY \mid A = a, S = s \right] \right] \quad (40)$$

$$= \arg \min_{Y^{\text{cf}} \in \mathcal{Y}^{\text{cf}}} \mathbb{E}_{P(S)} \left[\mathbb{E} \left[\text{KL} \left(\int P(A')P(\hat{Y}^{\text{eco}}(f(A', \epsilon))) dA' \mid \mid P(Y^{\text{cf}}(s)) \mid A = a, S = s \right) \right] \right] \quad (41)$$

$$= \int \int P(\hat{Y}^{\text{eco}}(f(A', \epsilon)))P(A')P(\epsilon | a, s) dA' d\epsilon \quad (42)$$

$$= \int \int P(\hat{Y}^{\text{eco}}(S(A')))P(S(A') | A = a, S = s)p(A') dS(A') dA' \quad (43)$$

$$= \hat{Y}^{\text{cf}}(a, s) \quad (44)$$

Eq. 36 is due to the definition of the kl divergence. Eq. 37 is due to \hat{Y}^{eco} being a given random variable. Eq. 38 is due to the observation that eco decision must only depend on $S = f(A, \epsilon)$ and cf decisions must only depend on ϵ . Eq. 39 switches the integral subject to conditions of the dominated convergence theorem. Eq. 40 is due to Fig. 2 and the tower property. Eq. 41 is due to the definition of the kl divergence and \hat{Y}^{eco} being given. Eq. 42 is because setting $P(Y^{\text{cf}}(\epsilon)) = \int P(A')P(\hat{Y}^{\text{eco}}(f(A', \epsilon))) dA'$ has $\text{KL}(\int P(A')P(\hat{Y}^{\text{eco}}(f(A', \epsilon))) dA' || P(Y^{\text{cf}}(\epsilon))) =$

0 for all ϵ . The expectation is also zero: $\mathbb{E} \left[\text{KL} \left(\int P(A') P(\hat{Y}^{\text{eco}}(f(A'), \epsilon)) dA' \middle| P(Y^{\text{cf}}(\epsilon)) \right) \right] = 0$. [Eq. 43](#) rewrites the ϵ in terms of A and S . [Eq. 44](#) is due to the definition of the cf decision-maker.

In addition, we show that \hat{Y}^{cf} recovers the marginal distribution of \hat{Y}^{eco} ,

$$\int P(\hat{Y}^{\text{cf}}(A, S)) P(A, S) dA dS \quad (45)$$

$$= \int P(\hat{Y}^{\text{eco}}(S(A'))) P(S(A') | A, S) P(A, S) P(A') dS(A') dA' dA dS \quad (46)$$

$$= \int P(\hat{Y}^{\text{eco}}(S)) P(A, S) dS(A) dA dS \quad (47)$$

$$= \int P(\hat{Y}^{\text{eco}}(S)) P(S) dS \quad (48)$$

$$= \int P(\hat{Y}^{\text{eco}}) P(A, S) dA dS. \quad (49)$$

The second equality is because $P(A) = P(A')$ by the construction of the cf decision-maker.

This calculation shows that the cf decision-maker \hat{Y}^{cf} preserves the marginal distribution of \hat{Y}^{eco} . \square

C The correctness of Algorithm 1

We prove the identifiability of the eco and cf decision given observational data, which establishes the correctness of [Alg. 1](#).

For the eco decision-maker, we utilize the backdoor adjustment formula ([Pearl, 2009](#)).

Proposition 4. Assume the causal graph [Fig. 2](#). The eco decision-maker can be computed using the observed data (A_i, S_i, Y_i) ,

$$P(\hat{Y}^{\text{ECO}}(a_{\text{new}}, s_{\text{new}})) = \int P(Y_i | A_i, S_i = s_{\text{new}}) P(A_i) dA_i, \quad (50)$$

if $P(S_i \in \mathbf{S} | A_i = a) > 0$ for all $a \in \mathcal{A}$ and sets \mathbf{S} such that $P(\mathbf{S}) > 0$ and $\mathbf{S} \subset \mathcal{S}$.

Proof. [Proposition 4](#) is a direct consequence of Theorem 3.3.2 of [Pearl \(2009\)](#). \square

[Eq. 50](#) reiterates the fact that \hat{Y}^{ECO} does not simply ignore the protected attribute A ; estimating $P(Y_i | A_i, S_i)$ relies on A in the training data. However, $\hat{Y}^{\text{ECO}}(a_{\text{new}}, s_{\text{new}})$ does not rely on the protected attribute a_{new} in predicting on the test data.

For the cf decision-maker, we adopt the abduction-action-prediction approach ([Pearl, 2009](#)).

Proposition 5. Assume the causal graph [Fig. 2](#). Write $S \stackrel{a,s}{=} f_S(A, \epsilon_S)$ for some function f_S and zero mean random variable ϵ_S satisfying $S \perp \epsilon_S$. The cf decision-maker can be computed using the observed data (A_i, S_i, Y_i) :

$$P(\hat{Y}^{\text{cf}}(a, s)) = \int P(\hat{Y}^{\text{eco}}(a, s')) \cdot P(S = s' | A = a', \epsilon_S) \quad (51)$$

$$\times P(\epsilon_S | A = a, S = s) \cdot P(A = a') ds' da' \quad (52)$$

$$\approx \int P(\hat{Y}^{\text{eco}}(\mathbb{E}[S_i | A_i = a'] + \hat{\epsilon}_S), a) \cdot P(A = a') da', \quad (53)$$

where $\hat{\epsilon}_S = s - \mathbb{E}[S_i | A_i = a]$.

Proof. [Eq. 52](#) is a direct consequence of Theorem 7.1.7 of [Pearl \(2009\)](#). [Eq. 53](#) is due to a linear approximation of f_S . \square

D Datasets in empirical studies

Each dataset in empirical studies consists of training and test sets, potentially with multiple protected and unprotected attributes.

The simulated college admissions data in [Fig. 1](#) comes from the structural model in [§ 4](#). We set $\beta_s = 2.0$ and $\beta_a = 1.0$ so that for the same test score the committee is more likely to admit males than females. We set $\lambda = 2.0$ to model that females perform less well on the test.

The case studies involve three public datasets.

- **Adult income data** ([Dua and Graff, 2017a](#)) The task is to decide whether an individual is loan worthy; the decision is binary. The data do not initially contain decisions about loan worthiness; instead, each individual’s income is provided. To evaluate the algorithms, we produce a decision about loan worthiness by thresholding income; individuals that make more than \$50K are decided as being loan worthy. The protected attributes are gender and race; other attributes include education and marital status. The dataset has 32,561 training samples and 16,282 test samples.
- **ProPublica’s COMPAS recidivism data.** The task is to decide an individual’s risk score of recidivating; the decision is real-valued. The protected attributes are gender and race; other attributes include the priors count, juvenile felonies count, and juvenile misdemeanor count. The dataset has 6,907 complete samples. We split them into 75% training and 25% testing.
- **German credit data.** ([Dua and Graff, 2017b](#)) The task is again to decide whether an individual is loan worthy. In this setting, each individual is initially labeled as having good or bad credit; we use this covariate to produce a decision. Individuals with good credit are decided as being loan worthy. The protected attributes are gender and marital status; other attributes include credit history, savings, and employment history. The dataset has 1,000 samples. We split them into 75% training and 25% testing.

E Implementation of decision-makers in the empirical studies

We implement classical ML decision-makers with linear/logistic regression using all attributes (including the protected ones). The ftu decision-makers perform regression but omit the protected attributes. The eco and cf decision-makers adjust the classical ML decision-makers using [Alg. 1](#). FairLearning infers background causal variables by fitting a linear regression of A against S and compute the residuals $\epsilon_i = S_i - \mathbb{E}[S_i | A_i]$; it then performs linear/logistic regression against the background variables. (This is the same procedure used in the abduction step for calculating f_{cf} .) With these decision-makers, we make algorithmic decisions \hat{y} about the test set units.

F Additional empirical studies

In the adult income data and COMPAS data, race and gender are the protected attributes. In the German credit data, gender and marital status are protected attributes. The decisions about loan worthiness are binary (as in the admissions example); the decisions in the COMPAS data are real-valued assessments of the recidivism score.

[Tables 1](#) and [2](#) present detailed results of the eco and cf decision-makers on the COMPAS and the German credit datasets.

Metrics ($\times 10^{-2}$) on COMPAS				
	ρ_{eco}	ρ_{cf}	KL	Prediction
ML predictor f_{ml}	-68.1(42.1)	-104.9(59.7)	17.1	28.0
ftu	0(0)	-52.9(30.1)	0.5	25.6
eco predictor f_{eco}	0(0)	-36.8(21.2)	0.5	25.6
FairLearning	-41.2(28.9)	0(0)	0.2	22.7
cf predictor f_{cf}	-41.2(29.4)	0(0)	0.2	22.7

Table 1: Both the eco predictor f_{eco} and ftu are eco-fair; they achieve zero in the eco metric (lower is better). cf predictor f_{cf} and FairLearning (Kusner et al., 2017) are cf-fair; they achieve zero in the cf metric (lower is better). cf predictor f_{cf} achieves demographic parity; it has close-to-zero kl divergence (Lower is better.) ML predictor f_{ml} predicts best; eco predicts best among the fair predictors (higher prediction scores are better.) We report mean values across individuals with the standard deviation in parentheses. eco and cf metric standard deviations are ≤ 0.42 and ≤ 0.6 , respectively. We also report the KL divergence and prediction accuracy, which are both distributional metrics.

Metrics ($\times 10^{-2}$) on German Credit				
	ρ_{eco}	ρ_{cf}	KL	Prediction
ML predictor f_{ml}	-5.4(2.5)	-3.9(1.8)	18.6	64.7
ftu	0(0)	-2.0(1.0)	15.6	63.4
eco predictor f_{eco}	0(0)	1.5(0.7)	13.2	64.5
FairLearning	-1.4(0.7)	0(0)	8.3	63.5
cf predictor f_{cf}	-1.3(0.6)	0(0)	7.8	64.3

Table 2: Both eco predictor f_{eco} and ftu are eco-fair; they achieve zero in the eco metric (lower is better). cf predictor f_{cf} and FairLearning (Kusner et al., 2017) are cf-fair; they achieve zero in the cf metric (lower is better). cf predictor f_{cf} achieves demographic parity; it has close-to-zero kl divergence (Lower is better.) ML predictor f_{ml} predicts best; eco predicts best among the fair predictors (higher prediction scores are better.) We report mean values across individuals with the standard deviation in parentheses. eco and cf metric standard deviations are ≤ 0.02 and ≤ 0.02 , respectively. We also report the KL divergence and prediction accuracy, which are both distributional metrics.

References

- Dheeru Dua and Casey Graff. 2017a. UCI Machine Learning Repository: Adult Data Set. <https://www.kaggle.com/wenruliu/adult-income-dataset>
- Dheeru Dua and Casey Graff. 2017b. UCI Machine Learning Repository: Statlog (German Credit Data) Data Set. <https://www.kaggle.com/uciml/german-credit>
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. 2017. Counterfactual fairness. *Neural Information Processing Systems* (2017).
- Judea Pearl. 2009. *Causality*. Cambridge University Press.