# Fair Policy Targeting

## Davide Viviano & Jelena Bradic

Taylor & Francis
Taylor & Francis Group

Check for updates

# Fair Policy Targeting

Davide Viviano[a] and Jelena Bradic[b]

[a]Stanford Graduate School of Business, Stanford, CA, and Department of Economics, Harvard University, Cambridge, MA; [b]Department of Mathematics and Haliciŏglu Data Science Institute, University of California at San Diego, La Jolla, CA

## ABSTRACT

One of the major concerns of targeting interventions on individuals in social welfare programs is discrimination: individualized treatments may induce disparities across sensitive attributes such as age, gender, or race. This article addresses the question of the design of fair and efficient treatment allocation rules. We adopt the nonmaleficence perspective of "first do no harm": we select the fairest allocation *within* the Pareto frontier. We cast the optimization into a mixed-integer linear program formulation, which can be solved using off-the-shelf algorithms. We derive regret bounds on the unfairness of the estimated policy function and small sample guarantees on the Pareto frontier under general notions of fairness. Finally, we illustrate our method using an application from education economics. Supplementary materials for this article are available online.

## 1. Introduction

Heterogeneity in treatment effects, widely documented in social sciences, motivates treatment allocation rules that assign treatments to individuals differently based on observable characteristics (Murphy 2003; Manski 2004). However, targeting individuals may induce disparities across sensitive attributes, such as age, gender, or race. Motivated by evidence for policymakers' preferences toward nondiscriminatory actions (Cowgill and Tucker 2019), this article designs fair and efficient targeting rules for applications in welfare and health programs. We construct treatment allocation rules using data from experiments or quasi-experiments, and we develop policies that trade off efficiency and fairness.

Fair targeting is a controversial task due to the lack of consensus on formulating the decision problem. Conventional approaches mostly developed in computer science consist in designing algorithmic decisions that maximize the *expected* utility across all individuals by imposing fairness constraints on the decision space of the policymaker (Nabi, Malinsky, and Shpitser 2019).[1] In contrast, the economic literature has outlined the importance of taking into account the welfare effects of such policies (Kleinberg et al. 2018). Fairness constraints on the policymaker's decision space may ultimately lead to suboptimal welfare for both sensitive groups. This is a significant limitation when the policymakers are concerned with the effects

of their decisions on each individual's utility: absent of legal constraints, we may not want to impose unnecessary constraints on the policy if such constraints are *harmful* to some or all individuals.

This article studies the design of fair and Pareto optimal treatment rules. We discuss targeting in a setting where decision-makers prefer allocations for which we cannot find any other policy that strictly improves the welfare of one of the two sensitive groups without decreasing the welfare of the opposite group. Within such a set, she then chooses the fairest allocation. The decision problem is conceived for applications in social welfare and health programs and motivated by the Hippocratic notion of "first do no harm" ("*primum non-nocere*") (Rotblat 1999): instead of imposing possibly harmful fairness constraints on the decision space, we restrict the set of admissible solutions to the Pareto optimal set, and among such, we choose the fairest one. For example, during a health-program campaign, the policymakers may not be willing to decrease all individuals' health status to gain fairness. Instead, they may be willing to *trade off* health status of different groups (e.g., young and old individuals) when considering fairness. Our framework has three desirable properties: (i) it applies to general notions of fairness which may reflect different decision-makers' preferences; (ii) it guarantees Pareto efficiency of the policy function, with the relative importance of each group solely chosen based on the notion of fairness adopted by the decision-maker; (iii) it also allows for arbitrary legal or ethical constraints, incorporating as a special case the presence of fairness constraints whenever such constraints are

---

[1] For a review, the reader may refer to Corbett-Davies and Goel (2018). Further discussion on the related literature is contained in Section 1.1.

**CONTACT** Davide Viviano ✉ *dviviano@fas.harvard.edu* 🖂 Stanford Graduate School of Business, Stanford, CA, and Department of Economics, Harvard University, Cambridge, MA.

This work was mostly conducted while at the Department of Economics, University of California at San Diego.

*binding* due to ethical or legal considerations.[2] We name our method Fair Policy Targeting.

This article contributes to the statistical treatment choice literature by introducing the notion, estimation procedure and studying the properties of Pareto optimal and fair treatment allocation rules. We allow for general notions of fairness, and as a contribution of independent interest, we study envy-freeness (Varian 1976) in the context of policy targeting.

The decision problem consists of lexicographic preferences of the policymaker of the following form: (i) Pareto-dominant allocations are preferred over dominated ones; (ii) Pareto optimal allocations are ranked based on fairness considerations. We identify the Pareto frontier as the set of maximizers over any weighted average of each group's welfare. Therefore, such an approach embeds as a special case maximizing a weighted combination of the welfare of each sensitive group such as in Athey and Wager (2021), Kitagawa and Tetenov (2018),[3] and in Rambachan et al. (2020). The above references take a specific weighted combination of welfare with weights as given. In contrast, in our case, weights are part of the decision problem and are directly selected to maximize fairness. This has important practical implications: our procedure is solely based on the notion of fairness adopted by the social planner. It does not require specific importance weights assigned to each sensitive group, which would be hard to justify to the general public.

Estimating the set of Pareto optimal allocations represents a fundamental challenge since (i) the set consists of maximizers over a continuum of weights between zero and one; (ii) each maximizer of the welfare (or a weighted combination of welfares) is often not unique (Elliott and Lieli 2013). To overcome these issues, we show that the Pareto frontier can be approximated using simple linear constraints. We use a discretization argument, and we evaluate weighted combinations of the objective functions separately to construct a polyhedron that contains Pareto allocations. Our approach drastically simplifies the optimization algorithm: instead of estimating the entire set of Pareto allocations, we maximize fairness under easy-to-implement linear constraints. Our theorems show that the distance between the Pareto frontier obtained via linear constraints and its population counterpart converges uniformly to zero at a rate $1/\sqrt{n}$.

We study regret guarantees, that is, the difference between the estimated policy function's expected unfairness against the minimal possible unfairness achieved by Pareto optimal allocations. We characterize the rate under high-level conditions for general notions of unfairness and derive upper bounds that scale at rate $1/\sqrt{n}$, in several examples, and a lower bound that matches the same rate. An application and a calibrated numerical study on targeting student awards illustrate the advantages of the proposed method compared to alternatives that ignore Pareto optimality.

The remainder of this article is organized as follows: we provide a brief overview of the literature in the following section; we introduce the decision problem in Section 2; Section 3 discusses estimation; Section 4 contains the theoretical analysis; Section 5 discusses counterfactual notions of fairness; Section 6 presents an empirical application and numerical studies, and Section 7 concludes. Derivations and extensions are in the online supplementary materials.

### 1.1. Related Literature

This article relates to a growing literature on statistical treatment rules (Manski 2004; Hirano and Porter 2009; Bhattacharya and Dupas 2012; Stoye 2012; Tetenov 2012; Armstrong and Shen 2015; Kitagawa and Tetenov 2018, 2019; Zhou, Athey, and Wager 2018; Viviano 2019; Sun 2021; Athey and Wager 2021; Mbakop and Tabord-Meehan 2021). Further connections are also related to the literature on classification (Elliott and Lieli 2013). However, none discuss the design of fair and Pareto optimal decisions.

Fairness is a rising concern in economics, see Cowgill and Tucker (2019), Kleinberg et al. (2018), and Rambachan et al. (2020). The authors provide economic insights on the characteristics of optimal decision rules when discrimination bias occurs. Here, we answer the different questions of the design and estimation of the optimal targeting rule within a statistical framework and derive the method's properties. A further difference from the references above is our focus on a multi-objective instead of a single-objective decision problem. Additional references include Kasy and Abebe (2021) that provide comparative statics on the impact of fairness on the individuals' welfare, focusing on the analysis of algorithms, while Narita (2021) motivates fairness based on incentive compatibility in the different context of the design of experiments.

In computer science, Pareto optimality has been considered in the context of binary predictions by Balashankar et al. (2019) and Martinez, Bertran, and Sapiro (2019). The authors propose semi-heuristic and computationally intensive procedures for estimating Pareto efficient classifiers. Xiao et al. (2017) discuss the different problems of estimation of a Pareto allocation that balance fairness and individual utilities for recommender systems, where the relative importance weights of the different objectives are selected a-priori. These references do not address the treatment choice problem discussed in the current article.

References in computer science include Chouldechova (2017), Dwork et al. (2012), Hardt, Price, and Srebro (2016) among others. Corbett-Davies and Goel (2018) contain a review. Additional work also includes Liu et al. (2017), who discuss fair bandits, and Ustun, Liu, and Parkes (2019), who propose decoupled estimation of tree classifiers without allowing for exogenous (legal or economic) constraints on the policy space. While the above references address the decision problem as a prediction problem, several papers discuss algorithmic fairness within a causal framework (Kilbertus et al. 2017; Nabi, Malinsky, and Shpitser 2019; Kusner et al. 2019; Coston et al. 2020). All such papers estimate decision rules under fairness constraints without discussing Pareto optimality. The different decision problem considered here is motivated by applications in social welfare and health programs. When not binding

---

[2] For binding fairness constraints, the policy that we propose is Pareto optimal within the set of constrained policies and achieves a lower unfairness compared to the policy that maximizes the utilitarian welfare under fairness constraints. See Section 2.3 for details.

[3] Under the utilitarian perspective considered in Kitagawa and Tetenov (2018), Athey and Wager (2021), the welfare maximization problem is equivalent to maximizing a weighted combination of the welfare of different groups with weights equal to the corresponding probabilities of such groups. See Section 2 for more discussion.

on policy-makers decisions, fairness constraints may lead to Pareto-dominated allocations and possibly harmful policies for advantaged and disadvantaged individuals. When fairness constraints are binding, the decision problem proposed in this article leads to fairer allocations than a constrained welfare maximization problem while not being Pareto-dominated.

## 2. Decision Making and Fairness

We start by introducing some notation. For each unit, we denote with $S \in \mathcal{S}$ a sensitive or protected attribute. For expositional convenience, we let $\mathcal{S} = \{0, 1\}$, with $S = 1$ denoting the disadvantaged group, and $X \in \mathcal{X} \subseteq \mathbb{R}^p$ individual characteristics. We define the post-treatment outcome with $Y \in \mathcal{Y} \subseteq \mathbb{R}$ realized only once the sensitive attribute, covariates, and the treatment assignment are realized. We define $Y(d)$, $d \in \{0, 1\}$ the potential outcomes under treatment $d$. The observed $Y$ satisfies the Single Unit Treatment Value Assumption (SUTVA) (Rubin 1990). Let

$$e(x, s) = P(D = 1 | X = x, S = s), \quad p_1 = P(S = 1) \quad (1)$$

be the propensity score and the probability of being assigned to the disadvantaged group. Here, treatments are independent of potential outcomes.

*Assumption 2.1 (Treatment Unconfoundedness).* For $d \in \{0, 1\}$, $Y(d) \perp D | X, S$.

### 2.1. Social Welfare

Given observables, $(Y_i, X_i, D_i, S_i)$ we seek to design a treatment assignment rule (i.e., policy function) $\pi : \mathcal{X} \times \mathcal{S} \mapsto \mathcal{T} \subseteq [0, 1], \pi \in \Pi$ that depends on the individual characteristics and protected attributes, and which can be either probabilistic or deterministic.[4] Here, $\Pi$ incorporates given and binding legal or economic constraints that restrict the decision space. The welfare generated by a policy $\pi$ on those individuals with sensitive attribute $S = s$ is defined as[5]

$$W_s(\pi) = \mathbb{E}\Big[(Y(1) - Y(0))\pi(X, S)\Big| S = s\Big]. \quad (2)$$

Under the utilitarian perspective (Manski 2004), the welfare maximization problem, that is, the population counterpart of the empirical welfare maximization (EWM) (Kitagawa and Tetenov 2018), solves

$$\max_{\pi \in \Pi} \Big\{ p_1 W_1(\pi) + (1 - p_1) W_0(\pi) \Big\}$$

where $p_1$ is defined as in Equation (1). However, whenever the sensitive group is a *minority* group, welfare maximization assigns a small weight to the welfare of the minority, disproportionally favoring the majority group. An alternative approach is to maximize the welfare separately for each possible sensitive group designing different policies for different groups (Ustun, Liu, and Parkes 2019). This approach may violate discriminatory laws, that is, the resulting policy function violates the constraint in $\Pi$. A simple example is when, due to legal reasons, the policy

$\pi(x, s)$ must be constant in the sensitive attribute $s$. Instead, we consider a framework where the policymaker simultaneously maximizes each group's welfare, imposing Pareto efficiency on the estimated policy under arbitrary legal or economic constraints encoded in $\Pi$. Given the set of efficient policies, the planner then selects the least unfair one. Our approach is designed for social and welfare programs where legal constraints naturally occur and where, given such constraints, the policymaker's preferences align with classical notions of "first do no harm."

### 2.2. Pareto Principle for Treatment Rules

The set of Pareto optimal choices is defined as $\Pi_o$, and it contains all such allocations $\pi \in \Pi$ for which the welfare for one of the two groups cannot be improved without reducing the welfare for the opposite group. We characterize $\Pi_o$ in the following lemma.

*Lemma 2.1 (Pareto Frontier).* The set $\Pi_o \subseteq \Pi$ is such that

$$\Pi_o = \Big\{\pi_\alpha : \pi_\alpha \in \arg\sup_{\pi \in \Pi} \alpha W_1(\pi) + (1 - \alpha) W_0(\pi), \ \alpha \in (0, 1)\Big\}. \quad (3)$$

The lemma follows from Negishi (1960), whose proof is in Appendix A, supplementary materials. It will be convenient to define

$$\bar{W}_\alpha = \sup_{\pi \in \Pi} \alpha W_1(\pi) + (1 - \alpha) W_0(\pi), \quad (4)$$

the largest value of the objective in Equation (3) for a fixed $\alpha$. In the following examples, we show that Pareto allocations *generalize* notions of treatment rules from previous literature.

*Example 2.1 (Welfare Maximization).* The population equivalent of the EWM problem belongs to the Pareto frontier. Namely, $\arg\max_{\pi \in \Pi} \Big\{ p_1 W_1(\pi) + (1 - p_1) W_0(\pi) \Big\} \subseteq \Pi_o$. An alternative approach consists in maximizing weighted combinations of the welfare with the weights for each group as given. For instance the allocation (Rambachan et al. 2020)

$$\check{\pi}_\omega \in \arg\max_{\pi \in \Pi} \Big\{ \omega W_1(\pi) + (1 - \omega) W_0(\pi) \Big\} \subseteq \Pi_o, \quad (5)$$

for some *specific* weight $\omega$ belongs to the Pareto frontier. $\square$

Pareto optimal allocations are often nonunique, allowing for flexibility in the choice of efficient policies. The policy-maker must appeal to some preferential ranking principle based on her preferences. We discuss those in the following lines.

### 2.3. Decision Problem

We start by defining $\mathcal{C}(\Pi)$ the *choice set* of the policy maker (Mas-Colell, Whinston, and Green 1995), where $\mathcal{C}$ is a choice function with $\mathcal{C}(\{\pi_1, \pi_2\}) = \pi_1$ if $\pi_1$ is strictly preferred to $\pi_2$. We let

$$\text{UnFairness} : \Pi \mapsto \mathbb{R} \quad (6)$$

an operator which quantifies the unfairness of a policy. We leave unspecified UnFairness and provide examples in Sections 4.3 and 5. We now state the planner's preferences.

---

[4]It is deterministic if $\mathcal{T} = \{0, 1\}$ and probabilistic if $\mathcal{T} = [0, 1]$.

[5]Welfare is interpreted from an *intention-to-treat* perspective similarly to Kitagawa and Tetenov (2018) and Athey and Wager (2021).

*Assumption 2.2 (Policy-maker's Preferences).* Preferences are rational,[6] and for each $\pi_1, \pi_2 \in \Pi$, (i) $\mathcal{C}(\{\pi_1, \pi_2\}) = \pi_1$ if $W_1(\pi_1) \geq W_1(\pi_2)$ and $W_0(\pi_1) \geq W_0(\pi_2)$ and either (or both) of the two inequalities hold strictly; (ii) if neither $\pi_1$ Pareto dominates $\pi_2$ nor $\pi_2$ Pareto dominates $\pi_1$, $\mathcal{C}(\{\pi_1, \pi_2\}) = \pi_1$ if UnFairness$(\pi_1) <$ UnFairness$(\pi_2)$; (iii) if neither Pareto dominates the other and with equal UnFairness, $\mathcal{C}(\{\pi_1, \pi_2\}) = \{\pi_1, \pi_2\}$.

Assumption 2.2 postulates lexicographic preferences of the following form: (i) an allocation is strictly preferred to another if it weakly improves welfare for both groups and strictly improves welfare for at least one group; (ii) given two allocations where none of the two Pareto dominates the other, allocations are ranked based on fairness.

While different applications reflect different planner's preferences, Assumption 2.2 is motivated by the applications for social welfare and health programs, where welfare depends on outcomes such as health status (Finkelstein et al. 2012), future earnings, or school achievements. Sacrificing each group's welfare (e.g., *health-status*) for fairness is undesirable in such applications. Conditional on achieving a Pareto efficient allocation, the planner minimizes UnFairness. Whenever, however, fairness constraints are binding (e.g., because of legal considerations), these can be directly incorporated into the function class $\Pi$. We can now characterize the decision problem.

*Proposition 2.2 (Decision Problem).* Under Assumption 2.2, $\pi^\star \in \mathcal{C}(\Pi)$ if and only if

$$\pi^\star \in \arg \inf_{\pi \in \Pi} \text{UnFairness}(\pi)$$
$$\text{subject to } \alpha W_1(\pi) + (1 - \alpha) W_0(\pi) \geq \bar{W}_\alpha, \text{ for some } \alpha \in (0, 1).$$
(7)

The proof is contained in Appendix A, supplementary materials. Proposition 2.2 formally characterizes the policy-makers decision problem, which consists of minimizing the policy's unfairness criterion under the condition that the policy is Pareto optimal. The policy-maker does not maximize a weighted combination of welfares, with some *pre-specified* and hard-to-justify weights. Instead, each group's importance (i.e., $\alpha$) is implicitly chosen within the optimization problem to maximize fairness. This approach allows for a transparent choice of policy based on the policy-makers definition of fairness.

*Example 2.2 (Why Pareto Efficiency? A simple example).* Let $X = 1$ for simplicity, take $\tau_s, \phi \in (0, 1), s \in \{0, 1\}$ and let $Y(d) = \tau_S d + \varepsilon(d)$, with $\mathbb{E}[\varepsilon(d)|S] = 0$. Consider a class of probabilistic decision rules

$$\Pi = \left\{ \pi(x, s) = \beta_s, \quad \beta_1, \beta_0 \in (0, 1), \quad \beta_0 p_0 + \beta_1 p_1 \leq \phi \right\},$$

with the share of treated units being at most $\phi$. Let UnFairness be the difference in the groups' welfares, namely $|\tau_1 \beta_1 - \tau_0 \beta_0|$. The smallest possible unfairness is zero, since we can choose $\beta_1 = \beta_0 = 0$ with one of the fairest allocation selecting none of the individuals to treatment. Consider now the Pareto frontier, defined as

$$\Pi_\circ = \left\{ \pi(x, s) = \beta_s^*, \quad \beta_0^* = \frac{\phi - p_1 \beta_1^*}{p_0}, \beta_1^* \in [0, 1] \right\} \subset \Pi.$$
(8)

The set of Pareto allocation rules out all those allocation for which the capacity constraint is attained with strict inequality, also excluding $\beta_1 = \beta_0 = 0$. The proposed policy assigns all benefits to individuals, and it tradeoffs *who* to treat to minimize $|\tau_1 \beta_1 - \tau_0 \beta_0|$.[7]    □

We conclude by comparing the properties of the policy in Proposition 2.2 with existing alternatives, stated as a corollary of Proposition 2.2. In particular, we compare our method with the policy that maximizes welfare with importance weights for different groups (Rambachan et al. 2020) in Equation (5) and the one with fairness constraints. For the latter, define $\Pi(\kappa) = \left\{ \pi \in \Pi : \text{UnFairness}(\pi) \leq \kappa \right\} \subseteq \Pi$, the set of policies with constraints, and

$$\tilde{\pi} \in \arg \max_{\pi \in \Pi(\kappa)} p_1 W_1(\pi) + (1 - p_1) W_0(\pi)$$
(9)

the policy that maximizes welfare by imposing fairness constraints (e.g., Nabi, Malinsky, and Shpitser 2019).

*Corollary 1 (Properties).* Let $\pi^\star$ be defined as in Equation (7) and $\pi_\omega, \tilde{\pi}$ and in Equations (5) and (9), respectively. Then UnFairness$(\pi^\star) \leq$ UnFairness$(\tilde{\pi}_\omega), \forall \omega \in (0, 1)$.

Suppose that either $\tilde{\pi} \in \Pi_\circ$ (i.e., it belongs to the Pareto frontier), or fairness constraints are binding to the policy-maker, that is, $\Pi(\kappa) = \Pi$. Then UnFairness$(\pi^\star) \leq$ UnFairness$(\tilde{\pi})$. Suppose instead that $\tilde{\pi} \notin \Pi_\circ$. Then UnFairness$(\pi^\star) \leq$ UnFairness$(\pi_\circ)$ for all $\pi_\circ \in \Pi_\circ$ that Pareto dominate $\tilde{\pi}$. In addition, $\pi_\omega$ and $\tilde{\pi}$ do not Pareto dominate $\pi^\star$.

Corollary 1 shows that UnFairness of $\pi^\star$ is Pareto optimal and smaller than UnFairness of the policy $\pi_\omega$ that maximizes a weighted combination of the welfares. It also shows that if $\tilde{\pi}$ *is* Pareto optimal, then its UnFairness is larger than UnFairness of $\pi^\star$. When instead $\tilde{\pi}$ is *not* Pareto optimal, its Pareto dominant allocations have larger UnFairness than $\pi^\star$. Further intuition can be gained under strong duality, which we discuss in Appendix B.1, supplementary materials. Intuitively, the constraint in Proposition 2.2 holding for *some* weighted combinations of welfares (instead of a particular choice of the weights) is key to achieve lower unfairness of $\pi^\star$ relative to $\tilde{\pi}$ when $\tilde{\pi}$ is Pareto efficient.[8] Finally, when fairness constraints are binding, the proposed procedure always leads to smaller UnFairness.

---

[6]Rational preferences imply transitivity and completeness (Mas-Colell, Whinston, and Green 1995).

[7]Observe that the level of unfairness with the frontier may or may not be potentially strictly larger than the unfairness obtained in an unconstrained scenario. Namely, to achieve zero unfairness for every $\pi \in \Pi_\circ$, we need that $\tau_1 \beta_1^* = \tau_0 \beta_0^*$. Substituting $\beta_0^* = \phi/p_0 - p_1 \beta_1^*/p_0$ this would require $\beta_1^* = \frac{\phi}{p_0}(\tau_1/\tau_0 + p_1/p_0)^{-1}$ which is not necessarily feasible (i.e., the expression is larger than one).

[8]Under strong duality, the dual of $\tilde{\pi}$ corresponds to minimize UnFairness for *one particular* weighted combination of welfare exceeding a certain threshold. In contrast, our decision problem imposes the constraint that *some* weighted combination of welfare exceeding a certain threshold. This difference reflects the difference between the lexicographic preferences we propose as opposed to an additive social planner's utility.

## 3. Fair Targeting: Estimation

We now construct an estimator of $\pi^\star$. We introduce some notation, and we define

$$m_{d,s}(x) = \mathbb{E}\Big[Y_i(d)\Big|X_i = x, S_i = s\Big],$$

$$\Gamma_{d,s,i} = \frac{1\{S_i = s\}}{p_s}\Big[\frac{1\{D_i = d\}}{e(X_i, S_i)}\Big(Y_i - m_{d,s}(X_i)\Big) + m_{d,s}(X_i)\Big] \tag{10}$$

the conditional mean of the group $s$ under treatment $d$, and the doubly robust score (Robins and Rotnitzky 1995), respectively. We let $\hat{\Gamma}_{d,s,i}$ the estimated counterpart of $\Gamma_{d,s,i}$. Define

$$\hat{W}_s(\pi) = \frac{1}{n}\sum_{i=1}^{n}\Big(\hat{\Gamma}_{1,s,i} - \hat{\Gamma}_{0,s,i}\Big)\pi(X_i, s). \tag{11}$$

the estimated welfare built upon semiparametric literature (Newey 1990; Robins and Rotnitzky 1995), with $\hat{m}_{d,s}(.), \hat{e}(.), \hat{p}_s$, constructed via cross-fitting (Chernozhukov, Newey, and Robins 2018). Details of the cross-fitting procedure are contained in Appendix B.2, supplementary materials. We consider first general notions of fairness, and introduce the corresponding estimator below.

*Definition 3.1 (Empirical UnFairness).* We define $\mathcal{V}_n(\pi, p_s, e, m)$ an unbiased estimate of UnFairness$(\pi)$ which depends on observables and the population propensity score and conditional mean. We write $\hat{\mathcal{V}}_n(\pi) = \mathcal{V}_n(\pi, \hat{p}_s, \hat{e}, \hat{m})$, the empirical counterpart.

We defer to Sections 4.3 and 5 explicit examples of $\hat{\mathcal{V}}_n(\pi)$.

### 3.1. (Approximate) Pareto Optimality

Next, we characterize the Pareto frontier using linear inequalities. To construct the Pareto frontier we use the constraint in Equation (7) after discretizing the set of weights $\alpha$. Namely, in the first step, we discretize the Pareto frontier, and construct a grid of equally spaced values $\alpha_j \in (0, 1), j \in \{1, \ldots, N\}$, with $N = \sqrt{n}$. We approximate the Pareto frontier using the set ($\hat{W}_0, \hat{W}_1$ are defined in Equation (11))

$$\hat{\Pi}_\circ = \Big\{\pi_\alpha \in \Pi : \pi_\alpha \in \arg\sup_{\pi \in \Pi}\Big\{\alpha\hat{W}_0(\pi) + (1-\alpha)\hat{W}_1(\pi)\Big\},$$

$$\text{s.t. } \alpha \in \{\alpha_1, \ldots, \alpha_N\}\Big\}. \tag{12}$$

The grid's choice is arbitrary, as long as values are *equally spaced*.

The set $\hat{\Pi}_\circ$ may be hard, if not impossible, to directly estimate, since we may have uncountably many solutions (Manski and Thompson 1989; Elliott and Lieli 2013). In particular, the solution to each optimization problem in Equation (12) may not be unique. Instead of directly estimating $\hat{\Pi}_\circ$, we characterize it through linear constraints. First, we find the largest empirical welfare achieved on the discretized Pareto Frontier defined as

$$\bar{W}_{j,n} = \sup_{\pi \in \Pi}\Big\{\alpha_j\hat{W}_0(\pi) + (1-\alpha_j)\hat{W}_1(\pi)\Big\}, \text{ for each } j \in \{1, \ldots, N\}, \tag{13}$$

which can be obtained through standard optimization routines (Kitagawa and Tetenov 2018; Zhou, Athey, and Wager 2018). Second, we observe that any $\pi \in \hat{\Pi}_\circ$, must satisfy $\alpha_j\hat{W}_0(\pi) + (1 - \alpha_j)\hat{W}_1(\pi) \geq \bar{W}_{j,n}$, for some $j \in \{1, \ldots, N\}$, since $\bar{W}_{j,n}$ defines the largest objective for a given $\alpha_j$. We impose such constraint up to a small slackness parameters $\lambda/\sqrt{n}$ and construct an approximate Pareto frontier as follows:

$$\hat{\Pi}_\circ(\lambda) = \Big\{\pi \in \Pi : \exists j \in \{1, \ldots, N\} \text{ such that}$$

$$\alpha_j\hat{W}_{0,n}(\pi) + (1 - \alpha_j)\hat{W}_{1,n}(\pi) \geq \bar{W}_{j,n} - \frac{\lambda}{\sqrt{n}}\Big\}, \tag{14}$$

where $\hat{\Pi}_\circ(0) = \hat{\Pi}_\circ$, and $\hat{\Pi}_\circ \subseteq \hat{\Pi}_\circ(\lambda)$ for any $\lambda \geq 0$.

Here, we introduced $-\frac{\lambda}{\sqrt{n}}$ which imposes that the resulting policy is "approximately" Pareto optimal. As shown in Section 4, $\lambda/\sqrt{n}$ guarantees that $\hat{\Pi}_\circ(\lambda)$ contains all Pareto optimal policies with high-probability, for $\lambda = \mathcal{O}(1)$. The estimated policy is defined as

$$\hat{\pi}_\lambda \in \arg\min_{\pi \in \hat{\Pi}_\circ(\lambda)} \hat{\mathcal{V}}_n(\pi). \tag{15}$$

*Remark 1 (The choice of the grid and $\lambda$).* The choice of $\lambda$ depends on the function class $\Pi$. In Theorems 4.2, 4.3, and 4.4 we discuss guarantees by imposing that $\lambda/\sqrt{n} \geq M\sqrt{v}/N$, for some finite constant $M$ with $\lambda$ increasing in the geometric complexity $v$ of $\Pi$, and where we choose $N = \sqrt{n}$.[9] In contrast, the function class complexity does not affect the choice of the grid (i.e., $N$). This is because the welfare loss due to the grid's approximation error is uniformly bounded by a constant independent of $\Pi$.[10] □

### 3.2. Optimization: Mixed Integer Quadratic Program

We provide a mixed-integer quadratic program (MIQP) for optimization. We define $\mathbf{z}_s = (z_{s,1}, \ldots, z_{s,n}), z_{s,i} = \pi(X_i, s), \pi \in \Pi$. Here, $z_{s,n}$ defines the treatment assignment under policy $\pi$, and sensitive attribute $s$ (see the example below); $\mathbf{z}_s$ have simple representation for general classes of policy functions, such as either probabilistic rules which we derive in Appendix B.2, supplementary materials or deterministic linear decision rules (Florios and Skouras 2008).

*Example 3.1 (Maximum score).* For the maximum score $\pi(X_i, s) = 1\{X_i\beta_x + S\mu \geq 0\}, \beta = (\beta_x, \mu) \in \mathcal{B}$, the indicators $z_{s,n}$ are defined via mixed-integer constraints of the form (Florios and Skouras 2008) $\frac{X_i^\top\beta + s\mu}{|C_i|} < z_{s,i} \leq \frac{X_i^\top\beta + s\mu}{|C_i|} + 1, C_i \geq \sup_{\beta \in \mathcal{B}}|(X_i, S_i)^\top\beta|, z_{s,i} \in \{0, 1\}$. Such constraint guarantees that $z_{s,i} = 1\{X_i^\top\beta_x + s\mu \geq 0\}$. □

---

[9]This guarantees that the estimated function class does not exclude Pareto optimal policies with high probability, while $\lambda = \mathcal{O}(1)$ guarantees uniform convergence of the Pareto frontier at $1/\sqrt{n}$ rate.

[10]Namely, take a grid of $N+1$ equally spaced $\alpha_j$. Then the approximation error reads as $\sup_{\pi \in \Pi}|\alpha W_1(\pi) + (1 - \alpha)W_0(\pi) - \max_{\alpha_j \in \{\alpha_1, \ldots, \alpha_N\}}\alpha_j W_1(\pi) - (1-\alpha)W_0(\pi)| \leq 2M/N$, which is uniformly bounded by $M$ where $M$ bounds the first moment of the potential outcomes, and independent of $\Pi$.

We now need to impose the constraint of Pareto optimality. To do so, we introduce an additional set of decision variables that guarantee the constraints in Equation (14) hold. The vector $\mathbf{u} = (u_1, \ldots, u_N) \in \{0,1\}^N$ encodes the locations on the grid of $\alpha$ for which the supremum in (14) is reached at; here, $u_j = 1$ whenever the constraint in Equation (14) holds for $\alpha_j$. The chosen policy must be Pareto optimal, that is, $u_j$ must be equal to one for at least one $j$. To ensure this, we impose the constraint $\sum_{j=1}^N u_j \geq 1$.

Combining such constraints, it directly follows that $\hat{\pi}_\lambda$ satisfies Equation (15) if and only if

$$\hat{\pi}_\lambda \in \arg\min_\pi \min_{\mathbf{z_0}, \mathbf{z_1}, \mathbf{u}} \widehat{\mathcal{V}}_n(\pi) \tag{16}$$

subject to
$$z_{s,i} = \pi(X_i, s), \quad 1 \leq i \leq n, \tag{A}$$

$$\begin{aligned} u_j \alpha_j \langle \hat{\mathbf{\Gamma}}_{1,0} - \hat{\mathbf{\Gamma}}_{0,0}, \mathbf{z_0} \rangle \\ + u_j(1 - \alpha_j) \langle \hat{\mathbf{\Gamma}}_{1,1} - \hat{\mathbf{\Gamma}}_{0,1}, \mathbf{z_1} \rangle \\ \geq u_j n \bar{W}_{j,n} - \sqrt{n}\lambda \end{aligned} \tag{B}$$

$$\langle \mathbf{1}, \mathbf{u} \rangle \geq 1 \tag{C}$$

$$\pi \in \Pi \tag{D}$$

$$u_j \in \{0,1\}, \quad 1 \leq j \leq N. \tag{E}$$

Here, $\mathbf{\Gamma}_{d,s}$ is the vector of $\Gamma_{d,s,i}$ defined in Equation (10). Constraints (B) and (C) state that the resulting policy is (approximately) Pareto optimal, or, equivalently, it maximizes a weighted combination of groups' welfare for *some* $\alpha_j$. Constraints (A), (C), (E) are (mixed-integer) linear constraints, while Constraint (B) is quadratic. Notice that we can further simplify (B) as a linear constraint at the expense of introducing additional $Nn$ binary variables and $2Nn$ additional constraints (e.g., see Wolsey and Nemhauser 1999; Viviano 2019). Finally, (D) is either linear or quadratic for deterministic assignments and linear probability models. Hence, the objective admits a MIQP representation whenever $\widehat{\mathcal{V}}_n(\pi)$ admits linear representation in $\pi$, as discussed in the following section. Note that the solution to the optimization problem might not be unique, depending on the function class. However, non-uniqueness does not affect theoretical properties in Section 4.

*Remark 2 (Computational complexity).* The complexity of the optimization problem depends on the policy function class. For discrete covariates, in Section 4.4 we show that the problem can be solved as a *sequence* of linear programs, for which algorithms that returns exact solutions in polynomial time exist; for example, Karmarkar 1984. For the maximum score and the optimal tree, researchers may rely on existing algorithms such as the branch and bound (Wolsey and Nemhauser 1999), efficiently computed by existing software, for example, GUROBI and CPLEX. However, their worst-case scalability may grow exponentially with the sample size, similarly to what was discussed in the policy learning literature (Kitagawa and Tetenov 2018; Zhou, Athey, and Wager 2018). One solution for the optimal tree is to use an exhaustive search method (Zhou, Athey, and Wager 2018), which, we show in Section 6.1 is feasible for a moderately large sample size. For the maximum score, researchers may instead use the early termination strategy, which we study in Section 4.4. □

## 4. Theoretical Analysis

Below we restrict the policy function class of interest (Condition (A), which holds for linear scores as in Manski 1975, and decision trees as in Zhou, Athey, and Wager 2018). We also impose measurability (Condition (B); Kosorok 2008; Rai 2018).

*Assumption 4.1.* Suppose that the following conditions hold: (A) $\Pi$ has finite VC-dimension, denoted as $v$; (B) $\Pi$ is pointwise measurable.

*Assumption 4.2.* Let: (i) $e(X_i, s), p_s \in (\delta, 1 - \delta)$, almost surely, for $\delta \in (0, 1)$, for all $s \in \{0, 1\}$; (ii) $Y_i(d) \in [-M, M]$, for some $M < \infty$, for all $d \in \{0, 1\}$ almost surely.

Condition (i) imposes the standard overlap assumption; Condition (ii) assumes uniformly bounded outcomes (e.g., Mbakop and Tabord-Meehan 2021, for related conditions). The following assumptions are imposed on the estimators.

*Assumption 4.3 (Nuisances' regularities).* For some $\xi_1 \geq 1/4, \xi_2 \geq 1/4$:

$$\begin{aligned} \mathbb{E}\left[\left(\hat{m}_{d,s}(X_i) - m_{d,s}(X_i)\right)^2\right] &= \mathcal{O}(n^{-2\xi_1}), \\ \mathbb{E}\left[\left(1\big/\hat{p}_s\hat{e}(X_i, s) - 1\big/p_s e(X_i, s)\right)^2\right] &= \mathcal{O}(n^{-2\xi_2}). \end{aligned} \tag{17}$$

for all $s, d \in \{0, 1\}$, where $X_i$ is out-of-sample. In addition, for a finite constant $M$ and $\delta \in (0, 1)$, $\sup_{d \in \{0,1\}, s \in \{0,1\}, x \in \mathcal{X}} |\hat{m}_{d,s}(x)| < M$, and $\hat{e}(X, S), \hat{p} \in (\delta, 1 - \delta)$ almost surely.

Assumption 4.3 states that the *product* of the mean-squared error of the estimated propensity score and conditional mean converges to zero at the parametric rate. This condition is standard in the doubly-robust literature (Farrell 2015; Chernozhukov, Newey, and Robins 2018). Assumption 4.3 also states that the conditional mean and the propensity score functions are uniformly bounded. The conditions can be stated asymptotically, in which case, under uniform consistency, the bound on estimated nuisance functions is not required, but results should be interpreted in the asymptotic sense only (Athey and Wager 2021).

### 4.1. Guarantees on the Pareto Frontier

It is interesting to study the behavior of the estimated frontier relative to its population counterpart. We do so in the following theorems.

*Theorem 4.1.* Under Assumptions 2.1, 4.1–4.3, for any $\gamma \in (0, 1), \lambda \geq 0$, a universal constant $c_0 < \infty$, with probability larger than $1 - \gamma$,

$$\begin{aligned} \sup_{\alpha \in (0,1), \pi \in \Pi} \Big| \alpha W_0(\pi) + (1 - \alpha) W_1(\pi) \\ - \max_{\alpha_j \in \{\alpha_1, \ldots, \alpha_N\}} \Big\{ \alpha_j \widehat{W}_0(\pi) + (1 - \alpha_j) \widehat{W}_1(\pi) - \tfrac{\lambda}{\sqrt{n}} \Big\} \Big| \\ \leq c_0 \sqrt{\tfrac{v}{n}} + c_0 \sqrt{\tfrac{\log(2/\gamma)}{n}} + \tfrac{\lambda}{\sqrt{n}}. \end{aligned} \tag{18}$$

Theorem 4.1 shows that the distance between the estimated Pareto frontier and its population counterpart converges to zero

at rate $1/\sqrt{n}$ for a choice of $\lambda = \mathcal{O}(1)$ where $\lambda$ is defined in Equation (14). The derivation uses properties of the double-robust estimator (Farrell 2015) and connects to the literature on empirical welfare maximization (Kitagawa and Tetenov 2018; Zhou, Athey, and Wager 2018; Athey and Wager 2021), while differently here we control the maximum deviation uniformly over a set of weights $\alpha$.

A natural question is whether the estimated Pareto frontier also contains all Pareto optimal allocations for a finite $\lambda$. We complement Theorem 4.1 by showing that with high probability the set of estimated allocations $\widehat{\Pi}_{\circ}(\lambda)$ contains the Pareto frontier for finite $\lambda$.

*Theorem 4.2.* Let Assumptions 2.1, 4.1–4.3 hold. For any $\gamma \in (0, 1), \lambda \geq \underline{b}(\sqrt{v} + \sqrt{\log(2/\gamma)} + 1)$, for a finite constant $\underline{b} > 0$ independent of $n$, $N = \sqrt{n}$, it follows that $\mathbb{P}\left(\Pi_{\circ} \subseteq \widehat{\Pi}_{\circ}(\lambda)\right) \geq 1 - \gamma$.

Theorem 4.2 complements Theorem 4.1 showing that it suffices $\lambda = \mathcal{O}(1)$ (and hence a slackness of order $\mathcal{O}(1/\sqrt{n})$) for the set of estimated allocations to contain the Pareto frontier. The proofs of Theorems 4.1, 4.2 are contained in Appendix A, supplementary materials. Theorem 4.2 uses *finite sample* properties of the estimated (discretized) frontier showing uniform concentration. The choice of $\lambda/\sqrt{n}$ matches the upper bound on the maximal deviations, and the choice $N = \sqrt{n}$ guarantees that the grid is coarse enough to control the estimation error.

*Remark 3 (Nonbinary policies).* While our framework considers a binary action space, our guarantees also generalize to multi-action policies. In such a case, the bound depends on the entropy integral of $\Pi$ and the derivation leverages concentration of $\left|\widehat{W}_d(\pi) - W_d(\pi)\right|$ for multi-action spaces (Zhou, Athey, and Wager 2018). See Appendix B.5, supplementary materials for a discussion. $\square$

### 4.2. General Fairness Bounds

Given the guarantees on the frontier, we next analyze guarantees on fairness. We start our discussion by introducing regret bounds for generic notions of unfairness under high-level assumptions and then provide examples of upper and lower bounds.

*Assumption 4.4 (High-level conditions on UnFairness).* For some $\eta > 0, \gamma > 0$,

$$\mathbb{P}\left(\sup_{\pi \in \Pi}\left|\widehat{\mathcal{V}}_n(\pi) - \text{UnFairness}(\pi)\right| \leq \mathcal{K}(\Pi, \gamma)n^{-\eta}\right) \geq 1 - \gamma$$

for some $\mathcal{K}(\Pi, \gamma) < \infty$. Also assume that $\text{UnFairness}(\pi)$ is uniformly bounded.

Assumption 4.4 states that the estimated unfairness converges with probability $1 - \gamma$ to population unfairness uniformly over $\Pi$ at rate $n^{-\eta}$ for some arbitrary $\eta$. The constant $\mathcal{K}(\Pi, \gamma)$ depends on the function class' complexity and the probability $\gamma$. We characterize the constant and the rate $\eta$ in examples in Section 4.3 and Appendix B.4, supplementary materials.

*Theorem 4.3.* Let Assumptions 2.1, 4.1–4.4 hold. Then for some constants $0 < c_0, \underline{b} < \infty$, independent of $n$, $\lambda \geq \underline{b}(\sqrt{v} + \sqrt{\log(2/\gamma)} + 1), N = \sqrt{n}$, with probability at least $1 - 2\gamma$,

$$\text{UnFairness}(\hat{\pi}_{\lambda}) - \inf_{\pi \in \Pi_{\circ}} \text{UnFairness}(\pi) \leq \frac{c_0}{\sqrt{n}} + \frac{c_0 \mathcal{K}(\Pi, \gamma)}{n^{\eta}}. \tag{19}$$

The proof is contained in Appendix A, supplementary materials and leverages Theorem 4.2 to show that the set of Pareto allocations is contained with high probability within the estimated allocations. Theorem 4.3 characterizes the convergence rate of the UnFairness of the estimated policy relative to the lowest unfairness within the class of Pareto allocations. To our knowledge, this is the first result of this type of fair policy. The rate depends on the convergence rate of the estimated UnFairness. In the following paragraphs, we provide examples and sufficient conditions for Assumption 4.4 to hold and formally characterize the rate of convergence $\eta$ and the constant $\mathcal{K}(\cdot)$.

### 4.3. Regret: Examples and Rate Characterization

Here we discuss three examples, one based on policy predictions, a second based on the welfare effect, and a third based on incentive compatibility.

*Definition 4.1 (Prediction disparity).* Prediction disparity and its empirical counterpart take the following form

$$C(\pi) = \mathbb{E}\left[\pi(X, S)|S = 0\right] - \mathbb{E}\left[\pi(X, S)|S = 1\right],$$
$$\hat{C}(\pi) = \frac{\sum_{i=1}^n \pi(X_i)(1 - S_i)}{n(1 - \hat{p}_1)} - \frac{\sum_{i=1}^n \pi(X_i)S_i}{n\hat{p}_1}.$$

Prediction disparity captures disparity in the treatment probability between groups. The second notion of UnFairness measures welfare disparities between the two groups.

*Definition 4.2 (Welfare disparity).* Define the welfare disparity and its empirical counterpart as

$$D(\pi) = W_0(\pi) - W_1(\pi), \quad \widehat{D}(\pi) = \widehat{W}_0(\pi) - \widehat{W}_1(\pi).$$

Between-groups disparity captures the difference in *welfare* between the advantaged group ($S = 0$) and the disadvantaged group ($S = 1$), relative to the baseline.[11]

The policymaker may also consider $|D(\pi)|$ or $|C(\pi)|$ as measures of UnFairness, in which case the policymaker treats the two groups symmetrically, whose regret bounds are discussed in Appendix A.10, supplementary materials. One last example is based on the notion of incentive compatibility, motivated by discussion in Narita (2021).

*Definition 4.3 (Incentive compatibility).* Incentive compatibility is defined as

$$\mathcal{I}(\pi) = I_1(\pi) + I_0(\pi),$$
$$I_s(\pi) = \mathbb{E}\left[\pi(X, 1 - s)(Y(1) - Y(0))|S = s\right]$$
$$\qquad - \mathbb{E}\left[\pi(X, s)(Y(1) - Y(0))|S = s\right]$$

---

[11]Recall the definition of welfare in Equation (2) where we only consider the effect under treatment the effect under control.

with estimator $\widehat{\mathcal{I}}(\pi) = \hat{I}_1(\pi) + \hat{I}_0(\pi)$, $\hat{I}_s(\pi) = \frac{1}{n} \sum_{i=1}^{n} (\hat{\Gamma}_{1,s,i} - \hat{\Gamma}_{0,s,i}) \pi(X_i, 1-s) - \widehat{W}_s(\pi)$.

Here $I_s(\pi)$ captures fairness based on the incentive of an individual in revealing her sensitive attribute: $I_s(\pi)$ is positive if the welfare of an individual generated from incorrectly reporting her sensitive attribute is larger than the welfare obtained if she reported it correctly. Additional notions, such as predictive parity, can also be considered and omitted for the sake of brevity; see Appendix B.4, supplementary materials for details. For each of the three definitions above, UnFairness is linear in $\pi$, and hence optimization can be performed via MIQP.

### 4.3.1. Upper and Lower Bounds: Rate Characterization

In the following theorem, we discuss the rate of the regret-bound.

*Theorem 4.4 (Regret bound).* Let Assumptions 2.1, 4.1–4.3 hold. Let either (i) UnFairness$(\pi) = D(\pi)$, and $\widehat{\mathcal{V}}_n(\pi) = \widehat{D}(\pi)$, (ii) or UnFairness$(\pi) = C(\pi)$, and $\widehat{\mathcal{V}}_n(\pi) = \widehat{C}(\pi)$, (iii) or UnFairness$(\pi) = \mathcal{I}(\pi)$, and $\widehat{\mathcal{V}}_n(\pi) = \widehat{\mathcal{I}}(\pi)$. Then for some constants $0 < \underline{b}, c_0 < \infty$ independent of the sample size, for any $\gamma \in (0, 1), \lambda \geq \underline{b}(\sqrt{\nu} + \sqrt{\log(2/\gamma)} + 1), N = \sqrt{n}$, with probability at least $1 - 2\gamma$,

$$\text{UnFairness}(\hat{\pi}_\lambda) - \inf_{\pi \in \Pi_o} \text{UnFairness}(\pi) \leq c_0 \sqrt{\frac{\nu}{n}} + c_0 \sqrt{\frac{\log(2/\gamma)}{n}}.$$

The proof is included in Appendix A, supplementary materials. Theorem 4.4 characterizes the regret bound for three different notions of UnFairness. The bound scales at rate $1/\sqrt{n}$. Here $\mathcal{K}(\Pi, \gamma) = \sqrt{\nu} + \sqrt{\log(2/\gamma)}, \eta = 1/2$ in Theorem 4.3. The lower bound depends, however, on the notion of unfairness. Below, we derive a lower bound for any data-dependent policy which achieves the same rate for the predictive disparity.

*Theorem 4.5 (Lower bound).* Let $\Pi$ be such that $\pi(x, s)$ is constant in its last argument $s$ for all $x \in \mathcal{X}, \pi \in \Pi$, and with finite VC-dimension $\nu \geq 3$. Let UnFairness$(\pi) = C(\pi)$, and $\widehat{\mathcal{V}}_n(\pi) = \widehat{C}(\pi)$. Let $\mathcal{U}$ be the set of distributions of $(X, S)$ and $\mathcal{P}(X, S) = \{P_{Y,D|(X,S)} : \text{ such that } |Y| < M \text{ a.s., and } P(D = 1|X, S) \in (\delta, 1 - \delta)\}$. Then, there exists a distribution $P_{X,S,Y,D} = P_{X,S} P_{Y,D|X,S}$ with $P_{X,S} \in \mathcal{U}, P_{Y,D|X,S} \in \mathcal{P}(X, S)$, such that for every rule $\pi_n \in \Pi_o$ based upon $(X_1, S_1, Y_1, D_1), \ldots, (X_n, S_n, Y_n, D_n)$, for finite constants constant $0 < c_0, \bar{C} < \infty$ independent of $n$, and any $\gamma \in (0, 1/4)$, $n \geq \max\{\bar{C} \log(1/(4\gamma)), \nu - 1\}$, with probability at least $\gamma$

$$\text{UnFairness}(\pi_n) - \inf_{\pi \in \Pi_o} \text{UnFairness}(\pi) \geq \sqrt{\frac{c_0 \log(\frac{1}{4\gamma})}{n}}.$$

The proof is contained in Appendix A, supplementary materials, and, to our knowledge, it is the first result of this type for fair and Pareto optimal policies. The lower bound states that we can find a distribution and some positive (nonvanishing) probability $\gamma$ such that any data-dependent policy $\pi_n$ achieves a regret which scales to zero at a rate no faster than $1/\sqrt{n}$. Observe that a direct corollary of such a result is that the rate of the lower bound is also achieved in expectation. The condition imposes a restriction on the set of policies $\Pi$: $\Pi$ does not contain policies that use the sensitive attribute as a covariate. This class of policies occurs if anti-discriminatory laws are enforced and incorporated over the set $\Pi$. The lower bound applies to prediction disparity, and we leave to future research a more comprehensive study of lower bounds under generic notions of fairness. The derivation modifies arguments in the empirical risk minimization literature (Devroye, Györfi, and Lugosi 2013) due to the dependence of the objective function with the *conditional* probability of treatment.

Throughout this section, we have considered *distributional* notions of fairness, that is, they depend on distributional statements relative to the sensitive attribute, often used in the literature (Donini et al. 2018; Kasy and Abebe 2021; Narita 2021). *Counterfactual* notions depend instead on counterfactual statements relative to the sensitive attribute (Kilbertus et al. 2017). We discuss one counterfactual notion in Section 5.
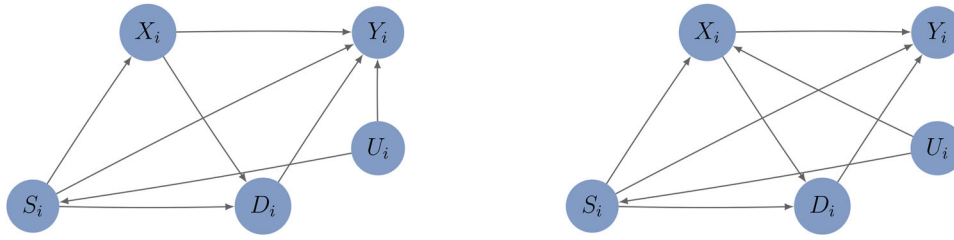
### 4.4. Computational Complexity

We conclude this section with a discussion on the computational complexity of the procedure. First, consider the case where $|\mathcal{X}| < \infty$, defines a *finite* number of strata, as often assumed in economic applications (e.g., Manski 2004). Since $|\mathcal{X}| < \infty$, let $X$ be a set of dummies $X \in \{0, 1\}^p, \sum_j X^{(j)} = 1$, and $\pi(X, S) = X\beta_S, \beta_0, \beta_1 \in [0, 1]^p$, where $\beta$ defines the treatment probability for a given individual type. Note that the result continues to hold if we require that $\pi$ assigns treatments on a finite number of strata, but $X$ is continuous.

*Proposition 4.6.* Let $|\mathcal{X}| < \infty$, and $\pi(X, S) = X\beta_S, \beta_0, \beta_1 \in [0, 1]^p$. Suppose that you can write $\mathcal{V}_n(\pi) = g(\sum_{i=1}^{n} \hat{F}_i \pi(X_i, S_i))$, for some arbitrary $\hat{F}_i$ and either $g(x) = x$ or $g(x) = |x|$. Let $||\hat{F}_i||_\infty, ||\hat{\Gamma}_i||_\infty < B < \infty$ be uniformly bounded. Then there exists an algorithm which solves Equation (16), with running time $\mathcal{O}(\sqrt{n}p^\omega)$, for some finite $\omega < \infty$.

The proof is contained in Appendix A.11, supplementary materials. Proposition 4.6 states that there exists an exact algorithm with a running time that is polynomial in the number of types and which scales at a rate $\sqrt{n}$ in the number of observations. The exponent $\omega$ depends on the algorithm.[12] The intuition is that we can represent the optimization problem in Equation (16) as a sequence of $\sqrt{n}$ many linear programs (see Appendix A.11, supplementary materials).

For generic function classes, the sequence of linear programs described above is not possible. Two examples are the maximum score with continuous variables and the optimal classification trees. These methods, however, admit a mixed-integer linear program (MILP), which can be solved *exactly* with, for example, Branch and Bound (BB) algorithms (Wolsey and Nemhauser 1999). In generic settings, MILP is known to be NP-hard in the *worst-case* scenario, hence, infeasible for large samples. Here, we characterize properties when an early termination is imposed. Namely, with an early termination, the BB algorithm reports

---

[12]Classic examples include Vaidya's and Karmarkar's algorithm (Karmarkar 1984; Vaidya 1990).

**Figure 1.** Directed acyclical graphs under which Assumption 2.1 does not hold in the presence of the confounder $U_i$, and holds in the absence of $U_i$.

an upper bound on the distance from the best objective (gap), informative for the regret.

*Proposition 4.7.* Define $\bar{W}_{j,n}^{\delta}$ the value function which maximizes $\alpha_j \hat{W}_0(\pi) + (1-\alpha_j)\hat{W}_1(\pi)$ with an early stopping criterion which stops when the estimated bound on the gap is $\delta$. Define $\hat{\pi}_{\lambda}^{\delta}$ the solution obtained after running the optimization algorithm in Equation (16) which stops whenever the estimated gap is $\delta$, and which replaces $\bar{W}_{j,n}$ in Constraint (B) with $\bar{W}_{j,n}^{\delta}$. Let the conditions in Theorem 4.4 hold. For some $0 < \underline{b}, c_0 < \infty$ independent of $n$, for any $\gamma \in (0,1)$, $\lambda \geq \underline{b}(\sqrt{v} + \sqrt{\log(2/\gamma)})$, $N = \sqrt{n}$, with probability at least $1 - 2\gamma$,

$$\text{UnFairness}(\hat{\pi}_{\lambda}^{\delta}) - \inf_{\pi \in \Pi_o} \text{UnFairness}(\pi) \leq c_0\sqrt{\frac{v}{n}} + c_0\sqrt{\frac{\log(2/\gamma)}{n}} + \delta.$$

The proof is in Appendix A.11, supplementary materials. Proposition 4.7 shows that the effect of early termination with gap $\delta$ is informative (and can be chosen appropriately) for the policy regret.

## 5. Counterfactual UnFairness

This section is of independent interest, and it discusses a novel notion of UnFairness which connects the literature on causal fairness (Kilbertus et al. 2017) and the economic literature on envy-freeness (Varian 1976). The notion is based on counterfactual statements relative to the *sensitive* attribute. We sketch the main intuition here and defer details to Appendix B.4.2, supplementary materials. This section defines $Y(d,s), X(s)$ the potential outcome and covariates as functions of the sensitive attribute $s$. The following causal model is considered.

*Assumption 5.1.* Let (A) $Y(d,s) \perp (D,S)|X(s)$, (B) $X(s) \perp S$.

Assumption 5.1 is required for estimation with counterfactual fairness and not for notions of fairness discussed in the previous sections. Condition (A) and (B) in Assumption 5.1 state that the sensitive attribute is independent of potential outcomes and covariates, while it allows for the dependence of *observed* covariates and outcomes with the sensitive attribute. Indexing potential outcomes and covariates captures this dependence by the sensitive attribute. Dependence can also occur through *unobserved* characteristics, which are dependent on both outcomes and sensitive attributes as long as observables do not *causally* affect the sensitive attribute. See Figure 1 for an

illustration. Assumption 5.1 holds when sensitive attributes do not have causal parents (e.g., Kilbertus et al. 2017).[13]

Let the conditional welfare, for the policy function being assigned to the opposite attribute, that is, the effect of $\pi(x, s_1)$, on the group $s_2$, conditional on covariates, be

$$V_{\pi(x,s_1)}(x, s_2)$$
$$= \mathbb{E}\left[\pi(x,s_1)Y_i(1,s_2) + (1 - \pi(x,s_1))Y_i(0,s_2) \Big| X_i(s_2) = x\right]. \quad (20)$$

Envy refers to the concept that "an allocation is equitable if and only if no agent prefers another agent's bundle to his own" (Varian 1976). We say that the agent with attribute $s_2$ *envies* the agent with attribute $s_1$, if her welfare (on the right-hand side of Equation (21)) exceeds the welfare she would have received had her covariate and policy been assigned the opposite attribute (left-hand side of Equation (21)), namely

$$\mathbb{E}_{X(s_1)}\left[V_{\pi(X(s_1),s_1)}\Big(X(s_1), s_2\Big)\right] > \mathbb{E}_{X(s_2)}\left[V_{\pi(X(s_2),s_2)}\Big(X(s_2), s_2\Big)\right]. \quad (21)$$

We then measure the unfairness toward an individual with attribute $s_2$ as

$$\mathcal{A}(s_1, s_2; \pi) = \mathbb{E}_{X(s_1)}\left[V_{\pi(X(s_1),s_1)}(X(s_1), s_2)\right] - \mathbb{E}_{X(s_2)}\left[V_{\pi(X(s_2),s_2)}\Big(X(s_2), s_2\Big)\right]. \quad (22)$$

Whenever we aim not to discriminate in either direction, we take the sum of the effects $\mathcal{A}(s_1, s_2; \pi)$ and $\mathcal{A}(s_2, s_1; \pi)$.[14] Equation (22) connects to previous notions of *counterfactual fairness* (Kilbertus et al. 2017), while, differently from previous references, (i) we provide formal justification to fairness using an envy-freeness argument; (ii) we construct the definition of fairness based on *distributional* impact of the treatment allocation rule on the welfare. It is complementary to Kusner et al. (2019), who compare the policy effects over individuals with the opposite sensitive attribute, lacking an envy-based justification. On the other hand, a shortcoming of the above notion is that, similarly to the above references, it does not capture notions of incentive compatibility differently from

---

[13]Whenever $S_i$ has not *causal parents*, such as, for instance, age and gender for our application (see Section 6), Assumption 2.1 holds. The case of race represents instead an exception under which Assumption 2.1 may fail since an individual's race depends on parents' characteristics. Assumption 2.1 can be stated after conditioning on baseline characteristics such as parents' observable characteristics to accommodate this latter case.

[14]Such an approach builds on the notion of "social envy" discussed in Feldman and Kirman (1974).

Definitions 4.3, discussed in the previous section. A second shortcoming is that it requires parametric estimation for the minimax rate of convergence. Namely, in Appendix B.4.2, supplementary materials we show that for a suitable choice of the estimator of $\mathcal{A}(\cdot)$, with probability at least $1 - 2\gamma$

$$\text{UnFairness}(\hat{\pi}) - \inf_{\pi \in \Pi_o} \text{UnFairness}(\pi) = \mathcal{O}\left(\sqrt{\frac{1}{n^{2\zeta}}} + \sqrt{\frac{\log(2/\gamma)}{n}}\right).$$

where here $\zeta$ denote the rate of convergence of the *conditional mean function* (see Appendix B.4.2, supplementary materials). The convergence rate is of order $n^{-1/2}$ for parametric estimators and slower for nonparametric estimators compared to the notions of UnFairness discussed in Section 4. The slower convergence rate is because counterfactual envy-freeness requires extrapolation on a different population. It opens new questions on the tradeoffs between counterfactual and predictive notions of fairness.

## 6. Empirical Application and Numerical Study

We now discuss the empirical application. This section designs a policy that assigns students to entrepreneurial programs while imposing fairness on gender. We use data that originated from Lyons and Zhang (2017). The article studies the effect of an entrepreneurship training and incubation program for undergraduate students in North America on subsequent entrepreneurial activity. We have 335 observations, of which 53% treated and the remaining under control, and 26% of applicants are women.[15] The population of interest is the pool of final applicants. We construct a targeting rule that assigns the award to the finalist based on the applicant's observable characteristics. We maximize subsequent entrepreneurial activity, which is captured using a dummy variable, indicating whether the participant worked in the startup once the program ended. The study is a quasi-experiment, and, as noted in Lyons and Zhang (2018) the focus on the pool of final applicants mitigates the selection on unobservables. Similarly to Lyons and Zhang (2018) we control for residual confounding through individual-level observable characteristics and an observable quality score of the final applicant. We estimate the nuisance functions through penalized regression, as discussed in Appendix C.1, supplementary materials.

We consider three notions of UnFairness: (i) *counterfactual envy*; (ii) *predictive disparity*, which minimizes the probability of treatment between the two groups as in Definition 4.1; (iii) *predictive disparity* with absolute value (i.e., it denotes the absolute difference between the probability of treatment between the two groups). While (ii) and (iii) do not impose conditions on the distribution of the sensitive attribute, counterfactual envy ((i)) assumes unconfoundedness also of the sensitive attribute. Such a condition is equivalent to assuming that the decision to change gender is exogenous. The reader may refer to Figure 1 for a graphical illustration. In case of failure of such an assumption, the reader should refer to results for (ii) and (iii) only.

We consider *linear* decision rules, given their large use in economics (Manski 1975)[16]

$$\Pi = \left\{\pi(x, \text{fem}) = 1\left\{\beta_0 + \beta_1 \text{fem} + x^\top \phi \geq 0\right\}, \right.$$
$$\left. (\beta_0, \beta_1, \phi) \in \mathcal{B}\right\}. \tag{23}$$

We allow covariates $x$ to be either (a) the years to graduation, years of entrepreneurship, the region of the startup, the major, the school rank, or (b) the score assigned to the candidate by the interviewer and the school rank. We refer to these two cases respectively as *Case 1* and *Case 2*. We consider in-sample capacity constraints imposed on the function class with at most 150 individuals selected for the treatment.[17]

We compare the proposed methodology to the method that maximizes the empirical welfare with the double robust score (Athey and Wager 2021). We consider three nested function classes for the welfare maximization method. The first does not impose any restriction except for the functional form in Equation (23). The second, imposes that $\beta_1 = 0$. The third class imposes that $\beta_1 = 0$ *and* that the average effect of the policy on females is at least as large as the one on males. The function classes are

$$\Pi_1 = \Pi, \quad \Pi_2 = \left\{\pi(x) = 1\left\{\beta_0 + x^\top \phi \geq 0\right\}\right\},$$
$$\Pi_3 = \left\{\pi(x) = 1\left\{\beta_0 + x^\top \phi \geq 0\right\},\right.$$
$$\mathbb{E}_n\left[(Y_i(1) - Y_i(0))\pi(X_i)\Big| S = 1\right]$$
$$\left. \geq \mathbb{E}_n\left[(Y_i(1) - Y_i(0))\pi(X_i)\Big| S = 0\right]\right\},$$

where $\mathbb{E}_n[\cdot]$ denote the empirical expectation, estimated using the doubly robust method.

Figure 2 reports the Pareto frontier over each function class.[18] The figure shows that restricting the function class leads to Pareto-dominated allocations. This outlines the limitations of maximizing welfare under fairness constraints: such constraints can be harmful to both groups. Instead, the proposed method enforces Pareto optimality in the least constrained environment (red line) and selects the policy based on fairness considerations.

In Table 1 we collect results[19] of the welfare of female and male students, as well as the relative importance weight assigned to each group for methods that maximize different UnFairness measures. In the table, we observe that minimizing Envy and Predictive Disparity leads to (weakly) larger welfare effects on the minority group. Envy leads to comparable results to welfare maximization for *Case*1 due to the discreteness of the frontier.[20] We observe an increase in the welfare of female students

---

[16]This is estimated solving Equation (16) with a small slackness parameter of order $10^{-6}$. The reader may refer to Appendix B.3, supplementary materials for details.

[17]The validity of the in-sample capacity constraints follows from a uniform concentration argument of the capacity constraint around its expectation.

[18]The value functions over the Pareto frontier can be exactly recovered as follows: we solve 2 optimization problems for each $\alpha_j, j \in \{1, \ldots, N\}$. For each of these problems, we impose constraints on the welfare of one of the two groups being larger than the other and vice-versa; we then select the subset of solutions that are not Pareto dominated by the other, and we plot the corresponding welfare in the figure.

[19]In computations, the competitors (Welfare Maximization) achieves the global optimum (dual gap equal to zero). For the proposed method, we impose a maximum time limit on the MIQP.

[20]Even if the weight $\alpha$ is larger for FTP Envy and FTP Parity Abs in *Case* 1 and 2 respectively, this does not lead to a different result than Welfare Max. 1 due to the discreteness of the frontier.
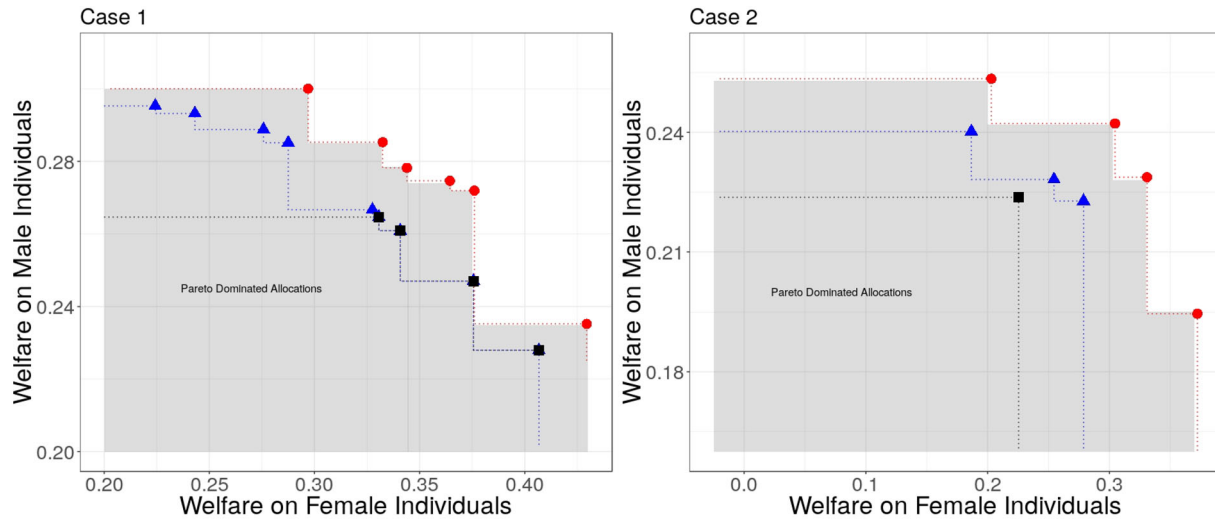
**Figure 2.** (Discretized) Pareto frontier under deterministic linear policy rule estimated through MIQP. Dots denote Pareto optimal allocations. Red dots (circle) correspond to $\Pi_1$, blue dots (triangle) to $\Pi_2$ and black dots (square) to $\Pi_3$.

**Table 1.** Empirical application.

| | Welfare female | | Welfare male | | Importance weight | |
| --- | --- | --- | --- | --- | --- | --- |
| | Case 1 | Case 2 | Case 1 | Case 2 | Case 1 | Case 2 |
| FTP Envy | 0.376 | 0.372 | 0.272 | 0.195 | 0.384 | 0.487 |
| FTP Pred | 0.432 | 0.374 | 0.224 | 0.180 | 0.847 | 0.924 |
| FTP Pred Abs | 0.433 | 0.351 | 0.208 | 0.235 | 0.924 | 0.487 |
| Welfare Max. 1 | 0.376 | 0.351 | 0.272 | 0.235 | 0.266 | 0.266 |
| Welfare Max. 2 | 0.288 | 0.307 | 0.285 | 0.238 | 0.266 | 0.266 |
| Welfare Max. 3 | 0.331 | 0.307 | 0.265 | 0.238 | 0.266 | 0.266 |

NOTE: The first two columns report the welfare improvement plus the baseline value. The last column reports the importance weights assigned by the method to the welfare of female students. FTP Envy refers to the Fair Targeting rule that minimizes envy-freeness unfairness; FTP Predictive Disp (Definition 4.1) refers to the Pareto allocation that minimizes the difference in probability of treatment (Abs indicate in absolute value); Welfare Max. 1 denotes the method that maximizes the empirical welfare considering $\Pi_1$, and similarly Welfare Max. 2, 3 for the function classes, respectively $\Pi_2, \Pi_3$.

when minimizing the *absolute* difference between probabilities of treatments for *Case* 1 and comparable results to the welfare maximization method for *Case* 2. The table shows that the proposed method finds importance weights assigned to each group solely based on the notion of fairness provided, without requiring any prior specification of relative weights assigned to each group. The method that maximizes the empirical welfare instead assigns the importance weight equal to its corresponding probability to the sensitive group, which is small for minorities. In two settings only, the results coincide with the proposed method due to the discreteness of the frontier.

Figure 3 reports the unfairness level for different sets of covariates, with unfairness measured as the difference in the probability of treatments between the two groups. Overall, Figure 3 shows that the level of the unfairness of the proposed method is uniformly smaller than the unfairness achieved by maximizing welfare, consistently with results in Section 2.

Finally, we compare also with probabilistic decision rules, which are allowed in our framework. Figure 3 collects result for a probabilistic policy function (in green) which is a super-set of $\Pi$ in Equation (23) and assigns different probabilities of treatments to groups below and above the hyperplane in Equation

(23) (see Appendix B.3, supplementary materials).[21] Results are mostly comparable across probabilistic or deterministic decisions. However, we find that a probabilistic decision enlarges the set of Pareto allocations in Appendix C.1, supplementary materials.
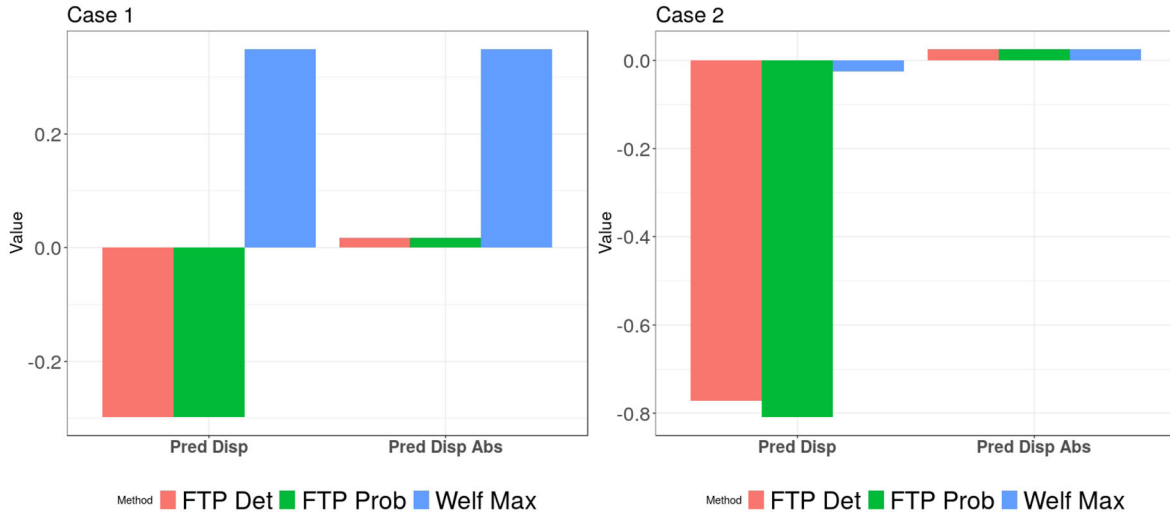
### 6.1. A Calibrated Experiment

Next, we conduct a calibrated experiment. We run the simulations calibrated to the same estimated model in the empirical application (using data from Lyons and Zhang 2017). Covariates and sensitive attributes are drawn with replacement from the empirical distribution. Formally, we draw $S_i \sim_{\text{iid}}$ Bern($\hat{p}_1$), where $\hat{p}_1$ is the probability of being female. We draw covariates $X|S = 1$ from the females' empirical distribution and similarly $X|S = 0$ for male applicants. We draw $D|X, S \sim$ Bern($\hat{e}(X, S)$), and $Y(d) = \hat{m}_{d,S}(X) + \varepsilon, \varepsilon \sim \mathcal{N}(0, 1)$, where $\hat{e}, \hat{m}$ are the estimated conditional mean and propensity score as in the application. We consider unfairness as prediction disparity in Definition 4.1.
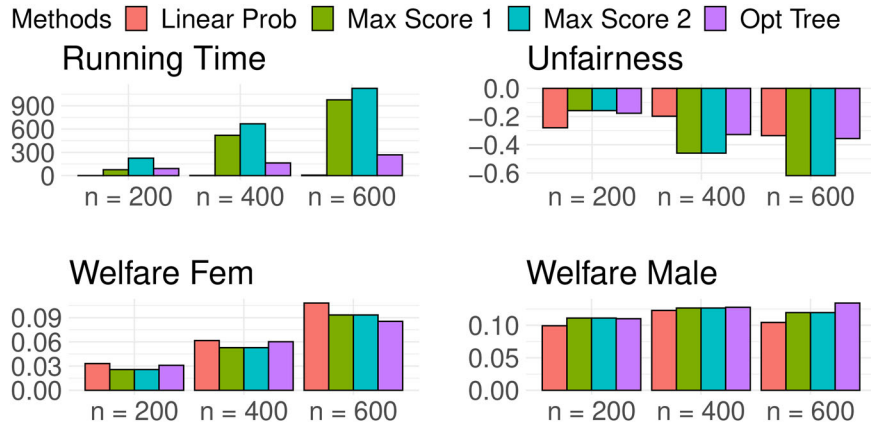
We consider three classes of policy functions: (i) probabilistic linear rule $x_1^\top \beta$, were $x_1$ is a set of binary variables ; (ii) maximum-score $\pi(x_2) = 1\{x_2^\top \beta > 0\}$; (iii) classification tree with depth equal to two. For each method, we compute results over three variables (whether the student is in a minority group, will graduate in more the one year, whether the average score exceeds the median score). We also include the average score as a continuous variable in the tree and the maximum score. We impose that the number of treated individuals does not exceed 150 individuals across each design and report welfare per share of treated individuals.[22] We estimate the probabilistic rule with a linear program, the maximum score with a mixed-integer linear program, and the classification tree via an exhaustive search.

---

[21]Formally, the function class is $\left\{ \pi_\beta(X, S) = p_1 1\{X_i^\top \beta + S\beta_0 > 0\} + p_0 1\{X^\top \beta + S\beta_0 \leq 0\}, p_1, p_0 \in [0, 1], \beta \in \mathcal{B} \right\}$.

[22]Welfare is scaled by the unconditional treatment probability since the number of treated units is fixed.

**Figure 3.** Empirical application. Unfairness level of the Fair Policy Targeting method with a deterministic allocation rule (in red), with a probabilistic decision rule (in green), and of the welfare maximization method (in blue). Pred disp refers to Definition 4.1 and Pred disp abs to Definition 4.1 in absolute value. Smaller values indicate smaller UnFairness.



**Figure 4.** Running time, UnFairness and welfare for $p = 4$. Here, Linear Prob is a linear probability rule estimated via linear programming, maximum score is estimated with MILP and optimal tree via exhaustive search with depth two. The maximum score algorithm presents two different stopping times denoted as Max Score 1 and Max Score 2 (with 50 and 200 sec stopping after estimating the Pareto frontier).

For the classification tree, we fix the number of possible splits to be four at equally spaced quantiles of each covariate distribution.[23] We run 100 replications, and over each replication, we correctly estimate the nuisance functions from the sampled observations.

In Figure 4 we report the running time (in seconds) of different function classes, with the maximum score having two different stopping times (see also Appendix C.2, supplementary materials for more results). Consistently with Section 4.4, the complexity of the linear rule scales much slower than the one of the maximum score. Also, the optimal tree is much faster than the maximum score, even for a larger sample size. The maximum score presents a relatively fast growth in terms of running time, which, however, is still feasible to handle for $n = 600$. Figure 4 also shows that a more stringent stopping time on the maximum score does not affect its performance either in terms of fairness or welfare. This is because by passing

as a starting point an "educated guess," most of the remaining optimization time is to discard dominated solutions. We obtain such a guess by taking the best solution estimated in the first step run to estimate the Pareto frontier. Finally, Figure 4 shows that different function classes mostly lead to nondominated male and female welfare comparisons.

Table 2 contrasts the unfairness and welfare of five different alternative approaches. Each competitor uses the same function class and estimation procedure as the proposed method. The first two competitors maximize a weighted average of female and male welfare with weights either $\alpha = 1/2$ (in the spirit of the planner's utility in Rambachan et al. 2020), or the empirical average $\mathbb{E}_n[S]$ as a welfare maximization problem. The third approach maximizes welfare with constraints in the form of UnFairness$_n(\pi) \leq \kappa/n$, where we choose $\kappa \in \{10, 1\}$ (Constrained Max and Constrained Max2, respectively).[24] The fourth approach maximizes welfare with constraints on disparate impact as in Definition 4.2 (in the spirit of optimization

---

[23]The choice of the exhaustive search follows in spirit to the discussion in Zhou, Athey, and Wager (2018), with differences due to the presence of multiple objectives and constraints here. Four splits facilitate computations.

[24]This is in the spirit of fairness constraints (e.g., Nabi, Malinsky, and Shpitser 2019), with constraints on statistical parity.

**Table 2.** Statistical disparity (Unfairness), welfare of male ($W_0$) and female ($W_1$) participants of the proposed method (Fair Targeting) and of the alternative procedures in *percentage points*.

| | Linear rule | | | Maximum score | | | Tree | | |
|---|---|---|---|---|---|---|---|---|---|
| | UnFair | $W_0$ | $W_1$ | UnFair | $W_0$ | $W_1$ | UnFair | $W_0$ | $W_1$ |
| Fair targeting | −33.5 | 10.4 | 10.8 | −62.4 | 12 | 9.4 | −35.5 | 13.4 | 8.5 |
| Weighted average | 13.2 | 14.7 | 8.5 | 14.3 | 16.6 | 7.7 | 2.4 | 15.2 | 8.3 |
| Utilitarian average | 25.6 | 16.3 | 6.4 | 14.3 | 16.6 | 7.7 | 30.9 | 17.0 | 6.3 |
| Constrained max | −13.2 | 14.4 | 8.1 | −32.2 | 14.2 | 8.3 | −13.8 | 15.0 | 7.1 |
| Disparate impact | 3.2 | 13.5 | 7.9 | −8.2 | 15.0 | 8.8 | 6.8 | 15.1 | 7.9 |
| Constrained Max2 | −16.3 | 14.1 | 8.3 | −34.5 | 13.9 | 8.4 | −19.3 | 14.7 | 7.3 |
| Disparate impact2 | −3.8 | 12.1 | 8.4 | −13.7 | 14.6 | 9.3 | 0.6 | 14.3 | 8.1 |

NOTE: Weighted average maximizes a weighted average of females and males' welfare with weight $\alpha = 1/2$; Utilitarian average uses instead $\alpha = \mathbb{E}_n[S]$; Constrained Max maximizes welfare under fairness constrain and Disparate Impact maximizes welfare under constrains on disparate welfare impact between the two groups. $n = 600, p = 4$. The constraint is $\kappa = 10$ for the methods in the fourth and fifth row and $\kappa = 1$ for the last two rows. The shaded row indicates the method proposed in the current paper.

in Donini et al. 2018). Interestingly, while the stricter constraint reduces the gap in males' and females' welfare for the competitor Disparate Impact, such a gap is due to the estimation error of the constraint (Appendix C.2, Figure C.4, supplementary materials presents details). We observe that the proposed method leads to the lowest UnFairness, and it is not Pareto-dominated. Our method favors the minority group, leading to larger welfare for female students. Appendix C.2, supplementary materials provides results for a smaller sample size.

## 7. Conclusion

In this article, we have introduced a novel method for estimating fair and optimal treatment allocation rules. We proposed a multi-objective decision problem, where the policymaker aims to select the least unfair policy in the set of Pareto optimal allocations. We discuss a set of theoretical guarantees on the estimated policy and provide an application. Theoretically, we open new questions on the tradeoffs between predictive and causal notions of fairness and its corresponding regret bound. Counterfactual notions require extrapolation, hence, possibly leading to a slower convergence rate. We leave a comprehensive study of the properties of different notions of fairness in terms of their implied regret to future research. From a practical perspective, an interesting new direction is estimation with nonutilitarian within-group welfare measures. Finally, the decision problem considered aims to balance efficiency and fairness, and a study of such tradeoffs in a decision theoretical framework remains an open research question.

## Supplementary Materials

The supplementary materials contain mathematical proofs omitted in the main text, extensions and additional numerical studies.

## Acknowledgments

We thank Graham Elliott, James Fowler, Ashesh Rambachan, Yixiao Sun, and Kaspar Wüthrich, the editor and anonymous referees for helpful comments. All mistakes are our own.

## Disclosure Statement

The authors report that there are no competing interests to declare.

## References

Armstrong, T., and Shen, S. (2015), "Inference on Optimal Treatment Assignments," Available at SSRN 2592479. [2]

Athey, S., and Wager, S. (2021), "Policy Learning with Observational Data," *Econometrica*, 89, 133–161. [2,3,6,7,10]

Balashankar, A., Lees, A., Welty, C., and Subramanian, L. (2019), "What is Fair? Exploring Pareto-Efficiency for Fairness Constrained Classifiers," arXiv preprint arXiv:1910.14120. [2]

Bhattacharya, D., and Dupas, P. (2012), "Inferring Welfare Maximizing Treatment Assignment under Budget Constraints," *Journal of Econometrics*, 167, 168–196. [2]

Chernozhukov, V., Newey, W. K., and Robins, J. (2018), "Double/de-biased Machine Learning using Regularized Riesz Representers," Technical report, cemmap working paper. [5,6]

Chouldechova, A. (2017), "Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments," *Big Data*, 5, 153–163. [2]

Corbett-Davies, S., and Goel, S. (2018), "The Measure and Mismeasure of Fairness: A Critical Review of Fair Machine Learning," arXiv preprint arXiv:1808.00023. [1,2]

Coston, A., Mishler, A., Kennedy, E. H., and Chouldechova, A. (2020), "Counterfactual Risk Assessments, Evaluation, and Fairness," in *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pp. 582–593. [2]

Cowgill, B., and Tucker, C. E. (2019), "Economics, Fairness and Algorithmic Bias," *Journal of Economic Perspectives* (in preparation). [1,2]

Devroye, L., Györfi, L., and Lugosi, G. (2013), *A Probabilistic Theory of Pattern Recognition* (Vol. 31), New York: Springer. [8]

Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J. S., and Pontil, M. (2018), "Empirical Risk Minimization under Fairness Constraints," in *Advances in Neural Information Processing Systems*, pp. 2791–2801. [8,13]

Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. (2012), "Fairness through Awareness," in *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pp. 214–226. [2]

Elliott, G., and Lieli, R. P. (2013), "Predicting Binary Outcomes," *Journal of Econometrics*, 174, 15–26. [2,5]

Farrell, M. H. (2015), "Robust Inference on Average Treatment Effects with Possibly More Covariates than Observations," *Journal of Econometrics*, 189, 1–23. [6,7]

Feldman, A., and Kirman, A. (1974), "Fairness and Envy," *The American Economic Review*, 64, 995–1005. [9]

Finkelstein, A., Taubman, S., Wright, B., Bernstein, M., Gruber, J., Newhouse, J. P., Allen, H., Baicker, K., and O. H. S. Group. (2012), "The Oregon Health Insurance Experiment: Evidence from the First Year," *The Quarterly Journal of Economics*, 127, 1057–1106. [4]

Florios, K., and Skouras, S. (2008), "Exact Computation of Max Weighted Score Estimators," *Journal of Econometrics*, 146, 86–91. [5]

Hardt, M., Price, E., and Srebro, N. (2016), "Equality of Opportunity in Supervised Learning," in *Advances in Neural Information Processing Systems*, pp. 3315–3323. [2]

Hirano, K., and Porter, J. R. (2009), "Asymptotics for Statistical Treatment Rules," *Econometrica*, 77, 1683–1701. [2]

Karmarkar, N. (1984), "A New Polynomial-Time Algorithm for Linear Programming," in *Proceedings of the Sixteenth Annual ACM Symposium on Theory of Computing*, pp. 302–311. [6,8]

Kasy, M., and Abebe, R. (2021), "Fairness, Equality, and Power in Algorithmic Decision-Making," in *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pp. 576–586. [2,8]

Kilbertus, N., Carulla, M. R., Parascandolo, G., Hardt, M., Janzing, D., and Schölkopf, B. (2017), "Avoiding Discrimination through Causal Reasoning," in *Advances in Neural Information Processing Systems*, pp. 656–666. [2,8,9]

Kitagawa, T., and Tetenov, A. (2018), "Who should be Treated? Empirical Welfare Maximization Methods for Treatment Choice," *Econometrica*, 86, 591–616. [2,3,5,6,7]

Kitagawa, T., and Tetenov, A. (2019), "Equality-Minded Treatment Choice," *Journal of Business & Economic Statistics*, 39, 561–574. [2]

Kleinberg, J., Ludwig, J., Mullainathan, S., and Rambachan, A. (2018), "Algorithmic Fairness," in *AEA Papers and Proceedings* (Vol. 108), pp. 22–27. [1,2]

Kosorok, M. R. (2008), *Introduction to Empirical Processes and Semiparametric Inference*, New York: Springer. [6]

Kusner, M., Russell, C., Loftus, J., and Silva, R. (2019), "Making Decisions that Reduce Discriminatory Impacts," in *International Conference on Machine Learning*, pp. 3591–3600. [2,9]

Liu, Y., Radanovic, G., Dimitrakakis, C., Mandal, D., and Parkes, D. C. (2017), "Calibrated Fairness in Bandits," arXiv preprint arXiv:1707.01875. [2]

Lyons, E., and Zhang, L. (2017), "The Impact of Entrepreneurship Programs on Minorities," *American Economic Review*, 107, 303–307. [10,11]

——— (2018), "Who does (not) Benefit from Entrepreneurship Programs?" *Strategic Management Journal*, 39, 85–112. [10]

Manski, C. F. (2004), "Statistical Treatment Rules for Heterogeneous Populations," *Econometrica*, 72, 1221–1246. [1,2,3,8]

——— (1975), "Maximum Score Estimation of the Stochastic Utility Model of Choice," *Journal of Econometrics*, 3, 205–228. [6,10]

Manski, C. F., and Thompson, T. S. (1989), "Estimation of Best Predictors of Binary Response," *Journal of Econometrics*, 40, 97–123. [5]

Martinez, N., Bertran, M., and Sapiro, G. (2019), "Fairness with Minimal Harm: A Pareto-Optimal Approach for Healthcare," arXiv preprint arXiv:1911.06935. [2]

Mas-Colell, A., Whinston, M. D., and Green, J. R. (1995), *Microeconomic Theory* (Vol. 1), New York: Oxford University Press. [3,4]

Mbakop, E., and Tabord-Meehan, M. (2021), "Model Selection for Treatment Choice: Penalized Welfare Maximization," *Econometrica*, 89, 825–848. [2,6]

Murphy, S. A. (2003), "Optimal Dynamic Treatment Regimes," *Journal of the Royal Statistical Society*, Series B, 65, 331–355. [1]

Nabi, R., Malinsky, D., and Shpitser, I. (2019), "Learning Optimal Fair Policies," *Proceedings of Machine Learning Research*, 97, 4674–4682. [1,2,4,12]

Narita, Y. (2021), "Incorporating Ethics and Welfare into Randomized Experiments," *Proceedings of the National Academy of Sciences*, 118, e2008740118. [2,7,8]

Negishi, T. (1960), "Welfare Economics and Existence of an Equilibrium for a Competitive Economy," *Metroeconomica*, 12, 92–97. [3]

Newey, W. K. (1990), "Semiparametric Efficiency Bounds," *Journal of Applied Econometrics*, 5, 99–135. [5]

Rai, Y. (2018), "Statistical Inference for Treatment Assignment Policies," Unpublished Manuscript. [6]

Rambachan, A., Kleinberg, J., Mullainathan, S., and Ludwig, J. (2020), "An Economic Approach to Regulating Algorithms," Technical report, National Bureau of Economic Research. [2,3,4,12]

Robins, J. M., and Rotnitzky, A. (1995), "Semiparametric Efficiency in Multivariate Regression Models with Missing Data," *Journal of the American Statistical Association*, 90, 122–129. [5]

Rotblat, J. (1999), "A Hippocratic Oath for Scientists," *Science*, 286, 1475–1475. [1]

Rubin, D. B. (1990), "Formal Mode of Statistical Inference for Causal Effects," *Journal of Statistical Planning and Inference*, 25, 279–292. [3]

Stoye, J. (2012), "Minimax Regret Treatment Choice with Covariates or with Limited Validity of Experiments," *Journal of Econometrics*, 166, 138–156. [2]

Sun, L. (2021), "Empirical Welfare Maximization with Constraints," arXiv preprint arXiv:2103.15298. [2]

Tetenov, A. (2012), "Statistical Treatment Choice based on Asymmetric Minimax Regret Criteria," *Journal of Econometrics*, 166, 157–165. [2]

Ustun, B., Liu, Y., and Parkes, D. (2019), "Fairness without Harm: Decoupled Classifiers with Preference Guarantees," in *International Conference on Machine Learning*, pp. 6373–6382. [2,3]

Vaidya, P. M. (1990), "An Algorithm for Linear Programming which Requires o (((m+ n) n 2+(m+ n) 1.5 n) l) Arithmetic Operations," *Mathematical Programming*, 47, 175–201. [8]

Varian, H. R. (1976), "Two Problems in the Theory of Fairness," *Journal of Public Economics*, 5, 249–260. [2,9]

Viviano, D. (2019), "Policy Targeting under Network Interference," arXiv preprint arXiv:1906.10258. [2,6]

Wolsey, L. A., and Nemhauser, G. L. (1999), *Integer and Combinatorial Optimization* (Vol. 55), New York: Wiley. [6,8]

Xiao, L., Min, Z., Yongfeng, Z., Zhaoquan, G., Yiqun, L., and Shaoping, M. (2017), "Fairness-Aware Group Recommendation with Pareto-Efficiency," in *Proceedings of the Eleventh ACM Conference on Recommender Systems*, pp. 107–115. [2]

Zhou, Z., Athey, S., and Wager, S. (2018), "Offline Multi-Action Policy Learning: Generalization and Optimization," arXiv preprint arXiv:1810.04778. [2,5,6,7,12]