# Unfairness Despite Awareness: Group-Fair Classification with Strategic Agents

Andrew Estornell[1], Sanmay Das[2], Yang Liu[3], Yevgeniy Vorobeychik[1]

[1] Washington University in Saint Louis, [2] George Mason University,
[3] University of California Santa Cruz

**Abstract**

The use of algorithmic decision making systems in domains which impact the financial, social, and political well-being of people has created a demand for these decision making systems to be "fair" under some accepted notion of equity. This demand has in turn inspired a large body of work focused on the development of fair learning algorithms which are then used in lieu of their conventional counterparts. Most analysis of such fair algorithms proceeds from the assumption that the people affected by the algorithmic decisions are represented as immutable feature vectors. However, strategic agents may possess both the ability and the incentive to manipulate this observed feature vector in order to attain a more favorable outcome. We explore the impact that strategic agent behavior could have on fair classifiers and derive conditions under which this behavior leads to fair classifiers becoming less fair than their conventional counterparts under the same measure of fairness that the fair classifier takes into account. These conditions are related to the the way in which the fair classifier remedies unfairness on the original unmanipulated data: fair classifiers which remedy unfairness by becoming more selective than their conventional counterparts are the ones that become less fair than their counterparts when agents are strategic. We further demonstrate that both the increased selectiveness of the fair classifier, and consequently the loss of fairness, arises when performing fair learning on domains in which the advantaged group is overrepresented in the region near (and on the beneficial side of) the decision boundary of conventional classifiers. Finally, we observe experimentally, using several datasets and learning methods, that this *fairness reversal* is common, and that our theoretical characterization of the fairness reversal conditions indeed holds in most such cases.

## 1 Introduction

The increasing deployment of algorithmic decision making systems in social, political, and economic domains has brought with it a demand that fairness of decisions be a central part of algorithm design. While the specific notion of

fairness appropriate to a domain is often a matter of debate, several have come to be commonly used in prior literature, such as positive (or selection) rate and false positive rate. A common goal in the design of fairness-aware (*group-fair*) algorithms is to balance predictive efficacy (such as accurate) with achieving near-equality on a chosen fairness measure among demographic categories, such as race or gender. A question that arises in many domains where such "fair" algorithms could be used is whether they are susceptible to, and create incentives for, manipulation by agents who may misrepresent themselves in order to achieve better outcomes.

Fair classifiers are often deployed in domains where feature manipulation is possible, under some constraints or potentially at some cost to agents. For example, in selection of individuals to receive assistance from social service programs (e.g homelessness services), or in admission to selective educational programs, it may be possible for applicants to misreport things, such as the number of dependents, income, or other self-reported characteristics. Importantly, we assume that agents cannot lie about group membership itself, and that the classifier cannot use group membership in making predictions (it can only do so during training). This is in keeping with most policy frameworks where fair machine learning might be used, since explicit use of protected categories in decision-making is often illegal.

We investigate the effects of such strategic manipulation of a binary *group-fair* classifier. In the context of the social services example, the classifier's job is to determine if an applicant should, or should not, be granted assistance, and the fairness guarantee of this classifier could be approximate equality of false positive rate between male and female applicants. Specifically, we aim to understand the circumstances under which a group-fair classifier becomes less fair than its conventional analog (i.e., a similar classifier that does not explicitly consider group fairness). Our main high-level observation is that if a group-fair classifier achieves greater fairness by becoming more selective than a fairness-agnostic classifier, it is also more unfair than the latter when agents are strategic. We first establish this theoretically, albeit in restricted settings, and subsequently show that this *fairness reversal* property obtains for several classifiers and datasets, and our theoretical analysis helps explain why.

**Summary of results:** We begin by examining threshold classifiers which operate on a single predictive feature (or an engineered score). Surprisingly, we show that in this setting strategic manipulation leads to a group-fair classifier being *less fair than its baseline counterpart* if and only if the group-fair classifier has a higher threshold (i.e., is more selective) than the baseline classifier. Furthermore, we provide conditions on the distribution of data for this to occur.

While the multivariate case is more complex, much of the intuition carries over. In particular, we demonstrate that for any hypothesis class from which the baseline and group-fair classifiers are sourced, there always exists a data distribution and a cost of manipulation such that fair classifiers become less fair under strategic manipulation. Additionally we show that when the fair classifier is more selective (aas formalized below), there exists a cost function

2

such that the group-fair classifier becomes less fair than the base classifier in the presence of strategic agents. Lastly, we experimentally evaluate two state-of-the-art group-fair learning algorithms on several standard datasets and demonstrate that in both the single- and multi-variable case, these algorithms often indeed produce models which are less fair than their baseline counterparts when agents act strategically.

**Selectivity and fairness:** Fairness reversal is a consequence of group-fair models becoming more selective than their conventional counterparts. In the case of threshold classifiers, higher selectivity means that the group-fair classifier has a higher selection threshold $\theta_F$ than the base (accuracy-maximizing) classifier ($\theta_C$). When does this happen? Figure 1 provides some intuition. De-
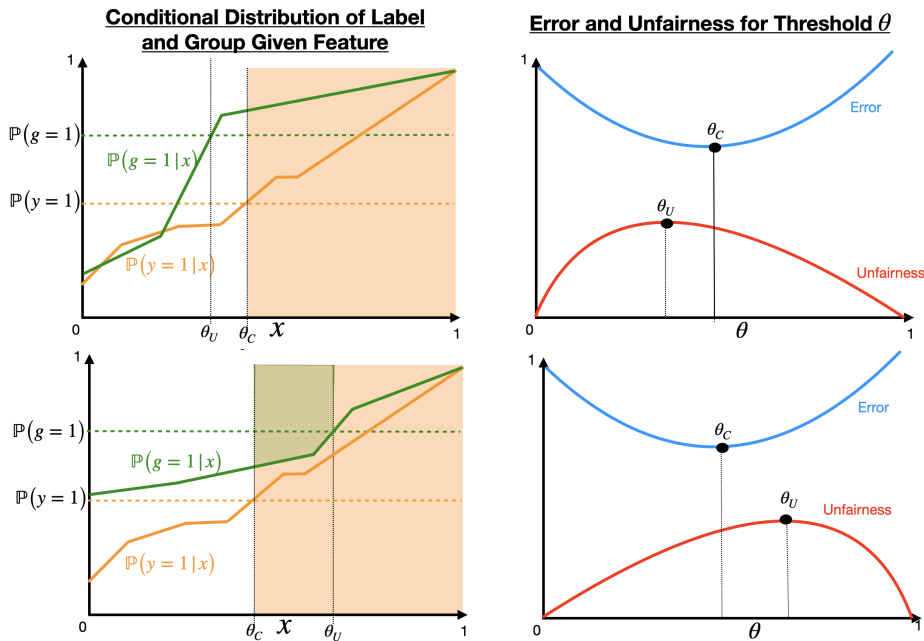


Figure 1: Stylized example of how the underlying conditional distribution leads to fair classifiers becoming more selective. The figures on the left hand correspond to the distribution, where $x$ is the predictive feature, $g$ indicates group membership, and $y$ indicates true label. The figures on the right correspond to the accuracy and unfairness of the threshold $\theta$ with respect to the distribution. Here $\theta_C$ corresponds to the threshold with the maximum accuracy, and $\theta_U$ is the threshold with the maximum unfairness.

fine $\theta_U$ as the threshold that maximizes unfairness. As discussed above, $\theta_C$ is the threshold that maximizes accuracy. Assuming that the informative feature $x$ is positively correlated with the desirable class $y = 1$, when $\theta_C > \theta_U$, the fair classifier will typically be more selective than the baseline, and the fairness reversal conditions will be met (top half of figure). This corresponds to a situa-

tion where the advantaged group is overrepresented in the region of the $x$-space just to the right of the decision boundary of $\theta_C$. Shifting this decision boundary to the right (excluding members of the advantaged group) increases fairness. However, doing so means that the newly excluded members of the advantaged group are now those who benefit most from strategic manipulation. Conversely, if $\theta_C < \theta_U$ this will likely not be the case, as pushing $\theta_C$ to the right increases unfairness. While this is a stylized example, we show in Section 6 that the selectivity condition often holds using several real-world examples.

## 2    Related Work

Our work is closely related to two major strands in literature: algorithmic fairness, and in particular, approaches for group-fair classification, and adversarial machine learning (also called strategic classification).

Broadly, algorithmic fairness literature aims to study the extent to which algorithmic decisions are perceived as unfair, for example, by being inequitable to historically disadvantaged groups [Buolamwini and Gebru, 2018, Corbett-Davies and Goel, 2018, Bolukbasi et al., 2016, Ajunwa et al., 2016]. In particular, many approaches have been introduced, particularly in machine learning, that investigate how to balance fairness and task-related efficacy, such as accurate [Agarwal et al., 2018, Feldman et al., 2015, Kearns et al., 2018, Xu et al., 2021, Zafar et al., 2019, Zemel et al., 2013, Hardt et al., 2016b]. Most of these impose hard constraints to ensure that pre-defined groups are near-equitable on some exogenously specified metric, such as selection (positive) rate [Agarwal et al., 2018, Kearns et al., 2018, Zafar et al., 2019], although alternative means, such as modifying the data to eliminate disparities, have also been proposed [Feldman et al., 2015, Chouldechova, 2017].

The adversarial machine learning literature was initially motivated by security considerations, such as spam and malware detection [Lowd and Meek, 2005, Huang et al., 2011, Vorobeychik and Kantarcioglu, 2018]. The primary issue of concern is that as we use machine learning techniques to identify malicious behavior, malicious actors change behavior characteristics to evade detection. It has, however, come to have a far broader scope, encompassing robustness of machine learning techniques in computer vision as well as social applications [Goodfellow et al., 2015, Hardt et al., 2016a, Björkegren et al., 2020, Dong et al., 2018, Chen et al., 2020]. In the latter context, this has come to be known as *strategic classification*, to indicate the concern that individuals impacted by algorithmic decisions change their features (e.g., by misreporting their household characteristics on surveys used to allocate housing to the homeless) and thereby undermine algorithms' efficacy. The intersection between strategic classification and fairness is particularly salient to our work, and has featured studies that highlight the inequity that results from strategic behavior by individuals [Hu et al., 2019], as well as inequity (social cost) resulting from making classifiers robust to strategic behavior [Milli et al., 2019, Xu et al., 2021]. Our goal, however, is quite distinct: we investigate the extent to which *group-fair* classification itself

leads to greater inequity compared to baseline approaches that do not include group-fairness constraints as a result of strategic behavior by individuals.

# 3   Preliminaries

We consider a setting with a population of agents, with each characterized by 1) a feature vector $\mathbf{x} \in \mathcal{X}$, 2) a group $g \in G \equiv \{0, 1\}$ to which it belongs (as is common in much prior literature, we treat it as binary here), and 3) a (true) binary label $y \in \mathcal{Y} \equiv \{0, 1\}$, denoting, for example, the agent's qualification (for a service, employment, bail, etc). Let $\mathcal{D}$ be the joint distribution over $G \times \mathcal{X} \times \mathcal{Y}$. We define $h(\mathbf{x})$ as the *marginal* pdf of $\mathbf{x}$, and assume that $h(\mathbf{x}) > 0$ for each $\mathbf{x} \in \mathcal{X}$.

Since using the sensitive group membership feature may pose a legal challenge, we assume that neither the conventional nor the group-fair classifier do so at prediction time (but may at training time). We denote the baseline, or conventional, classifier by $f_C$, while the group-fair classifier is denoted by $f_F$, and both map feature vectors $\mathbf{x}$ into a binary label $y \in \mathcal{Y}$. We assume that the baseline classifier aims to maximize accuracy, i.e., $f_C \in \arg\max_f \mathbb{P}_{(\mathbf{x},y)}\big(f(\mathbf{x}) = y\big)$, while $f_F$ aims to balance accuracy and fairness, solving

$$f_F = \mathrm{argmax}_f \ (1 - \alpha)\mathbb{P}_{(\mathbf{x},y)}(f(\mathbf{x}) = y)$$
$$- \alpha\big|\mathcal{M}(f_F; g = 0) - \mathcal{M}(f_F; g = 1)\big|,$$

where $\alpha \in [0, 1]$ specifies the relative weight of accuracy and fairness terms, while $\mathcal{M}(f; g)$ is a measure of efficacy (e.g., positive rate) of $f$ restricted to a group $g$.

In the literature fairness is sometimes defined with hard constraints, rather than the soft constraints of $\alpha$-fairness, for example

$$f_F = \mathrm{argmax}_f \ \mathbb{P}_{(\mathbf{x},y)}(f(\mathbf{x}) = y)$$
$$\text{subj. to } \big|\mathcal{M}(f_F; g = 0) - \mathcal{M}(f_F; g = 1)\big| \leq \beta$$

With hard constraints, decreasing $\beta$ can never increase the unfairness of $f_F$. In general soft constraints do not have the propriety that increasing $\alpha$ will never increase the unfairness of $f_F$. However, in the settings we study this is not an issue as there is a direct correspondence between $\alpha$ and $\beta$ fairness.

We consider the impact of strategic behavior of agents when they face a classifier $f$ (whether baseline or group-fair). Specifically, we suppose that each agent with features $\mathbf{x}$ can modify these, transforming them into another feature vector $\mathbf{x}'$ that is reported to the classifier. In doing so, the agent incurs a cost, captured by a cost function $c(\mathbf{x}, \mathbf{x}') \geq 0$. As in Hardt et al. [2016a], we define the agent's utility to be

$$u(\mathbf{x}, \mathbf{x}') = f(\mathbf{x}') - f(\mathbf{x}) - c(\mathbf{x}, \mathbf{x}').$$

Following the standard setting in strategic classification or adversarial machine learning, we assume any misreporting behavior would not change the true nature

of $\mathbf{x}$'s label $y$. We assume that all agents are rational utility maximizers. Thus, since $f(\mathbf{x}') - f(\mathbf{x}) \leq 1$, the agent will misreport its features only when $c(\mathbf{x}, \mathbf{x}') \leq 1$; additionally, the agent will not misreport if $f(\mathbf{x}) = 1$ (they are selected even with true values of features).

# 4 Classifiers on a Single Variable

First we consider the case in which there is a single continuous feature, i.e. $\mathcal{X} = [0, 1]$ and classifiers are thresholds on this feature. Recall that $f_C$ is the base classifier selected for maximum accuracy, and $f_F$ is an $\alpha$-fair classifier w.r.t. to fairness metric $\mathcal{M}$. We can express both classifiers as single parameter $\theta_C, \theta_F \in [0, 1]$ respectively where $f(x) = \mathbb{I}[x \geq \theta]$. We start by formalizing some standard notions we use extensively.

## 4.1 Preliminaries and preparations

**Definition 1. (Unimodal):** *A function $H : [a, b] \to \mathbb{R}$ is* negatively unimodal *(positively unimodal) on the interval $[a, b]$ if there exists a point $r \in [a, b]$ such that $H$ is monotone decreasing (increasing) on $[a, r]$ and monotone increasing (decreasing) on $[r, b]$.*
*(All convex functions are* negatively unimodal *and all concave functions are* positively unimodal*.)*

**Definition 2. (Single Crossing):** *A function $f$ is said to have a single crossing with the function $g$ if there exists $z$ s.t.*

$$\forall x \leq z : f(x) \geq g(x) \quad and \quad \forall x \geq z : f(x) \geq g(x)$$

This *single crossing* property is relevant to our work as we look at conditional distributions which have a single crossing with their unconditioned counterpart, e.g.. the functions $\mathbb{P}(y = 1|x)$ and $\mathbb{P}(y = 1)$ have a single crossing w.r.t. $x$. Experimentally we observe that this single crossing between $\mathbb{P}(y = 1|x)$ and $\mathbb{P}(y = 1)$, as well as $\mathbb{P}(g = 1|x)$ and $\mathbb{P}(g = 1)$ where $g$ indicates group membership, almost always holds in practice. Note that any monotone function has a single crossing with *every* constant function and thus monotone conditionals trivially satisfy this condition.

We begin our investigation by proving several lemmas (whose proofs are provided in the supplement) which will be useful in proving our main results, namely Theorems 1 and 2, later in the section.

**Lemma 1.** *Suppose that $\mathbb{P}(y = 1|x)$ has a single crossing with $\mathbb{P}(y = 1)$. Then error is negatively unimodal w.r.t. $\theta$ and the optimal base threshold is $\theta_C$ s.t. $\mathbb{P}(y = 1|\theta_C) = \mathbb{P}(y = 1)$.*

**Definition 3. (PR, TPR, FPR):** *Postive Rate (PR), True Positive Rate (TPR), and False Positive Rate (FPR) are defined, for classifier $f$ and distri-*

*bution $\mathcal{D}$ over $G \times \mathcal{X} \times \mathcal{Y}$, as*

$$PR_{\mathcal{D}}(f) = \mathbb{P}\big(f(x) = 1\big)$$
$$TPR_{\mathcal{D}}(f) = \mathbb{P}\big(f(x) = 1 | y = 1\big)$$
$$FPR_{\mathcal{D}}(f) = \mathbb{P}\big(f(x) = 1 | y = 0\big)$$

**Lemma 2.** *Suppose that fairness is defined in terms of Positive Rate (PR) and that $\mathbb{P}\big(g = 1 | x\big)$ has a single crossing with $\mathbb{P}\big(g = 1\big)$, then*

1. *$PR_{\mathcal{D}}(\theta | g = 1) \geq PR_{\mathcal{D}}(\theta | g = 0)$ for any $\theta \in [0, 1]$, (i.e. group 1 is advantaged under any threshold classifier), and*

2. *the unfairness term $\big| PR_{\mathcal{D}}(\theta | g = 1) - PR_{\mathcal{D}}(\theta | g = 0) \big|$ is positively unimodal w.r.t. $\theta$ and is maximized at any $\theta_U$ s.t. $\mathbb{P}\big(g = 1 | x = \theta_U\big) = \mathbb{P}\big(g = 1\big)$.*

**Lemma 3.** *Suppose that fairness is defined by either True Positive Rate or False Positive Rate and that $g, y$ are conditionally independent given $x$. Suppose further that $\mathbb{P}\big(g = 1 | x\big)$ has a single crossing with $\mathbb{P}\big(g = 1 | y = 1\big)$ in the TPR case and by $\mathbb{P}\big(g = 1 | y = 0\big)$ in the FPR case. Then when $\mathcal{M}$ is TPR or FPR,*

1. *$\mathcal{M}_{\mathcal{D}}(\theta | g = 1) \geq \mathcal{M}_{\mathcal{D}}(\theta | g = 0)$ for any $\theta \in [0, 1]$, (i.e. group 1 is advantaged under any threshold classifier), and*

2. *the unfairness term $\big| \mathcal{M}_{\mathcal{D}}(\theta | g = 1) - \mathcal{M}_{\mathcal{D}}(\theta | g = 0) \big|$ is positively unimodal w.r.t. $\theta$ and is maximized at any $\theta_U$ s.t. $\mathbb{P}\big(g = 1 | x = \theta_U\big) = \mathbb{P}\big(g = 1 | y = 1\big)$ in the TPR case and $\mathbb{P}\big(g = 1 | x = \theta_U\big) = \mathbb{P}\big(g = 1 | y = 0\big)$ in the FPR case.*

**Lemma 4.** *Suppose $f$ is of the form $f(x) = \mathbb{I}[x \geq \theta]$, and the cost of manipulating a feature from $x$ to $x'$ is given as $c(x, x')$, where an agent with true feature $x$ can submit any $x'$ subject to $c(x, x') \leq B$. Then for any cost function $c$ which is monotone in $|x' - x|$ there exists a classifier $f'(x) = \mathbb{I}[x \geq \theta']$ which makes identical predictions on the true distribution $\mathcal{D}$ as $f$ makes on manipulated data $\mathcal{D}_f^{(B)}$, i.e. when agents behave strategically $f(x') = f'(x)$ for all $x \in \mathcal{X}$. Moreover*

$$\theta' = argmin_x x \qquad s.t. \quad c(x, \theta) \leq B.$$

Lemma 4 implies that strategic agent behavior can be examined through both the perspective of the original classifier $f$ predicting on the modified distribution $\mathcal{D}_f^{(B)}$ or a modified classifier $f'$ on the original distribution $\mathcal{D}$. Since our investigation involves comparing two classifiers, $f_C$ and $f_F$, the latter perspective will prove particularly useful given that the distribution $\mathcal{D}$ remains invariant between classifiers.

## 4.2 Fair vs. baseline classifiers when agents are strategic

We now have the tools to compare the relative fairness of the fair classifier $f_F$ and the baseline classifier $f_C$ when agents are strategic. First we state a necessary and sufficient condition for which the fair classifier becomes less fair than the baseline classifiers when agents best respond to either classifier, namely that $f_F$ becomes less fair than $f_C$ if and only if $f_F$ is more selective than $f_C$ (i.e. the set of examples which $f_F$ classifies as a 1, is a subset of those classified as a 1 by $f_C$). Next we state a sufficient condition on the underlying distribution $\mathcal{D}$ that guarantees that $f_F$ will become less fair than $f_C$ when both are learned on $\mathcal{D}$. We observe that this sufficient condition is frequently satisfied in our experiments.

Note that if the budget of agents is so high that both $f_C$ and $f_F$ classify all agents as positive, then both classifiers have equal fairness and our results hold in an uninteresting manner. Thus we assume for the remainder of the paper that the available budgets do not induce such degenerate outcomes.

**Theorem 1.** *Suppose fairness is defined by PR, TPR, or FPR. Suppose further that $\mathbb{P}(y = 1|x)$ has a single crossing (Def 2) with $\mathbb{P}(y = 1)$, $\mathbb{P}(g = 1|x)$ has a single crossing with value given in Lemmas 2 and 3, $c(x, x')$ is monotone in $|x' - x|$, and $\theta_C$, $\theta_F$ are respectively the most accurate and optimal $\alpha$-fair thresholds. Then there exists a range of budgets $[B_1, B_2]$ such that strategic behavior agent behavior, with budget $B \in [B_1, B_2]$, leads to $f_F$ being less fair than $f_C$ if and only if $\theta_C < \theta_F$, (i.e. fair classifiers which are more* selective *than their baseline counterpart become less fair under strategic manipulation)*

*Proof Sketch.* The full proof is provided in the supplement. The unfairness of threshold $\theta$ w.r.t. to the distribution $\mathcal{D}$ and fairness metric $\mathcal{M} \in \{\mathrm{PR}, \mathrm{TPR}, \mathrm{FPR}\}$ is expressed as,

$$U_{\mathcal{D}}(\theta) = \left| \mathcal{M}_{\mathcal{D}}(\theta|g = 1) - \mathcal{M}_{\mathcal{D}}(\theta|g = 0) \right|,$$

For a given threshold $\theta$ and manipulation budget $B$ the best response of an agent with true feature $x$ is

$$x_\theta^{(B)} = \mathrm{argmax}_{x'} \left( \mathbb{I}[x' \geq \theta] - \mathbb{I}[x \geq \theta] \right) \ \ \mathrm{s.t.} \ \ c(x, x') \leq B,$$

When agents from $\mathcal{D}$ play this optimal response, let the resulting distribution be $\mathcal{D}_\theta^{(B)}$. The difference in unfairness, between classifiers, when agents are strategic is $U_{\mathcal{D}_{\theta_C}^{(B)}}(\theta_C)) - U_{\mathcal{D}_{\theta_F}^{(B)}}(\theta_F)$. By lemma 4 this unfairness can be expressed in terms of the true distribution $\mathcal{D}$, namely

$$U_{\mathcal{D}_{\theta_C}^{(B)}}(\theta_C) - U_{\mathcal{D}_{\theta_F}^{(B)}}(\theta_F) = U_{\mathcal{D}}(\theta_C^{(B)}) - U_{\mathcal{D}}(\theta_F^{(B)})$$

$$\text{where } \theta_C^{(B)} = \mathrm{argmin}_x x \ \mathrm{s.t.} \ c(x, \theta_C) \leq B \text{ and,}$$

$$\theta_F^{(B)} = \mathrm{argmin}_x x \ \mathrm{s.t.} \ c(x, \theta_F) \leq B$$

By the monotonicity of $c$ both $\theta_C, \theta_F$ are monotonically decreasing w.r.t. $B$ and $\theta_C^{(B)} \leq \theta_F^{(B)}$. By the unimodality of $U_{\mathcal{D}}(\theta)$, if there exists a $B'$ s.t. $\theta_F^{(B')} = \theta_U$, where $\theta_U = \operatorname{argmax}_\theta U_{\mathcal{D}}(\theta)$, then for any $B > B'$ we have $U_{\mathcal{D}}(\theta_C^{(B)}) \leq U_{\mathcal{D}}(\theta_F^{(B)})$. This must be true since the only setting in which this does not occur is $\theta_C < \theta_F < \theta_U$, which would imply that $f_C$ is at least as fair as $f_F$ on $\mathcal{D}$. Thus the forward direction holds.

Now assume that there exists a budget $B'$ such that $U_{\mathcal{D}}(\theta_C^{(B)}) \leq U_{\mathcal{D}}(\theta_F^{(B)})$. In this case, if $\theta_F < \theta_C$, then the unimodality of $U_{\mathcal{D}}(\theta)$ implies that $\theta_U < \theta_F < \theta_U$, which would imply again that $f_C$ is at least as fair as $f_F$ on the true distribution. Thus the converse direction also holds. $\qquad\square$

Theorem 1 indicates that the fair classifier becomes less fair than the base classifier (under PR, TPR, or FPR) iff $\theta_C < \theta_F$, i.e. only fair classifiers which are more *selective* than their baseline counter parts lead to greater unfairness under strategic behavior. With this observation, the question then becomes: under what conditions is $f_F$ actually more *selective* than $f_C$? We now provide a sufficient condition on the underlying distribution $\mathcal{D}$ such that the optimal $\alpha$-fair classifier is indeed more selective. In section 6 we demonstrate empirically that these sufficient conditions do in fact hold frequently in practice.

**Theorem 2.** *Suppose fairness is defined by PR, TPR, or FPR. Suppose further that $\mathbb{P}(y = 1|x)$ has a single crossing with $\mathbb{P}(y = 1)$, and that $\mathbb{P}(g = 1|x)$ has a single crossing with the respective value given in Lemmas 2 and 3, call this value $p_g$. Let $x_y$ and $x_g$ be defined by*

$$\mathbb{P}(y = 1|x_y) = \mathbb{P}(y = 1) \quad and \quad \mathbb{P}(g = 1|x_g) = p_g$$

*If $x_g < x_y$, then there exists a nonempty interval $[\alpha_0, \alpha_1]$ s.t. for any $\alpha \in [\alpha_0, \alpha_1]$ the optimal $\alpha$-fair classifier $f_F$, has the propriety that $\theta_C < \theta_F$ (implying that strategic agent behavior leads to $f_F$ becoming less fair than $f_C$ as outlined by Theorem 1).*

Intuitively this condition is saying that if the advantaged group (group 1) is overrepresented among features $x$ which have a slightly higher than base rate probability to be true positive examples (y=1), the optimal $\alpha$-fair classifier achieves its fairness by negatively classifying those "borderline" features on group 1 is overrepresented. This selectivity in turn leads to strategic agent behavior reversing the relative fairness of $f_C$ and $f_F$ since those newly rejected members of group 1 are now those who will most benefit from strategic manipulation. The size of this range of fairness coefficients will be a function of how great the over-representation of group 1 is on the features $x$. (See Figure 1)

*Proof Sketch.* The full proof is provided in the supplement. Note that $x_y = \theta_C$ and $x_g = \theta_U$, and thus $\theta_U < \theta_C$. For metric $\mathcal{M} \in \{\text{PR}, \text{TPR}, \text{FPR}\}$ the fair learning objective is

$$(1 - \alpha)\mathbb{P}(\mathbb{I}[x \geq \theta] = y) + \alpha U(\theta, \mathcal{D})$$

9

By Lemmas 1, 2, 3 the error term $\mathbb{P}\big(\mathbb{I}[x \geq \theta_F] = y\big)$ is negatively unimodal w.r.t. $\theta$ and the unfairness term $U(\theta, \mathcal{D})$ is positively unimodal. Therefore it cannot be the case that $\theta_U < \theta_F < \theta_C$, which implies that $\theta_F$ lives either on $[0, \theta_U)$ or $(\theta_C, 1]$. Further the optimal $\theta_F$ restricted to $[0, \theta_U)$ is monotonically decreasing, and the optimal $\theta_F$ restricted to $(\theta_C, 1]$ is monotonically increasing, w.r.t. $\alpha$. Thus unfairness is monotonically decreasing on both intervals w.r.t. to $\alpha$.

Since $\theta_U < \theta_C$ and both $\mathbb{P}\big(\mathbb{I}[x \geq \theta] = y\big)$ and $U(\theta, \mathcal{D})$ are smooth it must be the case that there exists a $\theta'$ s.t. any $\theta \in [\theta_C, \theta']$ has the property that

$$\mathbb{P}\big(\mathbb{I}[x \geq \theta'] = y\big) < \mathbb{P}\big(\mathbb{I}[x \geq \theta_U] = y\big) \text{ and}$$

$$U(\theta', \mathcal{D}) < U(\theta_U, \mathcal{D})$$

Thus for $\alpha \in (0, \alpha']$, where $\alpha'$ is the fairness coefficient corresponding $\theta'$, the optimal fair classifier has $\theta_C < \theta_F$. $\qquad\square$

We now turn our attention to a complementary observation. Previously we examined how strategic agent behavior can lead to a *fairness reversal* between the fair classifier $f_F$ and the baseline classifier $f_C$. This fairness reversal also comes with an interesting *accuracy reversal*. Not only is it the case that strategic behavior leads to $f_F$ taking on some of the unfairness of $f_C$, but also leads to $f_F$ obtaining some of the accuracy of $f_C$ as well. This is primarily due to the fact that $f_F$ becomes more selective and therefore more resilient to manipulation. This benefit overrides the accuracy drop on the original distribution when the budget $B$ is sufficiently large. To illustrate this concept we look at linearly separable data, i.e. there exists a $\theta^*$ such that for $x \geq \theta^*$, $y = 1$; otherwise $y = 0$.

Denote the following set:

$$\mathcal{X}(\theta, B) := \{x : x < \theta, \exists x' \geq \theta \text{ s.t. } c(x, x') \leq B\} \tag{1}$$

Note that for a fixed $\theta$, $\mathcal{X}(\theta, B)$ is monotone in $B$, i.e. a larger $B$ incurs a larger $\mathcal{X}(\theta, B)$.

**Theorem 3.** *When $\theta_F > \theta_C$ (the fair classifier is more selective), and $B$ is large enough such that*

$$\underbrace{\mathbb{P}\big(x \in \mathcal{X}(\theta_C, B)\big) + \mathbb{P}\big(x \in \mathcal{X}(\theta_F, B)\big)}_{Possible\ manipulations} \geq \underbrace{\mathbb{P}\big(x \in [\theta_C, \theta_F]\big)}_{Accuracy\ gap}$$

*Then $\theta_F$ is more accurate on $\mathcal{D}^{(B)}_{\theta_F}$ than $\theta_C$ on $\mathcal{D}^{(B)}_{\theta_C}$.*

*Proof.* The full proof is provided in the supplement. $\qquad\square$

**Fair classifiers when agents are strategic vs. when agents are truthful:**

Thus far we have considered fairness reversals between $f_F$ and $f_C$. However, it is worth noting that Theorem 1 also implies that the unfairness of $f_F$ must *strictly* increase for some range of manipulation budgets, and Theorem 2 is then also sufficient for this strict decrease of fairness to occur. Thus there always exists some range of manipulation budgets for which the fairness of the fair classifier strictly decreases in the presence of strategic agents.

# 5  Multivariate Classifiers

It is more challenging to show general results for multivariate classifiers. However, we can provide three main results that tie into both the univariate case and our subsequent empirical analyses. First we show existence: If the hypothesis class under consideration is at least as expressive as linear models, fairness reversal can occur under some distributions and cost functions. Second we show that when the fair classifier is, in a strict sense, more selective than the base classifier, there is always some cost function under which a fairness reversal occurs. Finally, to show that such cost functions are not necessarily esoteric, we provide an example of a natural cost function for linear classifiers that leads to fairness reversal.

**Theorem 4.** *Let $\mathcal{H}$ be a hypothesis class which is a super set of all linear models. If both $f_C, f_F$ are selected from $\mathcal{H}$, then there exists a distribution $\mathcal{D}$ and cost function c s.t. strategic behavior leads to a fairness reversal of $f_C$ and $f_F$.*

*Proof.* The full proof is provided in the supplement.  □

The next theorem states that if the set of examples which $f_F$ classifies as positive is partially a subset of those which $f_C$ classifies as positive, then there exists a cost function such that strategic agent behavior leads to $f_F$ being less fair than $f_C$.

**Theorem 5.** *Let $\mathcal{X}_1^{(f_F)} = \{\mathbf{x} \in \mathcal{X} : f_F(\mathbf{x}) = 1\}$ and $\mathcal{X}_1^{(f_C)} = \{\mathbf{x} \in \mathcal{X} : f_C(\mathbf{x}) = 1\}$, i.e. the set of examples predicted to be a 1 by the respective classifier when agents are truthful. Let $U_{\mathcal{D}}(f, c)$ be the difference in fairness between groups on $\mathcal{D}$ when best responding to f with cost function c. Then there exists a cost function c such that*

$$U_{\mathcal{D}}(f_F, c) - U_{\mathcal{D}}(f_C, c) > U_{\mathcal{D}}(f_C, \infty) - \mathbb{P}_{\mathbf{x} \sim \mathcal{D}}\big(\mathbf{x} \in \mathcal{X}_1^{(f_F)} \setminus \mathcal{X}_1^{(f_C)}\big)$$

*Proof Sketch.* The full proof is provided in the supplement. Strategic agent behavior causes the decisions of a classifier to change only in one direction: negative predictions become positive. Thus $f_C|_{(\mathcal{X}_1^{(f_C)} \setminus \mathcal{X}_1^{(f_F)})}$ can be constructed by a cost function c and $f_F$, yielding equal fairness between $f_C$ and $f_F$ on $\mathcal{X}_1^{(f_C)}$ (giving the positive dependence of $U_{\mathcal{D}}(f_C, \infty)$ in the bound). However, no cost

function can increase the unfairness of $f_F$ on $\mathcal{X}_1^{(f_F)} \setminus \mathcal{X}_1^{(f_C)}$ (giving the negative dependence of $\mathbb{P}\big(\mathbf{x} \in \mathcal{X}_1^{(f_F)} \setminus \mathcal{X}_1^{(f_C)}\big)$). $\qquad\square$

This result implies that the difference in unfairness between $f_F$ and $f_C$ grows as $\mathcal{X}_1^{(f_F)} \setminus \mathcal{X}_1^{(f_C)}$ shrinks. That is, as the positive predictions of $f_F$ are more likely to be subsumed by $f_C$, the difference in unfairness grows. This relationship between $f_C$ and $f_F$ holds particularly well when the cost function is *congruent* to the decision boundary of $f_C$, as is the case with monotone cost functions and threshold classifiers as well as the $l_2$-norm costs and linear classifiers.

While there may be edge cases which require this cost function to be obscure for the result to hold analytically, we demonstrate in our experimental section that the $l_2$ norm is often a sufficient cost function for this result to hold. Additionally we provide a concrete example where more selective linear classifiers exhibit this fairness reversal between $f_C$ and $f_F$ with $l_2$ manipulation cost.

**Linear classifiers:**

We now turn to linear classifiers as an illustrative example. Suppose both group membership and true label have the following relationship with $\mathbf{x}$: $\mathbb{P}\big(y = 1 | \mathbf{x}\big) = \varphi_y(\mathbf{v}_y^T \mathbf{x})$ and $\mathbb{P}\big(g = 1 | \mathbf{x}\big) = \varphi_g(\mathbf{v}_g^T \mathbf{x})$, where $\mathbf{v}_y, \mathbf{v}_g$ are fixed vectors and $\varphi_y, \varphi_g$ are a monotone functions describing a proper PDF. To model the "advantage" of group 1, let $\mathbf{v}_g^T \odot \mathbf{v}_y > 0$ (i.e. the Hamming product is elementwise nonnegative, implying that each feature $x$ is either negatively, or positively, correlated with both $y$ and $g$.). In the supplement we show this regularly arises on real data.

**Theorem 6.** *Let $f_C$ and $f_F$ be the optimal base and fair linear classifiers respectively. If $f_F$ is more selective than $f_C$, i.e. $\mathbf{w}_C \odot \mathbf{w}_F > 0$ and $\theta_F > \theta_g$, then there exists a range of manipulation budgets $[B_0, B_1]$ s.t. $c(\mathbf{x}, \mathbf{x}') = ||\mathbf{x} - \mathbf{x}'||$ lead $f_F$ to be less fair than $f_C$.*

*Proof.* The full proof is provided in the supplement. $\qquad\square$

# 6   Experiments

In this section we experimentally study the fairness reversal phenomenon that we so far considered theoretically. We use three datasets commonly found in the fairness literature,
**Recidivism:** The COMPAS dataset in which the objective is to predict re-offending.
**Law School:** Dataset of law students where the objective is to predict bar-exam passage.
**Community Crime:** Dataset of communities where the objective is to predict if the community has high crime.

All three datasets have binary outcomes, and we label the more desirable outcome for the individuals by $y = 1$ (e.g., *not* re-offending in the recidivism data), with the less desirable outcome labeled by $y = 0$. Consequently, higher

*positive rate (PR)*, , or *false positve rate (FPR)* is more desirable for individuals. Group membership in each dataset is determined by race, which in these datasets corresponds to a binary feature. In all cases, we refer to the "advantaged" group (e.g. the group with higher *PR* for *PR* based fairness) as group 1, or $G_1$, while the disadvantaged group is referred to as 0 or $G_0$.

We begin by considering single-variable threshold classifiers. For each dataset we look at thresholds over ordinal features. In both the Recidivism and Law School datasets there are 4 such features (excluding age). In the Community Crime dataset we select 4 features which have non-negligible unfairness and high correlation between the feature $x$ and label $y$. We then scale each feature to have domain $[0, 1]$ and for visual consistency across plots, if $\mathrm{Cor}(x, y) < 0$ we set $x := 1 - x$. In these three datasets most ordinal features satisfy the single crossing condition, and thus these featuers also satisfy the unimodality of error and unfairness w.r.t. to the threshold $\theta$ - this is discussed further in the supplement.

Recall that in Section 4 we showed that if error is *negatively unimodal* and unfairness is *positively unimodal* then $f_F$ becomes less fair than $f_C$ under strategic behavior if and only if $\theta_F \geq \theta_C$. We examine error and unfairness as a function of the selected features for each of the three datasets and each of the three fairness criteria. Due to space limitations, we present two of the figures here and defer the other seven, which are qualitatively similar, to the supplement. Figure 2 shows that both error and unfairness are approximately unimodal across each variable for each combination of dataset and fairness metric, with the exception of one example in which unfairness is negligible for any threshold. Moreover in these figures we see that $\theta_C < \theta_F$ holds in all but the aforementioned exception.

In the supplement, in addition to presenting graphs for other combinations of dataset and fairness metric, we also outline why the unimodality of error and unfairness arises so frequently. Specifically we show that when the conditional probabilities $\mathbb{P}(y = 1|x)$ and $\mathbb{P}(g = 1|x)$ meet the single crossing condition, as characterized in Section 4, unfairness, error, or both exhibit unimodality. Moreover, this is indeed the *typical* case in real data, with exceptions infrequent.

## 6.1 Multivariable Classifiers

Next, we study fairness reversal in settings where we make use of all ordinal features in the datasets, and consider three common baseline classifiers: logistic regression (LGR), support vector machines (SVM with an RBF kernel), and neural networks (NN). We consider two algorithms for group-fair classification: one due to Agarwal et al. [2018] (*Reductions*) and the second due to Kearns et al. [2018] (*GerryFair*). Both leverage a connection to cost-sensitive learning, but the specific techniques are different. Significant for our purposes is that both approaches turn the problem of learning with hard group-fairness constraints (that group averages on a given measure are within a specified tolerance level $\beta$) into a series of unconstrained optimization problems via an approximated Langrangian multiplier, which corresponds to $\alpha$ in our notation. Agents' manipulations are computed via *Projected Gradient Decent* (PGD) [Madry et al.,
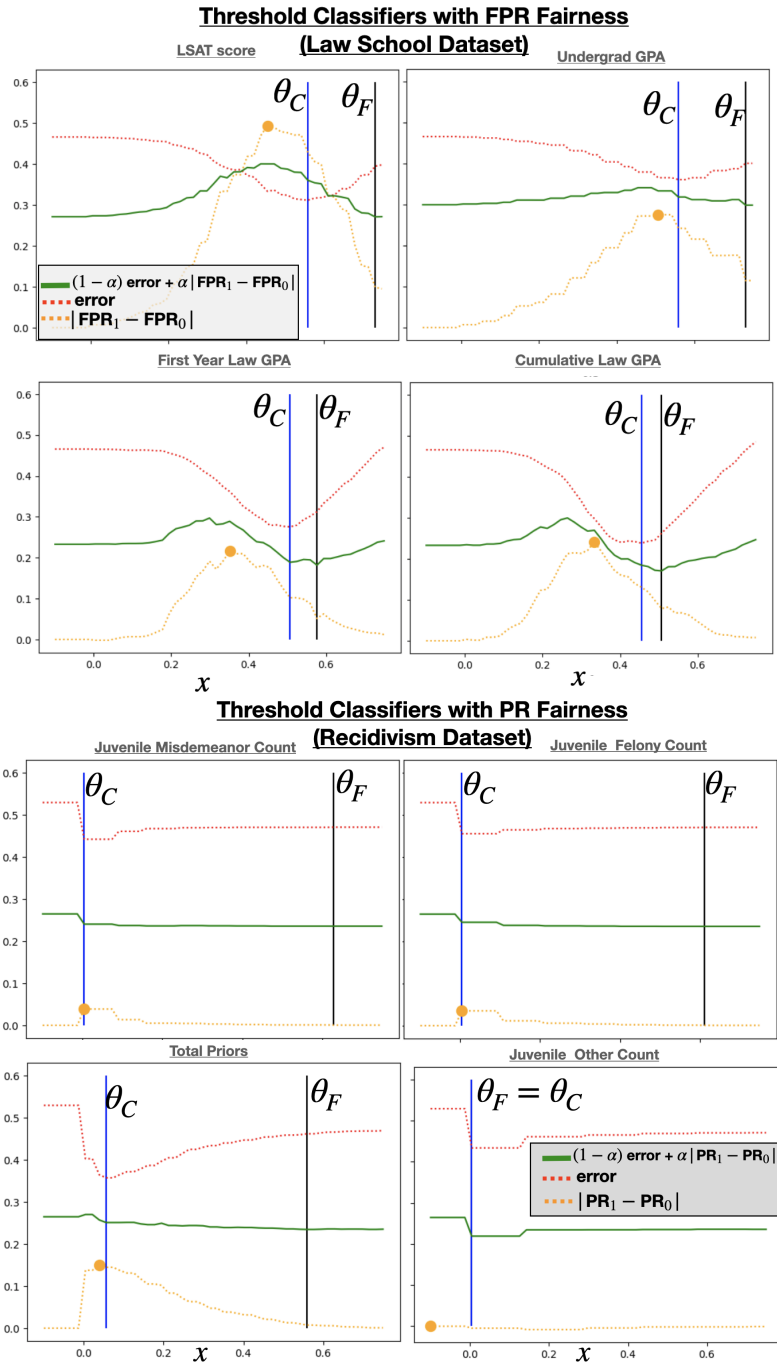
Figure 2: Optimal base and fair threshold classifiers $\theta_C, \theta_F$ respectively, on the Law School and Recidivism datasets. Recall that $\theta_F > \theta_C$ (observed in all examples, except for one - Juvenile Other Count + Recidivism) implies that strategic agent behavior will lead to $\theta_F$ becoming less fair than $\theta_C$. Error and unfairness are approximately unimodal.

14

2018]; this is discussed further in the supplement.

Figure 3 shows the unfairness and error of the fair $(f_F)$ and base $(f_C)$ classifiers, and presents three cases in which the fairness reversal occurs and one in which it does not. These are representative cases. Both error (dotted) and unfairness (solid) exhibit unimodality (w.r.t. the budget $B$, which is ultimately the parameter we care about – in the single variable case unimodality in $\theta$ implies unimodality in $B$ for any monotone cost function $c$) for both the base (blue) and fair (orange) classifiers. The shaded part indicates the region (and magnitude) of fairness reversal. In this region, the base classifier becomes more fair than than the fair classifier under strategic behavior, and there is a corresponding change where the fair classifier becomes more accurate than the base classifier, as predicted by Theorem 3. We observe experimentally that this fairness reversal is common on both the Law School and Community Crime datasets for any combination of fairness definition, base classifier, and fair learning scheme (see Supplement).

However, while the Recidivism dataset lends itself to unfairness in the single variable case, this reversal is quite infrequent in the multivariate setting, due in part to the fact that the particularly predictive features in the other two datasets are more directly correlated with group membership than those in the recidivism dataset (all though the correlation is in general high, it is lower than the other two datasets). This is discussed further in the supplement.

Figure 3: Error (dotted) and unfairness (solid) between $f_F$ (orange) and $f_C$ (blue) for several definitions of fairness, datasets, and learning schemes as a function of the manipulation budget $B$, with $l_2$-norm cost of manipulation. When evaluating fairness using PR, TPR, or FPR, we observe fairness reversal for a broad range of manipulation budget $B$.

16

# References

Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning*, pages 60–69, 2018.

Ifeoma Ajunwa, Sorelle A. Friedler, C. Scheidegger, and S. Venkatasubramanian. Hiring by algorithm: Predicting and preventing disparate impact, 2016.

Daniel Björkegren, Joshua E. Blumenstock, and Samsun Knight. Manipulation-proof machine learning. *arXiv preprint*, 2020.

Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *Advances in neural information processing systems*, 29:4349–4357, 2016.

Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *Conference on Fairness, Accountability and Transparency*, pages 77–91, 2018.

Yiling Chen, Yang Liu, and Chara Podimata. Learning strategy-aware linear classifiers. *Advances in Neural Information Processing Systems*, 33:15265–15276, 2020.

Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017.

Sam Corbett-Davies and Sharad Goel. The measure and mismeasure of fairness: A critical review of fair machine learning, 2018. arXiv preprint.

Jinshuo Dong, Aaron Roth, Zachary Schutzman, Bo Waggoner, and Zhiwei Steven Wu. Strategic classification from revealed preferences. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 55–70, 2018.

Michael Feldman, Sorelle A. Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, page 259–268, 2015.

Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.

Moritz Hardt, Nimrod Megiddo, Christos Papadimitriou, and Mary Wootters. Strategic classification. In *Innovations in Theoretical Computer Science*, 2016a.

Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29:3315–3323, 2016b.

Lily Hu, Nicole Immorlica, and Jennifer Wortman Vaughan. The disparate effects of strategic classification. In *Conference on Fairness, Accountability, and Transparency*, 2019.

Ling Huang, Anthony D Joseph, Blaine Nelson, Benjamin IP Rubinstein, and J Doug Tygar. Adversarial machine learning. In *ACM Workshop on Security and Artificial Intelligence*, pages 43–58, 2011.

Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing fairness gerrymandering: Auditing and learning for subgroup fairness. In *International Conference on Machine Learning*, pages 2564–2572, 2018.

Daniel Lowd and Christopher Meek. Adversarial learning. In *ACM SIGKDD International Conference on Knowledge Discovery in Data Mining*, pages 641–647, 2005.

Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018.

Smitha Milli, John Miller, Anca D. Dragan, and Moritz Hardt. The social cost of strategic classification. In *Conference on Fairness, Accountability, and Transparency*, page 230–239, 2019.

Yevgeniy Vorobeychik and Murat Kantarcioglu. *Adversarial machine learning*. Morgan & Claypool Publishers, 2018.

Han Xu, Xiaorui Liu, Yaxin Li, Anil K Jain, and Jiliang Tang. To be robust or to be fair: towards fairness in adversarial training. In *International Conference on Machine Learning*, 2021.

Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez-Rodriguez, and Krishna P. Gummadi. Fairness constraints: A flexible approach for fair classification. *Journal of Machine Learning Research*, 20(75):1–42, 2019.

Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International Conference on Machine Learning*, pages 325–333, 2013.
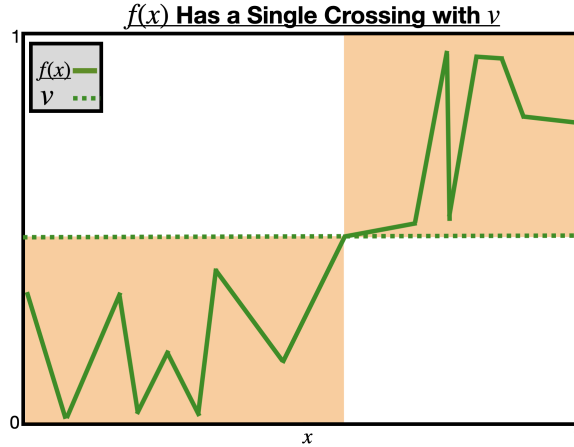
Figure 4: Example of a a function $f(x)$ (solid green) which has a *single crossing* (Def 2) with the constant function $v$ (dotted green). $f(x)$ can take on any values within the orange regions and maintaining the single crossing condition with $v$. That is, so long as $f(x)$ is upper bounded by $v$ prior to crossing $v$, and lower bounded by $v$ after crossing $v$, the single crossing condition holds.

## Supplement

### Single crossing and unimodality

**Lemma 5.** *A once-differentiable function is unimodal if its derivative has a single crossing with the constant function 0.*

*Proof.* Let $f(x) : \mathbb{R} \to \mathbb{R}$ be a once-differentiable function. Suppose that $f'(x)$ has a single crossing with 0. Then there exists a point $z$ such that $x \leq z$ implies $f'(x) \leq 0$ and $z \leq x$ implies $0 \leq f'(x)$. Thus $f$ is monotonically decreasing on the interval $(-\infty, z]$ and monotonically increasing on the interval $[z, \infty)$. Implying that $f$ is unimodal. $\square$

### Proofs

*Proof: (Lemma 1).* The error of a classifier $f(x) = \mathbb{I}[x \geq \theta]$ is given by,

$$1 - \mathbb{P}\big(\mathbb{I}[x \geq \theta] = y\big)$$
$$= 1 - \mathbb{P}\big(x \geq \theta, y = 1\big) - \mathbb{P}\big(x \leq \theta, y = 0\big)$$
$$= \mathbb{P}\big(y = 0\big) + \mathbb{P}\big(x \leq \theta, y = 1\big) - \mathbb{P}\big(x \leq \theta, y = 0\big)$$

Since $x$ is a continuous random variable and the terms involving $\theta$ are joint CDFs with well defined conditional PDFs, the derivative of the above expression w.r.t.

to $\theta$, exists and is equal to

$$h_{y,x}(y = 1, x = \theta) - h_{y,x}(y = 0, x = \theta)$$
$$= h_x(x = \theta)\big(\mathbb{P}(y = 1|x = \theta) - \mathbb{P}(y = 0|x = \theta)\big)$$
$$= h_x(x = \theta)\big(2\mathbb{P}(y = 1|x = \theta) - 1\big)$$

Since $\mathbb{P}(y = 1|x = \theta)$ is *split* by the value $1/2$, the above derivative is *split* by the value 0, thus by Lemma 5 error is *negatively unimodal* with global minima at any $\theta_C$ s.t. $\mathbb{P}(y = 1|x = \theta_C) = 1/2$. $\qquad\square$

*Proof: (Lemma 2).* For a classifier $f(x) = \mathbb{I}[x \geq \theta]$, we begin by demonstrating (1) the *unimodality* of $\mathbb{P}(x \geq \theta|g = 1) - \mathbb{P}(x \geq \theta|g = 0)$ and then use this propriety to show (2) the equivalence between $\mathbb{P}(x \geq \theta|g = 1) - \mathbb{P}(x \geq \theta|g = 0)$ and the unfairness term $\big|\mathbb{P}(x \geq \theta|g = 1) - \mathbb{P}(x \geq \theta|g = 0)\big|$. First, note that

$$= \mathbb{P}(x \geq \theta|g = 1) - \mathbb{P}(x \geq \theta|g = 1)$$
$$= \frac{\mathbb{P}(g = 1, x \geq \theta)}{\mathbb{P}(g = 1)} - \frac{\mathbb{P}(g = 0, x \geq \theta)}{\mathbb{P}(g = 0)}$$
$$= \frac{\mathbb{P}(g = 1) - \mathbb{P}(g = 1, x \leq \theta)}{\mathbb{P}(g = 1)} - \frac{\mathbb{P}(g = 0) - \mathbb{P}(g = 0, x \leq \theta)}{\mathbb{P}(g = 0)}$$
$$= 1 - \frac{\mathbb{P}(g = 1, x \leq \theta)}{\mathbb{P}(g = 1)} - 1 + \frac{\mathbb{P}(g = 0)\mathbb{P}(g = 0, x \leq \theta)}{\mathbb{P}(g = 0)}$$
$$= -\frac{\mathbb{P}(g = 1, x \leq \theta)}{\mathbb{P}(g = 1)} + \frac{\mathbb{P}(g = 0, x \leq \theta)}{\mathbb{P}(g = 0)}$$

Since each term involving $\theta$ is a joint CDF, the derivative of this term w.r.t to $\theta$ exists and is equal to

$$\frac{h_{g,x}(g = 0, x = \theta)}{\mathbb{P}(g = 0)} - \frac{h_{g,x}(g = 1, x = \theta)}{\mathbb{P}(g = 1)}$$
$$= \frac{\mathbb{P}(g = 0|x = \theta)h_x(x = \theta)}{\mathbb{P}(g = 0)} - \frac{\mathbb{P}(g = 1|x = \theta)h_x(x = \theta)}{\mathbb{P}(g = 1)}$$
$$= \frac{\big(1 - \mathbb{P}(g = 1|x = \theta)\big)h_x(x = \theta)}{\mathbb{P}(g = 0)}$$
$$\quad - \frac{\mathbb{P}(g = 1|x = \theta)h_x(x = \theta)}{\mathbb{P}(g = 1)}$$
$$= h_x(x = \theta)\frac{\mathbb{P}(g = 1) - \mathbb{P}(g = 1|x = \theta)}{\mathbb{P}(g = 1)\mathbb{P}(g = 0)}$$

Since $\mathbb{P}(g = 1|x)$ is *split* by the value $\mathbb{P}(g = 1)$ the above term is *split* by the value 0, thus by Lemma the term $\mathbb{P}(x \geq \theta|g = 1) - \mathbb{P}(x \geq \theta|g = 0)$ is *positively*

*unimodal*, and is maximized at any $\theta_U$ s.t.

$$h_x(x = \theta_U)\frac{\mathbb{P}(g = 1) - \mathbb{P}(g = 1 | x = \theta_U)}{\mathbb{P}(g = 1)\mathbb{P}(g = 0)} = 0$$

Since $h_x(x = \theta) > 0$ any such $\theta_U$ has the propriety that $\mathbb{P}(g = 1 | x = \theta_U) = \mathbb{P}(g = 1)$. Thus concluding the proof of (2).

We now use (2) to show that (1) immediately follows. Note that for $\theta \in \{0, 1\}$ we have $\mathbb{P}(x \geq \theta | g = 1) = \mathbb{P}(x \geq \theta | g = 0)$. Since the function is *positively unimodal* and $\mathbb{P}(g = 1) > 0$ neither $\theta = 0$ nor $\theta = 1$ can be points corresponding to local maximums, hence for any $\theta$ we have

$$\begin{aligned}
&\mathbb{P}(x \geq \theta | g = 1) - \mathbb{P}(x \geq \theta | g = 0) \\
&\geq \mathbb{P}(x \geq 1 | g = 1) - \mathbb{P}(x \geq 1 | g = 0) \\
&= 0
\end{aligned}$$

$\square$

*Proof: (Lemma 3.* This proof follows a similar argument to the the proof of 2. $\square$

*Proof: (Lemma 4).* When all agents prefer positive predictions to negative predictions, their manipulations will change the classifier in only a single direction, namely manipulations cause negatively predicted examples to become positively predicted. Thus, only agents with feature $x$, where $f(x) = 0$ need be considered.

Suppose $f$ is a threshold classifier with threshold $\theta$, then the agent's best response to $f$ is,

$$\begin{aligned}
x^* = \text{argmax}_x \mathbb{I}[x' \geq \theta] - \mathbb{I}[x \geq \theta] \\
\text{s.t. } c(x, x') \leq B
\end{aligned}$$

Since the cost function $c(x, x')$ is monotone w.r.t. $|x' - x|$ the above best response has solution

$$x^* = \begin{cases} \theta & \text{if } c(x, \theta) \leq B \text{ and } x < \theta \\ x & \text{otherwise} \end{cases}$$

Moreover, the monotonicity of $c(x, x')$ also implies that if an agent with feature $x$ has best response $x^* = \theta$, then so will any other agent with $x_1$ where $x \leq x_1 < \theta$.

Thus, the distribution shift of $\mathcal{D}$ caused by strategic behavior, can be quantified in terms of the agent with the smallest feature which is able to report a value of $\theta$, i.e. the feature

$$\begin{aligned}
x_{\min} = \text{argmin}_x x \\
\text{s.t. } c(x, \theta) \leq B
\end{aligned}$$

Thus when agents are strategic, any agent with feature $x \geq x_{\min}$ will be positively classified by $f$. Therefore, the threshold $\theta' = x_{\min}$ makes the same classifications on the unmanipulated distribution $\mathcal{D}$ as $\theta$ makes on the manipulated distribution $\mathcal{D}_\theta^{(B)}$. □

*Proof: (Theorem 1).* We first show that $\theta_C < \theta_F$ implies the existence of a budget interval $[B_1, B_2]$ s.t. strategic agent behavior under any $B \in [B_1, B_2]$ leads to $f_F$ being less fair than $f_C$. We then show that if $\theta_F < \theta_C$, no such budget interval exists.

The unfairness of threshold $\theta$ w.r.t. to the distribution $\mathcal{D}$ and fairness metric $\mathcal{M} \in \{\mathrm{PR}, \mathrm{TPR}, \mathrm{FPR}\}$ is expressed as,

$$U(\theta, \mathcal{D}) = \big| \mathcal{M}_\mathcal{D}(\theta | g = 1) - \mathcal{M}_\mathcal{D}(\theta | g = 0) \big|,$$

For a given threshold $\theta$ and manipulation budget $B$ the best response of an agent with true type $a = (g, x)$ is

$$x_\theta^{(B)} = \mathrm{argmax}_{x'} \big( \mathbb{I}[x' \geq \theta] - \mathbb{I}[x \geq \theta] \big) \ \text{ s.t. } \ c(x, x') \leq B,$$

When agents, originally distributed in accordance with $\mathcal{D}$, play this optimal responses w.r.t. $\theta$ and $B$, let the resulting distribution be $\mathcal{D}_\theta^{(B)}$. The difference in unfairness, between classifiers, when agents are strategic, can then be expressed as $U(\theta_C, \mathcal{D}_{\theta_C}^{(B)}) - U(\theta_F, \mathcal{D}_{\theta_F}^{(B)})$. Lemma 4 gives a way to express this difference in terms of the original distribution $\mathcal{D}$ by changing the thresholds, namely

$$U(\theta_C, \mathcal{D}_{\theta_C}^{(B)}) - U(\theta_F, \mathcal{D}_{\theta_F}^{(B)})$$
$$= U(\theta_C^{(B)}, \mathcal{D}) - U(\theta_F^{(B)}, \mathcal{D})$$

where

$$\theta_C^{(B)} = \mathrm{argmin}_x x \ \text{ s.t. } \ c(x, \theta_C) \leq B \text{ and,}$$
$$\theta_F^{(B)} = \mathrm{argmin}_x x \ \text{ s.t. } \ c(x, \theta_F) \leq B$$

By the monotonicity of $c(x, x')$ w.r.t. $x' - x$ we have that

$$\theta_C < \theta_F \implies \theta_C^{(B)} \leq \theta_F^{(B)} \quad \forall \, B \geq 0$$
$$\theta_C > \theta_F \implies \theta_C^{(B)} \geq \theta_F^{(B)} \quad \forall \, B \geq 0$$

i.e. the relative ordering of the thresholds is preserved under strategic behavior for any manipulation budget $B$, this fact will be of use later.

Let $\theta_U = \mathrm{argmax}_\theta U(\theta, \mathcal{D})$, we now proceed to prove the forward direction of our claim by three cases of the relative order of the thresholds $\theta_C, \theta_F, \theta_U$:

$$1.) \ \theta_C < \theta_F \leq \theta_U$$
$$2.) \ \theta_C \leq \theta_U \leq \theta_F$$
$$3.) \ \theta_U \leq \theta_C < \theta_F$$

First note that case (1) is infeasible. By Lemmas 2 and 3, we know that $U(\theta, \mathcal{D})$ is *positively unimodal* and maximized at $\theta_U$. Therefore, for $\theta \in [\theta_C, \theta_U]$ we have that $U(\theta, \mathcal{D})$ is monotonically increasing. Thus in case (1) we have $U(\theta_F, \mathcal{D}) \geq U(\theta_C, \mathcal{D})$, which is impossible since $f_F$ is assumed to be strictly more fair than $f_C$.

To prove that the claim holds in cases (2) and (3) we use the fact that the unfairness term $U(\theta, \mathcal{D})$, being *positively unimodal* implies that the term is also monotonically increasing on the interval $[0, \theta_U]$ and monotonically decreasing on $[\theta_U, 1]$. Hence, as stated previously, for any $\theta_1 \leq \theta_2 \leq \theta_U$ it must also be the case that $U(\theta_1, \mathcal{D}) \leq U(\theta_2, \mathcal{D}) \leq U(\theta_U, \mathcal{D})$. Therefore it suffices to show that there exists a budget interval $[B_1, B_2]$ s.t. for any $B \in [B_1, B_2]$ we have $\theta_C^{(B)} \leq \theta_F^{(B)} \theta_U$, and as stated previously, for any $B \geq 0$ $\theta_C < \theta_F$ implies that $\theta_C^{(B)} \leq \theta_F^{(B)}$. Hence we need only show that in cases (2), (3) having $B \in [B_1, B_2]$ implies $\theta_F^{(B)} \leq \theta_U$.

In both cases, (2) and (3), this follows immediately form Lemma 4, which gives the existence of a budget $B_U$, such that $\theta_F^{(B_U)} = \theta_U$, and implies that $\theta_F(B)$ is monotonically decreasing w.r.t. to $B$. Therefore, for any $B \in [B_U, \infty)$ it must be the case that $U(\theta_C^{(B)}, \mathcal{D}) \leq U(\theta_F^{(B)}, \mathcal{D})$.

To prove the reverse direction, we need to show that when $\theta_F < \theta_C$ it is the case that for any $B \geq 0$ we have $U(\theta_F^{(B)}, \mathcal{D}) \leq U(\theta_C^{(B)}, \mathcal{D})$. We again show this by three cases on the relative order of the thresholds $\theta_F, \theta_C, \theta_U$:

$$1.) \ \theta_F < \theta_C \leq \theta_U$$
$$2.) \ \theta_F \leq \theta_U \leq \theta_C$$
$$3.) \ \theta_U \leq \theta_F < \theta_C$$

Similar to the forward direction of the proof, one case is infeasible, namely case (3). This can be seen be a symmetric argument to the previous one, specifically that on the interval $[\theta_U, \theta_C]$ both error and unfairness are monotonically decreasing, and thus $\theta_F$ could not be an optimal fair threshold.

As shown previously, when $\theta_F < \theta_C$ it is also the case that for any $B \geq 0$ we have $\theta_F^{(B)} \leq \theta_C^{(B)}$, and if $\theta_C^{(B)} \geq \theta_U$ then $U(\theta_F^{(B)}, \mathcal{D}) \leq U(\theta_C^{(B)}, \mathcal{D})$. Thus the claim holds for case (1), leaving only case (2) left to prove.

In case (2) we have $\theta_F \leq \theta_U \leq \theta_C$. Let $B_U$ the budget s.t. $\theta_C(B_U) = \theta_U$, then for $B \in [0, B_U]$ the term $U(\theta_C^{(B)}, \mathcal{D})$ is monotone increasing, while $U(\theta_F^{(B)}, \mathcal{D})$ is monotone decreasing, and thus $U(\theta_F^{(B)}, \mathcal{D}) \leq U(\theta_C^{(B)}, \mathcal{D})$. Moreover for $B \in [B_U, \infty)$ we have already have show that $U(\theta_F^{(B)}, \mathcal{D}) \leq U(\theta_C^{(B)}, \mathcal{D})$.

Therefore the reverse direction of the claim holds, and thus there exists an interval $[B_1, B_2]$ s.t. $U(\theta_C^{(B)}, \mathcal{D}) \geq U(\theta_F^{(B)}, \mathcal{D})$ for $B \in [B_1, B_2]$ iff $\theta_C < \theta_F$. $\square$

*Proof: (Theorem 2).* Given $\alpha \in (0, 1)$, fairness metric $\mathcal{M} \in \{PR, TPR, FPR\}$, and data distribution $\mathcal{D}$, the objective of learning scheme is to find $\theta_F$ such that

$$\theta_F = \text{argmin}_\theta (1 - \alpha)\mathbb{P}\big(\mathbb{I}[x \geq \theta] \neq y\big) + \alpha U_\mathcal{D}(\theta) \tag{2}$$

where

$$U_{\mathcal{D}}(\theta) = \big| \mathcal{M}(\theta|g=1) - \mathcal{M}(\theta|g=0) \big|$$

By Lemma 1 the error term $\mathbb{P}\big(\mathbb{I}[x \leq \theta] = y\big)$ is negatively unimodal unimodal in $\theta$ and achieves a minimum at $\theta_C$ where $\mathbb{P}\big(y = 1|x = \theta_C\big) = \mathbb{P}(y = 1)$. Similarly, by Lemmas 2, 3 the unfairness term $U_{\mathcal{D}}(\theta)$ is positively unimodal in $\theta$ and achieves a maximum at $\theta_U$ where $\mathbb{P}\big(g = 1|x = \theta_U\big) = \mathbb{P}(g = 1)$. Thus for any $\alpha$ the fair learning objective (Equation 2) is monotonically increasing, thus $\theta_F \notin [\theta_U, \theta_C]$. Implying that either $\theta_F \in [0, \theta_U)$ or $\theta_F \in (\theta_C, 1]$. Moreover since $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ $\square$

*Proof: (Theorem 3).* Since the data is linearly separable, and the definition of $\theta_C$, we have $\theta_C = \theta^*$ which has an accuracy of 1. Due to the manipulation from $\mathcal{X}(\theta_C, B)$, the accuracy of $\theta_C$ on $\mathcal{D}'$ is then precisely $1 - \mathbb{P}\big(x \in \mathcal{X}(\theta_C, B)\big)$.

Define the following two quantities:

$$x_C := \inf_{x \in \mathcal{X}(\theta_C, B)} x, \quad x_F := \inf_{x \in \mathcal{X}(\theta_F, B)} x.$$

Consider the following two cases. First when $x_F \leq \theta_C$. Since all $x \geq \theta_C$ would either be classified by $\theta_F$ as 1 or be included in $\mathcal{X}(\theta_F, B)$ (so will report a $x'$ which will be classified as 1), the accuracy of $\theta_F$ is precisely $1 - \mathbb{P}\big(x \in [\mathbf{x}_F, \theta_C]\big)$. Since $c$ is monotonic in $|x - x'|$, and $\theta_F > \theta_C$, we have $x_F > x_C$, and therefore

$$1 - \mathbb{P}\big(x \in [x_F, \theta_C]\big) > 1 - \mathbb{P}\big(x \in [x_C, \theta_C]\big)$$

establishing $\theta_F$ is more accurate.

Now consider the case $x_F > \theta_C$. In this case, the error of $\theta_F$ is incurred by the $x \in [\theta_C, x_F]$ that has a $y = 1$ but has no incentive to deviate into $\mathcal{X}(\theta_F, B)$ and be misclassified. Therefore the accuracy of $\theta_F$ is:

$$\begin{aligned}
&1 - \mathbb{P}\big(x \in [\theta_C, x_F]\big) \\
=&1 - \big(\mathbb{P}\big(x \in [\theta_C, \theta_F]\big) - \mathbb{P}\big(x \in \mathcal{X}(\theta_F, B)\big)\big) \\
\geq&1 - \mathbb{P}\big(x \in \mathcal{X}(\theta_C, B)\big),
\end{aligned}$$

where the inequality is due to the condition we required. $\qquad\qquad\qquad$ $\square$

*Proof.* This proof follows a similar line of reasoning to its single dimensional counter part, namely that agents best responding to linear classifiers $f_C$ and $f_F$ can equivalently formulated as those same agents truthfully reporting to two modified linear classifiers $f_C$ and $f_F$. More specifically, agents sourced from $\mathcal{D}$ best responding to a classifier, say $f$ with cost function $c(\mathbf{x}, \mathbf{x}') = ||\mathbf{x} - \mathbf{x}'||$ and budget $B$, causes a shift in $D$. We call this shift $\mathcal{D}_f^{(B)}$, when agents strategically react to $f$, the classifier $f$ is now predicting over $\mathcal{D}_f^{(B)}$, rather than $\mathcal{D}$. However, there exists a modified linear classifier $f'$ whose predictions on $\mathcal{D}$ are precisely those of $f$ on $\mathcal{D}_f^{(B)}$. Thus our proof strategy is to prove that there exists a $B_0$

s.t. for any $B \geq B_0$, the corresponding modified base and fair classifier $f'_C$ and $f'_F$ have the propriety that $f'_F$ is less fair than $f'_C$ on $\mathcal{D}$.

To construct $f'_C$ and $f'_F$ we first compute the best response of agents. Under the cost function of $c(\mathbf{x}, \mathbf{x}') = ||\mathbf{x} - \mathbf{x}'||$ and a classifier linear classifier with parameters $\mathbf{w}, \theta$ the optimal response of an agent with true feature $\mathbf{x}$ is to play

$$\mathbf{x}' = \begin{cases} \mathbf{x} - \frac{\mathbf{w}^T\mathbf{x}+\theta}{||\mathbf{w}||_2^2}\mathbf{w} & \text{if } \mathbf{w}^T\mathbf{x} \leq \theta \text{ and } \left|\left|\frac{\mathbf{w}^T\mathbf{x}+\theta}{||\mathbf{w}||_2^2}\mathbf{w}\right|\right| \leq B \\ \mathbf{x} & \text{otherwise} \end{cases}$$

Using this best response and the fact that $\mathbf{w}_C, \mathbf{w}_F$ are unit vectors we can assume for the sake of analysis, that every agent plays $\mathbf{x}' = \mathbf{x} - B\mathbf{w}$ since doing so will not effect the decisions of $f$. Therefore we can write $f'_C$ and $f'_F$ as $\mathbf{w}_C, \theta_F - B$ and $\mathbf{w}_F, \theta_F - B$ respectively.

Thus, we need only show that for some $B$, thresholds $\theta_F - B$ and $\theta_C - B$ have "reversed" fairness compared to $\theta_F$ and $\theta_C$. Since $\mathbf{w}_C \odot \mathbf{w}_F > 0$, it must also be the case that $\mathbf{w}_F \odot \mathbf{v}_g > 0$ and $\mathbf{w}_C \odot \mathbf{v}_g > 0$, which in turn implies that $\mathbb{P}(g = 1|\mathbf{x})$ is also monotone in $\mathbf{w}_F^T\mathbf{x}$ and $\mathbf{w}_C^T\mathbf{x}$. As a result, if there are no agents which $f'_F$ predicts positively, but $f'_C$ predicts negativity, then $f'_C$ is at least as fair as $f'_F$. More specifically, let

$$S_{0,B} = \{\mathbf{x} \in \mathcal{X} : \mathbf{w}_F\mathbf{x} < \theta_F - B \text{ and } \mathbf{w}_C\mathbf{x} \geq \theta_F - B\},$$
$$S_{1,B} = \{\mathbf{x} \in \mathcal{X} : \mathbf{w}_F\mathbf{x} \geq \theta_F - B \text{ and } \mathbf{w}_C\mathbf{x} < \theta_F - B\}.$$

The fairness of $f'_C$ and $f'_F$ are equal on $\mathcal{X} \setminus (S_{0,B} \cup S_{1,B})$, since on this set their decisions agree. We show that if $S_{1,B'} = \emptyset$ for some $B'$, then there exists a budget $B_0 \geq B'$ s.t. for all $B \geq B_0$. The existence of $B'$ comes directly from the fact that $\theta_F > \theta_C$. Such a $B'$ must exist since $\theta_C < \theta_F$.

Intuitively this follows a similar argument to the single variable case, namely that $S_{1,B'} = \emptyset$ implies that the hyperplane induced by $f'_F$ is an upperbound for the hyperplane induced by $f'_C$. Meaning that unfairness is again unimodal w.r.t. the budget $B$.

More specifically, the PR fairness of a linear classifier can be expressed as

$$\mathbb{P}\big(f(x) = 1|g = 1\big) - \mathbb{P}\big(f(x) = 1|g = 0\big)$$
$$= \mathbb{P}\big(\mathbf{w}^T\mathbf{x} \geq \theta|g = 1\big) - \mathbb{P}\big(\mathbf{w}^T\mathbf{x} \geq \theta|g = 0\big)$$

Thus if $S_{1,B} = \emptyset$ then the difference in unfairness of $f'_F$ and $f'_C$ can be expressed

as

$$
\begin{aligned}
=&\mathbb{P}\big(\mathbf{w}_F^T\mathbf{x} \geq \theta_F - B|g=1, \mathbf{x} \in S_{0,B}\big)\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)\\
&- \mathbb{P}\big(\mathbf{w}_F^T\mathbf{x} \geq \theta_F - B|g=0, \mathbf{x} \in S_{0,B}\big)\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)\\
&- \mathbb{P}\big(\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B|g=1, \mathbf{x} \in S_{0,B}\big)\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)\\
&+ \mathbb{P}\big(\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B|g=0, \mathbf{x} \in S_{0,B}\big)\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)\\
=&\mathbb{P}\big(\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B, \mathbf{w}_F^T\mathbf{x} < \theta_F - B|g=0\big)\\
&- \mathbb{P}\big(\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B, \mathbf{w}_F^T\mathbf{x} < \theta_F - B|g=1\big)\\
=&\left(\frac{\mathbb{P}\big(g=0|\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B, \mathbf{w}_F^T\mathbf{x} < \theta_F - B\big)}{\mathbb{P}\big(g=0\big)}\right.\\
&\left.- \frac{\mathbb{P}\big(g=1|\mathbf{w}_C^T\mathbf{x} \geq \theta_C - B, \mathbf{w}_F^T\mathbf{x} < F - B\big)}{\mathbb{P}\big(g=1\big)}\right)\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)
\end{aligned}
$$

Since we only care that this quantity is non-negative, i.e. $f'_C$ is at least as fair as $f'_F$, the term $\mathbb{P}\big(\mathbf{x} \in S_{0,B}\big)$ can be dropped. The difference in conditionals can be written as

$$
\begin{aligned}
&\int_{\mathbf{x} \in S_{0,B}} \frac{1 - \varphi_g(\mathbf{v}_g^T\mathbf{x})}{\mathbb{P}\big(g=0\big)}d\mathbf{x} - \int_{\mathbf{x} \in S_{0,B}} \frac{\varphi_g(\mathbf{v}_g^T\mathbf{x})}{\mathbb{P}\big(g=1\big)}d\mathbf{x}\\
&= \int_{\mathbf{x} \in S_{0,B}} \frac{\mathbb{P}\big(g=1\big) - \varphi_g(\mathbf{v}_g^T\mathbf{x})}{\mathbb{P}\big(g=0\big)\mathbb{P}\big(g=1\big)}
\end{aligned}
$$

Since $\varphi_g$ is monotone in both $\mathbf{w}_C^T\mathbf{x}$ and $\mathbf{w}_F^T\mathbf{x}$ and decreasing $B$ can only introduce values of $\mathbf{x}$ into $S_{0,B}$ which have smaller values of $\mathbf{w}_C^T\mathbf{x}$ and remove values which have higher value of $\mathbf{w}_F^T\mathbf{x}$. Thus for $B$ such that any $\mathbf{x} \in S_{0,B}$ has $\varphi_g(\mathbf{v}_g^T\mathbf{x}) \leq \mathbb{P}(g=1)$ the integral is guaranteed to be positive, implying that for $B$ or any larger budget $f'_C$ is more fair than $f'_F$. Again due to the monotonicity of $\varphi_g$, such a sufficiently large $B$ must exist. $\qquad\square$

## Experiments

### Single Crossing

Figures 5, 6, 7 show the single crossing conditions between $\mathbb{P}\big(y=1|x\big)$, and $\mathbb{P}\big(g=1|x\big)$, and their respective constant functions given in Lemmas 1, 2, 3. We see that in all three datasets the single crossing conditions approximately holds in the sense that when the condition is violated, (i.e. crossing the respective horizontal line more than once) the violation is small in magnitude. Recall that the single crossing propriety implies the unimodality of the error and unfairness terms. Small violations (both in magnitude and duration) of the single crossing condition amount to small changes in the derivative of error or unfairness, which in term does not consequentially impact the unimodality of either term from an empirical perspective.

With this said we see that the Recidivism dataset breaks the single crossing assumptions on $\mathbb{P}(g = 1|x)$ more so than the other two datasets. However, we still observe both the unimodality of unfairness as well as a reversal of fairness between the base and fair classifier on this dataset (in the single variable case).

Moreover, we see that the variables in the Recidivism dataset have weaker relationships between $\mathbb{P}(y = 1|x)$ and $\mathbb{P}(g = 1|x)$ compared to those of the Law School and Community Crime datasets. Meaning that in the Recidivism dataset, label identification and group identification are easier to untangle, compared to the other two datasets. This is partially the reason we observe that the fairness reversal does not occur on the recidivism dataset in the multivariate setting.

**Fair learning schemes**

We make use of two fair learning algorithms to generate the fair models (denoted as $f_F$), namely GerryFair and Reductions. Each algorithm takes as input a base-learner (not to be confused with the baseline classifier which we denote as $f_C$). This base-learner is used solve the fair learning objective through cost sensitive learning. Each algorithm uses their respective base learner in a unique way, and the fair models produced by each learning scheme different considerably in terms of their structure. Reductions uses the base-learner to perform traditional cost sensitive learning and outputs a model which is of the same type of the base-learner. For example if the base-learner is Logistic Regression, then $f_F$ is also a Logistic Regression model. Thus the Projected Gradient Decent attack (PDG) is effective at computing an agents best response to $f_F$ when $f_F$ is learned via reductions and a differentiable base-learner (e.g. Logistic Regression, SVM, and Neural Networks). In the case of GerryFair the returned fair model $f_F$ has a different structure from the base learner, namely $f_F$ is an ensemble of models produced from the base learner. Thus the resulting model may not be smooth and PGD will not work to compute agents best response. However, of the models we examine (LRG, SVM, NN) GerryFair only supports LRG and SVM with a linear kernel. Meaning that each learner in the ensemble, produced by GerryFair, is linear and thus it is trivial to compute each agent's best response.

**Fairness Reversal**

Recall that in the single variable case, strategic manipulation leads to a fairness reversal between the base and fair thresholds $\theta_C$ and $\theta_F$ respectively, if and only if $\theta_C < \theta_F$. Figures 5-16 show the relationship between $\theta_C$ and $\theta_F$ for each of the variables, dataset, and fairness metrics we study. In these figures we see that $\theta_C < \theta_F$ is a common. Moreover, we see that the cases where this relationship does not hold are cases in which either $\theta_C < \theta_U$ (meaning the sufficient condition of Theorem 2 does not hold), the fair classifier is trivial (i.e. $\theta_F = 0$), or there is negligible unfairness regardless of the value selected for $\theta$. Moreover, we see that both error and unfairness are unimodal w.r.t. $\theta$, thus Lemma 4 implies that error and unfairness will remain unimodal w.r.t. the manipulation budget
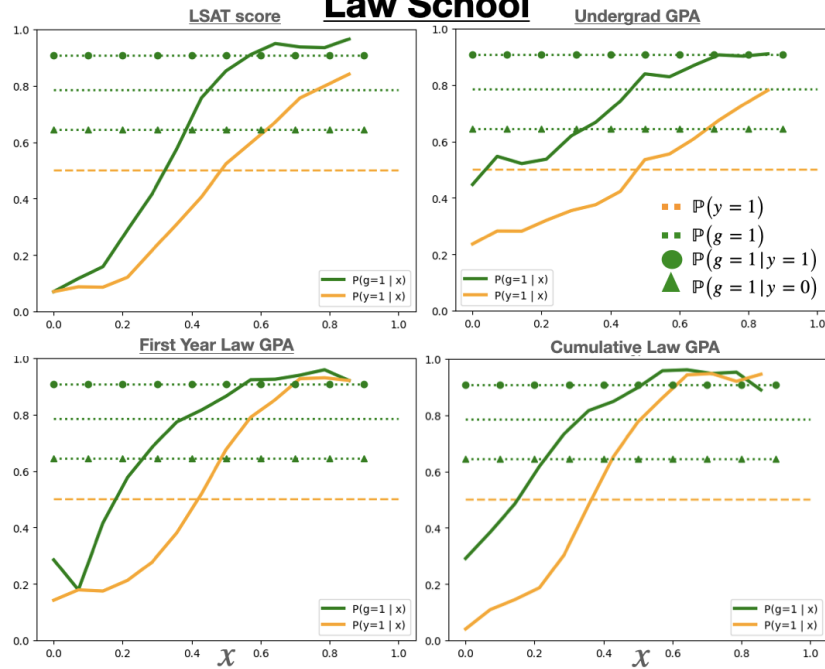
Figure 5: Probabilities of group membership $g$ (green) and true label $y$ (orange). Probabilities conditioned on the feature $x$ are given as solid lines, while those unconditioned are given as dotted, or dashed, lines. Recall that if the conditioned probabilities $\mathbb{P}(g = 1|x)$ and $\mathbb{P}(y = 1|x)$ having a single crossing with the respective unconditioned value (outlined in Lemmas 1, 2, 3) then error and unfairness will be unimodal w.r.t. to the threshold $\theta$. For example, in the case of PRfairness, if $\mathbb{P}(g = 1|x)$ has a single crossing with $\mathbb{P}(g = 1)$ and $\mathbb{P}(y = 1|x)$ has a single crossing with $\mathbb{P}(y = 1)$ then error and unfairness are unimodal w.r.t. to $\theta$.

# Conditional PDF of Group and Label Given Feature Recidivism



Figure 6:

Figure 7:

$B$ for *any* manipulation cost function $c(x, x')$ which is monotone in $|x' - x|$.

With respect to Figures 5-16, agent manipulation amounts to "moving" each threshold to the left. We can see that when $\theta_C < \theta_F$, moving $\theta_C$ to the left decreases unfairness, while moving $\theta_F$ to left increases unfairness, until the manipulated $\theta_F$ has been moved all the way to $\theta_U$ (the most unfair threshold). Additionally in these figures we see that not only does this leftward shift increase the unfairness of $\theta_F$, but also increased the accuracy of $\theta_F$: a phenomenon outlined by Theorem 3. That is, in the cases where $\theta_C < \theta_F$, strategic manipulation leads to both a fairness, and an accuracy, reversal between $\theta_C$ and $\theta_F$.

In the multivariate case, Figures 17-24, show that again the fairness reversal is common. In these figures the error (dotted) and unfairness (solid) are given for $f_C$ (blue) and $f_F$ (orange) as a function of the manipulation budget $B$ when agents have manipulation cost $c(\mathbf{x}, \mathbf{x}') = ||\mathbf{x} - \mathbf{x}'||$. The shaded orange region indicates both the duration (in terms of $B$) and magnitude of the fairness reversal of $f_F$ and $f_C$. Similar to the single variable case, these graphs again show a fundamental trade-off between fairness and accuracy in the presence of manipulation. In all settings where a fairness reversal between $f_F$ and $f_C$ occurs, an accuracy reversal also occurs. Namely, if $f_F$ becomes less fair than $f_C$, it also becomes more accurate than $f_C$. Moreover, as was the case in the single variable case, we see that in the multivariate case both error and unfairness exhibit unimodality w.r.t. to the budget $B$.

In the single variable case, we would expect that once $f_C$ and $f_F$ respectively hit the point with maximum unfairness (as a function of $B$) their unfairness would decrease at an equal rate from that point onward since both classifiers are effectively sharing the same unfairness curve, but sit at different points. In the multivariate case, we make this same observation. After reaching the most unfair $B$, both classifiers decreases at similar rates. However, $f_C$ requires a larger $B$, than $f_C$, to reach this point. Which ultimately leads to $f_F$ becoming less fair, since the unfairness of $f_F$ is still increasing while the unfairness of $f_C$ has already begun to fall.

The fairness reversal is common on the Law School and Community Crime datasets, but non-existent on the recidivism dataset. As mentioned previously there is a weaker correlation between group membership and true label on the Recidivism dataset compared to the other two datasets. We suspect that this weaker link between accuracy and unfairness results in fair classifiers which lose less of their fairness in the presence of strategic behavior, since the agents benefiting from strategic behavior should be more representative of the population as a whole, rather than one particular group.

We train a total of 9 linear models in our experiments (3 fairness definitions across 3 datasets using LRG and Reductions). We observe that 5/9 cases the propriety that $\mathbf{w}_C \odot \mathbf{w}_F > 0$ approximately holds. Here approximately holding means that if an element pair $w_{i,C} w_{i,F} < 0$ the respective magnitude of the product is small relative to those of the larger value weights.
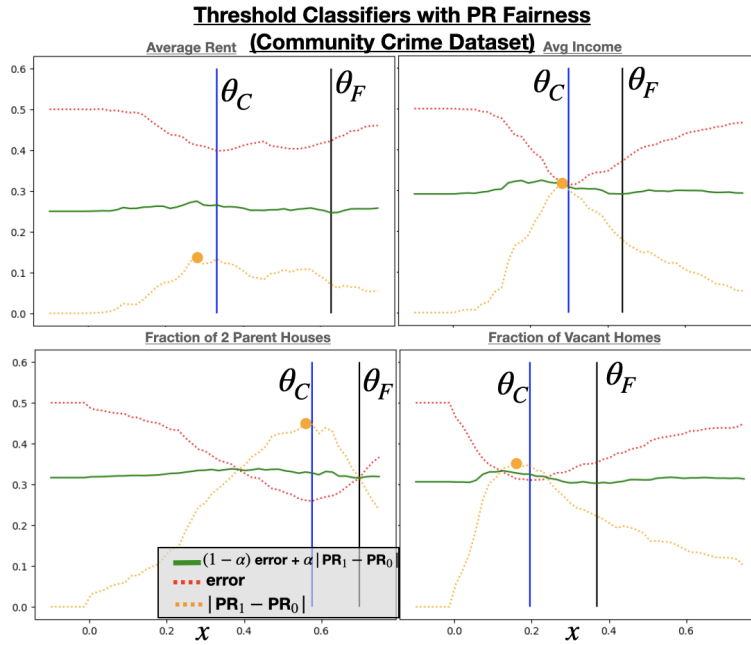
Figure 8: Unfairness and error of threshold classifiers. Both error and unfairness are approximately unimodal w.r.t. threshold $\theta = x$. Thus error and unfairness are also unimodal w.r.t. the manipulation budget $B$ for any manipulation cost function $c(x, x')$ which is monotone in $|x' - x|$. When this unimodality holds $\theta_C < \theta_F$ implies that strategic manipulation will lead to $\theta_C$ becoming more fair than $\theta_F$. This fairness reversal is due to the fact that strategic manipulation amounts to lowering (shifting to the left) the threshold. In this figure, as well as the subsequent figures, we see that $\theta_C < \theta_F$ is a common occurrence (namely 30 our of the 36 combinations of variable, fairness metric, and dataset studied).

Figure 9:



Figure 10:

**Threshold Classifiers with FPR Fairness**
**(Law School Dataset)**

Figure 11:

**Threshold Classifiers with FPR Fairness**
**(Community Crime Dataset)**

Figure 12:

34

Figure 13:



Figure 14:

**Threshold Classifiers with TPR Fairness (Recidivism Dataset)**



Figure 15:

**Threshold Classifiers with TPR Fairness (Community Crime Dataset)**



Figure 16:

36

Figure 17: Difference in fairness for Logistic Regression (blue) and Reduction (orange) when agents are strategic under $l_2$ cost with budget $B$. The shaded region indicates instances in which the fair classifier is less fair than its baseline counterpart. Note that when the fairness of $f_F$ and $f_C$ are reversed, there accuracy is reversed as well. Once hitting their maximally unfair point both $f_C$ and $f_F$ tend to decreased in unfairness at equal rates. However, $f_F$ typically hits this maximally unfair point for a $B$ larger than $f_C$.

Figure 18:

**Difference in FPR, TPR, PR for $f_C$ and $f_F$ on Law School (NN)**
**(Reductions)**



Figure 19:

**Difference in FPR and TPR for $f_C$ and $f_F$ on Law-School (LRG) (GerryFair)**
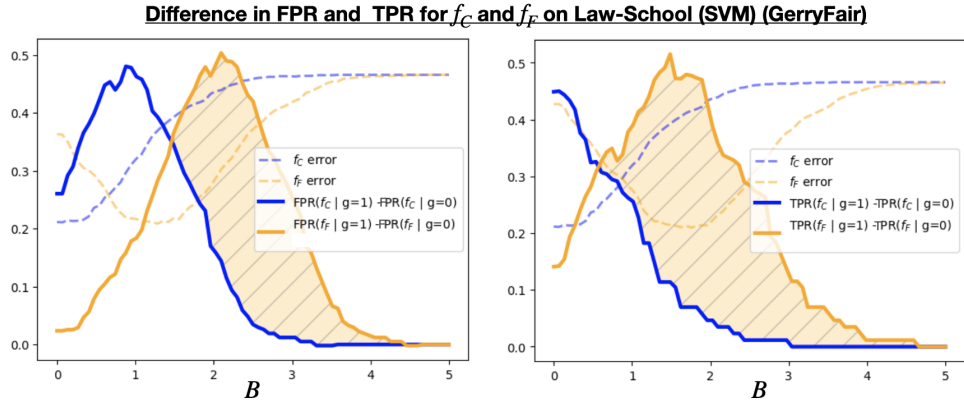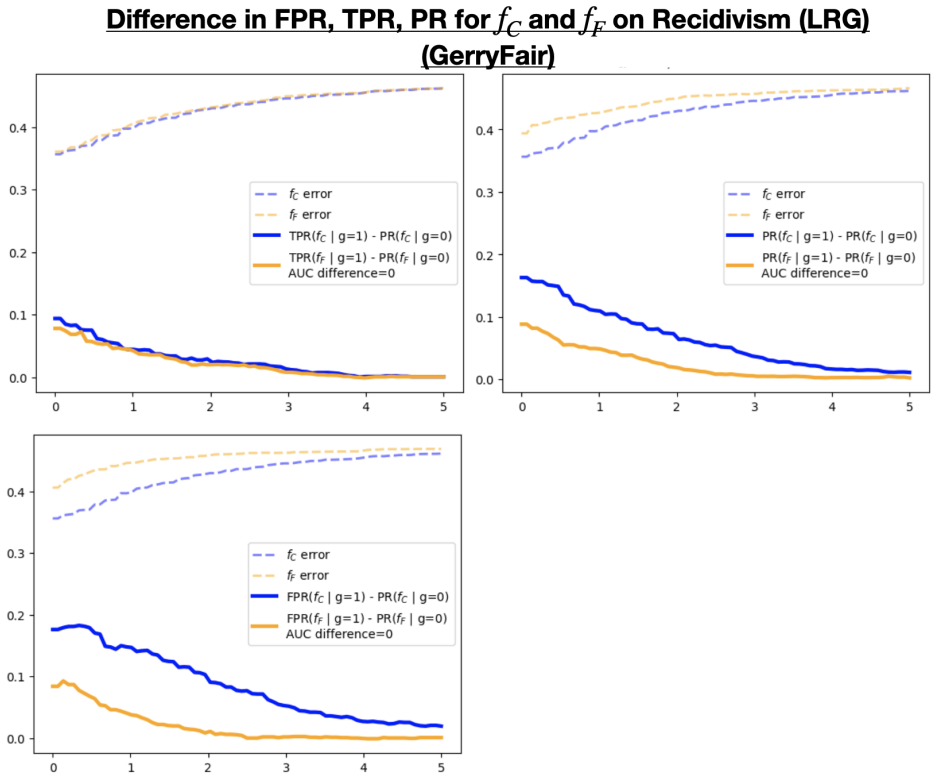


Figure 20:

Figure 21:



Figure 22:

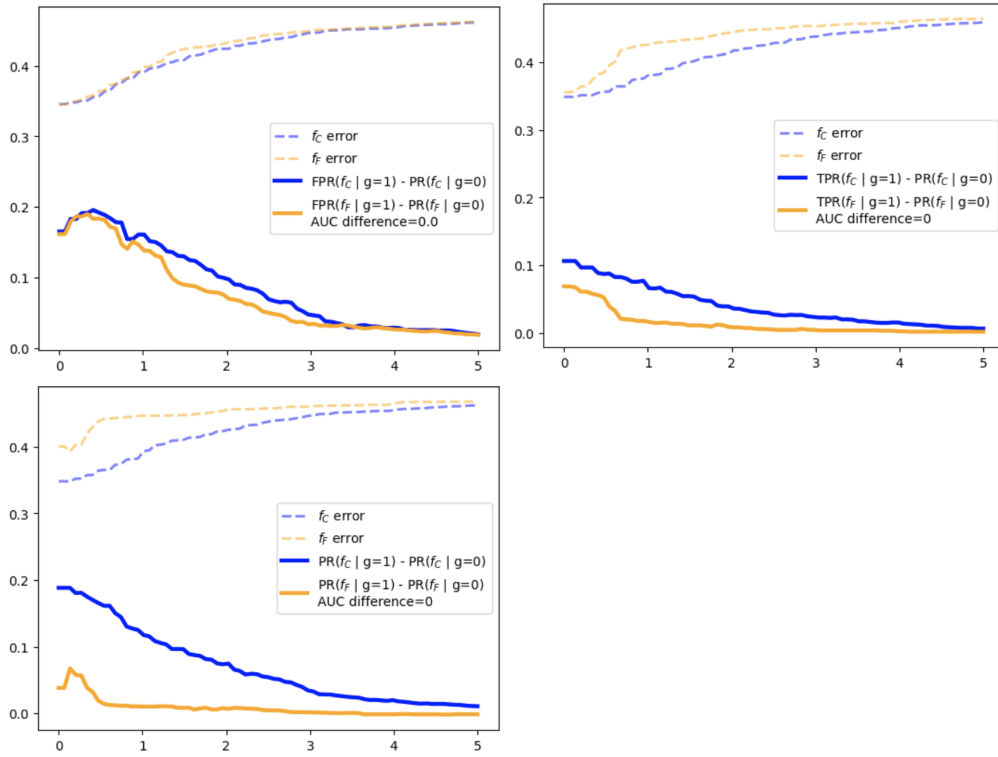**Difference in FPR, TPR, PR for $f_C$ and $f_F$ on Recidivism (NN) (Reduction)**



Figure 23:

**Difference in FPR, TPR for $f_C$ and $f_F$ on Community Crime (SVM) (GerryFair)**
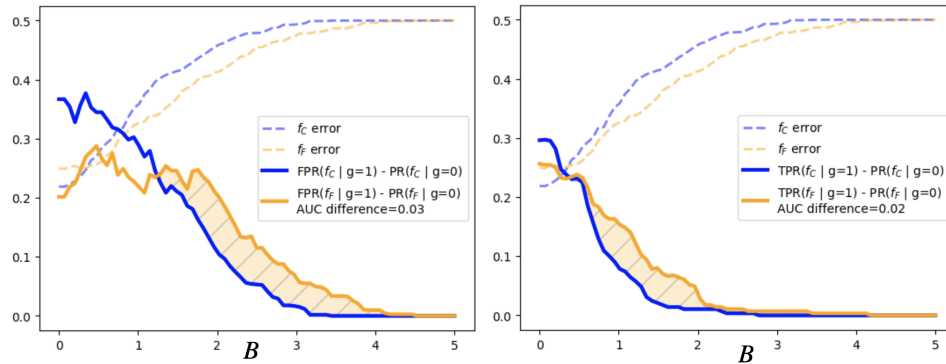


Figure 24: