

---

# Too Relaxed to Be Fair

---

Michael Lohaus<sup>1,2</sup> Michaël Perrot<sup>2,3,4</sup> Ulrike von Luxburg<sup>1,2</sup>

## Abstract

We address the problem of classification under fairness constraints. Given a notion of fairness, the goal is to learn a classifier that is not discriminatory against a group of individuals. In the literature, this problem is often formulated as a constrained optimization problem and solved using relaxations of the fairness constraints. We show that many existing relaxations are unsatisfactory: even if a model satisfies the relaxed constraint, it can be surprisingly unfair. We propose a principled framework to solve this problem. This new approach uses a strongly convex formulation and comes with theoretical guarantees on the fairness of its solution. In practice, we show that this method gives promising results on real data.

## 1. Introduction

Informally, a classifier is considered *unfair* when it unjustly promotes a group of individuals while being detrimental to others; it is considered *fair* when it is free of any unjust behavior. However, the details of what is fair and unfair can be vastly different from one application to another. For example, a college might want to admit a diverse student pool with respect to gender or race. This notion of fairness is called *demographic parity*. On the other hand, consider a bank giving out loans. If a group of individuals repays less frequently than others, it is normal that they receive fewer loans. However, it does not mean that all requests should be declined. In particular, any individual that is likely to repay a loan should be given the opportunity to get one, regardless of group membership. This is called *equality of opportunity*.

---

<sup>1</sup>University of Tübingen, Germany <sup>2</sup>Max Planck Institute for Intelligent Systems, Tübingen, Germany <sup>3</sup>Univ Lyon, UJM-Saint-Etienne, CNRS, IOGS, LabHC UMR 5516, F-42023, SAINT-ETIENNE, France <sup>4</sup>Most of the work was done when M.P. was affiliated with the Max Planck Institute. Correspondence to: Michael Lohaus <michael.lohaus@uni-tuebingen.de>, Michaël Perrot <michael.perrot@univ-st-etienne.fr>, Ulrike von Luxburg <luxburg@informatik.uni-tuebingen.de>.

The problem of learning fair classifiers has mainly been addressed in three ways. First, pre-processing approaches alter the labels of the examples or their representation to increase the intrinsic fairness of a dataset. A classifier learned on this modified data is then more likely to be fair (Feldman et al., 2015; Calmon et al., 2017; Kamiran & Calders, 2012; Dwork et al., 2012; Zemel et al., 2013). Second, post-hoc procedures transform existing accurate but unfair classifiers into fair classifiers (Chzhen et al., 2019; Hardt et al., 2016; Woodworth et al., 2017; Kamiran et al., 2010). Finally, direct methods learn a fair and accurate classifier in a single step (Kamishima et al., 2012; Zafar et al., 2017b;a; Calders & Verwer, 2010; Wu et al., 2019; Donini et al., 2018; Cotter et al., 2019; Agarwal et al., 2018; Goh et al., 2016). In this paper, we focus on the latter kind of approaches.

**Motivation: relaxations sometimes fail to produce fair solutions.** Recently, several direct methods have been proposed that use relaxed versions of the considered fairness constraint. These approaches work reasonably well for some applications. However, their relaxations are quite coarse and we demonstrate below that they can fail to find fair classifiers. In particular, there is typically no guarantee on the relationship between the relaxed fairness and the true fairness of a solution: a classifier that is perfectly fair in terms of relaxed fairness can be highly unfair in terms of true fairness (see Figure 1 for an illustration). In this paper, we study the limitations of a number of popular approaches (Zafar et al., 2017b;a; Wu et al., 2019; Donini et al., 2018).

**Algorithmic contributions.** We propose a new principled framework to tackle the problem of fair classification that is particularly relevant for application scenarios where formal fairness guarantees are required. Our approach is based on convex relaxations and comes with theoretical guarantees that ensure that the learned classifier is fair up to sampling errors. Furthermore, we use a learning theory framework for similarity-based classifiers to exhibit sufficient conditions that guarantee the existence of a fair and accurate classifier.

## 2. Problem Setting

Let  $\mathcal{X}$  be a feature space,  $\mathcal{Y} = \{-1, 1\}$  a space of binary class labels, and  $\mathcal{S} = \{-1, 1\}$  a space of binary sensitive attributes. Assume that there exists a distribution  $\mathcal{D}_{\mathcal{Z}}$  over  $\mathcal{Z} = \mathcal{X} \times \mathcal{S} \times \mathcal{Y}$  and that we can draw some examples

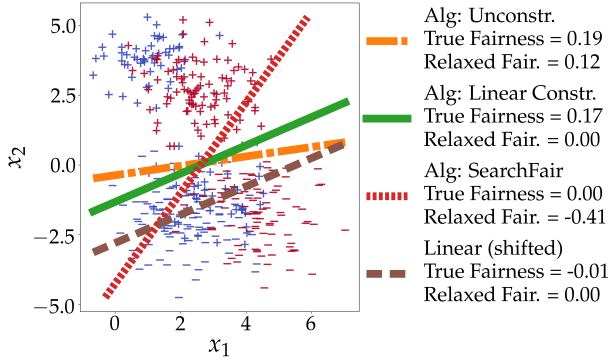


Figure 1. The goal is to separate the positive class (+) from the negative class (-) while remaining fair with respect to two sensitive groups: the blue and the red group. We evaluate the true fairness (DDP) and a linear relaxation of the fairness (Zafar, Section 3.1) of three linear classifiers learned with no fairness constraint (Unconstr., orange), a linear relaxation of the fairness constraint (Linear Constr., green), and our framework (SearchFair, red). We also plot the classifier obtained by translating Linear (Linear (shifted), brown). It has the same relaxed fairness as Linear but a different true fairness: the relaxation is oblivious to the intercept parameter. SearchFair finds the fairest classifier.

$(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}$ . Our goal in fair classification is to obtain a classifier, a mapping  $h : \mathcal{X} \rightarrow \mathcal{Y}$  defined as  $h(x) = \text{sign}(f(x))$  where  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a real valued function, that is fair with respect to the sensitive attribute while remaining accurate on the class labels. In this paper, we study two notions of fairness: *demographic parity* and *equality of opportunity*.

**Demographic Parity.** A classifier  $f$  is fair for demographic parity when its predictions are independent of the sensitive attribute (Calders et al., 2009; Calderys & Verwer, 2010). Formally, this can be written as

$$\mathbb{P}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [f(x) > 0 | s = 1] = \mathbb{P}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [f(x) > 0 | s = -1].$$

In practice, enforcing exact demographic parity might be too restrictive. Instead, we consider a fairness score (Wu et al., 2019) called Difference of Demographic Parity (DDP):

$$\text{DDP}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [\mathbb{I}_{f(x) > 0} | s = 1] - \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [\mathbb{I}_{f(x) > 0} | s = -1], \quad (1)$$

where  $\mathbb{I}_a$  is the indicator function that returns 1 when  $a$  is true and 0 otherwise. The DDP is positive when the favoured group is  $s = 1$  and negative when it is  $s = -1$ . Given a threshold  $\tau \geq 0$ , we say that a classifier  $f$  is  $\tau$ -DDP fair if  $|\text{DDP}(f)| \leq \tau$ . When  $\tau = 0$ , exact demographic parity is achieved and we say that the classifier is DDP fair.

**Equality of Opportunity.** A classifier  $f$  is fair for equality of opportunity when its predictions for positively labelled

examples are independent of the sensitive attribute (Hardt et al., 2016). Formally, it is

$$\mathbb{P}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [f(x) > 0 | y = 1, s = 1] = \mathbb{P}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [f(x) > 0 | y = 1, s = -1].$$

Again, instead of only considering exact equality of opportunity, we use a fairness score (Donini et al., 2018) called Difference of Equality of Opportunity (DEO):

$$\text{DEO}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [\mathbb{I}_{f(x) > 0} | y = 1, s = 1] - \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [\mathbb{I}_{f(x) > 0} | y = 1, s = -1]. \quad (2)$$

This quantity is positive when the favoured group is  $s = 1$  and negative when it is  $s = -1$ . Given a threshold  $\tau \geq 0$ , we say that a classifier  $f$  is  $\tau$ -DEO fair if  $|\text{DEO}(f)| \leq \tau$ . When  $\tau = 0$ , exact equality of opportunity is achieved and we say that the classifier is DEO fair.

It is worth noting that demographic parity and equality of opportunity are quite similar from a mathematical point of view. In the remainder of the paper, we focus our exposition on DDP as results that hold for DDP can often be readily extended to DEO by conditioning on the target label. We only provide details in the supplementary when these extensions are more involved.

**Learning a fair classifier.** Given a function class  $\mathcal{F}$ , a  $\tau$ -DDP fair and accurate classifier  $f^*$  is given as the solution of the following problem:

$$f^* = \arg \min_{\substack{f \in \mathcal{F} \\ |\text{DDP}(f)| \leq \tau}} L(f),$$

where  $L(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} [\ell(f(x), y)]$  is the true risk of  $f$  for a convex loss function  $\ell : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ . In practice, we only have access to a set  $\widehat{\mathcal{D}}_{\mathcal{Z}} = \{(x_i, s_i, y_i)\}_{i=1}^n$  of  $n$  examples drawn from  $\mathcal{D}_{\mathcal{Z}}$ . Hence, we consider the empirical version of this problem:

$$f^\beta = \arg \min_{\substack{f \in \mathcal{F} \\ |\text{DDR}(f)| \leq \tau}} \widehat{L}(f) + \beta \Omega(f), \quad (3)$$

where  $\Omega(f)$  is a convex regularization term used to prevent over-fitting,  $\beta$  is a trade-off parameter, and  $\widehat{L}(f) = \frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} \ell(f(x), y)$  is the empirical risk. The main difficulty involved in learning a fair classifier is to ensure that  $|\text{DDP}(f)| \leq \tau$ .

### 3. When Fairness Relaxations Fail

To obtain a  $\tau$ -DDP fair classifier, most approaches consider the fully empirical version of Optimization Problem 3:

$$\min_{f \in \mathcal{F}} \widehat{L}(f) + \beta \Omega(f) \quad \text{subject to } |\widehat{\text{DDP}}(f)| \leq \tau, \quad (4)$$

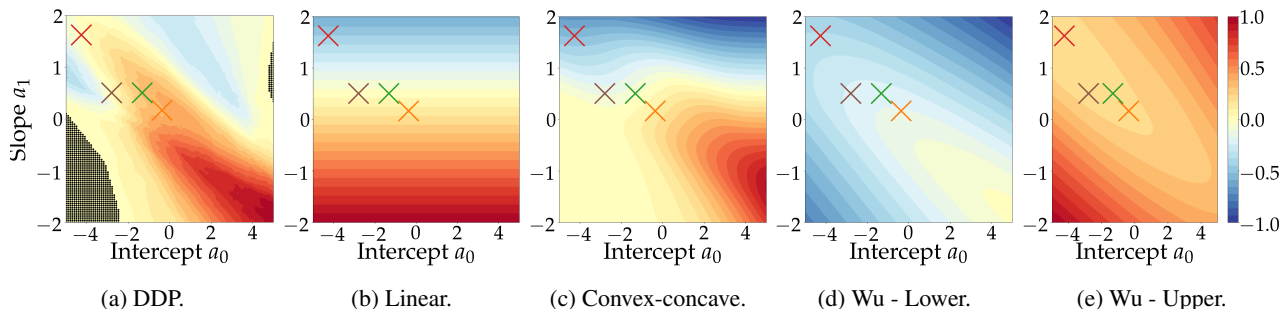


Figure 2. Consider linear classifiers for the dataset in Figure 1. The decision boundaries are of the form  $x_2 = a_1x_1 + a_0$  where  $a_1$  controls the slope and  $a_0$  the intercept. For given intercepts and slopes, we plot normalized values of (a) the DDP score (yellow is fair), (b) the linear relaxation (Section 3.1), (c) the convex-concave relaxation (Section 3.2), (d) the concave Wu lower bound, and (d) the convex Wu upper bound (Section 3.2). The black dotted area in (a) corresponds to trivial constant classifiers—the predicted class is the same for all points. The colored crosses correspond to the classifiers in Figure 1. A good relaxation should capture the true DDP reasonably well, in particular the yellow regions should match. However, none of the considered relaxations manage to achieve this.

where the empirical version of DDP is:

$$\widehat{\text{DDP}}(f) = \frac{1}{n} \sum_{\substack{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}} \\ s=1}} \mathbb{I}_{f(x)>0} - \frac{1}{n} \sum_{\substack{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}} \\ s=-1}} \mathbb{I}_{f(x)>0}.$$

The main issue with this optimization problem is the non-convexity of the constraints that makes it hard to find the optimal solution. A standard approach is then to first rewrite the DDP in an equivalent, but easier to handle form<sup>1</sup> and then replace the indicator functions with a relaxation. Zafar et al. (2017b) and Donini et al. (2018) use a linear relaxation to obtain a fully convex constraint. Zafar et al. (2017a) use a convex relaxation that leads to a convex-concave constraint. Wu et al. (2019) combine a convex relaxation with a concave one to obtain a fully convex problem. Below, we show that these approaches only loosely approximate the true constraint and might lead to suboptimal solutions (see Figure 2). Furthermore, when theoretical guarantees accompany the method, they are either insufficient to ensure that the learned classifier is fair (Wu et al., 2019) or they make assumptions that are hard to satisfy in practice (Donini et al., 2018).

### 3.1. Linear Relaxations

We first study methods that use a linear relaxation of the indicator function to obtain a convex constraint in Optimization Problem 4. First, Zafar et al. (2017b) rewrite the DDP:

$$\text{DDP}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} \frac{1}{p_1(1-p_1)} \left( \frac{s+1}{2} - p_1 \right) \mathbb{I}_{f(x)>0},$$

where  $p_1 = \mathbb{P}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} (s = 1)$  is the proportion of individuals in group  $s = 1$ . Then, they consider a linear

<sup>1</sup>In the supplementary, we provide the derivations for all the alternate formulations of DDP presented in this paper.

approximation of  $\mathbb{I}_{f(x)>0}$  and obtain the constraint:

$$\left| \frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} \frac{1}{\hat{p}_1(1-\hat{p}_1)} \left( \frac{s+1}{2} - \hat{p}_1 \right) f(x) \right| \leq \tau,$$

where  $\hat{p}_1$  is an empirical estimate of  $p_1$ . In their original formulation, Zafar et al. (2017b) get rid of the factor  $\frac{1}{\hat{p}_1(1-\hat{p}_1)}$  by replacing the right hand side of the constraint with  $c = \hat{p}_1(1 - \hat{p}_1)\tau$ .

Similarly, Donini et al. (2018) rewrite the DDP:

$$\text{DDP}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} \frac{s}{p_s} \mathbb{I}_{f(x)>0},$$

where  $p_s = \mathbb{P}_{(x',s',y') \sim \mathcal{D}_{\mathcal{Z}}} (s' = s)$  is the proportion of individuals in group  $s$ . Then, using the same linear relaxation as Zafar et al. (2017b) with  $\hat{p}_s$ , an empirical estimate of  $p_s$ , they obtain the constraint<sup>2</sup>

$$\left| \frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} \frac{s}{\hat{p}_s} f(x) \right| \leq \tau.$$

Both constraints are mathematically close and only differ in terms of the multiplicative factor in front of  $f(x)$  in the inner sum. Thus, they can be rewritten as

$$|\text{LR}_{\widehat{\text{DDP}}}(f)| = \left| \frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} C(s, \widehat{\mathcal{D}}_{\mathcal{Z}}) f(x) \right| \leq \tau.$$

<sup>2</sup>Donini et al. (2018) originally consider  $\tau$ -DEO fairness rather than DDP. In the constraint, instead of drawing the examples from  $\mathcal{D}_{\mathcal{Z}}$ , they use the conditional distribution  $\mathcal{D}_{\mathcal{Z}|y=1}$ . However, this does not change the intrinsic nature of the constraint, and the issues raised here remain valid.

where  $C(s, \widehat{\mathcal{D}}_{\mathcal{Z}})$  can be chosen to obtain any of the two constraints. In the following, we use this general formulation to show that both formulations have shortcomings that can lead to undesired behaviors.

**Linear relaxations are too loose.** In Figures 2(a) and 2(b) we illustrate the behaviors of  $\widehat{\text{DDP}}(f)$  and  $\text{LR}_{\widehat{\text{DDP}}}(f)$ . In the figures, we consider linear classifiers of the form  $f(x) = -x_2 + a_1x_1 + a_0$  where  $a_1$  controls the slope of the classifier and  $a_0$  the intercept. The underlying data is the same as in Figure 1. It shows that the linear relaxation of DDP can behave completely differently compared to the true DDP. It is particularly striking to notice that the intercept does not have any influence on the constraint. This behavior can be formally verified. Let  $f$  be a predictor of the form  $f(x) = g(x) + b$  where  $b$  is the intercept. Then,  $\text{LR}_{\widehat{\text{DDP}}}(f)$  is independent of changes in  $b$  since  $\frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} C(s, \widehat{\mathcal{D}}_{\mathcal{Z}}) = 0$  for both constraints presented above. The proofs are given in the supplementary.

**Theoretical guarantees for linear relaxations are not satisfactory.** Donini et al. (2018) study a sufficient condition under which the linear fairness relaxation  $\text{LR}_{\widehat{\text{DDP}}}(f)$  of a function  $f$  is close to its true fairness, that is it holds that  $|\widehat{\text{DDP}}(f)| \leq |\text{LR}_{\widehat{\text{DDP}}}(f)| + \hat{\Delta}$ . The condition that needs to be satisfied by  $f$  is

$$\frac{1}{2} \sum_{s' \in \{-1,1\}} \left| \frac{1}{2} \sum_{\substack{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}} \\ s=s'}} (\text{sign}(f(x)) - f(x)) \right| \leq \hat{\Delta}.$$

Unfortunately, the left hand side of this condition is non-convex and thus, it is difficult to use in practice. In particular, when they learn a classifier with their linear relaxation, Donini et al. (2018) do not ensure that it also has a small  $\hat{\Delta}$ . They only verify this condition when the learning process is over, that is when a classifier  $f$  has already been produced. However, at this time, it is also possible to compute  $\widehat{\text{DDP}}(f)$  directly, so the bound is not needed anymore.

If one could show that for a given function class  $\mathcal{F}$ , there exists a small  $\hat{\Delta}$  such that the condition holds for all  $f \in \mathcal{F}$ , then any classifier learned from this function class would be guaranteed to be fair when  $|\text{LR}_{\widehat{\text{DDP}}}(f)|$  is small. However, it is not clear whether such function class exists. Nevertheless, this argument hints that for linear relaxations of the fairness constraint, the complexity of the function class largely controls the DDP that can be achieved.

**Linear relaxations should not be combined with complex classifiers.** We demonstrate that, if the class of classifiers  $\mathcal{F}$  is complex, then the linear relaxation constraint has almost no influence on the outcome of the optimization problem. In Figure 3, we compare the performance, in terms

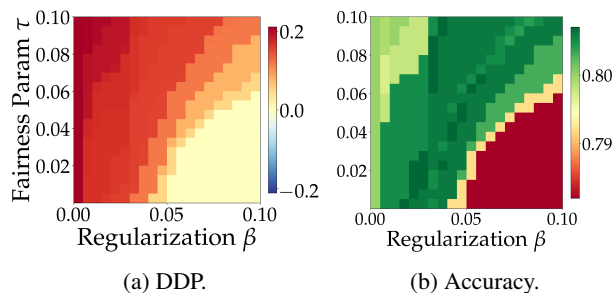


Figure 3. We consider a similarity-based classifier (Section 5) with rbf kernel and 1000 train and test points from the Adult dataset. Using a varying regularization parameter  $\beta$  and fairness parameter  $\tau$ , we train several classifiers using the linear fairness relaxation (Section 3.1). We plot the empirical test DDP of the learned models in Figure 3(a) (red and blue are bad, yellow is good) and their accuracy in Figure 3(b) (red is bad, green is good). We can see that, if  $\beta$  is small (complex model), the fairness relaxation parameter  $\tau$  has no influence on the DDP score. For higher values of  $\beta$  (simpler models), decreasing  $\tau$  improves the DDP. Best viewed in color.

of empirical DDP and accuracy, of several models learned by Optimization Problem 4 equipped with the linear relaxation for different parameters  $\beta$  (for regularization) and  $\tau$  (for fairness). Intuitively, one would expect that varying  $\tau$  leads to changes in the fairness level while varying  $\beta$  leads to changes in accuracy. However, this is not the case:  $\tau$  only has an effect on the result when  $\beta$  is sufficiently large. It means that the fairness of the model is mainly controlled by the regularization parameter rather than the fairness one.

This would not be an issue if the fairness of complex classifiers was small. Unfortunately, high-complexity models have a high capacity to alter their decision boundaries. It means that to achieve both high accuracy and high fairness at the same time, they tend to alter their prediction margin for a few selected examples. While this might not affect the accuracy by a lot, the linear relaxation is sensible to this kind of changes and thus can be largely improved—which is what the optimization aims for. However, altering labels of individual points does not have a big influence on the true DDP: it remains high. This effect is reduced when one learns models of low capacity, which have less freedom to deliberately change labels of individual points. Overall, linear relaxations are mainly relevant for simple classifiers and tend to have little effect on complex ones. We outline this undesirable behavior in the experiments.

### 3.2. Other Relaxations

In the previous section we demonstrated that linear relaxations are not sufficient to ensure fairness of the learned classifier. We now focus on two approaches that use non-linear relaxations of the indicator function to stay close to the original fairness definition.

**Convex-concave relaxation.** In a second paper, Zafar et al. (2017a) use the same fairness formulation as Zafar et al. (2017b), but, instead of a linear relaxation of the indicator function, they use a non-linear relaxation.<sup>3</sup> Hence, given  $\hat{p}_1$  defined as in Section 3.1, they obtain the constraint:

$$|\text{CCR}_{\widehat{\text{DDP}}}(f)| = \left| \frac{1}{n} \sum_{(x,s,y) \in \widehat{\mathcal{D}}_{\mathcal{Z}}} \frac{\left(\frac{s+1}{2} - \hat{p}_1\right)}{\hat{p}_1(1 - \hat{p}_1)} \min(0, f(x)) \right| \leq \tau.$$

In Figure 2(c) we give an illustration of  $\text{CCR}_{\widehat{\text{DDP}}}(f)$ . It more closely approximates the original  $\widehat{\text{DDP}}(f)$  than the linear relaxation. Nevertheless, it remains quite far from the original definition—in particular for classifiers that are not constant. Moreover, using such a convex relaxation leads to a convex-concave problem that turns out to be difficult to optimize without guarantees on the global optimality.

**Lower-upper relaxation with guarantees.** To derive their fairness constraint, Wu et al. (2019) propose to first equivalently rewrite the DDP as follows:

$$\text{DDP}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} \left[ \frac{\mathbb{I}_{s=1} \mathbb{I}_{f(x) > 0}}{p_1} + \frac{\mathbb{I}_{s=-1} \mathbb{I}_{f(x) < 0}}{1 - p_1} - 1 \right]$$

where  $p_1$  is defined as in Section 3.1. Replacing the indicator functions with a convex surrogate other than the linear one would lead to a convex-concave problem due to the absolute value in the constraint. Instead, Wu et al. (2019) propose to use a convex surrogate function  $\kappa$  for the requirement  $\text{DDP}(f) < \tau$  and a concave surrogate function  $\delta$  for  $\text{DDP}(f) > -\tau$ . The corresponding relaxation is

$$\text{DDP}_{\kappa}(f) = \mathbb{E}_{(x,s,y) \sim \mathcal{D}_{\mathcal{Z}}} \left[ \frac{\mathbb{I}_{s=1} \kappa(f(x))}{p_1} + \frac{\mathbb{I}_{s=-1} \kappa(-f(x))}{1 - p_1} - 1 \right],$$

and  $\text{DDP}_{\delta}(f)$  is defined analogously by simply replacing  $\kappa$  with  $\delta$ . It leads to the following convex problem:

$$\begin{aligned} \min_{f \in \mathcal{F}} \quad & \hat{L}(f) + \beta \Omega(f) \\ \text{subject to} \quad & \widehat{\text{DDP}}_{\kappa}(f) \leq \tau_{\kappa} \\ & -\widehat{\text{DDP}}_{\delta}(f) \leq \tau_{\delta}. \end{aligned} \quad (5)$$

Individually, the relaxations are far from the original fairness constraint (as illustrated in Figures 2(e) and 2(d)) but the idea is that combining the upper bound and the lower bound will help to learn a fair classifier. However, one needs to

<sup>3</sup>Zafar et al. (2017a) originally consider other notions of fairness than DPP, among them is the  $\tau$ -DEO fairness (Equation (5) in their paper). Instead of drawing the examples from  $\mathcal{D}_{\mathcal{Z}}$ , they consider the conditional distribution  $\mathcal{D}_{\mathcal{Z}|y=1}$ .

choose  $\tau_{\kappa}$  and  $\tau_{\delta}$  appropriately. To address this, Wu et al. (2019) show that choosing

$$\begin{aligned} \tau_{\kappa} &= \psi_{\kappa} \left( \tau_{\text{upper}} - \widehat{\text{DDP}}_{\kappa}^+ \right) + \widehat{\text{DDP}}_{\kappa}^-, \\ \tau_{\delta} &= \psi_{\delta} \left( \tau_{\text{lower}} + \widehat{\text{DDP}}_{\delta}^- \right) + \widehat{\text{DDP}}_{\delta}^+, \end{aligned}$$

guarantees that  $-\tau_{\text{lower}} \leq \widehat{\text{DDP}}(f) \leq \tau_{\text{upper}}$ . Here  $\widehat{\text{DDP}}_{\kappa}^+$  and  $\widehat{\text{DDP}}_{\delta}^-$  are the worst possible scores of  $\widehat{\text{DDP}}(f)$ : they are attained by those functions in the given function class that advantage either group  $s = -1$  or group  $s = 1$  the most. The values  $\widehat{\text{DDP}}_{\kappa}^-$  and  $\widehat{\text{DDP}}_{\delta}^+$  are defined in the same way for the relaxed scores. The functions  $\psi_{\kappa}$  and  $\psi_{\delta}$  are invertible functions that depend on the selected surrogate.

While this solution is appealing at a first glance, it turns out that Optimization Problem 5 is often infeasible for meaningful values of  $\tau_{\text{upper}}$  and  $\tau_{\text{lower}}$  as the constraints form disjoint convex sets. To illustrate this, consider  $\kappa(x) = \max\{0, 1 + x\}$  and  $\delta(x) = \min\{1, x\}$  as proposed by Wu et al. (2019) and the dataset used in Figure 1. Then, if  $\tau_{\text{upper}} = \tau_{\text{lower}} \leq 1.13$ , the problem is infeasible. If  $\tau_{\text{lower}} = 0$  and  $\tau_{\text{upper}} \leq 1.95$  the problem is also infeasible. Overall, the guarantees are often meaningless: they either make statements about the empty set (no feasible solution) or they are too loose to ensure meaningful levels of fairness.

## 4. New Approach with Guaranteed Fairness

In the previous section, we have seen that existing approaches use relaxations of the fairness constraint that lead to tractable optimization problems but have little control over the true fairness of the learned model. For this reason, we propose a new framework that solves the problem of finding *provably fair* solutions: given a convex approximation of the fairness constraint, our method is guaranteed to find a classifier with a good level of fairness.

We consider the following optimization problem:

$$f_{\widehat{\mathcal{D}}_{\mathcal{Z}}}^{\beta}(\lambda) = \arg \min_{f \in \mathcal{F}} \hat{L}(f) + \lambda \mathbf{R}_{\widehat{\text{DDP}}}(f) + \beta \Omega(f), \quad (6)$$

where  $\mathbf{R}_{\widehat{\text{DDP}}}(f)$  is a convex approximation of the signed fairness constraint, that is we do not consider the usual absolute value. In other words, we obtain a trade-off between accuracy and fairness that is controlled by two hyperparameters  $\lambda \geq 0$  and  $\beta > 0$  and, given  $\beta$  fixed, we can vary  $\lambda$  to move from strongly preferring one group to strongly preferring the other group. Our goal is then to find a parameter setting that is in the neutral regime and does not favor any of the two groups. The main theoretical ingredient for this procedure to succeed is the next theorem, which states that the function  $\lambda \mapsto \text{DDP}\left(f_{\widehat{\mathcal{D}}_{\mathcal{Z}}}^{\beta}(\lambda)\right)$  is continuous under reasonable assumptions on the data distribution, the candidate classifiers, and the convex relaxation.

**Theorem 1 (Continuity of DDP)**  $\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda)\right)$ . Let  $\mathcal{F}$  be a function space, we define the set of learnable functions as  $\mathcal{F}_{\Lambda} = \left\{f \in \mathcal{F} : \exists \lambda \geq 0, f = f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda)\right\}$ . Assume that the following conditions hold:

- (i) Optimization Problem 6 is  $m$ -strongly convex in  $f$ ,
- (ii)  $\forall f \in \mathcal{F}$ ,  $R_{\widehat{\text{DDP}}}(f)$  is bounded in  $[-B, B]$ ,
- (iii)  $\exists \rho$ , a metric, such that  $(\mathcal{F}_{\Lambda}, \rho)$  is a metric space,
- (iv)  $\forall x \in \mathcal{X}$ ,  $g(f) : f \mapsto f(x)$  is continuous,
- (v)  $\forall f \in \mathcal{F}_{\Lambda}$ ,  $f$  is Lebesgue measurable and the set  $\{x : (x, s, y) \in \mathcal{Z}, s = 1, f(x) = 0\}$  is a Lebesgue null set, as well as  $\{x : (x, s, y) \in \mathcal{Z}, s = -1, f(x) = 0\}$ ,
- (vi) the probability density functions  $f_{\mathcal{Z}|s=1}$  and  $f_{\mathcal{Z}|s=-1}$  are Lebesgue-measurable.

Then, the function  $\lambda \mapsto \text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda)\right)$  is continuous.

The proof of this theorem is given in the supplementary. The conditions (i) - (vi) are of a technical nature, but not very restrictive: Condition (i) can be satisfied by using a strongly convex regularization term, for example the squared  $L_2$  norm. Condition (ii) can be satisfied by assuming that  $\mathcal{X}$  is bounded. Condition (iii) is, for example, satisfied by any Hilbert Space equipped with the standard dot product. This includes, but is not restricted to, the set of linear classifiers. Condition (iv) ensures that small changes in the hypothesis, with respect to the metric associated to  $\mathcal{F}$ , also yield small changes in the predictions. Condition (v) ensures that the number of examples for which the predictions are zero is negligible, for example this happens when the decision boundary is sharp. Condition (vi) is satisfied by many usual distributions, for example the Gaussian distribution.

We demonstrate the continuous behavior of DDP on a real dataset in Figure 4. We plot the DDP score and the accuracy of classifiers learned with Optimization Problem 6 for varying parameters  $\lambda$  and  $\beta$ . Given a fixed  $\beta$ , the results support our theoretical findings: there is a smooth transition between favouring the group  $s = 1$  with small  $\lambda$  and favouring the group  $s = -1$  with higher  $\lambda$ . In between, there is always a region of perfect fairness. In the next corollary, we formally state the conditions necessary to ensure the existence of such a DDP-fair classifier.

**Corollary 1 (Existence of a DDP-fair classifier).** Assume that the conditions of Theorem 1 hold and that the convex approximation  $R_{\widehat{\text{DDP}}}(f)$  is chosen such that for Optimization Problem (6) there exist

- (i)  $\lambda_+$  such that  $\text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda_+)\right) > 0$ ,
- (ii)  $\lambda_-$  such that  $\text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda_-)\right) < 0$ .

Then, there exists at least one value  $\lambda_0$  in the interval  $[\min(\lambda_+, \lambda_-), \max(\lambda_+, \lambda_-)]$  such that  $\text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda_0)\right) = 0$ .

We prove this corollary in the supplementary. This sug-

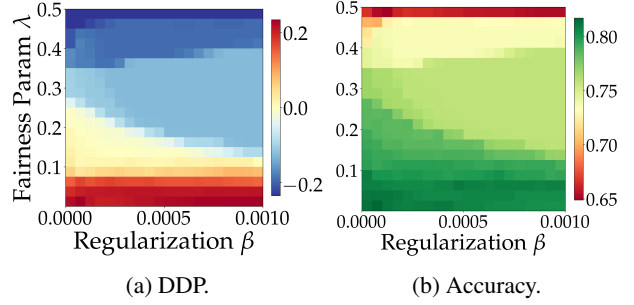


Figure 4. We consider a similarity-based classifier (Section 5) with rbf kernel and 1000 train and test points from the Adult dataset. Using a varying regularization parameter  $\beta$  and fairness parameter  $\lambda$ , we train several classifiers using Optimization Problem 6 with the same loss, convex relaxation, and regularization as SearchFair in the experiments. We plot the empirical test DDP of the learned models in (a) (red and blue are bad, yellow is good) and their accuracy in (b) (red is bad, green is good). We can see that, given a fixed regularization  $\beta$ , we can move from positive DDP (small  $\lambda$ , in red) to a negative DDP (large  $\lambda$ , in blue) with a region of perfect fairness in between (in yellow).

gests a very simple framework to learn provably fair models. First, we choose a convex fairness relaxation (e.g. the one proposed by Wu et al. (2019)) and search for two initial hyperparameters  $\lambda_+$  and  $\lambda_-$  that fulfill the assumptions of Corollary 1 (empirically,  $\lambda = 0$  and 1 are good candidates). Then, we use a binary search to find a  $\lambda_0$  between  $\lambda_+$  and  $\lambda_-$  such that  $\text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda_0)\right) = 0$ . We call this procedure *SearchFair* and summarize it in Algorithm 1 in the supplementary. Note that any convex approximation  $R_{\widehat{\text{DDP}}}(f)$  can be used as long as the conditions of Corollary 1 are respected. In Appendix A we give more details on how we choose this relaxation. Finally, SearchFair theoretically requires to evaluate the true population fairness  $\text{DDP}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda)\right)$  on the underlying distribution  $\mathcal{D}_{\mathcal{Z}}$ . In practice, we follow the example of existing fairness constraints (Woodworth et al., 2017) and simply approximate this quantity by its empirical counterpart  $\widehat{\text{DDP}}\left(f_{\mathcal{D}_{\mathcal{Z}}}^{\beta}(\lambda)\right)$ .

## 5. Towards Classifiers that are Fair and Accurate

In the last section, we presented a method that is guaranteed to find a DDP fair classifier. However, there is one important catch: we did not make any statement about the classification accuracy of this solution. Here, we take a step in this direction by proposing some sufficient conditions that ensure the existence of a classifier that is both fair and accurate. To this end, we focus on a particular set of classifiers: the family of similarity-based functions. Given a similarity function  $K : \mathcal{X} \times \mathcal{X} \rightarrow [-1, 1]$  and a set of points  $S = \{(x'_1, s'_1, y'_1), \dots, (x'_d, s'_d, y'_d)\}$ , we define a similarity

based classifier as  $f(x) = \sum_{i=1}^d \alpha_i K(x, x'_i)$ . The goal is then to learn the weights  $\alpha_i$ .

A theory of learning with such functions has been developed by Balcan et al. (2008). By defining a notion of good similarities, they provide sufficient conditions that ensure the existence of an accurate similarity-based classifier. Here, we build upon this framework and we introduce a notion of good similarities for both accuracy and fairness. Hence, in Definition 1 we give sufficient conditions that ensure the existence of a classifier that is—within a guaranteed margin—fair and accurate at the same time.

**Definition 1 (Good Similarities for Fairness).** A similarity function  $K$  is  $(\varepsilon, \gamma, \tau)$ -good for convex, positive, and decreasing loss  $\ell$  and  $(\mu, \nu)$ -fair for demographic parity if there exists a (random) indicator function  $R(x, s, y)$  defining a (probabilistic) set of “reasonable points” such that, given that  $\forall x \in \mathcal{X}, g(x) = \mathbb{E}_{(x', s', y') \sim \mathcal{D}_{\mathcal{Z}}} [y' K(x, x') | R(x', s', y')]$ , the following conditions hold:

- (i)  $\mathbb{E}_{(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}} \left[ \ell \left( \frac{yg(x)}{\gamma} \right) \right] \leq \varepsilon,$
- (ii)  $\left| \mathbb{P}_{\mathcal{D}_{\mathcal{Z}} | s=1} [g(x) \geq \gamma] - \mathbb{P}_{\mathcal{D}_{\mathcal{Z}} | s=-1} [g(x) \geq \gamma] \right| \leq \mu,$
- (iii)  $\mathbb{P}_{(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}} [|g(x)| \geq \gamma] \geq 1 - \nu,$
- (iv)  $\mathbb{P}_{(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}} [R(x, s, y)] \geq \tau.$

Roughly speaking, a similarity is good for classification if examples of the same class are on average closer to each other than examples of different classes up to a certain margin. Moreover, it is good for fairness if this margin is independent of group membership. Given such a similarity, we can prove the existence of a fair and accurate classifier as is summarized in the next theorem. The proof is given in the supplementary.

**Theorem 2 (Existence of a fair and accurate separator).**

Let  $K \in [-1, 1]$  be a  $(\varepsilon, \gamma, \tau)$ -good and  $(\mu, \nu)$ -fair metric for a given convex, positive and decreasing loss  $\ell$  with Lipschitz constant  $L$ . For any  $\varepsilon_1 > 0$  and  $0 < \delta < \frac{\gamma \varepsilon_1}{2(L + \ell(0))}$ , let  $S = \{(x'_1, s'_1, y'_1), \dots, (x'_d, s'_d, y'_d)\}$  be a set of  $d$  examples drawn from  $\mathcal{D}_{\mathcal{Z}}$  with

$$d \geq \frac{1}{\tau} \left[ \frac{L^2}{\gamma^2 \varepsilon_1^2} + \frac{3}{\delta} + \frac{4L}{\delta \gamma \varepsilon_1} \sqrt{\delta(1 - \tau) \log(2/\delta)} \right].$$

Let  $\phi^S : \mathcal{X} \rightarrow \mathbb{R}^d$  be a mapping with  $\phi_i^S(x) = K(x, x'_i)$ , for all  $i \in \{1, \dots, d\}$ . Then, with probability at least  $1 - \frac{\varepsilon}{2}\delta$  over the choice of  $S$ , the induced distribution over  $\phi^S(\mathcal{X}) \times \mathcal{S} \times \mathcal{Y}$  has a linear separator  $\alpha$  such that

$$\mathbb{E}_{(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}} \left[ \ell \left( \frac{y \langle \alpha, \phi^S(x) \rangle}{\gamma} \right) \right] \leq \varepsilon + \varepsilon_1,$$

and, with  $p_1 = \mathbb{P}_{(x, s, y) \sim \mathcal{D}_{\mathcal{Z}}} [s = 1]$ ,

$$|DDP(\alpha)| \leq \mu + (\nu + 2\delta) \max \left( \frac{1}{p_1}, \frac{1}{1 - p_1} \right).$$

## 6. Experiments

In this section, we empirically evaluate SearchFair by comparing it to 5 baselines on 6 real-world datasets. In all the experiments, SearchFair either reliably finds the fairest classifier and is comparable to a very recent non-convex optimization approach.

**Datasets.** We consider 6 different datasets: CelebA (Liu et al., 2015), Adult (Kohavi & Becker, 1996), Dutch (Zliobaite et al., 2011), COMPAS (Larson et al., 2016), Communities and Crime (Redmond & Baveja, 2002), and German Credit (Dua & Graff, 2017). In the supplementary, we give detailed descriptions of these datasets, how we preprocess the data, and the sizes of the train and test splits. Note that we remove the sensitive attribute  $s$  from the set of features  $x$  so that it is not needed at decision time.

**Baselines.** We compare SearchFair to 5 baselines. For 3 of them, we use Optimization Problem 4 with hinge loss and a squared  $\ell_2$  norm as the regularization term. As a function class  $\mathcal{F}$ , we use similarity-based classifiers presented in Section 5 with either the linear or the rbf kernel and with 70% (at most 1000) of the training examples as reasonable points. As a fairness constraint, we use either the linear relaxation of Zafar et al. (2017b) (Zafar), the linear relaxation of Donini et al. (2018) (Donini), or no constraint at all (Unconst). The fourth baseline is a recent method for non-convex constrained optimization by Cotter et al. (2019) (Cotter). Our last baseline is the constant classifier (Constant) that always predicts the same label but has perfect fairness.

**SearchFair.**<sup>4</sup> For SearchFair we also use the hinge loss, a squared  $\ell_2$  norm as the regularization term (it is strongly convex), and similarity-based classifiers. As a convex approximation of the fairness constraint, we use the bounds with hinge loss proposed by Wu et al. (2019) (see Section A in the supplementary for details).

**Metrics.** Our main goal is to learn fair classifiers. Hence, our main evaluation metrics are the empirical DDP and DEO scores on the test set (lower is better). As a secondary metric (in case of ties in the fairness scores), we consider the classification performance of the models and we report the errors on the test set (lower is better). All the experiments are repeated 10 times and we report the mean and standard deviation for all the metrics.

<sup>4</sup>The code is freely available online: [github.com/mlohaus/SearchFair](https://github.com/mlohaus/SearchFair).

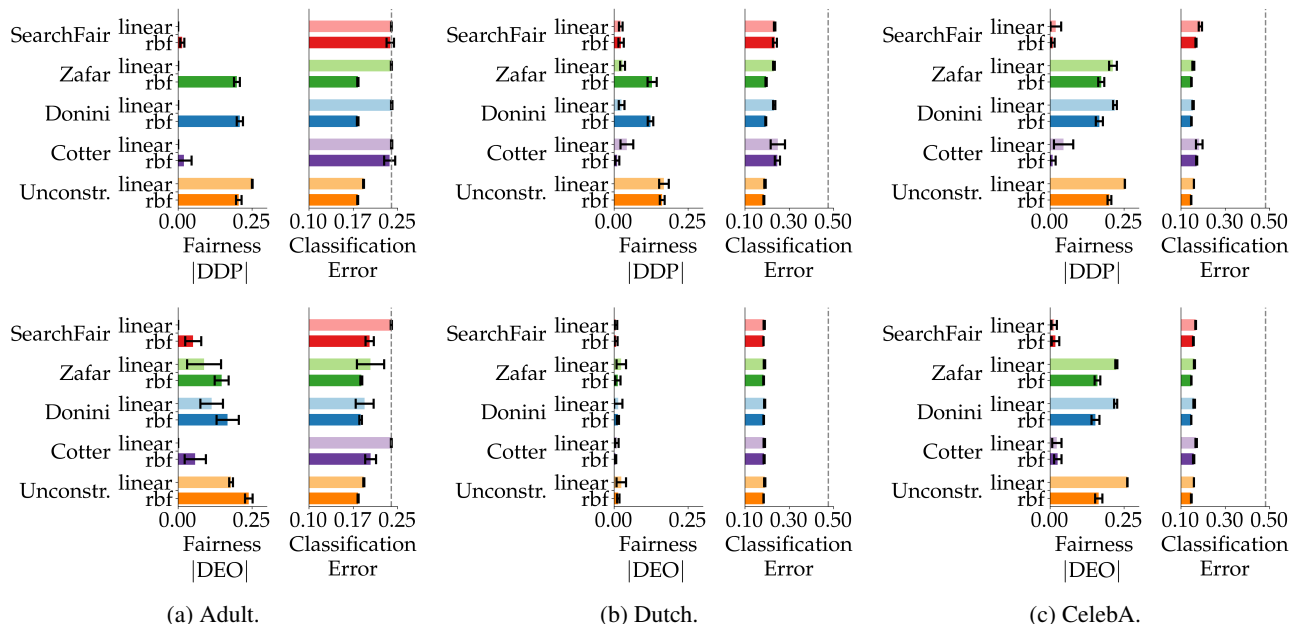


Figure 5. We report the average and standard deviation of classification error and absolute fairness scores DDP and DEO (closer to 0 is better) over 10 repetitions. The constant classifier is perfectly fair as it always predicts the same label. Its classification error is shown by the grey dashed vertical line. (a) To obtain good fairness on Adult, all DDP fair methods learn the constant classifier. Only SearchFair and Cotter reliably find fair classifiers for both kernels. (b) On Dutch, SearchFair obtains the lowest DDP with a slight loss in accuracy. Cotter performs comparably for both kernels, whereas the other methods only do well with a linear kernel but fail to learn fair classifiers with the rbf kernel. (c) For CelebA, SearchFair and Cotter are the only methods that obtain a low DDP and DEO with only a slight loss in accuracy. The other methods only provide little to no improvement.

**Hyperparameters.** Zafar, Donini and Cotter use a fairness parameter, that we call  $\tau$ , to control the fairness level. Since our goal is to learn classifiers that are fair, we set  $\tau = 0$  such that perfect fairness is required. For SearchFair, there is no fairness parameter since  $\lambda_0$  is automatically searched for between a lower bound  $\lambda_{\min}$  and an upper bound  $\lambda_{\max}$ . We set  $\lambda_{\min} = 0$  and  $\lambda_{\max} = 1$  as these values usually lead to classifiers with fairness scores of opposite sign (as needed). We use 10 iterations in the binary search.

We use 5-fold cross validation to choose other hyperparameters. For Cotter, only the width of the rbf kernel has to be tuned since we use the framework of the original paper with no regularization term. For all remaining methods we need to choose the regularization parameter  $\beta$  and the width of the rbf kernel. These values are respectively chosen in the sets  $\{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}\}$  and  $\{10^{-\lceil \log d \rceil - 1}, 10^{-\lceil \log d \rceil}, d^{-1}, 10^{-\lceil \log d \rceil + 1}, 10^{-\lceil \log d \rceil + 2}\}$ , with  $d$  the number of features. We select the set of parameters that lead to the most accurate classifier on average over the 5 folds. Indeed, the fairness level is automatically taken care of by the methods.

**Results.** We present the results for 3 out of 6 datasets in Figure 5. The other results are deferred to the supplementary as they follow the same trend. We make two main obser-

ations. First, SearchFair always obtains fairness values that are very close to zero. It learns the fairest classifiers out of all the methods and is only matched by Cotter, the non-convex approach. This sometimes comes with a small increase in terms of classification error. For example, in order to achieve perfect DDP fairness on the Adult dataset, SearchFair, and all the other fair methods, yield classifiers close to the trivial constant one. Second, the complexity of the model greatly influences the performances of the linear relaxations. For example, using the complex rbf kernel almost always results in an increase in the fairness score of Zafar and Donini. This is particularly striking for Adult and Dutch where the linear kernel yields reasonable fairness scores. Note that this trend is not always respected. For example, on CelebA, using an rbf kernel improves the fairness score compared to the linear kernel. However, neither of them obtain reasonable fairness levels in the first place.

**Discussion on hyperparameter selection.** Apart from the hyperparameter selection method used in our experiments, one can think of other cross validation procedures. For example, Donini et al. (2018) proposed NVP, a cross validation method where one selects the set of hyperparameters that gives the fairest classifier while obtaining an average accuracy above a given threshold. Similarly, one could select the set of hyperparameters that yields the most accurate



classifier under a given fairness threshold. In the supplementary, we report results that empirically show that these more complex procedures tend to improve the fairness of the baselines (SearchFair remains competitive on all the datasets). Unfortunately, they also blur the dividing line between hyperparameters that control the fairness of the model and the ones that control its complexity. In other words, it becomes unclear whether fairness is achieved thanks to the relaxation or thanks to the choice of hyperparameters (we already evoked this issue in Figure 3). We believe that it is better to have a method that is guaranteed to find a fair classifier for any given family of models and does not rely on a complex cross validation procedure.

## 7. Conclusion

In this paper, we have shown that existing approaches to learn fair and accurate classifiers have many shortcomings. They use loose relaxations of the fairness constraint and guarantees that relate the relaxed fairness to the true fairness of the solutions are either missing or not sufficient. We empirically demonstrated how these approaches can produce undesirable models. If “fair machine learning” is supposed to be employed in real applications in society, we need algorithms that actually find fair solutions, and ideally come with guarantees. We made a first step in this direction by proposing SearchFair, an approach that uses convex relaxations to learn a classifier that is guaranteed to be fair.

## Acknowledgements

This work has been supported by the German Research Foundation through the Cluster of Excellence “Machine Learning – New Perspectives for Science” (EXC 2064/1 number 390727645), the BMBF Tübingen AI Center (FKZ: 01IS18039A), and the International Max Planck Research School for Intelligent Systems (IMPRS-IS). It has also been supported by the ACADEMICS grant of the IDEXLYON, project of the Université de Lyon, PIA operated by ANR-16-IDEX-0005. The authors also thank Luca Rendsburg for helpful discussions.

## References

- Agarwal, A., Beygelzimer, A., Dudik, M., Langford, J., and Wallach, H. A reductions approach to fair classification. In *International Conference on Machine Learning*, 2018.
- Balcan, M.-F., Blum, A., and Srebro, N. Improved guarantees for learning via similarity functions. In *Conference on Learning Theory*, 2008.
- Calders, T. and Verwer, S. Three naive Bayes approaches for discrimination-free classification. In *International Conference on Data Mining and Knowledge Discovery*, 2010.
- Calders, T., Kamiran, F., and Pechenizkiy, M. Building classifiers with independency constraints. In *International Conference on Data Mining Workshops*, 2009.
- Calmon, F., Wei, D., Vinzamuri, B., Natesan Ramamurthy, K., and Varshney, K. R. Optimized pre-processing for discrimination prevention. In *International Conference on Neural Information Processing Systems*, 2017.
- Chzhen, E., Denis, C., Hebiri, M., Oneto, L., and Pontil, M. Leveraging labeled and unlabeled data for consistent fair binary classification. In *International Conference on Neural Information Processing Systems*, 2019.
- Cotter, A., Jiang, H., and Sridharan, K. Two-player games for efficient non-convex constrained optimization. In *International Conference on Algorithmic Learning Theory*, 2019.
- Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J., and Pontil, M. Empirical risk minimization under fairness constraints. In *International Conference on Neural Information Processing Systems*, 2018.
- Dua, D. and Graff, C. UCI machine learning repository, 2017.
- Dwork, C., Hardt, M., Pitassi, T., Reingold, O., and Zemel, R. Fairness through awareness. In *Innovations in Theoretical Computer Science Conference*, 2012.
- Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., and Venkatasubramanian, S. Certifying and removing disparate impact. In *International Conference on Knowledge Discovery and Data Mining*, 2015.
- Goh, G., Cotter, A., Gupta, M., and Friedlander, M. P. Satisfying real-world goals with dataset constraints. In *International Conference in Neural Information Processing Systems*, 2016.
- Hardt, M., Price, E., and Srebro, N. Equality of opportunity in supervised learning. In *International Conference on Neural Information Processing Systems*, 2016.
- Kamiran, F. and Calders, T. Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems*, 2012.
- Kamiran, F., Calders, T., and Pechenizkiy, M. Discrimination aware decision tree learning. In *International Conference on Data Mining*, 2010.
- Kamishima, T., Akaho, S., Asoh, H., and Sakuma, J. Fairness-aware classifier with prejudice remover regularizer. In *Machine Learning and Knowledge Discovery in Databases*, 2012.

- Kohavi, R. and Becker, B. Scaling up the accuracy of naive-bayes classifiers: a decision-tree hybrid, 1996.
- Larson, J., Mattu, S., Kirchner, L., and Angwin, J. How we analyzed the compas recidivism algorithm. *ProPublica*, 2016.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *International Conference on Computer Vision*, 2015.
- Redmond, M. A. and Baveja, A. A data-driven software tool for enabling cooperative information sharing among police departments. *European Journal of Operational Research*, 2002.
- Woodworth, B., Gunasekar, S., Ohannessian, M. I., and Srebro, N. Learning non-discriminatory predictors. In *Conference on Learning Theory*, 2017.
- Wu, Y., Zhang, L., and Wu, X. On convexity and bounds of fairness-aware classification. In *The World Wide Web Conference*, 2019.
- Zafar, M. B., Valera, I., Gomez Rodriguez, M., and Gummadi, K. P. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *International Conference on World Wide Web*, 2017a.
- Zafar, M. B., Valera, I., Rogriguez, M. G., and Gummadi, K. P. Fairness Constraints: Mechanisms for Fair Classification. In *International Conference on Artificial Intelligence and Statistics*, 2017b.
- Zemel, R., Wu, Y., Swersky, K., Pitassi, T., and Dwork, C. Learning fair representations. In *International Conference on Machine Learning*, 2013.
- Zliobaite, I., Kamiran, F., and Calders, T. Handling conditional discrimination. In *International Conference on Data Mining*, 2011.