
Efficient Post-Processing for Equal Opportunity in Fair Multi-Class Classification

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Fairness in machine learning is of growing concern as more instances of biased
2 model behavior are documented while their adoption continues to rise. The majority
3 of studies have focused on binary classification settings, despite the fact that many
4 real-world problems are inherently multi-class. This paper considers fairness
5 in multi-class classification under the notion of parity of true positive rates—an
6 extension of binary class equalized odds [23]—which ensures equal opportunity
7 to qualified individuals regardless of their demographics. We focus on algorithm
8 design and provide a post-processing method that derives fair classifiers from pre-
9 trained score functions. The method is developed by analyzing the representation
10 of the optimal fair classifier, and is efficient in both sample and time complexity,
11 as it is implemented by linear programs on finite samples. We demonstrate its
12 effectiveness at reducing disparity on benchmark datasets, particularly under large
13 numbers of classes, where existing methods fall short.

14 1 Introduction

15 Algorithmic fairness has emerged as a topic of significant concern in the field of machine learning,
16 due to the potential for models to exhibit discriminatory behavior towards historically disadvantaged
17 demographics [9, 4, 6], all while their adoption continues to rise in domains including high-stakes
18 areas such as criminal justice, healthcare, and finance [3, 7]. To address the concern, a variety of
19 fairness criteria have been proposed (e.g., demographic parity, equalized odds) along with mitigation
20 methods [10, 19, 23, 26]. On classification problems, the majority of work focuses on the binary class
21 setting [2, Table 1], where one class is typically considered to be more favorable (e.g., the approval
22 vs. rejection of a credit card application).

23 Yet, many real-world problems are multi-class in nature. In the case of credit card applications,
24 issuers may opt to assigning higher-tier interest rates to high-risk applicants rather than outright
25 rejecting them, which creates opportunities to applicants who would otherwise be denied credit
26 and also generates returns for the banks. Similarly, in online advertising, recruiting platforms can
27 employ machine learning models to match users to relevant job postings across multiple occupation
28 categories. There are evidences, however, for such systems to exhibit gender bias [8, 13, 44]; for
29 instance, models that are trained to identify occupation from biography tend to show higher accuracy
30 (recall) on male biographies than on their female counterparts in occupations that are historically
31 male-dominated [14].

32 In the example above, unfairness is manifested in a disparity of *true positive rates* (TPRs) across
33 demographic groups A (generalizing the true positive and negative rates in binary classification),

$$\text{TPR}_a(\hat{Y})_y := \mathbb{P}(\hat{Y} = y \mid Y = y, A = a), \quad \forall y \in [k], a \in [m].$$

34 A classifier satisfying parity of TPRs, i.e., $\text{TPR}_a = \text{TPR}_{a'}$ for all a, a' , ensures that individuals with
35 the same qualification (Y) will have *equal opportunity* of receiving their favorable outcome ($\hat{Y} = Y$)

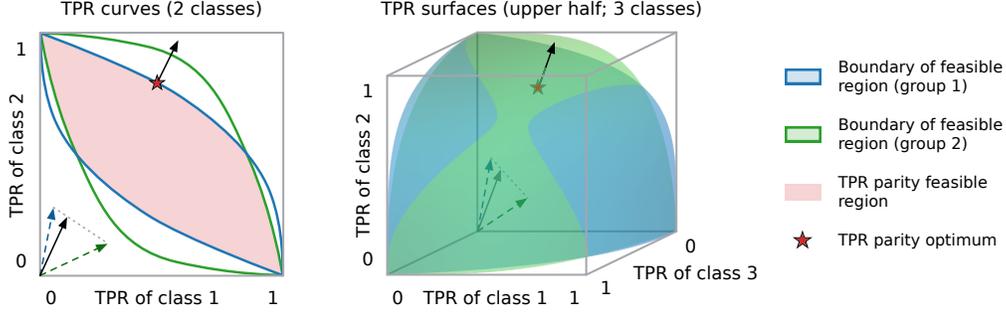


Figure 1: Feasible region of TPRs on a binary class (left) and a three-class problem (right). The black (resp. colored) arrow indicates the utility-maximizing direction (of each group).

36 regardless of demographics [20], e.g., being shown job postings on recruiting platforms for which the
 37 user is qualified. When the classes are binary, this fairness notion recovers *equalized odds* [23].

38 In this paper, we focus on the design of algorithm for mitigating TPR disparity and provide an efficient
 39 *post-processing* method that derives *attribute-aware* fair classifiers from (pre-trained) scoring models.
 40 Our method works on multi-class and multi-group classification problems, guarantees fairness by a
 41 sample complexity bound, can be implemented by linear programs, and achieves higher reductions in
 42 disparity compared to existing algorithms that are applicable to multi-class—a recently proposed post-
 43 processing method based on model projection [2], and adversarial debiasing [41], an in-processing
 44 method—especially when the number of classes is large.

45 **Organization.** We introduce the problem setup and objectives in Section 2, then describe our
 46 post-processing method for TPR parity in Section 3, along with suboptimality analyses; in particular,
 47 our method yields the optimal fair classifier when applied to the *Bayes optimal* score function.
 48 Our method is instantiated for finite sample estimation in Section 4, and we also provide sample
 49 complexity bounds to complete the analysis. Finally, in Section 5, we compare our algorithm with
 50 existing methods for disparity reduction on benchmark datasets.¹ A high-level summary of our results
 51 is provided in Section 1.1.

52 1.1 Summary of Results

53 One way to interpret and understand TPR parity is through visualizing the feasible regions of TPRs.
 54 In Fig. 1, we plot the feasible regions (achievable by probabilistic classifiers) of two groups on a
 55 (hypothetical) binary classification problem on the left, and those on a three-class problem on the
 56 right, where each axis represents the TPR of a class. Achieving optimal TPR parity amounts to
 57 first finding the TPR that maximizes the overall utility (e.g., accuracy) in the intersection of feasible
 58 regions, and subsequently an (attribute-aware) classifier attaining that target TPR on all groups. Note
 59 that the left figure is equivalent to the ROC curve (with a flip of the horizontal axis, because the TPR
 60 of class 1 equals one minus the false negative rate by treating class-1 as the negative class), which
 61 was used by Hardt et al. [23] for studying equalized odds. And thus, the TPR (hyper)surface plots in
 62 higher dimensions are a natural generalization of the ROC curve to multi-class settings.

63 Step one of finding the optimal fair TPR can be formulated as a linear program when estimating from
 64 finite samples. For the second step, our method derives a classifier attaining the target TPR from the
 65 score function; in particular, it yields the optimal fair classifier when the score is Bayes optimal:

Theorem 1.1. *Let $f_1^*, \dots, f_m^* : \mathcal{X} \rightarrow \Delta_k$ denote the Bayes score function on each group, $f_a^*(x) := \mathbb{E}[Y | X = x, A = a]$, and $q_1, \dots, q_m \in \Delta_k$ be arbitrary. Then under a continuity assumption (2.3), $\exists \beta_1, \dots, \beta_m \in [0, 1]$ and $\lambda_1, \dots, \lambda_m \in \mathbb{R}^k$ s.t. the probabilistic attribute-aware classifier*

$$(x, a) \mapsto \begin{cases} \arg \max_{y'} (\lambda_a)_{y'} \cdot f_a^*(x)_{y'} & \text{w.p. } 1 - \beta_a \\ y & \text{w.p. } \beta_a \cdot (q_a)_y, \quad \forall y \in [k] \end{cases} \quad (1a)$$

66 *achieves the maximum utility subject to TPR parity.*

¹Our code is provided in the supplemental material.

67 The post-processed classifier returned by our method is a mixture of two models (weighted by
68 β). Eq. (1a) returns the class with the highest likelihood after a class-wise rescaling, called a
69 *tilting* [2], which generalizes the concept of *thresholding* in binary classifiers. Eq. (1b) makes random
70 assignments sampled from a Multinoulli(q) distribution, which handles situations where the fair TPR
71 lies in the interior of the feasible region (see Fig. 1, where the optimum is located within the interior
72 of group 2 feasible region). To alleviate potential ethical concerns regarding this randomization, we
73 point out that the parameter q_a 's used in class sampling can be specified per-group by the practitioner
74 responsibly, e.g., uniform $1/k$, or $e_{y'}$ with y' being an advantaged outcome.

75 Among the possibly infinitely many fair classifiers derived from the score function f , our method
76 specifically seeks the simplistic representation in Eq. (1) because it can be estimated via linear
77 programs from finite samples. More importantly, it immediately extrapolates to unseen examples,
78 and provides good generalization performance at the rate of $\tilde{O}(\sqrt{k/n})$ thanks to its low function
79 complexity (Theorem 4.2).

80 When the score function being post-processed is not Bayes optimal, our method is still applicable, but
81 the resulting classifier may not be optimal nor exactly achieve TPR parity without access to labeled
82 data (the method itself only needs unlabeled data with the sensitive attribute) or additional knowledge
83 of the model. But these suboptimality are minimized if the model is *calibrated* (Theorem 3.5);
84 this answers the question raised in [2] about the effects of base model inaccuracies on downstream
85 post-processing.

86 1.2 Related Work

87 **Fairness Criteria.** The notion of TPR parity has appeared in the literature as *conditional procedure*
88 *accuracy equality* [7], *avoiding disparate mistreatment* [39], and (multi-class) *equal opportunity* [14,
89 29, 31] (to be distinguished from the fairness criterion with the same name in [23]). Other group
90 fairness notions that extend to multi-class include (but not limited to) *equalized odds* [23] (of which
91 TPR parity is a necessary condition), and *demographic parity* (DP) [10] (where Xian et al. [35]
92 recently proposed an optimal post-processing method). However, DP may be less desirable than
93 TPR parity in some use cases because the perfect classifier is not permitted under DP when the base
94 rates differ [42]. It is worth noting that TPR parity implies *accuracy parity* [9]. In addition to group
95 fairness, there are notions defined on the individual level [19].

96 **Mitigation Methods.** Our method is based on post-processing [25, 23]. There are also in-processing
97 methods via fair representation learning [40, 41, 43, 30] or solving zero-sum games [1, 36], and
98 pre-processing methods that debias the data prior to model training [11, 44]; see [4, 12] for a survey.

99 For multi-class TPR parity, the only applicable post-processing method to date, to our knowledge,
100 is due to Alghamdi et al. [2] (which is the primary baseline for our method in our experiments).
101 It is a general-purpose method that transforms the scores to satisfy fairness while minimizing the
102 distributional divergence (e.g., KL) between the transformed scores and the original. However, the
103 tradeoff between model performance and fairness is unclear as they did not relate the divergence to
104 utility. Furthermore, while the authors provided a sample complexity bound for their optimization
105 objective, it is not explicitly related to the violation of the fairness criteria.

106 2 Preliminaries

107 A k -class classification problem is defined by a joint distribution μ of input $X \in \mathcal{X}$, demographic
108 group membership $A \in [m] := \{1, \dots, m\}$ (a.k.a. the sensitive attribute), and class label $Y \in [k]$.
109 We denote the joint distribution of (X, A) by $\mu^{X,A}$, and, the $(k-1)$ -dimensional probability simplex
110 by $\Delta_k := \{z \in \mathbb{R}_{\geq 0}^k : \|z\|_1 = 1\}$.

111 Let $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$ be an attribute-aware (pre-trained) score function, whose outputs are probability
112 vectors that estimate the class probabilities as in $f(x, a)_y \approx \mathbb{P}_\mu(Y = y \mid X = x, A = a)$. We will
113 write $f_a : \mathcal{X} \rightarrow \Delta_k$ to denote the component of f associated with group a , i.e., $f_a(x) \equiv f(x, a)$. Our
114 goal is to find fair (probabilistic) post-processing maps $g_1, \dots, g_m : \Delta_k \rightarrow \mathcal{Y}$ to derive a classifier
115 $(x, a) \mapsto g_a \circ f_a(x)$ that satisfies TPR parity while maximizing utility (e.g., classification accuracy).

116 We allow for controllable tradeoffs between utility and fairness through the following relaxation of
117 TPR parity, and call a classifier α -fair if it satisfies α -TPR parity:

118 **Definition 2.1** (Approximate TPR Parity). Let $\alpha \in [0, 1]$. A predictor \hat{Y} is said to satisfy α -TPR
 119 parity if $\Delta_{\text{TPR}}(\hat{Y}) \leq \alpha$, where

$$\Delta_{\text{TPR}}(\hat{Y}) := \max_{a, a' \in \mathcal{A}} \left\| \text{TPR}_a(\hat{Y}) - \text{TPR}_{a'}(\hat{Y}) \right\|_{\infty}, \quad (2)$$

120 and $\text{TPR}_a(\hat{Y}) := \mathbb{P}(\hat{Y} | Y = y, A = a) \in [0, 1]^k$; \mathbb{P} includes the randomness of the predictor.

121 Beyond classification accuracy, we also allow for any utility functions that depend only on the TPRs:²

122 **Definition 2.2** (Utility). The utility function $u : [k] \times [k] \rightarrow \mathbb{R}$ is defined for some $v \in \mathbb{R}^k$ by

$$u(\hat{y}, y) := \sum_{y' \in [k]} v_{y'} \mathbb{1}[y = y', \hat{y} = y'].$$

123 E.g., accuracy, $\mathbb{1}[y = \hat{y}]$, is obtained by setting $v = \mathbf{1}_k$. The term v will appear in our analyses,
 124 and the significance of considering utilities of this form is that we could evaluate a classifier by a
 125 weighted sum of its TPRs. Define $p_{ay} := \mathbb{P}_{\mu}(A = a, Y = y)$, then

$$\mathcal{U}(\hat{Y}) = \mathbb{E} u(\hat{Y}, Y) = \sum_{a \in [m], y \in [k]} v_y p_{ay} \text{TPR}_a(\hat{Y})_y \equiv \mathcal{U}(\text{TPR}_1(\hat{Y}), \dots, \text{TPR}_m(\hat{Y})). \quad (3)$$

126 Finally, we make the following continuity assumption on the distributions of score to avoid technical
 127 complexities related to tie-breaking (on the atoms). This assumption has also appeared in prior work
 128 on fair post-processing [16, 21, 35]; it holds when the input distributions are continuous and the score
 129 function is injective, or can be satisfied by adding small random perturbations to the scores.

130 **Assumption 2.3.** The conditional distribution of score, $\mathbb{P}(f_a(X) | A = a)$, is (Lebesgue absolutely)
 131 continuous, $\forall a \in [m]$.

132 3 TPR Parity via Post-Processing

133 Given a score function $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$, and access to the (unlabeled) joint distribution $\mu^{X,A}$ (i.e.,
 134 no estimation error), we describe a method for deriving an attribute-aware α -fair classifier while
 135 maximizing utility, in the form of $(x, a) \mapsto g_a \circ f_a(x)$, where the g_a 's are (probabilistic) fair
 136 post-processing maps for each group. That is, we want to solve

$$\max_{g_1, \dots, g_m} \mathcal{U}(\hat{Y}) \quad \text{s.t.} \quad \Delta_{\text{TPR}}(\hat{Y}) \leq \alpha \quad \text{where} \quad \hat{Y} = g_A \circ f_A(X).$$

137 Although the method only returns classifiers derived from f as opposed to searching over the space of
 138 all classifiers $h : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{Y}$, it would yield the optimal fair classifier provided that the information
 139 of (A, Y) is preserved in the output of f ; this is the case when the score function is Bayes optimal.

140 3.1 Deriving Optimal Fair Classifier From Bayes Score Function

141 In this section, we explain how to obtain an optimal fair classifier by deriving from the Bayes score
 142 function f^* , thereby providing a *proof of Theorem 1.1* (omitted proof are deferred to the appendix).

143 **Step 1** (Finding Utility-Maximizing Fair TPRs). Let $D_a \subseteq [0, 1]^k$ denote the set of feasible TPRs
 144 on group a achieved by probabilistic classifiers. The first step is to find utility-maximizing fair TPRs
 145 contained in an ℓ_{∞} -ball of diameter α per Definition 2.1 of α -TPR parity (left figure of Fig. 2):

$$\max_{t_1 \in D_1, \dots, t_m \in D_m} \mathcal{U}(t_1, \dots, t_m) \quad \text{s.t.} \quad \|t_a - t_{a'}\|_{\infty} \leq \alpha, \quad \forall a, a' \in [m]. \quad (4)$$

146 When $\alpha = 0$, this reduces to finding a single $t \in \bigcap_a D_a$, and because each D_a is convex (since
 147 probabilistic classifiers are allowed), it can be found with ternary search as suggested in [23]. If
 148 instead the t_a 's are to be estimated from finite samples, then the empirical \hat{D}_a 's are described by
 149 polytopes and the problem can be formulated as a linear program (Section 4).

²This includes all possible utility/loss functions in binary classification, since $\text{TPR}(\hat{Y})_1$ (true negative rate)
 and $\text{TPR}(\hat{Y})_2$ (true positive rate) fully determine the 2×2 confusion matrix.

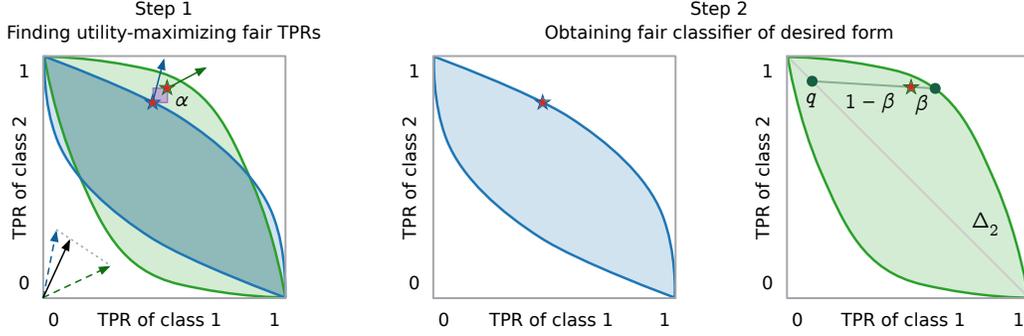


Figure 2: Achieving α -TPR parity on a binary class problem. First, the utility-maximizing TPRs residing in an ℓ_∞ -ball of diameter α are found (left). Then, classifiers achieving the fair TPRs are obtained: a tilting of the scores when the TPR lies on the boundary (middle), otherwise, a mixture of tilting and randomization (right). The simplex Δ_k is always inscribed in the feasible region.

150 The feasible regions of TPR generally differ across groups, due to uncertainties that are inherent to
 151 each group in the task of interest, or to inadequate and biased collection or sourcing of data. The more
 152 the D_a 's differ, the greater the tradeoff between fairness and utility; hence TPR parity incentivizes
 153 the practitioner to improve data collection and aspects of modeling that induces a balanced predictive
 154 capability on all groups [23].

155 Because $f^*(X, A)$ is *sufficient statistic* for Y , the fair TPRs we found above are always achievable
 156 by classifiers derived from f^* . Or more concretely,

157 **Proposition 3.1.** Let $f^* : \mathcal{X} \rightarrow \Delta_k$ denote the Bayes score function, then $D := \{\text{TPR}(h) \in [0, 1]^k \mid$
 158 $h : \mathcal{X} \rightarrow \mathcal{Y} \text{ (probabilistic)}\} = \{\text{TPR}(g \circ f^*) \in [0, 1]^k \mid g : \Delta_k \rightarrow \mathcal{Y} \text{ (probabilistic)}\}.$

159 **Step 2 (Obtaining Fair Classifier of Desired Form).** Having found the utility-maximizing fair TPR
 160 t_a 's, the next step is to derive a classifier that attains t_a on each group. This is provided by the
 161 following theorem:

162 **Theorem 3.2.** Let $f^* : \mathcal{X} \rightarrow \Delta_k$ denote the Bayes score function, and $q \in \Delta_k$ be arbitrary. Then
 163 under Assumption 2.3, $\forall t \in D$, there exists $\beta \in [0, 1]$ and $\lambda \in \mathbb{R}^k$ s.t. $\text{TPR}(h) = t$, where

$$h(x) = \begin{cases} \arg \max_{y'} \lambda_{y'} f^*(x)_{y'} & \text{w.p. } 1 - \beta \\ y & \text{w.p. } \beta q_y, \forall y \in [k]. \end{cases}$$

164 The construction uses the observation that the boundary of D , denoted by ∂D , is given by the set of
 165 TPRs attained by tiltings of the Bayes score:

166 **Proposition 3.3.** Let $f^* : \mathcal{X} \rightarrow \Delta_k$ denote the Bayes score function, then $h : \mathcal{X} \rightarrow \mathcal{Y}$ (probabilistic)
 167 satisfies $\text{TPR}(h) \in \partial D$ if and only if $\exists \lambda \in \mathbb{R}^k, \lambda \neq 0$ s.t. $h(x) \in \arg \max_y \lambda_y f^*(x)_y$.

168 *Proof of Theorem 3.2.* If the target TPR lies on the boundary of D , then by Proposition 3.3, it is
 169 achieved by a tilting of the Bayes score without any randomization (i.e., $\beta = 0$; center figure of
 170 Fig. 2). This holds due to Assumption 2.3, because we may break ties arbitrarily without affecting
 171 TPR, since the set of tied scores (finite union of $(k - 2)$ -d subspaces) has (Lebesgue) measure zero.

172 Otherwise, and generally, there must exist $t' \in \partial D$ and $\beta \in [0, 1]$ s.t. t can be written as a linear
 173 combination of $t = \beta q + (1 - \beta)t'$. This is simply because $q \in \Delta_k \subseteq D$, and the line connecting q
 174 and t must intersect ∂D at some point t' (right figure of Fig. 2). Since the TPR of the input-agnostic
 175 randomization according to $\text{Multinoulli}(q)$ equals q , and t' is achieved by a tilting of the score per
 176 Proposition 3.3, their β -mixture achieves the target TPR t by linearity. \square

177 3.2 Deriving From Any Score Function

178 The post-processing method described in the previous section, which only requires unlabeled data
 179 (X, A) , yields the optimal α -fair classifier when applied to Bayes scores f^* . Yet, in practice, there

Algorithm 1 Post-Process Score Function for α -TPR parity

- 1: **Input:** $\alpha \in [0, 1]$, $q_1, \dots, q_m \in \Delta_k$, score function $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$, distribution $\mu^{X,A}$
 - 2: $\tilde{D}_a := \{\widetilde{\text{TPR}}_a(h) \mid h : \mathcal{X} \rightarrow \mathcal{Y} \text{ (probabilistic)}\} \quad \triangleright \text{Eq. (5), induced TPR feasible region}$
 - 3: $\tilde{t}_1, \dots, \tilde{t}_m \leftarrow \arg \max_{\tilde{t}_1 \in \tilde{D}_1, \dots, \tilde{t}_m \in \tilde{D}_m} \mathcal{U}(\tilde{t}_1, \dots, \tilde{t}_m)$ s.t. $\|\tilde{t}_a - \tilde{t}_{a'}\|_\infty \leq \alpha, \forall a, a' \in [m]$
 $\triangleright \text{utility-maximizing fair TPRs}$
 - 4: **for** $a = 1$ **to** m **do**
 - 5: Find $h_a, \beta_a \in [0, 1]$ s.t. $\widetilde{\text{TPR}}_a(h_a) \in \partial \tilde{D}_a$ and $\tilde{t}_a = (1 - \beta_a)\widetilde{\text{TPR}}_a(h_a) + \beta_a q_a$
 - 6: Find $\lambda_a \in \mathbb{R}^k$ s.t. $h_a(x) \in \arg \max_{y'} (\lambda_a)_{y'} \cdot f_a(x)_{y'}, \forall x \in \text{supp}(\mu_a^X)$
 - 7: **end for**
 - 8: **Return:** $(x, a) \mapsto \arg \max_{y'} (\lambda_a)_{y'} \cdot f_a(x)$ w.p. $1 - \beta_a$, and y w.p. $\beta_a \cdot (q_a)_y$ for each $y \in [k]$
-

180 is the concern that Bayes score functions could be arbitrarily complex and are often not exactly
 181 learnable due to limited data or computational constraints [34].

182 Nonetheless, our method is still applicable to arbitrary (approximations to the Bayes) score functions
 183 $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$ for deriving classifiers that are approximately fair and optimal, by treating them
 184 *as if they were Bayes optimal* (Algorithm 1). Where, the only tweak we made is replacing the
 185 ground-truth TPRs and feasible regions (which are unknown without access to the Bayes score) by
 186 approximations *induced* by f , i.e.,

$$\tilde{D}_a := \left\{ \widetilde{\text{TPR}}_a(h) \in [0, 1]^k \mid h : \mathcal{X} \rightarrow \mathcal{Y} \text{ (probabilistic)} \right\}, \quad (5)$$

187 where

$$\widetilde{\text{TPR}}_a(h)_y := \frac{1}{\tilde{p}_{ay}} \int_{x \in \mathcal{X}} f_a(x)_y \mathbb{P}(h(x) = y) d\mu^{X,A}(x, a), \quad \tilde{p}_{ay} := \int_{x \in \mathcal{X}} f_a(x)_y d\mu^{X,A}(x, a). \quad (6)$$

188 It is not hard to show that they are equal to their ground-truth counterparts when $f = f^*$.

189 We may control and minimize the suboptimalities of the classifier returned from Algorithm 1 by
 190 performing *group-wise distribution calibration* to the score function f (using labeled data (X, A, Y)):

191 **Definition 3.4** (Distribution Calibration). A score R is said to be (group-wise) distribution calibrated
 192 if $\mathbb{P}(Y = y \mid R = s) = s_y, \forall s \in \Delta_k, y \in [k]$ (resp. $\mathbb{P}(Y = y \mid R = s, A = a) = s_y, \forall a \in [m]$).

193 Distribution calibration is a multi-class generalization of the original definition of calibration for
 194 binary predictors [15, 32], requiring the predicted score to match the underlying class distribution
 195 conditioned on the score across all classes, not just the most confident one [22]. Although this
 196 definition is convenient to work with mathematically, it could be difficult to achieve in practice. In the
 197 proof of Theorem 3.5, we relax it to a recently proposed notion of *decision calibration* [45] (w.r.t. the
 198 set of all tiltings; derived from *multicalibration* [24]), which could be achieved in polynomial time.

199 **Theorem 3.5.** *Let $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$ be a score function, and $h : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{Y}$ the (probabilistic)
 200 classifier derived from f using Algorithm 1. Then under Assumption 2.3, for any group-wise calibrated
 201 reference score function $\bar{f} : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$,*

$$|\bar{\mathcal{U}} - \mathcal{U}(h)| \leq \sum_{a \in [m], y \in [k]} 3v_y \epsilon_{ay}, \quad \Delta_{\text{TPR}}(h) \leq \alpha + \max_{a \in [m], y \in [k]} \frac{4\epsilon_{ay}}{p_{ay}},$$

202 where $p_{ay} := \mathbb{P}_\mu(A = a, Y = y)$, v is from the utility function in Definition 2.2, $\bar{\mathcal{U}}$ denotes the utility
 203 achieved by the optimal α -fair classifier derived from the calibrated reference \bar{f} , and

$$\epsilon_{ay} := \mathbb{E}[|\bar{f}_a(X)_y - f_a(X)_y| \cdot \mathbb{1}[A = a]]$$

204 is the $L^1(\mu)$ difference between f and the calibrated reference \bar{f} on group a and class y .

205 We draw two conclusions from this result. First, by using the Bayes score function f^* as the reference,
 206 it states that the suboptimality of the derived classifier when $f \neq f^*$ is upper bounded by the
 207 difference between the approximate scores and the ground-truth; this answers the question raised
 208 in [2] regarding the impact of base model inaccuracies. Second, if f satisfies calibration, then by
 209 using itself as the reference, the result guarantees that the classifier derived using Algorithm 1 exactly
 210 achieves the desired level of fairness, and is optimal among all fair classifiers derived from f (which
 211 cannot be improved without labeled data).

212 **4 Finite-Sample Algorithm and Guarantees**

213 We instantiate the post-processing method above for TPR parity to the case where we do not have
 214 access to the distribution $\mu^{X,A}$ but only samples drawn from it (i.e., to perform estimation), and
 215 analyze the sample complexity.

216 **Assumption 4.1.** We have n i.i.d. (unlabeled) samples of (X, A) , which are independent of the score
 217 function f being post-processed.

218 Denote the number of samples from group a by n_a , and the samples themselves by $(x_{a,i})_{i \in [n_a]}$.

219 **4.1 Algorithm**

220 We adapt Algorithm 1 to handle finite samples by replacing \tilde{D}_a and \mathcal{U} with their empirical counterparts
 221 (essentially calling it with the empirical distribution $\hat{\mu}^{X,A}$ formed by the samples as the argument),
 222 and implement the optimization problems on Lines 3, 5 and 6 using linear programs.

223 **Step 1** (Finding Utility-Maximizing Fair TPRs). The empirical induced feasible region of TPRs, \hat{D}_a ,
 224 can be computed via evaluating the TPRs of all (probabilistic) classifiers acting on the samples—by
 225 representing them using $n_a \times k$ lookup tables (each row gives the probabilities of the random class
 226 assignment on the corresponding sample):

$$\hat{D}_a := \left\{ \widehat{\text{TPR}}_a(\gamma_a) \mid \gamma_a \in \mathbb{R}_{\geq 0}^{n_a \times k}, \sum_{y \in [k]} (\gamma_a)_{i,y} = 1, \forall i \in [n_a] \right\},$$

227 where

$$\widehat{\text{TPR}}_a(\gamma)_y := \frac{1}{n \hat{p}_{ay}} \sum_{i \in [n_a]} f_a(x_{a,i})_y \cdot (\gamma_a)_{i,y}, \quad \hat{p}_{ay} := \frac{1}{n} \sum_{i \in [n_a]} f_a(x_{a,i})_y$$

228 (cf. Line 2 and Eqs. (5) and (6)). Note that \hat{D}_a is a polygon, as it is specified by linear constraints.

229 To obtain the utility-maximizing fair TPR \hat{t}_a 's, we take the empirical maximizer subject to the α -TPR
 230 constraint via solving a linear program (cf. Line 3 and Eqs. (3) and (4)):

$$\text{LP1}(\alpha) : \max_{\hat{t}_1 \in \hat{D}_1, \dots, \hat{t}_m \in \hat{D}_m} \widehat{\mathcal{U}}(\hat{t}_1, \dots, \hat{t}_m) \quad \text{s.t.} \quad \|\hat{t}_a - \hat{t}_{a'}\|_\infty \leq \alpha, \forall a, a' \in [m],$$

231 where $\widehat{\mathcal{U}}(\hat{t}_1, \dots, \hat{t}_m) := \sum_{a,y} v_y \hat{p}_{ay}(\hat{t}_a)_y$ is the empirical utility.

232 **Step 2** (Obtaining Fair Classifier of Desired Form). The next step is finding a classifier that achieves
 233 \hat{t}_a 's on the empirical distribution, i.e., Lines 5 and 6. To implement Line 5, note that another way of
 234 approaching this problem is to realize that among all eligible (β_a, h_a) -pairs, the h_a associated with
 235 the maximum β_a value must satisfy $\widehat{\text{TPR}}_a(h_a) \in \partial \hat{D}_a$ (otherwise, a contradiction can be reached
 236 using the fact that $\tilde{D}_a \subseteq [0, 1]^k$ is compact; also see the right figure of Fig. 2). Combined with the
 237 strategy above of representing classifiers using lookup tables, we get the following linear program:

$$\text{LP2}(t, q) : \max_{\gamma, \beta} \beta \quad \text{s.t.} \quad t = (1 - \beta) \widehat{\text{TPR}}(\gamma) + \beta q \quad \text{and} \quad \gamma \in \mathbb{R}_{\geq 0}^{n \times k}, \sum_{y \in [k]} \gamma_{i,y} = 1, \forall i \in [n].$$

238 Finally, on Line 6, we find a tilting λ_a s.t. after coordinate-wise multiplied by the scores, the argmax
 239 class assignment has nonzero probability according to the classifier γ_a found in the preceding step:

$$\text{LP3}(\gamma) : \min_{\lambda} 0 \quad \text{s.t.} \quad \lambda_y f(x_i)_y \geq \lambda_{y'} f(x_i)_{y'} \quad \forall i \in [n], y, y' \in [k], \gamma_{i,y} > 0.$$

240 The feasible set of this problem is nonempty by Proposition 3.1, because we are *treating* f as if it
 241 were the Bayes score function, and the empirical distribution $\hat{\mu}^{X,A}$ as the population.

242 All combined, our algorithm involves solving $(2m + 1)$ linear programs, where LP1 is the dominating
 243 one with $O(nk)$ variables and constraints; solving which (to near-optimality) takes, e.g., $\tilde{O}(\text{poly}(nk))$
 244 time using interior point methods [33].

245 **4.2 Sample Complexity**

246 Thanks to the low function complexity of post-processing maps used in our algorithm to derive
 247 classifiers (Eq. (1)), it enjoys the following efficient sample complexity:

248 **Theorem 4.2.** *Let $f : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$ be a score function, and $h : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{Y}$ the (probabilistic)
 249 classifier derived from f using Algorithm 1 with the empirical distribution formed by samples
 250 from Assumption 4.1 as the argument. Then under Assumption 2.3, for any group-wise calibrated
 251 (Definition 3.4) reference score function $\bar{f} : \mathcal{X} \times \mathcal{A} \rightarrow \Delta_k$, and $n \geq \Omega(\max_{a,y} \ln(mk/\delta)/p_{ay})$,*

$$|\bar{U} - U(h)| \leq O\left(\sum_{a \in [m], y \in [k]} v_y \left(\sqrt{\frac{kp_{ay}}{n}} \ln \frac{mk}{\delta} + \frac{k}{n} + \epsilon_{ay}\right)\right),$$

$$\Delta_{\text{TPR}}(h) \leq \alpha + O\left(\max_{a \in [m], y \in [k]} \left(\sqrt{\frac{k}{np_{ay}} \ln \frac{mk}{\delta}} + \frac{k}{np_{ay}} + \frac{\epsilon_{ay}}{p_{ay}}\right)\right),$$

252 where \bar{U} denotes the utility achieved by the optimal α -fair classifier derived from the calibrated
 253 reference \bar{f} , and $\epsilon_{ay} := \mathbb{E}[|\bar{f}_a(X)_y - f_a(X)_y| \cdot \mathbb{1}[A = a]]$.

254 The bound consists of a calibration error ϵ_{ay} as discussed in the remarks of Theorem 3.5, an estimation
 255 error from applying uniform convergence (the Natarajan dimension of the set of tiltings is $O(k)$), and
 256 a k/n term that comes from the disagreement over class assignments on the samples between the
 257 (deterministic) tilting found on Line 6 and the (probabilistic) classifier on Line 5 due to tie-breaking.

258 **5 Experiments**

259 We evaluate Algorithm 1 for reducing TPR disparity on benchmark datasets, and demonstrate its
 260 effectiveness compared to existing post-processing as well as in-processing bias mitigation methods.

261 **Datasets.** The first task is income prediction, for which, we use the ACSIncome dataset [18]—an
 262 extension of the UCI Adult dataset [27] with much more examples (1.6 million vs. 30,162), allowing
 263 us to compare methods confidently. We consider a binary setting where the sensitive attribute is
 264 gender and the target is whether the income is over \$50k, as well as a multi-group multi-class setting
 265 with five race categories and five income buckets. The second is text classification, of identifying
 266 occupations (28 in total) from biographies in the BiasBios dataset [14]; sensitive attribute is gender.

267 **Baselines and Setup.** The main baseline is FairProjection [2]—the only post-processing algo-
 268 rithm applicable for multi-class TPR parity to our knowledge.³ In the binary setting, we also compare
 269 to RejectOption [25]. To demonstrate the deficiencies of existing methods at reducing TPR dispar-
 270 ity, we additionally include in-processing results using Reductions [1] and Adversarial [41].^{4,5}

271 On each task, we first create a pre-training split from the dataset and train a linear logistic regression
 272 scoring model (with isotonic calibration and five-fold cross-validation as implemented in scikit-
 273 learn [37, 38, 28]), then randomly split the remaining data for post-processing and testing with
 274 10 different seeds and aggregate the results (the pre-trained model remains the same). For in-
 275 processing, we use the same splits but merge the pre-training and post-processing data for training.
 276 On BiasBios, linear logistic regression is performed on the embeddings of the biographies computed
 277 by a previously fine-tuned BERT model [17] (in other words, head-tuning). Additional details
 278 including hyperparameters are included in the appendix.

279 **Results.** In Fig. 3, we plot the tradeoff curves from varying the fairness tolerance (α for our
 280 method). Our method is consistently the most effective at minimizing TPR disparity, particularly
 281 under multi-class settings, where existing algorithms only manage to partially reduce Δ_{TPR} (and
 282 at a greater cost to accuracy when using FairProject and RejectOption). It also outperforms

³We use the authors’ code, where TPR parity is equivalent to the meo constraint. The results from using the KL divergence variant is included, which are better than the cross-entropy variant in our experiments.

⁴Although Reductions is extended to multi-class by Yang et al. [36], an implementation was not provided.

⁵The implementation (with minor modifications) in the AIF360 library is used for the latter methods [5].

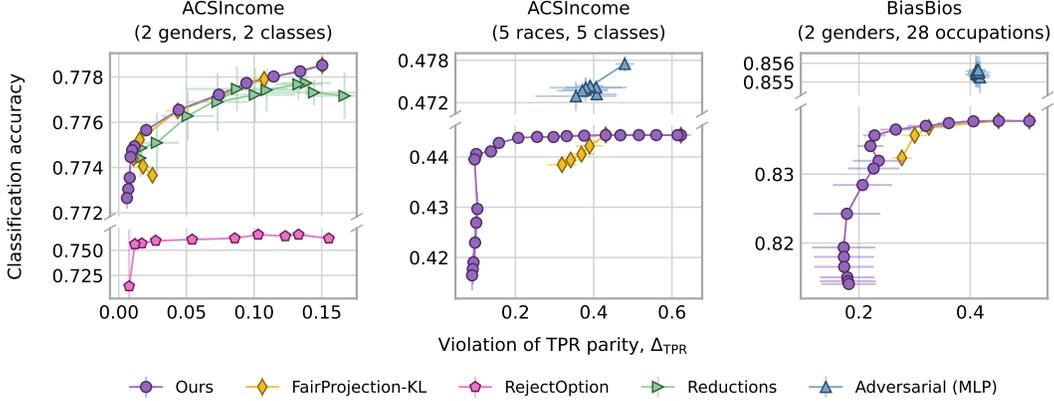


Figure 3: Tradeoff curves between accuracy and Δ_{TPR} (Eq. (2)). The base model is logistic regression (except for Adversarial, which uses a feedforward network). Error bars indicate the standard deviation over 10 runs with different random dataset splits. Running time is reported in the appendix.

283 the in-processing Reductions on binary ACSIncome, and Adversarial in terms of Δ_{TPR} , which,
 284 although enjoys higher accuracies because of the use of the more expressive feedforward networks
 285 as the prediction model, fails to reduce TPR parity. Sharper drops in accuracies are observed when
 286 applying our method with small α settings, e.g., 0.001 to 0.0001. We saw this happen when the
 287 randomized component in Eq. (1b) is activated (i.e., $\beta > 0$), meaning that Line 3 finds a fair TPR
 288 that lies in the interior of the feasible region of the better-performing group in order to match the
 289 feasible TPR on the worse-performing one(s). Hence the drop is expected because utility is being
 290 sacrificed to achieve TPR parity.

291 Although our method greatly reduces TPR disparity, there remains a gap to reaching $\Delta_{\text{TPR}} = 0$,
 292 especially on tasks with more classes (i.e., BiasBios, where a higher variance is also observed).
 293 While this could be due to miscalibration, or potentially a violation of Assumption 2.3, the main
 294 reason is suspected to be insufficient sample size. Recall from Theorem 4.2 that the sample complexity
 295 for Δ_{TPR} scales as $\tilde{O}(\sqrt{k/np_{ay}})$ in the worse-case (a, y) , which is itself at least $\tilde{O}(\sqrt{mk^2/n})$.
 296 Thus, learning generalizable classifiers that satisfy TPR parity under more groups and classes is much
 297 harder in terms of data requirement (and by extension, computing resource).

298 Lastly, we emphasize the necessity of group-wise calibration for achieving low Δ_{TPR} , as the
 299 definition of the criterion involves conditioning on the true label (it is also reflected by the calibration
 300 error term ϵ_{ay} in Theorem 4.2). In an ablation study in the appendix, a larger (minimum achievable)
 301 Δ_{TPR} is observed when no efforts are made to calibrate the scoring model. It is therefore necessary
 302 for model vendors to provide accurate uncertainty quantifications, and for practitioners building fair
 303 classifiers to verify and improve calibration.

304 6 Conclusions and Limitations

305 We described a post-processing method for reducing TPR disparity for equal opportunity in multi-class
 306 classification, and demonstrated its performance in comparison to existing algorithms on benchmarks
 307 datasets, especially when the number of classes is large. We analyzed the sample complexity of our
 308 method, and established its optimality under model calibration.

309 The effectiveness of our method at reducing TPR disparity is largely contributed to the tailored
 310 analysis, although it limits our method to this fairness notion only. Some use cases may demand
 311 equalized odds ($\hat{Y} \perp A \mid Y$) beyond TPR parity ($\mathbb{1}[\hat{Y} = Y] \perp A \mid Y$), which is a more stringent
 312 criterion: TPR parity only needs to match the main diagonal of the (conditional) confusion matrix
 313 across groups, whereas equalized odds requires matching all k^2 entries. The design of efficient
 314 algorithms for achieving equalized odds remains an open problem.⁶

⁶We note that most (general-purpose) fairness algorithms, e.g., [2], are only evaluated for TPR parity but not equalized odds.

315 References

- 316 [1] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. A
317 Reductions Approach to Fair Classification. In *Proceedings of the 35th International Conference*
318 *on Machine Learning*, pages 60–69, 2018.
- 319 [2] Wael Alghamdi, Hsiang Hsu, Haewon Jeong, Hao Wang, P. Winston Michalak, Shahab Asoodeh,
320 and Flavio P. Calmon. Beyond Adult and COMPAS: Fair Multi-Class Prediction via Information
321 Projection. In *Advances in Neural Information Processing Systems*, 2022.
- 322 [3] Solon Barocas and Andrew D. Selbst. Big Data’s Disparate Impact. *California Law Review*,
323 104(3):671–732, 2016.
- 324 [4] Solon Barocas, Moritz Hardt, and Arvind Narayanan. *Fairness and Machine Learning: Limita-*
325 *tions and Opportunities*. MIT Press, 2023.
- 326 [5] Rachel K. E. Bellamy, Kuntal Dey, Michael Hind, Samuel C. Hoffman, Stephanie Houde,
327 Kalapriya Kannan, Pranay Lohia, Jacquelyn Martino, Sameep Mehta, Aleksandra Mojsilovic,
328 Seema Nagar, Karthikeyan Natesan Ramamurthy, John Richards, Diptikalyan Saha, Prasanna
329 Sattigeri, Moninder Singh, Kush R. Varshney, and Yunfeng Zhang. AI Fairness 360: An
330 Extensible Toolkit for Detecting, Understanding, and Mitigating Unwanted Algorithmic Bias,
331 2018. *arxiv:1810.01943 [cs.AI]*.
- 332 [6] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On
333 the Dangers of Stochastic Parrots: Can Language Models Be Too Big? 🦜. In *Proceedings*
334 *of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623,
335 2021.
- 336 [7] Richard Berk, Hoda Heidari, Shahin Jabbari, Michael Kearns, and Aaron Roth. Fairness in
337 Criminal Justice Risk Assessments: The State of the Art. *Sociological Methods & Research*, 50
338 (1):3–44, 2021.
- 339 [8] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man
340 is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. In
341 *Advances in Neural Information Processing Systems*, volume 29, 2016.
- 342 [9] Joy Buolamwini and Timnit Gebru. Gender Shades: Intersectional Accuracy Disparities
343 in Commercial Gender Classification. In *Proceedings of the 2018 Conference on Fairness,*
344 *Accountability, and Transparency*, pages 77–91, 2018.
- 345 [10] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building Classifiers with Independency
346 Constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pages 13–18,
347 2009.
- 348 [11] Flavio P. Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and
349 Kush R. Varshney. Optimized Pre-Processing for Discrimination Prevention. In *Advances in*
350 *Neural Information Processing Systems*, volume 30, 2017.
- 351 [12] Simon Caton and Christian Haas. Fairness in Machine Learning: A Survey, 2020.
352 *arxiv:2010.04053 [cs.LG]*.
- 353 [13] Jeffrey Dastin. Amazon scraps secret AI recruiting tool that showed bias against women.
354 *Reuters*, oct 2018. URL [https://www.reuters.com/article/us-amazon-com-jobs-](https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G)
355 [automation-insight-idUSKCN1MK08G](https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G).
- 356 [14] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexan-
357 dra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. Bias in
358 Bios: A Case Study of Semantic Representation Bias in a High-Stakes Setting. In *Proceedings*
359 *of the 2019 ACM Conference on Fairness, Accountability, and Transparency*, pages 120–128,
360 2019.
- 361 [15] Morris H. DeGroot and Stephen E. Fienberg. The Comparison and Evaluation of Forecasters.
362 *Journal of the Royal Statistical Society. Series D (The Statistician)*, 32(1/2):12–22, 1983.

- 363 [16] Christophe Denis, Romuald Elie, Mohamed Hebiri, and François Hu. Fairness guarantee in
364 multi-class classification, 2023. *arxiv:2109.13642 [math.ST]*.
- 365 [17] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
366 Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019*
367 *Conference of the North American Chapter of the Association for Computational Linguistics:*
368 *Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- 369 [18] Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring Adult: New Datasets
370 for Fair Machine Learning. In *Advances in Neural Information Processing Systems*, 2021.
- 371 [19] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness
372 Through Awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science*
373 *Conference*, pages 214–226, 2012.
- 374 [20] Executive Office of the President. Big Data: A Report on Algorithmic Systems, Opportunity,
375 and Civil Rights. The White House, 2016. URL <https://www.fdlp.gov/GPO/gpo90618>.
- 376 [21] Solenne Gaucher, Nicolas Schreuder, and Evgenii Chzhen. Fair learning with Wasserstein
377 barycenters for non-decomposable performance measures. In *Proceedings of The 26th Interna-*
378 *tional Conference on Artificial Intelligence and Statistics*, pages 2436–2459, 2023.
- 379 [22] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On Calibration of Modern Neural
380 Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages
381 1321–1330, 2017.
- 382 [23] Moritz Hardt, Eric Price, and Nathan Srebro. Equality of Opportunity in Supervised Learning.
383 In *Advances in Neural Information Processing Systems*, volume 29, 2016.
- 384 [24] Úrsula Hébert-Johnson, Michael P. Kim, Omer Reingold, and Guy N. Rothblum. Multicalibra-
385 tion: Calibration for the (Computationally-Identifiable) Masses. In *Proceedings of the 35th*
386 *International Conference on Machine Learning*, pages 1939–1948, 2018.
- 387 [25] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. Decision Theory for Discrimination-
388 Aware Classification. In *2012 IEEE 12th International Conference on Data Mining*, pages
389 924–929, 2012.
- 390 [26] Michael Kearns, Seth Neel, Aaron Roth, and Zhiwei Steven Wu. Preventing Fairness Gerryman-
391 dering: Auditing and Learning for Subgroup Fairness. In *Proceedings of the 35th International*
392 *Conference on Machine Learning*, pages 2564–2572, 2018.
- 393 [27] Ron Kohavi. Scaling Up the Accuracy of Naive-Bayes Classifiers: A Decision-Tree Hybrid. In
394 *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*,
395 pages 202–207, 1996.
- 396 [28] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,
397 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vander-
398 plas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard
399 Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*,
400 12(85):2825–2830, 2011.
- 401 [29] Preston Putzel and Scott Lee. Blackbox Post-Processing for Multiclass Fairness. In *Proceedings*
402 *of the Workshop on Artificial Intelligence Safety 2022*, volume 3087, 2022.
- 403 [30] Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. Null It Out:
404 Guarding Protected Attributes by Iterative Nullspace Projection. In *Proceedings of the 58th*
405 *Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256, 2020.
- 406 [31] Julien Rouzot, Julien Ferry, and Marie-José Huguet. Learning Optimal Fair Scoring Systems for
407 Multi-Class Classification. In *2022 IEEE 34th International Conference on Tools with Artificial*
408 *Intelligence*, 2022.

- 409 [32] Hao Song, Tom Diethe, Meelis Kull, and Peter Flach. Distribution Calibration for Regression.
410 In *Proceedings of the 36th International Conference on Machine Learning*, pages 5897–5906,
411 2019.
- 412 [33] Pravin M. Vaidya. Speeding-up linear programming using fast matrix multiplication. In *30th*
413 *Annual Symposium on Foundations of Computer Science*, pages 332–337, 1989.
- 414 [34] Blake Woodworth, Suriya Gunasekar, Mesrob I. Ohannessian, and Nathan Srebro. Learning
415 Non-Discriminatory Predictors. In *Proceedings of the 2017 Conference on Learning Theory*,
416 pages 1920–1953, 2017.
- 417 [35] Ruicheng Xian, Lang Yin, and Han Zhao. Fair and Optimal Classification via Post-Processing
418 Predictors. In *Proceedings of the 40th International Conference on Machine Learning*, 2023.
- 419 [36] Forest Yang, Mouhamadou Cisse, and Sanmi Koyejo. Fairness with Overlapping Groups. In
420 *Advances in Neural Information Processing Systems*, volume 33, pages 4067–4078, 2020.
- 421 [37] Bianca Zadrozny and Charles Elkan. Obtaining calibrated probability estimates from decision
422 trees and naive Bayesian classifiers. In *Proceedings of the Eighteenth International Conference*
423 *on Machine Learning*, pages 609–616, 2001.
- 424 [38] Bianca Zadrozny and Charles Elkan. Transforming Classifier Scores into Accurate Multiclass
425 Probability Estimates. In *Proceedings of the Eighth ACM SIGKDD International Conference*
426 *on Knowledge Discovery and Data Mining*, pages 694–699, 2002.
- 427 [39] Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P. Gummadi.
428 Fairness Beyond Disparate Treatment & Disparate Impact: Learning Classification without
429 Disparate Mistreatment. In *Proceedings of the 26th International Conference on World Wide*
430 *Web*, pages 1171–1180, 2017.
- 431 [40] Richard Zemel, Yu Ledell Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning Fair
432 Representations. In *Proceedings of the 30th International Conference on Machine Learning*,
433 pages 325–333, 2013.
- 434 [41] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating Unwanted Biases with
435 Adversarial Learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and*
436 *Society*, pages 335–340, 2018.
- 437 [42] Han Zhao and Geoffrey J. Gordon. Inherent Tradeoffs in Learning Fair Representations. *Journal*
438 *of Machine Learning Research*, 23(57):1–26, 2022.
- 439 [43] Han Zhao, Amanda Coston, Tameem Adel, and Geoffrey J. Gordon. Conditional Learning of
440 Fair Representations. In *International Conference on Learning Representations*, 2020.
- 441 [44] Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender Bias
442 in Coreference Resolution: Evaluation and Debiasing Methods. In *Proceedings of the 2018*
443 *Conference of the North American Chapter of the Association for Computational Linguistics:*
444 *Human Language Technologies, Volume 2 (Short Papers)*, pages 15–20, 2018.
- 445 [45] Shengjia Zhao, Michael P. Kim, Roshni Sahoo, Tengyu Ma, and Stefano Ermon. Calibrating
446 Predictions to Decisions: A Novel Approach to Multi-Class Calibration. In *Advances in Neural*
447 *Information Processing Systems*, 2021.