

---

# Group-Aware Threshold Adaptation for Fair Classification

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The fairness in machine learning is getting increasing attention, as its applications  
2 in different fields continue to expand and diversify. To mitigate the discriminated  
3 model behaviors between different demographic groups, we introduce a novel post-  
4 processing method to optimize over multiple fairness constraints through group-  
5 aware threshold adaptation. We propose to learn adaptive classification thresholds  
6 for each demographic group by optimizing the confusion matrix estimated from  
7 the probability distribution of a classification model output. As we only need  
8 an estimated probability distribution of model output instead of the classification  
9 model structure, our post-processing model can be applied to a wide range of  
10 classification models and improve fairness in a model-agnostic manner to ensure  
11 privacy. This even allows us to post-process existing fairness methods to further  
12 improve the trade-off between accuracy and fairness. Moreover, our model is  
13 efficient with low computational cost by alternating optimization and flexible with  
14 the optimization over multiple fairness constraints. We provide Pareto frontier to  
15 characterize fairness-accuracy trade-off. Also, we provide a theoretical analysis  
16 of the optimal thresholds obtained from our model in terms of both accuracy  
17 and fairness in classification. Experimental results demonstrate that our method  
18 outperforms state-of-the-art methods and obtains the result that is closest to the  
19 theoretical accuracy-fairness trade-off boundary.

## 20 1 Introduction

21 Machine learning is broadening its impact in various fields including autonomous driving, credit  
22 analysis, and job application screening. As a consequence, the role and importance of fairness in  
23 machine learning are emerging. However, recent models have been found to behave differently  
24 between demographic groups in favorable predictions. For example, it has been discovered that  
25 COMPAS, the criminal risk assessment software currently used to help pretrial release decisions,  
26 has biases between different races [4]. Specifically, blacks got higher risk scores predicted from the  
27 model than whites with similar profiles. Therefore, discrimination truly exists and resolving it in  
28 machine learning is very important and urgent because its direct and potential impact is growing  
29 tremendously.

30 However, obtaining fairness is not a trivial problem, because the data set itself will be biased when it  
31 is accumulated artificially. Simply removing or manipulating sensitive features (such as *race*, *gender*)  
32 from the data does not solve the bias, because there is indirect discrimination [19] or disparate  
33 treatment [1] due to the feature redundancy and relevance, which means sensitive information can be  
34 inferred from other features.

35 In order to alleviate discrimination from different perspectives, various quantitative measurements  
36 of group equity [7, 11, 2, 13] have been proposed. It has been proven that the pursuit of fairness is

37 subject to a trade-off between fairness and accuracy [14, 10], i.e., if we want to improve fairness, we  
38 need to sacrifice accuracy.

39 Moreover, Pleiss et al. [20] studied the trade-offs between fairness notions that cannot be satisfied  
40 at the same time. Therefore, recent works usually target at a certain fairness notion in different  
41 approaches such as pre-processing [6], in-processing [24], and post-processing [7] methods. However,  
42 these approaches suffer from the *lack of flexibility*, since it is difficult to adapt a fair model that  
43 is trained w.r.t. one certain fairness criterion so as to optimize over other fairness measures. If  
44 the fairness constraints change under some circumstances, traditional fairness models need to be  
45 re-trained from scratch, which is computationally demanding and sometimes inapplicable due to  
46 model settings. To overcome the limitations above, we propose a novel post-processing method to  
47 improve fairness in a model-agnostic manner. Our GSTAR (Group Specific Threshold Adaptation  
48 for fair classification) model learns adaptive classification thresholds for each demographic group  
49 in classification task for improving the trade-off between fairness and accuracy. Given an existing  
50 classification model, GSTAR approximates the probability distribution of the model output via  
51 maximum likelihood estimation and utilizes confusion matrix to quantify accuracy and fairness w.r.t.  
52 the group-aware classification thresholds. This allows us to: 1) prevent from burdening additional  
53 complexity or deteriorate the stability of the training process of the classifier; 2) integrate different  
54 fairness notions into one unified objective function; 3) easily adapt one pre-trained model to other  
55 fairness constraints. We summarize our contributions of this paper as follows:

- 56 1. We propose a novel post-processing method, GSTAR, which can learn group-aware thresh-  
57 olds to optimize the trade-off between fairness and accuracy in classification. We derive  
58 rigorous theoretical analysis on the trade-off in our model, and empirically show that GSTAR  
59 outperforms state-of-the-art methods.
- 60 2. With GSTAR, we can simultaneously optimize over multiple fairness constraints with a low  
61 computational cost. GSTAR does not require multiple iterations over data, instead, it takes  
62 *at most* one pass of data in training for fast computation.
- 63 3. GSTAR can be adapted to a wide range of classification models in a model-agnostic manner  
64 and can adapt an existing classification model from one fairness criterion to another without  
65 re-training the classifier.
- 66 4. We derive Pareto frontiers of our model for the fairness-accuracy trade-offs that contextualize  
67 the quality of fair classification.

## 68 2 Related Works

69 In order to achieve group fairness, which quantifies the discrimination among different sensitive  
70 groups, a diverse notion of fairness has been introduced. Equalized odds [7] enforce equality of true  
71 positive rates and false positive rates between different demographic groups. Pleiss et al. [20] relaxed  
72 equalized odds to satisfy the calibration. Demographic parity or disparate impact [1] suggests that a  
73 model is unbiased if the model prediction is independent of the protected attribute.

74 Among different fairness methods, post-processing techniques propose to improve fairness by mod-  
75 ifying the output of a black-box classifier. Hardt et al. [7] propose to ensure equalized odds by  
76 constraining the model output. Kamiran et al. [9] propose to give a favorable outcome to unprivileged  
77 and an unfavorable outcome to the privileged group when the confidence of the prediction is beyond a  
78 certain range. However, such *static* confidence window keeps the same regardless of the demographic  
79 group and is determined by grid search, so it is less efficient.

80 Threshold adjustment (a.k.a. thresholding) was introduced to improve the performance of *static*  
81 thresholds. In the literature, Menon et al. [18] prove that instance-dependent thresholding of the  
82 predictive probability function is the optimal classifier in cost-sensitive fairness measures. Also,  
83 when considering immediate utility, Corbett-Davies et al. [3] show that optimal algorithm is achieved  
84 from group-specific threshold which is determined by group statistics. However, to the best of our  
85 knowledge, the threshold adjustment approach has not been deeply studied that neither encompasses  
86 broad group fairness metrics nor describes an explicit method to achieve the threshold.

87 Trade-off between fairness and accuracy exists when we impose fairness constraint to a model. Recent  
88 studies [2, 25] prove that models targeting at such fairness notions conform to an information theoretic

89 lower bound on the joint error across different sensitive groups. Therefore, our work presents a  
 90 practical upper bound of the best achievable accuracy given the fairness constraints.

91 Moreover, trade-offs between different fairness notions also exist if one has to consider multiple  
 92 fairness criteria. Some of them are theoretically proven to be incompatible [6, 18, 14]. To express  
 93 and formulate fairness, recent work [10] utilize confusion matrix and propose least-square accuracy-  
 94 fairness optimization problem on multiple fairness notions, and categorize the trade-offs between the  
 95 fairness notions.

96 Here, our work is the most related to the post-processing methods [7, 10]. Hardt et al. [7] propose  
 97 a post-processing method that utilizes the mixing rate to meet the equalized odds. Ours is similar  
 98 to Hardt et al. [7] in the manner that achieving group-wise threshold from the feasible region that  
 99 is geometrically generated by the intersection between the receiver operating characteristic (ROC)  
 100 curves conditioned on sensitive feature. Ours differ from [7] by generalizing the concept beyond  
 101 equalized odds to other multiple fairness constraints into consideration. FACT [10] utilizes a single  
 102 point (static) from the classifier to be post-processed as a reference which does not fully utilize the  
 103 classifier for the post-processing. In contrast, by approximating the distribution of the continuous  
 104 predicted logits, our GSTAR model enables a larger feasible region than [10] with a better fairness-  
 105 accuracy trade-off. We validate the improvement in trade-off via both theoretical and empirical results.  
 106 It is notable that these related methods [7, 10] can be considered as a special case of GSTAR.

### 107 3 GSTAR for Fair Classification

#### 108 3.1 Motivation

109 Consider a binary classification problem with a binary sensitive feature, such that the sensitive feature  
 110  $A \in \{0, 1\}$  and label  $Y \in \{0, 1\}$ . In general, for a given data  $X$ , a binary classification model  
 111 outputs an unnormalized logit  $h(X) \in \mathbb{R}$  with the class label probability  $p(X) = \sigma(h(X)) \in [0, 1]$ ,  
 112 where  $\sigma$  is an activation function (sigmoid function in logistic regression and neural network). It is  
 113 not necessary to calculate  $p$  in a classification model, e.g. support vector machines directly use the  
 114 positiveness/negativeness of logit  $h(X)$  to determine classification outcome. For traditional models,  
 115 we use a cut-off threshold  $\theta_h = 0$  for  $h(X)$  (i.e.,  $\theta_p = \sigma(0) = 0.5$  for  $p(X)$ ) in classification, such  
 116 that the predicted label is determined by  $\hat{Y} = \mathbb{I}\{h(X) \geq \theta_h\}$ . In the following context, unless  
 117 otherwise mentioned, we use  $\theta$  to refer to the threshold  $\theta_h$  on logit  $h$  since it is applicable to a  
 118 wider range of classification models, and the corresponding threshold on label probability  $\theta_p$  can be  
 119 easily inferred from the threshold on logit  $h$ . Traditional models use the same cut-off threshold  $\theta$  for  
 120 different demographic groups. However, since the distribution of logits  $h$  in different demographic  
 121 groups can be different, using the same threshold  $\theta$  brings biased classification.

122 In Figure 1, we show a real-world example of image classification on CelebA dataset with  
 123 ResNet50 [8] to show that the default setting of classification thresholds affects both accuracy  
 124 and fairness in classification. The goal of this classification example is to predict the image of a  
 125 person is whether attractive or not, and consider sensitive attribute as gender. This can be generalized  
 126 to different sensitive attributes such as age or race [22, 16]. We can observe an obvious difference in  
 127 the distribution of logit  $h$  between two gender groups. In this case, if we use a unified classification  
 128 threshold  $\theta_1 = \theta_0 = 0$ , it naturally brings a difference in the true positive rate and true negative  
 129 rate between two gender groups, thus renders bias in classification. Instead, we observe that the  
 130 optimal group-specific threshold obtained from GSTAR ( $\theta_1^* > \theta_1$ , and  $\theta_0^* < \theta_0$ ) can adapt to such  
 131 discrepancy in distribution between two demographic groups to improve both fairness and accuracy.

#### 132 3.2 Group-Aware Classification Thresholds

133 Given an existing classification model and a sensitive attribute  $a$ , we can denote true positive rate  
 134 ( $TP_a$ ), false positive rate ( $FP_a$ ), true negative rate ( $TN_a$ ), and false negative rate ( $FN_a$ ) in the confusion  
 135 matrix. Most fairness notions can be represented with entries in the confusion matrix. For instance,  
 136 Equal Opportunity (EOp) [7] requires  $TP_0 = TP_1$ , and Demographic Parity (DP) [1] requires

$$\frac{TP_1 n_{11} + FP_1 n_{01}}{N_1} = \frac{TP_0 n_{10} + FP_0 n_{00}}{N_0},$$

137 where  $n_{ya}$  denotes the number of samples in the subset  $\{Y = y, A = a\}$ ,  $N_a = \sum_y n_{ya}$  denotes the  
 138 number of samples in  $\{Y = y\}$ , and  $N = \sum_{y,a} n_{ya}$  is the total number of samples.

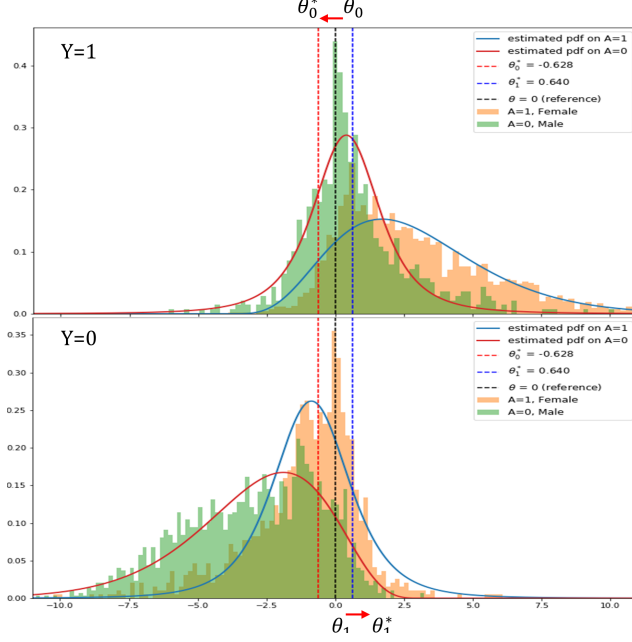


Figure 1: Histograms of logit  $h$  distribution from logistic regression on CelebA data, where logit  $h$  is used to determine the predicted label  $\hat{Y} = \mathbb{I}\{h(X) \geq \theta\}$ , and  $\theta$  is the classification threshold. The top plot is for positive samples ( $Y = 1$ , attractive), and the bottom plot for negative samples ( $Y = 0$ , unattractive). In each plot, yellow/green bars represent the distribution of logit  $h$  in different sensitive groups, and blue/red curves are estimated probability density functions of logit  $h$  in different sensitive groups.  $\theta_0 = \theta_1 = 0$  (black dashed line) are the default classification thresholds, that are identical for  $A = 0$  and  $A = 1$  groups. The default thresholds result in biased prediction towards the unprivileged group  $A = 0$  due to the different logit  $h$  distributions in different sensitive groups.  $\theta_0^*$  (red dashed line),  $\theta_1^*$  (blue dashed line) are group-aware thresholds from GSTAR for each sensitive group.

139 Consider the group-aware classification threshold  $\theta = (\theta_1, \theta_0)^\top$ , where  $\theta_a$  is the classification  
 140 threshold for sensitive group  $A = a$ . We can formulate the entries in the confusion matrix w.r.t.  $\theta$  as  
 141 below:

$$\begin{aligned}
 \text{TP}_a(\theta_a) &\approx 1 - \int_{-\infty}^{\theta_a} f_{1a}(x) dx, & \text{FN}_a(\theta_a) &\approx \int_{-\infty}^{\theta_a} f_{1a}(x) dx, \\
 \text{FP}_a(\theta_a) &\approx 1 - \int_{-\infty}^{\theta_a} f_{0a}(x) dx, & \text{TN}_a(\theta_a) &\approx \int_{-\infty}^{\theta_a} f_{0a}(x) dx,
 \end{aligned} \tag{1}$$

142 where  $f_{ya}(x)$  is an estimated parametric probability density function of the distribution of output  
 143 logit  $h$  in the subset  $\{Y = y, A = a\}$ . Here, we consider gamma, Student's t, and normal distribution  
 144 as the candidates for the estimated distribution, and select the one that has the maximum likelihood  
 145 with the output distribution. Without loss of generality, this can be generalized with other parametric  
 146 probability density function based on the needs or prior knowledge.

147 Then, we formulate the fairness-constrained classification problem with the objective of minimizing  
 148 classification error into a least-squared optimization problem. We denote our objective function  
 149 as  $\mathcal{L}(\theta)$  which consists of the performance loss  $\mathcal{L}_{per}(\theta)$  and fairness loss  $\mathcal{L}_{fair}(\theta)$ .  $\mathcal{L}_{per}(\theta)$  and  
 150  $\mathcal{L}_{fair}(\theta)$  measures the error in performance and fairness respectively that are represented with the  
 151 entries of the confusion matrix. In other words, our goal is to minimize the objective function  $\mathcal{L}(\theta)$   
 152 as below:

$$\mathcal{L}(\theta) = \mathcal{L}_{per}(\theta) + \lambda \mathcal{L}_{fair}(\theta), \tag{2}$$

153 where  $\lambda$  is a hyperparameter that determines how much fairness is enforced in the optimization.

154 The performance error  $\mathcal{L}_{per}(\theta)$  can be written as

$$\mathcal{L}_{per}(\theta) = \left( \frac{n_{01}}{N} \text{FP}_1(\theta_1) + \frac{n_{11}}{N} \text{FN}_1(\theta_1) + \frac{n_{00}}{N} \text{FP}_0(\theta_0) + \frac{n_{10}}{N} \text{FN}_0(\theta_0) \right)^2. \tag{3}$$

155 As for  $\mathcal{L}_{fair}(\theta)$ , it can be formulated to any fairness metrics that are expressible with confusion matrix.  
 156 For instance, when we impose EOp ( $TP_1 = TP_0$ ) and predictive equality (PE) ( $FP_1 = FP_0$ ) [2], we  
 157 can get the corresponding  $\mathcal{L}_{fair}(\theta)$  by summing over the least squared form of each constraint. Also,  
 158 satisfying EOp and PP is equivalent to satisfying Equalized Odds (EOd) [7], This can be formulated  
 159 in our  $\mathcal{L}_{fair}$  as

$$\begin{aligned}\mathcal{L}_{fair}^{EOd}(\theta) &= \mathcal{L}_{fair}^{EOp}(\theta) + \mathcal{L}_{fair}^{PP}(\theta) \\ &= (TP_1(\theta_1) - TP_0(\theta_0))^2 + (FP_1(\theta_1) - FP_0(\theta_0))^2.\end{aligned}\quad (4)$$

160 Note that a lower  $\mathcal{L}_{fair}$  value indicates a fairer threshold. When  $\mathcal{L}_{fair}^{EOd}(\theta) = 0$ , we can interpret as  
 161 the  $\theta$  satisfies the perfect EOd fairness.

162 Similar to (4), we can enforce multiple fairness constraints by summing over the least squared of  
 163 each metric with different weight constant  $\lambda$  to each fairness constraints if needed.

164 Also, it is notable that compared to the recent paper [10] that enforces fairness through confusion  
 165 tensor, our formulation of fairness in  $\mathcal{L}_{fair}(\theta)$  represents a direct notion of fairness metrics and  
 166 improves the measures that allows us to achieve better performance and Pareto frontiers that is shown  
 167 in Section 4.2 and Figure 2. For example,  $\mathbf{A}_{EOd}$  in the paper is calculated as  $M_1EOp + M_0PE$ , where  
 168  $M_y = n_{y0} + n_{y1}$ , such that EOd is a weighted sum of EOp and PE with weights being the number of  
 169 samples in each class. In this expression, the imbalance between the two fairness criteria will grow as  
 170 the degree of imbalance in the data increases. In contrast, our formulation expresses the constraints  
 171 as the exact notion of each metric that is not biased by the statistics of the dataset and we observe  
 172 improved Pareto frontier as in Figure 2.

173 We propose to optimize our threshold  $\theta$  with alternating optimization method. Here we take EOp  
 174 constraint as an example to show the alternating optimization steps, then  $\mathcal{L}_{fair}(\theta)$  can be written as

$$\mathcal{L}_{fair}^{EOp}(\theta) = (TP_1(\theta_1) - TP_0(\theta_0))^2. \quad (5)$$

175 **The first step** is to fix  $\theta_0$  and update  $\theta_1$ . We can approximate the terms that are related to  $\theta_1$  (e.g.,  
 176  $TP_1, FP_1, TN_1, FN_1$ ) in (1) with first-order Taylor expansion at  $\theta_1^{\tau-1}$ . For example,

$$TP_1(\theta_1) \approx TP_1(\theta_1^{\tau-1}) + \left. \frac{\partial TP_1}{\partial \theta_1} \right|_{\theta_1 = \theta_1^{\tau-1}} (\theta_1 - \theta_1^{\tau-1}) \quad (6)$$

177 From (1), we can easily derive that

$$\begin{aligned}TP_1(\theta_1^{\tau-1}) &= 1 - \int_{-\infty}^{\theta_1^{\tau-1}} f_{11}(x) dx, \\ \frac{\partial TP_1}{\partial \theta_1} &= -f_{11}(\theta_1^{\tau-1}).\end{aligned}\quad (7)$$

178 Similarly, we can find the first order Taylor expansion of  $FP_1, FN_1$ , and  $TN_1$ . Then, the update of  $\theta_1$   
 179 w.r.t. (2) can be approximated with the following minimization problem w.r.t.  $\Delta_1$

$$\Delta_1^\tau := \underset{\Delta_1}{\operatorname{argmin}} (\eta^\tau + \alpha^\tau \Delta_1)^2 + \lambda (\epsilon^\tau + \beta^\tau \Delta_1)^2, \quad (8)$$

180 where  $\Delta_1 = \theta_1 - \theta_1^{\tau-1}$  and

$$\begin{aligned}\alpha_1^\tau &= \frac{n_{11}}{N} f_{11}(\theta_1^{\tau-1}) - \frac{n_{01}}{N} f_{01}(\theta_1^{\tau-1}), \\ \beta_1^\tau &= -f_{11}(\theta_1^{\tau-1}), \\ \eta_1^\tau &= \int_{-\infty}^{\theta_1^{\tau-1}} \left( \frac{n_{11}}{N} f_{11}(x) + \frac{n_{01}}{N} (1 - f_{01}(x)) \right) dx + \int_{-\infty}^{\theta_0^{\tau-1}} \left( \frac{n_{10}}{N} f_{10}(x) + \frac{n_{00}}{N} (1 - f_{00}(x)) \right) dx, \\ \epsilon_1^\tau &= \int_{\infty}^{\theta_1^{\tau-1}} f_{11}(x) dx - \int_{\infty}^{\theta_0^{\tau-1}} f_{01}(x) dx.\end{aligned}\quad (9)$$

181 Taking the derivative of (8) w.r.t.  $\Delta_1$  and setting it to 0, we can easily obtain the closed-form solution  
 182 of  $\Delta_1^\tau$  as

$$\Delta_1^\tau = -\frac{\alpha^\tau \eta^\tau + \lambda \beta^\tau \epsilon^\tau}{(\alpha^\tau)^2 + \lambda (\beta^\tau)^2}. \quad (10)$$

---

**Algorithm 1** Optimization Algorithm of GSTAR Model

---

**Input** dataset  $\mathcal{X} \times \mathcal{A} \times \mathcal{Y} = \{(\mathbf{x}_i, \mathbf{a}_i, \mathbf{y}_i)\}_{i=1}^n$ , classification model  $h(X)$ , hyperparameter  $\lambda$ .

**Output** Group-specific threshold  $\theta = (\theta_1, \theta_0)$ .

**Initialize**  $\theta = (\theta_1, \theta_0) = (0, 0)$ .

1. Given a classifier  $H(x)$ , estimate probability density function  $f_{ya}, y, a \in \{0, 1\}$  by maximum likelihood estimation.

**while** not converge **do**

2. Calculate the optimal step  $\Delta_1$  as  $\Delta_1 = -\frac{\alpha_1\eta_1 + \lambda\beta_1\epsilon_1}{\alpha_1^2 + \lambda\beta_1^2}$ , with  $\alpha_1, \beta_1, \eta_1, \epsilon_1$  values shown in (9);

3. Update the threshold:  $\theta_1 \leftarrow \theta_1 + \Delta_1$ ;

4. Calculate the optimal step  $\Delta_0$  as  $\Delta_0 = -\frac{\alpha_0\eta_0 + \lambda\beta_0\epsilon_0}{\alpha_0^2 + \lambda\beta_0^2}$  with  $\alpha_0, \beta_0, \eta_0, \epsilon_0$  values calculated in a similar way as in (9);

5. Update the threshold:  $\theta_0 \leftarrow \theta_0 + \Delta_0$ .

**end while**

---

183 **The second step** is to fix  $\theta_1$  and update  $\theta_0$ , and this can be achieved in a similar way of updating  $\theta_1$ .  
184 Then we can finalize the alternating optimization as:

$$\begin{aligned}\theta_0^\tau &= \theta_0^{\tau-1} + \Delta_0^\tau, \\ \theta_1^\tau &= \theta_1^{\tau-1} + \Delta_1^\tau.\end{aligned}\tag{11}$$

185 It is notable that in each iteration we derive the optimal update step  $\theta$ , which eliminates the burden  
186 of tuning hyperparameter (such as learning rate) in iterative algorithm. The optimization step is  
187 summarized in Algorithm 1. The above algorithm can easily extend to multiple fairness constraints  
188 by adding corresponding squared-loss fairness terms to (2).

189 **Time Complexity:** The alternating optimization of GSTAR model is of low computational cost. We  
190 take at most one pass of the data for learning the estimated probability density functions  $f_{ya}$  in (1)  
191 (we do not even need to traverse the data if the parameters (such mean and variance in Gaussian  
192 distribution) for the estimated probability density functions  $f_{ya}$  can be provided). The optimization  
193 of  $\theta$  with alternating optimization is efficient since we only need  $f_{ya}$  as we have seen in (9) and (10).  
194  $\theta \in \mathbb{R}^2$  is a vector with fixed small size. Therefore, we need a constant time for each update. Overall,  
195 the time complexity of GSTAR is  $O(n + T)$ , where  $n$  is the number of samples, and  $T$  is the number  
196 of iterations in alternating optimization.

197 We further derive the theoretical analysis of our GSTAR model on the balance between fairness  
198 and accuracy, which indicates that the optimal solution provides guarantees on model accuracy  
199 under the optimal fairness constraint. Details of the theoretical analysis is in the Supplementary  
200 material. Besides, if a unified threshold is necessary [3], i.e.,  $\theta_1 = \theta_0$ , the optimization algorithm  
201 also applies and we only have one scalar variable in (2). When we have a unified threshold, we do not  
202 require sensitive information in the testing phase that we can conform more strict privacy regulations  
203 than group-aware thresholding. However, we have to sacrifice both fairness and accuracy as the  
204 thresholding is less flexible.

## 205 4 Experiments

206 In this section, we validate GSTAR model on four well-known fairness datasets and compare with  
207 other state-of-the-art methods. First, we plot Pareto frontiers of ours and FACT (MS) [10] to  
208 demonstrate the trade-offs between fairness and accuracy. Second, we evaluate the models with  
209 different fairness metrics and validate that our model is highly adaptive to any fairness metrics that  
210 are expressible with confusion matrix [7, 11, 2, 1]. Third, we use our model as a post-processing  
211 method to existing fair models and show that our model further improves existing fair models in an  
212 efficient and model-agnostic manner.

### 213 4.1 Experimental Setup

214 We compare with multiple fairness approaches in the experiments. For clear demonstration of results,  
215 we use different shapes of marker for each comparing methods in Figure 2 and Figure 4. The compar-  
216 ing methods include: **Learning fair representations for kernel models** (abbreviated as FGP) [23],

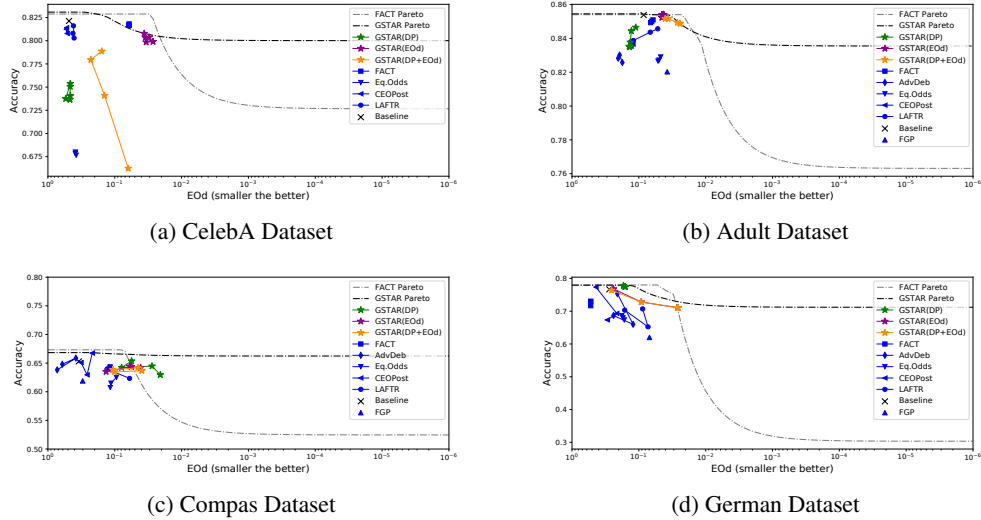


Figure 2: Model-specific Pareto frontiers of equalized odds to show the upper bound of best achievable accuracy under different fairness constraints. Upper right region under the boundary shows better fairness and higher accuracy. We plot three variations of GSTAR (star-shaped) with different fairness objectives. GSTAR is the closest to the Pareto frontier which indicates the best trade-offs.

217 **Fairness confusion tensor** (abbreviated as FACT) [10], **Disparate impact remover** (abbreviated  
 218 as DIR) [6], **Adversarial de-biasing** (abbreviated as AdvDeb) [24], **Calibrated equalized odds**  
 219 **post-processing** (abbreviated as CEOPost) [20], **Equality of opportunity in supervised learning**  
 220 (abbreviated as Odds) [7], **Learning adversarially fair and transferable representations** (abbreviated  
 221 as LAFTR) [17], and **Baseline**: For CelebA dataset, we use ResNet50 [8] as a reference, and  
 222 logistic regression for all other datasets. Our method is optimized with  $\lambda$  in the range of  $[10^{-1}, 10^4]$   
 223 with alternating optimization method. All experiments are implemented with Pytorch framework on  
 224 i9-9960X CPU and a Quadro RTX 6000 GPU.

225 We choose broadly used fairness metrics in evaluation including: **equal opportunity difference** and  
 226 **equalized odds difference** (abbreviated as EOp, and EOd respectively) [7] ; **1-disparate impact**  
 227 (abbreviated as 1-DIMP) [1]; **balanced accuracy difference** (abbreviated as BD).

228 We evaluate the methods on four fairness datasets: **CelebA** image dataset<sup>1</sup> [15], **Adult** dataset from  
 229 the UCI repository [12], **COMPAS**<sup>2</sup> (Correctional Offender Management Profiling for Alternative  
 230 Sanctions) dataset, and **German** credit dataset from the UCI repository [5]. All data is split as 70%  
 231 for training and 30% for testing. More details of the comparing methods, evaluation metrics, and  
 232 datasets are provided in the Supplementary material.

## 233 4.2 Performance and Fairness-Accuracy Trade-Offs

234 In this subsection, we look into the performance evaluation of GSTAR comparing with other state-of-  
 235 the-art methods. We consider Pareto frontier to visualize the trade-offs between fairness and accuracy  
 236 to demonstrate the measure of performance.

237 In Figure 2, we plot Pareto frontier, which is the upper bound for the accuracy-fairness trade-offs,  
 238 desired output locates at the upper right region under the boundary which corresponds to higher  
 239 values in accuracy and lower values in fairness discrepancy. With the same fairness constraints  
 240 are given, we achieve a better frontier than the FACT [10] as we equally weigh on demographic  
 241 statistics and have a better feasible region. To obtain our results (star points), we first estimate the  
 242 logit distribution from the output of the baseline model, and then we get optimal adaptive thresholds  
 243 with corresponding fairness metric by updating w.r.t. the objective function in (2). Here we have three  
 244 combinations of fairness imposed to GSTAR: demographic parity (DP), equalized odds (EOd), and

<sup>1</sup><http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

<sup>2</sup><https://github.com/propublica/compas-analysis>

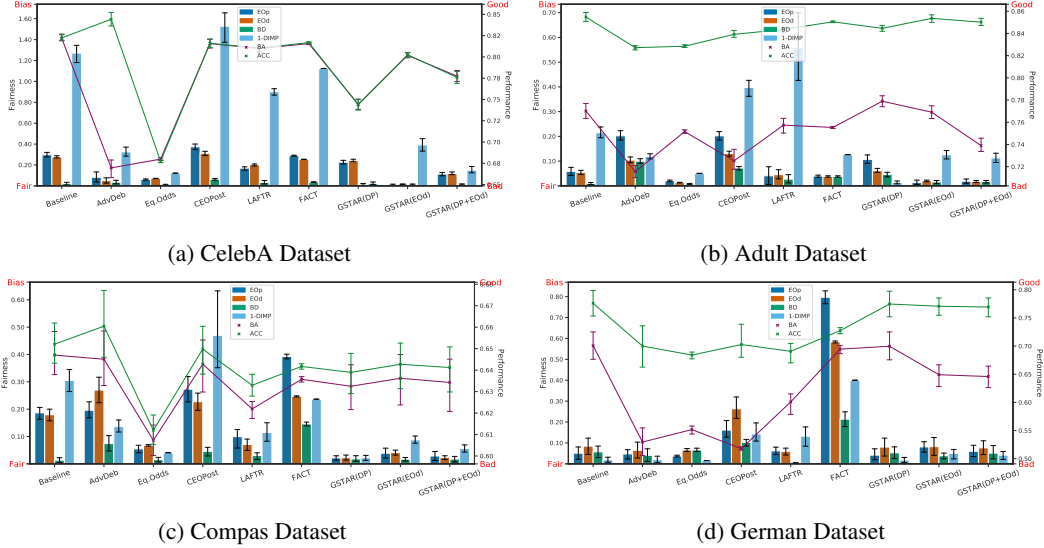


Figure 3: Quantitative evaluation on fairness and performance metrics. The bar plots indicate fairness measures (EOp, EOd, BD, 1-DISP) of each model. Lower fairness values in the left y-axis shows better fairness. The line plots indicate the performance measure (balanced accuracy (BA) and accuracy (ACC) of each model. Higher performance values in the right y-axis shows better classification performance. We consider three variations of GSTAR models (DP, EOd, DP+EOd).

245 with both constraints (DP+EOd). By post-processing on a simple baseline, we achieved significantly  
 246 better fairness with small or no sacrifice in accuracy. In all datasets, GSATR got competitive or better  
 247 results than other state-of-the-art methods on both fairness and accuracy.

248 For example, we got  $\theta_{EOd}^* = (0.640, -0.627)^\top$  for the CelebA dataset. This shows that we have  
 249 a higher threshold for the privileged group and a lower threshold for the unprivileged group. This  
 250 optimal thresholding from GSTAR allows more samples from the privileged group to be correctly  
 251 predicted as unattractive that would compensate for the discrimination of the original model. In other  
 252 words, this improves predictive equality [2] with a huge amount from 0.235 to 0.014. Also, true  
 253 positive rate difference (also known as equality of opportunity [7]) got reduced from 0.282 to 0.018.  
 254 It is notable that GSTAR only sacrificed 2.2% of accuracy to bring the big improvement in fairness.

255 Since the objective function of our model is independent to data dimensionality, our model is much  
 256 more efficient especially for high dimensional data. We mostly outperform the computational cost  
 257 comparing to the other methods. The comparison of computational time on the datasets can be found  
 258 in the Supplementary material.

### 259 4.3 Flexibility and Multiple Fairness Constraints

260 Since each fairness metric has different interests, it has been theoretically proven that they cannot be  
 261 perfectly satisfied all together [20, 2, 11]. Because of this inherent trade-offs between fairness metrics,  
 262 most of the recent works focus on a single metric at a time to achieve fairness. However with GSTAR,  
 263 we have the flexibility to optimize on multiple fairness constraints that can be represented in the  
 264 confusion matrix format. Moreover, given the estimated distribution  $f_{y_a}$  of a black-box classification  
 265 model, we can adjust the optimal  $\theta$  based on the needs by accommodating different fairness criteria.

266 Figure 3 demonstrates the result of the methods with fairness metrics and accuracy trade-off eval-  
 267 uations. Overall, the variations of GSTAR achieve the best fairness on each target fairness while  
 268 preserving the performance. For example in Figure 3(a), GSTAR with EOd constraint has outstanding  
 269 performance in most fairness metrics with comparable accuracy (80.3%). Comparing with GSTAR  
 270 (EOd), when we introduce EOd and DP together (DP+EOd), we achieve significantly better w.r.t. DP  
 271 fairness with sacrificing a small amount of accuracy and EOd.



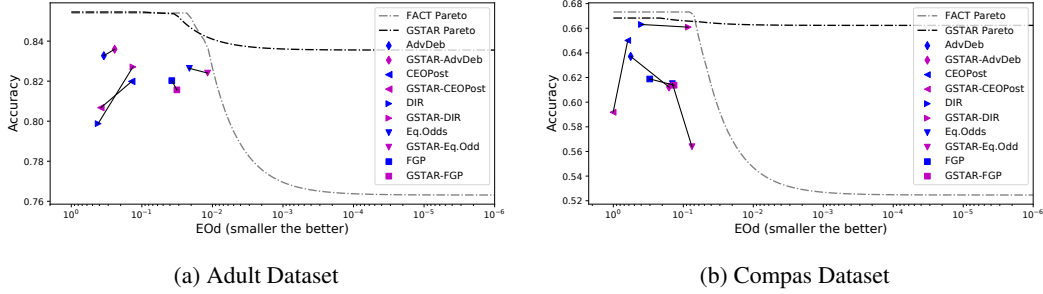


Figure 4: Illustration of post-processing (magenta colored points) on existing fairness models (blue colored points). Given the outputs of each model, we efficiently improve existing fairness models with optimized group-aware thresholds from GSTAR.

272 In general, by sacrificing individual fairness performance, we could introduce multiple constraints.  
 273 Also, we implicitly observe that the more fairness constraints are introduced, the more accuracy is  
 274 sacrificed. We empirically found that in some cases (e.g. Figure 3(c)), introducing multiple fairness  
 275 is complementary to each other that improves both conditions.

#### 276 4.4 Post-Processing on an Existing Fair Model

277 For a binary classifier that has a single fixed classification threshold (0 for out logit, and 0.5 for  
 278 label probability), we can improve the trade-off between fairness and accuracy via GSTAR post-  
 279 processing. Given the logit/probability of the dataset from a black-box model, we can improve the  
 280 fairness as illustrated in Figure 4. In most cases, we observe improvement in fairness after GSTAR  
 281 post-processing. It is also interesting to note that by optimizing the different thresholds for each  
 282 protected group, we even obtain better performance on both fairness and accuracy, which indicates  
 283 that the threshold optimization can not only improve fairness but also accuracy.

284 However, when the distribution of the logits/probability is highly extreme (such as the results of using  
 285 GSTAR to post-process CEOPost), it is difficult to estimate the distribution and thus causes erroneous  
 286 optimization in GSTAR. We empirically found that when the dataset is extremely imbalanced such  
 287 that we do not have enough samples to estimate the logit/probability distribution, or black-box model  
 288 is too certain to the prediction that samples are concentrated to certain output, this problem arises.

## 289 5 Conclusion and Discussion

290 In this paper, we propose a group-aware threshold adaptation method (GSTAR) to post-process a black-  
 291 box model and optimize over multiple fairness constraints. We directly optimize the classification  
 292 threshold for each demographic group w.r.t. the classification error and multiple fairness constraints  
 293 in a unified objective function, such that we can practically achieve an optimal trade-off between  
 294 accuracy and fairness in fair classification. Our method is applicable to diverse notions of group  
 295 fairness as the majority of fairness notions can be expressed as a linear or quadratic equation through  
 296 confusion matrix. We empirically show that GSTAR is *flexible* with fairness regularization, *efficient*  
 297 with low computational cost. We also notice that the adaptive thresholds benefit accuracy in some  
 298 cases. GSTAR agrees to protect *privacy* such as article 17 of EU’s GDPR [21] with model-agnostic  
 299 post-processing. We only require the estimated distribution of the output from a black-box model  
 300 i.e., our post-processing method is oblivious to features. Thus training data is no longer needed and  
 301 allowed to be discarded after training the black-box model.

302 Further, we empirically find that GSTAR is not applicable to post-process some classification models  
 303 in the following situations: 1) the model does not provide logit/probability as the outcome; 2) The  
 304 model provides an extreme distribution of the output logit/probability. For example, when the model  
 305 is too certain about its prediction, it will be difficult to perform probability density estimation. In our  
 306 future work, we will study possible strategies to solve the above limitations, and extend GSTAR to  
 307 multi-class, multi-sensitive group problems and improve the fairness-accuracy trade-off in a more  
 308 general scheme.

309 **References**

- 310 [1] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016.
- 311 [2] Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism  
312 prediction instruments. *Big data*, 5(2):153–163, 2017.
- 313 [3] Sam Corbett-Davies, Emma Pierson, Avi Feller, Sharad Goel, and Aziz Huq. Algorithmic  
314 decision making and the cost of fairness. In *KDD*, pages 797–806, 2017.
- 315 [4] Julia Dressel and Hany Farid. The accuracy, fairness, and limits of predicting recidivism. *Sci.*  
316 *Adv*, 4(1):eaao5580, 2018.
- 317 [5] Dheeru Dua and Casey Graff. UCI machine learning repository, 2019.
- 318 [6] Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkata-  
319 subramanian. Certifying and removing disparate impact. In *KDD*, pages 259–268, 2015.
- 320 [7] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In  
321 *NeurIPS*, pages 3315–3323, 2016.
- 322 [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
323 recognition. In *CVPR*, pages 770–778, 2016.
- 324 [9] Faisal Kamiran, Asim Karim, and Xiangliang Zhang. Decision theory for discrimination-aware  
325 classification. In *ICDM*, pages 924–929. IEEE, 2012.
- 326 [10] Joon Sik Kim, Jiahao Chen, and Ameet Talwalkar. Model-agnostic characterization of fairness  
327 trade-offs. *arXiv preprint arXiv:2004.03424*, 2020.
- 328 [11] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair  
329 determination of risk scores. *arXiv preprint arXiv:1609.05807*, 2016.
- 330 [12] Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In  
331 *KDD*, volume 96, pages 202–207, 1996.
- 332 [13] Lydia T Liu, Sarah Dean, Esther Rolf, Max Simchowitz, and Moritz Hardt. Delayed impact of  
333 fair machine learning. *arXiv preprint arXiv:1803.04383*, 2018.
- 334 [14] Lydia T Liu, Max Simchowitz, and Moritz Hardt. The implicit fairness criterion of unconstrained  
335 learning. In *ICML*, pages 4051–4060, 2019.
- 336 [15] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the  
337 wild. In *ICCV*, pages 3730–3738, 2015.
- 338 [16] Vishnu Suresh Lokhande, Aditya Kumar Akash, Sathya N Ravi, and Vikas Singh. Fairalm:  
339 Augmented lagrangian method for training fair models with little regret. In *ECCV*, pages  
340 365–381. Springer, 2020.
- 341 [17] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair  
342 and transferable representations. *arXiv preprint arXiv:1802.06309*, 2018.
- 343 [18] Aditya Krishna Menon and Robert C Williamson. The cost of fairness in binary classification.  
344 In *ACM FAccT*, pages 107–118, 2018.
- 345 [19] Dino Pedreshi, Salvatore Ruggieri, and Franco Turini. Discrimination-aware data mining. In  
346 *KDD*, pages 560–568, 2008.
- 347 [20] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On fairness  
348 and calibration. In *NeurIPS*, pages 5680–5689, 2017.
- 349 [21] General Data Protection Regulation. Regulation eu 2016/679 of the european parliament and of  
350 the council of 27 april 2016. *Official Journal of the European Union*. Available at: [http://ec.europa.eu/justice/data-protection/reform/files/regulation\\_oj\\_en.pdf](http://ec.europa.eu/justice/data-protection/reform/files/regulation_oj_en.pdf) (accessed 20 September  
351 2017), 2016.  
352

- 353 [22] Hee Jung Ryu, Hartwig Adam, and Margaret Mitchell. Inclusivefacenet: Improving face  
354 attribute detection with race and gender diversity. *arXiv preprint arXiv:1712.00193*, 2017.
- 355 [23] Zilong Tan, Samuel Yeom, Matt Fredrikson, and Ameet Talwalkar. Learning fair representations  
356 for kernel models. In *AISTATS*, pages 155–166, 2020.
- 357 [24] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with  
358 adversarial learning. In *AIES*, pages 335–340, 2018.
- 359 [25] Han Zhao and Geoff Gordon. Inherent tradeoffs in learning fair representations. In *NeurIPS*,  
360 pages 15675–15685, 2019.

## 361 Checklist

- 362 1. For all authors...
- 363 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
364 contributions and scope? [Yes]
- 365 (b) Did you describe the limitations of your work? [Yes] See Section 5.
- 366 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See  
367 Section 5.
- 368 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
369 them? [Yes]
- 370 2. If you are including theoretical results...
- 371 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section 1  
372 in the Supplementary material.
- 373 (b) Did you include complete proofs of all theoretical results? [Yes] See Section 1 in the  
374 Supplementary material.
- 375 3. If you ran experiments...
- 376 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
377 perimental results (either in the supplemental material or as a URL)? [Yes] We will  
378 provide the code and instructions at request. The data used in the experiments are  
379 public available. See Section 2.3 in the Supplementary material.
- 380 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
381 were chosen)? [Yes] See Section 4.1.
- 382 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
383 ments multiple times)? [Yes] See Section 4.3.
- 384 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
385 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4.1.
- 386 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 387 (a) If your work uses existing assets, did you cite the creators? [N/A]
- 388 (b) Did you mention the license of the assets? [N/A]
- 389 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]  
390
- 391 (d) Did you discuss whether and how consent was obtained from people whose data you’re  
392 using/curating? [N/A]
- 393 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
394 information or offensive content? [N/A]
- 395 5. If you used crowdsourcing or conducted research with human subjects...
- 396 (a) Did you include the full text of instructions given to participants and screenshots, if  
397 applicable? [N/A]
- 398 (b) Did you describe any potential participant risks, with links to Institutional Review  
399 Board (IRB) approvals, if applicable? [N/A]
- 400 (c) Did you include the estimated hourly wage paid to participants and the total amount  
401 spent on participant compensation? [N/A]