

# Group-Aware Threshold Adaptation for Fair Classification

Taeuk Jang<sup>1</sup>, Pengyi Shi<sup>2</sup>, Xiaoqian Wang<sup>1\*</sup>

<sup>1</sup>School of Electrical and Computer Engineering, Purdue University, West Lafayette, USA, 47907

<sup>2</sup>Krannert School of Management, Purdue University, West Lafayette, USA, 47907  
{jang141, shi178, joywang}@purdue.edu

## Abstract

The fairness in machine learning is getting increasing attention, as its applications in different fields continue to expand and diversify. To mitigate the discriminated model behaviors between different demographic groups, we introduce a novel post-processing method to optimize over multiple fairness constraints through group-aware threshold adaptation. We propose to learn adaptive classification thresholds for each demographic group by optimizing the confusion matrix estimated from the probability distribution of a classification model output. As we only need an estimated probability distribution of model output instead of the classification model structure, our post-processing model can be applied to a wide range of classification models and improve fairness in a model-agnostic manner and ensure privacy. This even allows us to post-process existing fairness methods to further improve the trade-off between accuracy and fairness. Moreover, our model has low computational cost. We provide rigorous theoretical analysis on the convergence of our optimization algorithm and the trade-off between accuracy and fairness. Our method theoretically enables a better upper bound in near optimality than previous method under the same condition. Experimental results demonstrate that our method outperforms state-of-the-art methods and obtains the result that is closest to the theoretical accuracy-fairness trade-off boundary.

## Introduction

Machine learning is broadening its impact in various fields including credit analysis, job screening and *etc.* Consequently, the importance of fairness in machine learning is emerging. However, recent models have been found to behave differently between demographic groups in favorable predictions. For example, it has been discovered that COMPAS, the criminal risk assessment software currently used to help pretrial release decisions, has biases between different races (Dressel and Farid 2018). Specifically, blacks got higher risk scores predicted from the model than whites with similar profiles. Therefore, discrimination truly exists and resolving it is critical as its direct and potential impact is growing tremendously.

However, obtaining fairness is not a trivial problem, as the dataset itself will be biased when it is accumulated artificially (Jang, Zheng, and Wang 2021). Simply modifying sensitive features (such as *race*, *gender*) from the data does not solve the bias, because there is indirect discrimination (Pedreshi, Ruggieri, and Turini 2008) caused by the feature relevance, which means sensitive information can be inferred from other features.

In order to alleviate discrimination from different perspectives, various quantitative measurements of group equity (Hardt, Price, and Srebro 2016; Kleinberg, Mullainathan, and Raghavan 2016; Chouldechova 2017) have been proposed. It has been proven that the pursuit of fairness is subject to a trade-off between fairness and accuracy (Liu et al. (2019), Kim et al. (2020)).

Moreover, Pleiss et al. (2017) studied the trade-offs between fairness notions that cannot be satisfied at the same time. Therefore, recent works (Feldman et al. 2015; Zhang, Lemoine, and Mitchell 2018; Hardt, Price, and Srebro 2016) usually target at a certain fairness notion. However, these approaches suffer from the *lack of flexibility*, *i.e.*, target fairness cannot be adjusted to meet the needs. If the fairness constraints change under some circumstances, traditional fairness models need to be re-trained from scratch, which is computationally demanding and sometimes inapplicable due to model settings.

To overcome the limitations, we propose a novel post-processing method to improve fairness in a model-agnostic manner *i.e.*, we only need the prediction of an unknown model. Our GSTAR (Group Specific Threshold Adaptation for fair classification) model learns adaptive classification thresholds for each demographic group in classification task to improve the trade-off between fairness and accuracy. Given an existing classification model, GSTAR approximates the probability distribution of the model output and utilizes confusion matrix to quantify accuracy and fairness w.r.t. the group-aware classification thresholds. This allows us to: 1) prevent from burdening additional complexity or deteriorate the stability of the training process of the classifier; 2) integrate different fairness notions into one unified objective function; 3) easily adapt one pre-trained model to other fairness constraints.

We summarize our contributions of this paper as follows:

1. We propose a novel post-processing method, named

\*Corresponding author.

GSTAR, which can learn group-aware thresholds to optimize the fairness-accuracy trade-off in classification. We empirically show that GSTAR outperforms state-of-the-art methods.

2. With GSTAR, we can simultaneously optimize over multiple fairness constraints with a low computational cost. GSTAR does not require multiple iterations over data, instead, it takes *at most* one pass of data in training for fast computation.
3. We conduct extensive rigorous theoretical analysis on our method, in terms of convergence analysis and fairness-accuracy trade-off. We introduce theoretical improvement in terms of near optimality.
4. We derive Pareto frontiers of our model for the fairness-accuracy trade-offs that contextualize the quality of fair classification.

## Related Works

In order to achieve group fairness, which quantifies the discrimination among different sensitive groups, a diverse notion of fairness has been introduced. Equalized odds (Hardt, Price, and Srebro 2016) enforce equality of true positive rates and false positive rates between different demographic groups. Pleiss et al. (2017) relaxed equalized odds to satisfy group-wise calibration. Demographic parity or disparate impact (Barocas and Selbst 2016) suggests that a model is unbiased if the model prediction is independent of the protected attribute.

Among different fairness methods, post-processing techniques propose to improve fairness by modifying the output of a given classifier. Hardt et al. (2016) propose to ensure equalized odds by constraining the model output. Kim et al. (2020) utilize confusion matrix and propose least-square accuracy-fairness optimization problem. Kamiran et al. (2012) propose to give a favorable outcome to unprivileged and an unfavorable outcome to the privileged group when the confidence of the prediction is in a certain range. However, such *static* confidence window keeps the same regardless of the demographic group and is determined by grid search, so it is less efficient.

Threshold adjustment (a.k.a. thresholding) was introduced to improve the performance of *static* thresholds. In the literature, Menon et al. (2018) prove that instance-dependent thresholding of the predictive probability function is the optimal classifier in cost-sensitive fairness measures. Also, when considering immediate utility, Corbett-Davies et al. (2017) show that optimal algorithm is achieved from group-specific threshold which is determined by group statistics. However, to the best of our knowledge, the threshold adjustment approach has not been deeply studied that neither encompasses broad group fairness metrics nor describes an explicit method to achieve the threshold.

Trade-off between fairness and accuracy exists when we impose fairness constraint to a model. Recent studies (Chouldechova 2017; Zhao and Gordon 2019) prove that models targeting at such fairness notions conform to an information theoretic lower bound on the joint error across different sensitive groups. Therefore, our work presents a prac-

tical upper bound of the best achievable accuracy given the fairness constraints.

Here, our work is the most related to the post-processing methods (Hardt, Price, and Srebro 2016; Kim, Chen, and Talwalkar 2020). However, ours differ from theirs in several aspects. First, we theoretically prove that GSTAR achieves a better upper bound of near optimality than Hardt et al. (2016) as we directly operate on ROC curve instead of linear intersections in Hardt et al. (2016). Also, GSTAR corrects the predicted label by the confidence of the prediction from a given model instead of randomly flipping the output to achieve equalized odds, which is more reliable in post-processing. FACT (Kim, Chen, and Talwalkar 2020) utilizes a single point (static) from the classifier to be post-processed as a reference which does not fully utilize the classifier for the post-processing. In contrast, by approximating the distribution of the continuous predicted logits, GSTAR model enables a larger feasible region than Kim et al. (2020) with a better fairness-accuracy trade-off. We validate the improvement in this trade-off via both theoretical and empirical results. It is notable that these related methods can be considered as a special case of GSTAR.

## GSTAR for Fair Classification

### Motivation

Consider a binary classification problem with a binary sensitive feature, such that the sensitive feature  $A \in \{0, 1\}$  and label  $Y \in \{0, 1\}$ . In general, for a given data  $X$ , a binary classification model outputs an unnormalized logit  $h(X) \in \mathbb{R}$  with the class label probability  $R(X) = \sigma(h(X)) \in [0, 1]$ , where  $\sigma$  is an activation function (*e.g.*, sigmoid function). It is not necessary to calculate  $R$  in a classification model, *e.g.*, support vector machines directly use the positiveness/negativeness of logit  $h(X)$  to determine classification outcome.

For traditional models, we use a cut-off threshold  $\theta_h = 0$  for  $h(X)$  (*i.e.*,  $\theta_R = \sigma(0) = 0.5$  for  $R(X)$ ) in classification, such that the predicted label is determined by  $\hat{Y} = \mathbb{I}\{h(X) \geq \theta_h\}$ . In the following context, unless otherwise mentioned, we use  $\theta$  to refer to the threshold  $\theta_h$  on logit  $h$  since it is applicable to a wider range of classification models, and the corresponding threshold on label probability  $\theta_R$  can be easily inferred from the threshold on logit  $h$ . Traditional models use the same cut-off threshold  $\theta$  for different demographic groups. However, since the distribution of logits  $h$  in different demographic groups can be different, using the same threshold  $\theta$  brings biased classification.

In Fig. 1, we use a real-world example of image classification on CelebA dataset with ResNet50 (He et al. 2016) to show that the default setting of classification thresholds affects both accuracy and fairness in classification. The goal is to predict whether the image of a person is attractive or not, and consider sensitive attribute as gender. This can be generalized to different sensitive attributes in image classification task, *e.g.*, age or race (Lokhande et al. 2020). We can observe an obvious difference in the distribution of logit  $h$  between two gender groups. If we use a unified classification threshold  $\theta_1 = \theta_0 = 0$ , it naturally brings a difference in the true positive rate and true negative rate between two

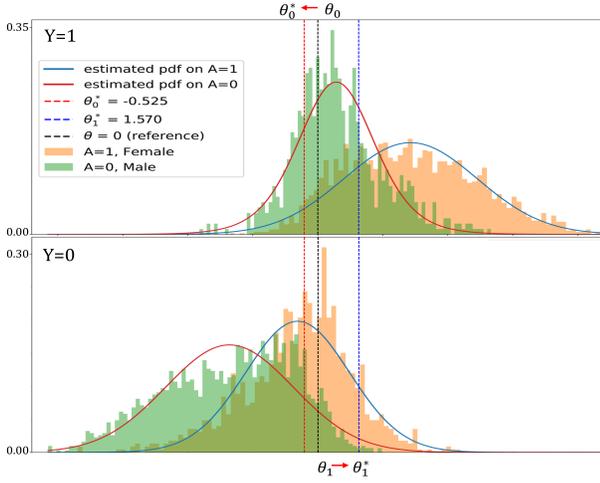


Figure 1: Histograms of logit  $h$  distribution from logistic regression on CelebA data, where  $\theta$  is the threshold to assign predicted label based on  $h$ . The top and bottom plot is for positive samples ( $Y = 1$ , attractive), and negative samples ( $Y = 0$ , unattractive). Bars represent the distributions of logit  $h$  of sensitive groups, and curves are estimated probability density functions of logit  $h$  of sensitive groups as in the legend.  $\theta = 0$  (black dashed line) is the default classification thresholds. The default thresholds result in biased prediction towards the unprivileged group ( $A = 0$ ) due to the different logit  $h$  distributions in different sensitive groups.  $(\theta_0^*, \theta_1^*)$  (colored dashed line) are group-aware thresholds for each sensitive group achieved by GSTAR.

gender groups, thus it behaves as a biased classification. Instead, we observe that the optimal group-specific threshold obtained from GSTAR ( $\theta_1^* > \theta_1$ , and  $\theta_0^* < \theta_0$ ) can adapt to such discrepancy in distribution between two demographic groups to improve both fairness and accuracy.

### Group-Aware Classification Thresholds

Given an existing classification model and a sensitive attribute  $a$ , we can denote true positive rate ( $TP_a$ ), false positive rate ( $FP_a$ ), true negative rate ( $TN_a$ ), and false negative rate ( $FN_a$ ) in the confusion matrix. Most fairness notions can be represented with entries in the confusion matrix. For instance, Equal Opportunity (EOp) (Hardt, Price, and Srebro 2016) requires  $TP_0 = TP_1$ , and Demographic Parity (DP) (Barocas and Selbst 2016) requires

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

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

Consider the group-aware classification threshold  $\theta = (\theta_1, \theta_0)^T$ , where  $\theta_a$  is the classification threshold for sensitive group  $A = a$ . We can formulate the entries in the

confusion matrix w.r.t.  $\theta$  as below:

$$\begin{aligned} TP_a(\theta_a) &\approx 1 - \int_{-\infty}^{\theta_a} f_{1a}(x) dx, & FN_a(\theta_a) &\approx 1 - TP_a(\theta_a) \\ FP_a(\theta_a) &\approx 1 - \int_{-\infty}^{\theta_a} f_{0a}(x) dx, & TN_a(\theta_a) &\approx 1 - FP_a(\theta_a) \end{aligned} \quad (1)$$

where  $f_{ya}(x)$  is an estimated probability density function of the distribution of output logit  $h$  in the subset  $\{Y = y, A = a\}$ .

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

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

where  $\lambda$  is a hyperparameter that determines how much fairness is enforced in the optimization. The performance error  $\mathcal{L}_{per}(\theta)$  can be written as

$$\begin{aligned} \mathcal{L}_{per}(\theta) &= \left( \frac{n_{01}}{N} FP_1(\theta_1) + \frac{n_{11}}{N} FN_1(\theta_1) \right. \\ &\quad \left. + \frac{n_{00}}{N} FP_0(\theta_0) + \frac{n_{10}}{N} FN_0(\theta_0) \right)^2. \end{aligned}$$

As for  $\mathcal{L}_{fair}(\theta)$ , it can be formulated to any fairness metrics that are expressible with confusion matrix. For instance, when we impose EOp ( $TP_1 = TP_0$ ) and predictive equality (PE) ( $FP_1 = FP_0$ ) (Chouldechova 2017), we can get the corresponding  $\mathcal{L}_{fair}(\theta)$  by summing over the least squared form of each constraint. Also, satisfying EOp and PP is equivalent to satisfying Equalized Odds (EOd) (Hardt, Price, and Srebro 2016). This can be formulated in our  $\mathcal{L}_{fair}$  as

$$\begin{aligned} \mathcal{L}_{fair}^{EOd}(\theta) &= \mathcal{L}_{fair}^{EOp}(\theta) + \mathcal{L}_{fair}^{PP}(\theta) \\ &= (TP_1(\theta_1) - TP_0(\theta_0))^2 + (FP_1(\theta_1) - FP_0(\theta_0))^2. \end{aligned} \quad (3)$$

Note that a lower  $\mathcal{L}_{fair}$  value indicates a fairer threshold. When  $\mathcal{L}_{fair}^{EOd}(\theta) = 0$ , we can interpret as the  $\theta$  satisfies the perfect EOd fairness. Similar to (3), we can enforce multiple fairness constraints by summing over the least square of each metric with different weight constant  $\lambda$  to each fairness constraints if needed.

Also, it is notable that compared to FACT (Kim, Chen, and Talwalkar 2020) that enforces fairness through confusion tensor, our formulation of fairness in  $\mathcal{L}_{fair}(\theta)$  represents a direct notion of fairness metrics and improves the measures that allows us to achieve better performance and Pareto frontiers that is shown in Section and Fig. 2. For example, FACT integrates multiple constraints as a weighted sum with the weights being the number of samples in each class. In this expression, the imbalance between the two fairness criteria will grow as the degree of imbalance in the data increases. In contrast, our formulation expresses the constraints as the exact notion of each metric that is not biased by the statistics of the dataset and we observe improved Pareto frontier as in Fig. 2.

## Optimization of GSTAR

Our GSTAR objective in (2) lies in the family of Non-linear Least Squares Problem (NLSP) (Gratton, Lawless, and Nichols 2007). To optimize objective (2) and find the threshold  $\theta$ , we adopt the Gaussian-Newton optimization method (Gratton, Lawless, and Nichols 2007). Here we take EOp constraint as an example to show the alternating optimization steps, then  $\mathcal{L}_{fair}(\theta)$  can be written as

$$\mathcal{L}_{fair}^{EOp}(\theta) = (\text{TP}_1(\theta_1) - \text{TP}_0(\theta_0))^2. \quad (4)$$

To solve NLSP with the Gauss-Newton method, we first convert the nonlinear optimization problem to a linear least square problem using Taylor expansion. That is, the parameter values are calculated in an iterative fashion with

$$\theta_a \approx \theta_a^{k+1} = \theta_a^k + \Delta_a, \quad (5)$$

in the  $k$ -th iteration number, with the vector of increments  $\Delta = \{\Delta_a\} = \{\theta_a^{k+1} - \theta_a^k\}$  (also known as the shift vector).

We rewrite our objective function as a real vector function  $r(\theta) = (r_1(\theta), r_2(\theta)) = (\mathcal{L}_{per}, \lambda \mathcal{L}_{fair})$ . We linearize each component in the loss function to a first-order Taylor polynomial expansion as

$$r_i(\theta) \approx r_i(\theta^k) + \sum_a \frac{\partial r_i(\theta^k)}{\partial \theta_a} \Delta_a \quad (6)$$

with  $\theta^k = (\theta_0^k, \theta_1^k)$ . Plugging this linearized equation into the objective function, we get the usual least square problem. Then, the optimal solution can be obtained as

$$\Delta = -(J^T J)^{-1} J^T f(\theta^k), \quad (7)$$

where  $J = \{J_{ia}\}$  with  $J_{ia} = \{\frac{\partial r_i(\theta)}{\partial \theta_a}\}$  is the Jacobian. Each entry of the jacobian can be expressed with linear combination of pdf and cdf of  $f_{ya}$  for  $i, a, y \in \{0, 1\}$ . we can finalize the alternating optimization as

$$\theta_0^\tau = \theta_0^{\tau-1} + \Delta_0^\tau, \quad \theta_1^\tau = \theta_1^{\tau-1} + \Delta_1^\tau. \quad (8)$$

It is notable that in each iteration we derive the optimal update step  $\Delta_a$ , which eliminates the burden of tuning hyperparameter (such as learning rate) in iterative algorithm. See the supplementary for detailed optimization process.

The alternating optimization of GSTAR model is of low computational cost. We take at most one pass of the data for learning the estimated probability density functions  $f_{ya}$  in (1) (we do not even need to traverse the data if the parameters (such mean and variance in Gaussian distribution) for the estimated probability density functions  $f_{ya}$  can be provided). The optimization of  $\theta$  with alternating optimization is efficient since we only need  $f_{ya}$ . Therefore, we need a constant time for each update. Overall, the time complexity of GSTAR is  $O(n + T)$ , where  $n$  is the number of samples, and  $T$  is the number of iterations in alternating optimization.

Besides, if a unified threshold is necessary (Corbett-Davies et al. 2017), *i.e.*,  $\theta_1 = \theta_0$ , the optimization algorithm also applies and we only have one scalar variable in (2). When we have a unified threshold, we do not require sensitive information in the testing phase that we can conform more strict privacy regulations than group-aware thresholding. However, we have to sacrifice both fairness and accuracy as the thresholding is less flexible.

## Theoretical Analysis

**Upper Bounds on FPR/FNR Gap between Groups** We first state the assumptions we need to make for Theorem 1 and 2.

**Assumption 1** For any given classifier  $h$  and its induced PDF  $f_{ya}$  and CDF  $F_{ya}$ , we assume the following holds:

- The PDF  $f_{ya}(x)$  is uniformly bounded, *i.e.*, there is an  $\hat{f}_{ya}(x) = \max_x f_{ya}(x)$ .
- The inverse CDF  $F_{ya}^{-1}(x)$  is Lipschitz continuous with Lipschitz constant  $M_{ya}$ .
- The difference in the CDF between two groups is uniformly bounded, *i.e.*,

$$|F_{y1}(x) - F_{y0}(x)| \leq u_y, \quad \forall x.$$

**Theorem 1** For any given classifier that satisfies Assumption 1 and any given pair of thresholds  $(\theta_0, \theta_1)$  that satisfies the perfect EOp condition, the gap between false-positive rates of the two group is upper bounded by

$$|\epsilon_1| = |FP_0(\theta_0) - FP_1(\theta_1)| \leq u_0 + C_1 u_1, \quad (9)$$

where  $C_1 = \hat{f}_{01} M_{10}$ .

**Theorem 2** For any given classifier that satisfies Assumption 1 and any given pair of thresholds  $(\theta_0, \theta_1)$  that satisfies the perfect PE condition, the gap between false-negative rates of the two group is upper bounded by

$$|\epsilon_2| = |FN_0(\theta_0) - FN_1(\theta_1)| \leq u_1 + C_0 u_0, \quad (10)$$

where  $C_0 = \hat{f}_{11} M_{00}$ .

Theorem 1 and 2 characterize the upper bound of false positive/negative rate gap between two groups when the false negative/positive rate gap is 0. At the same time, it captures the upper bound of additional accuracy loss due to the two different thresholds for different groups under a perfect fairness (EOp or PE) condition.

**Trade-off between Accuracy and Fairness** Now we prove a theorem to characterize the trade-off between accuracy and fairness. Let  $\theta_a^* = \text{argmin}_{\theta_a} \mathcal{L}_{per}(\theta_a)$ , and its perturbed value  $\tilde{\theta}_a$  as

$$\begin{aligned} |FN_1(\theta_1^*) - FN_1(\tilde{\theta}_1)| &\leq \gamma/2, \\ |FN_0(\theta_0^*) - FN_0(\tilde{\theta}_0)| &\leq \gamma/2, \end{aligned} \quad (11)$$

for some perturbation coefficient  $\gamma$ . Then for optimal perturbed version  $\tilde{\theta}_a^* = \text{argmin}_{\tilde{\theta}_a} \mathcal{L}_{per}(\tilde{\theta}_a)$ , we state the theorem below:

**Theorem 3** Under Assumption 1 and condition (11),

$$\mathcal{L}_{per}(\theta_1^*) - \mathcal{L}_{per}(\tilde{\theta}_1^*) \leq C\gamma,$$

where

$$C = 2L^* \left( \frac{r_1}{2} + r_0 \frac{\hat{f}_{01} M_{11}}{2} + \frac{n_{00}}{N} \left( \hat{f}_{00} M_{10} + \frac{\hat{e}'_1 M_{11}}{2} \right) + \frac{n_{10}}{N} \right)$$

and  $\hat{e}'_1 = \max \tilde{e}'_1$  is the maximum of the derivative of  $\tilde{e}_1$ .

Theorem 3 quantifies the decrease in accuracy loss (*i.e.*, the improvement in accuracy) when we allow a gap of true positive rates between two groups, *i.e.*, relaxation from the perfect fairness cases in Theorem 1 and 2.

**Convergence Analysis of GSTAR** Our objective function and the optimization solution algorithm belong to the family of Gauss-Newton algorithm. Given the assumptions A1 and A2 below,

- A1. There exists  $\theta^*$  such that  $J^T(\theta^*)r(\theta^*) = 0$ ,
- A2. The Jacobian at  $\theta^*$  has full rank,

we state the following theorem of convergence:

**Theorem 4** Assume that the estimated density function  $f(\cdot)$  satisfy assumptions A1 and A2. Further,  $f(\cdot)$  satisfies that

$$\|Q(\theta^k)(J^T J)^{-1}(\theta^k)\|_2 \leq \eta$$

for some constant  $\eta \in [0, 1)$  for each iteration  $k$ , where  $Q(\theta)$  denotes the second order terms  $\sum_i r_i(\theta)\nabla^2 r_i(\theta)$ . Then as long as the initial solution is sufficiently close to the true optimal with  $\|\theta^0 - \theta^*\|_2 \leq \epsilon$ , the sequence of Gauss-Newton iterates  $\{\theta^k\}$  converges to  $\theta^*$ .

**Near Optimality of GSTAR** Following the proof of Theorem 5.6 of Hardt et al. (2016), we provide the following near optimality theorem for our GSTAR model.

**Theorem 5** With a bounded loss function  $\ell$  and a given estimated density function  $f(x)$ , let  $\hat{R}_h \in [0, 1]$  be the induced random variable from the density  $f(x)$  of logit  $h(x)$ . Then the equalized odds predictor  $\hat{Y}_h$  derived from  $(\hat{R}_h, A)$  using the method in our paper can achieve near optimality in the following sense:

$$\mathbb{E}[\ell(\hat{Y}_h, Y)] \leq \mathbb{E}[\ell(Y^*, Y)] + 2d_K(\hat{R}_h, R^*).$$

Here,  $Y$  is the true label,  $Y^*$  is the optimal equalized odds predictor derived from the Bayes optimal regressor  $R^*$  as given in Hardt et al. (Hardt, Price, and Srebro 2016), and  $d_K(\hat{R}_h, R^*)$  is the conditional Kolmogorov distance.

Theorem 5 provides that GSTAR has tighter bound of near optimality than Hardt et al. (2016) under the same condition. See the supplementary for the proof of Theorem 1 - 5.

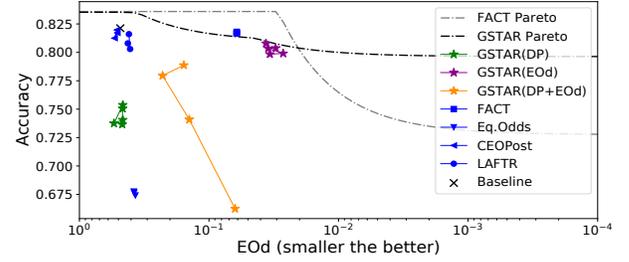
## Experiments

In this section, we validate GSTAR model on four well-known fairness datasets and compare with other state-of-the-art methods.

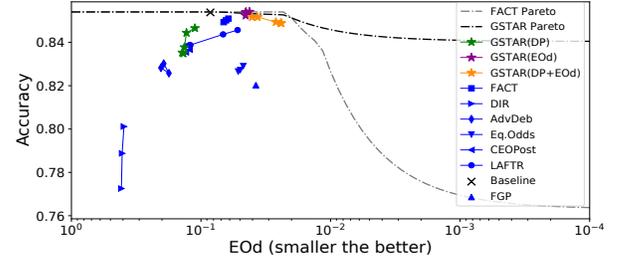
### Experimental Setup

We compare with multiple fairness approaches in the experiments. For clear demonstration of results, we use different shapes of marker for each comparing methods in Fig. 2 and Fig. 4. The comparing methods include: FGP (Tan et al. 2020), FACT (Kim, Chen, and Talwalkar 2020), DIR (Feldman et al. 2015), AdvDeb (Zhang, Lemoine, and Mitchell 2018), CEOPost (Pleiss et al. 2017), Eq.Odds (Hardt, Price, and Srebro 2016), LAFTR (Madras et al. 2018), and Baseline: For CelebA dataset, we use ResNet50 (He et al. 2016) as a reference, and logistic regression for all other datasets.

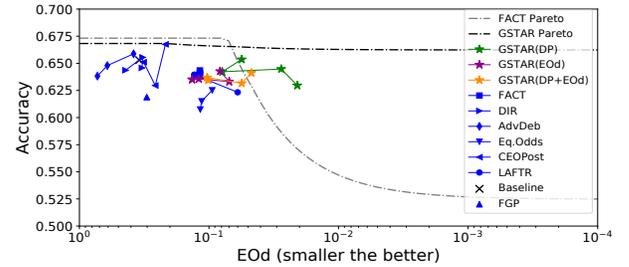
We choose broadly used fairness metrics in evaluation including: equal opportunity difference (EOp) and equalized



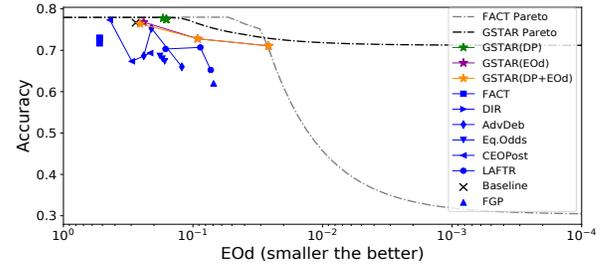
(a) CelebA Dataset



(b) Adult Dataset



(c) Compas Dataset



(d) German Dataset

Figure 2: Pareto frontiers of equalized odds to show the upper bound of best achievable accuracy under different fairness constraints. Upper right region under the boundary is desired. The variations of GSTAR generally achieve the best trade-offs as they are the closest to the Pareto frontier.

odds difference (EOd) (Hardt, Price, and Srebro 2016); 1-disparate impact (1-DIMP) (Barocas and Selbst 2016); balanced accuracy difference (BD). We use balanced accuracy (BA) and accuracy (ACC) as performance metrics.

We evaluate the methods on four fairness datasets:

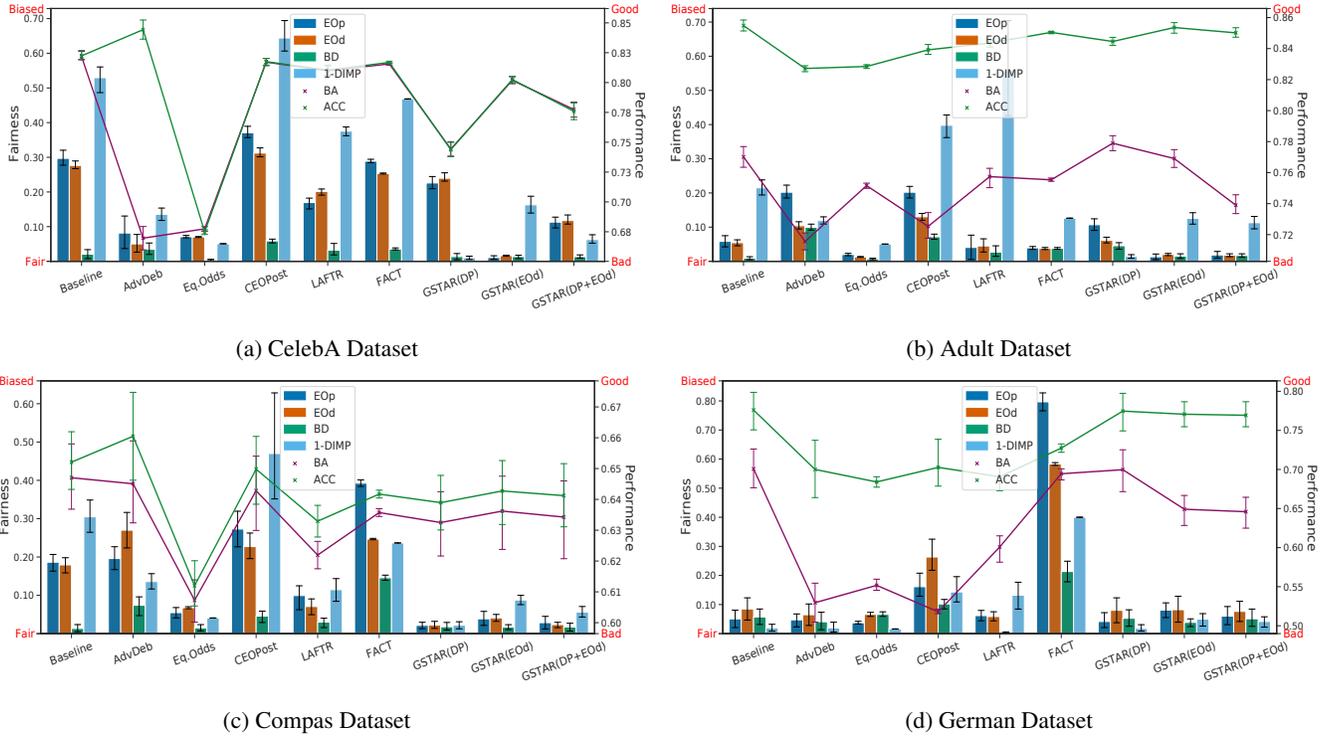


Figure 3: Evaluation on fairness and performance metrics. The bar plots indicate fairness measures of each model. The line plots indicate the performance measure of each model. Lower fairness values (left y-axis) and higher performance values (right y-axis) show better fairness and performance respectively. We consider three variations of GSTAR models (DP, EOd, DP+EOd).

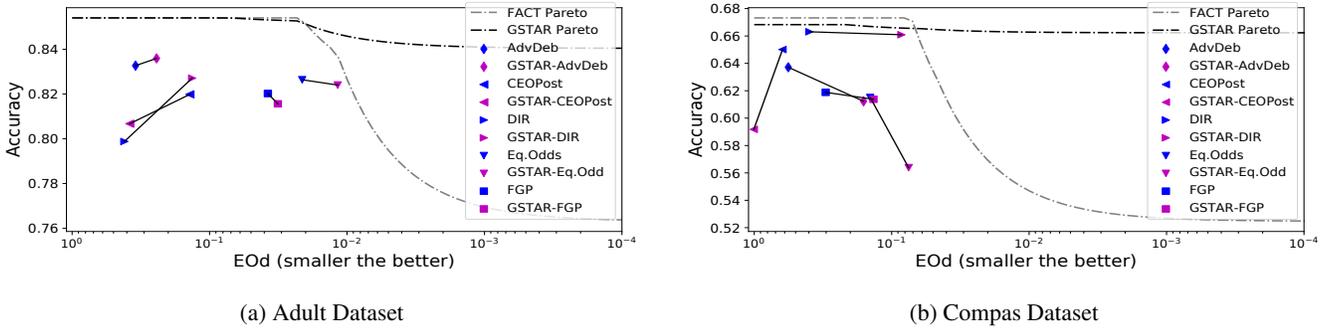


Figure 4: Illustration of post-processing (magenta colored points) on existing fairness models (blue colored points). Given the outputs of each model, GSTAR efficiently improves most existing fairness models with optimized group-aware thresholds.

CelebA dataset (Liu et al. 2015), Adult dataset (Kohavi 1996), COMPAS<sup>1</sup> dataset, and German dataset (Dua and Graff 2019). More details of the comparing methods, evaluation metrics, and datasets are provided in the Supplementary.

### Performance and Fairness-Accuracy Trade-Offs

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

performance.

In Fig. 2, we plot Pareto frontier, which is the upper bound for the accuracy-fairness trade-offs, desired output locates at the upper right region under the boundary which corresponds to higher values in accuracy and lower values in fairness discrepancy. With the same fairness constraints are given, we achieve a better frontier than the FACT (Kim, Chen, and Talwalkar 2020) as we equally weigh on demographic statistics and have a better feasible region. To obtain our results (star points), we first estimate the logit distribution from the output of the baseline model, and then we get optimal adaptive thresholds with corresponding fairness

<sup>1</sup><https://github.com/propublica/compas-analysis>

metric by updating w.r.t. the objective function in (2). Here we have three combinations of fairness imposed to GSTAR: demographic parity (DP), equalized odds (EOd), and with both constraints (DP+EOd). By post-processing on a simple baseline, we achieved significantly better fairness with small or no sacrifice in accuracy. In all datasets, GSATR got competitive or better results than other state-of-the-art methods on both fairness and accuracy.

For example, we got  $\theta_{EOd}^* = (1.570, -0.525)^\top$  for the CelebA dataset. This shows that we have a higher threshold for the privileged group and a lower threshold for the unprivileged group. This optimal thresholding from GSTAR allows more samples from the privileged group to be correctly predicted as unattractive that would compensate for the discrimination of the original model. In other words, this improves false positive rate difference (also known as predictive equality (Chouldechova 2017)) with a huge amount from 0.235 to 0.014. Also, true positive rate difference (also known as equality of opportunity (Hardt, Price, and Srebro 2016)) got reduced from 0.282 to 0.018. It is notable that GSTAR only sacrificed 2.2% of accuracy to bring the big improvement in fairness.

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

### Flexibility and Multiple Fairness Constraints

Since each fairness metric has different interests, it has been theoretically proven that they cannot be perfectly satisfied all together (Pleiss et al. 2017; Chouldechova 2017; Kleinberg, Mullainathan, and Raghavan 2016). Because of this inherent trade-offs between fairness metrics, most of the recent works focus on a single metric at a time to achieve fairness. However with GSTAR, we have the flexibility to optimize on multiple fairness constraints that can be represented in the confusion matrix format. Moreover, given the estimated distribution  $f_{ya}$  of a arbitrary classification model, we can adjust the optimal  $\theta$  based on the needs by accommodating different fairness criteria.

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

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

### Post-Processing on an Existing Fair Model

For a binary classifier that has a single fixed classification threshold (0 for out logit, and 0.5 for label probability), we can provide better trade-off between fairness and accuracy with GSTAR. Given the logit/probability in the model-agnostic manner, we can improve the fairness as illustrated in Fig. 4. In most cases, we observe improvement in fairness after GSTAR post-processing. It is also interesting to note that by optimizing the different thresholds for each protected group, we even obtain better performance on both fairness and accuracy, which indicates that the threshold optimization can not only improve fairness but also accuracy.

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

### Conclusion and Discussion

In this paper, we propose a group-aware threshold adaptation method (GSTAR) to post-process in model-agnostic manner and optimize over multiple fairness constraints. We directly optimize the classification threshold for each demographic group w.r.t. the classification error and multiple fairness constraints in a unified objective function, such that we can practically achieve an optimal trade-off between accuracy and fairness in fair classification. Our method is applicable to diverse notions of group fairness as the majority of fairness notions can be expressed as a linear or quadratic equation through confusion matrix. We empirically show that GSTAR is *flexible* with fairness regularization, *efficient* with low computational cost. We also notice that the adaptive thresholds benefit accuracy in some cases. GSTAR agrees to protect *privacy* such as article 17 of EU’s GDPR (Regulation 2016). We only require the estimated distribution of the output from a given model i.e., our post-processing method is oblivious to features. Thus training data is no longer needed and allowed to be discarded after training the model that to be post-processed. Thus, GSTAR can be applied to relaxed scenarios where practitioners cannot access individual-level sensitive information but have estimated distributions of logits for each sensitive group.

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

## Acknowledgements

This work was partially supported by NSF IIS #1955890, Purdue's Elmore ECE Emerging Frontiers Center.

## References

- Barocas, S.; and Selbst, A. D. 2016. Big data's disparate impact. *Calif. L. Rev.*, 104: 671.
- Chouldechova, A. 2017. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2): 153–163.
- Corbett-Davies, S.; Pierson, E.; Feller, A.; Goel, S.; and Huq, A. 2017. Algorithmic decision making and the cost of fairness. In *KDD*, 797–806.
- Dressel, J.; and Farid, H. 2018. The accuracy, fairness, and limits of predicting recidivism. *Sci. Adv.*, 4(eaao5580): 1–5.
- Dua, D.; and Graff, C. 2019. UCI Machine Learning Repository. *University of California, Irvine, School of Information and Computer Sciences*.
- Feldman, M.; Friedler, S. A.; Moeller, J.; Scheidegger, C.; and Venkatasubramanian, S. 2015. Certifying and removing disparate impact. In *KDD*, 259–268.
- Gratton, S.; Lawless, A. S.; and Nichols, N. K. 2007. Approximate Gauss–Newton methods for nonlinear least squares problems. *SIAM*, 18(1): 106–132.
- Hardt, M.; Price, E.; and Srebro, N. 2016. Equality of opportunity in supervised learning. In *NIPS*, 3315–3323.
- He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *CVPR*, 770–778.
- Jang, T.; Zheng, F.; and Wang, X. 2021. Constructing a fair classifier with generated fair data. In *AAAI*, volume 35, 7908–7916.
- Kamiran, F.; Karim, A.; and Zhang, X. 2012. Decision theory for discrimination-aware classification. In *ICDM*, 924–929. IEEE.
- Kim, J. S.; Chen, J.; and Talwalkar, A. 2020. Model-Agnostic Characterization of Fairness Trade-offs. *arXiv preprint arXiv:2004.03424*.
- Kleinberg, J.; Mullainathan, S.; and Raghavan, M. 2016. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*.
- Kohavi, R. 1996. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In *KDD*, volume 96, 202–207.
- Liu, L. T.; Simchowitz, M.; and Hardt, M. 2019. The implicit fairness criterion of unconstrained learning. In *ICML*, 4051–4060.
- Liu, Z.; Luo, P.; Wang, X.; and Tang, X. 2015. Deep learning face attributes in the wild. In *ICCV*, 3730–3738.
- Lokhande, V. S.; Akash, A. K.; Ravi, S. N.; and Singh, V. 2020. FairALM: Augmented Lagrangian Method for Training Fair Models with Little Regret. In *ECCV*, 365–381. Springer.
- Madras, D.; Creager, E.; Pitassi, T.; and Zemel, R. 2018. Learning adversarially fair and transferable representations. *arXiv preprint arXiv:1802.06309*.
- Menon, A. K.; and Williamson, R. C. 2018. The cost of fairness in binary classification. In *ACM FAccT*, 107–118.
- Pedreshi, D.; Ruggieri, S.; and Turini, F. 2008. Discrimination-aware data mining. In *KDD*, 560–568.
- Pleiss, G.; Raghavan, M.; Wu, F.; Kleinberg, J.; and Weinberger, K. Q. 2017. On fairness and calibration. In *NIPS*, 5680–5689.
- Regulation, G. D. P. 2016. Regulation EU 2016/679 of the European Parliament and of the Council of 27 April 2016. *OJEU*, 43–44.
- Tan, Z.; Yeom, S.; Fredrikson, M.; and Talwalkar, A. 2020. Learning fair representations for kernel models. In *AISTATS*, 155–166.
- Zhang, B. H.; Lemoine, B.; and Mitchell, M. 2018. Mitigating unwanted biases with adversarial learning. In *AIES*, 335–340.
- Zhao, H.; and Gordon, G. 2019. Inherent tradeoffs in learning fair representations. In *NeurIPS*, 15675–15685.