
Supplementary Material for “Group-Aware Threshold Adaptation for Fair Classification”

Anonymous Author(s)

Affiliation

Address

email

1 Upper Bounds on False-Positive/Negative Rate Gap Between Groups

2 1.1 Notations

3 We start from defining notations. We denote $f_{ya}(x)$ for the estimated parametric probability density
4 function (PDF) of the distribution of output logit h in the subset $\{Y = y, A = a\}$. Correspondingly,
5 we denote the corresponding cumulative distribution function (CDF) as

$$F_{ya}(x) = \int_{-\infty}^x f_{ya}(x)dx.$$

6 We use $F_{ya}^{-1}(x)$ to denote the inverse of the CDF.

7 Then, following the definitions given in the main paper, we have

$$\begin{aligned} \text{TP}_a(\theta_a) &= 1 - F_{1a}(\theta_a), & \text{FN}_a(\theta_a) &= F_{1a}(\theta_a), \\ \text{FP}_a(\theta_a) &= 1 - F_{0a}(\theta_a), & \text{TN}_a(\theta_a) &= F_{0a}(\theta_a). \end{aligned} \quad (1)$$

8 1.2 Characterizing the accuracy loss function under perfect EOp condition

9 Before stating the theorem, we illustrate the difference between $\mathcal{L}_{per}(\theta)$ used in our paper versus loss
10 function one would use in a population-wise classification problem (without considering group-aware
11 thresholds). That is, one would only consider the loss function on accuracy

$$\bar{\mathcal{L}}_{per}(\theta) = (r_1\bar{\text{FN}}(\theta) + r_0\bar{\text{FP}}(\theta))^2, \quad (2)$$

12 where only one threshold θ (for both groups) needs to be decided, $r_y = (n_{y0} + n_{y1})/N$ is the
13 population ratio of data samples with label y , $\bar{\text{FN}}(\theta)$, $\bar{\text{FP}}(\theta)$ are the population-wise false-negative
14 and false-positive rate. $\bar{\text{FN}}(\theta)$, $\bar{\text{FP}}(\theta)$ are defined in a similar way as in (1) except that we just use
15 the population-wise pdf $\bar{f}_y(x)$ in the integral for label y . (2) will be our benchmark to compare with
16 $\mathcal{L}_{per}(\theta)$ used in our paper.

17 We start from considering the case that we achieve perfect EOp condition, that is

$$\text{TP}_1(\theta_1) = \text{TP}_0(\theta_0), \quad \text{or equivalently} \quad \text{FN}_1(\theta_1) = \text{FN}_0(\theta_0). \quad (3)$$

18 This means that θ_0 and θ_1 satisfies the following condition

$$F_{11}(\theta_1) = F_{10}(\theta_0). \quad (4)$$

19 Equivalently, we have

$$\theta_0 = F_{10}^{-1}(F_{11}(\theta_1)). \quad (5)$$

20 Under any given pair of (θ_0, θ_1) that satisfies (5), recall that the performance error $\mathcal{L}_{per}(\boldsymbol{\theta})$ is defined
 21 as

$$\mathcal{L}_{per}(\boldsymbol{\theta}) = \left(\frac{n_{01}}{N} \text{FP}_1(\theta_1) + \frac{n_{11}}{N} \text{FN}_1(\theta_1) + \frac{n_{00}}{N} \text{FP}_0(\theta_0) + \frac{n_{10}}{N} \text{FN}_0(\theta_0) \right)^2. \quad (6)$$

22 From (3), we get

$$\frac{n_{11}}{N} \text{FN}_1(\theta_1) + \frac{n_{10}}{N} \text{FN}_0(\theta_0) = \frac{n_{11} + n_{10}}{N} \text{FN}(\theta_1) = r_1 \text{FN}(\theta_1),$$

23 where r_1 denotes, over the entire population (across different groups), proportion of samples with
 24 positive labels. In other words, $r_1 \text{FN}(\theta_1)$ represents the proportion of data samples (from both
 25 groups) with positive label but falsely classified as negative out of the entire dataset.

26 Next, we look at the other two terms:

$$\frac{n_{01}}{N} \text{FP}_1(\theta_1) + \frac{n_{00}}{N} \text{FP}_0(\theta_0).$$

27 This sum can be written as

$$\begin{aligned} \frac{n_{01}}{N} \text{FP}_1(\theta_1) + \frac{n_{00}}{N} \text{FP}_0(\theta_0) &= \frac{n_{01} + n_{00}}{N} \text{FP}_1(\theta_1) + \frac{n_{00}}{N} (\text{FP}_0(\theta_0) - \text{FP}_1(\theta_1)) \\ &= r_0 \text{FP}(\theta_1) + \frac{n_{00}}{N} (\text{FP}_0(\theta_0) - \text{FP}_1(\theta_1)). \end{aligned}$$

28 We denote $\epsilon_1 = (\text{FP}_0(\theta_0) - \text{FP}_1(\theta_1))$. Hence,

$$\mathcal{L}_{per}(\boldsymbol{\theta}) = \mathcal{L}_{per}(\theta_1) = \left(r_1 \text{FN}(\theta_1) + r_0 \text{FP}(\theta_1) + \frac{n_{00}}{N} \epsilon_1 \right)^2. \quad (7)$$

29 Comparing (2) with (7), we can see that, when $\text{FP}_0(\theta_0) > \text{FP}_1(\theta_1)$, the term $\frac{n_{00}}{N} \epsilon_1$ captures the
 30 additional accuracy loss due to that we have chosen two different thresholds even though that
 31 condition (4) is satisfied. Next, we characterize an upper bound for ϵ_1 .

32 1.3 Theorem 1 and its Proof

33 We first state the assumptions we need to make for Theorem 1.

34 **Assumption 1.** For any given classifier h and its induced parametric PDF f_{y_a} and CDF F_{y_a} , we
 35 assume the following holds:

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

$$|F_{y_1}(x) - F_{y_0}(x)| \leq u_y, \quad \forall x.$$

38 **Theorem 1.** For any given classifier that satisfies Assumption 1 and any given pair of thresholds
 39 (θ_0, θ_1) that satisfies the perfect EOp condition, the gap between false-positive rates of the two group
 40 is upper bounded by

$$|\epsilon_1| = |\text{FP}_0(\theta_0) - \text{FP}_1(\theta_1)| \leq u_0 + C_1 u_1, \quad (8)$$

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

42 *Proof.* Recall that $\text{FP}_1(\theta_1) = 1 - F_{01}(\theta_1)$ and $\text{FP}_0(\theta_0) = 1 - F_{00}(\theta_0)$. Hence,

$$\begin{aligned} |\text{FP}_0(\theta_0) - \text{FP}_1(\theta_1)| &= |F_{01}(\theta_1) - F_{00}(\theta_0)| \\ &\leq |F_{01}(\theta_1) - F_{01}(\theta_0)| + |F_{01}(\theta_0) - F_{00}(\theta_0)|. \end{aligned}$$

43 To bound ϵ , we just need to bound $|F_{01}(\theta_1) - F_{01}(\theta_0)|$ and $|F_{01}(\theta_0) - F_{00}(\theta_0)|$.

For the second one, we note that from Assumption 1 that

$$|F_{01}(\theta_0) - F_{00}(\theta_0)| \leq u_0.$$

44 For the first one, we note that

$$|F_{01}(\theta_1) - F_{01}(\theta_0)| \leq \hat{f}_{01}|\theta_1 - \theta_0|,$$

45 where $\hat{f}_{01} = \max_x f_{01}(x)$.

46 Next, we bound $|\theta_1 - \theta_0|$. Note that from (5),

$$\begin{aligned} |\theta_1 - \theta_0| &= |F_{10}^{-1}(F_{11}(\theta_1)) - \theta_1| \\ &= \left| F_{10}^{-1}(F_{11}(\theta_1)) - F_{10}^{-1}(F_{10}(\theta_1)) \right| \\ &\leq M_{10}|F_{11}(\theta_1) - F_{10}(\theta_1)| \\ &\leq M_{10}u_1. \end{aligned}$$

47

□

48 Theorem 1 provides an upper bound on the difference in the false positive rate between the two
49 groups, for any given pair of (θ_0, θ_1) such that the false negative rates are the same for the two
50 groups (i.e., satisfies the perfect EOp condition). As discussed in Section 1.2, this upper bound also
51 characterizes the additional accuracy loss due to that we have group-dependent thresholds compared
52 to the case with only one threshold for both groups.

53 1.4 Under perfect PE condition

54 For predictive equality (PE) condition, we prove a similar result. That is, assuming we achieve perfect
55 PE condition with

$$FP_1(\theta_1) = FP_0(\theta_0), \quad \text{or equivalently} \quad TN_1(\theta_1) = TN_0(\theta_0). \quad (9)$$

56 This means that θ_0 and θ_1 satisfies the following condition

$$F_{01}(\theta_1) = F_{00}(\theta_0). \quad (10)$$

57 Equivalently, we have

$$\theta_0 = F_{00}^{-1}(F_{01}(\theta_1)). \quad (11)$$

58 Under any given pair of (θ_0, θ_1) that satisfies (11), the performance error $\mathcal{L}_{per}(\theta)$ can be written as

$$\begin{aligned} \mathcal{L}_{per}(\theta) &= \left(\frac{n_{01}}{N}FP_1(\theta_1) + \frac{n_{11}}{N}FN_1(\theta_1) + \frac{n_{00}}{N}FP_0(\theta_0) + \frac{n_{10}}{N}FN_0(\theta_0) \right)^2 \\ &= \left(r_1FN(\theta_1) + r_0FP(\theta_1) + \frac{n_{10}}{N}\epsilon_2 \right)^2, \end{aligned}$$

59 where

$$\epsilon_2 = (FN_0(\theta_0) - FN_1(\theta_1)).$$

60 Similar to Theorem 1, we can provide an upper bound on ϵ_2 under Assumption 1.

61 **Theorem 2.** For any given classifier that satisfies Assumption 1 and any given pair of thresholds
62 (θ_0, θ_1) that satisfies the perfect PE condition, the gap between false-negative rates of the two group
63 is upper bounded by

$$|\epsilon_2| = |FN_0(\theta_0) - FN_1(\theta_1)| \leq u_1 + C_0u_0, \quad (12)$$

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

65 *Proof.* The proof is similar to that of Theorem 1. We provide the main steps and omit details that
66 repeat with the proof of Theorem 1. We have

$$\begin{aligned} |FN_0(\theta_0) - FN_1(\theta_1)| &= |F_{11}(\theta_1) - F_{10}(\theta_0)| \\ &\leq |F_{11}(\theta_1) - F_{11}(\theta_0)| + |F_{11}(\theta_0) - F_{10}(\theta_0)| \\ &\leq \hat{f}_{11}|\theta_1 - \theta_0| + u_1 \\ &\leq \hat{f}_{11}M_{00}u_0 + u_1. \end{aligned}$$

67

□

68 Theorem 2 provides an upper bound on the difference in the false negative rate between the two
 69 groups, for any given pair of (θ_0, θ_1) such that the false positive rates are the same for the two groups
 70 (i.e., satisfies the perfect PE condition).

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

75 2 Characterizing the Tradeoff between Accuracy and Fairness

76 In this section, we prove a theorem to characterize the tradeoff between accuracy and fairness. That
 77 is, we start from the perfect EOP or PE conditions and perturb the solution by a small amount. We
 78 then bound the difference in the accuracy loss by comparing the perturbed solution with the original
 79 solution that satisfies the perfect fairness conditions.

80 2.1 Perturbed EOP condition

81 To start with, let us consider solutions (θ_0, θ_1) that satisfy the perfect EOP condition (5). Under this
 82 condition, the optimization problem becomes one dimensional, that is,

$$\theta_1^* = \operatorname{argmin}_{\theta_1} \mathcal{L}_{per}(\theta_1),$$

83 where

$$\mathcal{L}_{per}(\theta_1) = \left(r_1 \text{FN}_1(\theta_1) + r_0 \text{FP}_1(\theta_1) + \frac{n_{00}}{N} \epsilon_1(\theta_1) \right)^2 \quad (13)$$

84 and

$$\epsilon_1(\theta_1) = \text{FP}_0(\theta_0) - \text{FP}_1(\theta_1) = F_{01}(\theta_1) - F_{00}(F_{10}^{-1}(F_{11}(\theta_1))).$$

85 From θ_1^* , we can get the corresponding $\theta_0^* = F_{10}^{-1}(F_{11}(\theta_1^*))$. We further denote this optimal accuracy
 86 loss value as

$$L^* = \mathcal{L}_{per}(\theta_1^*).$$

87 Now with the optimal solution (θ_0^*, θ_1^*) , we investigate the changes in $\mathcal{L}_{per}(\theta_1^*)$ when we perturb the
 88 perfect EOP condition and allow a small difference. That is, now consider solution $(\tilde{\theta}_0, \tilde{\theta}_1)$ such that

$$|\text{FN}_1(\theta_1^*) - \text{FN}_1(\tilde{\theta}_1)| \leq \gamma/2, \quad |\text{FN}_0(\theta_0^*) - \text{FN}_0(\tilde{\theta}_0)| \leq \gamma/2. \quad (14)$$

89 Consequently, the solution $(\tilde{\theta}_0, \tilde{\theta}_1)$ satisfy the following perturbed EOP condition:

$$|\text{TP}_1(\tilde{\theta}_1) - \text{TP}_0(\tilde{\theta}_0)| = |\text{FN}_1(\tilde{\theta}_1) - \text{FN}_0(\tilde{\theta}_0)| \leq \gamma. \quad (15)$$

90 Without loss of generality, we assume that (i) the true positive rate of group 1 is higher than that
 91 of group 0, and (ii) the above inequality is binding (because if not binding, then we can always
 92 choose a smaller γ to make it binding). Thus, we have $\text{TP}_1(\tilde{\theta}_1) - \text{TP}_0(\tilde{\theta}_0) = \gamma$, or equivalently,
 93 $\text{FN}_0(\tilde{\theta}_0) - \text{FN}_1(\tilde{\theta}_1) = \gamma$. This gives us

$$\tilde{\theta}_0 = F_{10}^{-1}(F_{11}(\tilde{\theta}_1) + \gamma). \quad (16)$$

94 Next, we analyze $\mathcal{L}_{per}(\tilde{\theta}_1)$ by substituting $(\tilde{\theta}_0, \tilde{\theta}_1)$ in (6), which gives us

$$\mathcal{L}_{per}(\tilde{\theta}_1) = \left(r_1 \text{FN}_1(\tilde{\theta}_1) + r_0 \text{FP}_1(\tilde{\theta}_1) + \frac{n_{10}}{N} \gamma + \frac{n_{00}}{N} \tilde{\epsilon}_1(\tilde{\theta}_1) \right)^2, \quad (17)$$

95 where

$$\tilde{\epsilon}_1(\tilde{\theta}_1) = \text{FP}_0(\tilde{\theta}_0) - \text{FP}_1(\tilde{\theta}_1) = F_{01}(\tilde{\theta}_1) - F_{00}(F_{10}^{-1}(F_{11}(\tilde{\theta}_1) + \gamma)).$$

96 We denote the optimal value for this perturbed version as $\tilde{\theta}_1^*$, and its corresponding loss value as

$$\tilde{L}^* = \mathcal{L}_{per}(\tilde{\theta}_1^*).$$

97 Furthermore, from (14), we have

$$|\text{FN}_1(\theta_1^*) - \text{FN}_1(\tilde{\theta}_1^*)| = |F_{11}(\theta_1^*) - F_{11}(\tilde{\theta}_1^*)| \leq \gamma/2. \quad (18)$$

98 Under Assumption 1, we have

$$\begin{aligned} |\theta_1^* - \tilde{\theta}_1^*| &= |F_{11}^{-1}(F_{11}(\theta_1^*)) - F_{11}^{-1}(F_{11}(\tilde{\theta}_1^*))| \\ &\leq M_{11}|F_{11}(\theta_1^*) - F_{11}(\tilde{\theta}_1^*)| \\ &= M_{11}\gamma/2. \end{aligned}$$

99 2.2 Theorem 3 and its proof

100 We are ready to compare $\mathcal{L}_{per}(\theta_1^*)$ and $\mathcal{L}_{per}(\tilde{\theta}_1^*)$. The latter loss should be no larger than the former
101 since we relaxed the perfect EOp condition (constraint) in the optimization, i.e., $L^* \geq \tilde{L}^*$.

102 **Theorem 3.** *Under Assumption 1 and condition (14),*

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

103 where $C = 2L^* \left(\frac{r_1}{2} + r_0 \frac{\hat{f}_{01}M_{11}}{2} + \frac{n_{00}}{N} \left(\hat{f}_{00}M_{10} + \frac{\hat{\epsilon}'_1 M_{11}}{2} \right) + \frac{n_{10}}{N} \right)$, and $\hat{\epsilon}'_1 = \max \tilde{\epsilon}'_1$ is the maxi-
104 mum of the derivative of $\tilde{\epsilon}_1$.

105 *Proof.* We have that

$$\begin{aligned} &\mathcal{L}_{per}(\theta_1^*) - \mathcal{L}_{per}(\tilde{\theta}_1^*) \\ &\leq 2L^* \left| r_1 \text{FN}_1(\theta_1^*) + r_0 \text{FP}_1(\theta_1^*) + \frac{n_{00}}{N} \epsilon_1(\theta_1^*) \right. \\ &\quad \left. - \left(r_1 \text{FN}_1(\tilde{\theta}_1^*) + r_0 \text{FP}_1(\tilde{\theta}_1^*) + \frac{n_{10}}{N} \gamma + \frac{n_{00}}{N} \tilde{\epsilon}_1(\tilde{\theta}_1^*) \right) \right| \\ &\leq 2L^* \left(r_1 \gamma/2 + r_0 |\text{FP}_1(\theta_1^*) - \text{FP}_1(\tilde{\theta}_1^*)| + \frac{n_{00}}{N} |\epsilon_1(\theta_1^*) - \tilde{\epsilon}_1(\tilde{\theta}_1^*)| + \frac{n_{10}}{N} \gamma \right), \end{aligned}$$

106 where we further have that

$$\begin{aligned} |\text{FP}_1(\theta_1^*) - \text{FP}_1(\tilde{\theta}_1^*)| &= |F_{01}(\theta_1^*) - F_{01}(\tilde{\theta}_1^*)| \\ &\leq \hat{f}_{01} |\theta_1^* - \tilde{\theta}_1^*| \\ &\leq \hat{f}_{01} M_{11} \gamma/2, \end{aligned}$$

107 and

$$\begin{aligned} |\epsilon_1(\theta_1^*) - \tilde{\epsilon}_1(\tilde{\theta}_1^*)| &\leq |\epsilon_1(\theta_1^*) - \tilde{\epsilon}_1(\theta_1^*)| + |\tilde{\epsilon}_1(\theta_1^*) - \tilde{\epsilon}_1(\tilde{\theta}_1^*)| \\ &\leq |F_{00}(F_{10}^{-1}(F_{11}(\theta_1^*))) - F_{00}(F_{10}^{-1}(F_{11}(\theta_1^*) + \gamma))| + \hat{\epsilon}'_1 M_{11} \gamma/2 \\ &= (\hat{f}_{00} M_{10} + \hat{\epsilon}'_1 M_{11}/2) \gamma. \end{aligned}$$

108 Here, $\hat{\epsilon}'_1 = \max \tilde{\epsilon}'_1$ is the maximum of the derivative of $\tilde{\epsilon}_1$. Combining all the terms in front of γ
109 gives us the desired upper bound. \square

110 Theorem 3 quantifies the decrease in accuracy loss (i.e., the improvement in accuracy) when we allow
111 a gap of true positive rates between two groups (i.e., relaxation from the perfect EOp condition).

112 Repeating the analysis for the perturbed PE condition, we can obtain a similar bound for the changes
113 in the accuracy loss function. We omit the details here in the interest of space.

114 3 Experimental Details

115 3.1 Comparing Methods

116 We compared our method with multiple state-of-the-art methods to verify our work. The details about
117 the comparing methods are as below:

- 118 • **Learning fair representations for kernel models** (abbreviated as FGP) [11]: a pre-
119 processing method to learn representation focusing on kernel-based models. The fair
120 model that satisfies certain fairness criterion is obtained by Bayesian learning from fair
121 Gaussian process (FGP) prior.
- 122 • **Fairness confusion tensor** (abbreviated as FACT) [5]: a post-processing model that mini-
123 mize the least-squares accuracy-fairness optimality problem based on confusion tensor.
- 124 • **Adversarial de-biasing** (abbreviated as AdvDeb) [12]: an in-processing model that miti-
125 gates the conflicting gradient directions in utility and fairness objectives by projecting one
126 gradient to another to remove the opposite direction.
- 127 • **Calibrated equal odds post-processing** (abbreviated as CEOPost) [9]: a post-processing
128 method that minimizes the disparity in the predicted probability to the preferred class among
129 different sensitive groups, while maintaining the calibration condition in a relaxed condition.
- 130 • **Equality of opportunity in supervised learning** (abbreviated as Eq.Odds) [3]: a post-
131 processing method that learns the threshold to yield the equalized odds/opportunity between
132 different demographic by exploring the intersection of achievable true positive rates and
133 false positive rates.
- 134 • **Learning adversarially fair and transferable representations** (abbreviated as
135 LAFTR) [8]: a fair representation learning model that adopts fairness metrics as the adver-
136 sarial objectives and analyze the balance between utility and fairness.
- 137 • **Baseline**: For CelebA dataset, we use ResNet50 [4] as a reference because we input second
138 last layer(2048 features) of ResNet to all methods. For other tabular datasets, logistic
139 regression is used as all other methods except for FGP and LAFTR are based on logistic
140 regression.

141 If the hyperparameter is adjustable in the listed methods, we report the result with the fairness
142 coefficient that has an accuracy closest to the average accuracy in the coefficient sweep to balance
143 utility and fairness.

144 3.2 Evaluation Metrics

145 In the experiments, we evaluate the methods on four fairness and two performance measures. Four
146 fairness metrics are as below:

- **Equal Opportunity** (abbreviated as EOp) [3] : This measures absolute difference of favor-
able prediction given positive label.

$$|P(\hat{Y} = 1|Y = 1, A = 1) - P(\hat{Y} = 1|Y = 1, A = 0)|.$$

- 147 • **Equalized Odds** (abbreviated as EOd) [3] : This measures the difference between the
148 probability given the true labels.

$$|P(\hat{Y} = 1|Y = 1, A = 1) - P(\hat{Y} = 1|Y = 1, A = 0)| + \\ |P(\hat{Y} = 1|Y = 0, A = 1) - P(\hat{Y} = 1|Y = 0, A = 0)|.$$

- 149 • **Balanced Accuracy Difference** (abbreviated as BD) : This measures difference between
150 balanced accuracy between the groups.

$$|P(\hat{Y} = 1|Y = 1, A = 1) + P(\hat{Y} = 0|Y = 0, A = 1)| \\ - |P(\hat{Y} = 1|Y = 1, A = 0) + P(\hat{Y} = 0|Y = 0, A = 0)|.$$

151 If BD and EOd has the same value, it indicates that both TPR and TNR are higher in a certain
152 sensitive group. However, if the gap between the two terms is large, we can interpret as the
153 classifier is more biased as a group with higher TPR has lower TNR. This is more unfair as
154 a sample from the privileged group is more likely to be falsely and correctly predicted as
155 positive output. EOp is a partial measure of EOd as it only measures the difference from a
156 favorable class.

CelebA				
Model	GSTAR	FGP	FACT	CEOPost
Time	0.287	-	0.067	0.077
Model	DIR	Eq.Odds	LAFTR	AdvDeb
Time	183.20	0.062	107.04(min)	303.15
Adult				
Model	GSTAR	FGP	FACT	CEOPost
Time	0.29	51.28	0.055	25.61
Model	DIR	Eq.Odds	LAFTR	AdvDeb
Time	168	0.037	53.04(min)	102.00
Compas				
Model	GSTAR	FGP	FACT	CEOPost
Time	0.292	43.74	0.035	8.3
Model	DIR	Eq.Odds	LAFTR	AdvDeb
Time	123.20	0.034	57.04(min)	15.45
German				
Model	GSTAR	FGP	FACT	CEOPost
Time	0.271	7.08	0.0257	2.64
Model	DIR	Eq.Odds	LAFTR	AdvDeb
Time	1.68	0.034	56.51(min)	2.17

Table 1: Computational time (in seconds) for all comparing fairness methods for each dataset.

- 157 • **Absolute (1 - Disparate Impact)** (abbreviated as 1-DIMP) [1] : This measures ratio of the
158 probability of the favorable prediction given a protected group.

$$\left| 1 - \frac{P(\hat{Y} = 1|A = 1)}{P(\hat{Y} = 1|A = 0)} \right|.$$

159 We evaluate performance of the methods with two metrics.

- 160 • **Balanced Accuracy** (abbreviated as BA) : This measures average between true positive rate
161 and true negative rate. Compared to the traditional accuracy, this measure effectively shows
162 the whether the classifier is focusing on the performance of a certain class when the dataset
163 is unbalanced.

$$\frac{1}{2} \left(P(\hat{Y} = 1|Y = 1) + P(\hat{Y} = 0|Y = 0) \right).$$

- 164 • **Accuracy** (abbreviated as ACC) : This measures traditional classification accuracy of the
165 method.

166 3.3 Dataset Description

167 We evaluate the methods on four fairness datasets. The goal for all datasets is binary classification on
168 binary sensitive feature. The details of the datasets are as below:

- 169 • **CelebA image dataset**¹ [7]: The data consists of 202,599 face images in diverse demo-
170 graphics. The images are annotated with 40 attributes (face shape, skin tone, smiling, etc.).
171 Similar to Quadrianto *et al.* [10], the goal is to predict whether a person in the image is
172 attractive or not. The feature *sex* is used as the sensitive feature.
- 173 • **Adult** dataset from the UCI repository [6] contains 48,842 instances described by 14 features
174 (workclass, age, education, sex, race, *etc*) with the goal of the income prediction whether a
175 person’s income exceeds 50K USD per year. The feature *sex* is used as the sensitive feature.
- 176 • **COMPAS**²(Correctional Offender Management Profiling for Alternative Sanctions) dataset
177 includes 6,167 samples described by 401 features with the target of recidivism prediction

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

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

178 with the label showing if each person gets rearrested within two years. The feature *race* is
179 used as the sensitive feature for this dataset.
180 • **German** credit dataset from the UCI repository [2] contains 1,000 samples described by
181 20 features. The goal is to predict the credit risks. The feature *sex* is used as the sensitive
182 feature.

183 3.4 Computational Cost

184 In Table 1, we describe the computational time for each method on each dataset. By introducing
185 estimated PDF functions for post-processing, we outperform other methods except Eq.Odds [3] and
186 FACT [5]. As they both only utilize the entries of the confusion matrix to find optimal mixing rate in
187 their methods, they have less computation than ours. However, as we discussed in the main paper, we
188 explore better feasible region than theirs by group-specific thresholding that results better in both
189 fairness and performance by sacrificing little efficiency, yet outperforms most of the other works.

190 References

- 191 [1] Solon Barocas and Andrew D Selbst. Big data’s disparate impact. *Calif. L. Rev.*, 104:671, 2016.
- 192 [2] Dheeru Dua and Casey Graff. UCI machine learning repository, 2019.
- 193 [3] Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. In
194 *NeurIPS*, pages 3315–3323, 2016.
- 195 [4] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
196 recognition. In *CVPR*, pages 770–778, 2016.
- 197 [5] Joon Sik Kim, Jiahao Chen, and Ameet Talwalkar. Model-agnostic characterization of fairness
198 trade-offs. *arXiv preprint arXiv:2004.03424*, 2020.
- 199 [6] Ron Kohavi. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In
200 *KDD*, volume 96, pages 202–207, 1996.
- 201 [7] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the
202 wild. In *ICCV*, pages 3730–3738, 2015.
- 203 [8] David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair
204 and transferable representations. *arXiv preprint arXiv:1802.06309*, 2018.
- 205 [9] Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On fairness
206 and calibration. In *NeurIPS*, pages 5680–5689, 2017.
- 207 [10] Novi Quadrianto, Viktoriia Sharmanska, and Oliver Thomas. Discovering fair representations
208 in the data domain. In *CVPR*, pages 8227–8236, 2019.
- 209 [11] Zilong Tan, Samuel Yeom, Matt Fredrikson, and Ameet Talwalkar. Learning fair representations
210 for kernel models. In *AISTATS*, pages 155–166, 2020.
- 211 [12] Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with
212 adversarial learning. In *AIES*, pages 335–340, 2018.