Supplement: Scalable and Stable Surrogates for Flexible Classifiers with Fairness Constraints

Anonymous Author(s) Affiliation Address email

1 A Experimental Protocols

For each experiment, we begin by observing the latent fairness discrepancy in an unconstrained
model, and use that to assign group labels *a* and *b*. In our results, we obscure these assignments and

refer to the sensitive groups as a and b in order to reduce implicit biases associated with identity
groups, as in Denton et al. [7].

6 A.1 Fair Logistic Regression for Tabular Data

⁷ For fairness-regularized logistic regression, we optimize Eq. (11) with negative log-likelihood loss ⁸ function, linear model class $f_{\theta}(\mathbf{x}) = \theta^T \mathbf{x}$, and an added L_2 regularizer.

⁹ To preprocess each data set, we make a random $\frac{2}{3} - \frac{1}{3}$ train-test split and apply min-max normalization to improve regularization performance. To get an L_2 penalty coefficient, we perform 5-fold crossvalidation on the training data and search to find a coefficient which maximizes cross-validated AUC.

¹² We searched over the set of powers of 2 from 2^{-15} to 2^{10} . For Fig. 2 only, we use a penalty coefficient

13 of 2^{-20} to better demonstrate the degeneracy of the surrogates.

¹⁴ To set the fairness penalizer λ , we searched in the range [0, 1] to find a λ^* for which Δ was slightly ¹⁵ less than 0 on the training data set. Note that this λ^* varies for each relaxation and fairness criterion, ¹⁶ even on the same data set. We then ran experiments on a dense uniform grid of λ from 0 to λ^* .

All relaxations are optimized via our Lagrangian framework. All code was implemented using PyTorch, and optimized using L-BFGS. $\lambda = 0$ models were initialized at the all-0 parameter vector, and each subsequent model was initialized starting from the solution to the previous λ value.

20 We set the initial learning rate 0.1, which was

²¹ chosen to achieve quick convergence on the unconstrained model. Training was terminated ²³ when every component of θ changed by $< 10^{-8}$ ²⁴ in a single iteration, which took less than 1 ²⁵ minute for every λ on both data sets. Our code

²⁶ is publicly available online.¹

Logistic vs Hinge Convergence. Both Lohaus 27 et al. [13] and Wu et al. [15] use the hinge func-28 tion as their surrogates. The hinge and logistic 29 30 functions have the same asymptotic behavior, but in Fig. 1 we show that Lagrangian optimiza-31 tion of the logistic function is quicker and more 32 consistent due to its smoothness. We compare 33 our logistic upper bound formulation to a hinge 34 upper bound $q(r) = \max(0, 1+r)$, and our lo-35 gistic difference formulation to the rectified lin-36 ear relaxation $q(r) = \max(0, r)$ of Zafar et al. 37



Figure 1: Comparison of convergence rates for logistic and piecewise linear relaxations on tabular data sets. On the left, both use upper bound regularizer R_g^+ to achieve demographic parity on the Adult data set. On the right, the difference framework is used to achieve equality of opportunity on COMPAS. Error bars show the standard deviation of number of iterations until convergence (changes in all parameters drop below 10^{-10}), across 5 independent train-test splits, for the L-BFGS quasi-Newton method. For both relaxations, standard gradient descent is orders of magnitude slower.

Submitted to 35th Conference on Neural Information Processing Systems (NeurIPS 2021). Do not distribute.

¹anonymous.4open.science/r/FairSurrogates-F3EF

³⁸ [16]. We compare over 5 random train-test splits, with each constructed as described above. We ³⁹ found the optimal L_2 penalty coefficient independently for each split.

40 Adult. The Adult data set [12] is one of the most popular in the fair classification literature [1, 4, 6,

8, 10, 11, 13, 15, 17]. The goal of the prediction task is to predict whether an individual is earning

⁴² more or less than \$50,000 per year. We use sex as a sensitive attribute with values *male* and *female*.

43 We use the preprocessed data compiled by [10], which has 30,162 total data points and 100 attributes. 44 After removing the target attribute and sensitive attributes *race*, *sex*, and *race-sex*, we further remove

45 *capital-gain* and *capital-loss* as was done in [13]. This leaves us with p = 94 predictor attributes.

⁴⁶ Here we define equality of opportunity on false negative rates, i.e. predicting that someone earns

under \$50,000 per year when actually they earn more.

48 COMPAS. The COMPAS data set was compiled by Angwin et al. [2], to investigate racial bias in 49 recidivism prediction. The goal of the classification task is to predict whether the defendant will 50 commit another crime or *recidivate* within two years. Importantly, the data set contains information 51 compiled by Freedom of Information Act (FOIA) requests, and does *not* contain all data used by 52 Northpointe in building their predictive model.

As with the Adult data set, we use the preprocessed data of [10]. We drop the target and sensitive attributes to form a data set with p = 399 predictive attributes and 5,273 data points. In this case, equality of opportunity is concerned with false positive rates, where a defendant who will not recidivate is incorrectly labeled high-risk. We use only the subset of the data with sensitive attribute *Caucasian* or *African-American*.

Toy Model. We use binary search [13] to find a model with $\Delta = 0.03$ for each relaxation. Adding an

outlier with y = 1, s = a does not change the decision boundaries for any relaxation. Setting s = b,

60 however, causes the linear relaxation to degenerate.

61 A.2 Fair Deep Learning

For our deep learning experiments, we used the approach of Sec. A.1 to construct a range of λ values. All results are reported on test data. All computation was done on NVIDIA GeForce RTX 2080Ti 11GB GPUs, and GPU times are reported with respect to that hardware

⁶⁴ 11GB GPUs, and GPU times are reported with respect to that hardware.

CelebA. We used the pre-split data provided in torchvision which has 162,770 training images, 19,867 validation images and 19,962 testing images. For our architecture, we used a wide residual network (WRN-50-2) [18] initialized with random low-noise parameter values. The last layer is a soft-max layer, which is mathematically equivalent to logistic regression performed on the attributes learned from the previous layers. Thus in order to enforce fairness, the only change we make is to add the scaled fairness surrogate to the loss function. For this data set, we defined equality of opportunity on false negative rates.

The network was trained on training mini-batches of size 32 for 3 epochs. We used Adam to perform stochastic optimization with an initial learning rate of 0.01, and a scheduler which reduced the learning rate by a factor of 10 when validation loss plateaus for 2000 batches. Each epoch took 25 minutes of GPU time. With 8 relaxations and 21 λ -values per relaxation, the total GPU time was 210 hours.

Faces of the World. The Looking at People CVPR Challenge Track 2 [9] required participants to, given an image, return the bounding box around the face, the subject's gender, and whether or not the subject is smiling. The data set has 6,171 training images, 3,087 validation images and 8,506 test images. The participants were allowed to train on any additional data. The Faces of the World data set shows people from more varied angles than CelebA and is not limited to celebrities, a group that is not representative of the broader population in many physical or sensitive attributes.

We crop the images according to the bounding boxes provided, and resize to 224 by 224 pixels as 83 expected by the WRN. Because the data set is small, we first trained a WRN-50-2 on CelebA using 84 the scheme described above for our CelebA results, and then froze the first two layers of the network 85 to prevent overfitting on Faces of the World. All experiments were initialized at this same baseline. 86 The network was trained on mini-batches of size 32 for 30 epochs using Adam. The initial learning 87 rate was set to 0.01 and a scheduler reduced the learning rate by a factor of 10 when validation loss 88 plateaued for 2000 batches. Each epoch took 1 minute of GPU time. With 3 relaxations and 21 λ 89 values per relaxation, the total GPU time was 32 hours. 90

On top of having performance costs, kernel-based methods have quadratic memory requirements 91 in the number of data points. In order to get around this, Lohaus et al. [13] construct a small set of 92 "reasonable points" and perform learning on those. We use their publicly available code², converting 93 the input images to 150,528-dimensional input vectors (224 by 224 by 3 channels) and run with as 94 many reasonable points as we can hold in 125 gigabytes of RAM. 95

Yelp. We took the subset of reviews from the 5,000 most prolific reviewers, totaling 337,723 reviews. 96 To estimate those reviewers' genders, we use Gender API [3] as in [14]. When no gender can be 97 confidently inferred, we set the sensitive attribute to unknown. 98

The reweighting baseline of Calders & Verwer [5] does not propose a way weighting data with no 99 labelled sensitive attribute. For this baseline, we simply set those points to have weight 1. 100

We modified publicly-available code from Onepoint Consulting³. Our model was initialized to 101 the pretrained BertForSequenceClassification from the Pytorch Transformers library. No 102 parameters were frozen during our training. 103

The learning rate was initialized to $2 \cdot 10^{-5}$ and decreased linearly to 0 over the 3 training epochs. 104 Weight decay with parameter 10^{-3} was used to improve regularization. Each epoch took 200 minutes 105 of GPU time. To save computation time we use the same $\lambda = 0$ training session for every relaxation. 106 With 3 relaxations and 7 λ values per relaxation, as well as the unconstrained training session, the 107 total GPU time was 220 hours. 108

Conditions for Surrogate Degeneracy B 109

For the more general case, we must consider averages of $\Phi(\mathbf{x})_j$ where we zero out negative values, as 110

well as averages where we zero out the positive values. Define: 111

$$\gamma_{jsy} = \frac{1}{N_{sy}} \sum_{\substack{(\mathbf{x}, s', y') \in \mathcal{D} \\ s' = s, y' = y}} \max(0, \Phi(\mathbf{x})_j), \tag{1}$$

$$\eta_{jsy} = \frac{1}{N_{sy}} \sum_{\substack{(\mathbf{x}, s', y') \in \mathcal{D} \\ s' = s, y' = y}} \min(0, \Phi(\mathbf{x})_j).$$
(2)

Note that the normalizing N_{sy} include all data with sensitive label s and target attribute y, even those 112

which the thresholding functions replace with 0. Further let $\gamma_{j,y}$ be the mean of $\Phi(\mathbf{x})_j$ (with negative 113

114

values zeroed out) for data with y' = y but any group, and $\gamma_{js.}^{j.y}$ be the mean of $\Phi(\mathbf{x})_j$ (with negative values zeroed out) for data with s' = s but any target value. Define $\eta_{js.}, \eta_{j.y}$ similarly. Note that 115

when attribute $\Phi(\mathbf{x})_j$ is non-negative, all η_{jsy} are 0 and $\gamma_{jsy} = \mu_{jsy}$. 116

We again assume that g(r) is continuous and monotonically increasing, and that $\delta_q^+ = \lim_{r \to \infty} g'(r)$ 117

118 and
$$\delta_q^- = \lim_{r \to -\infty} g'(r)$$
 exist.

Theorem 4 (General case for degeneracy in demographic parity). *Consider a feature j*. 119

(1) If $\delta_g^-(\gamma_{ja.} - \gamma_{jb.}) + \delta_g^+(\eta_{ja.} - \eta_{jb.}) > 0$ and $\lambda > \lambda_j^* = \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{\delta_g^-(\gamma_{ja.} - \gamma_{jb.}) + \delta_g^+(\eta_{ja.} - \eta_{jb.})}$, then 120 $\lim \mathcal{L}_{a,\lambda}(\theta) = -\infty.$ 121

$$\lim_{\substack{w_{j} \to -\infty \\ y_{j} \to +\infty \\ y_{j} \to +\infty \\ w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ y_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ y_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ y_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ y_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\ } \sum_{w_{j} \to +\infty \\ } \sum_{\substack{w_{j} \to +\infty \\$$

If either (1) or (2) hold and
$$g(r)$$
 is linear, then for this fixed $\Phi(\mathbf{x})$, $\mathcal{L}_{g,\lambda}$ has no stationary points.

²https://github.com/mlohaus/SearchFair, GNU General Public License v3.0

³https://github.com/onepointconsulting/yelp_bert/blob/master/bert_training.ipynb

Proof of Theorems 1 and 4. We begin by proving Theorem 4. Consider the two cases separately.

(1) It is sufficient to show that as $\mathbf{w}_j \to -\infty$, $\frac{\partial}{\partial \mathbf{w}_j} \mathcal{L}_{g,\lambda}(\theta) \to \kappa > 0$.

$$\frac{\partial}{\partial \mathbf{w}_{j}} \mathcal{L}_{g,\lambda}(\theta) = \frac{\partial}{\partial \mathbf{w}_{j}} \left(\frac{1}{N} \sum_{(\mathbf{x},s,y)\in\mathcal{D}} \left(\mathcal{L}(\mathbf{w}^{T}\Phi(\mathbf{x}), y) + \lambda R_{g}(\Phi(\mathbf{x}), y, s) \right) \right) \\
= \frac{\partial}{\partial \mathbf{w}_{j}} \left(\frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ y=1}} -\log \sigma(\mathbf{w}^{T}\Phi(\mathbf{x})) + \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ y=0}} -\log \sigma(-\mathbf{w}^{T}\Phi(\mathbf{x})) \right) \\
+ \frac{\lambda}{N_{a}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ s=a}} g(\mathbf{w}^{T}\Phi(\mathbf{x})) - \frac{\lambda}{N_{b}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ s=b}} g(\mathbf{w}^{T}\Phi(\mathbf{x})) \right) \qquad (3)$$

$$= \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ s=a}} -\sigma(-\mathbf{w}^{T}\mathbf{x})\Phi(\mathbf{x})_{j} + \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ y=0}} \sigma(\mathbf{w}^{T}\Phi(\mathbf{x}))\Phi(\mathbf{x})_{j} \\
+ \frac{\lambda}{N_{a}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ s=a}} g'(\mathbf{w}^{T}\Phi(\mathbf{x}))\Phi(\mathbf{x})_{j} - \frac{\lambda}{N_{b}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D} \\ s=b}} g'(\mathbf{w}^{T}\Phi(\mathbf{x}))\Phi(\mathbf{x})_{j}$$

127 Note that as $\mathbf{w}_j \to -\infty$, $\sigma(-\mathbf{w}^T \Phi(\mathbf{x})) \to 0$ for data with $\Phi(\mathbf{x})_j < 0$ and $\sigma(-\mathbf{w}^T \Phi(\mathbf{x})) \to 1$ for 128 data with $\Phi(\mathbf{x})_j > 0$. Similarly, $\sigma(\mathbf{w}^T \Phi(\mathbf{x})) \to 0$ for data with $\Phi(x)_j > 0$ and $\sigma(\mathbf{w}^T \Phi(\mathbf{x})) \to 1$ for 129 data with $\Phi(\mathbf{x})_j < 0$. Finally, $g'(\mathbf{w}^T \Phi(\mathbf{x})) \to \delta_g^+$ for data with $\Phi(\mathbf{x})_j < 0$ and $g'(\mathbf{w}^T \Phi(\mathbf{x})) \to \delta_g^-$ 130 for data with $\Phi(\mathbf{x})_j > 0$. Thus

$$\lim_{\mathbf{w}_{j}\to-\infty} \frac{\partial}{\partial \mathbf{w}_{j}} \mathcal{L}_{g,\lambda}(\theta) = \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\y=1}} -\Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} + \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\y=0}} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} \\
+ \frac{\lambda}{N_{a}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\s=a}} \left(\delta_{g}^{+} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} + \delta_{g}^{-} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} \right) \\
- \frac{\lambda}{N_{b}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\s=b}} \left(\delta_{g}^{+} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} + \delta_{g}^{-} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} \right) \\
= -\gamma_{j.1} p_{1} + \eta_{j.0} p_{0} + \lambda \left(\delta_{g}^{+}(\eta_{ja.} - \eta_{jb.}) + \delta_{g}^{-}(\gamma_{ja.} - \gamma_{jb.}) \right) \\
> - \gamma_{j.1} p_{1} + \eta_{j.0} p_{0} + \lambda_{j}^{*} \left(\delta_{g}^{+}(\eta_{ja.} - \eta_{jb.}) + \delta_{g}^{-}(\gamma_{ja.} - \gamma_{jb.}) \right) = 0$$
(4)

131 Where $\lambda_j^* = \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{\delta_g^-(\gamma_{ja.} - \gamma_{jb.}) + \delta_g^+(\eta_{ja.} - \eta_{jb.})}$ is the lower bound on λ . The derivative converges to 132 $\kappa = -\gamma_{j.1}p_1 + \eta_{j.0}p_0 + \lambda \left(\delta_g^+(\eta_{ja.} - \eta_{jb.}) + \delta_g^-(\gamma_{ja.} - \gamma_{jb.})\right) > 0.$

For condition (2) it is sufficient to show that as $\mathbf{w}_j \to \infty$, $\frac{\partial}{\partial \mathbf{w}_j} \mathcal{L}_{g,\lambda}(\theta) \to \kappa < 0$. Equation (3) still holds, but we instead note that as $\mathbf{w}_j \to \infty$, $\sigma(-\mathbf{w}^T \Phi(\mathbf{x})) \to 0$ for data with $\Phi(\mathbf{x})_j > 0$ and $\sigma(-\mathbf{w}^T \Phi(\mathbf{x})) \to 1$ for data with $\Phi(\mathbf{x})_j < 0$. Similarly, $\sigma(\mathbf{w}^T \Phi(\mathbf{x})) \to 0$ for data with $\Phi(\mathbf{x})_j < 0$ and $\sigma(\mathbf{w}^T \Phi(\mathbf{x})) \to 1$ for data with $\Phi(\mathbf{x})_j > 0$. Finally, $g'(\mathbf{w}^T \Phi(\mathbf{x})) \to \delta_g^+$ for data with $\Phi(\mathbf{x})_j > 0$ 137 and $g'(\mathbf{w}^T \Phi(\mathbf{x})) \to \delta_g^-$ for data with $\Phi(\mathbf{x})_j < 0$. Thus

$$\lim_{\mathbf{w}_{j}\to\infty} \frac{\partial}{\partial \mathbf{w}_{j}} \mathcal{L}_{g,\lambda}(\theta) = \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\y=1}} -\Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} + \frac{1}{N} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\y=0}} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} \\
+ \frac{\lambda}{N_{a}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\s=a}} \left(\delta_{g}^{+} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} + \delta_{g}^{-} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} \right) \\
- \frac{\lambda}{N_{b}} \sum_{\substack{(\mathbf{x},s,y)\in\mathcal{D}\\s=b}} \left(\delta_{g}^{+} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}>0\}} + \delta_{g}^{-} \Phi(\mathbf{x})_{j} \mathbf{1}_{\{\Phi(\mathbf{x})_{j}<0\}} \right) \\
= -\eta_{j.1} p_{1} + \gamma_{j.0} p_{0} + \lambda \left(\delta_{g}^{+}(\gamma_{ja.} - \gamma_{jb.}) + \delta_{g}^{-}(\eta_{ja.} - \eta_{jb.}) \right) \\
< - \eta_{j.1} p_{1} + \gamma_{j.0} p_{0} + \lambda_{j}^{*} \left(\delta_{g}^{+}(\gamma_{ja.} - \gamma_{jb.}) + \delta_{g}^{-}(\eta_{ja.} - \eta_{jb.}) \right) = 0$$
(5)

138 Where $\lambda_j^* = \frac{\gamma_{j.0}p_0 - \eta_{j.1}p_1}{\delta_g^+(\gamma_{jb.} - \gamma_{ja.}) + \delta_g^-(\eta_{jb.} - \eta_{ja.})}$ is the lower bound on λ . The derivative converges to

139
$$\kappa = -\eta_{j.1}p_1 + \gamma_{j.0}p_0 + \lambda \Big(\delta_g^+(\gamma_{ja.} - \gamma_{jb.}) + \delta_g^-(\eta_{ja.} - \eta_{jb.})\Big) < 0, \text{ and thus } \mathcal{L}_{g,\lambda}(\theta) \to -\infty.$$

When g(r) is linear and $\Phi(\mathbf{x})$ is considered fixed, the entire objective is a sum of convex functions and therefore convex with respect to **w**. A convex function's only stationary point is its global minimum, which does not exist for an unbounded function. Thus there are no stationary points.

Theorem 1 is a special case of Theorem 4 where all $\eta_{jsy} = 0$, $\gamma_{jsy} = \mu_{jsy}$. The conditions in Theorem 1 make use of the fact that $\delta_g^+, \delta_g^- > 0$ from the monotonicity assumption.

146 (1) If
$$\delta_g^-(\gamma_{ja0} - \gamma_{jb0}) + \delta_g^+(\eta_{ja0} - \eta_{jb0}) > 0$$
 and $\lambda > \lambda_j^* = \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{\delta_g^-(\gamma_{ja0} - \gamma_{jb0}) + \delta_g^+(\eta_{ja0} - \eta_{jb0})}$, then

$$\begin{array}{l} & \mathcal{L}_{g,\lambda}(b) = -\infty. \\ & \text{if } \delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0}) > 0 \text{ and } \lambda > \lambda_{j}^{*} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0})}, \text{ then} \\ & \text{if } \lambda_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0}) > 0 \text{ and } \lambda > \lambda_{j}^{*} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0})}, \text{ then} \\ & \text{if } \lambda_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0}) > 0 \text{ and } \lambda > \lambda_{j}^{*} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0})}, \end{array}$$

149
$$\lim_{w_i \to +\infty} \mathcal{L}_{g,\lambda}(\theta) = -\infty.$$

150 If either (1) or (2) hold and g(r) is linear, then for this fixed $\Phi(\mathbf{x})$, $\mathcal{L}_{g,\lambda}$ has no stationary points.

Proof of Theorems 2 and 5. The proof is identical to the proof for Theorems 1 and 4, except with the fairness metrics defined only on the negative instances. \Box

- ¹⁵³ Finally, we present a generalization of Theorem 3:
- **Theorem 6.** For any surrogate g(r) such that $\delta_g^- = 0$, both thresholds in Theorem 5 are bounded below by p_{b0}/δ_q^+ :

156 (1) If
$$\delta_{g}^{-}(\gamma_{ja0} - \gamma_{jb0}) + \delta_{g}^{+}(\eta_{ja0} - \eta_{jb0}) > 0$$
, then $\lambda_{j}^{*} = \frac{\gamma_{j.1}p_{1} - \eta_{j.0}p_{0}}{\delta_{g}^{-}(\gamma_{ja0} - \gamma_{jb0}) + \delta_{g}^{+}(\eta_{ja0} - \eta_{jb0})} \ge p_{b0}/\delta_{g}^{+}$.
157 (2) If $\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0}) > 0$ then $\lambda_{j}^{*} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0})} \ge p_{b0}/\delta_{g}^{+}$.

Proof of Theorems 3 and 6. We start by proving condition (1). Applying the assumption that $\delta_g^- = 0$, we see that $\delta_g^+(\eta_{ja0} - \eta_{jb0}) > 0$. Because all $\eta_{jsy} \le 0$ and $\delta_g^+ \ge 0$, we know $-\eta_{jb0}\delta_g^+ \ge$ 160 $\delta_g^+(\eta_{ja0} - \eta_{jb0}) > 0$. Thus

$$\lambda_j^* = \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{\delta_g^-(\gamma_{ja0} - \gamma_{jb0}) + \delta_g^+(\eta_{ja0} - \eta_{jb0})} = \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{\delta_g^+(\eta_{ja0} - \eta_{jb0})} \ge \frac{\gamma_{j.1}p_1 - \eta_{j.0}p_0}{-\eta_{jb0}\delta_g^+} \ge \frac{-\eta_{j.0}p_0}{-\eta_{jb0}\delta_g^+}$$

161 The last inequality uses $\gamma_{j,1}p_1 \ge 0$. Further note that,

$$\eta_{j,0} = \frac{1}{N_0} \sum_{\substack{(\mathbf{x},s,y) \in \mathcal{D} \\ y=0}} -\min(0, \Phi(\mathbf{x})_j) \ge \frac{1}{N_0} \sum_{\substack{(\mathbf{x},s,y) \in \mathcal{D} \\ s=b,y=0}} -\min(0, \Phi(\mathbf{x})_j) = -\frac{N_{b0}}{N_0} \eta_{jb0}$$

162 This follow from the fact that every term in these sums is non-negative. Thus

$$\lambda_{j}^{*} \geq \frac{-\eta_{j.0}p_{0}}{-\eta_{jb0}\delta_{g}^{+}} \geq \frac{-\eta_{jb0}\frac{N_{b0}}{N_{0}}\frac{N_{0}}{N}}{-\eta_{jb0}\delta_{g}^{+}} = \frac{p_{b0}}{\delta_{g}^{+}}$$

- Next we prove condition (2). Applying the assumption that $\delta_g^- = 0$, we see that $\delta_g^+(\gamma_{jb0} \gamma_{ja0}) > 0$. Because all $\gamma_{jsy} \ge 0$ and $\delta_g^+ \ge 0$, we know $\gamma_{jb0}\delta_g^+ \ge \delta_g^+(\gamma_{jb0} \gamma_{ja0}) > 0$. Thus 163
- 164

$$\lambda_{j}^{*} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0}) + \delta_{g}^{-}(\eta_{jb0} - \eta_{ja0})} = \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}(\gamma_{jb0} - \gamma_{ja0})} \ge \frac{\gamma_{j.0}p_{0} - \eta_{j.1}p_{1}}{\delta_{g}^{+}\gamma_{jb0}} \ge \frac{\gamma_{j.0}p_{0}}{\delta_{g}^{+}\gamma_{jb0}}$$

165 The last inequality uses $-\eta_{j,1}p_1 \ge 0$. Further note that

$$\gamma_{j,0} = \frac{1}{N_0} \sum_{\substack{(\mathbf{x},s,y) \in \mathcal{D} \\ y=0}} \max(0, \Phi(\mathbf{x})_j) \ge \frac{1}{N_0} \sum_{\substack{(\mathbf{x},s,y) \in \mathcal{D} \\ s=b,y=0}} \max(0, \Phi(\mathbf{x})_j) = \frac{N_{b0}}{N_0} \gamma_{jb0}$$

This follows from the fact that every term in these sums is non-negative. Thus 166

$$\lambda_j^* \geq \frac{\gamma_{j.0} p_0}{\delta_g^+ \gamma_{jb0}} \geq \frac{\gamma_{jb0} \frac{N_{b0}}{N_0} \frac{N_0}{N}}{\delta_g^+ \gamma_{jb0}} = \frac{p_{b0}}{\delta_g^+}$$

167

Theorem 3 is a special case of Theorem 6, where $\Phi(\mathbf{x})_j \ge 0$ so all $\eta_{jsy} = 0$ and all $\gamma_{jsy} = \mu_{jsy}$. Thus condition (1) of Theorem 6 cannot hold, and given condition (2), $\lambda_j^* \ge p_{b0}/\delta_g^+$. Note that for 168 all relaxations we consider where $\delta_g^+ \neq 0$, $\delta_g^+ = 1$ and thus $\lambda_j^* \geq p_{b0}$. 169

References 170

- [1] Agarwal, A., Beygelzimer, A., Dudík, M., Langford, J., and Wallach, H. A reductions approach 171 to fair classification. arXiv preprint arXiv:1803.02453, 2018. 172
- [2] Angwin, J., Larson, J., Mattu, S., and Kirchner, L. Machine bias. propublica, may 23, 2016, 173 2016. 174
- [3] API, G., 2016. URL https://gender-api.com/. 175
- [4] Bechavod, Y. and Ligett, K. Penalizing unfairness in binary classification. arXiv preprint 176 177 arXiv:1707.00044, 2017.
- [5] Calders, T. and Verwer, S. Three naive bayes approaches for discrimination-free classification. 178 Data Mining and Knowledge Discovery, 21(2):277–292, 2010. 179
- [6] Cotter, A., Jiang, H., and Sridharan, K. Two-player games for efficient non-convex constrained 180 181 optimization. In Algorithmic Learning Theory, pp. 300–332. PMLR, 2019.
- [7] Denton, E., Hutchinson, B., Mitchell, M., and Gebru, T. Detecting bias with generative 182 counterfactual face attribute augmentation. arXiv preprint arXiv:1906.06439, 2019. 183
- [8] Donini, M., Oneto, L., Ben-David, S., Shawe-Taylor, J. S., and Pontil, M. Empirical risk 184 minimization under fairness constraints. In Advances in Neural Information Processing Systems, 185 pp. 2791-2801, 2018. 186
- [9] Escalera, S., Baró, X., Escalante, H. J., and Guyon, I. Chalearn looking at people: A review of 187 events and resources. In 2017 International Joint Conference on Neural Networks (IJCNN), pp. 188 1594-1601. IEEE, 2017. 189
- [10] Friedler, S. A., Scheidegger, C., Venkatasubramanian, S., Choudhary, S., Hamilton, E. P., and 190 Roth, D. A comparative study of fairness-enhancing interventions in machine learning. In 191 Proceedings of the Conference on Fairness, Accountability, and Transparency, pp. 329–338, 192 2019. 193
- [11] Goh, G., Cotter, A., Gupta, M., and Friedlander, M. P. Satisfying real-world goals with dataset 194 constraints. In Advances in Neural Information Processing Systems, pp. 2415–2423, 2016. 195
- [12] Kohavi, R. Scaling up the accuracy of naive-bayes classifiers: A decision-tree hybrid. In Kdd, 196 volume 96, pp. 202-207, 1996. 197
- [13] Lohaus, M., Perrot, M., and von Luxburg, U. Too relaxed to be fair. In International Conference 198 on Machine Learning, 2020. 199
- [14] Mansoury, M., Mobasher, B., Burke, R., and Pechenizkiy, M. Bias disparity in collaborative rec-200 ommendation: algorithmic evaluation and comparison. In 2019 Workshop on Recommendation 201 in Multi-Stakeholder Environments, RMSE 2019, pp. 6. CEUR-WS. org, 2019. 202

- [15] Wu, Y., Zhang, L., and Wu, X. On convexity and bounds of fairness-aware classification. In
 The World Wide Web Conference, pp. 3356–3362, 2019.
- [16] Zafar, M. B., Valera, I., Gomez Rodriguez, M., and Gummadi, K. P. Fairness beyond disparate treatment & disparate impact: Learning classification without disparate mistreatment. In *Proceedings of the 26th international conference on world wide web*, pp. 1171–1180, 2017.
- [17] Zafar, M. B., Valera, I., Gomez-Rodriguez, M., and Gummadi, K. P. Fairness constraints: A
 flexible approach for fair classification. *Journal of Machine Learning Research*, 20(75):1–42,
 2019. URL http://jmlr.org/papers/v20/18-262.html.
- 2019. OKL http://jmii.org/papers/v20/10-202.html
- [18] Zagoruyko, S. and Komodakis, N. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.