

ATTAINABILITY AND OPTIMALITY: THE EQUALIZED-ODDS FAIRNESS REVISITED

Anonymous authors

Paper under double-blind review

ABSTRACT

Fairness of machine learning algorithms has been of increasing interest. In order to suppress or eliminate discrimination in prediction, various notions as well as approaches to impose fairness have been proposed. However, in different scenarios, whether or not the chosen notion of fairness can always be attained, even if with unlimited amount of data, is not well addressed. In this paper, focusing on the Equalized Odds notion of fairness, we consider the attainability of this criterion, and furthermore, if attainable, the optimality of the prediction performance under various settings. In particular, for classification with a deterministic prediction function of the input, we give the condition under which Equalized Odds can hold true; if randomized prediction is acceptable, we show that under mild assumptions, fair classifiers can always be derived. Moreover, we prove that compared to enforcing fairness by post-processing, one can always benefit from exploiting all available features during training and get better prediction performance while remaining fair. However, for regression tasks, Equalized Odds is not always attainable if certain conditions on the joint distribution of the features and the target variable are not met. This indicates the inherent difficulty in achieving fairness in certain cases and suggests a broader class of prediction methods might be needed for fairness.

1 INTRODUCTION

As machine learning models become widespread in automated decision making systems, apart from the efficiency and accuracy of the prediction, their potential social consequence also gains increasing attention. To date, there is ample evidence that machine learning models have resulted in discrimination against certain groups of individuals under many circumstances, for instance, the discrimination in ad delivery when searching for names that can be predictive of the race of individual (Sweeney, 2013); the gender discrimination in job-related ads push (Datta et al., 2015); stereotypes associated with gender in word embeddings (Bolukbasi et al., 2016); the bias against certain ethnicities in the assessment of recidivism risk (Angwin et al., 2016).

The call for accountability and fairness in machine learning has motivated various (statistical) notions of fairness. The *Demographic Parity* criterion (Calders et al., 2009) requires the independence between prediction (e.g., of a classifier) and the protected feature (sensitive attributes of an individual, e.g., gender, race). *Equalized Odds* (Hardt et al., 2016), also known as *Error-rate Balance* (Chouldechova, 2017), requires the output of a model be conditionally independent of protected feature(s) given the ground truth of the target. *Predictive Rate Parity* (Zafar et al., 2017a), on the other hand, requires the actually proportion of positives (negatives) in the original data for positive (negative) predictions should match across groups (well-calibrated).

On the theoretical side, results have been reported regarding relationships among fairness notions. It has been independently shown that if base rates of true positives differ among groups, then Equalized Odds and Predictive Rate Parity cannot be achieved simultaneously for non-perfect predictors (Kleinberg et al., 2016; Chouldechova, 2017). Any two out of three among Demographic Parity, Equalized Odds, and Predictive Rate Parity are incompatible with each other (Barocas et al., 2017). At the interface of privacy and fairness, the impossibility of achieving both *Differential Privacy* (Dwork et al., 2006) and *Equal Opportunity* (Hardt et al., 2016) while maintaining non-trivial accuracy is also established (Cummings et al., 2019).

In practice, one can broadly categorize computational procedures to derive a fair predictor into three types: pre-processing approaches (Calders et al., 2009; Dwork et al., 2012; Zemel et al., 2013; Zhang et al., 2018; Madras et al., 2018; Creager et al., 2019; Zhao et al., 2020), in-processing approaches (Kamishima et al., 2011; Pérez-Suay et al., 2017; Zafar et al., 2017a;b; Donini et al., 2018; Song et al., 2019; Mary et al., 2019; Baharlouei et al., 2020), and post-processing approaches (Hardt et al., 2016; Fish et al., 2016; Dwork et al., 2018). In accord with the fairness notion of interest, a pre-processing approach first maps the training data to a transformed space to remove discriminatory information between protected feature and target, and then pass on the data to make prediction. In direct contrast, a post-processing approach treats the off-the-shelf predictor(s) as uninterpretable black-box(es), and imposes fairness by outputting a function of the original prediction. For in-processing approaches, various kinds of regularization terms are proposed so that one can optimize the utility function while suppressing the discrimination at the same time. Approaches based on estimating/bounding causal effect between the protected feature and final target have also been proposed (Kusner et al., 2017; Russell et al., 2017; Zhang et al., 2017; Nabi & Shpitser, 2018; Zhang & Bareinboim, 2018; Chiappa, 2019; Wu et al., 2019).

Focusing on the Equalized-Odds criterion, although various approaches have been proposed to impose the fairness requirement, whether or not it is always attainable is not well addressed. The attainability of Equalized Odds, namely, the existence of the predictor that can score zero violation of fairness in the large sample limit, is an asymptotic property of the fairness criterion. This characterizes a completely different kind of violation of fairness compared to the empirical error bound of discrimination in finite-sample cases. If utilizing a “fair” predictor which is actually biased, the discrimination would become a snake in the grass, making it hard to detect and eliminate. Actually, as we illustrate in this paper, Equalized Odds is not always attainable for regression and even classification tasks, if we use deterministic prediction functions. This calls for alternative definitions in the same spirit as Equalized Odds that can always be achieved under various circumstances. Our contributions are mainly:

- For regression and classification tasks with deterministic prediction functions, we show that Equalized Odds is not always attainable if certain (rather restrictive) conditions on the joint distribution of the features and the target variable are not met.
- Under mild assumptions, for binary classification we show that if randomized prediction is taken into consideration, one can always derive a non-trivial Equalized Odds classifier.
- Considering the optimality of performance under fairness constraint(s), when exploiting all available features, we show that the predictor derived via an in-processing approach would always outperform the one derived via a post-processing approach (unconstrained optimization followed by a post-processing step).

2 PRELIMINARIES

In this section, we first illustrate the difference between *prediction fairness* and *procedure fairness*, and then, we present the formal definition of Equalized Odds (Hardt et al., 2016).

2.1 HIERARCHY OF FAIRNESS

Before presenting the formulation of fairness, it is important to see the distinction between different levels of fairness when discussing fair predictors. When evaluating the performance of the proposed fair predictor, it is a common practice to compare the loss (with respect to the utility function of choice, e.g., accuracy for binary classification) computed on target variable and the predicted value. There is an implicit assumption lying beneath this practice: the generating process of the data, which is just describing a real-world procedure, is *not* biased in any sense (Danks & London, 2017). Only when we treat the target variable (recorded in the dataset) as unbiased can we justify the practice of loss evaluation and the conditioning on target variable when imposing fairness (as we shall see in the definition of Equalized Odds in Equation 1).

One may consider a music school admission example. The music school committee would decide if they admit a student to the violin performance program based on the applicant’s personal information, educational background, instrumental performance, and so on. When evaluating whether or not

the admission is ‘‘fair’’, there are actually two levels of fairness. First, based on the information at hand, did the committee evaluate the qualification of applicants without bias (How committee evaluate the applicants)? And second, is committee’s procedure of evaluating applicants’ qualification reasonable (How other people view the evaluation procedure used by the committee)?

In this paper, we consider *prediction fairness*, namely, assuming the data recorded is unbiased, the prediction (made with respect to current reality) itself should not include any biased utilization of information. The fairness with respect to the data generating procedure as well as the potential future influence of the prediction are beyond the scope of this paper.

2.2 EQUALIZED-ODDS FAIRNESS

Hardt et al. (2016) proposed *Equalized Odds* which requires conditional independence between prediction and protected feature(s) given ground truth of the target. Let us denote the protected feature by A , with domain of value \mathcal{A} , additional (observable) feature(s) by X , with domain of value \mathcal{X} , target variable by Y , with domain \mathcal{Y} , (not necessarily fair) predictors by \tilde{Y} , and fair predictors by \hat{Y} . Equalized-Odds fairness requires

$$\tilde{Y} \perp\!\!\!\perp A \mid Y. \quad (1)$$

For classification tasks, one can conveniently use the probability distribution form:

$$\forall a \in \mathcal{A}, t, y \in \mathcal{Y} : P(\tilde{Y} = t \mid A = a, Y = y) = P(\tilde{Y} = t \mid Y = y), \quad (2)$$

$$\text{or more concisely, } P_{\tilde{Y}|AY}(t \mid a, y) = P_{\tilde{Y}|Y}(t \mid y). \quad (3)$$

For better readability, we also use the formulation in Equation 3 in cases without ambiguity. In the context of binary classification ($\mathcal{Y} = \{0, 1\}$), Equalized Odds requires that the True Positive Rate (TPR) and False Positive Rate (FPR) of each certain group match population positive rates. Throughout the paper, without loss of generality we assume there is only one protected feature for the purpose of simplifying notation. However, considering the fact that the protected feature can be discrete (e.g., race, gender) or continuous (e.g., the ratio of ethnic group in the population for certain district of a city), we do not assume discreteness of the protected feature. Due to the space limit, we will focus on the illustration and implication of our results and defer all the proofs to the appendix.

3 FAIRNESS IN REGRESSION MAY NOT BE ATTAINED

In this section we consider the attainability of Equalized Odds for regression tasks, namely, whether or not it is possible to find a predictor that is conditionally independent from the protected feature given true value of the target. For linearly Gaussian cases, one can attain Equalized Odds by constraining zero partial correlation between the prediction and the protected feature given target variable (Woodworth et al., 2017). Various regularization terms have also been proposed to suppress discrimination when predicting a continuous target (Berk et al., 2017; Mary et al., 2019). However, whether or not one can always achieve 0-discrimination for regression, even if with an unlimited amount of data, is not clear yet.

If ‘‘fair’’ predictors are deployed without carefully checking the attainability of fairness, the discrimination would become a hidden hazard, making it hard to detect and eliminate. Actually as we will show in this section, even in the simple setup of linearly correlated continuous data, Equalized Odds is not always attainable.

3.1 UNATTAINABILITY OF EQUALIZED ODDS IN LINEAR NON-GAUSSIAN REGRESSION

As stated in Section 2.1, in this paper we consider *prediction fairness*, and therefore any possible bias introduced by the data generating procedure itself is beyond the scope of the discussion. Consider the situation where the data is generated as following (H is not measured in the dataset):

$$\begin{aligned} X &= qA + E_X, \\ H &= bA + E_H, \\ Y &= cX + dH + E_Y, \end{aligned} \quad (4)$$

where (A, E_X, E_H, E_Y) are mutually independent. In fact, if at most one of E_X and $E := E_Y + dE_H$ is Gaussian, then any linear combination of A and X with non-zero coefficients will not be conditionally independent from A given Y , meaning that it is not possible to achieve Equalized-Odds fairness. Let Z be a linear combination of A and X , i.e., $Z = \alpha A + \beta X = (\alpha + q\beta)A + \beta E_X$, with linear coefficients α and β , where $\beta \neq 0$. In Theorem 3.1, we present the general result in linear non-Gaussian cases, where one cannot achieve the conditional independence between Z and A given Y .

Theorem 3.1. (Unattainability of Equalized Odds in the Linear Non-Gaussian Case)

Assume that feature X has a causal influence on Y , i.e., $c \neq 0$ in Equation 4, and that the protected feature A and Y are not independent, i.e., $qc + bd \neq 0$. Assume p_{E_X} and p_E are positive on \mathbb{R} . Let $f_1 := \log p_A$, $f_2 := \log p_{E_X}$, and $f_3 := \log p_E$. Further assume that f_2 and f_3 are third-order differentiable. Then if at most one of E_X and E is Gaussian, Z is always conditionally dependent on A given Y .

From Theorem 3.1, we see that in linear non-Gaussian cases, any non-zero linear combination of the feature (which is a deterministic function of the input) will not satisfy Equalized Odds. One may wonder whether Equalized Odds can be achieved by nonlinear regression, instead of a linear model. Although a proof with general nonlinear models is rather complicated, our simulation results in Section 5.1 strongly suggest that the unattainability of Equalized Odds persists in nonlinear regression cases.

In light of the unattainability of Equalized Odds for prediction with deterministic functions of A and X , it is desirable to develop general, nonlinear prediction algorithms to produce a probabilistic prediction (i.e., with a certain type of randomness in the prediction). One possible way follows the framework of Generative Adversarial Networks (GANs) (Goodfellow et al., 2014): we use random standard Gaussian noise, in addition to A and X , as input, such that the output will have a specific type of randomness. The parameters involved are learned by minimizing prediction error and enforcing Equalized Odds on the ‘‘randomized’’ output at the same time. Given that this approach is not essential to illustrate the claims made in this paper and that theoretical properties of such nonlinear regression algorithms with randomized output are not straightforward to establish, this is left as future work.

4 FAIRNESS IN CLASSIFICATION

In this section, we consider the attainability of Equalized Odds for binary classifiers (with a deterministic or randomized prediction function), and furthermore, if attainable, the optimality of performance under the fairness criterion. Admittedly, as is already pointed out by Woodworth et al. (2017), we generally cannot have 0-discriminatory predictors with a finite number of samples; instead, one should consider imposing δ -discrimination in practice (δ is the violation of Equalized Odds). However, this does not guarantee nor rule out the possibility of attaining 0-discrimination on population when the sample size goes to infinity.

4.1 CLASSIFICATION WITH DETERMINISTIC PREDICTION

We begin with considering cases when the classification is performed by a deterministic function of the input. In particular, we derive the condition under which Equalized Odds can possibly hold true.

Theorem 4.1. *Assume that the protected feature A and Y are dependent and that their joint probability $P(A, Y)$ (for discrete A) or joint probability density $p(A, Y)$ (for continuous A) is positive for every combination of possible values of A and Y . Further assume that Y is not fully determined by A , and that there are additional features X that are not independent of Y . Let the output of the classifier \tilde{Y} be a deterministic function $f : \mathcal{A} \times \mathcal{X} \rightarrow \mathcal{Y}$. Let $S_A^{(t)} := \{a \mid \exists x \in \mathcal{X} \text{ s.t. } f(a, x) = t\}$, and $S_{X|a}^{(t)} := \{x \mid f(a, x) = t\}$. Equalized Odds holds true if and only if the following two conditions are satisfied:*

- (i) $\forall t \in \mathcal{Y} : S_A^{(t)} = \mathcal{A}$,
- (ii) $\forall t \in \mathcal{Y}, \forall a, a' \in \mathcal{A}$ (for continuous X , replace summation with integration accordingly):

$$\sum_{x \in S_{X|a}^{(t)}} P_{X|AY}(x \mid a, y) = \sum_{x \in S_{X|a'}^{(t)}} P_{X|AY}(x \mid a', y).$$

Let us take a look at the two conditions. Condition (i) says that within each class determined by the classification function f , A should be able to take all possible values in \mathcal{A} . While condition (i) is already pretty restrictive, condition (ii) specifies an even stronger constraint on the relation between the conditional probability $P_{X|AY}(x|a, y)$ (or the conditional probability density $p_{X|AY}(x|a, y)$ for continuous X) and the set $S_{X|a}^{(t)}$ (which is determined by the function f). By definition $S_{X|a}^{(t)}$ has following properties: (1) for any fixed value of $a \in \mathcal{A}$, if $t \neq t'$, then $S_{X|a}^{(t)} \cap S_{X|a}^{(t')} = \emptyset$; (2) for any fixed value of $a \in \mathcal{A}$, $\bigcup_{t \in \mathcal{Y}} S_{X|a}^{(t)} = \mathcal{X}$. Condition (ii) says that the set $S_{X|a}^{(t)}$ and the conditional distribution $P_{X|AY}(x|a, y)$ are coupled in some specific way so that they happen to satisfy the specified equality.

In special cases when $X \perp\!\!\!\perp A \mid Y$, if f is a function of only X , condition (ii) would always hold true. In general situations, if there does not exist any subsets $K_a, K_{a'} \subseteq \mathcal{X}$ for different values of $a, a' \in \mathcal{A}$ such that $\sum_{x \in K_a} P_{X|AY}(x|a, y) = \sum_{x \in K_{a'}} P_{X|AY}(x|a', y)$, then condition (ii) can never hold true (i.e., we cannot find a deterministic function $f(A, X)$ that satisfies Equalized Odds). Generally speaking, in order to score a better classification accuracy, one would like to make $P_{\hat{Y}|A, X}(t|a, x)$ as close as possible to $P_{Y|A, X}(y|a, x)$, and if the set $S_{X|a}^{(t)}$ and $P_{X|AY}(x|a, y)$ are not strictly coupled, condition (ii) would be violated.

4.2 CLASSIFICATION WITH RANDOMIZED PREDICTION

In this section, we consider cases when randomized prediction is acceptable, namely, the classifier would output class labels with certain probabilities. We first derive the relation between positive rates (TPR and FPR) of binary classifiers before and after the post-processing step, i.e., \hat{Y}_{opt} (the unconstrainedly optimized classifier) and \tilde{Y}_{post} (the fair classifier derived by post-processing \hat{Y}_{opt}), and show that under mild assumptions, one can always derive a non-trivial Equalized-Odds (on population level) \tilde{Y}_{post} via a post-processing step. Then, from the ROC feasible area perspective, we prove that post-processing approaches are actually equivalent to in-processing approaches but with additional “pseudo” constraints enforced. Therefore, using the same loss function, post-processing approaches can perform no better than in-processing approaches.

4.2.1 THE POST-PROCESSING STEP

The post-processing step of a predictor \hat{Y} (here we drop the subscript if without ambiguity) only utilizes the information in the joint distribution (A, Y, \hat{Y}) . A fair predictor \tilde{Y}_{post} derived via a post-processing step, for instance, the *shifted decision boundary* (Fish et al., 2016), the *derived predictor* (Hardt et al., 2016), or the *(monotonic) joint loss optimization over decoupled classifiers* (Dwork et al., 2018), is then fully specified by a (possibly randomized) function of (A, \hat{Y}) . This implies the conditional independence $\tilde{Y}_{\text{post}} \perp\!\!\!\perp Y \mid A, \hat{Y}$. Since we can denote the positive rates of \hat{Y} as $P_{\hat{Y}|AY}(1|a, y)^1$, positive rates of \tilde{Y} (here we drop the subscript for readability) as $P_{\tilde{Y}|AY}(1|a, y)$, the relation between positive rates of binary classifiers before and after a post-processing step would satisfy (for every $a \in \mathcal{A}, u, y \in \mathcal{Y}$):

$$P_{\tilde{Y}|AY}(1|a, y) = \sum_{u \in \mathcal{Y}} \beta_a^{(u)} P_{\hat{Y}|AY}(u|a, y), \text{ where } \beta_a^{(u)} := P(\tilde{Y} = 1 \mid A = a, \hat{Y} = u). \quad (5)$$

Notice that Equation 5 is just a factorization of probability under the conditional independence (between \tilde{Y}_{post} and Y given A and \hat{Y}). Therefore, post-processing an existing predictor boils down to optimizing parameters (for discrete A) or functions (for continuous A) $\beta_a^{(u)}$.

4.2.2 ROC FEASIBLE AREA

On the Receiver Operator Characteristic (ROC) plane, a two-dimensional plane with horizontal axis denoting FPR and vertical axis denoting TPR, the performance of any binary predictor \hat{Y} (not

¹Recall that $P_{\hat{Y}|AY}(u|a, y) = P(\hat{Y} = u \mid A = a, Y = y)$. When $u = 1$, $P_{\hat{Y}|AY}(1|a, y) = P(\hat{Y} = 1 \mid A = a, Y = y)$ represents positive rates of \hat{Y} ; When $u = 0$, $P_{\hat{Y}|AY}(0|a, y) = P(\hat{Y} = 0 \mid A = a, Y = y)$ represents positive rates of $1 - \hat{Y}$ (the classifier that flips the prediction of \hat{Y}).

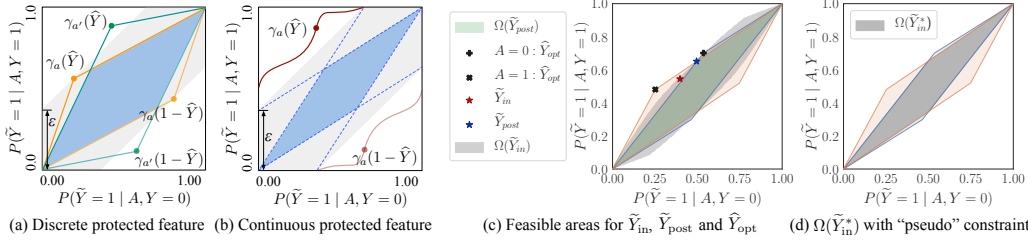


Figure 1: ROC feasible area illustrations. Panels (a)-(b): Attainability of Equalized Odds for binary classifiers with discrete or continuous protected feature. Panels (c)-(d): ROC feasible areas comparison between $\Omega(\tilde{Y}_{in})$, $\Omega(\tilde{Y}_{post})$, $\Omega(\tilde{Y}_{opt})$, and $\Omega(\tilde{Y}_{in}^*)$.

necessarily a fair one) with certain value of protected feature $A = a$ corresponds to a point $\gamma_a(\hat{Y}) = (\text{FPR}, \text{TPR})$ on the plane. Denote each coordinate according to the value of Y as $\gamma_{ay}(\hat{Y})$:

$$\gamma_a(\hat{Y}) = (\gamma_{a0}(\hat{Y}), \gamma_{a1}(\hat{Y})) := (P_{\hat{Y}|AY}(1|a, 0), P_{\hat{Y}|AY}(1|a, 1)). \quad (6)$$

Further denote the corresponding convex hull of \hat{Y} on the ROC plane as $\mathcal{C}_a(\hat{Y})$ using vertices:

$$\mathcal{C}_a(\hat{Y}) := \text{convhull}\{(0, 0), \gamma_a(\hat{Y}), \gamma_a(1 - \hat{Y}), (1, 1)\}, \quad (7)$$

and then, as already stated in Hardt et al. (2016), the (FPR, TPR) pair corresponding to a post-processing predictor falls within (including the boundary of) $\mathcal{C}_a(\hat{Y})$.

Definition 4.1. (ROC feasible area) The feasible area of a predictor $\Omega(\hat{Y})$, specified by the hypothesis space of available predictors \hat{Y} , is the set containing all attainable (FPR, TPR) pairs by the predictor on the ROC plane satisfying Equalized Odds.

In Hardt et al. (2016) it is proposed that the post-processing fair predictor can be derived by solving a linear programming problem on the ROC plane. However, it is not clearly stated whether or not such problem always has a non-trivial solution. Following Hardt et al. (2016), we analyze the relation between the (FPR, TPR) pair of predictors on the ROC plane and formally establish the existence of the non-trivial Equalized-Odds predictor. Actually as we shall see in Theorem 4.2, under mild assumptions, an Equalized-Odds predictor \tilde{Y}_{post} derived via post-processing \hat{Y} (a predictor optimized without fairness concern) always has non-empty ROC feasible area.

Theorem 4.2. (Attainability of Equalized Odds)

Assume that the feature X is not independent from Y , and that \hat{Y} is a function of A and X . Then for binary classification, if \hat{Y} is a non-trivial predictor for Y , there is always at least one non-trivial (possibly randomized) predictor \tilde{Y}_{post} derived by post-processing \hat{Y} that can attain Equalized Odds:

$$\Omega(\tilde{Y}_{post}) \neq \emptyset.$$

Here \tilde{Y}_{post} is a possibly randomized function of only A and \hat{Y} , trading off TPR with FPR across groups with different value of protected feature. From the panels (a) and (b) of Figure 1 we can also see that $\Omega(\tilde{Y}_{post})$, the ROC feasible area of \tilde{Y}_{post} , is the intersection of $\mathcal{C}_a(\hat{Y})$, indicating that although Equalized Odds is attained, the performance of \tilde{Y}_{post} is always worse than the weakest performance across different groups, which is obviously suboptimal.

4.2.3 OPTIMALITY OF PERFORMANCE AMONG FAIR CLASSIFIERS

In this subsection we discuss the optimality of performance of fair classifiers derived via different approaches. Considering the fact that recent efforts to impose Equalized Odds in the pre-processing manner (Madras et al., 2018; Zhao et al., 2020) approach the problem from a representation learning perspective, where the main focus is to learn fair representations that at the same time preserve sufficient information from the original data, we omit pre-processing approaches from the discussion and compare the performance of post-processing and in-processing fair classifiers.

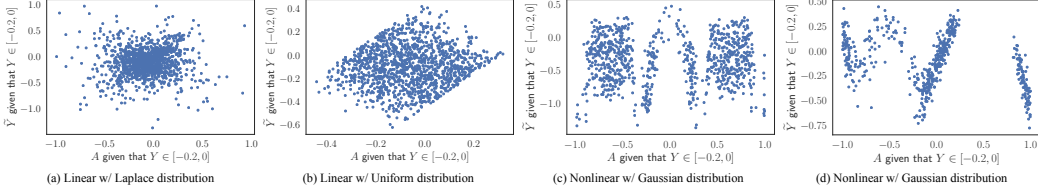


Figure 2: Illustration of unattainable Equalized Odds for regression tasks. Panel (a)-(b): Linear regression on the data generated with linear transformations and non-Gaussian distributed exogenous terms (following Laplace, Uniform distribution respectively). Panel (c)-(d): Nonlinear regression with a neural net regressor (Mary et al., 2019) on the data generated with nonlinear transformations and Gaussian exogenous terms. We can observe obvious dependencies between \tilde{Y} and A on a small interval of Y . This indicates the conditional dependency between \tilde{Y} and A given Y , i.e., the Equalized Odds is not achieved.

Admittedly, when only the information about joint distribution of (A, Y, \hat{Y}) is available, post-processing is the best we can do. However, this is not the case when we have access to additional available features during training. For any predictor specified by parameters $\theta \in \Theta$, the derivation of the in-processing fair predictor \tilde{Y}_{in} and the unconstrained statistical optimal predictor \hat{Y}_{opt} take following forms respectively:

$$\begin{aligned}
 & \min_{\theta \in \Theta} \mathbb{E}[l(f(A, X; \theta), Y)] & \min_{\theta \in \Theta} \mathbb{E}[l(f(A, X; \theta), Y)] \\
 & \text{s.t. } P_{\tilde{Y}_{\text{in}}|AY}(t | a, y) = P_{\tilde{Y}_{\text{in}}|Y}(t | y) & \text{where } \hat{Y}_{\text{opt}} = f(A, X; \theta). \\
 & \text{where } \tilde{Y}_{\text{in}} = f(A, X; \theta), &
 \end{aligned} \tag{9}$$

It is natural to wonder, now that one can always directly solve for \tilde{Y}_{in} from Equation 8, how is it related to \tilde{Y}_{post} , which is derived by post-processing the \hat{Y}_{opt} solved from Equation 9? Interestingly, although \tilde{Y}_{in} and \tilde{Y}_{post} are solved separately using different constrained optimization schemes, one can draw a connection between them by utilizing \hat{Y}_{opt} as a bridge and reason about the relation between their ROC feasible areas $\Omega(\tilde{Y}_{\text{in}})$ and $\Omega(\tilde{Y}_{\text{post}})$, as we summarize in the following theorem.

Theorem 4.3. (Equivalence between ROC feasible areas)

Let $\Omega(\tilde{Y}_{\text{post}})$ denote the ROC feasible area specified by the constraints enforced on \tilde{Y}_{post} . Then $\Omega(\tilde{Y}_{\text{post}})$ is identical to the ROC feasible area $\Omega(\tilde{Y}_{\text{in}}^*)$ that is specified by the following set of constraints:

- (i) constraints enforced on \tilde{Y}_{in} ;
- (ii) additional “pseudo” constraints: $\forall a \in \mathcal{A}, \beta_{a0}^{(0)} = \beta_{a1}^{(0)}, \beta_{a0}^{(1)} = \beta_{a1}^{(1)}$, where $\beta_{ay}^{(u)} = \sum_{x \in \mathcal{X}} P(\tilde{Y}_{\text{in}} = 1 | A = a, X = x) P(X = x | A = a, Y = y, \hat{Y}_{\text{opt}} = u)$.

As we can see from panels (c) and (d) of Figure 1, if the additional “pseudo” constraints are introduced when optimizing \tilde{Y}_{in}^* , we have $\Omega(\tilde{Y}_{\text{in}}) \supseteq \Omega(\tilde{Y}_{\text{post}}) = \Omega(\tilde{Y}_{\text{in}}^*)$. The ROC feasible area is fully specified by the hypothesis class and the fairness constraint. Therefore, with the same objective function and fairness constraint, the fair classifier derived from an in-processing approach always outperforms the one derived from a post-processing approach. We can see that when we have access to additional features and choose a post-processing approach, we lose performance (compared to \tilde{Y}_{in}) by unintentionally introducing “pseudo” constraints during optimization. These “pseudo” constraints actually offset the benefit of utilizing additional features (in the hope to score a better performance while remaining fair).

5 EXPERIMENTS

In order to intuitively illustrate the claims, we provide numerical results for various settings. We first present the result for (linear non-Gaussian and nonlinear) regression tasks when Equalized Odds is not attained. We demonstrate the dependence between the prediction and the protected feature given true value of the target variable. Then for classification tasks we compare the performance of several

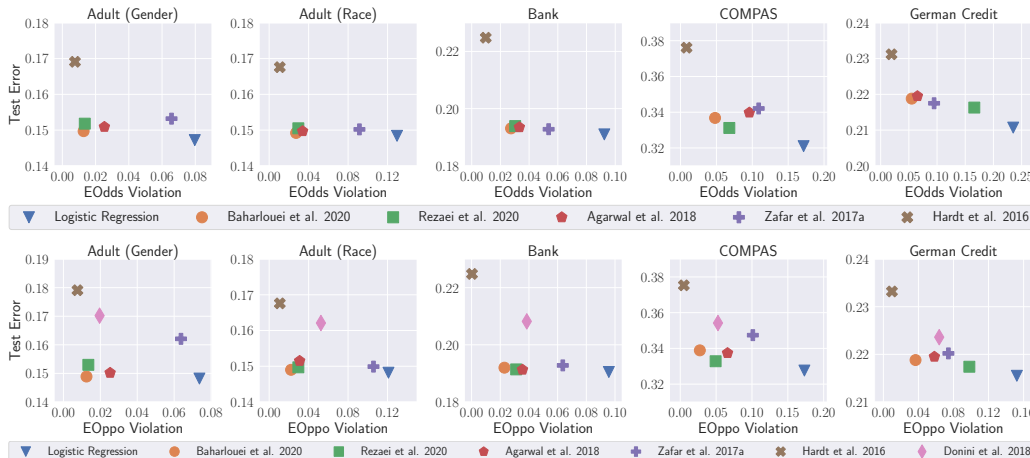


Figure 3: Results for classification with Equalized Odds/Equal Opportunity criterion.

existing methods in the literature on the *Adult*, *Bank*, *COMPAS*, and *German Credit* data sets. The detailed description of the data sets is available in the appendix.

5.1 REGRESSION WITH LINEAR NON-GAUSSIAN AND NONLINEAR DATA

In Section 3.1 we showed the unattainability of Equalized Odds for regression with linear non-Gaussian data. Although a proof for similar results in nonlinear cases does not seem straightforward, as strongly suggested by our numerical illustrations, the unattainability of Equalized Odds persists in nonlinear regression cases. In Figure 2 we present scatter plots of \tilde{Y} versus A for Y in a small (compared to its support) interval, for linear non-Gaussian as well as nonlinear regression cases.

For linear cases, the data is generated as stated in Equation 4, with non-Gaussian distributed exogenous terms (E_X , E_H , and E_Y). We use linear regression with the *Equalized Correlations* constraint (Woodworth et al., 2017), a weaker notion of Equalized Odds for linearly correlated data, as the predictor. For nonlinear cases, the data is generated using a similar scheme but with nonlinear transformations (e.g., combinations of $\sin(\cdot)$, $\log(\cdot)$, and polynomials) and Gaussian distributed exogenous terms. We use a neural net regressor with an Equalized Odds regularization term (Mary et al., 2019) to perform nonlinear fair regression. As we can see in Figure 2, for nonlinear regression tasks, Equalized Odds may not be attained even if every exogenous term is Gaussian distributed.

5.2 FAIR CLASSIFICATION

In Figure 3, we compare the performance under Equalized Odds of multiple methods proposed in the literature. Hardt et al. (2016) propose a post-processing approach where the prediction is randomized to minimize violation of fairness; Zafar et al. (2017a) use a covariance proxy measure as the regularization term when optimizing classification accuracy; Agarwal et al. (2018) take the reductions approach and reduce fair classification into solving a sequence of cost-sensitive classification problems; Rezaei et al. (2020) minimize the worst-case log loss using an approximated regularization term; Baharlouei et al. (2020) propose to use Rényi correlation as the regularization term to account for nonlinear dependence between variables. To measure the violation of the fairness criterion, we use Equalized Odds (EOdds) violation, defined as $\max_{y \in \mathcal{Y}} \max_{a, a' \in \mathcal{A}} |P_{\tilde{Y}|AY}(1|a, y) - P_{\tilde{Y}|AY}(1|a', y)|$. Following Agarwal et al. (2018), we pick 0.01 as the default violation bound that the EOdds violation does not exceed (if practically achievable for the method) during training. For each method we plot the testing accuracy versus the violation of Equalized Odds. Although a probabilistic classification model is used across each method (here is logistic regression), if an algorithm output the class label where the prediction likelihood is maximized, the prediction is in essence performed by a deterministic function of input features (e.g., Rezaei et al. (2020); Baharlouei et al. (2020)). As we have shown in Section 4.1, for classification with a deterministic function, in general cases the conditions specified in Theorem 4.1 are easily violated, i.e., Equalized Odds may not be

attained even if there is an unlimited amount of data. Therefore, although here we are considering finite data cases, we can still anticipate a lower level of fairness violation with a randomized prediction. This is validated by the numerical experiment: while the approach by Hardt et al. (2016) does not score the lowest test error, the violation of Equalized Odds is the lowest compared to other approaches. The benefit of introducing randomization can also be witnessed by the Pareto frontier presented in Agarwal et al. (2018), where the approach can potentially achieve any desired fairness-accuracy trade-off between that of the post-processing approach and that of the unconstrainedly optimized classifier². In some scenarios people tend to only care about equal TPR (e.g., the rate of acceptance/admission) across groups, i.e., the Equal Opportunity (Hardt et al., 2016) notion of fairness. The related numerical result on real-world data sets is also presented.

6 CONCLUSION AND FUTURE WORK

In this paper, we focus on the Equalized-Odds criterion and consider the attainability of fairness, and furthermore, if attainable, the optimality of the prediction performance under various settings. We first show that, for fair regression, one can only achieve Equalized Odds when certain conditions on the joint distribution of the features and the target variable are met. Then for classification tasks with deterministic classifiers, we give the condition under which Equalized Odds can hold true; we also show that under mild assumptions, one can always find a non-trivial Equalized-Odds (randomized) predictor, even with a continuous protected feature; in terms of the optimality of performance, one can always (if conditions permit) benefit from exploiting all available features during training. Future work would naturally consider nonlinear regression algorithms with randomized output and fairness guarantees, and the attainability of more fine-grained (compared to group fairness) criteria of fairness (e.g., individual fairness) as well as the *procedure fairness* in the fairness hierarchy.

REFERENCES

- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *International Conference on Machine Learning*, pp. 60–69, 2018.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. Machine bias: There’s software used across the country to predict future criminals, and it’s biased against blacks. *ProPublica*, 2016.
- Sina Baharlouei, Maher Nouiehed, Ahmad Beirami, and Meisam Razaviyayn. Rényi fair inference. In *International Conference on Learning Representations*, 2020.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness in machine learning. *NIPS Tutorial*, 2017.
- Richard Berk, Hoda Heidari, Shahin Jabbari, Matthew Joseph, Michael Kearns, Jamie Morgenstern, Seth Neel, and Aaron Roth. A convex framework for fair regression. *arXiv preprint arXiv:1706.02409*, 2017.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Advances in neural information processing systems*, pp. 4349–4357, 2016.
- Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pp. 13–18. IEEE, 2009.
- Silvia Chiappa. Path-specific counterfactual fairness. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 7801–7808, 2019.

²In the approach proposed by Agarwal et al. (2018), the randomization can come in two folds: the first kind of randomization comes from picking a classifier from the distribution of multiple available classifiers; the second kind of randomization comes from the probabilistic prediction (if the hypothesis class contains probabilistic prediction models).

- Adrián Pérez-Suay, Valero Laparra, Gonzalo Mateo-García, Jordi Muñoz-Marí, Luis Gómez-Chova, and Gustau Camps-Valls. Fair kernel learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pp. 339–355. Springer, 2017.
- Ashkan Rezaei, Rizal Fathony, Omid Memarrast, and Brian Ziebart. Fairness for robust log loss classification. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- Chris Russell, Matt J Kusner, Joshua Loftus, and Ricardo Silva. When worlds collide: integrating different counterfactual assumptions in fairness. In *Advances in Neural Information Processing Systems*, pp. 6414–6423, 2017.
- Jiaming Song, Pratyusha Kalluri, Aditya Grover, Shengjia Zhao, and Stefano Ermon. Learning controllable fair representations. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2164–2173, 2019.
- Latanya Sweeney. Discrimination in online ad delivery. *Queue*, 11(3):10–29, 2013.
- Blake Woodworth, Suriya Gunasekar, Mesrob I Ohannessian, and Nathan Srebro. Learning non-discriminatory predictors. In *Conference on Learning Theory*, pp. 1920–1953, 2017.
- Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. Pc-fairness: A unified framework for measuring causality-based fairness. In *Advances in Neural Information Processing Systems*, pp. 3399–3409, 2019.
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rodriguez, and Krishna P Gummadi. 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, 2017a.
- Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P Gummadi. Fairness constraints: Mechanisms for fair classification. In *Artificial Intelligence and Statistics*, pp. 962–970, 2017b.
- Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International Conference on Machine Learning*, pp. 325–333, 2013.
- Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 335–340, 2018.
- Junzhe Zhang and Elias Bareinboim. Fairness in decision-making—the causal explanation formula. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- Lu Zhang, Yongkai Wu, and Xintao Wu. A causal framework for discovering and removing direct and indirect discrimination. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pp. 3929–3935, 2017.
- Han Zhao, Amanda Coston, Tameem Adel, and Geoffrey J. Gordon. Conditional learning of fair representations. In *International Conference on Learning Representations*, 2020.

A APPENDIX

A.1 PROOF FOR THEOREM 3.1

To prove the unattainability of Equalized Odds in regression, we will need the following lemma, which provides a way to characterize conditional independence/dependence with conditional or joint distributions.

Lemma A.1. *Variables V_1 and V_2 are conditionally independent given variable V_3 if and only if there exist functions $h(v_1, v_3)$ and $g(v_2, v_3)$ such that*

$$p_{V_1, V_2 | V_3}(v_1, v_2 | v_3) = h(v_1, v_3) \cdot g(v_2, v_3). \quad (10)$$

Proof. First, if V_1 and V_2 are conditionally independent given variable V_3 , then Equation 10 holds:

$$p_{V_1, V_2 | V_3}(v_1, v_2 | v_3) = p_{V_1 | V_3}(v_1 | v_3) \cdot p_{V_2 | V_3}(v_2 | v_3).$$

We then let $\tilde{h}(v_3) := \int h(v_1, v_3) dv_1$ and $\tilde{g}(v_3) := \int g(v_2, v_3) dv_2$. Take the integral of Equation 10 w.r.t. v_1 and v_2 , we have:

$$\begin{aligned} p_{V_2 | V_3}(v_2 | v_3) &= \tilde{h}(v_3) \cdot g(v_2, v_3), \\ p_{V_1 | V_3}(v_1 | v_3) &= \tilde{g}(v_3) \cdot h(v_1, v_3), \end{aligned}$$

respectively. Bearing in mind Equation 10, one can see that the product of the two equations above is

$$\begin{aligned} p_{V_2 | V_3}(v_2 | v_3) \cdot p_{V_1 | V_3}(v_1 | v_3) &= \tilde{h}(v_3) \cdot g(v_2, v_3) \cdot \tilde{g}(v_3) \cdot h(v_1, v_3) \\ &= \tilde{h}(v_3) \cdot \tilde{g}(v_3) \cdot p_{V_1, V_2 | V_3}(v_1, v_2 | v_3). \end{aligned}$$

Take the integral of the equation above w.r.t. v_1 and v_2 gives $\tilde{h}(v_3) \cdot \tilde{g}(v_3) \equiv 1$. The above equation then reduces to

$$p_{V_2 | V_3}(v_2 | v_3) \cdot p_{V_1 | V_3}(v_1 | v_3) = p_{V_1, V_2 | V_3}(v_1, v_2 | v_3).$$

That is, V_1 and V_2 are conditionally independent given V_3 . \square

Now we are ready to prove the unattainability of Equalized Odds in linear non-Gaussian regression:

Theorem. (Unattainability of Equalized Odds in the Linear Non-Gaussian Case)

Assume that feature X has a causal influence on Y , i.e., $c \neq 0$ in Equation 4, and that the protected feature A and Y are not independent, i.e., $qc + bd \neq 0$. Assume p_{E_X} and p_E are positive on \mathbb{R} . Let $f_1 := \log p_A$, $f_2 := \log p_{E_X}$, and $f_3 := \log p_E$. Further assume that f_2 and f_3 are third-order differentiable. Then if at most one of E_X and E is Gaussian, Z is always conditionally dependent on A given Y .

Proof. According to Equation 4, we have

$$\begin{bmatrix} A \\ Z \\ Y \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ \alpha + q\beta & \beta & 0 \\ qc + bd & c & 1 \end{bmatrix} \cdot \begin{bmatrix} A \\ E_X \\ E \end{bmatrix}. \quad (11)$$

The determinant of the above linear transformation is β , which relates the probability density function of the variables on the LHS and that of the variables on the RHS of the equation. Therefore, according to Equation 11, we can rewrite the joint probability density function by making use of the Jacobian determinant and factor the joint density into marginal density functions (A , E_X , E are mutually independent according to the data generating process). Further let

$$\tilde{\alpha} := \frac{\alpha + q\beta}{\beta}, \quad \tilde{r} := bd - \frac{c\alpha}{\beta}, \quad \text{and} \quad \tilde{c} := \frac{c}{\beta}. \quad (12)$$

Then we have $E_X = \frac{1}{\beta}Z - \tilde{\alpha}A$, $E = Y - \tilde{r}A - \tilde{c}Z$, and

$$\begin{aligned} p_{A, Z, Y}(a, z, y) &= p_{A, E_X, E}(a, e_x, e) / |\beta| \\ &= \frac{1}{|\beta|} p_A(a) p_{E_X}(e_x) p_E(e) \\ &= \frac{1}{|\beta|} p_A(a) p_{E_X}\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) p_E(y - \tilde{r}a - \tilde{c}z). \end{aligned}$$

On its support, the log-density can be written as

$$\begin{aligned}
J &:= \log p_{A,Z,Y}(a, z, y) \\
&= \log p_A(a) + \log p_{E_X}\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) + \log p_E(y - \tilde{r}a - \tilde{c}z) - \log|\beta| \\
&= f_1(a) + f_2\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) + f_3(y - \tilde{r}a - \tilde{c}z) - \log|\beta|.
\end{aligned} \tag{13}$$

According to Lemma A.1, $A \perp\!\!\!\perp Z \mid Y$ if and only if $p_{A,Z|Y}(a, z \mid y)$ is a product of a function of a and y and a function of z and y . $p_{A,Z,Y}(a, z, y)$ is further a product of the above function and a function of only y . This property, under the conditions in Theorem 3.1, is equivalent to the constraint

$$\frac{\partial^2 J}{\partial A \partial Z} \equiv 0. \tag{14}$$

According to Equation 13, we have

$$\begin{aligned}
\frac{\partial J}{\partial z} &= \frac{1}{\beta} \cdot f_2'\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) - \tilde{c} \cdot f_3'(y - \tilde{r}a - \tilde{c}z) \\
\Rightarrow \frac{\partial^2 J}{\partial a \partial z} &= -\frac{\tilde{\alpha}}{\beta} \cdot f_2''\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) + \tilde{r}\tilde{c} \cdot f_3''(y - \tilde{r}a - \tilde{c}z).
\end{aligned} \tag{15}$$

Combining Equations 14 and 15 gives

$$\tilde{r}\tilde{c} \cdot f_3''(y - \tilde{r}a - \tilde{c}z) = \frac{\tilde{\alpha}}{\beta} \cdot f_2''\left(\frac{1}{\beta}z - \tilde{\alpha}a\right). \tag{16}$$

Further taking the partial derivative of both sides of the above equation w.r.t. y yields

$$\tilde{r}\tilde{c} \cdot f_3'''(y - \tilde{r}a - \tilde{c}z) \equiv 0. \tag{17}$$

There are three possible situations where the above equation holds:

- (i) $\tilde{c} = 0$, which is equivalent to $c = 0$ and contradicts with the theorem assumption.
- (ii) $\tilde{r} = 0$. Then according to Equation 16, we have $\frac{\tilde{\alpha}}{\beta} \cdot f_2''\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) \equiv 0$, implies either $\tilde{\alpha} = 0$ or $f_2''\left(\frac{1}{\beta}z - \tilde{\alpha}a\right) \equiv 0$. If the latter is the case, then f_2 is a linear function and, accordingly, $\exp(f_2)$ is not integrable and does not correspond to any valid density function. If the former is true, i.e., $\tilde{\alpha} = 0$, then according to Equation 12, we have $\alpha = -q\beta$, which further implies $\tilde{r} = bd - \frac{c\alpha}{\beta} = bd + qc$. Therefore, in this situation, $bd + qc = 0$, which again contradicts with the theorem assumption.
- (iii) $f_3'''(y - \tilde{r}a - \tilde{c}z) \equiv 0$. That is, f_3 is a quadratic function with a nonzero coefficient for the quadratic term (otherwise f_3 does not correspond to the logarithm of any valid density function). Thus E follows a Gaussian distribution.

Only situation (iii) is possible, i.e., $\tilde{r}\tilde{c} \neq 0$ and E follows a Gaussian distribution. This further tells us that the RHS of Equation 16 is a nonzero constant. Hence f_2 is a quadratic function and E_X also follows a Gaussian distribution. Therefore if $A \perp\!\!\!\perp Z \mid Y$ were to be true, then E_X and E are both Gaussian. Its contrapositive gives the conclusion of this theorem. \square

Corollary. *Suppose that both E_X and E are Gaussian, with variances $\sigma_{E_X}^2$ and σ_E^2 , respectively. (The protected feature A is not necessarily Gaussian.) Then $Z \perp\!\!\!\perp A \mid Y$ if and only if*

$$\frac{\alpha}{\beta} = \frac{bdc \cdot \sigma_{E_X}^2 - q \cdot \sigma_E^2}{c^2 \cdot \sigma_{E_X}^2 + \sigma_E^2}. \tag{18}$$

Proof. Under the condition that E_X and E are Gaussian, their log-density functions are third-order differentiable. Then according to the proof of Theorem 3.1, the Equalized Odds condition $A \perp\!\!\!\perp Z \mid Y$ is equivalent to Equation 16, which, together with Equation 12 as well as the fact that $f_2'' = \frac{1}{\sigma_{E_X}^2}$ and $f_3'' = \frac{1}{\sigma_E^2}$, yields Equation 18. \square

A.2 PROOF FOR THEOREM 4.1

Theorem. Assume that the protected feature A and Y are dependent and that their joint probability $P(A, Y)$ (for discrete A) or joint probability density $p(A, Y)$ (for continuous A) is positive for every combination of possible values of A and Y . Further assume that Y is not fully determined by A , and that there are additional features X that are not independent of Y . Let the output of the classifier \tilde{Y} be a deterministic function $f : \mathcal{A} \times \mathcal{X} \rightarrow \mathcal{Y}$. Let $S_A^{(t)} := \{a \mid \exists x \in \mathcal{X} \text{ s.t. } f(a, x) = t\}$, and $S_{X|a}^{(t)} := \{x \mid f(a, x) = t\}$. Equalized Odds holds true if and only if the following two conditions are satisfied:

- (i) $\forall t \in \mathcal{Y} : S_A^{(t)} = \mathcal{A}$,
- (ii) $\forall t \in \mathcal{Y}, \forall a, a' \in \mathcal{A}, a \neq a' :$

$$\sum_{x \in S_{X|a}^{(t)}} P(X = x \mid A = a, Y = y) = \sum_{x \in S_{X|a'}^{(t)}} P(X = x \mid A = a', Y = y).$$

Proof. We begin by considering the case when A and X are discrete (for the purpose of readability). The Equalized Odds criterion can be written in terms of the conditional probabilities:

$$\forall a \in \mathcal{A}, t, y \in \mathcal{Y} : P(\tilde{Y} = t \mid A = a, Y = y) = P(\tilde{Y} = t \mid Y = y). \quad (19)$$

Expand the LHS of Equation 19:

$$P(\tilde{Y} = t \mid A = a, Y = y) = \sum_{x \in \mathcal{X}} P(\tilde{Y} = t \mid A = a, X = x, Y = y)P(X = x \mid A = a, Y = y),$$

and bear in mind that $\tilde{Y} := f(A, X)$ is a deterministic function of (A, X) , we have:

$$P(f(A, X) = t \mid A = a, X = x, Y = y) = P(f(A, X) = t \mid A = a, X = x) \in \{0, 1\}. \quad (20)$$

From Equation 20 we can see that the conditional probability $P(X = x \mid A = a, Y = y)$ can contribute to the summation only when $f(a, x) = t$. We can rewrite the LHS of Equation 19:

$$P(\tilde{Y} = t \mid A = a, Y = y) = \sum_{x \in S_{X|a}^{(t)}} P(X = x \mid A = a, Y = y) := Q^{(t)}(a, y).$$

Similarly, for the RHS of Equation 19, we have:

$$\begin{aligned} P(\tilde{Y} = t \mid Y = y) &= \sum_{a \in \mathcal{A}} \sum_{x \in \mathcal{X}} P(\tilde{Y} = t \mid A = a, X = x, Y = y)P(X = x, A = a \mid Y = y) \\ &= \sum_{a \in S_A^{(t)}} \sum_{x \in S_{X|a}^{(t)}} P(X = x \mid A = a, Y = y)P(A = a \mid Y = y) \\ &= \sum_{a \in S_A^{(t)}} Q^{(t)}(a, y)P(A = a \mid Y = y). \end{aligned}$$

Since Equalized Odds holds true if and only if Equation 19 holds true, then the LHS of the equation does not involve a (as is the case for the RHS), i.e., $Q^{(t)}(a, y)$ does not change with a . Then Equation 19 becomes:

$$Q^{(t)}(a, y) = \sum_{a \in S_A^{(t)}} Q^{(t)}(a, y)P(A = a \mid Y = y) = Q^{(t)}(a, y) \sum_{a \in S_A^{(t)}} P(A = a \mid Y = y),$$

which gives the condition (i) that $S_A^{(t)}$ contains all possible values of A , i.e., $\mathcal{A} = S_A^{(t)}$ (otherwise $\sum_{a \in S_A^{(t)}} P(A = a \mid Y = y) < 1$). Since $Q^{(t)}(a, y)$ does not change with a , we have:

$$\forall a, a' \in \mathcal{A}, a \neq a' : \sum_{x \in S_{X|a}^{(t)}} P(X = x \mid A = a, Y = y) = \sum_{x \in S_{X|a'}^{(t)}} P(X = x \mid A = a', Y = y),$$

which gives the condition (ii). Therefore, Equalized Odds implies conditions (i) and (ii). On the other hand, it is easy to see that when conditions (i) and (ii) are satisfied, Equation 19 holds true, i.e., Equalized Odds holds true.

When A and X are continuous, one can replace the summation with integration accordingly. \square

A.3 PROOF FOR THEOREM 4.2

Theorem. (Attainability of Equalized Odds)

Assume that the feature X is not independent from Y , and that \widehat{Y} is a function of A and X . Then for binary classification, if \widehat{Y} is a non-trivial predictor for Y , there is always at least one non-trivial predictor \widetilde{Y}_{post} derived by post-processing \widehat{Y} that can attain Equalized Odds, i.e.,

$$\Omega(\widetilde{Y}_{post}) \neq \emptyset.$$

Proof. Since \widehat{Y} is a function of (A, X) and $X \not\perp Y$, \widehat{Y} is not conditionally independent from Y given protected feature A . Furthermore, since \widehat{Y} is a non-trivial estimator of the binary target Y , there exists a positive constant $\epsilon > 0$, such that:

$$|P(\widehat{Y} = 1 | A = a, Y = 1) - P(\widehat{Y} = 1 | A = a, Y = 0)| \geq \epsilon, \forall a \in \mathcal{A}. \quad (21)$$

Equation 21 implies that for each value of A , the corresponding true positive rate of the non-trivial predictor is always strictly larger than its false positive rate³. As illustrated in panels (a) and (b) of Figure 1, the (FPR, TPR) pair of the predictor \widehat{Y} when $A = a$, i.e., the point $\gamma_a(\widehat{Y})$ on ROC plane, will never fall in the gray shaded area, and its coordinates are bounded away from the diagonal by at least ϵ . Therefore, the intersection of all $\mathcal{C}_a(\widehat{Y})$ would always form a parallelogram with non-empty area, which corresponds to attainable non-trivial post-processing fair predictors \widetilde{Y}_{post} . \square

A.4 PROOF FOR THEOREM 4.3

Theorem. (Equivalence between ROC feasible areas)

Let $\Omega(\widetilde{Y}_{post})$ denote the ROC feasible area specified by the constraints enforced on \widetilde{Y}_{post} . Then $\Omega(\widetilde{Y}_{post})$ is identical to the ROC feasible area $\Omega(\widetilde{Y}_{in}^*)$ that is specified by the following set of constraints:

- (i) constraints enforced on \widetilde{Y}_{in} ;
- (ii) additional ‘‘pseudo’’ constraints: $\forall a \in \mathcal{A}$, $\beta_{a0}^{(0)} = \beta_{a1}^{(0)}$, $\beta_{a0}^{(1)} = \beta_{a1}^{(1)}$, where $\beta_{ay}^{(u)} = \sum_{x \in \mathcal{X}} P(\widetilde{Y}_{in} = 1 | A = a, X = x)P(X = x | A = a, Y = y, \widehat{Y}_{opt} = u)$.

Proof. Since the post-processing predictor \widetilde{Y}_{post} is derived by optimizing over parameters or functions (of A) $\beta_a^{(u)}$. Therefore, considering the fact that $P_{\widetilde{Y}_{post}|AY}(1|a, y) = \gamma_{ay}(\widetilde{Y}_{post})$, $P_{\widehat{Y}_{opt}|AY}(1|a, y) = \gamma_{ay}(\widehat{Y}_{opt})$, we have the relation between $\gamma_{ay}(\widetilde{Y}_{post})$ and $\gamma_{ay}(\widehat{Y}_{opt})$:

$$\begin{aligned} \gamma_{ay}(\widetilde{Y}_{post}) &= \beta_a^{(0)} \gamma_{ay}(1 - \widehat{Y}_{opt}) + \beta_a^{(1)} \gamma_{ay}(\widehat{Y}_{opt}), \\ \beta_a^{(0)} &= P(\widetilde{Y}_{post} = 1 | A = a, \widehat{Y}_{opt} = 0), \\ \beta_a^{(1)} &= P(\widetilde{Y}_{post} = 1 | A = a, \widehat{Y}_{opt} = 1). \end{aligned} \quad (22)$$

Similarly, consider the relation between positive rates of \widetilde{Y}_{in} and those of \widehat{Y}_{opt} , i.e., $P_{\widetilde{Y}_{in}|AY}(1|a, y)$ and $P_{\widehat{Y}_{opt}|AY}(1|a, y)$, by factorizing $P_{\widetilde{Y}_{in}|AY}(1|a, y)$ over X and \widehat{Y}_{opt} :

$$P_{\widetilde{Y}_{in}|AY}(1|a, y) = \sum_{u \in \mathcal{Y}} \left[\sum_{x \in \mathcal{X}} P_{\widetilde{Y}_{in}|AX}(1|a, x) P_{X|AY\widehat{Y}_{opt}}(x|a, y, u) \right] P_{\widehat{Y}_{opt}|AY}(u|a, y). \quad (23)$$

Therefore, we have the relation between $\gamma_{ay}(\widetilde{Y}_{in})$ and $\gamma_{ay}(\widehat{Y}_{opt})$:

$$\begin{aligned} \gamma_{ay}(\widetilde{Y}_{in}) &= \beta_{ay}^{(0)} \gamma_{ay}(1 - \widehat{Y}_{opt}) + \beta_{ay}^{(1)} \gamma_{ay}(\widehat{Y}_{opt}), \\ \beta_{ay}^{(0)} &= \sum_{x \in \mathcal{X}} P(\widetilde{Y}_{in} = 1 | A = a, X = x)P(X = x | A = a, Y = y, \widehat{Y}_{opt} = 0), \\ \beta_{ay}^{(1)} &= \sum_{x \in \mathcal{X}} P(\widetilde{Y}_{in} = 1 | A = a, X = x)P(X = x | A = a, Y = y, \widehat{Y}_{opt} = 1). \end{aligned} \quad (24)$$

³If the TPR of the predictor is always smaller than its FPR, one can simply flip the prediction (since the target is binary) and then Equation 21 holds true.

If there is more than one variable in X in Equation 24, one can expand the summation if needed; if some variables are continuous, one may also substitute the summation with integration accordingly.

From Equation 24, $\beta_{ay}^{(0)}$ and $\beta_{ay}^{(1)}$ depend on the value of Y :

$$\begin{aligned}\beta_{ay}^{(0)} &= P(\tilde{Y}_{\text{in}} = 1 \mid A = a, Y = y, \hat{Y}_{\text{opt}} = 1), \\ \beta_{ay}^{(1)} &= P(\tilde{Y}_{\text{in}} = 1 \mid A = a, Y = y, \hat{Y}_{\text{opt}} = 0).\end{aligned}\tag{25}$$

Apart from Equalized Odds constraints (which are shared by \tilde{Y}_{in} and \tilde{Y}_{post}), when enforcing additional “pseudo” constraints $\beta_{a0}^{(0)} = \beta_{a1}^{(0)}$ and $\beta_{a0}^{(1)} = \beta_{a1}^{(1)}$, conditional independence $\tilde{Y}_{\text{in}} \perp\!\!\!\perp Y \mid A, \hat{Y}_{\text{opt}}$ is enforced, making $\beta_{ay}^{(0)}$ and $\beta_{ay}^{(1)}$ no longer depend on Y . This is exactly the inherent constraint \tilde{Y}_{post} satisfies. Therefore the stated equivalence between ROC feasible areas $\Omega(\tilde{Y}_{\text{post}})$ (specified by the constraints enforced on \tilde{Y}_{post}) and $\Omega(\tilde{Y}_{\text{in}}^*)$ (specified by the constraints enforced on \tilde{Y}_{in} together with the additional “pseudo” constraints) hold true. \square

A.5 DESCRIPTION OF THE DATA SETS

- (1) **Adult**⁴: The UCI Adult data set contains 14 features for 45,222 individuals (32,561 samples for training and 12,661 samples for testing). The census information includes gender, marital status, education, capital gain, etc. The classification task is to predict whether a person’s annual income exceeds 50,000 USD. We use the provided testing set for evaluations and present the result with gender and race (consider white and black people only) set as the protected feature respectively.
- (2) **Bank**⁵: The UCI Bank Marketing data set is related with marketing campaigns of a banking institution, containing 16 features of 45,211 individuals. The assigned classification task is to predict if a client will subscribe (yes/no) to a term deposit. The original data set is very unbalanced with only 4,667 positives out of 45,211 samples. Therefore, we combine “yes” points with randomly subsampled “no” points and perform experiments on the down-sampled data set with 10,000 data points. The protected feature is the marital status of the client.
- (3) **COMPAS** (Angwin et al., 2016): The COMPAS data set contains records of over 11,000 defendants from Broward County, Florida, whose risk (of recidivism) was assessed using the COMPAS tool. Each record contains multiple features of the defendant, including demographic information, prior convictions, degree of charge, and the ground truth for recidivism within two years. Following Zafar et al. (2017a); Nabi & Shpitser (2018), we limit our attention to the subset consisting of African-Americans and Caucasians defendants. The features we use include age, gender, race, number of priors, and degree of charges. The task is to predict the recidivism of the defendant and we choose race as the protected feature.
- (4) **German Credit**⁶: The UCI German Credit data contains 20 features (7 numerical, 13 categorical) describing the social and economical status of 1,000 customers. The prediction task is to classify people as good or bad credit risks. We use the provided numerical version of the data and choose gender as the protected feature.

A.6 ADDITIONAL DISCUSSION

For classification, while randomization can ensure group level of fairness, there is still some inherent shortcoming of the criterion that we should pay attention to. For example, in the FICO case study in Hardt et al. (2016), for a specific client from certain demographic group, the decision of approve/deny the loan actually comes in two folds: if his/her credit score is above (below) the upper (lower) threshold, the bank approve (deny) the application for sure; if the score falls in the interval between two thresholds, the bank would flip a coin to make a decision. Then we can imagine the following situation when a client whose credit score falls within the interval between the upper and

⁴<http://archive.ics.uci.edu/ml/datasets/Adult>

⁵<https://archive.ics.uci.edu/ml/datasets/bank+marketing>

⁶[https://archive.ics.uci.edu/ml/datasets/statlog+\(german+credit+data\)](https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data))

lower thresholds goes to a bank to apply a loan. He/she can ask (if conditions permit) the bank to run the model multiple times until the decision is approval. This would make the randomization that was built into the system for the sake of fairness no longer effective, and the system in essence only has one fixed threshold (i.e., the original lower threshold).