
Washing The Unwashable : On The (Im)possibility of Fairwashing Detection

Ali Shahin Shamsabadi*
The Alan Turing Institute
Vector Institute

Mohammad Yaghini*
University of Toronto
Vector Institute

Natalie Dullerud*
University of Toronto
Vector Institute

Sierra Wyllie
University of Toronto
Vector Institute

Ulrich Aïvodji
ÉTS Montréal

Aisha Alaagib
University of Toronto
Vector Institute

Sébastien Gambs
Université du Québec à Montréal

Nicolas Papernot
University of Toronto
Vector Institute

Abstract

The use of black-box models (*e.g.*, deep neural networks) in high-stakes decision-making systems, whose internal logic is complex, raises the need for providing explanations about their decisions. Model explanation techniques mitigate this problem by generating an interpretable and high-fidelity surrogate model (*e.g.*, a logistic regressor or decision tree) to explain the logic of black-box models. In this work, we investigate the issue of fairwashing, in which model explanation techniques are manipulated to rationalize decisions taken by an unfair black-box model using deceptive surrogate models. More precisely, we theoretically characterize and analyze fairwashing, proving that this phenomenon is difficult to avoid due to an irreducible factor—the unfairness of the black-box model. Based on the theory developed, we propose a novel technique, called FRAUD-Detect (FaiRness AUDit Detection), to detect fairwashed models by measuring a divergence over subpopulation-wise fidelity measures of the interpretable model. We empirically demonstrate that this divergence is significantly larger in purposefully fairwashed interpretable models than in honest ones. Furthermore, we show that our detector is robust to an informed adversary trying to bypass our detector. The code implementing FRAUD-Detect is available at <https://github.com/cleverhans-lab/FRAUD-Detect>.

1 Introduction

The wide applicability of machine learning models has recently increased their usage in high-stakes decision systems such as credit scoring [43], insurance risk [10] and predictive justice [30]. The consequences of erroneous decisions that are based on predictions of machine learning models (*e.g.*, people being wrongly denied parole [46]) have increased the demand—from both the public and government—to provide an explanation to humans about model decisions. For instance, this appears as an explanation requirement in the European General Data Protection Regulation [26]. However, various widely-used model architectures, such as deep neural networks, are considered black-boxes due to their complex and hidden internal logic, impeding the ability to explain their decisions in terms that are understandable by a human. To address this issue, black-box model

*Contributed equally.

using the corresponding fairwashed interpretable model. In this paper, we are the first to propose a method for detecting such dishonest model explanations. This is a challenging problem. Fairwashing cannot be detected through computing the explanation violation of the explainable model, such as fidelity of explainable models, as we demonstrate that the explainable model cannot achieve perfect fidelity with respect to the black-box model. In addition to this, measuring the differences between the fairness of the black-box model and the explainable model cannot help to detect fairwashing as there are several design choices (including fairwashing) that could lead to different fairness, thus non-intentional fairwashing.

Our primary contributions include (1) an extensive theoretical analysis of fairwashing and the fairness limitations in explanation techniques and (2) a novel approach for detecting, and thus deterring, fairwashing. Our proposed method, FRAUD-Detect (Fairness AUDit Detection), formulates the problem of detecting fairwashing as a non-cooperative test; both to overcome possible dishonest behaviors in entities that are actively being audited and for communication efficiency. FRAUD-Detect operates in a realistic scenario: FRAUD-Detect does not require access to the black-box model provided the users calling for the audit have provided their data and received decisions, which is a realistic setting as often the model is not provided by the entity due to intellectual property and trade-secret concerns [39, 37]. Thus, FRAUD-Detect only relies on the predictions of the interpretable and black-box models. Both our theoretical and empirical analysis quantify the per-subpopulation, per-label fidelity of the interpretable model with respect to the black-box model.

Figure 1 shows confusion matrix distributions for two sensitive subpopulations before and after fairwashing. These matrices are constructed by comparing the interpretable to the black-box model predictions. Fairwashing induces a divergence between subpopulation values to conceal the unfairness of the black-box model. FRAUD-Detect leverages the Kullback–Leibler (KL) divergence between subpopulation confusion matrices to distinguish between honest and fairwashed interpretable models. For a comprehensive empirical study, we examine logistic regressors and decision trees as interpretable models with respect to black-box Deep Neural Networks, AdaBoost [24], Gradient Boosted Decision Trees [17] and Random Forests [13] on three benchmark datasets: COMPAS [5], Adult Income [22] and Bank Marketing [34].

We evaluate the strength of our detector in the presence of an *informed adversary* representing an informed dishonest entity with knowledge of FRAUD-Detect that seeks to fairwash *and* bypass detection by constraining the divergence over subpopulation-wise confusion matrices. Our empirical experiments quantify the capacity of this informed adversary by computing the achievable fairness gap and fidelity under the additional detection constraint. We provide additional theoretical support for the challenge of jointly satisfying fairwashing and evasion in appendices.

In summary, our main contributions are as follows:

- We characterize fairwashing via a theoretical analysis on the difference in subpopulation gap between the black-box model and the interpretable model, *fairwashing violation* (Section 3).
- Based on our fairwashing theory, we establish requirements for a *sufficient* fairwashing detection method and prove that fairwashing is impossible to avoid completely due to an *irreducible* component in the violation term.
- We introduce the first method for fairwashing detection (Section 4.2). Informed by our theoretical results, we observe that fairwashing causes disparate effect on the subpopulation-wise fidelity distributions of the interpretable model with respect to the black-box predictions. We leverage this observation to propose a non-cooperative black-box access method, FRAUD-Detect, which utilizes KL divergence on per-subpopulation confusion matrices to distinguish between fairwashed and honest interpretable models.
- Our empirical results demonstrate that FRAUD-Detect successfully detects fairwashing based on the KL divergence between subpopulation-wise confusion matrices (Section 6). We show that the divergence over subpopulation confusion matrices can vary by over 0.6 between an honest and a fairwashed interpretable models.
- We illustrate the robustness of FRAUD-Detect against an informed adversary (*i.e.*, a dishonest entity who attempts to fairwash while evading detection). Our empirical results show that evading our detector comes at the cost of a significant increase in subpopulation gap, negating fairwashing. Specifically, a dishonest entity that jointly attempts to fairwash while evading detection only achieves parity gaps greater than 10%.

2 Preliminaries and Problem Formulation

2.1 Related Work

Fairness. Many different formal definitions of algorithmic fairness have been proposed in the machine learning community [35]. One of the main challenges of algorithmic fairness is the lack of consensus on a universally applicable fairness definition. More precisely, it has been formally demonstrated that many such measures are incompatible with each other (*i.e.*, they cannot be achieved jointly in some situations) [31]. In addition, their differences are also rooted in their underlying philosophical and moral assumptions [28]. Thus, the choice of a particular fairness metric is usually context dependent. Nonetheless broadly speaking, there are two main families of fairness definitions: individual fairness and group fairness [23]. Individual fairness posits that “similar individuals should be treated similarly”. In contrast, group fairness strives to equalize statistical properties of classification outcomes across subpopulations created by partitioning the population based on a sensitive attribute A such as gender. In our work, we focus on a particular notion of group fairness, demographic parity, also called statistical group parity, which refers to observing equal probability of positive label prediction over subpopulations. One of the reasons we rely on demographic parity as the fairness metric is that true labels will not be available in a realistic scenario (especially in fairwashing scenario described in Section 2.3). In general, the existence of true labels conflicts with the assumption that when requesting decisions (or explanations), users (or auditors) do not have access to these true labels [42]. Furthermore, using the predictions of a black-box model as the true labels is paradoxical in the sense that the black-box model is assumed to be unfair. Demographic parity enforces independence between the class predicted by the model and inclusion in a particular subpopulation.

Definition 1 (Demographic parity [15]). *A classifier \hat{Y} satisfies demographic parity with respect to sensitive attribute A if*

$$\Pr[\hat{Y} = 1|A = a] = \Pr[\hat{Y} = 1|A = b] \quad \forall a, b \in A.$$

Global black-box explanation methods focus on explaining the whole logic of black-box models by training an inherently interpretable surrogate model. We refer the interested reader to Appendix A for the description of other explanation methods. In terms of abusing explanations via fairwashing, Fukuchi et al., similarly to us, try to detect fairwashing but with a different setting [25]. Slack et al. shows an alternate method of fairwashing that evades local model explanation techniques in a setting that is comparable to ours but they do not attempt to detect fairwashing [44]. We explore these additional fairwashing settings in Appendix B.

2.2 Desiderata for Global Black-box Explanations

We consider the setting in which an entity learns a complex black-box model $B(\cdot)$ on a training set $(X_{\text{Tr}}, Y_{\text{Tr}}, A)$, in which A represents a sensitive attribute. We maintain the setup from recent fairwashing literature [3, 2] and consider binary classifiers mapping M features into a binary label $B: \mathbb{R}^M \mapsto \{0, 1\}$. We also consider binary-valued sensitive attributes, $A \in \{0, 1\}$. Note that we do not assume that A is necessarily used during the training of $B(\cdot)$. Given user queries, the entity uses the predictions of $B(\cdot)$ as part of a high-stakes decision system. Later, a group of these users may request an explicit explanation from the entity due to concerns over improper use of the sensitive attribute. Such a scenario frequently arises in high-impact domains, such as bank loans and credit scoring. Let X_{sg} be a suing set comprising the unlabeled data of users demanding an explanation for their particular outcomes (see Figure 1). Consequently, an external auditing entity requires the company to provide explanations in terms that are understandable to humans about the predictions of the black-box. Here, we will assume that the explanation will be provided in the form of a global explanation. This means that a simple interpretable model $I(\cdot)$ will be trained on a dataset X labeled by querying the black-box model $B(\cdot)$ trained on X_{Tr} , such that $I(\cdot)$ accurately reflects and explains the logic of $B(\cdot)$ in terms that are understandable to humans.

Formalizing further the auditing desiderata for black-box model explanations, the interpretable model $I(\cdot)$ must accurately mirror: 1) the output predictions of the black-box model (*i.e.*, fidelity) and 2) the fairness (according to a pre-defined metric) of the black-box model.

Fidelity. The interpretable model should satisfy the fidelity criterion for any $X \sim \mathcal{D}$. Fidelity (defined in Appendix C) is difficult to perfectly achieve and challenging to measure for an auditor over all $X \sim \mathcal{D}$. Thus, we additionally introduce the notion of empirical fidelity.

Definition 2 (Empirical Fidelity). *The empirical fidelity on a dataset $X = \{x_i\}_{i=1}^N$ including N data points x_i is defined as the relative accuracy of $I(\cdot)$ with respect to $B(\cdot)$ on X :*

$$\text{EmpiricalFidelity}(I, B; X) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(I(x_i) = B(x_i)), \quad (1)$$

where the indicator $\mathbb{1}(\cdot, \cdot)$ outputs 1 if the output label of $I(\cdot)$ and $B(\cdot)$ are the same and 0 otherwise.

Fairness. As the audit of $B(\cdot)$ is expressly requested due to fairness concerns with respect to A , the interpretable model should reflect the fairness violations or adherences of the black-box model.

2.3 Fairwashing Definitions

To evade legal consequences of decision-making based on improper use of a sensitive attribute, a dishonest company may perform *fairwashing*. Fairwashing with respect to a pre-defined fairness metric permits the dishonest company to learn an interpretable model that hides the unfair behaviour of their black-box model. In practice, fairwashing could also occur due to the fact that current regulations [26] do not define what constitutes a *valid explanation*, thus leaving the possibility for the model provider to choose the interpretable model that meets their needs [45]. Fairwashing can be quantified with respect to a particular fairness measure.

Definition 3 (Fairwashing). *Let Γ_I and Γ_B define the fairness gaps of the interpretable model and black-box model, respectively, with respect to a pre-defined fairness metric. For example, in terms of demographic parity (Definition 1):*

$$\begin{aligned} \Gamma_I &:= \Pr \left[\hat{Y}_I = 1 \mid A = 0 \right] - \Pr \left[\hat{Y}_I = 1 \mid A = 1 \right], \\ \Gamma_B &:= \Pr \left[\hat{Y}_B = 1 \mid A = 0 \right] - \Pr \left[\hat{Y}_B = 1 \mid A = 1 \right]. \end{aligned} \quad (2)$$

We define the fairwashing violation $\gamma > 0$ as the difference between Γ_I and Γ_B :

$$\gamma := \Gamma_B - \Gamma_I. \quad (3)$$

Assumption 1. *The fairwashing violation γ is non-negative. Remark that the case of $\gamma < 0$ is the opposite of fairwashing since the interpretable model is displaying a bigger fairness gap than that of the black-box model, thus a model owner has no incentive to employ such an interpretable model.*

3 Fairness Gap of the Black-Box Breeds Fairwashing

We analyze whether it is possible to eliminate the risk of fairwashing altogether. We characterize the fairwashing violation γ and show that completely eliminating fairwashing is impossible: an interpretable model explaining an unfair black-box model always has a non-zero fairwashing violation.

Theorem 1. *Assume \hat{Y}_B are the black-box model $B(\cdot)$ predictions, \hat{Y}_I are the predictions of a surrogate (interpretable) model $I(\cdot)$ trained on $B(\cdot)$ outputs and A is a sensitive attribute: 1) If $B(\cdot)$ does not satisfy demographic parity, completely eliminating fairwashing is impossible (i.e., $\gamma > 0$); 2) A detector for fairwashing measuring false-positive rates and true-positive rates of $I(\cdot)$ with respect to $B(\cdot)$ is sufficient.*

We provide a proof sketch here (see Appendix D for full proof).

Proof Sketch. Define T_a^+ , the true-positive rate of $I(\cdot)$ with respect to $B(\cdot)$ on $X \sim \mathcal{D}_a$, in which \mathcal{D}_a denotes the distribution over data with attribute $a \in A$. Define F_a^+ , the false-positive rate of $I(\cdot)$ with respect to $B(\cdot)$ on $X \sim \mathcal{D}_a$. Denote the differences between T_0^+ and T_1^+ , and F_0^+ and F_1^+ as $\tilde{\delta}$, δ' , respectively. We can express Γ_I in terms of $\tilde{\delta}$, δ' , T_1^+ , F_1^+ , Γ_B and $Y_B|A$ by expanding terms and using Bayes' formula:

$$\Gamma_I = \Gamma_B (T_1^+ - F_1^+) + \left(\tilde{\delta} - \delta' \right) \Pr \left[\hat{Y}_B = 1 \mid A = 0 \right] + \delta' \quad (4)$$

Now that we have derived Γ_I in terms of our desired terms, we can eliminate Γ_I from the equation for γ . The resulting formula for γ in terms of $\tilde{\delta}, \delta', T_1^+, F_1^+, \Gamma_B$ and $Y_B|A$ allows us to demonstrate **correctness, irreducibility** and **sufficiency**.

$$\gamma = \Gamma_B(1 + F_1^+ - T_1^+) - \delta' \Pr[\hat{Y}_B = 0 | A = 0] - \tilde{\delta} \Pr[\hat{Y}_B = 1 | A = 0]. \quad (5)$$

Correctness. Note that if $\Gamma_B = 0$ (i.e., black-box model is fair) then no-fairwashing is present, and no detector should mistakenly report fairwashing. We can verify this with Equation (15). Under Assumption 1, we know that $\gamma \geq 0$ and as below terms are, by definition non-negative, thus:

$$\gamma = - \left(\delta' \Pr[\hat{Y}_B = 0 | A = 0] + \tilde{\delta} \Pr[\hat{Y}_B = 1 | A = 0] \right) \geq 0 \implies \tilde{\delta} = \delta' = 0.$$

Irreducibility. If $\Gamma_B \neq 0$, even if $\tilde{\delta} = \delta' = 0$, i.e. when true-positive and false-positive rates are equal over subpopulations ($T_0^+ = T_1^+$ and $F_0^+ = F_1^+$), there exists an irreducible fairwashing violation:

$$\gamma = \Gamma_B(1 + F_1^+ - T_1^+) > 0, \quad (6)$$

since $1 + F_1^+ - T_1^+ > 0$ except in trivial cases (see Appendix D for the proof). We further point out that $1 + F_1^+ - T_1^+$ is not a function of disparity in the interpretable model as the value does not measure disparity between subpopulations. Rather, $1 + F_1^+ - T_1^+$ is a function of the difference in true-positive and false-positive rates (T_1^+ and F_1^+ , respectively) within a single subpopulation. The sole measure of disparity in the gap above is black-box Γ_B , which is a given constant.

Sufficiency. The remaining factors of the gap are only function of $\tilde{\delta}$ and δ' as $\Pr[\hat{Y}_B = 1 | A = 0]$ and $\Pr[\hat{Y}_B = 0 | A = 0]$ are also given constants. Therefore, a detector that measures $\tilde{\delta}$ and δ' is sufficient to detect fairwashing. \square

4 Detecting Fairwashing

4.1 Fairwashing In Practice

Eliminating the fairness gap of the black-box model would be the ideal outcome as it would also remove the possibility for a dishonest model provider to perform fairwashing. However, in practice absolute fairness is often unrealistic and even regulatory bodies tolerate small violations of fairness laws [18]. In Section 3, we showed that fairness violations result in fairwashing violations unless the interpretable model has a higher fairness gap than the black-box model, which would not be desirable from the point of view of the model provider. Therefore, not all fairwashing violations may have been intended by the model owners. These *involuntary* violations are in fact side-effects of the imperfect fidelity, which is inherent to optimizing the empirical fidelity (Definition 2) of the interpretable model on a training set $X_{tr} \sim \mathcal{D}$,

$$\max \text{EmpiricalFidelity}(I, B; X_{tr}).$$

In contrast, deliberate fairwashing is the result of a similar optimization problem purposefully constrained to provide better fairness on a target suing set $X_{sg} \sim \mathcal{D}$:

Definition 4 (Fairwashing Optimization [2]). *Given a black-box model $B(\cdot)$ and a suing set X_{sg} , fairwashing optimization is defined as learning an interpretable model $I(\cdot)$ from $B(\cdot)$ on X_{sg} such that the interpretable model has 1) high fidelity with respect to the black-box model and 2) is less unfair than this black-box model:*

$$\begin{aligned} \max \quad & \text{EmpiricalFidelity}(I, B; X_{sg}) \\ \text{subject to} \quad & \text{FairnessGap}(I; X_{sg}) \leq \epsilon, \end{aligned} \quad (7)$$

in which ϵ is an upper bound on the fairness gap of the interpretable model on the suing set $\text{FairnessGap}(I; X_{sg})$, and $\text{FairnessGap}(\cdot, \cdot)$ constitutes a measure of fairness gap according to some fairness metric – generically measured by the subpopulation gap in a fairness metric. Note that $\epsilon < \text{FairnessGap}(B; X_{sg})$ must be less than the fairness gap of the black-box model on the suing set $\text{FairnessGap}(B; X_{sg})$ to cause fairwashing violation $\gamma > 0$ (see Definition 3).

4.2 Proposed method: FRAUD-Detect

As demonstrated in Section 3, a fairwashing detection method that relies on the difference over subpopulations in true-positive and false-positive rates of the interpretable model w.r.t the black-box model ($\tilde{\delta}$ and δ' respectively) is sufficient. Inspired by this theoretical analysis, we propose our fairwashing detector, FRAUD-Detect, which locates fairwashing violations. FRAUD-Detect aims to distinguish between honest and fairwashed interpretable models approximating the unfair black-box model $B(\cdot)$ as described in Section 2.2 under the additional following constraints: 1) without white-box access to the black-box and interpretable models as well as; 2) independently of the cooperation of the model provider being audited. These constraints are introduced to conform to realistic scenarios in which entities desire to retain the confidentiality of their black-box model due to intellectual property rights and because they considered it as a company's asset.

Similarly to definitions of true-positive T_a^+ and false-positive F_a^+ rates of $I(\cdot)$ w.r.t $B(\cdot)$ on data with attribute $a \in A = \{0, 1\}$, we define the true-negative rate, T_a^- , and false-negative rate, F_a^- as:

$$T_a^- = \Pr[\hat{Y}_I = 0 \mid \hat{Y}_B = 0, A = a], \quad F_a^- = \Pr[\hat{Y}_I = 0 \mid \hat{Y}_B = 1, A = a] \quad (8)$$

for $a \in A = \{0, 1\}$. Note that $T_a^- = 1 - F_a^+$ and $F_a^- = 1 - T_a^+$.

FRAUD-Detect computes the $T_a^+, F_a^+, T_a^-, F_a^-$ for each subpopulation $a \in A = \{0, 1\}$ using the predictions of the interpretable model and black-box model on the suing set X_{sg} . Let

$$C_0 = [T_0^+, F_0^+, T_0^-, F_0^-], \quad C_1 = [T_1^+, F_1^+, T_1^-, F_1^-], \quad (9)$$

be the (flattened) confusion matrices of the interpretable model w.r.t the black-box model on X_{sg} with attribute 0 and 1, respectively. Then, FRAUD-Detect computes the divergence between C_0 and C_1 using Kullback–Leibler (KL) as:

$$\mathcal{C}_{\text{KL}} = \text{KL}(C_0, C_1). \quad (10)$$

Recall that we demonstrated in Section 3 that the divergence in T_0^+ and T_1^+ ($\tilde{\delta}$) and F_0^+ and F_1^+ (δ') were sufficient for a detection method. We can complement our detection method with additional divergences in T_a^- and F_a^- and measure the dissimilarity via KL divergence. This choice is natural as the KL divergence is commonly used to quantify the divergence over probability distributions— C_a functions as a simplified representation of the probability distribution $Y_I|Y_B, A = a$. For an honest interpretable model, \mathcal{C}_{KL} is generally relatively low ($\tilde{\delta}, \delta' \approx 0$) as the only divergence arises from general error in fidelity optimization and the fairness gap of the black-box model (as shown in Section 3). In other words, the honest interpretable model approximates the black-box model equivalently across subpopulations. However, for a fairwashed interpretable model, \mathcal{C}_{KL} grows significantly, as the interpretable model is explicitly manipulated to optimize the fairness across subpopulations and improve over the black-box model. We quantify the distinction between honest and fairwashed interpretable models via a threshold $\Delta > 0$ on \mathcal{C}_{KL} such that if $\mathcal{C}_{\text{KL}} > \Delta$, the interpretable model is considered fairwashed. Note here that due to the irreducibility result stated in Section 3, Δ must be chosen such that Δ is tightly greater than the irreducible term in order to properly distinguish between accidental fairwashing that occurs as a result of the irreducibility and malicious fairwashing. This procedure is detailed in Algorithm 1.

5 Evading FRAUD-Detect

In this section, we investigate on whether a dishonest entity could evade FRAUD-Detect while performing fairwashing. We assume that this dishonest entity, which we call the *adversary*, is informed about FRAUD-Detect. Namely, the adversary is aware of the fairwashing detection method and desires to fairwash while evading detection by the auditor. To achieve this goal, the informed adversary has the objective of finding a fairwashed interpretable model that satisfies an additional constraint on the Kullback–Leibler (KL) divergence of the confusion matrices among subpopulations, \mathcal{C}_{KL} in the previous section. To probe the robustness of our detector to an informed adversary empirically, we explore the range of fairness gap given a fixed value of fidelity and a fixed value of \mathcal{C}_{KL} via solving the informed adversary optimization problem.

²We performed experiments in two settings: 1) considering all related true-positive, false-positive, true-negative and false-negative rates; 2) considering only independent true-positive and false-positive rates.

Algorithm 1 FRAUD-Detect. $\{\text{predicate}\} = \{(x, a, \hat{y}_B, \hat{y}_I) \in (X_{sg}, A, \hat{Y}_B, \hat{Y}_I) : \text{predicate}\}$

Input: Query access to the interpretable model $I(\cdot)$, suing dataset X_{sg} , Black-box model predictions on the suing set $B(X_{sg})$, sensitive attribute $a \in A$ and a threshold $\Delta > 0$.

Output: $\begin{cases} T & \text{fairwashing is detected} \\ F & \text{fairwashing is not detected} \end{cases}$

<p>1: $\hat{Y}_I \leftarrow I(X_{sg})$ \triangleright Query interpretable model</p> <p>2: for $i \leftarrow 0, 1$ do</p> <p>3: $T_i^+ \leftarrow \frac{ \{\hat{y}_B=1, \hat{y}_I=1, a=i\} }{ \{\hat{y}_B=1, a=i\} }$ \triangleright TPR</p> <p>4: $F_i^+ \leftarrow \frac{ \{\hat{y}_B=0, \hat{y}_I=1, a=i\} }{ \{\hat{y}_B=0, a=i\} }$ \triangleright FPR</p>	<p>5: $F_i^- \leftarrow \frac{ \{\hat{y}_B=1, \hat{y}_I=0, a=i\} }{ \{\hat{y}_B=1, a=i\} }$ \triangleright FNR</p> <p>6: $T_i^- \leftarrow \frac{ \{\hat{y}_B=0, \hat{y}_I=0, a=i\} }{ \{\hat{y}_B=0, a=i\} }$ \triangleright TNR</p> <p>7: $C_0 \leftarrow [T_0^+, F_0^+, T_0^-, F_0^-]$</p> <p>8: $C_1 \leftarrow [T_1^+, F_1^+, T_1^-, F_1^-]$</p> <p>9: $\mathcal{C}_{KL} \leftarrow \text{KL}(C_0, C_1)$ \triangleright Kullback–Leibler</p> <p>10: if $\mathcal{C}_{KL} > \Delta$ then</p> <p>11: return T</p> <p>12: else</p> <p>13: return F</p>
---	---

Definition 5 (Informed Adversary Optimization). *Given a black-box model $B(\cdot)$, a suing set X_{sg} , a sensitive attribute A , a loss threshold v and a fairwashing detection threshold Δ , an informed adversary aims to learn an interpretable $I(\cdot)$ such that*

$$\begin{aligned} & \text{minimize} && \text{FairnessGap}(I; X_{sg}) \\ & \text{subject to} && L(I; X_{sg}) \leq v \quad \text{and} \quad \mathcal{C}_{KL} \leq \Delta, \end{aligned} \tag{11}$$

where \mathcal{C}_{KL} is the KL divergence between the subpopulation confusion matrices of X_{sg} .

The differences between the fairwashing optimization in Definition 4 and the informed adversary optimization of Definition 5 are two fold. First, the order of the maximand and the constraint is reversed. Both formulations are valid since each indicates priorities of the adversary, namely, whether to put a hard limit on the fidelity loss or on the fairness gap of the interpretable model. Second, in Definition 5 we replaced the non-differentiable $\text{EmpiricalFidelity}(I, B; X_{tr})$ with a differentiable measure of loss $L(I; X_{sg})$. In our empirical evaluations, we use the logistic regression loss.

To solve the above optimisation problem, we consider the Rashomon set of high-fidelity interpretable models and compute the range of the fairness gap of interpretable models. In short, given a classification task, the Rashomon Set [41] is defined as the set of almost-equally-accurate models (interpretable models in our case). More precisely, given a model class \mathcal{F} , a loss function $L_{\mathcal{D}}(\cdot)$ over a dataset \mathcal{D} of interest, a reference model f^* (e.g., optimal model) and a performance threshold $\tau \in [0, 1]$, the Rashomon set $R_s(\mathcal{F}, f^*, \tau) = \{f \in \mathcal{F} \mid L_{\mathcal{D}}(f) \leq L_{\mathcal{D}}(f^*) + \tau\}$.

We extend the Fairness in the Rashomon Set (FaiRS) algorithm [19] to compute the range of the fairness gap of interpretable models that can be generated over the set of high-fidelity models satisfying a constraint on the KL divergence. FaiRS exploits the so-called Rashomon Effect [14], which is an empirical phenomenon resulting in multiple models displaying the same performance overall (e.g., w.r.t their global accuracy) but have significant differences in terms of their individual predictions [19]. In our setting, the Rashomon effect implies that several interpretable models can achieve the same high fidelity w.r.t the black-box model while displaying different values of fairness gap. Thus, computing the range of the fairness gap of high-fidelity interpretable models that satisfy a constraint on \mathcal{C}_{KL} quantifies the robustness of FRAUD-Detect to the informed adversary.

To the best of our knowledge, the FaiRS algorithm of [19] is the only work proposing a practical solution to efficiently characterize the range of the fairness gap over the Rashomon set. Alternative approaches would require measuring the range of the fairness gap over approximations of the Rashomon set obtained by either generating models using several hyperparameter values (e.g., seeds [33, 21]) or brute-forcing over a particular class of model (e.g., depth seven decision trees [41]). This is not practical due to the high computational cost and may lead to sub-optimal results.

6 Empirical validation

FRAUD-Detect aims to detect fairwashing by distinguishing between a fairwashed interpretable model and honest interpretable one during an audit. Additionally, FRAUD-Detect seeks to be robust

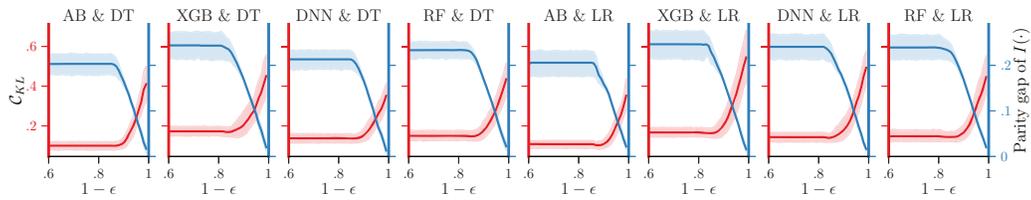


Figure 2: KL divergence between the confusion matrices of subpopulations, C_{KL} , and parity gap as a function of ϵ fairwashing in interpretable models (Logistic Regression (LR) and Decision Tree (DT)) explaining black-box models (AB, Deep Neural Network (DNN), Random Forest (RF) and Gradient Boosted Decision Trees (XgBoost)) using COMPAS. Values before $1 - \epsilon = .6$ are constant and not shown. See Appendix F for the results of other two datasets.

even against an informed adversary, *i.e.*, demonstrate the ability to detect fairwashing even when a dishonest entity is aware of FRAUD-Detect and attempts evasion. Therefore, we empirically validate the performance of FRAUD-Detect in: 1) capturing a relationship between C_{KL} and demographic parity gap of the interpretable model $I(\cdot)$: C_{KL} increases as the demographic parity gap of the interpretable model decreases (*i.e.* as fairwashing becomes more severe); 2) restricting the fairwashing capabilities of the informed adversary: incorporation of C_{KL} constraint into fairwashing problem prevents favorable demographic parity gaps for the interpretable model, precluding fairwashing.

We assess the performance of FRAUD-Detect using a diverse set of black-box architectures, interpretable models and datasets. More precisely, we consider four architectures of black-box models: Deep Neural Networks (DNN), AdaBoost (AB) [24], Gradient Boosted Decision Trees (XGBs) [17] and Random Forests (RFs) [13]. We evaluate the approach on three real-world datasets corresponding to critical decision systems: Adult Income [22], Bank Marketing [34] and COMPAS [5]. We rely on Logistic Regression (LR) and Decision Trees (DTs) as our interpretable models due to their high performance and ease of explainability. We refer to Appendix E for details on these datasets, black-box models, interpretable models and implementation of FRAUD-Detect.

FRAUD-Detect successfully detects fairwashing. Figure 2 illustrates both the demographic parity gap and C_{KL} of interpretable models as a function of ϵ . The bigger the value of ϵ , the stronger the fairwashing (see Definition 4). Any graph within the figures may be read from left to right as the degree of fairwashing increases. As fairwashing becomes more severe, the demographic parity gap decreases while C_{KL} increases such that we observe a relationship between the demographic parity gap of interpretable models and the associated C_{KL} values. For example, decreasing the demographic parity gap of Decision Trees trained on a COMPAS dataset labelled by a RF black-box model by half, increases C_{KL} by a factor of 4. The amount of change in C_{KL} and in the demographic parity gap of fairwashed interpretable models differ across the dataset and black-box models due to the different original fairness of the black-box models. As C_{KL} is sensitive to changes in fairwashing, FRAUD-Detect can be used to detect fairwashing by not permitting C_{KL} to pass a certain threshold value. The C_{KL} thresholds corresponding to 5% and 50% fairwashing are reported in Appendix G. Per-seed results are provided in Appendix H.

FRAUD-Detect is robust to the informed adversary. Figure 3 shows the range of the demographic parity gap of high-fidelity fairwashed interpretable models subjected to a constraint on the C_{KL} by solving the optimization problem of Equation (11). We use four types of black-box models (AB, XGB, DNN and RF). However, for simplicity, we performed the experiment using logistic regression as the interpretable model. The fairwashing detection thresholds were chosen from a wide range, *i.e.* $\Delta \in \{0.01, 0.03, 0.05, 0.07, 0.10, 0.20\}$ (the constraint on C_{KL} in Equation (11)). A direct implication for FRAUD-Detect is that using low fairwashing detection thresholds makes fairwashing difficult. Consistently over all these results, imposing a constraint on the C_{KL} significantly narrows the range of demographic parity gap of fairwashed interpretable models. For instance, on COMPAS with $\Delta = 0.2$, for an AB model that has demographic parity gap of 0.214, fairwashing can produce a completely fair (*i.e.*, demographic parity gap of 0.0) high-fidelity interpretable model while remaining undetected (*i.e.*, $C_{KL} < 0.2$). However, when C_{KL} decreases to 0.1, the fairest high-fidelity interpretable model obtained by the adversary exhibits a demographic parity gap of 0.107.

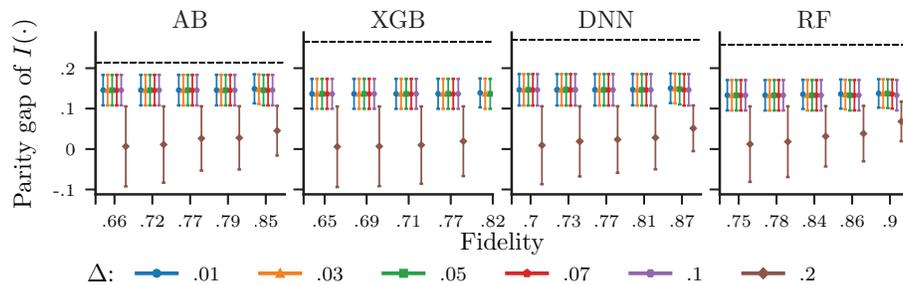


Figure 3: Range of the demographic parity gap of fairwashed logistic regression interpretable models subjected to a constraint (Δ) on the \mathcal{C}_{KL} ($\Delta \in \{0.01, 0.03, 0.05, 0.07, 0.10, 0.20\}$), explaining AB, Deep Neural Network (DNN), RF (RF) and Gradient Boosted Decision Trees (XgBoost) black-box models trained on using COMPAS. Horizontal lines denote the parity gap of the black-box models. See Appendix F for the results of other two datasets.

7 Conclusion and Future Work

Incorporating fairness and explainability in decision systems can help foster trust [29] in them by characterizing failure modes (*e.g.*, unfairness) and providing assurances (*e.g.*, fairness constraints) [16]. Auditing provides an oversight component and helps avoid *first-order failures* such as unfairness. In this sense, fairwashing can be seen as failure of auditing—a *second-order failure* that reduces trust in both fairness and explainability. Second-order failures help characterize the limits of techniques designed to reduce the risks of first-order failures. Our work is the first to delineate the theoretical limits of what is possible in auditing fairness and in containing the risk of fairwashing. In practice, we demonstrated the possibility of performing fairwashing through joint optimization of fidelity and fairness constraints when providing explanations. This could thwart auditing processes that are solely based on an analysis of the delivered interpretable model. We address this issue by proposing an additional auditing protocol that queries the interpretable model. Finally, through both a theoretical framework and an experimental evaluation, we demonstrate that our fairwashing detector cannot be evaded by an attacker, informed of our detector, without significantly degrading the fidelity of model explanations. Future directions include the extension of the detector to other fairness measures (*e.g.*, measures incorporating ground-truth information) as well as extending the analysis to multi-class setups. However, these would both require additional assumptions on the data manifold (to account for class correlations) to establish similar results as in Section 3.

Acknowledgments

This work was supported by CIFAR (through a Canada CIFAR AI Chair and a Catalyst grant), by NSERC (under the Discovery Program, and COHESA strategic research network), by the Ontario Early Researcher Award, and by gifts from Intel and Meta. We also thank the Vector Institute’s sponsors. Ali Shahin Shamsabadi was also partially supported by The Alan Turing Institute. Sierra Wyllie was also partially supported by Engineering Science at University of Toronto through ESROP grant. Ulrich Aïvodji was supported by the NSERC Discovery Grants program (2022-04006). Sébastien Gambs is supported by the Canada Research Chair program, a Discovery Grant from NSERC as well as the NSERC-RDC DEEL project.

References

- [1] Alekh Agarwal, Alina Beygelzimer, Miroslav Dudík, John Langford, and Hanna Wallach. 2018. A reductions approach to fair classification. In *Proceedings of the 35th International Conference on Machine Learning (ICML)*. Stockholm, Sweden, (July 2018).
- [2] Ulrich Aïvodji, Hiromi Arai, Olivier Fortineau, Sébastien Gambs, Satoshi Hara, and Alain Tapp. 2019. Fairwashing: the risk of rationalization. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*. Long Beach, California, USA, (June 2019).

- [3] Ulrich Aïvodji, Hiromi Arai, Sébastien Gambis, and Satoshi Hara. 2021. Characterizing the risk of fairwashing. In *Proceedings of the 35th Conference on Neural Information Processing Systems (NeurIPS)*. Sydney, Australia, (December 2021).
- [4] Christopher Anders, Plamen Pasliev, Ann-Kathrin Dombrowski, Klaus-Robert Müller, and Pan Kessel. 2020. Fairwashing explanations with off-manifold detergent. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*. Online, (July 2020).
- [5] Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. Machine bias. *ProPublica*, (May 2016). Retrieved January 15, 2022.
- [6] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. 2020. Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115.
- [7] Solon Barocas and Andrew D. Selbst. 2016. Big data’s disparate impact. *California Law Review*, 104, 671.
- [8] Adrien Bibal, Michael Lognoul, Alexandre De Streeel, and Benoît Frènay. 2021. Legal requirements on explainability in machine learning. *Artificial Intelligence and Law*, 29, 2, 149–169.
- [9] Sarah Bird, Miro Dudík, Richard Edgar, Brandon Horn, Roman Lutz, Vanessa Milan, Mehrnoosh Sameki, Hanna Wallach, and Kathleen Walker. 2020. Fairlearn: A toolkit for assessing and improving fairness in AI. Technical report MSR-TR-2020-32. Microsoft, (May 2020). <https://www.microsoft.com/en-us/research/publication/fairlearn-a-toolkit-for-assessing-and-improving-fairness-in-ai/>.
- [10] Noorhannah Boodhun and Manoj Jayabalan. 2018. Risk prediction in life insurance industry using supervised learning algorithms. *Complex & Intelligent Systems*, 4, 2, 145–154.
- [11] Olcay Boz. 2002. Extracting decision trees from trained neural networks. In *Proceedings of the 8th International Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*. Edmonton, Canada, (July 2002).
- [12] Leo Breiman. 2017. *Classification and regression trees*. Routledge.
- [13] Leo Breiman. 2001. Random forests. *Machine Learning*, 45, 1, 5–32.
- [14] Leo Breiman. 2001. Statistical modeling: the two cultures (with comments and a rejoinder by the author). *Statistical Science*, 16, 3, 199–231.
- [15] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building classifiers with independency constraints. In *2009 IEEE International Conference on Data Mining Workshops (ICDM)*. Miami, Florida, USA, (December 2009).
- [16] Varun Chandrasekaran, Hengrui Jia, Anvith Thudi, Adelin Travers, Mohammad Yaghini, and Nicolas Papernot. 2021. Sok: machine learning governance. *arXiv preprint arXiv:2109.10870*.
- [17] Tianqi Chen and Carlos Guestrin. 2016. XgBoost: a scalable tree boosting system. In *Proceedings of the 22nd International Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*. San Francisco, California, USA, (August 2016).
- [18] U.S. Equal Employment Opportunity Commission. 1979. Questions and answers to clarify and provide a common interpretation of the uniform guidelines on employee selection procedures. (March 1979). Retrieved 01/15/2022 from <https://www.eeoc.gov/laws/guidance/questions-and-answers-clarify-and-provide-common-interpretation-uniform-guidelines>.
- [19] Amanda Coston, Ashesh Rambachan, and Alexandra Chouldechova. 2021. Characterizing fairness over the set of good models under selective labels. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*. Online, (July 2021).
- [20] Mark Craven and Jude Shavlik. 1995. Extracting tree-structured representations of trained networks. *Advances in Neural Information Processing Systems (NIPS)*, 8, 24–30.
- [21] Alexander D’Amour, Katherine Heller, Dan Moldovan, Ben Adlam, Babak Alipanahi, Alex Beutel, Christina Chen, Jonathan Deaton, Jacob Eisenstein, Matthew D. Hoffman, et al. 2020. Underspecification presents challenges for credibility in modern machine learning. *arXiv preprint arXiv:2011.03395*.
- [22] Dheeru Dua and Casey Graff. 2017. UCI machine learning repository. <http://archive.ics.uci.edu/ml>.

- [23] Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. 2012. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference (ITCS)*. Cambridge, Massachusetts, USA, (August 2012).
- [24] Yoav Freund and Robert E. Schapire. 1997. A decision-theoretic generalization of on-line learning and an application to boosting. *Journal of Computer and System Sciences*, 55, 1, 119–139. ISSN: 0022-0000. DOI: <https://doi.org/10.1006/jcss.1997.1504>. <https://www.sciencedirect.com/science/article/pii/S002200009791504X>.
- [25] Kazuto Fukuchi, Satoshi Hara, and Takanori Maehara. 2020. Faking fairness via stealthily biased sampling. In *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI) number 01*. Volume 34. New York, New York, USA, (February 2020).
- [26] Bryce Goodman and Seth Flaxman. 2017. European union regulations on algorithmic decision-making and a “right to explanation”. *AI Magazine*, 38, 3, 50–57.
- [27] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM Computing Surveys (CSUR)*, 51, 5, 1–42.
- [28] Hoda Heidari, Michele Loi, Krishna P. Gummadi, and Andreas Krause. 2019. A moral framework for understanding of fair ML through economic models of equality of opportunity. In *Proceedings of the 2019 Conference on Fairness, Accountability, and Transparency (ACM FAT*)*. Atlanta, Georgia, USA, (January 2019).
- [29] Alon Jacovi, Ana Marasović, Tim Miller, and Yoav Goldberg. 2021. Formalizing trust in artificial intelligence: prerequisites, causes and goals of human trust in ai. In *Proceedings of the 2021 Conference on Fairness, Accountability, and Transparency (ACM FAccT)*. Online, (March 2021). ISBN: 9781450383097.
- [30] Jon Kleinberg, Himabindu Lakkaraju, Jure Leskovec, Jens Ludwig, and Sendhil Mullainathan. 2018. Human decisions and machine predictions. *The Quarterly Journal of Economics*, 133, 1, 237–293.
- [31] Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. 2017. Inherent trade-offs in the fair determination of risk scores. In *8th Innovations in Theoretical Computer Science Conference (ITCS 2017) (Leibniz International Proceedings in Informatics (LIPIcs))*. Christos H. Papadimitriou, editor. Volume 67. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany. DOI: [10.4230/LIPIcs.ITCS.2017.43](https://doi.org/10.4230/LIPIcs.ITCS.2017.43). Retrieved 07/16/2019 from.
- [32] Zachary C. Lipton. 2018. The mythos of model interpretability: in machine learning, the concept of interpretability is both important and slippery. *ACM Queue*, 16, 3, 31–57.
- [33] Charles Marx, Flavio Calmon, and Berk Ustun. 2020. Predictive multiplicity in classification. In *Proceedings of the 37th International Conference on Machine Learning (ICML)*. Online, (July 2020).
- [34] Sérgio Moro, Paulo Cortez, and Paulo Rita. 2014. A data-driven approach to predict the success of bank telemarketing. *Decision Support Systems*, 62, 22–31.
- [35] Arvind Narayanan. 2018. Tutorial: 21 fairness definitions and their politics. FAT* 2018, (March 2018). <https://www.youtube.com/watch?v=jIXIuYdnyyk>.
- [36] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. 2011. Scikit-learn: machine learning in Python. *Journal of Machine Learning Research*, 12, 2825–2830.
- [37] Inioluwa Deborah Raji, Andrew Smart, Rebecca N. White, Margaret Mitchell, Timnit Gebru, Ben Hutchinson, Jamila Smith-Loud, Daniel Theron, and Parker Barnes. 2020. Closing the AI accountability gap: defining an end-to-end framework for internal algorithmic auditing. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency (ACM FAT*)*. Association for Computing Machinery, New York, New York, USA, (January 2020). DOI: [10.1145/3351095.3372873](https://doi.org/10.1145/3351095.3372873). Retrieved 01/19/2022 from.
- [38] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2016. "why should i trust you?" explaining the predictions of any classifier. In *Proceedings of the 22nd International Conference on Knowledge Discovery and Data Mining (ACM SIGKDD)*. San Francisco, California, USA, (August 2016).

- [39] Christian Sandvig, Kevin Hamilton, Karrie Karahalios, and Cedric Langbort. 2014. Auditing algorithms: research methods for detecting discrimination on internet platforms. In *Data and Discrimination: Converting Critical Concerns into Productive Inquiry, a preconference at the 64th Annual Meeting of the International Communication Association*. Seattle, Washington, USA, (May 2014).
- [40] Ramprasaath R. Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: visual explanations from deep networks via gradient-based localization. In *Proceedings of the 2017 IEEE International Conference on Computer Vision (ICCV)*. Venice, Italy, (October 2017).
- [41] Lesia Semenova, Cynthia Rudin, and Ronald Parr. 2019. A study in rashomon curves and volumes: a new perspective on generalization and model simplicity in machine learning. *arXiv preprint arXiv:1908.01755*.
- [42] Reza Shokri, Martin Strobel, and Yair Zick. 2021. On the privacy risks of model explanations. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society (AIES)*. Online, (May 2021).
- [43] Naeem Siddiqi. 2012. *Credit risk scorecards: developing and implementing intelligent credit scoring*. Volume 3. John Wiley & Sons.
- [44] Dylan Slack, Sophie Hilgard, Emily Jia, Sameer Singh, and Himabindu Lakkaraju. 2020. Fooling lime and shap: adversarial attacks on post hoc explanation methods. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society (AIES)*. New York, New York, USA, (February 2020).
- [45] Sandra Wachter, Brent Mittelstadt, and Luciano Floridi. 2017. Why a right to explanation of automated decision-making does not exist in the general data protection regulation. *International Data Privacy Law*, 7, 2, 76–99.
- [46] Rebecca Wexler. 2017. When a computer program keeps you in jail. *New York Times*. Retrieved January 15, 2022.