FAIRNESS VIA INDEPENDENCE: A GENERAL REGULAR-IZATION FRAMEWORK FOR MACHINE LEARNING

Anonymous authors

000

001

002003004

010 011

012

013

014

016

017

018

019

021

024

025

026

027

028

029

031

033

035

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Fairness in machine learning has emerged as a central concern, as predictive models frequently inherit or even amplify biases present in training data. Such biases often manifest as unintended correlations between model outcomes and sensitive attributes, leading to systematic disparities across demographic groups. Existing approaches to fair learning largely fall into two directions: incorporating fairness constraints tailored to specific definitions, which limits their generalizability, or reducing the statistical dependence between predictions and sensitive attributes, which is more flexible but highly sensitive to the choice of distance measure. The latter strategy in particular raises the challenge of finding a principled and reliable measure of dependence that can perform consistently across tasks. In this work, we present a general and model-agnostic approach to address this challenge. The method is based on encouraging independence between predictions and sensitive features through an optimization framework that leverages the Cauchy-Schwarz (CS) Divergence as a principled measure of dependence. Prior studies suggest that CS Divergence provides a tighter theoretical bound compared to alternative distance measures used in earlier fairness methods, offering a stronger foundation for fairness-oriented optimization. Our framework, therefore, unifies prior efforts under a simple yet effective principle and highlights the value of carefully chosen statistical measures in fair learning. Through extensive empirical evaluation on four tabular datasets and one image dataset, we show that our approach consistently improves multiple fairness metrics while maintaining competitive accuracy.

1 Introduction

Fairness in machine learning has garnered growing concern, as machine learning (ML) models are playing key roles in many high-stakes decision-making scenarios, such as credit scoring (Petrasic et al., 2017), the job market (Hu & Chen, 2018), healthcare (Grote & Keeling, 2022), and education (Bøyum, 2014; Kizilcec & Lee, 2022). Among the various fairness notions, group fairness is one of the most extensively studied ones as it addresses the prediction disparities across demographic groups, including gender, age, skin color, and region (Mehrabi et al., 2021; Dwork et al., 2012; Barocas et al., 2017). While many group fairness ML algorithms are proposed, they have challenges in their applications, especially the *generalizability*, i.e., their adaptation to different fairness notions, and *robustness*, i.e., the stability of the fairness when they encounter a slight change of model parameter.

Existing group fairness approaches can be intrinsically categorized into two main approaches based on their debiasing objectives: *i)* directly integrate the fairness notion into the training objective, *ii)* minimizing the correlation between predictions and sensitive attributes. Methods belongs to *i)* such as a demographic parity (DP) regularizer, and an equality of opportunity (EO) regularizer. The benefit of this approach is that the model trained by the target fairness objective can perform well on specific fairness notions. For example, the machine learning model trained at demographic parity has a high possibility of achieving good demographic parity in testing. However, such methods limited their generalizability to other fairness notions (shown in Fig. 1). Method *ii)* solves from a more fundamental way that can deal with generalizability, including using information theory, or an adversarial approach to minimize the correlation between the prediction and the sensitive attribute. The most straightforward way is to use a distance measurement that assesses the relationship between the sensitive attribute and the prediction, thus minimizing this distance during the training. This

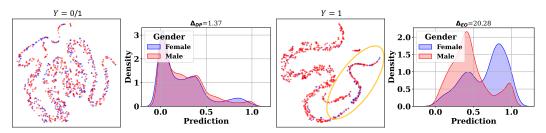


Figure 1: From left to right: (1) Prediction distribution of *all classes*; (2) T-SNE plot of embeddings for samples from *all classes*; (3) Prediction distribution of *class 1*; (4) T-SNE plot of embeddings for samples from Adult, and the sensitive attribute is gender. The blue points represent samples with sensitive attribute 0, while the red points represent samples with sensitive attribute 1.

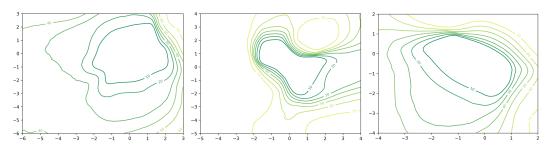


Figure 2: Fairness loss landscapes evaluated using three functions, presented from left to right: Kullback-Leibler (KL) divergence, Hilbert-Schmidt Independence Criterion (HSIC), and Cauchy-Schwarz (CS) divergence. A smaller inner circle indicates greater robustness. Among these methods, the CS divergence achieves the smallest inner circle, ranging from -2 to 1, while the inner circles of KL and HSIC divergences both span from -2 to 2.

enables the generalizability, but ascribes the pressure of the fairness performance to the quality of the distance measurement.

Existing fairness regularizers mainly assessed this correlation using gap parity, MMD, Kullback-Leibler (KL) divergence, and the Hilbert-Schmidt Independence Criterion (HSIC). However, the current fairness regularizers are sensitive to the model parameter change, making them less robust in maintaining fairness, responding to a small change in the model parameters (shown in Fig. 2). Theoretical studies have shown that the Cauchy-Schwarz divergence provides a tighter bound compared to the Kullback-Leibler divergence and gap parity, suggesting its potential to improve fairness in machine learning models. Motivates by this, we would like to know if using CS divergence can result in more generalizable and more robust fairness due to the benefit from a tighter upper bound of CS divergence. In light of this, we propose a new fairness regularizer based on the Cauchy-Schwarz divergence for fair machine learning. To evaluate the generalizability, we tested the fairness under a wide range of fairness notions proposed by previous studies, and to evaluate the robustness, we visualize if a small change of the learned model parameters can influence the fairness. We summarize our contributions as follows:

- We introduce the Cauchy-Schwarz divergence to fair machine learning and present a novel regularization method.
- We elucidate the relationships between the Cauchy-Schwarz regularizer and other fairness regularizers, emphasizing its superior effectiveness in debiasing.
- Our experimental results, obtained from four tabular datasets and one image dataset, validate the
 efficacy of the proposed Cauchy-Schwarz regularizer in achieving fairness across multiple fairness
 notions simultaneously.

2 Preliminaries

In this section, we establish the foundational concepts for our study. We start by exploring the notion of fairness in machine learning, including the relevant notations. Next, we provide an overview

of general fairness-aware machine learning methods. Finally, we introduce the Cauchy-Schwarz divergence and discuss its benefits in reducing bias.

Problem scope. This paper focuses on in-process group fairness in the context of binary classification with binary-sensitive attributes. While group fairness seeks to ensure that machine learning models treat different demographic groups equitably, where groups are defined based on sensitive attributes such as gender, race, and age (Feldman et al., 2015; Zemel et al., 2013).

Notations. Under this setting, we consider a dataset $\mathcal{D} = \{(\mathbf{x}_i, y_i, s_i)\}_{i=1}^M$, where M is the number of samples, $\mathbf{x}_i \in \mathbb{R}^d$ represents the features excluding the sensitive attribute, $y_i \in \{0,1\}$ is the label of the downstream task, and $s_i \in \{0,1\}$ is the sensitive attribute of the i-th sample. The predicted probability for the i-th sample is denoted as $z_i \in [0,1]$, computed by the machine learning model as $z_i = f(\mathbf{x}_i, s_i) : \mathbb{R}^d \to [0,1]$. The binary prediction is represented as $\hat{y}_i \in \{0,1\}$, defined by $\hat{y}_i = \mathbf{1}_{\{\}} \geq t[z_i]$, where $\mathbf{1}_{\{\}} \geq t(\cdot)$ is the indicator function that evaluates whether its input is greater than or equal to the threshold t. Finally, X, Y, S, and \hat{Y} denote random variables corresponding to \mathbf{x}_i, y_i, s_i , and \hat{y}_i , respectively.

Problem Formulation. Generally, the fairness objective can be summarized as follows:

$$\min_{f} \quad \mathcal{L}_{utility} + \lambda \mathcal{L}_{fairness}, \tag{1}$$

where the term $\mathcal{L}_{utility}$ denotes the loss function that measures the utility of the model, often a binary entropy loss, for our binary classification problem, while $\mathcal{L}_{fairness}$ (also shown in Fig. 2) indicates the fairness constraint applied in the model. The parameter λ is used to control the trade-off between utility and fairness.

2.1 GROUP FAIRNESS

There are many ways to define and measure the group fairness. Each definition focuses on distinct statistical measures aimed at achieving balance among subgroups within the data. Among these, Demographic Parity and Equal Opportunity are the most popular ones, other popular fairness definitions are summarized in Appendix L.2.

Demographic Parity (DP). Demographic Parity (Zafar et al., 2017; Feldman et al., 2015; Dwork et al., 2012) mandates that the predicted outcome \hat{Y} be independent of the sensitive attribute S, expressed mathematically as $\hat{Y} \perp S$. Most of the existing literature primarily addresses binary classification and binary attributes, where $Y \in \{0,1\}$ and $S \in \{0,1\}$. Similar to the concept of equal opportunity, the metric evaluating the DP fairness is defined by:

$$\Delta_{DP} = |P(\hat{Y}|S=0) - P(\hat{Y}|S=1)|. \tag{2}$$

A lower value of \triangle_{DP} signifies a fairer classifier.

Equal Opportunity (EO). Equal Opportunity (Hardt et al., 2016) mandates that a classifier achieves equal true positive rates across various subgroups, striving towards the ideal of a perfect classifier. The corresponding fairness measurement for EO can be articulated as follows:

$$\Delta_{EO} = |P(\hat{Y}|Y=1, S=0) - P(\hat{Y}|Y=1, S=1)|. \tag{3}$$

A low \triangle_{EO} indicates that the difference in the probability of an instance in the positive class being assigned a positive outcome is relatively small for both subgroup members. Both DP and EO can be effectively extended to problems involving multi-class classifications and multiple sensitive attribute categories. Note that in the binary classification task, \triangle_{DP} and \triangle_{EO} are sometimes calculated after binarizing $P(\hat{Y})$. Specifically, \triangle_{DP} is defined as $\triangle_{DP} = |P(\hat{Y}=1|S=0) - P(\hat{Y}=1|S=1)|$, and \triangle_{EO} is defined as $\triangle_{EO} = |P(\hat{Y}=1|Y=1,S=0) - P(\hat{Y}=1|Y=1,S=1)|$ (Beutel et al., 2017; Dai & Wang, 2021; Dong et al., 2022).

2.2 CAUCHY-SCHWARZ DIVERGENCE

Motivated by the well-known Cauchy-Schwarz (CS) inequality for square-integrable functions¹, which holds with equality if and only if $p(\mathbf{x})$ and $q(\mathbf{x})$ are linearly dependent, we can define a measure

$${}^{1}\left(\int p(\mathbf{x})q(\mathbf{x})\,dx\right)^{2} \leq \int p(\mathbf{x})^{2}\,dx\int q(\mathbf{x})^{2}\,dx$$

of the distance between $p(\mathbf{x})$ and $q(\mathbf{x})$. This measure is referred to as the CS divergence (Principe et al., 2000; Yu et al., 2023), given by:

$$D_{CS}(p;q) = -\log\left(\frac{\left(\int p(\mathbf{x})q(\mathbf{x})dx\right)^2}{\int p(\mathbf{x})^2dx\int q(\mathbf{x})^2dx}\right). \tag{4}$$

The CS divergence, denoted as D_{CS} , is symmetric for any two probability density functions (PDFs) p and q, satisfying $0 \le D_{\text{CS}} < \infty$. The minimum divergence is achieved if and only if $p(\mathbf{x}) = q(\mathbf{x})$.

3 WHAT MAKES A GOOD FAIRNESS REGULARIZER?

Current fair machine learning algorithms adopt a variety of approaches, prompting us to explore the essential properties that make for effective fairness regularizers. By categorizing these algorithms into three types and conducting preliminary experiments, we aim to identify the key characteristics that contribute to mitigating bias in machine learning models.

Reg.	Fairness Objective ($\mathcal{L}_{fairness}$)
DP	$ \mathbb{E}(\hat{Y} S=0) - \mathbb{E}(\hat{Y} S=1) $
EO	$ \mathbb{E}(\hat{Y} Y=1, S=0) - \mathbb{E}(\hat{Y} Y=1, S=1) $
MMD	$D_{\text{MMD}}(Z_{S=0}, Z_{S=1})$
HSIC	$D_{ ext{HSIC}}(\hat{Y},S)$
	$D_{ extsf{PR}}(\hat{Y},S)$

to mitigating bias in machine learning models. Table 1: Fairness regularizers (Reg.) and objectives.

3.1 BALANCING THE PREDICTION ACROSS DIFFERENT SENSITIVE GROUPS

The first is to directly integrate the fairness notions, such as DP and EO, into the fairness objective.

$$\mathcal{L}_{\text{fairness}} = D(\mathbb{P}, \mathbb{Q}) \quad \text{where}$$

$$\begin{cases}
\mathbb{P} = P(\hat{Y} \mid S = 0), & \mathbb{Q} = P(\hat{Y} \mid S = 1) \text{ for DP,} \\
\mathbb{P} = P(\hat{Y} \mid Y = 1, S = 0), & \mathbb{Q} = P(\hat{Y} \mid Y = 1, S = 1) \text{ for EO.}
\end{cases}$$
(5)

Calculating the distance between \mathbb{P} and \mathbb{Q} has many ways by using difference distance measurement D, and the most used one is to calculate the absolute distance between the mean empirical estimations:

$$\mathcal{L}_{fairness} = |\mathbb{E}(\mathbb{P}) - \mathbb{E}(\mathbb{Q})|, \tag{6}$$

where the expected values are calculated as the mean of summation since \mathbb{P} and \mathbb{Q} are discrete distributions.

Previous fair machine learning studies have shown that the fairness loss of the testing can be upper bounded by the loss of training. Therefore, a distance function having a tighter generalization error bound used in training will lead to a better fairness guarantee for testing.

$$\mathcal{L}_{fairness} = |\mathbb{E}(\hat{Y}|S=0) - \mathbb{E}(\hat{Y}|S=1)|. \tag{7}$$

We can see that the basic idea is to balance the prediction distribution between two sensitive groups. Therefore, for this type of approach, we can also use other distance measurements: Therefore, except for the absolute distance between the two prediction distributions.

3.2 BALANCING THE LATENT REPRESENTATION ACROSS DIFFERENT SENSITIVE GROUPS

Distance measures and minimization:

$$\mathcal{L}_{fairness} = D(Z_{S=0}, Z_{S=1}), \tag{8}$$

where Z is the latent representation from the neural networks, and $Z_{S=0}$ and $Z_{S=1}$ are the representation when the sensitive attribute is 1 or 0. The distance metric $D(\cdot)$ here can be a Mean Maximum Discrepancy (Louizos et al., 2016).

3.3 MINIMIZING THE RELATIONSHIP BETWEEN PREDICTIONS AND SENSITIVE ATTRIBUTES

The goal of this category is to ensure that a fair machine learning algorithm's predictions retain minimal sensitive information.

$$\mathcal{L}_{fairness} = D(\hat{Y}, S), \tag{9}$$

The HSIC and PR in Table 1 belong to this category. Specifically, the $D_{PR}(\hat{Y}, S)$ is defined as:

$$D_{PR}(\hat{Y}, S) = \sum_{\hat{y} \in \hat{Y}} \sum_{s \in S} p(\hat{y}, s) \log \left(\frac{p(\hat{y}, s)}{p(\hat{y})p(s)} \right). \tag{10}$$

Note that, adversarial debiasing methods (Zhang et al., 2018) fall into this category because they employ discriminators to predict sensitive group membership from the learned encoded representations. These methods aim to make sensitive attributes difficult to deduce from the encoded representations. However, our study does not include adversarial methods, as they require training additional discriminator models, thereby adding complexity to the framework. Instead, we focus our comparisons on simple distance-based fairness regularizers.

4 CAUCHY-SCHWARZ FAIRNESS REGULARIZER

In this section, we first introduce three prominent fairness regularizers that assess distribution distance using different metrics: Mean Maximum Discrepancy, Kullback-Leibler divergence, and Hilbert-Schmidt Independence Criterion (HSIC). For each metric, we explore its relationship with CS divergence. Subsequently, we explain how CS divergence can be utilized to achieve fairness.

4.1 HOW CAN THE CAUCHY-SCHWARZ DIVERGENCE BE APPLIED TO MITIGATE BIAS?

Given samples $\{\mathbf{x}_i^p\}_{i=1}^m$ and $\{\mathbf{x}_i^q\}_{i=1}^n$ drawn independently and identically distributed (i.i.d.) from $p(\mathbf{x})$ and $q(\mathbf{x})$ respectively, we can estimate the empirical CS divergence. This estimation can be performed using the kernel density estimator (KDE) as described in (Parzen, 1962) and follows the empirical estimator formula in (Jenssen et al., 2006).

Proposition 4.1. Given two sets of observations $\{\mathbf{x}_i^p\}_{i=1}^{N_1}$ and $\{\mathbf{x}_j^q\}_{j=1}^{N_2}$, let p and q denote the distributions of two groups. The empirical estimator of the CS divergence $D_{CS}(p;q)$ is then given by:

$$\tilde{D}_{CS}(p;q) = \log\left(\frac{1}{N_1^2} \sum_{i,j=1}^{N_1} \kappa(\mathbf{x}_i^p, \mathbf{x}_j^p)\right) + \log\left(\frac{1}{N_2^2} \sum_{i,j=1}^{N_2} \kappa(\mathbf{x}_i^q, \mathbf{x}_j^q)\right) - 2\log\left(\frac{1}{N_1 N_2} \sum_{i=1}^{N_2} \sum_{j=1}^{N_2} \kappa(\mathbf{x}_i^p, \mathbf{x}_j^q)\right).$$
(11)

The proof of this proposition is detailed in Appendix B.1. where κ represents a kernel function, such as the Gaussian kernel defined as $\kappa_{\sigma}(x, x') = \exp(-\|x - x'\|_2^2/2\sigma^2)$. In the following sections, we will explore the relationship between this kernel function and the existing fairness regularizer.

As mentioned earlier, the goal of fairness is to ensure an equal distribution of predictions across sensitive attributes. To achieve this, fairness-aware algorithms focus on minimizing the dependency of predictions on these sensitive attributes. Therefore, effectively modeling the relationship between the outcome variable Y and the sensitive attribute S becomes crucial. The prediction distribution over the sensitive attribute S is defined as follows:

$$\mathbb{P} = P(\hat{Y} \mid S = 0); \quad \mathbb{Q} = P(\hat{Y} \mid S = 1).$$
 (12)

By substituting the distribution of predictions over the sensitive attribute into Eq. (20), where $p = \mathbb{P}$ and $q = \mathbb{Q}$, we can define the objective we aim to solve as follows:

$$\min_{\theta} \mathcal{L}_{BCE} + \alpha \tilde{D}_{CS} \left(\mathbb{P}, \mathbb{Q} \right) + \frac{\beta}{2} \|\theta\|_{2}^{2}, \tag{13}$$

where \mathcal{L}_{BCE} is the binary cross-entropy loss, which measures the classifier's accuracy. It is defined as:

$$\mathcal{L}_{BCE} = \frac{1}{M} \sum_{i=1}^{M} -Y_i \log \hat{Y}_i, \tag{14}$$

where \hat{Y}_i is the predicted output obtained from the training model parameterized by θ . This model can be a Multi-Layer Perceptron for tabular data or a ResNet for image data. Additionally, $\|\theta\|_2^2$ serves as an L_2 regularizer.

Meth	nods		Uti	lity			Fa	irness	
		ACC (%)	\uparrow	AUC (%)	<u></u>	Δ_{DP} (%)	\downarrow	$ \Delta_{EO}$ (%)	\downarrow
	MLP	85.63±0.34	_	90.82±0.23	_	16.52±0.91	_	8.43±3.20	
Gender t t t d n 1 t	DP MMD HSIC PR CS	81.90 ± 0.68 82.89 ± 0.23 81.81 ± 0.52	-4.36% -3.20% -4.46%	$\begin{array}{c} 86.91\pm0.80 \\ 85.27\pm0.52 \\ \underline{87.25}\pm0.41 \\ \overline{85.38}\pm0.82 \\ \textbf{90.15}\pm0.49 \end{array}$	-6.11% -3.93% -5.99%	2.47 ± 0.52 2.66 ± 0.54 0.71 ± 0.40	85.05% 83.90% 95.70%	$\begin{array}{ c c c }\hline 20.15\pm1.13\\ 17.53\pm1.36\\ 18.47\pm1.22\\ \underline{12.45}\pm2.38\\ \hline{\textbf{2.27}}\pm1.04\\ \end{array}$	-107.95%
A	MLP	$ 84.42\pm0.31 $	_	90.15±0.36	_	13.47±0.83	_	$ 9.25\pm3.86 $	_
Race	DP MMD HSIC PR CS	83.12±0.82 84.98±0.17 82.13±1.16	-1.54% 0.66% -2.71%	$\begin{array}{c} 88.45 \pm 0.32 \\ 88.36 \pm 0.67 \\ \textbf{90.90} \pm 0.19 \\ 87.44 \pm 0.33 \\ \underline{90.26} \pm 0.47 \end{array}$	-1.99% 0.83% -3.01%	2.58 ± 0.75 7.90 ± 0.72 1.53 ± 0.83	80.85% 41.35% 88.64%	$ \begin{array}{c} 2.16 \pm 1.06 \\ 3.33 \pm 0.93 \\ 2.11 \pm 0.18 \\ 0.86 \pm 0.60 \\ \hline{\textbf{0.44}} \pm 0.12 \end{array} $	76.65% 64.00% 77.19% 90.70% 95.24%
	MLP	$ 66.85\pm0.72 $	_	72.10±0.94	_	13.22±3.32	_	11.41±5.83	_
Gender	DP MMD HSIC PR CS	64.82 ± 1.62 63.17 ± 3.46 64.95 ± 0.15	-3.04% -5.50% -2.84%	$\begin{array}{c} 70.64 \!\pm\! 1.05 \\ 70.72 \!\pm\! 0.92 \\ 71.17 \!\pm\! 0.84 \\ \textbf{72.12} \!\pm\! 0.75 \\ \underline{71.53} \!\pm\! 0.61 \end{array}$	-1.91% -1.29% 0.03%	$ \begin{vmatrix} 3.09 \pm 0.92 \\ 1.84 \pm 0.43 \\ 3.85 \pm 0.60 \end{vmatrix} $	76.63% 86.08% 70.88%		40.58% 72.39% 77.21% 65.73% 96.14%
8	MLP	$ 66.99\pm1.05 $	_	$ 72.46\pm0.88 $		17.24±4.15		19.44±4.63	
Race	DP MMD HSIC PR CS	64.41±2.04 64.52±2.20 67.22 ±0.90	-3.85% -3.69% 0.34%	72.16 ± 0.94	0.50% 0.41% -0.55%	$ \begin{array}{c} 4.42 \pm 2.11 \\ 2.21 \pm 0.68 \\ \hline 5.60 \pm 1.12 \end{array} $	74.36% 87.18% 67.52%	$\begin{array}{c} 7.04{\pm}2.13\\ 5.60{\pm}1.25\\ \underline{2.72}{\pm}0.87\\ 6.52{\pm}1.30\\ \textbf{1.48}{\pm}1.64 \end{array}$	63.79% 71.19% 86.01% 66.46% 92.39%
	MLP	$ 82.04\pm0.27$	_	$ 90.16\pm0.18$	_	$ 10.26\pm4.68 $		$ 2.13\pm3.64 $	
Gender	DP MMD HSIC PR CS	80.93 ± 0.55 81.40 ± 0.12 80.03 ± 0.30	-1.35% -0.78% -2.45%	$\begin{array}{c} 89.33 \pm 0.15 \\ \overline{88.44} \pm 1.71 \\ \textbf{89.53} \pm 0.10 \\ 88.10 \pm 0.26 \\ 89.15 \pm 0.60 \end{array}$	-1.91% -0.70% -2.28%	2.45 ± 0.65 1.54 ± 0.18 0.35 ± 0.20	76.12% 84.99% 96.59%	$\begin{array}{c} 5.37{\pm}0.32 \\ 4.91{\pm}1.48 \\ 4.95{\pm}0.39 \\ 4.54{\pm}0.41 \\ \hline{\textbf{0.90}}{\pm}0.46 \end{array}$	-152.11% -130.52% -132.39% -113.15% 57.75%
4	MLP	$ 81.23\pm0.14 $	_	$ 90.16\pm0.18 $	_	$ 10.06\pm1.84 $	_	7.42±0.66	
Race	DP MMD HSIC PR CS	80.22 ± 1.22 81.41 ± 0.15 80.27 ± 0.26	-1.24% 0.22% -1.18%	$\begin{array}{ } 89.45 \pm 0.11 \\ \hline 88.42 \pm 1.63 \\ \textbf{89.67} \pm 0.12 \\ 88.45 \pm 0.21 \\ 89.14 \pm 0.94 \end{array}$	-1.93% -0.54% -1.90%	1.45 ± 0.89 1.04 ± 0.53 0.37 ± 0.30	85.59% 89.66% 96.32%	$ \begin{vmatrix} 4.53 \pm 0.48 \\ 4.01 \pm 0.54 \\ 2.77 \pm 0.35 \\ \hline 4.25 \pm 0.49 \\ \textbf{1.35} \pm 0.64 \end{vmatrix} $	38.95% 45.96% 62.67% 42.72% 81.81%

Table 2: Fairness performance of existing fair models on the tabular datasets, considering race and gender as sensitive attributes. \uparrow indicates accuracy improvement **compared to MLP**, with higher accuracy reflecting better performance, and \downarrow denotes fairness improvement **compared to MLP**, where lower values indicate better fairness. All results are based on 10 runs for each method. The best results for each metric and dataset are highlighted in **bold** text.

4.2 Why is the CS Divergence More Effective for Ensuring Fairness?

The CS Divergence is particularly well-suited for promoting fairness due to several key reasons:

(1) Closed-form solution for the mixture of Gaussians. The CS divergence has several advantageous properties, one of which is that it provides a *closed-form solution for the mixture of Gaussians* (Kampa et al., 2011). This particular property has facilitated its successful application in various tasks, including deep clustering (Trosten et al., 2021), disentangled representation learning (Tran et al., 2022), and point-set registration (Sanchez Giraldo et al., 2017).

(2) CS Divergence has a tighter error bound than the KL divergence.

Proposition 4.2. For any d-variate Gaussian distributions $p \sim \mathcal{N}(\mu_p, \Sigma_p)$ and $q \sim \mathcal{N}(\mu_q, \Sigma_q)$, where Σ_p and Σ_q are positive definite, the following inequality holds:

$$D_{\rm CS}(p;q) \le D_{\rm KL}(p;q)$$
 and $D_{\rm CS}(p;q) \le D_{\rm KL}(q;p)$. (15)

The proof can be found in Appendix B.3. It is important to note that the divergences are being compared under the same model parameter θ .

(3) CS divergence can provide tighter bounds than MMD and DP when the distributions are far apart or when the scale of the embeddings varies significantly. Based on Remark A.1, we

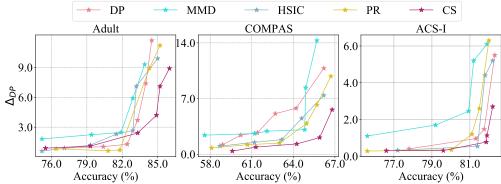


Figure 3: Fairness-accuracy trade-off curves on the test sets for (left) Adult, (middle) COMPAS, and (bottom) ACS-I. Ideally, results should be positioned in the bottom-right corner.

know that CS divergence employs cosine distance, while MMD relies on Euclidean distance. In addition, DP Eq. (5) utilizes a mean disparity, which is a Manhattan distance for the mean estimations of two distributions. CS divergence measures the angle between two distributions in the feature space, focusing on the difference in direction rather than magnitude. In cases where the distributions have significantly different variances or scales, MMD and DP may yield a large distance even if the distributions are aligned in the feature space. In contrast, CS divergence normalizes this comparison, resulting in a more accurate measure of similarity and thereby providing a tighter generalization bound. This normalization enhances the robustness of CS divergence, preventing MMD and DP from overestimating the discrepancy due to their reliance on an unnormalized distance measure.

5 EXPERIMENTS

In this section, we evaluate the effectiveness of the CS fairness regularizer from several perspectives: (1) utility and fairness performance, (2) the tradeoff between utility and fairness, (3) prediction distributions across different sensitive groups, (4) T-SNE plots for these sensitive groups, and (5) the sensitivity of parameters in Eq. (13). Our evaluation encompasses five datasets with diverse sensitive attributes, including four tabular datasets: Adult, COMPAS, ACS-I, and ACS-T, as well as one image dataset, CelebA-A. Utility performance is assessed based on accuracy and the area under the curve (AUC), while fairness performance is measured using \triangle_{DP} Eq. (2) and \triangle_{EO} Eq. (3). Detailed information about the datasets and baselines can be found in the Appendix. We denote an observation drawn from the results as **Obs.**.

5.1 Fairness and Utility Performance

We conducted experiments on five datasets along with their corresponding baselines, as previously mentioned. For each dataset, we performed 10 different splits to ensure robustness in our results. We calculated the mean and standard deviation for each metric across these splits. The accuracy and fairness performance of the downstream tasks is in Table 2. Our observations are as follows:

Obs. 1: CS consistently achieves the best \triangle_{EO} and ranks among the top four for \triangle_{DP} across the Adult, COMPAS, and ACS-I datasets, with only a small margin behind the best results on the remaining datasets. Notably, CS demonstrates exceptional fairness performance on the image dataset, CelebA-A, where the disparity in the 'Young' and 'Non-Young' groups sees a \triangle_{DP} reduction of 97.36% and a \triangle_{EO} reduction of 98.58%. Furthermore, in the Adult and ACS-I datasets, which include gender groups, traditional methods such as DP, MMD, HSIC, and PR do not effectively optimize for EO fairness. In contrast, the proposed CS achieves significant reductions in \triangle_{EO} by 72.12% and 63.85%, respectively, compared to MLP.

Obs. 2: CS achieves good fairness performance with a small sacrifice in utility. Specifically, CS exhibits a decrease of less than 3.1% in accuracy and less than 2.2% in AUC. The only exception is observed with COMPAS when gender is treated as a sensitive attribute, resulting in a slightly higher accuracy loss of 3.6%. Notably, CS demonstrates either equivalent or improved AUC performance,

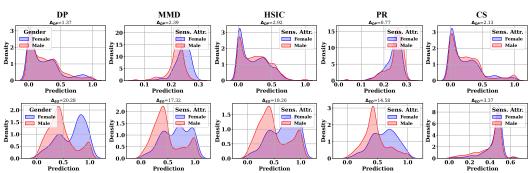


Figure 4: Prediction distributions for female and male groups in the Adult dataset.

with increases of 0.02% and 0.58% on Adult for the gender and race groups, respectively, as well as a 0.35% increase on COMPAS for the race group. Among the baselines, HSIC ranks highest in utility, achieving the best performance on ACS-I for the race group and on ACS-T for both the gender and race groups. This is followed by PR, which shows the best utility on COMPAS for both the gender and race groups, as well as on CelebA-A for the gender group.

5.2 How Do Accuracy and Fairness Trade-Off in Baseline Models and CS?

We evaluate the trade-off between accuracy and \triangle_{DP} for the baselines by varying the fairness hyperparameters (Yao et al., 2023; Deka & Sutherland, 2023). The results are presented in Fig. 3, where the x-axis represents the target accuracy, while the y-axis shows the average Demographic Parity (DP) across both positive and negative target classes. It is important to note that the figure in the bottom right corner represents the optimal result.

Obs. 3: At the same utility level, CS is the most effective method in promoting fairness. Analyzing the results, we find that CS consistently achieves the lowest \triangle_{DP} across most accuracy levels, with this effect becoming more pronounced at higher accuracy levels. This is evidenced by the significant gap in \triangle_{DP} between CS and other baselines. It is important to note that while all baselines can demonstrate good fairness when the optimization prioritizes fairness over task objectives, the task objective remains critical for the practical application of these models.

Obs. 4: High accuracy can sometimes lead to worse fairness compared to MLP, as the fairness objective becomes more challenging to optimize when there is a stronger focus on task-specific objectives. As shown in Table 2, the \triangle_{DP} for MMD is over 14.0, which is greater than the average \triangle_{DP} of 13.22 for MLP. However, these fairness regularizers generally prove effective in controlling bias in representations, especially when more emphasis is placed on the task-specific objective. Notably, some datasets with particular sensitive attributes pose greater challenges for achieving fairness. For instance, the COMPAS dataset, which includes gender as a sensitive attribute, demonstrates this difficulty. One possible explanation is the relatively small sample size of COMPAS, which contains only 6,172 samples, significantly fewer than other datasets where fairness is easier to achieve. For example, the ACS-I dataset has 195,995 samples, approximately 31.7 times that of COMPAS, and features a more balanced gender distribution.

Obs. 5: CS displays a significant increase in \triangle_{DP} at a slower rate than other baselines as accuracy increases. We analyze the slope of the lines representing the increase in \triangle_{DP} with rising accuracy. Many methods, such as PR and DP, demonstrate strong fairness performance at low accuracy levels; however, they quickly lose control over fairness as accuracy begins to increase. This is evident from the abrupt rise in \triangle_{DP} observed at around 82.0% on Adult, 63.0% on COMPAS, and 81.0% on ACS-I. In contrast, CS only exhibits a sudden increase at 85.0%, 65.5%, and 81.5% for the same datasets, respectively.

5.3 HOW CAN THE CS FAIRNESS REGULARIZER PERFORM WELL ON BOTH DP AND EO?

We visualize the kernel density estimate plot 2 of the predictions \hat{Y} across different sensitive groups to analyze how CS achieves a better balance of various fairness definitions compared to other baselines.

https://seaborn.pydata.org/generated/seaborn.kdeplot.html

 The first row displays the predictions for all target classes, specifically Y=0 and Y=1, grouped by sensitive attributes. In this row, the blue areas represent the prediction density for S=0, while the red areas indicate the prediction density for S=1. The second row illustrates the prediction density for the positive target class, Y=1, across two different sensitive groups. Fig. 4 presents the results for Adult based on gender and race groups, with additional results for other datasets available in Appendix C.3.

Obs. 6: CS effectively optimizes the prediction distributions for the two sensitive groups, specifically $\hat{Y}|S=0$ and $\hat{Y}|S=1$. Additionally, it optimizes the prediction distributions for these groups within the positive target group, i.e., $\hat{Y}|S=0, Y=1$ and $\hat{Y}|S=1, Y=1$. Achieving DP and EO fairness requires different objectives. For instance, DP directly optimizes the Δ_{DP} , which results in reduced effectiveness for achieving EO fairness. This is evident across all datasets, as DP ranks among the worst, achieving 7/10 of the lowest EO fairness scores on Δ_{EO} when tested on five datasets with two types of sensitive attributes. The distribution plots for DP further illustrate this, showing a generally larger gap between the two sensitive groups in the EO plots compared to other methods. In contrast, CS consistently minimizes the prediction density gap between the two sensitive groups.

5.4 PARAMETER SENSITIVITY ANALYSIS

For all models, we tune the hyperparameters using cross-validation on the training set. The hyperparameters for these variants are determined through grid search during cross-validation. Specifically, we vary the parameters α and β in Eq. (13) across the ranges (1e-6,150) and (1e-3,10), respectively. In this experiment, we specifically visualize the values of α in the range (1e-4,1e-1) for CS.

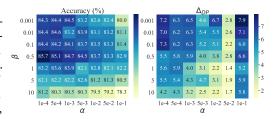


Figure 5: Parameter sensitivity study on Adult.

The heatmap in Fig. 5 illustrates the accuracy and \triangle_{DP} across various combinations of α and β values for the Adult. In the accuracy plots, darker colors indicate higher values, which are preferable, while lighter colors in the \triangle_{DP} plots represent better fairness performance.

Obs. 7: The highest accuracy is achieved when α is set to its smallest value, 1e-4, while the best fairness is obtained with $\alpha=5e-2$. Notably, fairness drops significantly when α increases from 5e-2 to 1e-1. Generally, smaller values of α can still yield satisfactory fairness performance when paired with an appropriate range of β , specifically around 5-10.

Obs. 8: The fairness performance is more sensitive to changes in α than in β . For instance, adjusting β from 1e-3 to 10, which represents a $10,000\times$ increase, results in only a slight decrease in \triangle_{DP} from 7.2 to 4.2. In contrast, increasing α from 1e-2 to 5e-2, a $5\times$ change, leads to a significant drop in \triangle_{DP} from 6.7 to 2.8, when keeping β fixed at 1e-3.

6 CONCLUSION

In this paper, we introduce a novel fair machine learning method called the Cauchy-Schwarz (CS) fairness regularizer. Our approach achieves more robust and generalizable fairness by minimizing the Cauchy-Schwarz divergence between the prediction distribution and the sensitive attributes. We demonstrate that the CS divergence provides a tighter bound compared to both the Kullback-Leibler divergence and the Maximum Mean Discrepancy, as well as the mean disparity used in Demographic Parity regularization. This superiority is particularly evident when the distributions are significantly different or when there is substantial variation in the scale of the embeddings. As a result, our CS fairness regularizer delivers improved fairness performance in practical scenarios. While our work currently only evaluates on general machine learning tasks, and thus leave future work to other tasks such as graph learning.

ETHICS STATEMENT

Our work aims to improve fairness in machine learning by introducing a general and model-agnostic regularization method that reduces statistical dependence between predictions and sensitive attributes. By encouraging independence through Cauchy–Schwarz (CS) Divergence, our framework helps mitigate systematic biases that often lead to disparate treatment of demographic groups. We use publicly available datasets containing sensitive attributes (e.g., gender, race) solely for the purpose of evaluating fairness interventions. We acknowledge that fairness is a multidimensional concept and that no single method can fully eliminate all forms of bias. While our method improves several group fairness metrics, we caution against deploying fairness-enhancing techniques without thorough domain-specific evaluation and stakeholder engagement.

REPRODUCIBILITY STATEMENT

We ensure the reproducibility of our results by providing the following: (1) a comprehensive description of the CS Fairness Regularization method, including the formulation of the objective function and optimization procedure; (2) details on hyperparameters, model architectures, and training procedures used in all experiments; (3) evaluation on five benchmark datasets using multiple fairness metrics and baselines; and (4) open-sourcing our code and data preprocessing scripts upon publication. We follow standard experimental protocols and report average results across multiple runs to account for randomness. These practices facilitate independent verification and encourage future research on robust fairness methods.

REFERENCES

- Tameem Adel, Isabel Valera, Zoubin Ghahramani, and Adrian Weller. One-network adversarial fairness. In *Proceedings of the AAAI*, pp. 2412–2420, 2019.
- Chirag Agarwal, Himabindu Lakkaraju, and Marinka Zitnik. Towards a unified framework for fair and stable graph representation learning. In *UAI 2021: Uncertainty in Artificial Intelligence*, 2021.
- Wael Alghamdi, Hsiang Hsu, Haewon Jeong, Hao Wang, Peter Michalak, Shahab Asoodeh, and Flavio Calmon. Beyond adult and compas: Fair multi-class prediction via information projection. *Advances in Neural Information Processing Systems*, 35:38747–38760, 2022.
- Sina Baharlouei, Maher Nouiehed, Ahmad Beirami, and Meisam Razaviyayn. Rényi fair inference. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=HkgsUJrtDB.
- Solon Barocas, Moritz Hardt, and Arvind Narayanan. Fairness in machine learning. NeurIPS tutorial, 1:2017, 2017.
- Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. Data decisions and theoretical implications when adversarially learning fair representations. *arXiv preprint arXiv:1707.00075*, 2017.
- Steinar Bøyum. Fairness in education—a normative analysis of oecd policy documents. *Journal of Education Policy*, 29(6):856–870, 2014.
- Maarten Buyl and Tijl De Bie. Optimal transport of classifiers to fairness. In *Advances in Neural Information Processing Systems*, 2022.
- Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Optimized pre-processing for discrimination prevention. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017a. URL https://proceedings.neurips.cc/paper_files/paper/2017/file/9a49a25d845a483fae4be7e341368e36-Paper.pdf.
- Flavio Calmon, Dennis Wei, Bhanukiran Vinzamuri, Karthikeyan Natesan Ramamurthy, and Kush R Varshney. Optimized pre-processing for discrimination prevention. *NeurIPS*, 30, 2017b.

- Alexandra Chouldechova. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. *Big data*, 5(2):153–163, 2017.
- 543 Ching-Yao Chuang and Youssef Mroueh. Fair mixup: Fairness via interpolation. In *ICLR*, 2020.
 - Enyan Dai and Suhang Wang. Say no to the discrimination: Learning fair graph neural networks with limited sensitive attribute information. In *Proceedings of the 14th ACM WSDM*, pp. 680–688, 2021.
- Namrata Deka and Danica J Sutherland. Mmd-b-fair: Learning fair representations with statistical testing. In *International Conference on Artificial Intelligence and Statistics*, pp. 9564–9576. PMLR, 2023.
 - Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair machine learning. *NeurIPS*, 2021.
 - Yushun Dong, Ninghao Liu, Brian Jalaian, and Jundong Li. Edits: Modeling and mitigating data bias for graph neural networks. In *Proceedings of the ACM web conference* 2022, pp. 1259–1269, 2022.
 - Dheeru Dua and Casey Graff. UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.
 - Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Conference on Innovations in Theoretical Computer Science (ITCS)*, 2012.
 - Harrison Edwards and Amos Storkey. Censoring representations with an adversary. *arXiv preprint arXiv:1511.05897*, 2015.
 - Michael Feldman, Sorelle A Friedler, John Moeller, Carlos Scheidegger, and Suresh Venkatasubramanian. Certifying and removing disparate impact. In *proceedings of the 21th ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 259–268, 2015.
 - Arthur Gretton, Olivier Bousquet, Alex Smola, and Bernhard Schölkopf. Measuring statistical dependence with hilbert-schmidt norms. In *International conference on algorithmic learning theory*, pp. 63–77. Springer, 2005.
 - Arthur Gretton, Kenji Fukumizu, Choon Teo, Le Song, Bernhard Schölkopf, and Alex Smola. A kernel statistical test of independence. *Advances in neural information processing systems*, 20, 2007.
 - Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13(1):723–773, 2012.
 - Thomas Grote and Geoff Keeling. Enabling fairness in healthcare through machine learning. *Ethics and Information Technology*, 24(3):39, 2022.
 - Xiaotian Han, Zhimeng Jiang, Hongye Jin, Zirui Liu, Na Zou, Qifan Wang, and Xia Hu. Retiring \$\delta \text{DP}\$: New distribution-level metrics for demographic parity. *Transactions on Machine Learning Research*, 2023. ISSN 2835-8856. URL https://openreview.net/forum?id=LjDFIWWVVa.
 - Moritz Hardt, Eric Price, and Nati Srebro. Equality of opportunity in supervised learning. *Advances in neural information processing systems*, 29, 2016.
 - Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, June 2016.
 - Roger A Horn and Charles R Johnson. *Matrix analysis*. Cambridge university press, 2012.
 - Lily Hu and Yiling Chen. A short-term intervention for long-term fairness in the labor market. In *WWW*, pp. 1389–1398, 2018.
 - Robert Jenssen, Jose C Principe, Deniz Erdogmus, and Torbjørn Eltoft. The cauchy–schwarz divergence and parzen windowing: Connections to graph theory and mercer kernels. *Journal of the Franklin Institute*, 343(6):614–629, 2006.

- Ray Jiang, Aldo Pacchiano, Tom Stepleton, Heinrich Jiang, and Silvia Chiappa. Wasserstein fair classification. In *Uncertainty in Artificial Intelligence*, pp. 862–872. PMLR, 2020.
 - Faisal Kamiran and Toon Calders. Data preprocessing techniques for classification without discrimination. *Knowledge and information systems*, 33(1):1–33, 2012.
 - Toshihiro Kamishima, Shotaro Akaho, Hideki Asoh, and Jun Sakuma. Fairness-aware classifier with prejudice remover regularizer. In *Joint European conference on machine learning and knowledge discovery in databases*. Springer, 2012.
 - Kittipat Kampa, Erion Hasanbelliu, and Jose C Principe. Closed-form cauchy-schwarz pdf divergence for mixture of gaussians. In *IJCNN*, pp. 2578–2585. IEEE, 2011.
 - Jian Kang, Jingrui He, Ross Maciejewski, and Hanghang Tong. Inform: Individual fairness on graph mining. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 379–389, 2020.
 - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
 - René F Kizilcec and Hansol Lee. Algorithmic fairness in education. In *The ethics of artificial intelligence in education*, pp. 174–202. Routledge, 2022.
 - Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. Inherent trade-offs in the fair determination of risk scores. *arXiv preprint arXiv:1609.05807*, 2016.
 - Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. *NeurIPS*, 30, 2017.
 - Jeff Larson, Surya Mattu, Lauren Kirchner, and Julia Angwin. Propublica compas analysis—data and analysis for 'machine bias'. https://github.com/propublica/compas-analysis, 2016. Accessed: 2023-03-13.
 - Tai Le Quy, Arjun Roy, Vasileios Iosifidis, Wenbin Zhang, and Eirini Ntoutsi. A survey on datasets for fairness-aware machine learning. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 12(3):e1452, 2022.
 - Peizhao Li, Yifei Wang, Han Zhao, Pengyu Hong, and Hongfu Liu. On dyadic fairness: Exploring and mitigating bias in graph connections. In *International Conference on Learning Representations*, 2020.
 - Zhu Li, Adrian Perez-Suay, Gustau Camps-Valls, and Dino Sejdinovic. Kernel dependence regularizers and gaussian processes with applications to algorithmic fairness. *arXiv* preprint *arXiv*:1911.04322, 2019.
 - Hongyi Ling, Zhimeng Jiang, Youzhi Luo, Shuiwang Ji, and Na Zou. Learning fair graph representations via automated data augmentations. In *ICLR*, 2023.
 - Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
 - Christos Louizos, Kevin Swersky, Yujia Li, Max Welling, and Richard S Zemel. The variational fair autoencoder. In *International Conference on Learning Representations (ICLR)*, 2016.
 - Gilles Louppe, Michael Kagan, and Kyle Cranmer. Learning to pivot with adversarial networks. *NeurIPS*, 30, 2017.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning adversarially fair and transferable representations. *International Conference on Machine Learning*, 2018.
 - Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6):1–35, 2021.
 - Anay Mehrotra and Nisheeth Vishnoi. Fair ranking with noisy protected attributes. *Advances in Neural Information Processing Systems*, 35:31711–31725, 2022.

- Debarghya Mukherjee, Mikhail Yurochkin, Moulinath Banerjee, and Yuekai Sun. Two simple ways to learn individual fairness metrics from data. In *ICML*, pp. 7097–7107. PMLR, 2020.
 - Debarghya Mukherjee, Felix Petersen, Mikhail Yurochkin, and Yuekai Sun. Domain adaptation meets individual fairness. and they get along. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=XSNfXG9HBAu.
 - Emanuel Parzen. On estimation of a probability density function and mode. *The annals of mathematical statistics*, 33(3):1065–1076, 1962.
 - Kevin Petrasic, Benjamin Saul, James Greig, Matthew Bornfreund, and Katherine Lamberth. Algorithms and bias: What lenders need to know. *White & Case*, 2017.
 - JC Principe, Dongxin Xu, J Fisher, and S Haykin. Information theoretic learning. unsupervised adaptive filtering. *Unsupervised Adapt Filter*, 1, 2000.
 - Luis G Sanchez Giraldo, Erion Hasanbelliu, Murali Rao, and Jose C Principe. Group-wise point-set registration based on renyi's second order entropy. In *CVPR*, pp. 6693–6701, 2017.
 - Changjian Shui, Gezheng Xu, Qi Chen, Jiaqi Li, Charles X Ling, Tal Arbel, Boyu Wang, and Christian Gagné. On learning fairness and accuracy on multiple subgroups. *Advances in Neural Information Processing Systems*, 35:34121–34135, 2022.
 - Linh Tran, Maja Pantic, and Marc Peter Deisenroth. Cauchy–schwarz regularized autoencoder. *JMLR*, 23(115):1–37, 2022.
 - Daniel J Trosten, Sigurd Lokse, Robert Jenssen, and Michael Kampffmeyer. Reconsidering representation alignment for multi-view clustering. In *CVPR*, pp. 1255–1265, 2021.
 - Ioannis Tsaousis and Mohammed H Alghamdi. Examining academic performance across gender differently: Measurement invariance and latent mean differences using bias-corrected bootstrap confidence intervals. *Frontiers in Psychology*, 13:896638, 2022.
 - Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
 - Sahil Verma and Julia Rubin. Fairness definitions explained. In 2018 ieee/acm international workshop on software fairness (fairware), pp. 1–7. IEEE, 2018.
 - Yao Yao, Qihang Lin, and Tianbao Yang. Stochastic methods for auc optimization subject to aucbased fairness constraints. In *International Conference on Artificial Intelligence and Statistics*, pp. 10324–10342. PMLR, 2023.
 - Wenzhe Yin, Shujian Yu, Yicong Lin, Jie Liu, Jan-Jakob Sonke, and Efstratios Gavves. Domain adaptation with cauchy-schwarz divergence. *arXiv preprint arXiv:2405.19978*, 2024.
 - Shujian Yu, Hongming Li, Sigurd Løkse, Robert Jenssen, and José C Príncipe. The conditional cauchy-schwarz divergence with applications to time-series data and sequential decision making. *arXiv* preprint arXiv:2301.08970, 2023.
 - Shujian Yu, Xi Yu, Sigurd Løkse, Robert Jenssen, and Jose C Principe. Cauchy-schwarz divergence information bottleneck for regression. *arXiv* preprint arXiv:2404.17951, 2024.
 - Mikhail Yurochkin and Yuekai Sun. Sensei: Sensitive set invariance for enforcing individual fairness. *arXiv preprint arXiv:2006.14168*, 2020.
 - Mikhail Yurochkin, Amanda Bower, and Yuekai Sun. Training individually fair ml models with sensitive subspace robustness. *arXiv preprint arXiv:1907.00020*, 2019.
 - Muhammad Bilal Zafar, Isabel Valera, Manuel Gomez Rogriguez, and Krishna P. Gummadi. Fairness constraints: Mechanisms for fair classification. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2017.

Rich Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning fair representations. In *International conference on machine learning*, pp. 325–333. PMLR, 2013.

Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial learning. In AAAI/ACM Conference on AI, Ethics, and Society (AIES), 2018.

Guanhua Zhang, Yihua Zhang, Yang Zhang, Wenqi Fan, Qing Li, Sijia Liu, and Shiyu Chang. Fairness reprogramming. *arXiv preprint arXiv:2209.10222*, 2022.

Ji Zhao and Deyu Meng. Fastmmd: Ensemble of circular discrepancy for efficient two-sample test. *Neural computation*, 27(6):1345–1372, 2015.

Aoqi Zuo, Susan Wei, Tongliang Liu, Bo Han, Kun Zhang, and Mingming Gong. Counterfactual fairness with partially known causal graph. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022. URL https://openreview.net/forum?id=9albntHz1Uq.

LLM DISCLOSURE

We used an LLM (OpenAI ChatGPT) only for copyediting—streamlining phrasing and correcting grammar, spelling, and stylistic inconsistencies. The model played no role in conceiving the research, designing methods, implementing code, running experiments, analyzing results, or shaping claims. All edits were reviewed by the authors, and only manuscript text was provided to the tool.

A WHAT IS THE RELATIONSHIP BETWEEN CS DIVERGENCE AND EXISTING DISTRIBUTION DISTANCE MEASURES?

To illustrate the advantages of the CS fairness regularizer, we begin by summarizing the commonly used distribution distance metrics: Maximum Mean Discrepancy (MMD), Kullback-Leibler divergence (KL), and Hilbert-Schmidt Independence Criterion (HSIC).

Demographic Parity Regularizer. The demographic parity regularizer is widely utilized in fairness-focused machine learning studies (Chuang & Mroueh, 2020). It aims to optimize the mean disparity between two *prediction distributions*. This regularizer can be formally expressed as:

$$DP(p;q) = \left| \frac{1}{N_1} \sum_{i=1}^{N_1} p(\mathbf{x}_i) - \frac{1}{N_2} \sum_{j=1}^{N_2} q(\mathbf{x}_j) \right|, \tag{16}$$

where \mathbf{x}_i and \mathbf{x}_j are data points from S=0 and S=1, in the context of fairness. In the following, we represent \mathbf{x}_i with distribution p and \mathbf{x}_j with distribution q as \mathbf{x}_i^p and \mathbf{x}_i^q for simplicity. However, only optimizing on the mean disparity of two distributions cannot always generate an optimized DP or EO, as the Eq. (16) equals 0 is a necessary but not sufficient condition for achieving DP and EO.

Mean Maximum Discrepancy. One of the most widely used distance metrics is the Mean Maximum Discrepancy (MMD) (Gretton et al., 2012). In the context of fairness, previous studies have employed MMD as a regularizer to enforce statistical parity among the *embeddings* of different sensitive groups within a machine learning model (Deka & Sutherland, 2023; Louizos et al., 2016). This approach aims to facilitate fair representation learning.

$$\widetilde{\text{MMD}}^{2}(p;q) = \frac{1}{N_{1}^{2}} \sum_{i,j=1}^{N_{1}} \kappa(\mathbf{x}_{i}^{p}, \mathbf{x}_{j}^{p}) + \frac{1}{N_{2}^{2}} \sum_{i,j=1}^{N_{2}} \kappa(\mathbf{x}_{i}^{q}, \mathbf{x}_{j}^{q}) - \frac{2}{N_{1}N_{2}} \sum_{i=1}^{N_{1}} \sum_{j=1}^{N_{2}} \kappa(\mathbf{x}_{i}^{p}, \mathbf{x}_{j}^{q}).$$
(17)

By comparing with Eq. (20), we observe that the CS divergence introduces a logarithmic term for each component of the MMD. Through simple transformations, we can deduce the following:

Remark A.1. CS divergence measures the cosine distance between empirical mean embedding $\mu_p = \frac{1}{N_1} \sum_{i=1}^{N_1} f(\mathbf{x}_i^p)$ and $\mu_q = \frac{1}{N_2} \sum_{j=1}^{N_2} f(\mathbf{x}_j^q)$ in a Reproducing Kernel Hilbert Space, while MMD utilizes Euclidean distance.

Kullback-Leibler Divergence. Kullback-Leibler (KL) Divergence is a key concept in information bottleneck theory, where it is used to quantify the mutual information between two probability distributions. This metric has gained popularity across various domains, including fair machine learning (Kamishima et al., 2012).

$$D_{KL} = \int p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right)$$
 (18)

Hilbert-Schmidt Independence Criterion (HSIC). Let K and L denote the Gram matrices for the variables x and y, respectively. Specifically, K is defined such that $K_{ij} = \kappa(\mathbf{x}_i, \mathbf{x}_j)$, and L is defined as $L_{ij} = \kappa(\mathbf{y}_i, \mathbf{y}_j)$, where κ is the Gaussian kernel function given by $\kappa = \exp\left(-\frac{\|\cdot\|^2}{2\sigma^2}\right)$. The Hilbert-Schmidt Independence Criterion (HSIC) can be estimated using the following expression (Gretton et al., 2007):

$$\widetilde{HSIC}(p;q) = \frac{1}{N^2} \sum_{i,j}^{N} K_{ij} L_{ij} + \frac{1}{N^4} \sum_{i,j,q,r}^{N} K_{ij} L_{qr} - \frac{2}{N^3} \sum_{i,j,q}^{N} K_{ij} L_{iq} = \frac{1}{N^2} \operatorname{tr}(KHLH),$$
(19)

where $H = I - \frac{1}{N} \mathbb{1} \mathbb{1}^T$ represents a centering matrix of size $N \times N$. In this expression, I is the identity matrix, $\mathbb{1}$ is a vector of ones, and $\frac{1}{N} \mathbb{1} \mathbb{1}^T$ computes the average across the columns, effectively centering the data by subtracting the mean from each entry.

Compared to Eq. (17), The HSIC can be interpreted as the MMD between the joint distribution $p(\mathbf{x}, \mathbf{y})$ and the product of their marginal distributions $p(\mathbf{x})p(\mathbf{y})$.

B DETAILS ON THE RELATION OF CS AND EXISTING FAIRNESS REGULARIZERS

B.1 Proof of Proposition 4.1

Proposition 1. Given two sets of observations $\{\mathbf{x}_i^p\}_{i=1}^{N_1}$ and $\{\mathbf{x}q_j\}_{j=1}^{N_2}$, let p and q denote the distributions of two groups. The empirical estimator of the CS divergence $D_{CS}(p;q)$ is given by:

$$\tilde{D}_{CS}(p;q) = \log\left(\frac{1}{N_1^2} \sum_{i,j=1}^{N_1} \kappa(\mathbf{x}_i^p, \mathbf{x}_j^p)\right) + \log\left(\frac{1}{N_2^2} \sum_{i,j=1}^{N_2} \kappa(\mathbf{x}_i^q, \mathbf{x}_j^q)\right) - 2\log\left(\frac{1}{N_1 N_2} \sum_{i=1}^{N_1} \sum_{j=1}^{N_2} \kappa(\mathbf{x}_i^p, \mathbf{x}_j^q)\right).$$
(20)

Proof. The CS divergence is defined as:

$$D_{CS}(p;q) = -\log\left(\frac{(\int p(\mathbf{x})q(\mathbf{x}) d\mathbf{x})^2}{\int p(\mathbf{x})^2 d\mathbf{x} \int q(\mathbf{x})^2 d\mathbf{x}}\right),\tag{21}$$

where $\hat{p}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^{M} \kappa_{\sigma}(\mathbf{x} - \mathbf{x}_{j}^{p})$ and $\hat{q}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^{N} \kappa_{\sigma}(\mathbf{x} - \mathbf{x}_{j}^{q})$ are kernel density estimation.

Then we can obtain:

$$\int \hat{p}^2(\mathbf{x}) \, d\mathbf{x} = \frac{1}{M^2} \sum_{i=1}^M \sum_{j=1}^M \kappa_{\sqrt{2}\sigma}(\mathbf{x}_i^p - \mathbf{x}_j^p). \tag{22}$$

By a similar approach,

$$\int \hat{q}(\mathbf{z})^2 d\mathbf{x} = \frac{1}{N^2} \sum_{i=1}^N \sum_{j=1}^N \kappa_{\sqrt{2}\sigma}(\mathbf{x}_i^q - \mathbf{x}_j^q), \tag{23}$$

and

$$\int \hat{p}(\mathbf{x})\hat{q}(\mathbf{x}) \, d\mathbf{x} = \frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \kappa_{\sqrt{2}\sigma}(\mathbf{x}_{i}^{q} - \mathbf{x}_{j}^{p}). \tag{24}$$

Substituting (Eq. (22))-(Eq. (24)) into Eq. (21), we obtain:

$$\widetilde{D}_{CS}(p;q) = \log \left(\frac{1}{M^2} \sum_{i,j=1}^{M} \kappa_{\sqrt{2}\sigma} (\mathbf{x}_i^p - \mathbf{x}_j^p) \right) + \log \left(\frac{1}{N^2} \sum_{i,j=1}^{N} \kappa_{\sqrt{2}\sigma} (\mathbf{x}_i^q - \mathbf{x}_j^q) \right) - 2 \log \left(\frac{1}{MN} \sum_{i=1}^{M} \sum_{j=1}^{N} \kappa_{\sqrt{2}\sigma} (\mathbf{x}_i^q - \mathbf{x}_j^p) \right).$$
(25)

B.2 Proof of Remark A.1

Remark 1. CS divergence measures the cosine distance between μ_p and μ_q in a Reproducing Kernel Hilbert Space, while MMD utilizes Euclidean distance.

Proof. Let \mathcal{H} be a Reproducing Kernel Hilbert Space (RKHS) associated with a kernel $\kappa(\mathbf{x}_i^p, \mathbf{x}_j^q) = \langle f(\mathbf{x}_i^p), f(\mathbf{x}_j^q) \rangle_{\mathcal{H}}$ (Yu et al., 2024). The mean embeddings of two distributions p and q in \mathcal{H} are denoted by $\mu_p = \frac{1}{N_1} \sum_{i=1}^{N_1} f(\mathbf{x}_i^p)$ and $\mu_q = \frac{1}{N_2} \sum_{j=1}^{N_2} f(\mathbf{x}_j^q)$ in \mathcal{H} , respectively. The CS divergence defined by Eq. (20) can thus be written as:

$$\widetilde{D}_{CS}(p;q) = -2\log\frac{\langle \boldsymbol{\mu}_p, \boldsymbol{\mu}_q \rangle_{\mathcal{H}}}{\|\boldsymbol{\mu}_p\|_{\mathcal{H}} \|\boldsymbol{\mu}_q\|_{\mathcal{H}}} = -2\log D_{COS}(\boldsymbol{\mu}_p, \boldsymbol{\mu}_q)$$

Here, $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ denotes the inner product in the RKHS, and $\| \cdot \|_{\mathcal{H}}$ represents the norm induced by the inner product. The mean embeddings μ_p and μ_q are elements of \mathcal{H} . Thus, the CS divergence is computed based on the cosine distance D_{COS} between μ_p and μ_q .

Similarly, the Maximum Mean Discrepancy (MMD) between distributions p and q defined in Eq. (17) can be written as:

$$\mathrm{MMD}^2(p,q) = \|\boldsymbol{\mu}_p - \boldsymbol{\mu}_q\|_{\mathcal{H}}^2 = D_{\mathrm{EUC}}(\boldsymbol{\mu}_p, \boldsymbol{\mu}_q).$$

Thus, the MMD measures the Euclidean distance between the mean embeddings of p and q in the RKHS \mathcal{H} , i.e., the μ_p and μ_q .

B.3 Proof of Proposition 4.2

Proposition 2. For any d-variate Gaussian distributions $p \sim \mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p)$ and $q \sim \mathcal{N}(\boldsymbol{\mu}_q, \Sigma_q)$ with positive definite Σ_p and Σ_q , the following inequality holds:

$$D_{\rm CS}(p;q) \le D_{\rm KL}(p;q)$$
 and $D_{\rm CS}(p;q) \le D_{\rm KL}(q;p)$. (26)

Proof. The KL divergence for p and q is given by:

$$D_{\mathrm{KL}}(p;q) = \frac{1}{2} \left(\operatorname{tr}(\Sigma_q^{-1} \Sigma_p) - d + (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^{\top} \Sigma_q^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) + \log \left(\frac{|\Sigma_q|}{|\Sigma_p|} \right) \right). \tag{27}$$

The CS divergence is expressed as (Kampa et al., 2011):

$$D_{CS}(p;q) = -\log(d_{xy}) + \frac{1}{2}\log(d_{xx}) + \frac{1}{2}\log(d_{yy}), \tag{28}$$

$$d_{pq} = \frac{\exp\left(-\frac{1}{2}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)^{\top}(\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_q)^{-1}(\boldsymbol{\mu}_p - \boldsymbol{\mu}_q)\right)}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_q|}},$$
(30)

$$d_{pp} = \frac{1}{\sqrt{(2\pi)^d |2\Sigma_p|}}, \quad d_{qq} = \frac{1}{\sqrt{(2\pi)^d |2\Sigma_q|}}.$$
 (31)

We simplify:

$$D_{\text{CS}}(p;q) = \frac{1}{2} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^{\top} (\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_q)^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) + \frac{1}{2} \log \left(\frac{|\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_q|}{2^d \sqrt{|\boldsymbol{\Sigma}_p||\boldsymbol{\Sigma}_q|}} \right).$$
(32)

When the mean vectors differ, based on the property (Horn & Johnson, 2012), $\Sigma_q^{-1} - (\Sigma_p + \Sigma_q)^{-1}$ is positive semi-definite given $\Sigma_p = \Sigma_q$, we have:

$$2(D_{\text{CS}}(p;q) - D_{\text{KL}}(p;q))$$

$$= (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^{\top} (\boldsymbol{\Sigma}_p + \boldsymbol{\Sigma}_q)^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)$$

$$- (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^{\top} \boldsymbol{\Sigma}_q^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p) \le 0.$$
(33)

When the covariance matrices differ, let I be the d-dimensional identity matrix (Yin et al., 2024):

$$2(D_{CS}(p;q) - D_{KL}(p;q)) = \log\left(\frac{|\Sigma_p + \Sigma_q|}{2^d \sqrt{|\Sigma_p||\Sigma_q|}}\right)$$

$$-\log\left(\frac{|\Sigma_q|}{|\Sigma_p|}\right) - \operatorname{tr}(\Sigma_q^{-1}\Sigma_p) + d$$

$$= -d\log 2 + \log\left(|\Sigma_q^{-1}\Sigma_p + I|\right)$$

$$+ \frac{1}{2}\log\left(|\Sigma_q^{-1}\Sigma_p|\right) - \operatorname{tr}(\Sigma_q^{-1}\Sigma_p) + d.$$
(34)

We have $|\Sigma_q^{-1}\Sigma_p| \leq \left(\frac{1}{d}\operatorname{tr}(\Sigma_q^{-1}\Sigma_p)\right)^d$, and $|\Sigma_q^{-1}\Sigma_p + I| \leq \left(1 + \frac{1}{d}\operatorname{tr}(\Sigma_q^{-1}\Sigma_p)\right)^d$. Thus, based on Eq. (34), we can obtain:

$$2(D_{CS}(p;q) - D_{KL}(p;q))$$

$$\leq -d\log 2 + d\log\left(1 + \frac{1}{d}\operatorname{tr}(\Sigma_q^{-1}\Sigma_p)\right)$$

$$+ \frac{d}{2}\log\left(\frac{1}{d}\operatorname{tr}(\Sigma_q^{-1}\Sigma_p)\right) - \operatorname{tr}(\Sigma_q^{-1}\Sigma_p) + d.$$
(35)

The combined Eq. (33) and Eq. (35), we can obtain:

$$2(D_{CS}(p;q) - D_{KL}(p;q)) \le 0, (36)$$

Similarly, we can obtain $2(D_{CS}(q; p) - D_{KL}(q; p)) \le 0$. In conclusion, we conclude:

$$D_{\rm CS}(p;q) \le D_{\rm KL}(p;q)$$
 and $D_{\rm CS}(p;q) \le D_{\rm KL}(q;p)$. (37)

	Methods			Uti	lity			Fair	ness	
			ACC (%)	<u></u>	AUC (%)	<u></u>	Δ_{DP} (%)	\downarrow	Δ_{EO} (%)	\downarrow
		MLP	66.21±0.95	_	73.78±0.25		8.32±2.67	_	5.11±3.55	_
		DP	65.38±0.29		72.40±0.38		0.29 ± 0.15		1.83±0.26	64.19%
	Gender		64.48 ± 0.27		$\frac{72.92}{2} \pm 0.31$		1.22±0.36		2.11±0.49	58.71%
		HSIC			73.16 ±0.32		0.98 ± 0.26		1.00 ± 0.28	80.43%
		PR	62.72±1.01	-5.27%	69.36 ± 0.85		0.78 ± 0.50		1.07 ± 0.36	79.06%
\vdash		CS	$ \underline{65.70}\pm0.42 $	-0.77%	72.83 ± 0.58	-1.29%	0.17 ±0.08	97.96%	0.75 ±0.22	85.32%
ACS-		MLP	66.38±0.42	_	73.69±0.63	_	$9.28{\pm}1.63$		6.21±1.63	_
Ä		DP	$ 64.96\pm0.23 $	-2.14%	71.86±0.23	-2.48%	0.82 ± 0.33	91.16%	1.30 ± 0.26	79.07%
	Race	MMD	65.71 ± 0.65	-1.01%	70.57±0.52	-4.23%	$\overline{3.97} \pm 0.97$	57.22%	1.55±0.79	75.04%
		HSIC	$\overline{65.81} \pm 0.24$	-0.86%	72.92 ±0.23	-1.04%	1.75 ± 0.31	81.14%	0.43 ±0.23	93.08%
		PR	64.25±0.87	-3.21%	70.25±0.30	-4.67%	1.56 ± 0.87	83.19%	1.21±0.74	80.52%
		CS	65.16±0.45	-1.84%	72.56 ± 0.72	-1.41%	0.55 ±0.19	94.07%	1.38 ± 0.46	77.78%
		RN	78.14±0.47	_	86.58±0.55		51.66±0.97	_	35.67±1.11	
		DP	62.42±4.79	-20.12%	66.86±3.19	-22.78%	0.46 ±0.25	99.11%	4.84±2.37	86.43%
	Gender	MMD	62.54±4.26	-19.96%	66.47±3.85	-23.23%	1.39 ± 0.64	97.31%	5.89 ± 3.12	83.49%
Ø		HSIC	63.39±3.63	-18.88%	69.33±3.25	-19.92%	2.24 ± 0.36	95.66%	3.83 ± 2.22	89.26%
Ą		PR	65.51 ±3.52	-16.16%	71.70 ±2.88	-17.19%	4.00 ± 0.52	92.26%	5.05 ± 2.57	85.84%
ep		CS	65.05 ± 3.80	-16.75%	71.42 ± 2.46	-17.51%	0.98 ± 0.62	98.10%	1.53 ±1.05	95.71%
CelebA-A		RN	78.14±0.47	_	86.67±0.53	_	41.74±1.17	_	18.35±1.56	
You		DP	66.78 ±3.61	-14.54%	73.95 ±3.44	-14.68%	2.43±0.83	94.18%	0.91±1.77	95.04%
	Young	MMD	65.82±4.87	-15.77%	72.84±3.61	-15.96%	3.49 ± 0.83	91.64%	1.60 ± 0.71	91.28%
	υ	HSIC			73.08±2.69			95.23%	1.04 ± 0.60	94.33%
		PR	62.98 ± 4.69	-19.40%	69.63±4.02	-19.66%	1.32 ± 0.49	96.84%	1.82 ± 0.53	90.08%
		CS			73.15 ± 3.84			96.93%	0.30 ±0.12	98.37%

Table 3: The fairness performance on the tabular dataset for existing fair models, and we consider race and gender as sensitive attributes. A higher accuracy metric indicates better performance. \uparrow represents the accuracy improvement compared to MLP. A lower fairness metric indicates better fairness. \downarrow represents the improvement of fairness compared to MLP. The results are based on 10 runs for all methods.

C MORE EXPERIMENTAL RESULTS

C.1 EXPERIMENTS ON IMAGE DATASET

In this section, we present the experimental results on the CelebA-A image dataset. The CelebA-A face attributes dataset (Liu et al., 2015) contains over 200, 000 face images, where each image has 40 human-labeled attributes. Among the attributes, we select 'Attractive' as a binary classification task and consider 'Gender' and 'Young' as sensitive attributes. The results are presented in Table 3. The results show a similar finding with the tabular dataset, demonstrating that 1) DP method always achieves a lower Δ_{DP} but a relatively high Δ_{EO} . 2) HSIC is a more promising fair model to achieve equal opportunity.

C.2 IS THE REPRESENTATION LEARNED BY APPLYING CS VIEWED AS FAIR?

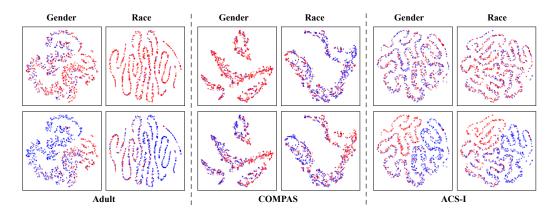


Figure 6: T-SNE visualizations of the latent representations on Adult, COMPAS, and ACS-I, colored by the target attribute (top) and the sensitive attribute (bottom).

To further validate that CS can learn fair representations, we visualize the T-SNE embeddings of the latent space from the last layer before the prediction layer (Van der Maaten & Hinton, 2008)³. Fig. 6 displays the representations learned from the last embedding layer on the Adult, COMPAS, and ACS-I datasets, while Fig. 12 presents the results for ACS-T and CelebA-A. Based on these visualizations, we make the following observations:

Obs. 7: The CS can learn representations that are indistinguishable between sensitive groups. This observation validates the effectiveness of CS in learning fair representations. Specifically, the plots in the first row of Fig. 6 illustrate the embedding visualization of two sensitive groups: blue for S=0 and red for S=1. Overall, the points are uniformly dispersed, with no clear clusters of nodes sharing the same color. This indicates that the embeddings are learned independently of the sensitive attribute. Although some groups have a greater number of data points—such as in the Adult dataset with the sensitive attribute race, where the ratio of S=0:S=1 is 1:9.20, and in the COMPAS dataset with gender, where the ratio is 1:4.17 (as shown in Table 4)—the distribution of points in both colors remains even.

Obs. 8: The CS can learn distinguishable representations for different target attributes. Observing the second row of Fig. 6, we can identify a distinct pattern in the distribution of the blue and red points across different locations in the plot. Among these, the embedding for ACS-I exhibits the clearest pattern, followed by Adult. This observation is consistent with the utility results presented in Table 2, which show a decrease in accuracy and AUC as the degree of negativity increases, particularly evident in the \(\gamma\) columns compared to the MLP. In contrast, COMPAS presents a greater challenge in ensuring utility while considering fairness, as indicated by the less distinct pattern in the learned embeddings, corroborated by the most significant utility drops in Table 2.

 $^{^3} https://scikit-learn.org/stable/modules/generated/sklearn\protect\penalty\z@.manifold.TSNE.html https://scikit-learn.org/stable/modules/generated/sklearn.manifold.TSNE.html$

C.3 More Prediction Distributions over the Sensitive Groups

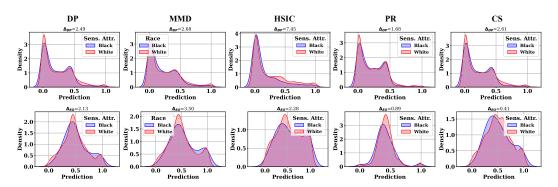
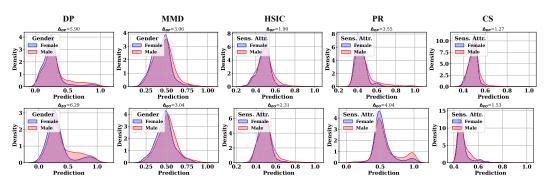
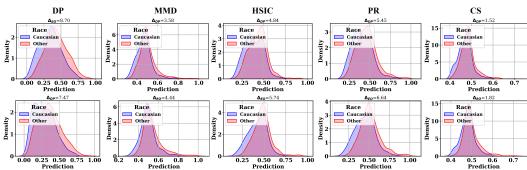


Figure 7: Prediction distributions for black and white groups in the Adult dataset.

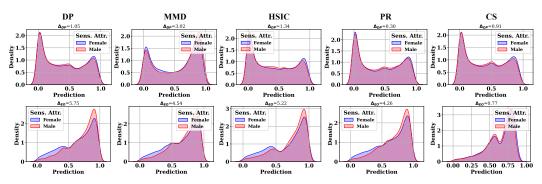


(a) Prediction distributions for female and male groups in the COMPAS dataset.

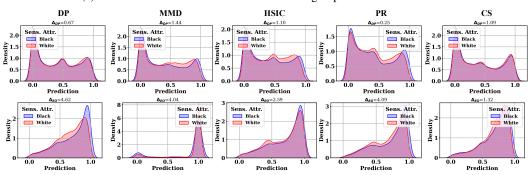


(b) Prediction distributions for Caucasian and (all) other groups in the COMPAS dataset.

Figure 8: Accuracy and \triangle_{DP} trade-off on COMPAS with sensitive attribute gender and race. Results located in the bottom-right corner are preferable.

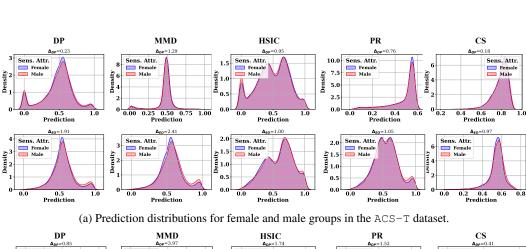


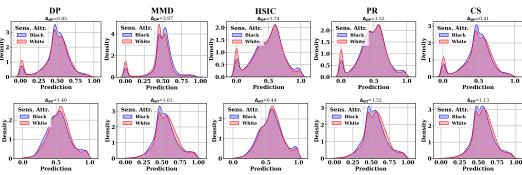
(a) Prediction distributions for female and male groups in the ACS-I dataset.



(b) Prediction distributions for black and white groups in the ACS-I dataset.

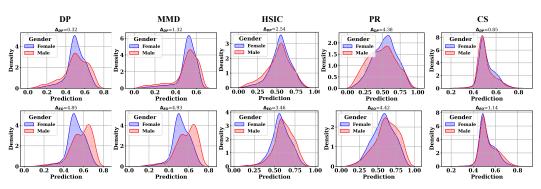
Figure 9: Accuracy and \triangle_{DP} trade-off on ACS-I with sensitive attribute gender and race. Results located in the bottom-right corner are preferable.



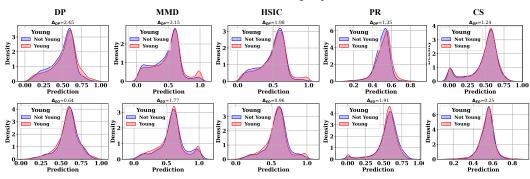


(b) Prediction distributions for black and white groups in the ACS-T dataset.

Figure 10: Accuracy and \triangle_{DP} trade-off on ACS-T with sensitive attribute gender and race. Results located in the bottom-right corner are preferable.



(a) Prediction distributions for female and male groups in the CelebA-A dataset.



(b) Prediction distributions for young and non-yong groups in the CelebA-A dataset.

Figure 11: Accuracy and \triangle_{DP} trade-off on CelebA-A with sensitive attribute gender and race. Results located in the bottom-right corner are preferable.

C.4 MORE T-SNE PLOTS

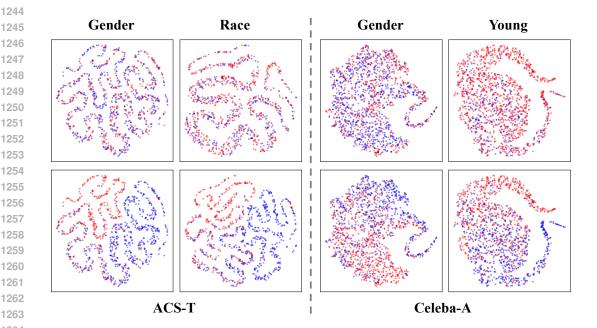


Figure 12: Accuracy and \triangle_{DP} trade-off on ACS-T and CelebA-A. Results located in the bottomright corner are preferable.

In addition to the T-SNE plots shown in Fig. 6, which show the results on three datasets, we also include the T-SNE plots on the two remaining datasets ACS-T and CelebA-A in Fig. 12.

D DATASET DESCRIPTIONS AND DETAILS

We conducted experiments on five datasets, including four tabular datasets and one image data. The introduction of these datasets is as follows:

- Adult (Dua & Graff, 2017) The Adult dataset includes data from 45, 222 individuals based on the 1994 US Census. The primary task is to predict whether an individual's income exceeds \$50k USD, using various personal attributes. In this analysis, we focus on gender and race as sensitive
- COMPAS (Larson et al., 2016) The COMPAS dataset contains records of criminal defendants and is designed to predict the likelihood of recidivism within two years. It encompasses various attributes related to the defendants, including their criminal history, gender, and race.
- ACS-I and ACS-T⁶ (Ding et al., 2021) The ACS dataset is derived from the American Community Survey (ACS) Public Use Microdata Sample and encompasses several prediction tasks. These tasks include predicting whether an individual's income exceeds \$50k and whether an individual is employed, with features such as race, gender, and other relevant characteristics tailored to each
- CelebA-A⁷ (Liu et al., 2015) The CelebFaces Attributes dataset comprises 20,000 face images of 10,000 distinct celebrities. Each image is annotated with 40 binary labels representing various facial attributes, including gender, hair color, and age. In this study, we focus on the 'attractive' label for a binary classification task, while considering 'young' and 'gender' as sensitive attributes.

The detailed statistics for the aforementioned datasets are summarized as follows:

https://archive.ics.uci.edu/ml/datasets/adult

⁵https://github.com/propublica/compas-analysis

⁶https://github.com/zykls/folktables

⁷https://mmlab.ie.cuhk.edu.hk/projects/CelebA.html

Dataset	Task	Sen. Attr. (S)	#Samples	#Feat.	Class <i>Y</i> 0:1	1st S 0:1	$2nd\ S\ 0\!:\!1$
Adult	Income	Gender, Race	45,222	101	1:0.33	1:2.08	1:9.20
COMPAS	Credit	Gender, Race	6,172	405	1:0.83	1:4.17	1:0.52
ACS-I	Income	Gender, Race	195,665	908	1:0.70	1:0.89	1:1.62
ACS-T	Travel Time	Gender, Race	172,508	1,567	1:0.94	1:0.89	1:1.61
CelebA-A	Attractive	Gender, Young	202,599	48×48	1:0.95	1:0.71	1:3.45

Table 4: The table presents the statistics of the datasets. #Feat. refers to the total number of features after preprocessing. The ratio 0:1 represents the proportion between the two categories of the target label or sensitive attributes.

E BASELINES DETAILS

We consider four widely used fairness methods: DP, MMD, HSIC, and PR. Specifically, DP and HSIC minimize the demographic parity and Hilbert-Schmidt Independence Criterion, correspondingly. MMD learns a classifier that optimizes the Mean Maximum Discrepancy. We also include base models MLP and RN for tabular data and image data, correspondingly.

- DP: It is a gap regularization method for demographic parity (Chuang & Mroueh, 2020). As these
 fairness definitions cannot be optimized directly, gap regularization is differentiable and can be
 optimized using gradient descent.
- MMD: The Maximum Mean Discrepancy (MMD) (Gretton et al., 2012) is a metric used to measure
 the distance between probability distributions. Previous research has leveraged MMD to enhance
 fairness in machine learning models, specifically in variational autoencoders (Louizos et al., 2016)
 and MLPs (Deka & Sutherland, 2023). In this paper, we build on the methodologies from earlier
 works (Zhao & Meng, 2015) to compute the MMD baseline.
- HSIC: It minimizes the Hilbert-Schmidt Independence Criterion between the prediction accuracy and the sensitive attributes (Gretton et al., 2005; Baharlouei et al., 2020; Li et al., 2019).
- Prejudice Remover (PR) (Kamishima et al., 2012) (Prejudice Remover) minimizes the prejudice index which is the mutual information between the prediction accuracy and the sensitive attributes.

F MORE FAIRNESS DEFINITIONS

G DETAILS OF THE GROUP FAIRNESS

In this section, we provide the details of the group fairness. We first introduce the definition of group fairness. Then, we introduce the existing group fairness metrics and algorithms.

- DP (Demographic Parity or Statistical Parity) (Zemel et al., 2013). A classifier satisfies demographic parity if the predicted outcome \(\hat{Y}\) is independent of the sensitive attribute S, i.e., \(P(\hat{Y} | S = 0) = P(\hat{Y} | S = 1)\).
- prule (Zafar et al., 2017). A classifier satisfies p%-rule if the ratio between the probability of subjects having a certain sensitive attribute value assigned the positive decision outcome and the probability of subjects not having that value also assigned the positive outcome should be no less than p/100, i.e., $|P(\hat{Y} = 1 \mid S = 1)/P(\hat{Y} = 1 \mid S = 0)| < p/100$.
- prule (Equality of Opportunity) (Hardt et al., 2016). A classifier satisfies equalized opportunity if the predicted outcome Y is independent of the sensitive attribute S when the label Y = 1, i.e., $P(\hat{Y} \mid S = 0, Y = 1) = P(\hat{Y} \mid S = 1, Y = 1)$.
- prule (Equalized Odds) (Hardt et al., 2016). A classifier satisfies equalized odds if the predicted outcome Y is independent of the sensitive attribute S conditioned on the label Y, i.e., $P(\hat{Y} \mid S = 0, Y = y) = P(\hat{Y} \mid S = 1, Y = y), y \in \{0, 1\}.$
- ACC (Accuracy Parity). A classifier satisfies accuracy parity if the error rates of different sensitive attribute values are the same, i.e., $P(\hat{Y} \neq Y \mid S = 0) = P(\hat{Y} \neq Y \mid S = 1), y \in \{0, 1\}.$
- aucp (ROC AUC Parity). A classifier satisfies ROC AUC parity if its area under the receiver operating characteristic curve of w.r.t. different sensitive attribute values are the same.

- ppv (Predictive Parity Value Parity) A classifier satisfies predictive parity value parity if the probability of a subject with a positive predictive value belonging to the positive class w.r.t. different sensitive attribute values are the same, i.e., $P(Y = 1 \mid \hat{Y}, S = 0) = P(Y = 1 \mid \hat{Y}, S = 1)$.
- bnegc (Balance for Negative Class). A classifier satisfies balance for the negative class if the average predicted probability of a subject belonging to the negative class is the same w.r.t. different sensitive attribute values, i.e., $\mathbb{E}[f(X) \mid Y = 0, S = 0] = \mathbb{E}[f(X) \mid Y = 0, S = 1]$.
- bposc (Balance for Positive Class). A classifier satisfies balance for the negative class if the average predicted probability of a subject belonging to the positive class is the same w.r.t. different sensitive attribute values, i.e., $\mathbb{E}[f(X) \mid Y = 1, S = 0] = \mathbb{E}[f(X) \mid Y = 1, S = 1]$.
- abcc (Area Between Cumulative density function Curves) (Han et al., 2023) is proposed to
 precisely measure the violation of demographic parity at the distribution level. The new fairness
 metrics directly measure the difference between the distributions of the prediction probability for
 different demographic groups

H ADDITION EXPERIMENTS ON MORE FAIRNESS METRICS

we provide additional results comparing our framework with baselines under the following fairness notions: Predictive Parity (PPV) (Chouldechova, 2017), p%-Rule (PRULE) (Zafar et al., 2017), Balance for Positive Class (BFP) (Kleinberg et al., 2016), and Balance for Negative Class (BFN) (Kleinberg et al., 2016). The dataset is Adult, using gender as the sensitive attribute. All other experimental settings are consistent with Table 1 in the paper.

Method	$\Delta_{PPV}\left(\downarrow\right)$	PRULE (†)	$\Delta_{BFP}\left(\downarrow\right)$	$\Delta_{BFN} (\downarrow)$
DP	27.35 ± 5.64	81.21 ± 9.04	11.25 ± 2.75	5.15 ± 0.44
MMD	35.19 ± 6.33	85.83 ± 7.15	18.32 ± 3.74	3.49 ± 0.25
HSIC	37.25 ± 3.19	96.18 ± 2.12	16.47 ± 1.21	4.04 ± 0.32
PR	25.46 ± 3.17	89.57 ± 7.39	21.45 ± 2.37	3.46 ± 0.28
CS	31.59 ± 4.35	97.75 ± 3.24	15.25 ± 2.58	3.18 ± 0.36

Table 5: Fairness performance comparison on the Adult dataset, with gender as the sensitive attribute under Δ_{PPV} , PRULE, Δ_{BFP} , and Δ_{BFN} .

We observe the following:

- CS generally achieves the best fairness trade-off performance across the four tested fairness notions. - On the Adult, BFN is generally minimized more effectively than BFP. - Since BFN is related to EO, the ranking of Δ_{BFN} aligns with Δ_{EO} in Table 1 of the paper. Note that, as stated in previous studies (Kleinberg et al., 2016), there is an inherent trade-off between BFP and BFN in practice.

I ADDITIONAL EXPERIMENTS ON COMBINING MULTIPLE REGULARIZER TERMS SIMULTANEOUSLY

We conducted additional experiments where we combined both KL divergence and CS divergence as regularizers. The experiments were performed on the Adult, with gender as the sensitive attribute.

Actually, combining multiple fairness objectives has several drawbacks, which is why most existing studies avoid using multiple regularizers. Instead, they often choose to add simple constraint terms. The key drawbacks of combining fairness regularizers are summarized as follows:

- The CS divergence is upper-bounded by the KL divergence. Therefore, adding KL as an additional fairness objective is theoretically redundant and will not provide further benefits. - Adding KL or

⁸We adopt the preprocessing in previous studies (Le Quy et al., 2022; Mehrabi et al., 2021) involving identifying the target labels and sensitive attributes, and then selecting the relevant features for the analysis.

Method	$\Delta_{DP}(\downarrow)$	$\Delta_{EO} \left(\downarrow \right)$
DP	1.29 ± 0.95	20.15 ± 1.13
CS	2.42 ± 0.85	2.27 ± 1.04
KL	2.77 ± 0.86	10.42 ± 4.34
CS+KL	2.46 ± 1.25	13.42 ± 6.12
$\mathtt{CS+}0.5\mathtt{KL}$	2.25 ± 1.14	9.33 ± 6.36

Table 6: The fairness performance on Adult (gender).

other fairness metrics increases computational complexity, making the optimization process more challenging.

These experimental results further shows the significance of our contribution: identifying a suitable, tighter-bounded fairness regularizer that balances effectiveness and computational efficiency.

J ADDITIONAL EXPERIMENTS ON PRE-PROCESSING AND POST-PROCESSING BASELINES

we have added a post-processing method, PostEO (Hardt et al., 2016) on the Adult dataset (with gender as the sensitive attribute).

Method	ACC (↑)	AUC (↑)	$\triangle_{DP}\left(\downarrow\right)$	$\triangle_{EO} (\downarrow)$
DP	82.42 ± 0.39	86.91 ± 0.80	1.29 ± 0.95	20.15 ± 1.13
CS	83.31 ± 0.47	90.15 ± 0.49	2.42 ± 0.85	2.27 ± 1.04
PR	81.81 ± 0.52	85.38 ± 0.82	0.71 ± 0.40	12.45 ± 2.38
PostPro	80.25 ± 0.83	84.35 ± 0.98	5.75 ± 1.67	2.12 ± 1.44

Table 7: Comparison of methods on various metrics.

the PostPro method is specifically designed to optimize for EO (Hardt et al., 2016), which explains its lower Δ_{EO} .

However, both pre-processing and post-processing methods share a common limitation: they result in lower utility (ACC or AUC). Considering the need for a balanced trade-off between fairness and utility, CS emerges as the most favorable option in our comparison.

K More Experimental Details

In this section, we describe the details of the experimental setup. In this work, we adopted a straightforward stopping strategy. We employ a linear decay strategy for the learning rate, halving it every 50 training step. The model training is stopped when the learning rate decreases to a value below $1e^{-5}$. Across all datasets, we use a weight decay of 0.0, StepLR with a step size of 50 and a gamma value of 0.1, and train for 150 epochs using the Adam Optimizer (Kingma & Ba, 2014). The batch size and learning rate vary depending on the dataset, with specific values provided below. Additionally, Table 8 lists the range of the control hyperparameter β for each fairness approach. The experiments were executed using NVIDIA RTX A4000 GPUs with 16GB GDDR6 Memory.

K.1 Hyperparameter Settings

1. Training Hyperparameters:

- Tabular data (Adult, COMPAS, ACS-I, and ACS-T):
 - Learning rate: $1e^{-2}$ - Weight decay: 0.0

- StepLR_step: 50 - StepLR_gamma: 0.1 - Training epochs: 150 - Batch sizes: 1,024 on Adult, 32 on COMPAS, 4,096 on ACS-I, 4,096 on ACS-T Image data (CelebA-A): - Learning rate: $1e^{-3}$ - Weight decay: 0.0 - StepLR_step: 50 StepLR_gamma: 0.1 - Training epochs: 150

2. Architecture Hyperparameters:

Batch sizes: 256.

• Multilayer perceptron:

- Number of layers: 3
- Number of hidden neurons: {512, 256, 64}
- ResNet-18 (He et al., 2016):
 - Model: https://github.com/pytorch/vision/blob/main/torchvisio n/models/resnet.py

K.2 Hyperparameter Selection

To implement CS and the baseline methods, we adjust the hyperparameter β by tuning it within a specified range. The details of the hyperparameter selection process and the specific range for β are provided below:

Method	Fairness Control Hyperparameter β
DP	0.5, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.5, 3.0, 3.5, 4
HSIC	0.1, 1, 5, 10, 50, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1,000
PR	0.05, 0.2, 0.3, 0.40, 0.50, 0.7, 0.9, 1.0
ADV	0.5, 1.0, 1.2, 1.4, 1.6, 1.8, 2.0, 2.5, 3.0, 3.5
CS	$1e^{-6}, 1e^{-5}, 1e^{-4}, 1e^{-3}, 1e^{-2}, 2e^{-2}, 5e^{-2}, 0.1, 0.5, 1.0, 2.0, 3.0, 4.0, 50, 150$

Table 8: The selections of fairness control hyperparameter blue β .

L RELATED WORK

In this section, we first review relevant prior studies, beginning with an overview of algorithmic fairness in machine learning. We then narrow our focus to regularization-based in-processing methods, which are central to our approach.

L.1 ALGORITHMIC FAIRNESS IN MACHINE LEARNING

The importance of fairness in machine learning has grown significantly as the demand for unbiased decision-making models for individuals and groups increases. This is especially critical in high-stakes applications where the consequences of biased decisions can be severe. Fairness is commonly categorized into three main types: *Individual fairness* (Yurochkin et al., 2019; Mukherjee et al., 2020; Yurochkin & Sun, 2020; Kang et al., 2020; Mukherjee et al., 2022), which aims to ensure that similar individuals are treated similarly; *Group fairness* (Hardt et al., 2016; Verma & Rubin, 2018; Li et al., 2020; Ling et al., 2023), which focuses on achieving fairness across predefined subgroups, often defined by sensitive attributes such as gender or race; *Counterfactual fairness* (Kusner et al., 2017; Agarwal et al., 2021; Zuo et al., 2022), which seeks to ensure fairness by considering how decisions would hold under alternative scenarios. Given the widespread adoption of group fairness

 metrics in real-world applications and the increasing development of in-processing techniques for deep neural network models, we focus on benchmarking these methods to ensure group fairness in neural networks, particularly for tabular and image data.

Various techniques for mitigating bias in machine learning models can be categorized into three main approaches: *pre-processing, in-processing*, and *post-processing*. *Pre-processing* methods focus on addressing biases present in the dataset itself to ensure that the trained model exhibits fairness (Kamiran & Calders, 2012; Calmon et al., 2017a). For instance, these techniques may involve rebalancing the dataset or modifying the data collection process (Calmon et al., 2017b). *In-processing* methods, on the other hand, adjust the training objectives by incorporating fairness constraints directly into the learning process (Kamishima et al., 2012; Zhang et al., 2018; Madras et al., 2018; Zhang et al., 2022; Buyl & De Bie, 2022; Alghamdi et al., 2022; Shui et al., 2022; Mehrotra & Vishnoi, 2022). This approach aims to ensure that the model learns fair representations during training. Finally, *post-processing* methods modify the predictions made by classifiers after the model has been trained, with the goal of promoting fairness across different groups (Hardt et al., 2016; Jiang et al., 2020; Tsaousis & Alghamdi, 2022). By categorizing these techniques, we can better understand the different strategies available for mitigating bias in machine learning systems.

L.2 REGULARIZATION-BASED IN-PROCESSING METHODS

In this paper, we explore three types of regularization-based in-processing methods. First, *Gap Regularization* (Chuang & Mroueh, 2020) streamlines the optimization process by offering a smooth approximation of real-world loss functions, which are typically non-convex and difficult to optimize directly. This category includes methods such as DP, EO, and EOD. Second, the *Independence* approach integrates fairness constraints into the optimization, aiming to mitigate the influence of protected attributes on model predictions while maintaining overall performance. Notable examples of this approach include PR (Kamishima et al., 2012) and HSIC (Li et al., 2019). Lastly, *adversarial debiasing* seeks to minimize utility loss while hindering an adversary's ability to accurately predict the protected attributes. This approach encompasses methods like ADV (Zhang et al., 2018; Louppe et al., 2017; Beutel et al., 2017; Edwards & Storkey, 2015; Adel et al., 2019) and LAFTR (Madras et al., 2018).