

Minority Collective Action for User-Side Fairness

Anonymous authors

Paper under double-blind review

Abstract

Machine learning models often preserve biases present in training data, leading to unfair treatment of certain minority groups. Despite an array of existing firm-side bias-mitigation techniques, these methods typically incur utility costs and require organizational buy-in. Recognizing that many models rely on user-contributed data, end-users can induce fairness through the framework of Algorithmic Collective Action. In this setting, a coordinated minority group strategically relabels its own data to enhance fairness without altering the firm’s training process. We design practical, model-agnostic, minority-only collective-action methods and validate them on real-world datasets. Our findings show that a subgroup of the minority can substantially reduce unfairness with little impact on overall prediction error.

1 Introduction

As machine learning (ML) tools become increasingly accessible, more firms deploy them for decision-making. However, ML models often perpetuate societal biases present in their training data, leading to unfair outcomes across demographic groups (Barocas & Selbst, 2016). Moreover, most fairness-preserving learning algorithms incur a non-negligible cost in accuracy or computational resources (Menon & Williamson, 2018; Zhao & Gordon, 2019; Dehdashtian et al., 2024; Sadeghi et al., 2022), which can discourage practical adoption. Since firms control the ML pipeline, end-users lack access to these algorithms and cannot directly enforce fair treatment. Yet, affected users routinely generate and share data, through clicks, ratings, or other contributions, that is used to train the firm’s models. Such cases, where users collaborate to influence what firms learn, are not uncommon and are well-documented (Sigg et al., 2025). Consequently, this paper asks whether underrepresented minority groups can collaboratively alter the data they share to steer learned models toward fairer behavior, even without access to the firm’s training pipeline.

For example, consider a human resources company that fills vacancies and trains ML models on resumes to predict candidates’ skills. While the majority may have more formal education and college degrees, disadvantaged groups may have informal training or internships. As a result, the ML model may fail to assign the correct skills to minority candidates because they lack formal education, despite their practical experience. Minority members can react to this injustice by collectively submitting their resumes while reframing their reported skills, such as sales or management. Appendix A describes other examples where such *collective action* is applicable. This idea is reminiscent of *pre-processing* fairness techniques (Kamiran & Calders, 2009; Luong et al., 2011; Zemel et al., 2013; Madras et al., 2018), which modify the data before model training. Unlike these prior approaches, which assume centralized control over the data, we consider the setting of *algorithmic collective action* (ACA) (Hardt et al., 2023; Ben-Dov et al., 2024; Baumann & Mendler-Dünner, 2024; Sigg et al., 2025; Gauthier et al., 2025), in which a small group of users strategically modifies their own data to influence the correlations learned by the model.

1.1 Our contributions

Shifting focus from firm-side fairness methods to the user side. After formally introducing fairness in ML and ACA in Section 2, we adapt the *erasure strategy* from Hardt et al. (2023) to reduce predictive correlation between group membership and the target label by relabeling minority samples. By erasing a counterfactual signal, the success metric of the collective aligns with the definition of counterfactual fairness, which under certain conditions can translate to group fairness.

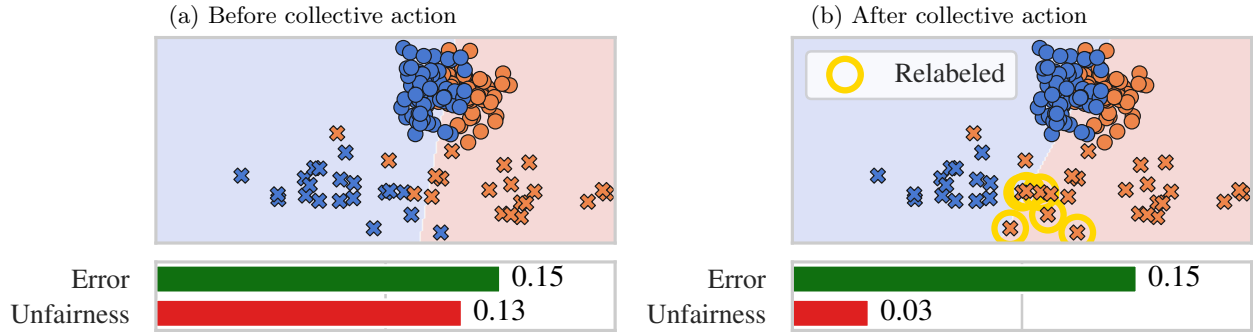


Figure 1: Minority-only collective action can substantially improve fairness. When only 6 minority members change their labels, the EqOd violation (Eq. (2)) of logistic regression drops by over 75% with a negligible rise in prediction error (Eq. (1)). Circles represent the majority and crosses represent the minority.

Designing minority-only collective action with counterfactual estimation. The erasure strategy faces two major challenges in our setting. First, the collective is limited to only minority members, since minority members are more motivated to join collective action (Saleem et al., 2021; Begeny et al., 2022) and can be efficiently mobilized (McAdam, 1999; Michelson, 2005), while majority-group users may be less inclined to disrupt the status quo. Second, the optimal strategy requires counterfactual labels which are generally intractable. To overcome these obstacles, in Section 3 we propose a strategy in which only a specific number of candidates flip their labels. The candidates are ranked by one of three simple methods that estimate how likely a candidate is to have a positive counterfactual label. To test this strategy, in Section 4 we conduct experiments on common real-world fairness datasets. As we cannot evaluate counterfactual fairness, we instead measure group fairness metrics as surrogates, i.e., equalized odds (EqOd) and statistical parity (SP). These experiments show that the proposed strategy substantially decreases group fairness violations with only a small impact on prediction error, as illustrated in Fig. 1.

Limitations of minority-only collectives. Section 5 investigates the fundamental limitations of minority-only collectives, such as the impossibility of reaching zero EqOd violation, and provides theoretical results showing that learned representations and accurate counterfactual approximation methods can improve on the basic strategy we propose.

2 Collective Action for Fairness

To establish the connection between fairness and ACA, Section 2.1 first defines the problem setting and how unfairness can be measured. Then, Section 2.2 describes the theoretical framework of ACA and how it can be utilized to mitigate bias. Finally, Section 2.3 formally relates ACA to group fairness metrics through counterfactual fairness.

2.1 Group fairness for classification

In our setting a firm uses ML to predict a binary label $y \in \{0, 1\}$. The firm collects data from its users, forming a dataset $\mathcal{D} = \{(x_i, a_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^m$ denotes user i 's feature vector, $a_i \in \{0, 1\}$ is a sensitive attribute indicating binary group membership ($a_i = 0$ for the majority group, $a_i = 1$ for the minority), and $y_i \in \{0, 1\}$ is the true label. We assume the users are drawn independently and identically distributed (i.i.d.) from a distribution \mathbb{P}_0 over $\mathbb{R}^m \times \{0, 1\} \times \{0, 1\}$. The firm trains a classifier $h : \mathbb{R}^m \rightarrow \{0, 1\}$ to minimize the prediction error, defined

$$\text{Error}(h) = \mathbb{P}[h(x) \neq y]. \quad (1)$$

To do so, the firm minimizes the empirical error on \mathcal{D} via Empirical Risk Minimization (ERM).

In the group-fairness paradigm, the sensitive attribute $a \in \{0, 1\}$ partitions the data into subgroups, and fairness criteria seek to ensure similar outcomes across these groups. Common metrics include statistical

parity (SP) (Calders et al., 2009; Dwork et al., 2012) and equalized odds (EqOd) (Hardt et al., 2016). In this work, we focus primarily on violations of EqOd, formally defined as

$$\text{EqOd}(h) = \frac{1}{2} \sum_{z=0,1} |\mathbb{P}[h(x) = 1|a = 1, y = z] - \mathbb{P}[h(x) = 1|a = 0, y = z]|, \quad (2)$$

which measures the differences between true positive and false positive rates. Appendix B.1 provides formal definitions and further discussion of these metrics.

ERM-trained models tend to achieve low predictive error, but often at the cost of SP and EqOd violations (Menon & Williamson, 2018; Zhao & Gordon, 2019; Bardenhagen et al., 2021; Sanyal et al., 2022). Despite significant progress in fairness research, most solutions have traditionally focused on *firm-side* solutions: pre-processing the dataset, in-processing modifications to the training algorithm, or post-processing the classifier’s predictions. These approaches almost always incur errors or additional pipeline complexity, discouraging firms from deploying them in practice.

While most prior work has focused on firm-side solutions, this work shifts the focus to *user-side* methods that do not require the firm’s participation. Since users generate the training data, they can collectively influence the learned model by strategically modifying their own behavior. For instance, consider a digital platform that recommends content to a user based on a classifier that predicts engagement labels $y_i \in \{\text{will engage, will not engage}\}$. The classifier, trained on historical user interactions, may unintentionally rely on group membership rather than individual preferences when making recommendations for minority members. In response, users can coordinate to alter interaction patterns, such as clicking on or avoiding certain items. This ACA affects the dataset in a way that steers the learned classifier toward fairer outcomes, and is generally studied under the field of algorithmic collective action (Hardt et al., 2023).

2.2 Algorithmic collective action

In social sciences, *collective action* refers to the coordinated efforts of individuals working together to pursue a shared goal (Olson, 1989; Marwell & Oliver, 1993). Hardt et al. (2023) adapt this notion to ML, proposing that a group of users, termed a collective, can strategically modify their data to align the behavior of a trained classifier h with the collective’s goals. In this formulation, the training distribution is a mixture distribution $\mathcal{D} \sim \mathbb{P}_\alpha = \alpha\mathbb{P}^* + (1 - \alpha)\mathbb{P}_0$, where \mathbb{P}^* and \mathbb{P}_0 are the collective and base distributions, and $\alpha \in [0, 1]$ denotes the proportion of the collective.

Relation to fair representation learning. With user agency over the data, one possible form of ACA for fairness is to modify their features to increase correlation with the label $y = 1$. An analogous firm-side approach is fair representation learning (FRL), which learns a transformation from the input space to a representation space such that ERM leads to a classifier that is both accurate and fair (Zemel et al., 2013; Jovanović et al., 2023). However, a hindrance of FRL in the context of ACA is that the transformation must be applied consistently at inference time, requiring active cooperation from each minority member to transform their features. In contrast, our setting assumes users have control only over the labels and cannot intervene in other parts of the machine learning pipeline.

Erasing a signal. Suppose the collective seeks a classifier that is invariant under a transformation $g : \mathbb{R}^m \rightarrow \mathbb{R}^m$ applied to the features. The success of the collective is

$$S(\alpha) = \mathbb{P}_0[h(g(x)) = h(x)], \quad (3)$$

the probability, under the base distribution, that the classifier’s prediction is unchanged after applying g to the features. In words, the collective’s goal is to *erase the signal* g : to ensure the classifier behaves identically regardless if g is applied. Intuitively, if g embeds a pattern correlated with group membership (i.e., minority or majority), then achieving invariance under g promotes fairness by reducing the classifier’s dependence on group-identifying information.

To achieve signal erasure, Hardt et al. (2023) propose the collective relabels itself with the most likely label under the transformation g . Formally, the strategy is defined as

$$x, y \rightarrow x, \operatorname{argmax}_{y' \in \{0,1\}} \mathbb{P}_0(y'|g(x)). \quad (4)$$

Since this strategy leaves the features unchanged, it is well-suited for settings where the minority is limited to modify only their labels, such as ours. For ϵ -optimal Bayes classifiers (Definition B.1 in Appendix B.2), Hardt et al. (2023) prove the following lower bound for its success

$$S(\alpha) \geq 1 - \frac{2(1-\alpha)}{\alpha} \cdot \tau - \frac{\epsilon}{(1-\epsilon)\alpha}, \quad (5)$$

where $\tau = \mathbb{E}_{x \sim \mathbb{P}_0} \left[\max_{y' \in \{0,1\}} |\mathbb{P}_0(y'|x) - \mathbb{P}_0(y'|g(x))| \right]$ captures the sensitivity of y under g .

Note that the strategy in Eq. (4) may require some majority members to relabel themselves with the label $y = 0$. Such a change might deter them from participating in the collective action, either because majority members are unwilling to give up their advantage or prefer to maintain the status quo. To avoid this conflict, we restrict the collective to include only minority members. We discuss the implications of this restriction in Section 5.

2.3 Counterfactual fairness

The concept of *counterfactual fairness* (CF) (Kusner et al., 2017; Garg et al., 2019; Wu et al., 2019) connects signal-erasure success to group fairness. To introduce this idea, assume that a sample x is generated by a causal model, in which the group membership A is a causal parent. Then a classifier h is counterfactually fair if its predictions are invariant to interventions on the group membership, i.e., $h(x) = h(x_{A \leftarrow a'})$ for any a' , where $x_{U \leftarrow u}$ denotes an intervention on a causal parent U of a sample x . In certain causal contexts, CF implies or aligns with group fairness criteria such as SP or EqOd (Anthis & Veitch, 2023). Therefore, if ACA induces a counterfactually fair classifier, it may also induce a fair classifier under SP or EqOd. As our focus is on fairness for the minority, we relax the original definition of CF (Kusner et al., 2017).

Definition 2.1. A classifier h is **minority-focused counterfactually fair** if, under the causal distribution induced by \mathbb{P}_0 ,

$$\mathbb{P}_0(h(x_{A \leftarrow 1}) = y | X = x, A = 1) = \mathbb{P}_0(h(x_{A \leftarrow a'}) = y | X = x, A = 1), \quad (6)$$

for all y and for any value a' attainable by A .

By this definition, changing the group membership of a minority individual, in a counterfactual sense, has no effect on the classifier’s prediction. ACA can theoretically enforce such fairness by applying the erasure strategy from Eq. (4) with the counterfactual signal $g(x) = x_{A \leftarrow 0}$, which replaces a minority individual with its majority-group counterfactual. This ACA aligns the signal erasure success from Eq. (3) with minority-focused counterfactual fairness from Definition 2.1. The following informal proposition, proved in Appendix C.1, formalizes this alignment.

Proposition 2.2 (Informal; see Proposition C.1). *For a deterministic Bayes classifier, perfect success of the minority collective is equivalent to (almost-sure) minority-focused counterfactual fairness.*

This result directly connects ACA theory to fairness. Thus, perfect success of the collective is equivalent to achieving minority-focused counterfactual fairness.

3 Minority Collective Action With Approximated Counterfactuals

This section describes how a minority collective can approximate a signal-erasure strategy to promote fairness in practice. While the theory of signal erasure has been studied before (Hardt et al., 2023; Gauthier et al., 2025), prior work lacks empirical evaluation. In this paper, we present the first practical algorithm for signal

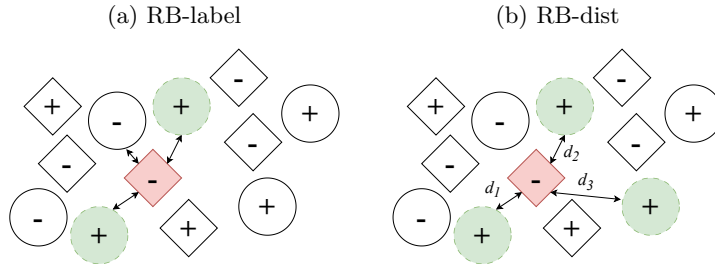


Figure 2: Visualization of KNN scoring methods in Section 3 with $k=3$. Squares represent the minority and circles the majority, marked with a positive “+” or a negative “-” label. **(a) RB-label**: Two of the nearest majority neighbors have a positive label, resulting in the score $s=2$. **(b) RB-dist**: The average distance to the nearest positive majority neighbors results in the score $s=-(d_1+d_2+d_3)/3$.

erasure and provide experimental results in Section 4. As discussed in Section 2.3, a suitable signal to erase is $g(x) = x_{A \leftarrow 0}$, where each collective member relabels themselves according to Eq. (4).

However, end-users lack access to the causal model and cannot compute the counterfactual labels directly. To address this limitation, we propose to assign each collective member i a score s_i , which serves as a proxy for the likelihood that they would receive the label $y = 1$ if they belonged to the majority. Given a budget of M label flips, the collective selects the M members with the highest scores; these individuals flip their labels from $y = 0$ to $y = 1$. The budget M controls the accuracy–fairness tradeoff: a higher budget typically improves fairness but increases error.

We suggest three scoring functions, each capturing a different notion of similarity to majority users:

1. **Rank by probability (RB-prob)**: Train a regressor $f : \mathbb{R}^m \rightarrow \mathbb{R}$ exclusively on majority data ($a = 0$) to estimate the probability $\mathbb{P}(Y = 1 | X = x, A = 0)$ of having the label $y = 1$ under the majority distribution. Each collective member i receives a score based on the model’s prediction:

$$s_i = f(x_i). \quad (7)$$

2. **Rank by label (RB-label)**: For each collective member i , identify the set K_i of their k nearest majority neighbors using Euclidean distance. The score is the number of neighbors with the label $y = 1$:

$$s_i = \sum_{j \in K_i} \mathbf{1}\{y_j = 1\}. \quad (8)$$

3. **Rank by distance (RB-dist)**: Restrict the neighbors set K_i to only majority users with the label $Y = 1$. The score is the negative mean Euclidean distance to these neighbors:

$$s_i = -\frac{1}{k} \sum_{j \in K_i} \|x_i - x_j\|_2. \quad (9)$$

Intuitively, RB-prob acts as a naive counterfactual estimator by training only on majority data. RB-label assumes that a sample with more positive-labeled majority neighbors is more likely to have a positive counterfactual label, and RB-dist assumes that proximity to positive-labeled majority samples is evidence for a positive counterfactual label. Note that RB-label and RB-dist are k nearest neighbors (KNN) based methods, which are common in counterfactual estimation (Stuart, 2010; Li et al., 2016; Verma et al., 2024).

4 Experimental Results

This section evaluates the performance of our methods. However, since we cannot compute the counterfactual fairness gap in Definition 2.1, we instead measure surrogate group fairness metrics. The main text presents

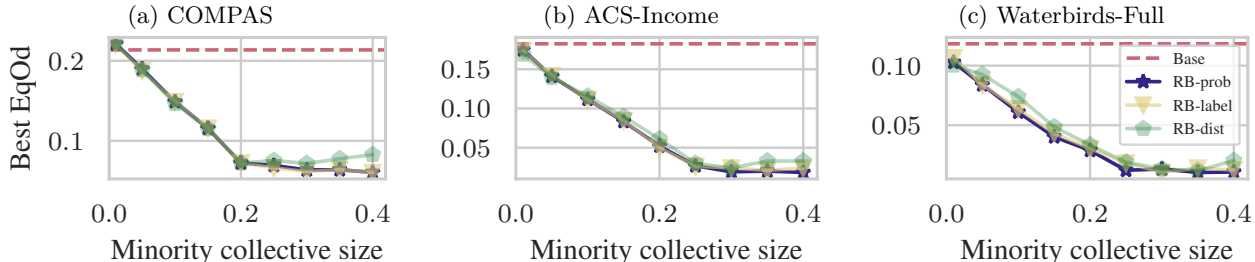


Figure 3: The lowest EqOd violation a collective can achieve decreases as the collective size increases, up to a certain point. Each point is the mean of 10 runs, with the standard deviation being smaller than the markers. Across the datasets we experimented on, the lowest EqOd violation stabilizes around $\alpha = 0.3$. Additional results are presented in Fig. 10 in the appendix.

results for EqOd, while Appendix E.2 presents results for both EqOd and SP, along with additional datasets. We compare the three methods, RB-label, RB-dist, RB-prob, against a random baseline that flips $y = 0$ labels to $y = 1$ for M randomly selected collective members. We conducted experiments on the tabular datasets COMPAS (Mattu et al., 2016), Adult (Becker & Kohavi, 1996), HSLs (Jeong et al., 2022), and ACS-Income (Ding et al., 2021), the image dataset Waterbirds (Sagawa et al., 2020), and the text dataset CivilComments (Borkan et al., 2019). For Waterbirds, we use features extracted from a pre-trained *ResNet-18* (denoted Waterbirds-Full), and for CivilComments, we use features extracted from Hugging Face’s pre-trained *bert-base-uncased* model (denoted CivilComments-Full). In addition to the complete features of Waterbirds and CivilComments, we also include experiments on the PCA features, with 85 components for Waterbirds (denoted Waterbirds-PCA) and 100 components for CivilComments (denoted CivilComments-PCA). Details on the datasets and preprocessing are provided in Appendix D.1.

All reported metrics are computed on a fixed test set without any ACA and averaged over 10 independent runs for each method described in Section 3. In each run, we randomly selected a minority collective to apply the method. For the KNN-based methods, we tuned the neighborhood size k using a 15% validation split from the train set, optimizing for EqOd and SP. Finally, we trained a gradient-boosted decision tree on each modified train set. More technical details are given at Appendix D.2.

Importance of collective size While the number of label flips M is the primary factor for balancing between accuracy and fairness, the size of the collective, α , also plays a role. In addition to bounding the possible number of flips, increasing α also expands the candidate pool from which the most effective labels to flip can be selected. To measure this effect, our experiments include a range of α values, each tested with multiple values of M . For each α , we define the best achievable EqOd as the minimum EqOd across all tested values of M . As shown in Fig. 3, increasing α decreases the lowest achievable EqOd violation until saturating around $\alpha = 0.3$. We fix this value for all remaining experiments.

Flipping cost Since each method scores candidates differently, they may also vary in efficiency, that is, the number of label flips required to achieve a given level of EqOd violation. To evaluate efficiency, Fig. 4 plots EqOd as a function of number of label flips M , where lower curves indicate more efficient methods. The random baseline consistently yields the worst EqOd across all values of M , highlighting the value of informed relabeling algorithms. However, no single method dominates the others in all settings: While RB-prob and RB-label often outperform the other methods, RB-dist can surpass them in specific cases (e.g., Fig. 4a), or perform comparably to the random baseline in others (Fig. 4c). These results suggest that a well-chosen scoring function enables the collective to achieve a desired level of EqOd violation with fewer label flips, reducing the “cost” of ACA and mitigating the accuracy loss from excessive relabeling.

Interestingly, beyond a certain number of flips, EqOd begins to increase, indicating that excessive flipping can shift unfairness from the minority to the majority. This upturn reflects the fundamental limits of minority ACA for fairness, a point we elaborate on in Section 5.

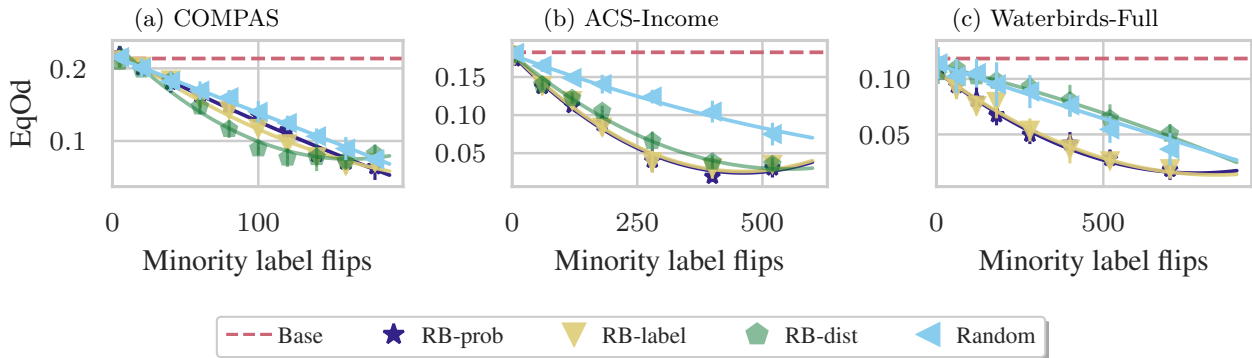


Figure 4: Our proposed methods are generally more efficient than randomly flipping labels, requiring fewer label flips to attain the same EqOd violation level. Each marker is the mean of 10 random runs with a specific number of label flips. Error bars show the standard deviation. The dashed line shows the mean EqOd for a classifier trained on the dataset without collective action.

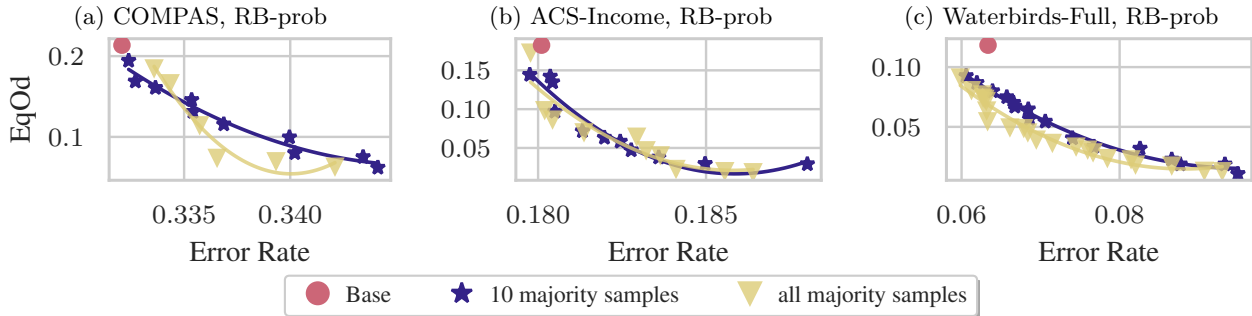


Figure 5: Limiting the collective’s knowledge of the majority does not significantly harm the Pareto front. Each point is the mean of 10 runs and the curves are fitted to guide the eye.

Partial knowledge of the majority In the previous experiments, we assumed that the collective has full access to the majority data. Here we investigate the performance of our methods when this knowledge is limited. To visualize the fairness-error tradeoff, we measure the error and EqOd for a range of label flips, yielding a set of pairs (Error, EqOd). This set forms a Pareto front, representing the tradeoff. A Pareto front is said to *dominate* another if it lies entirely to the left (lower error) and below (lower unfairness) of the other. The Pareto fronts in Fig. 5 show that a collective employing RB-prob, when restricted to only 10 majority members, performs similarly to a collective with full knowledge. While the Pareto fronts remain similar, limited majority knowledge can increase the number of required flips. This is evident when comparing to the zero-knowledge scenario, designated as random in Fig. 4. This finding implies that the fewer flips the collective is allowed, the more important it is to have access to the majority data.

5 Limitations of Minority Collective Action

Previous work on ACA assumes that the collective is uniformly sampled from the distribution \mathbb{P}_0 and that the collective has a perfect oracle for the conditional distribution $\mathbb{P}_0(Y|X)$. Yet, our method restricts collective participation to minority members and approximates this conditional distribution. These differences introduce limitations to the existing theory, which we analyze in this section.

Collective restricted to the minority. As mentioned above, we focus on minority collectives, unlike prior work. This restriction captures scenarios in which majority members lack incentives to support changes that would benefit the minority and instead prefer to preserve the status quo. Naturally, this limitation

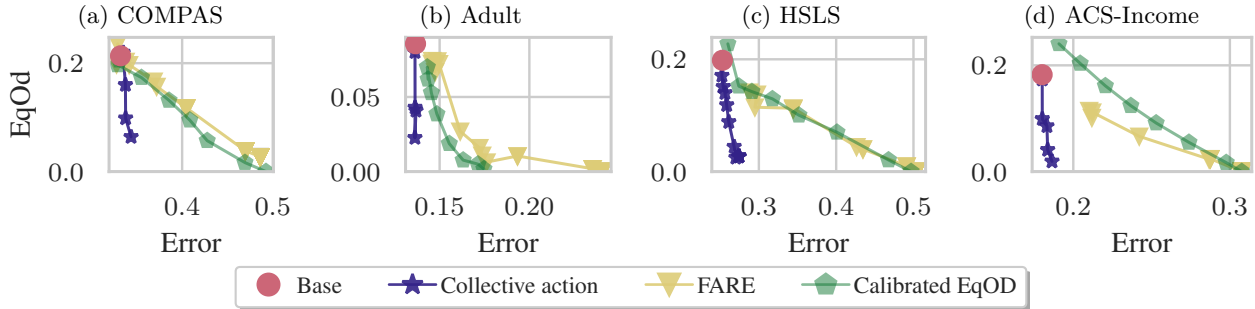


Figure 7: The user-side method cannot achieve zero EqOd violation, while the firm-side pre-processing method FARE Jovanović et al. (2023) and the post-processing method calibrated equalized odds Pleiss et al. (2017) attain 0 EqOd with large error. However, RB-prob has lower EqOd violation than the base classifier, with a smaller error than the firm-side methods.

reduces the collective’s impact. Consider a binary classification task on the two-dimensional 4-Gaussian mixture model $\mathbb{P}_{4\text{GMM}}$ where each Gaussian belongs to a distinct combination of label and group membership, as illustrated in Fig. 6. Each label consists of a large majority subgroup and a significantly smaller minority subgroup. We can then state the following informal result about the EqOd violation of ERM.

Proposition 5.1 (Informal). *Consider $\mathbb{P}_{4\text{GMM}}$, where every minority point participates in the ACA by flipping all $y=0$ labels to $y=1$. Then, under sufficiently separable clusters, with high probability, the EqOd of the ERM classifier minimizing the logistic loss will asymptotically approach 0.5.*

The formal Proposition C.3, which holds for a broader family of distributions, is provided in Appendix C.2 along with all necessary assumptions. It can be extended to any dimensionality \mathbb{R}^d using techniques similar to those in Chaudhuri et al. (2023). Although Proposition 5.1 is not a formal lower bound, it emphasizes an important limitation: ACA restricted to the minority cannot generally achieve zero EqOd violation, even under very advantageous conditions involving a maximum-sized collective, a strong strategy, and a disregard for accuracy. This limitation stands in contrast to standard firm-side bias mitigation methods, which can in principle achieve zero EqOd violation. There are several reasons why relabeling alone may not be enough to get zero EqOd violation. For one, relabeling according to the counterfactual implicitly assumes that the label is determined by the same features across the majority and minority, but this assumption may fail under distribution shift between the groups. To illustrate, consider a firm training a classifier to screen candidates for a managerial position. Majority members may be more educated, while minority members may have more hands-on experience rather than formal education. In this case, a counterfactual label associated with the majority is disjoint from the features associated with the minority, rendering the signal erasure strategy irrelevant. Future work could study this problem and determine when it is beneficial to change the features as well.

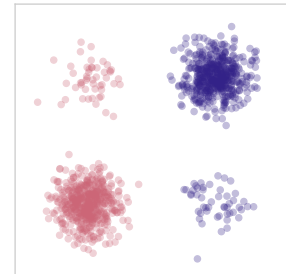


Figure 6: The distribution $\mathbb{P}_{4\text{GMM}}$ used in Proposition 5.1. The color signifies the label, and the density shows the group membership.

We empirically corroborate the findings of Proposition 5.1 on real-world datasets by examining the EqOd–accuracy tradeoff of several fair ML methods. Fig. 7 compares the Pareto fronts of RB-prob with established firm-side methods. We observe that the lowest EqOd violation achievable by RB-prob is greater than that of the firm-side approaches. However, the firm-side methods arrive at zero EqOd violation only at the cost of prohibitively high prediction error. In the region where the error is small compared to the base classifier, the EqOd violation of RB-prob is comparable to that of the firm-side methods.

Estimating counterfactuals. The methods from Section 3 estimate which individuals would receive a counterfactual label that is different than their original label. However, the success lower bound in Eq. (5) assumes perfect knowledge of \mathbb{P}_0 and of the underlying causal model. To account for estimation error, we

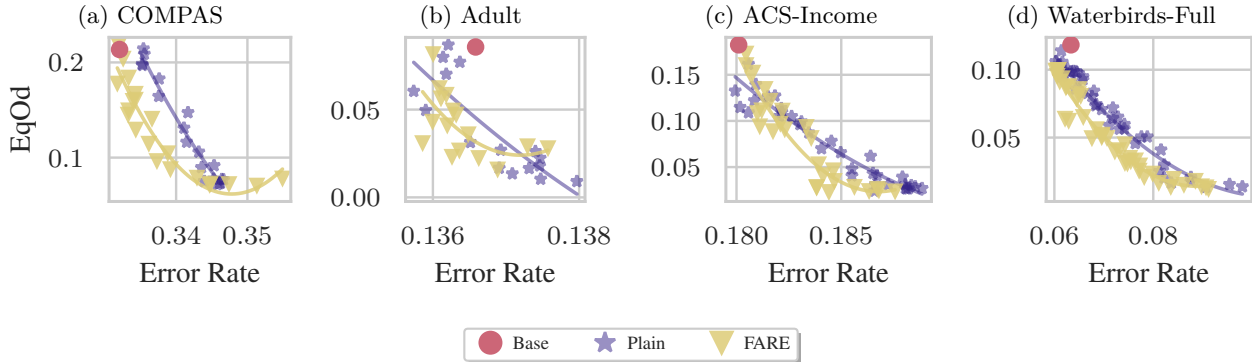


Figure 8: The Pareto fronts for using a fair representation when computing the KNN for RB-dist dominate the Pareto fronts for KNN computed on untransformed features. The blue stars represent the KNN without transforming the data, and the yellow triangles represent the KNN when the data is transformed using FARE Jovanović et al. (2023). The lines are fitted to guide the eye.

model the collective’s prediction as the output of an algorithm $\mathcal{A}(x) \approx \arg \max_y \mathbb{P}_0(y|x_{A \leftarrow 0})$ with error rate ρ , defined as

$$\rho := \mathbb{P}_0(\mathcal{A}(x) \neq \arg \max_y \mathbb{P}_0[y'|g(x)]). \quad (10)$$

Given this definition, we derive the following lower bound on success, proved in Appendix C.3.

Proposition 5.2. *With an algorithm $\mathcal{A}(x)$ with label error $\rho < 1/2$, the success of the collective is bounded by*

$$S(\alpha) \geq 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha} \tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}. \quad (11)$$

This bound recovers Eq. (5) when $\rho=0$, but higher ρ worsens the bound. Next, we show how to use FRL to reduce the error ρ , thereby improving the lower bound.

Impact of feature representations Since the methods RB-label and RB-dist rely on KNN, they are sensitive to the choice of distance metric and feature representation. In Section 4, we used Euclidean distance in the original feature space, which is convenient but could be suboptimal. Here, we explore whether FRL can learn a more suitable representation space for KNN. A *fair representation* maps the data into a space where the group-based bias is removed, while preserving informative features. Intuitively, such representations may help RB-label and RB-dist to better estimate the counterfactual labels. To formalize this intuition, we consider predicting the counterfactual label of minority points using a 1-NN classifier on majority data, i.e., assigning each minority point the label of its nearest neighbor in the majority. In settings where the minority and majority are differently distributed (e.g., \mathbb{P}_{4GMM}), this task can be challenging. The following informal result compares the error of 1-NN in the original feature space to its error in FRL.

Proposition 5.3 (Informal). *Let data be drawn from \mathbb{P}_{4GMM} , and ρ_{plain} denote the error of a 1-NN classifier that assigns the label of the nearest majority neighbor in the original feature space. Then there exists a fair representation in which a 1-NN classifier achieves error ρ_{FRL} such that, asymptotically with respect to the dataset size, $\rho_{FRL} \leq \rho_{plain}$.*

The formal statement, Theorem C.6, suggests that FRL can reduce the counterfactual label error ρ of KNN methods, raising the lower bound of the success according to Proposition 5.2. Empirically, Fig. 8 indicates that applying FARE (Jovanović et al., 2023) before the KNN step improves the Pareto front for RB-dist. Yet, methods that rely purely on predictive information, such as RB-prob, can perform worse, due to inadvertent removal of features.

6 Related work

Fairness in ML often comes at the cost of reduced classification accuracy, leading to the well-documented accuracy–fairness tradeoff (Menon & Williamson, 2018; Zhao & Gordon, 2019; Dehdashtian et al., 2024; Sadeghi et al., 2022). In response, previous work has proposed fairness interventions at different stages of the ML pipeline: pre-processing methods modify the training data before learning (Kamiran & Calders, 2009; Luong et al., 2011; Zemel et al., 2013; Jovanović et al., 2023), in-processing methods adjust the learning algorithm itself (Agarwal et al., 2018; Nam et al., 2020; Sagawa et al., 2020; Liu et al., 2021), and post-processing methods correct the predictions of a trained (unfair) classifier (Hardt et al., 2016; Alghamdi et al., 2022; Tifrea et al., 2024; Cruz & Hardt, 2024). A firm can introduce any of these categories into its pipeline, while users, who control only their data, can only partially implement pre-processing methods. However, as mentioned in Section 2.2, feature-changing pre-processing methods such as fair representation learning (Zemel et al., 2013; Jovanović et al., 2023) demand changing features at inference time.

Still, some pre-processing methods change only the labels, similarly to our proposed collective action. The method by Luong et al. (2011) compares minority KNN and majority KNN and flips labels according to the difference in positive labels between the two groups of neighbors. This method resembles RB-label, with the difference that RB-label examines only the majority KNN in order to approximate the counterfactual. The approach of Kamiran & Calders (2009) trains a regressor to predict $y = 1$ outcome probabilities, flips minority members with $y = 0$ labels and high predicted probabilities to $y = 1$, and similarly flips majority $y = 1$ labels to $y = 0$. Flipping from both groups is intended to preserve the error of the classifier. Our method RB-prob differs by training the regressor only on the majority to better approximate the counterfactuals. Since this approach requires flipping the labels of majority members as well, it cannot be completely adopted by the collective. In Appendix E.1 we compare RB-prob to CND and KDP, and find that our method is more efficient in terms of the number of label flips. However, as Fig. 4 demonstrates, each dataset has different sensitivity to the counterfactual approximation method, and it is not guaranteed that our proposed methods will always perform better. This uncertainty is a result of several factors, including the effect of the data representation (as illustrated in Fig. 8) and the fact that the collective modifies only the labels, even though the counterfactual features may also be different.

7 Conclusion

This work demonstrates that user-side methods can effectively reduce unfairness in machine learning. While much of the existing fairness research focused on firm-side methods, these often come at a cost that may not be worth to the firm. This tradeoff emphasizes the importance of studying user-side approaches for bias mitigation. We show empirically that ACA can considerably reduce unfairness in a variety of datasets, though not completely. Importantly, we also examine the limitations of collectives composed only of minority members, and how success is affected by approximating the counterfactual labels. Our proposed methods require collective members to relabel themselves, which often comes with a price, as users have to report labels that differ from their uncoordinated labels. However, real-world studies show that minority members can willingly participate in collective action to benefit their demographic after being encouraged by their community (Begeny et al., 2022) or by the media (Saleem et al., 2021) and efficiently mobilized at large enough scales (McAdam, 1999; Michelson, 2005). Overall, this paper shows a practical use case of collective action in the hopes of sparking further research into applications of ACA and user-side methods for social good.

Broader Impact Statement. Our paper discusses how ACA can contribute to social good by improving fairness in ML systems. However, ACA is a collaboration among many individuals sharing the same goal, and that goal may be malicious, such as gaining self-favor at the cost of the greater good. While our paper utilizes ACA for fairness in ML, it does not provide guarantees against misuse. As seen in our motivating problems in Appendix A, we advocate responsible use of ACA, and believe that any organized collective should be transparent and have clear goals.

References

- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A reductions approach to fair classification. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 60–69. PMLR, 2018.
- 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. In *Advances in Neural Information Processing Systems*, volume 35, pp. 38747–38760. Curran Associates, Inc., 2022.
- Jacy Anthis and Victor Veitch. Causal context connects counterfactual fairness to robust prediction and group fairness. In *Advances in Neural Information Processing Systems*, volume 36, pp. 34122–34138. Curran Associates, Inc., 2023.
- Vincent Bardenhagen, Alexandru Tifrea, and Fan Yang. Boosting worst-group accuracy without group annotations. In *NeurIPS 2021 Workshop on Distribution Shifts: Connecting Methods and Applications*. OpenReview, 2021.
- Solon Barocas and Andrew D. Selbst. Big data’s disparate impact. *California Law Review*, 104(3):671–732, 2016.
- Joachim Baumann and Celestine Mendler-Dünger. Algorithmic Collective Action in Recommender Systems: Promoting Songs by Reordering Playlists. In *Advances in Neural Information Processing Systems*, volume 37, pp. 119123–119149. Curran Associates, Inc., 2024.
- Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996.
- Christopher T. Begeny, Jolien van Breen, Colin Wayne Leach, Martijn van Zomeren, and Aarti Iyer. The power of the Ingroup for promoting collective action: How distinctive treatment from fellow minority members motivates collective action. *Journal of Experimental Social Psychology*, 101:104346, 2022.
- Omri Ben-Dov, Jake Fawkes, Samira Samadi, and Amartya Sanyal. The Role of Learning Algorithms in Collective Action. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pp. 3443–3461. PMLR, 2024.
- Daniel Borkan, Lucas Dixon, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Nuanced Metrics for Measuring Unintended Bias with Real Data for Text Classification. In *Companion Proceedings of The 2019 World Wide Web Conference*, pp. 491–500. Association for Computing Machinery, 2019.
- Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building Classifiers with Independency Constraints. In *2009 IEEE International Conference on Data Mining Workshops*, pp. 13–18, 2009.
- Kamalika Chaudhuri, Kartik Ahuja, Martin Arjovsky, and David Lopez-Paz. Why does throwing away data improve worst-group error? In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 4144–4188. PMLR, 2023.
- André Cruz and Moritz Hardt. Unprocessing seven years of algorithmic fairness. In *The Twelfth International Conference on Learning Representations*, 2024.
- Sepehr Dehdashtian, Bashir Sadeghi, and Vishnu Naresh Boddeti. Utility-fairness trade-offs and how to find them. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 12037–12046, 2024.
- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring Adult: New Datasets for Fair Machine Learning. In *Advances in Neural Information Processing Systems*, volume 34, pp. 6478–6490. Curran Associates, Inc., 2021.
- Cynthia Dwork, Moritz Hardt, Toniann Pitassi, Omer Reingold, and Richard Zemel. Fairness through awareness. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pp. 214–226. Association for Computing Machinery, 2012.

- Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H. Chi, and Alex Beutel. Counterfactual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 219–226. Association for Computing Machinery, 2019.
- Etienne Gauthier, Francis Bach, and Michael I. Jordan. Statistical collusion by collectives on learning platforms. In *Forty-Second International Conference on Machine Learning*, 2025.
- Moritz Hardt, Eric Price, and Nathan Srebro. Equality of Opportunity in Supervised Learning. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 3323–3331, 2016.
- Moritz Hardt, Eric Mazumdar, Celestine Mendler-Dünnler, and Tijana Zrnica. Algorithmic Collective Action in Machine Learning. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202, pp. 12570–12586, 2023.
- Haewon Jeong, Hao Wang, and Flavio P. Calmon. Fairness without Imputation: A Decision Tree Approach for Fair Prediction with Missing Values. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(9):9558–9566, 2022.
- Ziwei Ji and Matus Telgarsky. The implicit bias of gradient descent on nonseparable data. In *Proceedings of the Thirty-Second Conference on Learning Theory*, volume 99, pp. 1772–1798. PMLR, 2019.
- Nikola Jovanović, Mislav Balunovic, Dimitar Iliev Dimitrov, and Martin Vechev. FARE: Provably Fair Representation Learning with Practical Certificates. In *Proceedings of the 40th International Conference on Machine Learning*, pp. 15401–15420. PMLR, 2023.
- Faisal Kamiran and Toon Calders. Classifying without discriminating. In *Control and Communication 2009 2nd International Conference on Computer*, pp. 1–6, 2009.
- Matt J Kusner, Joshua Loftus, Chris Russell, and Ricardo Silva. Counterfactual fairness. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Sheng Li, Nikos Vlassis, Jaya Kawale, and Yun Fu. Matching via dimensionality reduction for estimation of treatment effects in digital marketing campaigns. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pp. 3768–3774. AAAI Press, 2016.
- Evan Z Liu, Behzad Haghgoo, Annie S Chen, Aditi Raghunathan, Pang Wei Koh, Shiori Sagawa, Percy Liang, and Chelsea Finn. Just Train Twice: Improving Group Robustness without Training Group Information. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 6781–6792, 2021.
- Binh Thanh Luong, Salvatore Ruggieri, and Franco Turini. K-NN as an implementation of situation testing for discrimination discovery and prevention. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 502–510. Association for Computing Machinery, 2011.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. Learning Adversarially Fair and Transferable Representations. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 3384–3393. PMLR, 2018.
- Gerald Marwell and Pamela Oliver. *The Critical Mass in Collective Action*. Cambridge University Press, 1993.
- Jeff Mattu, Julia Larson, Lauren Angwin, and Surya Kirchner. How We Analyzed the COMPAS Recidivism Algorithm. <https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm>, 2016.
- Doug McAdam. *Political Process and the Development of Black Insurgency, 1930-1970*. University of Chicago Press, 1999.
- Aditya Krishna Menon and Robert C. Williamson. The cost of fairness in binary classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency*, pp. 107–118. PMLR, 2018.

- Melissa R. Michelson. Meeting the Challenge of Latino Voter Mobilization. *The ANNALS of the American Academy of Political and Social Science*, 601(1):85–101, 2005.
- Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684, 2020.
- Mancur Olson. Collective action. In *The Invisible Hand*, pp. 61–69. Palgrave Macmillan UK, 1989.
- Geoff Pleiss, Manish Raghavan, Felix Wu, Jon Kleinberg, and Kilian Q Weinberger. On fairness and calibration. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- Bashir Sadeghi, Sepehr Dehdashtian, and Vishnu Boddeti. On Characterizing the Trade-off in Invariant Representation Learning. *Transactions on Machine Learning Research*, 2022.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally Robust Neural Networks. In *Eighth International Conference on Learning Representations*, 2020.
- Muniba Saleem, Ian Hawkins, Magdalena E. Wojcieszak, and Jessica Roden. When and how negative news coverage empowers collective action in minorities. *Communication Research*, 48(2):291–316, 2021.
- Amartya Sanyal, Yaxi Hu, and Fanny Yang. How unfair is private learning? In *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180, pp. 1738–1748. PMLR, 2022.
- Dorothee Sigg, Moritz Hardt, and Celestine Mendler-Dünner. Decline now: A combinatorial model for algorithmic collective action. In *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems*. Association for Computing Machinery, 2025.
- Daniel Soudry, Elad Hoffer, Mor Shpigel Nacson, Suriya Gunasekar, and Nathan Srebro. The implicit bias of gradient descent on separable data. *Journal of Machine Learning Research*, 19(70):1–57, 2018.
- Elizabeth A. Stuart. Matching methods for causal inference: A review and a look forward. *Statistical science : a review journal of the Institute of Mathematical Statistics*, 25(1):1–21, 2010.
- Alexandru Tifrea, Preethi Lahoti, Ben Packer, Yoni Halpern, Ahmad Beirami, and Flavien Prost. FRAPPÉ: A group fairness framework for post-processing everything. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235, pp. 48321–48343. PMLR, 2024.
- Sahil Verma, Varich Boonsanong, Minh Hoang, Keegan Hines, John Dickerson, and Chirag Shah. Counterfactual explanations and algorithmic recourses for machine learning: A review. *Acm Computing Surveys*, 56(12), 2024.
- Yongkai Wu, Lu Zhang, and Xintao Wu. Counterfactual fairness: Unidentification, bound and algorithm. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 1438–1444. International Joint Conferences on Artificial Intelligence Organization, 2019.
- Richard Zemel, Yu Wu, Kevin Swersky, Toni Pitassi, and Cynthia Dwork. Learning Fair Representations. In *Proceedings of the 30th International Conference on Machine Learning*, volume 28, pp. 325–333. PMLR, 2013.
- Han Zhao and Geoff Gordon. Inherent tradeoffs in learning fair representations. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.

A Motivating examples

To recognize real-world problems that lend themselves well to collective action for fairness one needs to look for the following few characteristics:

- Firm and goal: A firm trains a predictive model primarily to minimize average error, with little incentive to protect minority groups.
- End-users: People who use the platform and whose behavior generates data for the firm’s dataset.
- Disadvantaged group: A subgroup of end-users who is treated unfairly.
- Relabeling possibility: How the minority can relabel themselves to make the trained classifier fairer.

Here we suggest several cases that answer these characteristics.

1. Content moderation

- Firm and goal: A global social-media company optimizes a high-recall harmful-content detector measured on its largest user pools.
- End-users: Everyday users of the platform who can flag offensive content.
- Disadvantaged group: Slurs, insults, or cultural references specific to minority communities are not flagged often enough, so the model fails to detect harmful content in those groups’ languages.
- Relabeling possibility: The minority flags borderline content from their community that the platform’s global guidelines ignore.

2. Resume screening

- Firm and goal: A multi-national HR firm trains a classifier to extract skills from resumes.
- End-users: Job applicants submitting resumes.
- Disadvantaged group: Applicants from a disadvantaged minority may lack formal education and degrees compared to the majority, but may have informal training which the classifier ignores.
- Relabeling possibility: Applicants can reframe their work experience, e.g. framing working at a store as being a salesperson, or managing shifts as managerial experience.

3. Medical treatment prediction

- Firm and goal: A nationwide insurer builds a treatment-recommendation model to minimize average costs and adverse events.
- End-users: Patients who report their treatment outcomes (pain levels, recovery time, side effects).
- Disadvantaged group: Minority groups may experience different side effects or recovery rates than the majority, so the model recommends suboptimal treatments for them.
- Relabeling possibility: Individual patients record more detailed outcomes rather than underreporting, e.g., consistently marking “still in pain” instead of “fine”.

4. Credit scoring

- Firm and goal: A lender trains a credit-risk model to predict defaults and set loan terms, using historical repayment data.
- End-users: Borrowers whose repayment or default becomes training labels.
- Disadvantaged group: Disadvantaged groups may not have credit cards or have never taken loans, and only deal with cash but still pay their bills. These actions are “credit-invisible”.
- Relabeling possibility: A borrower can report their paid bills, such as rent or utilities, as repaid loans. These become additional positive repayment labels.

5. Recommender systems

- Firm and goal: A streaming platform trains recommender system to maximize engagement, heavily weighted toward mainstream content Baumann & Mendler-Dünner (2024).
- End-users: Users who like, skip, or re-listen to songs.
- Disadvantaged group: Niche genres or local musicians get suppressed, as engagement data mostly comes from the majority’s preferences.
- Relabeling possibility: Users can promote underrepresented content by repeatedly listening, liking, or playlisting it.

B Preliminaries

B.1 Statistical parity and equalized odds

Among the various ways fairness can be defined in machine learning, group fairness is one of the most studied. Group fairness requires that a model’s predictions should not systematically differ between protected groups. One standard measure of this is statistical parity (SP), which captures the difference in the probability of a positive prediction across groups. Formally, it is defined as

$$\text{SP}(h) = |P[h(x) = 1|a = 1] - P[h(x) = 1|a = 0]|, \quad (12)$$

where a smaller SP value indicates fairer treatment across groups. However, SP does not account for the ground-truth labels y , and thus optimizing for SP can degrade the overall accuracy. For example, a classifier that always predicts $\hat{y} = 1$ will have perfect SP but a high prediction error. Alternatively, a stricter notion called equalized odds (EqOd) Hardt et al. (2016) requires that both the true positive rate and false positive rate be equal across groups. Here the EqOd difference is defined as

$$\text{EqOd}(h) = \frac{1}{2} \sum_{z=0,1} |P[h(x) = 1|a = 1, y = z] - P[h(x) = 1|a = 0, y = z]|. \quad (13)$$

B.2 Suboptimal Bayes classifier

Definition B.1 (ϵ -suboptimal classifier). A classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$ is ϵ -suboptimal on a set $\mathcal{X}' \subseteq \mathcal{X}$ under the distribution \mathbb{P} if there exists a \mathbb{P}' with $\text{TV}(\mathbb{P}_{Y|X=x}, \mathbb{P}'_{Y|X=x}) \leq \epsilon$ such that for all $x \in \mathcal{X}'$

$$f(x) = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{P}'(y|x).$$

$\text{TV}(\cdot, \cdot)$ is the total variation distance between two distributions. The definition is discussed more in Hardt et al. (2023).

C Theoretical Results and Proofs

C.1 Counterfactual fairness as success

Proposition C.1. [Formal version of Proposition 2.2] Let $h_\alpha(x) = \operatorname{argmax}_y P_\alpha(Y = y | X = x)$ be a deterministic Bayes classifier, and let $g(x) = x_{A \leftarrow 0}$. Suppose that $g(x) = x$ for majority individuals with $A = 0$, and let $\beta = P_0(A = 1) > 0$. Then $S(\alpha) = 1$ if and only if

$$h_\alpha(X) = h_\alpha(X_{A \leftarrow 0}) \quad P_0\text{-almost-surely conditional on } A = 1.$$

Proof. By definition of the success $S(\alpha)$ we have,

$$S(\alpha) = P_0[h_\alpha(X) = h_\alpha(g(X))].$$

For majority individuals, $A = 0$, the intervention $A \leftarrow 0$ leaves the group membership unchanged, so $g(X) = X$. Hence the event $h_\alpha(X) = h_\alpha(g(X))$ holds with probability one conditional on $A = 0$. Therefore,

$$S(\alpha) = 1 - \beta + \beta P_0[h_\alpha(X) = h_\alpha(X_{A \leftarrow 0}) \mid A = 1].$$

Since $\beta > 0$, we have $S(\alpha) = 1$ if and only if

$$P_0[h_\alpha(X) = h_\alpha(X_{A \leftarrow 0}) \mid A = 1] = 1,$$

which is exactly the claimed almost-sure counterfactual invariance for minority individuals. \square

C.2 Impossibility of fairness under ERM

The following proposition follows the structure of Theorem 6 in Chaudhuri et al. (2023). For a vector $x \in \mathbb{R}^d$, let $D(x)$ denote a distribution on \mathbb{R}^d with mean x . Let p and m be the number of majority and minority sample, respectively with $p \gg m$.

Assumption C.2 (Concentration Condition, Assumption 2 from Chaudhuri et al. (2023)). Let $x_1, \dots, x_n \stackrel{\text{i.i.d.}}{\sim} D(0)$ in \mathbb{R}^d . There exist maps $X_{\max}, c, C : \mathbb{Z}_+ \times [0, 1] \times \mathbb{Z}_+ \rightarrow \mathbb{R}$ such that for all $n \geq n_0$, all $\delta \in (0, 1)$, and all unit vectors $v \in \mathbb{R}^d$, with probability at least $1 - \delta$

$$\max_{i \in \{1, \dots, n\}} \{v^\top x_i\} \in [X_{\max}(n, \delta, d) - c(n, \delta, d), X_{\max}(n, \delta, d) + C(n, \delta, d)]$$

and $\lim_{n \rightarrow \infty} C(n, \delta, d) = 0$, $\lim_{n \rightarrow \infty} c(n, \delta, d) = 0$.

Data Model Labels $y \in \{-1, 1\}$ and protected attribute $a \in \{-1, 1\}$ define four groups whose class-conditional distributions share the same shape $D(\cdot)$ but have different means:

$$x \mid (y, a) \sim D(y\mu + ya\psi),$$

where $\mu, \psi \in \mathbb{R}^2$ and $\mu \perp \psi$, with $\hat{\mu} = \mu / \|\mu\|$ and $\hat{\psi} = \psi / \|\psi\|$. For concreteness, take $\mu = \|\mu\| (0, 1)^\top$ and $\psi = \|\psi\| (1, 0)^\top$. Without loss of generality, let the majority attribute be $a_M = +1$ (the minority is $a_m = -1$). Thus the two majority means lie on the positive diagonal $\pm(\mu + \psi)$ and the two minority means on the negative diagonal $\pm(\mu - \psi)$.

Let A^M, B^M, A^m, B^m denote centered samples from $D(0)$, where A indexes examples with original label $+1$, and B indexes examples with original label -1 . We analyze the SVM objective in signed coordinates. Thus, before relabeling, majority signed examples are centered at $\mu + \psi$, while minority signed examples are centered at $\mu - \psi$.

The majority examples are not relabeled. For the minority examples, write

$$A^m = A^{m,+} \cup A^{m,-}, \quad B^m = B^{m,+} \cup B^{m,-},$$

where the superscript denotes the training label after relabeling. Hence $A^{m,+}$ and $B^{m,-}$ keep their original labels, while $A^{m,-}$ and $B^{m,+}$ are flipped.

If the minority were absent, the ERM SVM converges to the spurious direction

$$w_{\text{spu}}^{\text{maj}} \propto \mu + \psi.$$

Proposition C.3. *Suppose $D(0)$ satisfies Assumption C.2 and*

$$X_{\max}(p, \delta, 2) - X_{\max}(m, \delta, 2) \geq 2 \max\{\|\mu\|, \|\psi\|\} + c(p, \delta, 2) + C(m, \delta, 2). \quad (14)$$

Then, for any (possibly adversarial) relabeling of minority (which preserves separability) of training examples, if $p \rightarrow \infty$, with probability at least $1 - 4\delta$, gradient descent of ERM with logistic loss converges to the same spurious solution $w_{\text{spu}}^ \propto \mu + \psi$. Under a centrally symmetric $D(0)$ and when $\|\mu\| = \|\psi\|$, this limit satisfies $\text{EqOd}(w_{\text{spu}}^*) \rightarrow 0.5$.*

Proof. As in Chaudhuri et al. (2023), the limiting direction of gradient descent on the logistic loss over separable data is the hard-margin SVM direction. It is therefore enough to characterize the SVM direction.

Let

$$w_{\alpha,\sigma} = \alpha \hat{\mu} + \sigma \gamma \hat{\psi}, \quad \gamma = \sqrt{1 - \alpha^2}, \quad \alpha \in [-1, 1], \quad \sigma \in \{-1, +1\}.$$

We write γ instead of β to avoid confusing this coefficient with a population fraction. For fixed (α, σ) , define the terms

$$\begin{aligned} M_A^\sigma(\alpha) &:= \sup_{u \in A^M} -w_{\alpha,\sigma}^\top u, & M_B^\sigma(\alpha) &:= \sup_{u \in B^M} w_{\alpha,\sigma}^\top u, \\ R_{A,+}^\sigma(\alpha) &:= \sup_{u \in A^{m,+}} -w_{\alpha,\sigma}^\top u, & R_{A,-}^\sigma(\alpha) &:= \sup_{u \in A^{m,-}} w_{\alpha,\sigma}^\top u, \\ R_{B,+}^\sigma(\alpha) &:= \sup_{u \in B^{m,+}} -w_{\alpha,\sigma}^\top u, & R_{B,-}^\sigma(\alpha) &:= \sup_{u \in B^{m,-}} w_{\alpha,\sigma}^\top u. \end{aligned}$$

We use the convention that the supremum over an empty relabeled subset is $-\infty$.

The negative-label side of the SVM objective contains the majority negative examples, the flipped $A^{m,-}$ examples, and the kept $B^{m,-}$ examples. The positive-label side contains the majority positive examples, the kept $A^{m,+}$ examples, and the flipped $B^{m,+}$ examples. Therefore the SVM objective can be written as

$$\begin{aligned} F_\sigma(\alpha) = \max & \left\{ M_B^\sigma(\alpha) - \alpha \|\mu\| - \sigma \gamma \|\psi\|, \right. \\ & R_{A,-}^\sigma(\alpha) + \alpha \|\mu\| - \sigma \gamma \|\psi\|, \\ & \left. R_{B,-}^\sigma(\alpha) - \alpha \|\mu\| + \sigma \gamma \|\psi\| \right\} \\ & + \max \left\{ M_A^\sigma(\alpha) - \alpha \|\mu\| - \sigma \gamma \|\psi\|, \right. \\ & R_{A,+}^\sigma(\alpha) - \alpha \|\mu\| + \sigma \gamma \|\psi\|, \\ & \left. R_{B,+}^\sigma(\alpha) + \alpha \|\mu\| - \sigma \gamma \|\psi\| \right\}. \end{aligned}$$

By Assumption C.2, for the majority samples of size p , there exist X_p, c_p, C_p such that, with probability at least $1 - 4\delta$, uniformly over α and σ ,

$$M_A^\sigma(\alpha), M_B^\sigma(\alpha) \in [X_p - c_p, X_p + C_p].$$

For the minority samples, every relabeled subset is contained in the original minority sample, so

$$R_{A,+}^\sigma(\alpha), R_{A,-}^\sigma(\alpha), R_{B,+}^\sigma(\alpha), R_{B,-}^\sigma(\alpha) \leq X_m + C_m,$$

where

$$X_m := X_{\max}(m, \delta, 2), \quad C_m := C(m, \delta, 2).$$

Using Equation (14), we obtain, uniformly over α and σ ,

$$M_A^\sigma(\alpha) - R^\sigma(\alpha) \geq 2 \max\{\|\mu\|, \|\psi\|\}, \quad M_B^\sigma(\alpha) - R^\sigma(\alpha) \geq 2 \max\{\|\mu\|, \|\psi\|\},$$

for every term R^σ appearing above.

We now show that the majority entries attain both inner maxima. Consider the first maximum. Comparing the majority term with the flipped $A^{m,-}$ term,

$$\begin{aligned} & (M_B^\sigma(\alpha) - \alpha \|\mu\| - \sigma \gamma \|\psi\|) - (R_{A,-}^\sigma(\alpha) + \alpha \|\mu\| - \sigma \gamma \|\psi\|) \\ &= M_B^\sigma(\alpha) - R_{A,-}^\sigma(\alpha) - 2\alpha \|\mu\| \\ &\geq 2 \|\mu\| - 2\alpha \|\mu\| \geq 0. \end{aligned}$$

Comparing the majority term with the kept $B^{m,-}$ term,

$$\begin{aligned} & (M_B^\sigma(\alpha) - \alpha \|\mu\| - \sigma \gamma \|\psi\|) - (R_{B,-}^\sigma(\alpha) - \alpha \|\mu\| + \sigma \gamma \|\psi\|) \\ &= M_B^\sigma(\alpha) - R_{B,-}^\sigma(\alpha) - 2\sigma \gamma \|\psi\| \\ &\geq 2 \|\psi\| - 2\gamma \|\psi\| \geq 0. \end{aligned}$$

Thus the first maximum is attained by the majority term. The same comparisons show that the second maximum is also attained by the majority term. Hence

$$F_\sigma(\alpha) = M_A^\sigma(\alpha) + M_B^\sigma(\alpha) - 2\alpha \|\mu\| - 2\sigma\gamma \|\psi\|.$$

Write

$$E_\sigma(\alpha) := M_A^\sigma(\alpha) + M_B^\sigma(\alpha) - 2X_p.$$

The concentration condition gives

$$E_\sigma(\alpha) \in [-2c_p, 2C_p]$$

uniformly over α and σ . Therefore

$$F_+(\alpha) = 2X_p + E_+(\alpha) - 2(\alpha \|\mu\| + \gamma \|\psi\|),$$

whereas

$$F_-(\alpha) = 2X_p + E_-(\alpha) - 2(\alpha \|\mu\| - \gamma \|\psi\|).$$

For the $\sigma = +1$ branch, $F_+(\alpha)$ is minimised at

$$\alpha_g = \frac{\|\mu\|}{\sqrt{\|\mu\|^2 + \|\psi\|^2}},$$

By Lemma C.4, the minimizer of F_+ converges to α_g as $p \rightarrow \infty$.

Lemma C.4 (Approximate Maximization Lemma - I, Lemma 14 from Chaudhuri et al. (2023)). *Let $F(\alpha) = f(\alpha) + g(\alpha)$ where $g(\alpha) = \alpha u + \sqrt{1 - \alpha^2}v$, $u, v > 0$, and $f(\alpha) \in [-L, U]$. Let $\alpha_F \in \operatorname{argmax}_\alpha F(\alpha)$, and let $\alpha_g = \frac{u}{\sqrt{u^2 + v^2}} \in \operatorname{argmax}_\alpha g(\alpha)$.*

Then, the angle between $(\alpha_F, \sqrt{1 - \alpha_F^2})$ and $(\alpha_g, \sqrt{1 - \alpha_g^2})$ is at most $\cos^{-1}\left(1 - \frac{L+U}{\sqrt{u^2 + v^2}}\right)$, and $\max_\alpha F(\alpha) \geq \sqrt{u^2 + v^2} - L$.

For the $\sigma = -1$ branch, minimizing $F_-(\alpha)$ is asymptotically equivalent to maximizing

$$\alpha \|\mu\| - \sqrt{1 - \alpha^2} \|\psi\|.$$

This quantity is at most $\|\mu\|$. By contrast,

$$\max_\alpha \left[\alpha \|\mu\| + \sqrt{1 - \alpha^2} \|\psi\| \right] = \sqrt{\|\mu\|^2 + \|\psi\|^2} > \|\mu\|$$

whenever $\psi \neq 0$. Since $c_p, C_p \rightarrow 0$, the $\sigma = +1$ branch has the strictly smaller asymptotic SVM objective. Hence the optimal sign is $\sigma = +1$. Therefore,

$$w^* \rightarrow \alpha_g \hat{\mu} + \gamma_g \hat{\psi} = \frac{\mu + \psi}{\sqrt{\|\mu\|^2 + \|\psi\|^2}} \propto \mu + \psi.$$

Thus, independently of how the minority samples were relabeled in training, the limiting ERM direction is the same majority spurious direction

$$w_{\text{spu}}^* \propto \mu + \psi.$$

Under a centrally symmetric $D(0)$ and if $\|\mu\| = \|\psi\|$, the majority group ($a = +1$) is separated in the limiting spurious direction, while the minority signed means lie along $\mu - \psi$. Since

$$(\mu + \psi)^\top (\mu - \psi) = \|\mu\|^2 - \|\psi\|^2 = 0,$$

the limiting classifier cuts the two minority classes symmetrically. Hence

$$\text{TPR}_{a=+1} \rightarrow 1, \quad \text{FPR}_{a=+1} \rightarrow 0, \quad \text{TPR}_{a=-1} \rightarrow \frac{1}{2}, \quad \text{FPR}_{a=-1} \rightarrow \frac{1}{2}.$$

Consequently,

$$\text{EqOd}(w_{\text{spu}}^*) \rightarrow 0.5.$$

Finally, since the data are linearly separable under the stated assumptions, gradient descent on the logistic loss converges in direction to the hard-margin SVM solution (Soudry et al., 2018; Ji & Telgarsky, 2019). Hence the same limiting direction holds for ERM with logistic loss. \square

This result can also be extended to \mathbb{R}^d using techniques similar to those in Chaudhuri et al. (2023). This result also encompasses the 4-Gaussian mixture model $\mathbb{P}_{4\text{GMM}}$ used in Section 5 as a special case, leading to the following.

Proposition 5.1 (Informal). *Consider $\mathbb{P}_{4\text{GMM}}$, where every minority point participates in the ACA by flipping all $y=0$ labels to $y=1$. Then, under sufficiently separable clusters, with high probability, the EqOd of the ERM classifier minimizing the logistic loss will asymptotically approach 0.5.*

C.3 Success Bound With Label Error

The following proof uses Lemma 11 from Hardt et al. (2023).

Lemma C.5 (Lemma 11 from Hardt et al. (2023)). *Suppose that P, P' are two distributions such that $\text{TV}(P, P') \leq \epsilon$. Take any two events E_1, E_2 measurable under P, P' . If $P(E_1) > P(E_2) + \frac{\epsilon}{1-\epsilon}$, then $P'(E_1) > P'(E_2)$.*

Proposition 5.2. *With an algorithm $\mathcal{A}(x)$ with label error $\rho < 1/2$, the success of the collective is bounded by*

$$S(\alpha) \geq 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha} \tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}. \quad (11)$$

Proof. This proof follows closely the proof of Theorem 5 by Hardt et al. (2023). We start under the assumption of an optimal Bayes classifier, setting $\epsilon = 0$.

When the new label y' is wrong with probability ρ , then we can think of the collective as being union of two sub-collectives: one with the correct label and one with the incorrect label. In the binary case this can be formulated with correct subcollective P^+ as having label $y' = \arg \max_y P_0(y|g(x))$ and the incorrect subcollective P^- as with label $y' = \arg \min_y P_0(y|g(x))$. Then we can write the train distribution as

$$\begin{aligned} P_\alpha &= \alpha (\rho P^- + (1-\rho) P^+) + (1-\alpha) P_0 \\ &= \alpha \rho P^- + (1-\rho) \alpha P^+ + (1-\alpha) P_0. \end{aligned} \quad (15)$$

Denote $y^*(x) = \arg \max_y P_0(y|g(x))$, then the probability to get prediction y^* is

$$\begin{aligned} P_\alpha(y^*|x) &= \alpha \rho P^-(y^*|x) + (1-\rho) \alpha P^+(y^*|x) + (1-\alpha) P_0(y^*|x) \\ &= (1-\rho) \alpha + (1-\alpha) P_0(y^*|x), \end{aligned} \quad (16)$$

and the probability to get the prediction $y \neq y^*$ is

$$\begin{aligned} P_\alpha(y|x) &= \alpha \rho P^-(y|x) + (1-\rho) \alpha P^+(y|x) + (1-\alpha) P_0(y|x) \\ &= \alpha \rho + (1-\alpha) P_0(y|x), \end{aligned} \quad (17)$$

where $P^+(y^*|x) = 1$, $P^-(y^*|x) = 0$, $P^+(y \neq y^*|x) = 0$, $P^-(y \neq y^*|x) = 1$ by definition.

A Bayes classifier h returns the most probable label $h(x) = \arg \max_y P(y|x)$. Therefore, a Bayes classifier will output y^* if the probability is greater, which can be written as the condition

$$\begin{aligned} P_\alpha(y^*|x) &> P_\alpha(y|x) \\ (1-\rho)\alpha + (1-\alpha)P_0(y^*|x) &> \alpha\rho + (1-\alpha)P_0(y|x) \\ (1-2\rho)\alpha &> (1-\alpha)(P_0(y|x) - P_0(y^*|x)). \end{aligned} \quad (18)$$

Let $\tau(x) = \max_y [P_0(y|x) - P_0(y|g(x))]$, then

$$\begin{aligned} P_0(y|x) - P_0(y^*|x) &\leq P_0(y|x) - P_0(y|g(x)) + P_0(y^*|g(x)) - P_0(y^*|x) \\ &\leq 2\tau(x). \end{aligned} \quad (19)$$

With that, the condition in Eq. (18) can be written as

$$(1-2\rho)\alpha > 2(1-\alpha)\tau(x). \quad (20)$$

With that, the success can be bounded as

$$\begin{aligned} S &= P_0[f(x) = f(g(x))] \\ &= P_0[f(x) = y^*(x)] \\ &\geq P_0[(1-2\rho)\alpha > 2(1-\alpha)\tau(x)] \\ &= P_0\left[1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau(x) > 0\right] \\ &= \mathbb{E}_{x \sim P_0}\left[\mathbf{1}\left\{1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau(x) > 0\right\}\right] \\ &\geq \mathbb{E}_{x \sim P_0}\left[1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau(x)\right] \\ &= 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau \end{aligned} \quad (21)$$

With sub-optimality $\epsilon > 0$ A result of Lemma C.5 is to write the condition in Eq. (20) as

$$(1-2\rho)\alpha > 2(1-\alpha)\tau(x) + \frac{\epsilon}{1-\epsilon}, \quad (22)$$

which by following the same steps as with $\epsilon = 0$ results in the final bound

$$S(\alpha) \geq 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}. \quad (23)$$

□

C.4 Label Error With Better Representation

For the following we assume a similar setting as in Appendix C.2, visualised as a 2D distribution in Fig. 6. We are given the majority data, and tasked with labeling the minority data. Assume all labels are distributed equally $\mathbb{P}[Y = 1] = \mathbb{P}[Y = -1] = \frac{1}{2}$. The minority features X_{\min} are distributed as $X_{\min} \sim \mathcal{N}(y\mu_{\min}, \Sigma_{\min})$ with $X_{\min} \in \mathbb{R}^d$. The label $\hat{y}_{1NN}^{(n)}$ is predicted according to a 1NN classifier from n majority samples $\mathcal{D}_n = (x_i, y_i)_{i=0}^n$. Majority samples with $y = +1$ are distributed as $X_+ \sim \mathcal{N}(\mu, \Sigma)$, and with $y = -1$ are distributed as $X_- \sim \mathcal{N}(-\mu, \Sigma)$.

Theorem C.6. *Assume that $\mu_{\min}^\top \Sigma^{-1} \mu < 0$. Further, consider the setting with $\Sigma_{\min} = I$, and the minority (i.e. test) distribution introduced above with $\mathbb{P}[Y = 1] = \mathbb{P}[Y = -1] = 0.5$ and $X_{\min} \sim \mathcal{N}(y\mu_{\min}, \Sigma_{\min})$.*

Then, there exists a projection $P \in \mathbb{R}^{d \times d}$ such that asymptotically for $n \rightarrow \infty$, $err_{1NN}^{rep} < err_{1NN}^{raw}$.

Proof. Consider the projection on the hyperplane perpendicular to w , where $w = \frac{\mu - \mu_{\min}}{2}$. The assumption $\mu_{\min}^\top \Sigma^{-1} \mu < 0$ implies $\mu \neq \mu_{\min}$, so $w \neq 0$ and the projection matrix is well-defined as $P = I - \frac{ww^\top}{w^\top w}$. Since P projects onto the subspace orthogonal to $\mu - \mu_{\min}$, it satisfies

$$P\mu = P\mu_{\min}. \quad (24)$$

We apply Lemma C.7 to obtain closed forms for the asymptotic error of 1NN applied to the initial representation and to the features after the projection P . Namely, using the notation $v := \Sigma^{-1}\mu$ we have:

$$\text{err}_{1\text{NN}} = \frac{1}{2} \mathbb{P}_{X_{\min}|y=1}[\hat{y}_{1\text{NN}} = -1] + \frac{1}{2} \mathbb{P}_{X_{\min}|y=-1}[\hat{y}_{1\text{NN}} = 1] \quad (25)$$

$$= \frac{1}{2} \left(1 - \Phi \left(\frac{v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}} \right) \right) + \frac{1}{2} \Phi \left(\frac{-v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}} \right) \quad (26)$$

$$= 1 - \Phi \left(\frac{v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}} \right) \quad (27)$$

$$= 1 - \Phi(\text{SNR}), \quad (28)$$

where we used the fact that $\Phi(-z) = 1 - \Phi(z)$ and we denote $\text{SNR} := \frac{v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}}$.

Since the numerator is negative by assumption and the denominator is positive, $\text{SNR} < 0$, and therefore $\text{err}_{1\text{NN}}^{\text{raw}} > 1/2$.

We now analyze the projected representation using coordinates on the image of P . Let $U \in \mathbb{R}^{d \times r}$ be an orthonormal basis for $\text{Im}(P)$, where $r = \text{rank}(P)$, so $U^\top U = I_r$ and $UU^\top = P$. The projected majority distributions in coordinates $z = U^\top x$ are

$$Z_+ \sim \mathcal{N}(\eta, \Sigma_U), \quad Z_- \sim \mathcal{N}(-\eta, \Sigma_U), \quad (29)$$

where $\eta := U^\top \mu$ and $\Sigma_U := U^\top \Sigma U$. The projected minority distribution is

$$Z_{\min} | y \sim \mathcal{N}(y\eta, I_r), \quad (30)$$

because $U^\top \mu_{\min} = U^\top \mu$ and $\Sigma_{\min} = I$.

If $\eta = 0$, the projected representation contains no label signal and the balanced-label error is $\text{err}_{1\text{NN}}^{\text{rep}} = 1/2$, which is strictly smaller than $\text{err}_{1\text{NN}}^{\text{raw}}$. Otherwise, applying Lemma C.7 in the projected coordinate system gives

$$\text{err}_{1\text{NN}}^{\text{rep}} = 1 - \Phi(\text{SNR}_{\text{proj}}), \quad \text{SNR}_{\text{proj}} = \frac{\eta^\top \Sigma_U^{-1} \eta}{\sqrt{(\Sigma_U^{-1} \eta)^\top (\Sigma_U^{-1} \eta)}}. \quad (31)$$

Since Σ_U is positive definite and $\eta \neq 0$, the numerator is positive and the denominator is positive, so $\text{SNR}_{\text{proj}} > 0$. Therefore $\text{err}_{1\text{NN}}^{\text{rep}} < 1/2 < \text{err}_{1\text{NN}}^{\text{raw}}$. □

Lemma C.7. *For a unimodal minority distribution $X_{\min} \sim \mathcal{N}(\mu_{\min}, \Sigma_{\min})$ it holds that:*

$$\lim_{n \rightarrow \infty} \mathbb{P}_{X_{\min}}[\hat{y}_{1\text{NN}}^{(n)} = -1] = 1 - \Phi \left(\frac{v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}} \right),$$

where $v := \Sigma^{-1}\mu$ and Φ is the CDF of a standard Gaussian.

Proof. Let us denote $\hat{y}_{1NN} := \lim_{n \rightarrow \infty} \hat{y}_{1NN}^{(n)}$ and let p_+ and p_- be the densities of two class-conditional distribution. Notice that the two class conditional training distributions are supported on the entire domain of \mathbb{R}^d . Therefore, in the asymptotic regime, the label \hat{y}_{1NN} at a test point x is given according to the class-conditional distribution that has higher density. Namely, we have:

$$\hat{y}_{1NN} = \begin{cases} -1 & \text{if } p_+(x) < p_-(x), \\ 1 & \text{otherwise.} \end{cases}$$

Given $X_{\min} \sim \mathcal{N}(\mu_{\min}, \Sigma_{\min})$, we can then write the probability of predicting $\hat{y}_{1NN} = -1$ as:

$$\mathbb{P}_{X_{\min}}[\hat{y}_{1NN} = -1] = \mathbb{P}_{X_{\min}}[p_+(x) < p_-(x)].$$

Using the closed forms for the pdf of a Gaussian, we write the corresponding log-probabilities as follows:

$$\log p_+(x) = -\frac{1}{2}(x - \mu)^\top \Sigma^{-1}(x - \mu) + \text{const.}$$

$$\log p_-(x) = -\frac{1}{2}(x + \mu)^\top \Sigma^{-1}(x + \mu) + \text{const.}$$

Using the fact that log is monotonically increasing and Σ (and by extension Σ^{-1}) is a symmetric matrix, we can write after some simple calculations:

$$\mathbb{P}_{X_{\min}}[\hat{y}_{1NN} = -1] = \mathbb{P}_{X_{\min}}[\mu^\top \Sigma^{-1} x < 0].$$

Let us denote the random variable $Z := (\Sigma^{-1}\mu)^\top X$. Since Z is a linear transformation of Gaussian random variable, it is itself Gaussian and we can write its mean and variance as follows:

$$\mu_Z := v^\top \mu_{\min}, \text{ and } \sigma_Z^2 := v^\top \Sigma_{\min} v, \text{ where } v := \Sigma^{-1}\mu.$$

After this change of variable, we can rewrite the probability of predicting $\hat{y}_{1NN} = -1$ as:

$$\begin{aligned} \mathbb{P}_{X_{\min}}[\hat{y}_{1NN} = -1] &= \mathbb{P}_Z[Z < 0] \\ &= \Phi\left(\frac{0 - \mathbb{E}[Z]}{\sqrt{\text{Var}[Z]}}\right) \\ &= \Phi\left(\frac{-v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}}\right) \\ &= 1 - \Phi\left(\frac{v^\top \mu_{\min}}{\sqrt{v^\top \Sigma_{\min} v}}\right). \end{aligned}$$

□

Note that the error from Theorem C.6 is defined the same as ρ (Eq. (10)). This leads to the following.

Proposition 5.3 (Informal). *Let data be drawn from \mathbb{P}_{4GMM} , and ρ_{plain} denote the error of a 1-NN classifier that assigns the label of the nearest majority neighbor in the original feature space. Then there exists a fair representation in which a 1-NN classifier achieves error ρ_{FRL} such that, asymptotically with respect to the dataset size, $\rho_{\text{FRL}} \leq \rho_{\text{plain}}$.*

D Technical Details

D.1 Datasets

COMPAS The Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) dataset contains the data of criminal defendants in Broward county sheriff’s office in Florida with the task of predicting the recidivism risk. The label in this dataset represents whether the person re-offended and the sensitive attribute is the race. We follow the same data cleaning and pre-processing as Alghamdi et al. (2022).

Adult The Adult dataset Becker & Kohavi (1996) contains demographic features of US citizens and is tasked with predicting the income level of an individual. The label represents if the individual has income higher than \$50,000 and the sensitive attribute we use is the race. We follow the same data cleaning and pre-processing as Alghamdi et al. (2022).

HLSL The High School Longitudinal Study of 2009 (HLSL) Jeong et al. (2022) contains details of high-school students across the US and the task is to predict the academic success of the students. The label represents the exam score and the sensitive attribute is the race. We follow the same data cleaning and pre-processing as Alghamdi et al. (2022).

ACS-Income Ding et al. (2021) offer different classification tasks derived by US census data. In our work we used the pre-defined task of predicting level of income denoted as *ACSIncome*, where the data is already pre-processed. The label represents if the individual has income higher than \$50,000 and the sensitive attribute is the race.

Waterbirds The waterbirds dataset Sagawa et al. (2020) contains images of landbirds and waterbirds super-imposed on either land or water backgrounds, with the task of classifying the image as of a landbird or a waterbird. The label represents the type of bird, and the sensitive attribute is whether the background is land or water. To obtain the features, we used the output of the penultimate layer of a pre-trained ResNet-18 network from *PyTorch*¹. We report the results on those features as Waterbirds-Full. We also performed PCA (using *scikit-learn*) and kept the first 85 principal components which retain about 75% of the variance, and report the results of these components as Waterbirds-PCA.

CivilComments The CivilComments dataset Borkan et al. (2019) is a collection of text comments found on the internet, with the goal of training a classifier to fairly detect toxicity. For this paper, we modified the dataset to keep only the comments that include either *christian* or *muslim* (but not both), with a label 0 meaning toxic and 1 meaning safe. To obtain the features, we used the word embeddings given by Hugging Face’s *bert-base-uncased* model². We report the results on those features as CivilComments-Full. We also performed PCA (using *scikit-learn*) and kept the first 100 principal components which retain about 75% of the variance, and report the results of these components as CivilComments-PCA.

D.2 Training

All classification experiments were trained with *scikit-learn*’s histogram-based gradient boosting classification tree with the default parameters³. When there was not a pre-defined test set, we set the train-test split as 80-20 before applying the collective action.

The probabilities for RB-prob were inferred by training *scikit-learn*’s histogram-based gradient boosting classification tree on the majority data with the default parameters, and using its *predict_proba* function. For LFR Zemel et al. (2013) we used the implementation in *Holistic AI*’s open source library⁴ with the default parameters. For FARE Jovanović et al. (2023) we used the official implementation⁵ with hyperparameters

¹<https://pytorch.org/vision/main/models/generated/torchvision.models.resnet18.html>

²<https://huggingface.co/google-bert/bert-base-uncased>

³scikit-learn.org/stable/modules/generated/sklearn.ensemble.HistGradientBoostingClassifier.html

⁴<https://github.com/holistic-ai/holisticai>

⁵<https://github.com/eth-sri/fare>

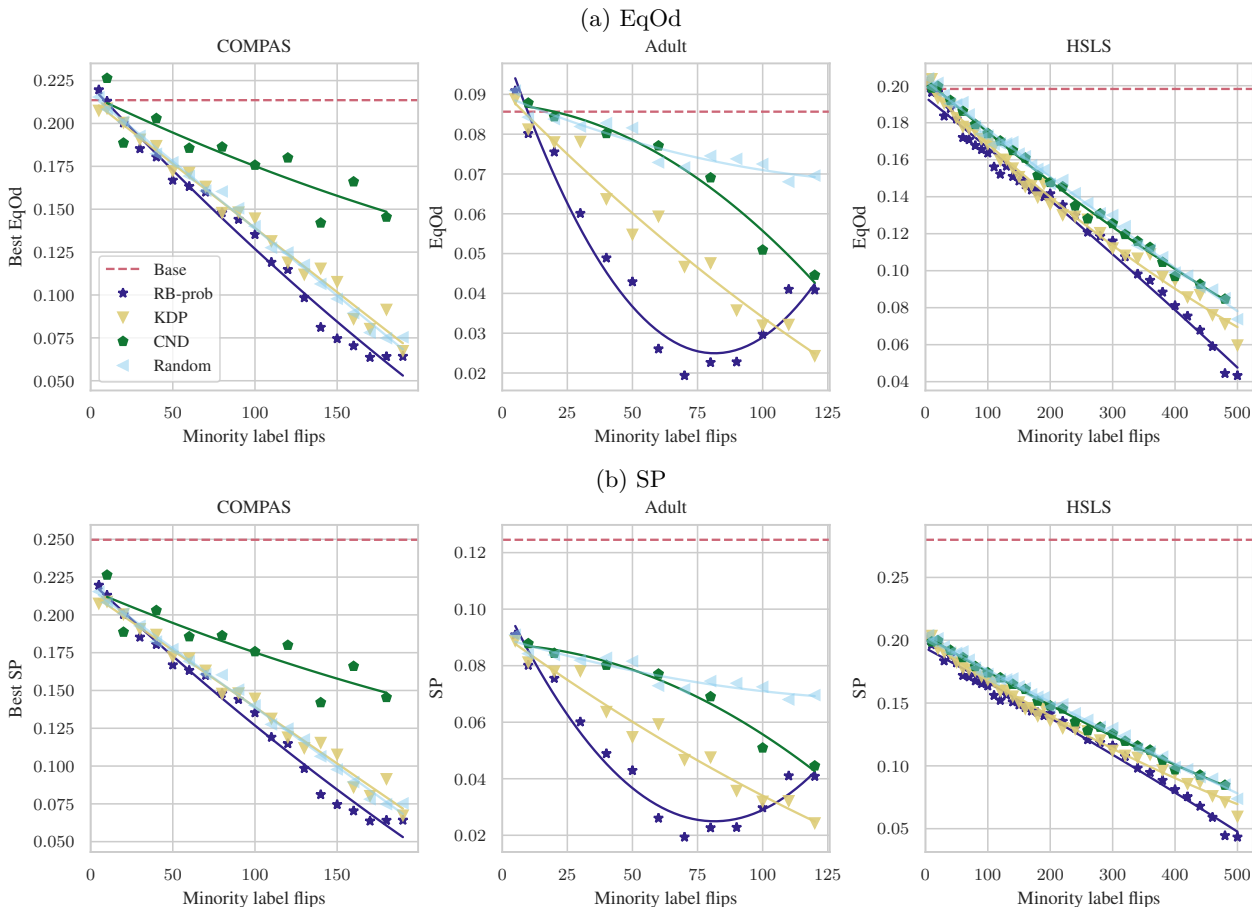


Figure 9: EqOd and SP violations per number of label flips for the random baseline, our method RB-prob, and the existing methods KDP (Luong et al., 2011) and CND (Kamiran & Calders, 2009). Our method is more efficient than prior work, requiring fewer flips to achieve comparable violation levels. Note that in this experiment CND could flip any label, while all other methods were restricted to the labels of 30% of the minority.

$\gamma = 0.85$, $k = 200$ and $n = 100$. For all distance computation we used the Euclidean norm ℓ^2 -norm as $d(v, u) = \|v - u\|_2 = \sqrt{\sum_i (v_i - u_i)^2}$.

E Additional Results

E.1 Comparison with prior work

We compare our method RB-prob with the existing methods KDP Luong et al. (2011) and CND Kamiran & Calders (2009) in Fig. 9. Note that CND requires flipping labels for both majority and minority members, and we report the total number of label flips. Fig. 9 shows that our method, motivated by the counterfactual labeling, is more efficient in terms of required number of label flips, than the existing works.

E.2 Expanded results

The following figures include the results of the experiments reported in the main text using all methods on all dataset, both with EqOd (Eq. (2)) and SP (Eq. (12)) as a measure of unfairness

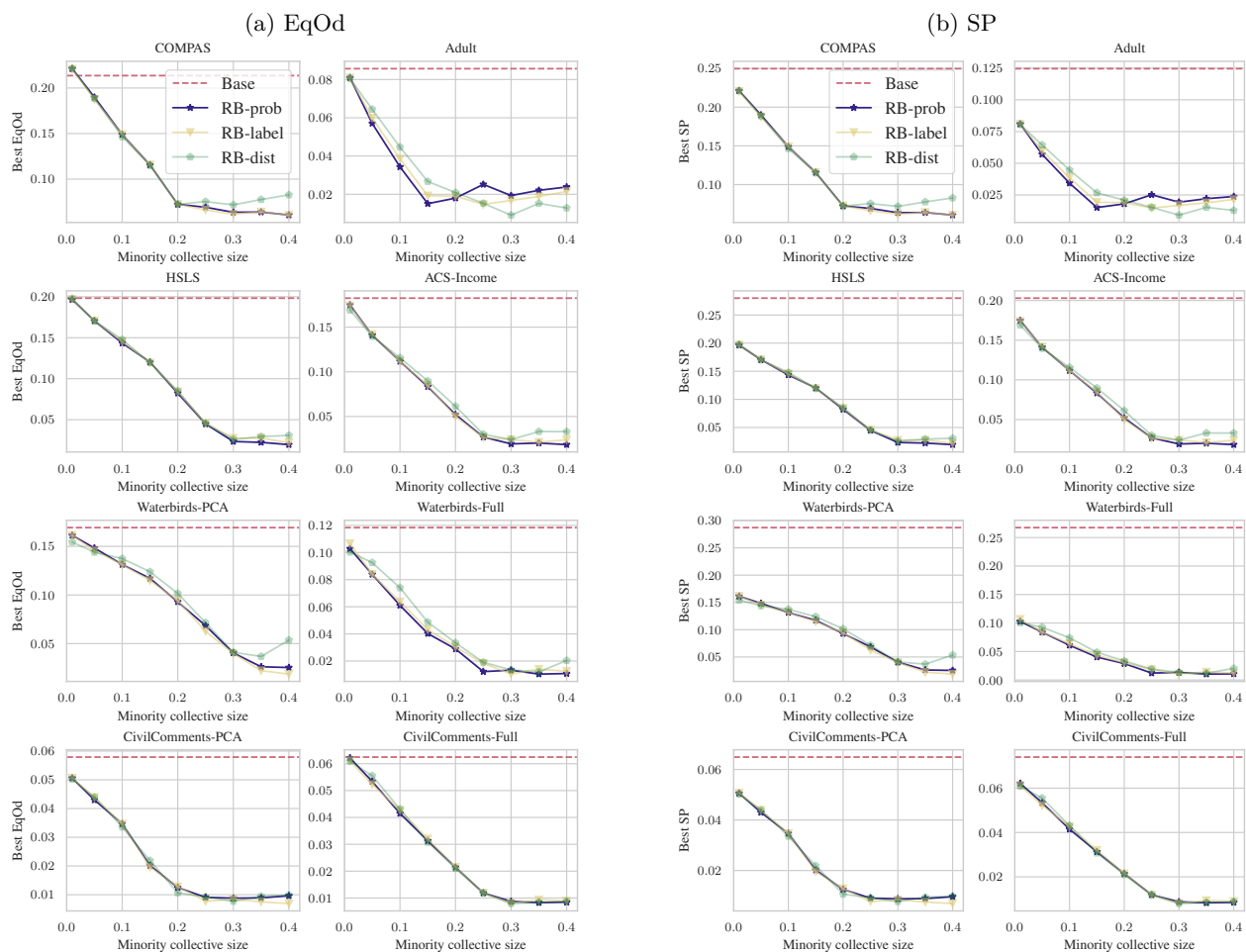


Figure 10: The lowest EqOd and SP violations a collective can achieve decrease as the collective size increases, up to a certain point. Each point is a mean of 10 runs, with the standard deviation being smaller than the markers. In all the datasets we experimented on, both violation metrics stabilize around $\alpha = 0.3$.

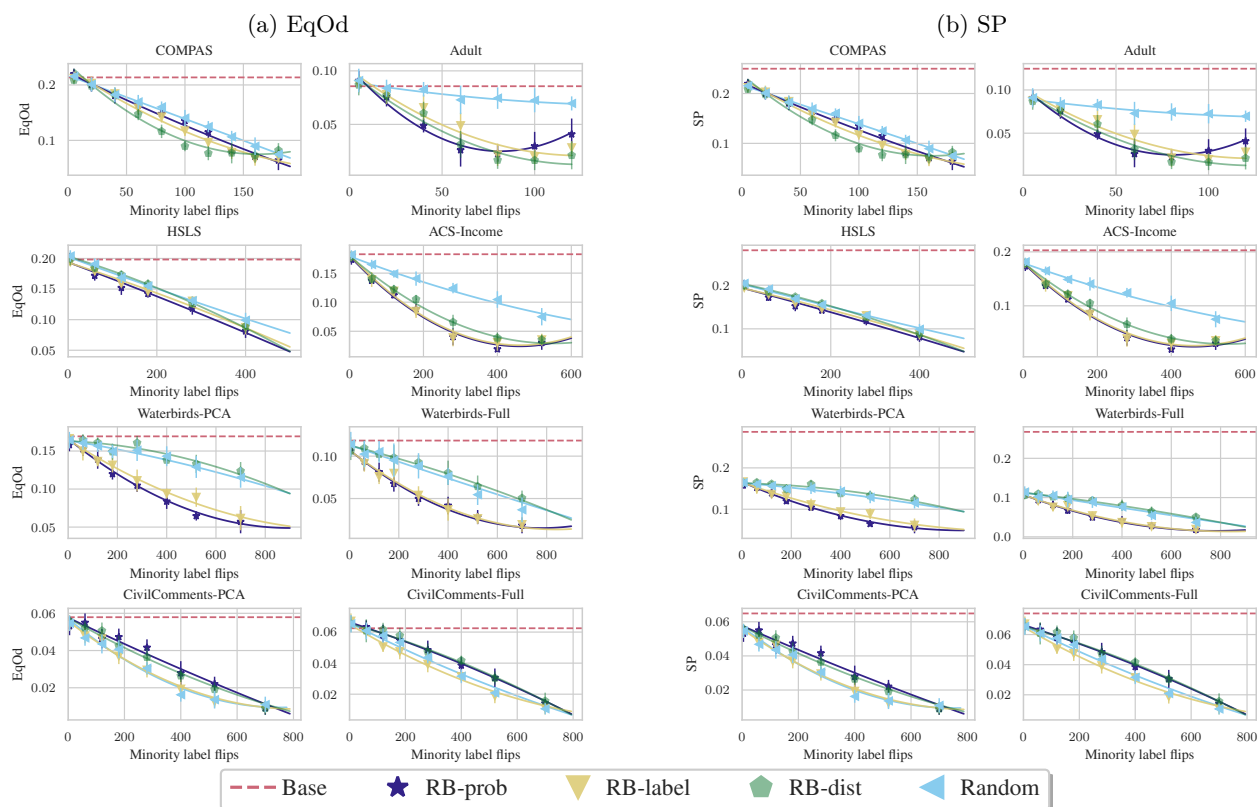


Figure 11: Our proposed methods are consistently more efficient than randomly flipping labels, requiring fewer label flips to attain comparable EqOd or SP violation levels. Each marker is the mean of 10 random runs with a specific number of label flips. The dashed line shows the corresponding mean violation for a classifier trained on the dataset without collective action.

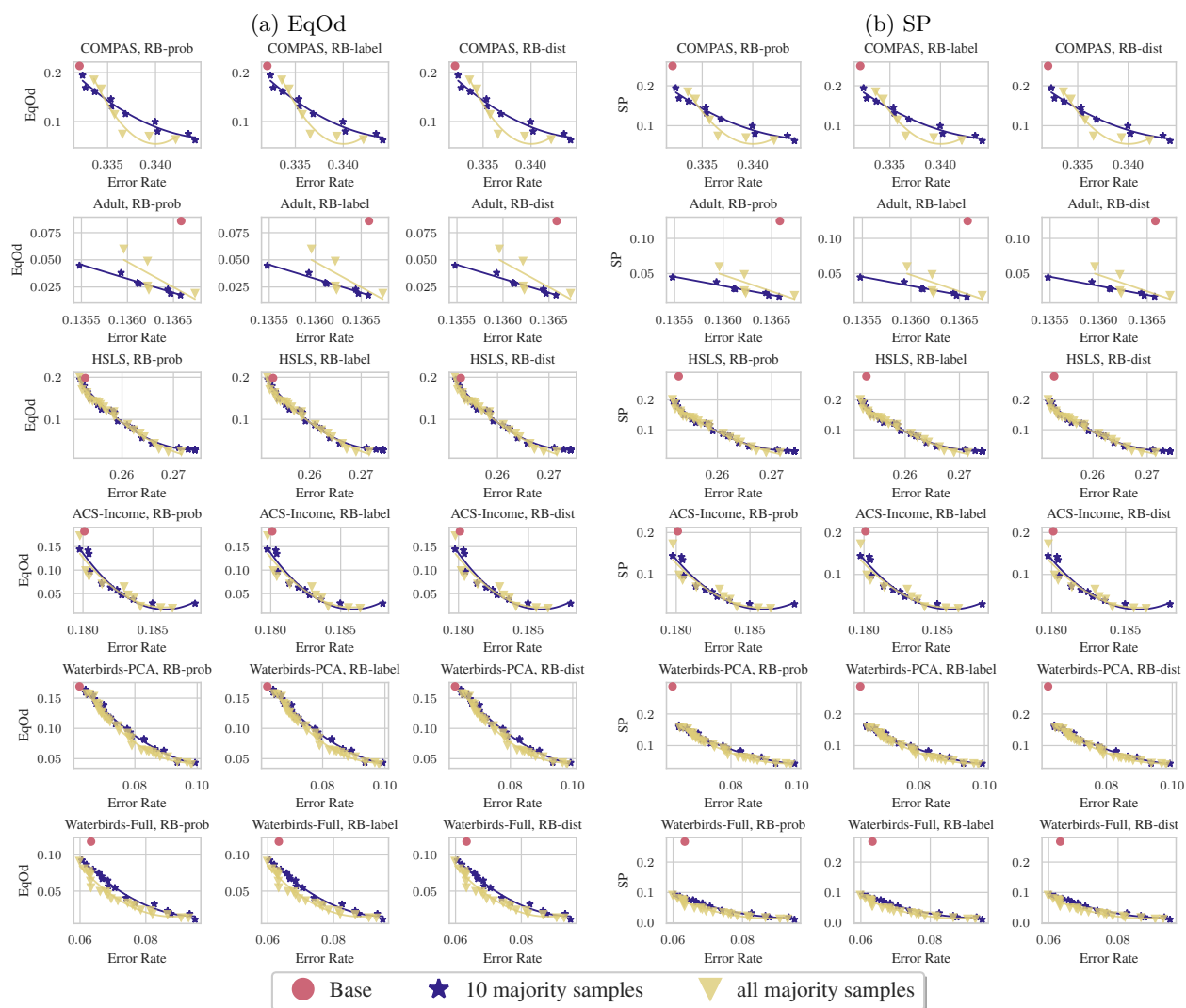


Figure 12: Limiting the knowledge of the collective about the majority does not significantly harm the Pareto front. Each point is the mean of 10 runs and the curves are fitted to guide the eye.

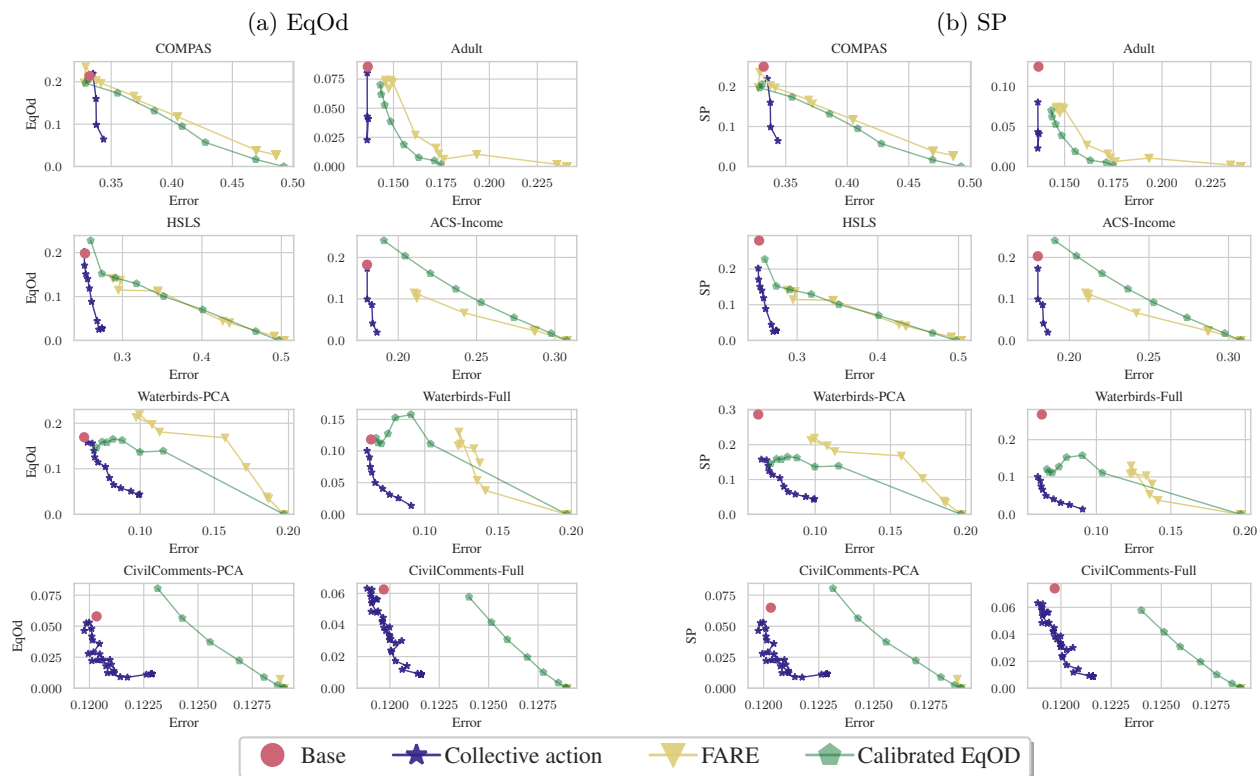


Figure 13: The EqOd panel shows that the firm-side pre-processing method FARE Jovanović et al. (2023) and the post-processing method calibrated equalized odds Pleiss et al. (2017) attain 0 EqOd with large error. Across both EqOd and SP, RB-prob with $\alpha = 0.3$ (Section 3) has much smaller error and lower violations than the base classifier, but does not reach zero violation.

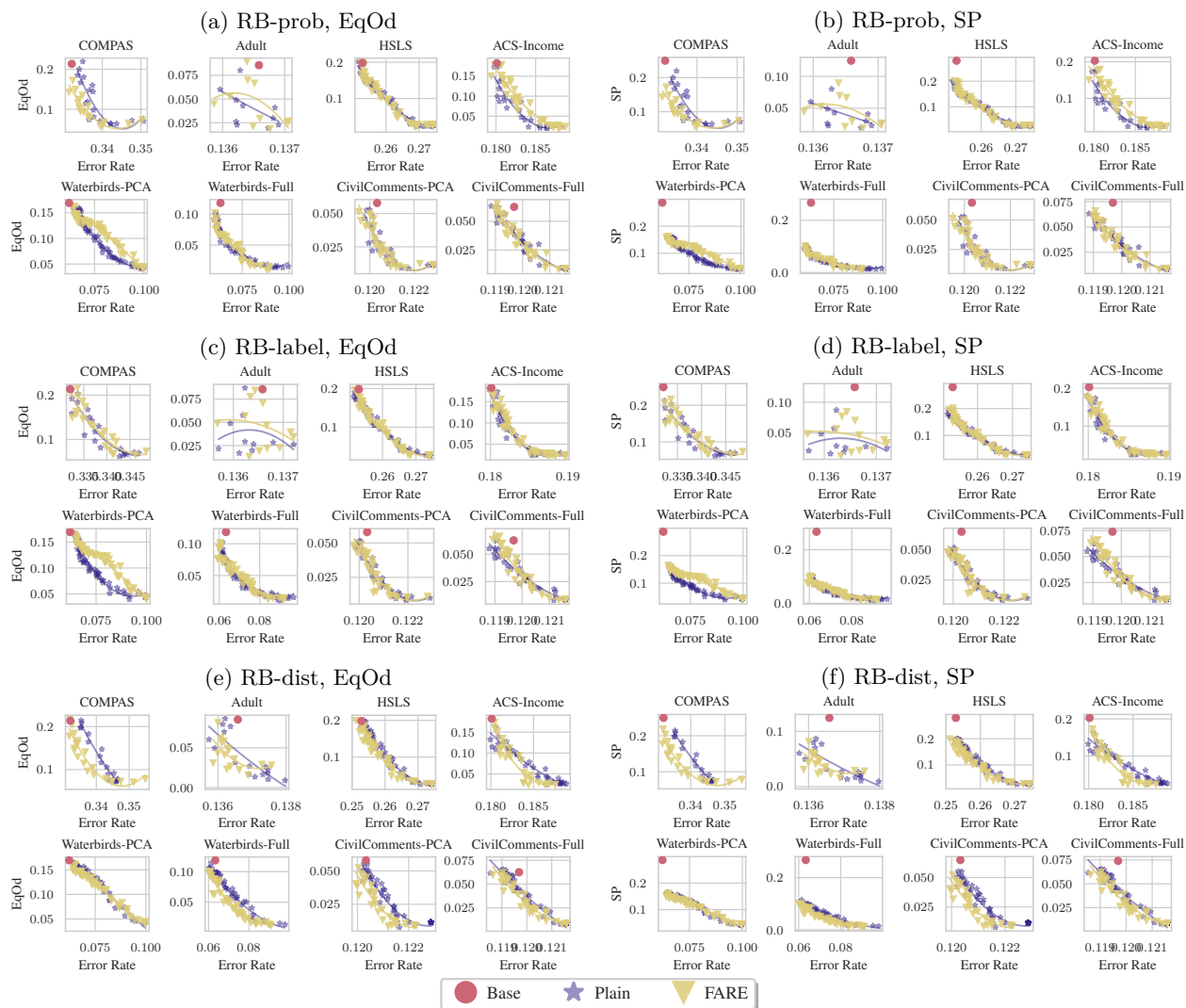


Figure 14: Using a fair representation generally shifts the Pareto fronts lower and left across RB-prob, RB-label, and RB-dist. The blue stars represent each method without transforming the data, and the yellow triangles represent each method after transforming the data using FARE (Jovanović et al., 2023). The lines are fitted by a polynomial of degree 2 to guide the eye.