Fairness for the People, by the People: Minority Collective Action

Anonymous Author(s)

Affiliation Address email

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, they 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, where a coordinated minority group strategically relabels its own data to enhance fairness, without altering the firm's training process. We propose three practical, model-agnostic methods to approximate ideal relabeling and validate them on real-world datasets. Our findings show that a subgroup of the minority can substantially reduce unfairness with a small impact on the overall prediction error.

1 Introduction

5

6

8

9

10

25

26

27 28

29 30

31

As machine learning (ML) tools become increasingly accessible, more firms deploy them for decisionmaking. However, ML models often perpetuate biases present in their data, leading to unfair outcomes across demographic groups [1]. Moreover, most fair-learning algorithms incur a non-negligible cost in accuracy or computational resources [2, 3, 4, 5], which can discourage practical adoption.

Since firms control the training pipeline, end-users lack access to these algorithms and cannot directly 16 enforce fair treatment. Yet, affected users routinely generate and share data — through clicks, ratings, 17 or other contributions — that is used to train the firm's models. Consequently, if underrepresented 18 minority groups collaboratively alter the data they share, they might be able to steer the learned model 19 towards fairer behavior, even without access to the firm's training pipeline. This idea is reminiscent of pre-processing fairness techniques [6, 7, 8, 9], which modify the data before model training. Unlike 21 22 these prior approaches, which assume centralized control over the data, we consider the setting 23 of algorithmic collective action [10, 11, 12, 13, 14], in which a small group of users strategically modifies their own data to influence the correlations learned by the model. 24

We adapt the *erasure strategy* from Hardt et al. [10] to reduce correlation between group membership and label by relabeling minority samples. The collective is restricted to the minority, because minority members are more motivated to join minority collective action [15, 16], and majority-group users may be less inclined to disrupt the status quo. We show that when a classifier is trained on data affected by this form of collective action, standard fairness metrics (e.g., demographic parity, equalized odds) improve substantially. This improvement is illustrated in Figure 1, where a small collective of minority samples significantly reduces unfairness with minimal impact on prediction error.

The key obstacle in implementing the erasure strategy is that it requires knowledge of each user's label under a counterfactual group membership. Computing such counterfactual labels exactly would require access to an underlying causal model, which is typically infeasible in practice. To overcome this challenge, we propose three *model-agnostic* methods to estimate the counterfactual labels.

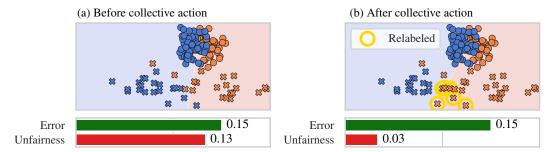


Figure 1: Minority-only collective action can substantially improve fairness. With only 6 label flips, the fairness violation of logistic regression goes down by over 75% with only a negligible increase in prediction error. Circles and crosses represent majority and minority points, respectively.

2 Collective Action for Fairness

We consider a setting in which 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 $h: \mathbb{R}^m \to \{0,1\}$ to minimize the prediction error, defined as

$$\operatorname{Error}(h) = \mathbb{P}\left[h\left(x\right) \neq y\right]. \tag{1}$$

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) [17, 18] and equalized odds (EqOd) [19]. In this work, we focus primarily on violations of EqOd, formally defined as

$$EqOd(h) = \frac{1}{2} \sum_{z=0.1} |\mathbb{P}[h(x) = 1 | a = 1, y = z] - \mathbb{P}[h(x) = 1 | a = 0, y = z]|, \qquad (2)$$

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

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. Appendix A suggests real-world scenarios where collective action can contribute to fairness. These collectives and their influential abilities in ML are studied as the field of algorithmic collective action [10].

In social sciences, *collective action* refers to the coordinated efforts of individuals working together to pursue a shared goal [20, 21]. Hardt et al. [10] adapt this notion to machine learning, 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 population that belongs to the collective.

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

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

the probability, under the base distribution, that the classifier's prediction remains unchanged after applying q to the features.

To achieve signal erasure, Hardt et al. [10] propose the collective relabels itself with the most likely label under the transformation g. Formally, the strategy is defined as

$$(x,y) \to \left(x, \arg\max_{y' \in \{0,1\}} \mathbb{P}_0\left(y'|g\left(x\right)\right)\right).$$
 (4)

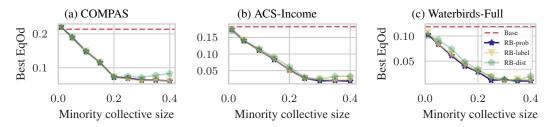


Figure 2: The lowest EqOd violation a collective can achieve greatly improves as the collective size increases, up to a certain point. Each point is a mean of 10 runs, with the standard deviation being smaller than the markers. In all the datasets we experimented on, the lowest EqOd violation converges around $\alpha = 0.3$. Additional results are presented in Figure 9 in the appendix.

Intuitively, if g is a feature pattern correlated with group membership (e.g., minority vs. majority), then achieving invariance under g promotes fairness by reducing the classifier's dependence on group-identifying information. We define q to be the counterfactual features a minority member would have had they belonged to the majority. Appendix B.3 describes in more detail the success of signal erasure and Appendix B.4 connects this action to fairness through counterfactual fairness.

Approximating the Counterfactual Label

70

71

72

85

86

87

88

89

90

91

92

93

94

This section describes how a minority collective can approximate a signal-erasure strategy to promote 73 fairness in practice. While the theory of signal erasure has been studied before [10, 14], prior work 74 lacks empirical evaluation. In this paper, we present the first practical algorithm for signal erasure 75 and provide experimental results in Section 4. As discussed in Appendix B.4, a suitable signal to 76 erase is $g(x) = x_{A \leftarrow 0}$, where each collective member relabels themselves according to Equation (4). 77 However, end-users lack access to the true causal model and cannot compute the counterfactual 78 labels directly. To address this limitation, we propose to assign each collective member i a score s_i , 79 which serves as a proxy for the likelihood that they would receive the label y=1 if they belonged 80 to the majority. Given a budget of M label flips, the collective selects the M members with the 81 highest scores; these individuals flip their labels from y=0 to y=1. The budget M controls the 82 accuracy-fairness tradeoff, where a higher budget typically leads to better fairness, but higher error 83 (Figure 6).

We introduce three model-agnostic scoring functions, each capturing a different notion of similarity to majority users:

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

$$s_i = f\left(x_i\right). \tag{5}$$

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:

$$s_i = \sum_{j \in K_i} \mathbf{1} \{ y_j = 1 \}.$$
 (6)

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:

$$s_i = -\frac{1}{k} \sum_{j \in K_i} \|x_i - x_j\|_2. \tag{7}$$

Intuitively, RB-prob assigns a higher score where a classifier trained solely on majority data predicts a higher likelihood of the label y = 1. RB-label scores collective members according to the frequency of y = 1 among their majority neighbors, while RB-dist prioritizes those who are closer majority

users with y=1. The similarity to prior work is discussed in Appendix D

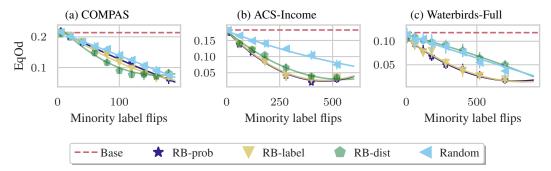


Figure 3: Our proposed methods are consistently more efficient than randomly flipping labels, requiring less label flips to attain the same level of EqOd. Each marker is the mean of 10 random runs with a specific number of label flips. The standard deviation is presented by the error bars. The dashed line shows the mean EqOd for a classifier trained on the dataset without collective action.

99 4 Experimental Results

115

116

117 118

119

120

121

122

123

124

125

126

127

128

129

We compare our methods against a random baseline that flips y=0 labels to y=1 for M 101 randomly selected collective members. We experimented on the tabular datasets COMPAS [22], Adult [23], HSLS [24], ACS-Income [25], the image dataset Waterbirds [26] and the text dataset 102 CivilComments [27]. For Waterbirds, we use features extracted from a pre-trained ResNet-18 (denoted 103 Waterbirds-Full) and for CivilComments, we used the extracted features from Hugging Face's pre-104 trained bert-base-uncased model (denoted CivilComments-Full). In addition to the complete features 105 of Waterbirds and CivilComments, we also include experiments on the PCA features, with 85 106 components for Waterbirds (denoted Waterbirds-PCA) and 100 components for CivilComments 107 (denoted CivilComments-PCA). Details on the datasets are provided in Appendix F.1. 108

All reported metrics are computed on a fixed test set, without any collective action, and averaged over 10 independent runs. In each run, we randomly selected a minority collective, which then applies one of the methods described in Section 3. 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. The complete set of results can be found in Appendix G.2, including an experiment with limited knowledge of the majoritry (Figure 11).

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, the experiments included 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 Figure 2, increasing α improves the best achievable EqOd 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 fairness. To evaluate efficiency, Figure 3 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., Figure 3a), or perform comparably to the random baseline in others (Figure 3c).

These results suggest that a well-chosen scoring function enables the collective to achieve a desired level of fairness with fewer label flips, reducing the "cost" of collective action 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 collective action for fairness, a point we elaborate on in Appendix C.

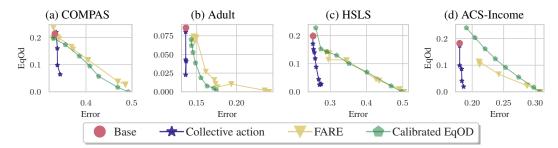


Figure 4: User-side method cannot achieve perfect fairness, while the firm-side pre-processing method FARE [28] and the post-processing method calibrated equalized odds [29] attain 0 EqOd with large error. However, RB-prob's fairness is better than the base classifier, with a smaller error than the firm-side methods.

5 Limitations of Minority Collective Action

Previous work on collective action 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. Those differences introduce limitations to the existing theory, which we analyze and theoretically quantify in this section and in Appendix C.

This restriction expresses 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 reduces the collective's impact. In Appendix E.2 we formally prove that there exists a case where a minority-only collective is unable to acheive perfect fairness.

We empirically corroborate this claim on real world datasets by examining the fairness-accuracy 146 tradeoff of several fair learning methods. Most of these methods include a hyperparameter that 147 controls this trade-off, yielding a set of pairs (Error, EqOd) as it varies. This set forms a Pareto front, 148 representing the best attainable trade-offs. A Pareto front is said to dominate another if it lies entirely 149 to the left (lower error) and below (lower unfairness) of the other. Figure 4 compares the Pareto fronts 150 of RB-prob, one of our minority collective action methods, with established firm-side methods. We 151 observe that the lowest fairness violation achievable by RB-prob is greater than that of the firm-side 152 approaches. However, the firm-side methods are able to arrive at perfect fairness only at a cost of 153 prohibitively high prediction error. But, inspecting the region where the error is small compared to 154 the base classifier, the fairness of RB-prob is comparable to that of the firm-side methods. 155

Appendix C additionaly discusses the limitations given that the counterfactual labels is onyl approximated and not exact and how the success bound is affected, and also how the same estimation error can be decreased by using representation learning.

6 Conclusion

159

160

161

162

163

164

165

166

This work demonstrates that user-side methods, specifically minority collective action, can effectively reduce unfairness in machine learning. While much of the existing fairness research focused on firm-side methods, paradoxically these often come at a cost that may not be worth to the firm. This catch emphasizes the importance of studying user-side approaches for bias mitigation. We also note that in general, collective action methods can be exploited by malicious parties seeking self-gain or harming other communities, and it is important to be discussing these limitations and possibly regulate them.

We introduce three practical methods that a collective can easily implement to relabel itself, and show empirically that collective action can considerably reduce unfairness in a variety of datasets, though not completely. Importantly, we also examine the limitations of a minority being composed of only minority members, and how the success is affected by approximating the counterfactual labels.

Overall, this paper shows a practical use case of collective action in the hopes of sparking further research into applications of collective action and user-side methods for social good.

References

173

- [1] Solon Barocas and Andrew D. Selbst. Big data's disparate impact. *California Law Review*, 104 (3):671–732, 2016.
- [2] 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*, pages 107–118. PMLR, 2018.
- [3] Han Zhao and Geoff Gordon. Inherent tradeoffs in learning fair representations. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [4] 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)*, pages 12037–12046, 2024.
- [5] Bashir Sadeghi, Sepehr Dehdashtian, and Vishnu Boddeti. On Characterizing the Trade-off in
 Invariant Representation Learning. Transactions on Machine Learning Research, 2022.
- [6] Faisal Kamiran and Toon Calders. Classifying without discriminating. In *Control and Communication 2009 2nd International Conference on Computer*, pages 1–6, 2009.
- [7] 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*, pages 502–510.
 Association for Computing Machinery, 2011. ISBN 978-1-4503-0813-7.
- [8] 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, pages 325–333. PMLR, 2013.
- [9] 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, pages 3384–3393. PMLR, 2018.
- [10] Moritz Hardt, Eric Mazumdar, Celestine Mendler-Dünner, and Tijana Zrnic. Algorithmic
 Collective Action in Machine Learning. In *Proceedings of the 40th International Conference* on Machine Learning, volume 202, pages 12570–12586, 2023.
- 201 [11] Omri Ben-Dov, Jake Fawkes, Samira Samadi, and Amartya Sanyal. The Role of Learning
 202 Algorithms in Collective Action. In *Proceedings of the 41st International Conference on*203 *Machine Learning*, volume 235, pages 3443–3461. PMLR, 2024.
- [12] Joachim Baumann and Celestine Mendler-Dünner. Algorithmic Collective Action in Recommender Systems: Promoting Songs by Reordering Playlists. In *Advances in Neural Information Processing Systems*, volume 37, pages 119123–119149. Curran Associates, Inc., 2024.
- [13] 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. ISBN 979-8-4007-1394-1.
- 211 [14] Etienne Gauthier, Francis Bach, and Michael I. Jordan. Statistical collusion by collectives on learning platforms. In *Forty-Second International Conference on Machine Learning*, 2025.
- 213 [15] 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.
- [16] 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.

- [17] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. Building Classifiers with Independency
 Constraints. In 2009 IEEE International Conference on Data Mining Workshops, pages 13–18,
 2009.
- [18] 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*, pages 214–226. Association for Computing Machinery, 2012. ISBN 978-1-4503-1115-1.
- [19] 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,
 pages 3323–3331, 2016.
- [20] Mancur Olson. Collective action. In *The Invisible Hand*, pages 61–69. Palgrave Macmillan UK,
 1989. ISBN 978-1-349-20313-0.
- [21] Gerald Marwell and Pamela Oliver. The Critical Mass in Collective Action. Cambridge
 University Press, 1993.
- [22] Jeff Mattu, Julia Larson, Lauren Angwin, and Surya Kirchner. How We Analyzed the COM PAS Recidivism Algorithm. https://www.propublica.org/article/how-we-analyzed-the-compas recidivism-algorithm, 2016.
- 237 [23] Barry Becker and Ronny Kohavi. Adult. UCI Machine Learning Repository, 1996.
- [24] 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.
- [25] 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,
 pages 6478–6490. Curran Associates, Inc., 2021.
- [26] Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally
 Robust Neural Networks. In *Eighth International Conference on Learning Representations*,
 2020.
- [27] 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, pages 491–500. Association for
 Computing Machinery, 2019. ISBN 978-1-4503-6675-5.
- [28] 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*, pages 15401–15420. PMLR, 2023.
- [29] 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.
- [30] 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.
- Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H. Chi, and Alex Beutel. Counter-factual fairness in text classification through robustness. In *Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society*, pages 219–226. Association for Computing Machinery, 2019. ISBN 978-1-4503-6324-2.
- 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*, pages 1438–1444. International Joint Conferences on Artificial Intelligence Organization, 2019.

- [33] 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, pages 34122–34138. Curran Associates, Inc., 2023.
- [34] 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, pages 4144–4188. PMLR, 2023.
- [35] 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, pages 60–69. PMLR, 2018.
- [36] 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.
- 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, pages 6781–6792, 2021.
- [38] 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, pages 38747–
 38760. Curran Associates, Inc., 2022.
- [39] 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, pages 48321–48343. PMLR,
 2024.
- ²⁹¹ [40] André Cruz and Moritz Hardt. Unprocessing seven years of algorithmic fairness. In *The Twelfth*²⁹² International Conference on Learning Representations, 2024.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: We describe the goals and contributions of the paper in the abstract and in the introduction.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The paper mentions the limitations in the discussion section.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach.
 For example, a facial recognition algorithm may perform poorly when image resolution
 is low or images are taken in low lighting. Or a speech-to-text system might not be
 used reliably to provide closed captions for online lectures because it fails to handle
 technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: In our theory and proofs we refer to the previous work and assumptions we rely on

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We describe the technical details of the experiments in the experiments section and appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
- (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The experiments we report are not difficult to recreate, and we provide the exact dataset and python packages we use. We are also planning to include the code in the supplementary material.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be
 possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not
 including code, unless this is central to the contribution (e.g., for a new open-source
 benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how
 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We describe the technical details of the experiments in the experiments section and appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: Each point in our figures is a mean of 10 random runs. Some of the figures report the error bars, but in all cases the averaged points present clear trends that support our claims with small standard deviation. For this reason we did not find it necessary to include error bars in most figures in the main text since in most cases they are smaller than the markers and would have made the figures less readable.

Guidelines:

The answer NA means that the paper does not include experiments.

- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how
 they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [No]

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

479

480

482

483

484 485

486

487

488

489

490

491

492

493

494

495 496

497

498

499

500

Justification: Our paper does not focus on algorithmic efficiency, but on classifier behavior. For this reason, we do not think incluidng the compute power is relevant.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our experiments used public datasets. The societal impact is discussed in the discussion section.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a
 deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: The societal impact is discussed in the discussion section.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal
 impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: This paper offers no new data or models.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with
 necessary safeguards to allow for controlled use of the model, for example by requiring
 that users adhere to usage guidelines or restrictions to access the model or implementing
 safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We use publicly available python packages, mention them by name and provide a URL to the functions and objects.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the
 package should be provided. For popular datasets, paperswithcode.com/datasets
 has curated licenses for some datasets. Their licensing guide can help determine the
 license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
 - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597 598

599

600

601

602

603

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: The paper does not release new assets.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

 The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

16. Declaration of LLM usage

Question: Does the paper describe the usage of LLMs if it is an important, original, or non-standard component of the core methods in this research? Note that if the LLM is used only for writing, editing, or formatting purposes and does not impact the core methodology, scientific rigorousness, or originality of the research, declaration is not required.

Answer: [NA]

Justification: The core method development in this research does not involve LLMs as any important, original, or non-standard components.

Guidelines:

- The answer NA means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
- Please refer to our LLM policy (https://neurips.cc/Conferences/2025/LLM) for what should or should not be described.

5 7 Appendix

626

629 630

631 632

633

634

635

637

639

640

643 644

645

646

647

648

650

651

652

653

654

655

656

657

658

659

662

663

664

665

666

667

668

669

670

671

672

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.

636 Here are a few concrete examples that we will incorporate in the paper:

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 detects 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 "creditinvisible".
- Relabeling possibility: A borrower can report their payed 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 [12].
- 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

$$SP(h) = |P[h(x) = 1|a = 1] - P[h(x) = 1|a = 0]|, \tag{8}$$

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) [19] requires that both the true positive rate and false positive rate be equal across groups. Here the EqOd difference is defined as

$$EqOd(h) = \frac{1}{2} \sum_{z=0,1} |P[h(x) = 1|a = 1, y = z] - P[h(x) = 1|a = 0, y = z]|.$$
 (9)

692 B.2 Suboptimal Bayes classifier

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

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

TV (\cdot, \cdot) is the total variation distance between two distributions. The definition is discussed more in Hardt et al. [10].

B.3 Algorithmic collective action

In social sciences, *collective action* refers to the coordinated efforts of individuals working together to pursue a shared goal [20, 21]. Hardt et al. [10] adapt this notion to machine learning, 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 population that belongs to the collective.

Relation to fair representation learning. When users have agency over the training data, one possible form of collective action 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 [8, 28]. However, a hindrance of FRL in the context of collective action 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.

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

$$S(\alpha) = \mathbb{P}_0 \left[h\left(g\left(x\right)\right) = h\left(x\right) \right],\tag{10}$$

the probability, under the base distribution, that the classifier's prediction remains 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 the g is applied. Intuitively, if g removes a feature pattern correlated with group membership (e.g., minority vs. 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. [10] propose the collective relabels itself with the most likely label under the transformation q. Formally, the strategy is defined as

$$(x,y) \to \left(x, \arg\max_{y' \in \{0,1\}} \mathbb{P}_0\left(y'|g\left(x\right)\right)\right).$$
 (11)

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 1), Hardt et al. [10] prove the following lower bound for its success

$$S(\alpha) \ge 1 - \frac{2(1-\alpha)}{\alpha} \cdot \tau - \frac{\epsilon}{(1-\epsilon)\alpha},$$
 (12)

where $\tau = \underset{x \sim \mathbb{P}_0}{\mathbb{E}} \left[\max_{y' \in \{0,1\}} \left| \mathbb{P}_0\left(y'|x\right) - \mathbb{P}_0\left(y'|g\left(x\right)\right) \right| \right]$ measures the sensitivity of the true label distribution to the transformation g.

Note that the strategy in Equation (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 Appendix C.

731 B.4 Counterfactual fairness

The concept of counterfactual fairness (CF) [30, 31, 32] bridges between signal erasure success 732 to group fairness. To introduce this idea, assume that a sample x is generated by a causal model, 733 in which the group membership A is a causal parent. Then a classifier h is counterfactually fair if 734 its predictions are invariant to interventions on the group membership, i.e., $h(x) = h(x_{A \leftarrow a'})$ for 735 any a', where $x_{U\leftarrow u}$ denotes an intervention on a causal parent U of a sample x. In certain causal 736 contexts, CF implies or aligns with group fairness criteria such as SP or EqOd [33]. Therefore, if 737 collective action induces a counterfactually fair classifier, it may also induce a fair classifier under SP 738 739 or EqOd.

Since our focus is on fairness for the minority group, we relax the original definition of CF [30].

Definition 2. A classifier h is minority-focused counterfactually fair if under any context X = x,

$$\mathbb{P}_0(h(x_{A \leftarrow a}) = y | X = x, A = 1) = \mathbb{P}_0(h(x_{A \leftarrow a'}) = y | X = x, A = 1), \tag{13}$$

742 for any value a' attainable by A.

By this condition, changing the group membership of a minority individual, in a counterfactual sense, has no effect on the classifier's prediction. Collective action can theoretically enforce such fairness by applying the erasure strategy from Equation (4) with the counterfactual signal $g(x) = x_{A\leftarrow 0}$, which replaces a minority individual with its majority-group counterfactual. This collective action aligns the signal erasure success from Equation (3) with minority-focused counterfactual fairness from Definition 2. The following proposition, proved in Appendix E.1, formalizes this alignment.

Proposition 1. A Bayes classifier trained on \mathbb{P}_{α} is minority-focused counterfactually fair if and only

Proposition 1. A Bayes classifier trained on \mathbb{P}_{α} is minority-focused counterfactually fair if and only if the success of a minority collective is S=1.

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

C Limitations of Minority Collective Action

Previous work on collective action 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. Those differences introduce limitations to the existing theory, which we analyze and theoretically quantify in this section.

Collective restricted to the minority. As mentioned above, we focus on collectives composed solely of minority members, unlike prior work. This restriction expresses 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 reduces the collective's impact, as demonstrated in the following example.

Consider a binary classification task on the two-dimensional 4-Gaussian mixture model \mathbb{P}_{4GMM} where each Gaussian belongs to a distinct combination of label and group membership, as illustrated in Figure 5. 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 fairness violation of ERM.

Proposition 2 (Informal). Consider a dataset sampled from the distribution \mathbb{P}_{4GMM} described above, where every minority point participates in the collective action 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.

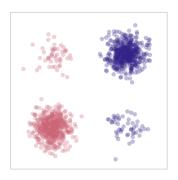


Figure 5: The distribution \mathbb{P}_{4GMM} used in Proposition 2. The color signifies the label, and the density shows the group membership.

The formal Proposition 5 is provided in Appendix E.2 along with all necessary assumptions, which holds for a broader family of distributions and can be extended to any dimensionality \mathbb{R}^d using techniques similar to those in Chaudhuri et al. [34]. Although Proposition 2 is not a formal lower bound, it emphasizes an important limitation: collective action restricted to the minority cannot generally achieve perfect fairness, even under very advantageous conditions involving a maximum-sized collective, a strong strategy, and a complete disregard for accuracy. This limitation stands in contrast to standard firm-side bias mitigation methods, which can, in principle, achieve perfect fairness.

We empirically corroborate the findings of Proposition 2 on real world datasets by examining the fairness–accuracy tradeoff of several fair learning methods. Most of these methods include a hyperparameter that controls this trade-off, yielding a set of pairs (Error, EqOd) as it varies. This set forms a Pareto front, representing the best attainable trade-offs. A Pareto front is said to *dominate* another if it lies entirely to the left (lower error) and below (lower unfairness) of the other.

Figure 4 compares the Pareto fronts of RB-prob, one of our minority collective action methods, with established firm-side methods. We observe that the lowest fairness violation achievable by RB-prob is greater than that of the firm-side approaches. However, the firm-side methods are able to arrive at perfect fairness only at a cost of prohibitively high prediction error. But, inspecting the region where the error is small compared to the base classifier, the fairness of RB-prob is comparable to that of the firm-side methods.

Approximating the class-conditional $\mathbb{P}_0(Y|X)$. In Section 3 we proposed methods to estimate which individuals would receive a different counterfactual label than their original label. However, the success lower bound in Equation (12) assumes perfect knowledge of \mathbb{P}_0 and its causal model. To account for approximation error, we model the collective's prediction as the output of algorithm $\mathcal{A}(x) \approx \mathbb{P}_0 \max_y (y|x_{A\leftarrow 0})$ that has an error rate

$$\rho := \mathbb{P}_{0}\left(\mathcal{A}\left(x\right) \neq \arg\max_{y'} \mathbb{P}_{0}\left[y'|g\left(x\right)\right]\right). \tag{14}$$

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

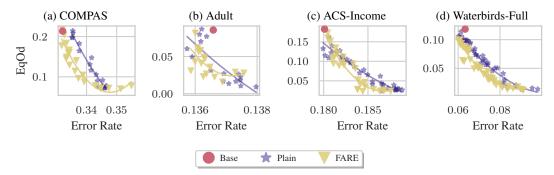


Figure 6: The Pareto fronts for using a fair representation when computing the KNN for RB-dist dominate the Pareto fronts for KNN computed on untransformed features. The blue stars represent the KNN without transforming the data, and the yellow triangles represent the KNN when the data is transformed using FARE [28]. The lines are fitted by a polynomial of degree 2 to guide the eye.

Proposition 3. With algorithm A(x) with label error ρ , the success of the collective is bounded by

$$S(\alpha) \ge 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}.$$
 (15)

This bound recovers Equation (12) when $\rho = 0$, but higher values of the error ρ worsen 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, their performance is sensitive to the choice of distance metric and feature representation. In our main experiments, 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 is distributed differently than the majority (e.g., \mathbb{P}_{4GMM}), this task can be challenging. The following informal result compares the error of 1-NN in the original features space to its error in a learned fair representation.

Proposition 4 (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 1, with the proof and assumptions can be found in Appendix E.4. The result suggests that FRL can reduce the counterfactual label error ρ of RB-label and RB-dist, consequently improving the lower bound of the collective's success according to Proposition 3. Empirically, Figure 6 indicates that applying FARE [28] before the KNN step improves the Pareto front for RB-dist. On the other hand, methods that rely purely on predictive information, such as RB-prob, can perform worse, due to FRL inadvertently removing features predictive of the class label. This behavior, and additional results, are provided in Figure 13 in the appendix.

D Related work

807

808

811

812

813

819 820

821

822

830

Optimizing for fairness metrics often comes at the cost of reduced classification accuracy, leading to the well-documented accuracy–fairness tradeoff [2, 3, 4, 5]. In response, previous work has proposed fairness interventions at different stages of the ML pipeline: pre-processing methods modify the training data before learning [6, 7, 8, 28], in-processing methods adjust the learning algorithm itself [35, 36, 26, 37], and post-processing methods correct the predictions of a trained (unfair)

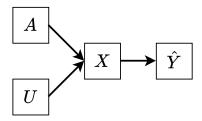


Figure 7: Assumed causal model for data generation and prediction. The group membership A and the other latent variables U are the causal parents of the observable features X. The classifier outputs a predicted label \hat{Y} that depends on the features X.

classifier [19, 38, 39, 40]. 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 837 mentioned in Appendix B.3, using feature-changing pre-processing methods such fair representation 838 learning [8, 28] demand changing those features during inference time as well. 839 Still, a couple of pre-processing methods suggest changing only the labels, similarly to our proposed 840 collective action. The method by Luong et al. [7] compares between the minority KNN and majority 841 KNN and flip the labels according to the difference of positive labels between the two groups of 842 neighbors. This method resembles RB-label, with the difference that RB-label examines only the 843 majority KNN in order to approximate the counterfactual. Similarly, the approach of Kamiran and 844 Calders [6] trains a regressor to predict y = 1 outcome probabilities, and flip the label of minority 845 members with y=0 labels and high probability according to the regressor to have y=1, and 846 similarly flip majority y = 1 labels to y = 0. Flipping from both groups is done to preserve the 847 error of the classifier. Our method RB-prob differs by training the regressor only on the majority to 848 better approximate the counterfactuals. Since this approach requires flipping the labels of majority members as well, it cannot be completey adopted by the collective. In Appendix G.1 we compare the between RB-prob to CND and KDP, and find that our method, based on the counterfactual search, is 851 more efficient in terms of the required number of label flips. 852

E Theoretical Results and Proofs

836

853

854 E.1 Counterfactual fairness as success

Proposition 1. A Bayes classifier trained on \mathbb{P}_{α} is minority-focused counterfactually fair if and only if the success of a minority collective is S=1.

Proof. For this proof, we assume the data is generated according to the causal model presented in Figure 7, where the features X are conditioned on the group membership A and other latent causal parent U. The features X are then used by a classifier to compute a predicted label $h\left(x\right)\hat{Y}$. In our case, the predicted label is the output of an optimal Bayes classier that predicts the most probable label as $h\left(x\right) = \arg\max_{y} P\left(y|x\right)$.

The data distribution is a mixture distribution between the majority distribution $\mathbb{P}_{A=0}$ and the minority distribution $\mathbb{P}_{A=1}$, which is defined as

$$\mathbb{P}_0 = (1 - \beta) \, \mathbb{P}_{A=0} + \beta \mathbb{P}_{A=1}, \tag{16}$$

where β is the proportion of the minority in the data.

The collective is employing the signal erasure strategy from Equation (4), where the erased signal is the counterfactual of x if they were a member of the majority group A=0, or formally as

$$g(x) = x_{A \leftarrow 0} \sim \mathbb{P}\left(X_{A \leftarrow 0}\right). \tag{17}$$

The training distribution is a mixture distribution of the data distribution \mathbb{P}_0 and the collective distribution \mathbb{P}^* , which is defined as

$$\mathbb{P}_{\alpha} = \alpha \mathbb{P}^* + (1 - \alpha) \mathbb{P}_0. \tag{18}$$

We now write the success of the collective (Equation (3)) in terms of the Bayes classifier as

$$S = \mathbb{P}_{0} \left[h\left(x \right) = h\left(g\left(x \right) \right) \right]$$

$$= \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(y | x \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(y | g\left(x \right) \right) \right].$$
(19)

- To compute this probability, we split it into two cases, conditioning on the group membership A.
- When conditioning the success on the majority group A = 0, then g(x) = x as the intervention on A,
- which converts to the majority, does not change the value of A, which is already the majority. This
- 873 trivially leads to

$$S_{A=0} = \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(y | x, A = 0 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(y | g \left(x \right), A = 0 \right) \right]$$

$$= \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(y | x, A = 0 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(y | x, A = 0 \right) \right]$$

$$= 1.$$
(20)

- For conditioning the success on the minority, recall that the data is generated according to the causal
- model in Figure 7 which means that intervention on the group membership A can be passed down to
- the features X as

$$\mathbb{P}\left(h\left(x_{A\leftarrow 0}\right) = y|X, A = 1\right) = \mathbb{P}\left(h\left(x\right) = y|X_{A\leftarrow 0}, A = 1\right) = \mathbb{P}\left(h\left(x\right) = y|g\left(X\right), A = 1\right). \tag{21}$$

877 This can be used to write the success conditioned on the minority as

$$S_{A=1} = \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(y | x, A = 1 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(y | g \left(x \right), A = 1 \right) \right]$$

$$= \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(h \left(x_{A \leftarrow 1} \right) = y | X, A = 1 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(h \left(x_{A \leftarrow 0} \right) = y | X, A = 1 \right) \right]. \tag{22}$$

- The first term is rewritten to use the intervention notation even though the intervened variable is
- 879 unchanged.
- As the proportion of the minority is known to be β , the success can be written by combining
- 881 Equations (20) and (22) using the law of total probability as

$$S = 1 - \beta + \beta \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 1}\right) = y | X, A = 1 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 0}\right) = y | X, A = 1 \right) \right]$$

$$= 1 - \beta \left(1 - \mathbb{P}_{0} \left[\arg \max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 1}\right) = y | X, A = 1 \right) = \arg \max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 0}\right) = y | X, A = 1 \right) \right] \right). \tag{23}$$

- This equality can be examined under two scenarios: when the success is perfect S=1 and when the classifier is minority-focused counterfactually fair.
- When the success is S=1 If the success of the collective is S=1, then Equation (23) leads to

$$\mathbb{P}_{0}\left[\arg\max_{y}\mathbb{P}_{\alpha}\left(h\left(x_{A\leftarrow1}\right)=y|X,A=1\right)=\arg\max_{y}\mathbb{P}_{\alpha}\left(h\left(x_{A\leftarrow0}\right)=y|X,A=1\right)\right]=1. (24)$$

885 This means that it is certain that

$$\arg\max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 1} \right) = y | X, A = 1 \right) = \arg\max_{y} \mathbb{P}_{\alpha} \left(h\left(x_{A \leftarrow 0} \right) = y | X, A = 1 \right), \tag{25}$$

Since the label is binary, then it follows that the same applies to using $\arg \min$. Therefore, for all $y \in \{0, 1\}$ we have

$$\mathbb{P}_{\alpha} (h(x_{A \leftarrow 1}) = y | X, A = 1) = \mathbb{P}_{\alpha} (h(x_{A \leftarrow 0}) = y | X, A = 1), \tag{26}$$

which is the definition of a minority-focused counterfactually fair classifier (Definition 2).

When the classifier is one-sided counterfactually fair If the classifier is one-sided counterfactually fair (Definition 2), then by definition 890

$$\mathbb{P}_0\left[\arg\max_{y}\mathbb{P}_{\alpha}\left(h\left(x_{A\leftarrow 1}\right)=y|X,A=1\right)=\arg\max_{y}\mathbb{P}_{\alpha}\left(h\left(x_{A\leftarrow 0}\right)=y|X,A=1\right)\right]=1 \quad (27)$$

and plugging that in Equation (23) results in S=1. 891

E.2 Impossibility of fairness under ERM 892

- The following proposition follows the structure of Theorem 6 in Chaudhuri et al. [34]. For a vector 893 $x \in \mathbb{R}^d$, let D(x) denote a distribution on \mathbb{R}^d with mean x. Let p and m be the number of majority 894 and minority sample, respectively with $p \gg m$. 895
- Assumption 1 (Concentration Condition, Assumption 2 from Chaudhuri et al. [34]). Let 896 $x_1, \ldots, x_n \overset{i.i.d.}{\sim} D(0)$ in \mathbb{R}^d . There exist maps $X_{\max}, c, C : \mathbb{Z}_+ \times [0, 1] \times \mathbb{Z}_+ \to \mathbb{R}$ 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 $1 - \delta$ 897

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

- and $\lim_{n\to\infty} C(n,\delta,d) = 0$, $\lim_{n\to\infty} c(n,\delta,d) = 0$. 899
- **Data Model** Labels $y \in \{-1, 1\}$ and protected attribute $a \in \{-1, 1\}$ define four groups whose 900 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 $\mu = 0$ 902 $\|\mu\|$ $(0,1)^{\top}$ and $\psi=\|\psi\|$ $(1,0)^{\top}$. Without loss of generality, let the majority attribute be $a_M=+1$ 903
- (the minority is $a_m = -1$). Thus the two majority means lie on the positive diagonal $\pm (\mu + \psi)$ and 904
- the two minority means on the negative diagonal $\pm (\mu \psi)$. 905
- Let B, A be two sets of points that are sampled from D(0). We will always associate B with 906 negatively labelled points and A with positive, as will be clear below. Following Chaudhuri et al. 907
- [34], define the sets 908

$$A_{\mu} = \{x + \mu : x \in A\}, \qquad -B_{\mu} = \{x + \mu : x \in -B\}.$$

- We split by attribute and (for the minority) allow arbitrary relabeling before training. Write A_{μ}^{M} , 909 B_{μ}^{M} for the majority parts and A_{μ}^{m} , B_{μ}^{m} for the minority subsets used with positive/negative labels in training after centering. Incorporating the attribute shifts, set 910

$$A^{M}_{\mu,\psi} = -\psi + A^{M}_{\mu}, \quad B^{M}_{\mu,\psi} = +\psi + B^{M}_{\mu}, \qquad A^{m}_{\mu,\psi} = +\psi + A^{m}_{\mu}, \quad B^{m}_{\mu,\psi} = -\psi + B^{m}_{\mu},$$

- and similarly for the relabeled minority pieces $A_{\mu,\psi}^{m,\pm}$ and $B_{\mu,\psi}^{m,\pm}$ (these are subsets of $A_{\mu,\psi}^m$ and $B_{\mu,\psi}^m$, 912 respectively).
- 913
- If the minority were absent, the ERM SVM converges to the spurious direction 914

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

- However, we assume that the minority is performing some relabeling. As a result, we denote the set
- 916
- $A_{\mu}^{m,+}$ as the samples relabeled with y=1 and the set $A_{\mu}^{m,-}$ as the samples keeping the original label y=0. Similarly we denote the set $B_{\mu}^{m,+}$ as the positive minority keeping their labels and $B_{\mu}^{m,-}$ as
- the positive minority who flip to y = 0. 918
- **Proposition 5.** Suppose D(0) satisfies Assumption 1 and 919

$$X_{max}(p,\delta,2) - X_{max}(m,\delta,2) \ge 2 \|\psi\| + c(p,\delta,2) + C(m,\delta,2).$$
 (28)

- Then, for any (possibly adversarial) relabeling of minority training examples, if $p \to \infty$, with 920
- probability at least $1-4\delta$, the SVM ERM solution converges to the same spurious solution $w_{spu}^* \propto$
- $\mu + \psi$. Under a centrally symmetric D(0) and when $\|\mu\| = \|\psi\|$, this limit satisfies $EqOd\left(w_{spu}^*\right) \to 0$
- 923 0.5.

Proof. As Chaudhuri et al. [34] shows, the ERM solution can be written as $w^* = \alpha^* \hat{\mu} + \sigma \beta^* \hat{\psi}$,

where 925

$$\alpha^* = \arg\min_{\alpha \in [-1,1], \sigma \in \{-1,1\}} \sup_{x \in \{x|y=0\}} \left(\alpha \hat{\mu} + \sigma \beta \hat{\psi} \right)^{\mathsf{T}} (x-\mu) + \sup_{x \in \{x|y=1\}} \left(\alpha \hat{\mu} + \sigma \beta \hat{\psi} \right)^{\mathsf{T}} (x-\mu)$$

and $\beta = \sqrt{1 - \alpha^2}$. 926

With the shorthand 927

$$f_1(\alpha) := \sup_{x \in A_{\mu,\psi}^M} (\alpha \hat{\mu} + \sigma \beta \hat{\psi})^\top x, \qquad f_{2,\pm}(\alpha) := \sup_{x \in A_{\mu,\psi}^{m,\pm}} (\alpha \hat{\mu} + \sigma \beta \hat{\psi})^\top x,$$

$$f_3(\alpha) := \sup_{x \in -B_{\mu,\psi}^M} (\alpha \hat{\mu} + \sigma \beta \hat{\psi})^\top x, \qquad f_{4,\pm}(\alpha) := \sup_{x \in -B_{\mu,\psi}^{m,\pm}} (\alpha \hat{\mu} + \sigma \beta \hat{\psi})^\top x,$$

the SVM objective is

$$F(\alpha) = \min_{\alpha} \Big\{ \max \big(f_1(\alpha) - \alpha \|\mu\| + \sigma \beta \|\psi\|, \ f_{2,-}(\alpha) - \alpha \|\mu\| - \sigma \beta \|\psi\|, \ f_{4,-}(\alpha) - \alpha \|\mu\| - \sigma \beta \|\psi\| \big) + \max \big(f_3(\alpha) - \alpha \|\mu\| + \sigma \beta \|\psi\|, \ f_{2,+}(\alpha) - \alpha \|\mu\| - \sigma \beta \|\psi\|, \ f_{4,+}(\alpha) - \alpha \|\mu\| - \sigma \beta \|\psi\| \big) \Big\}.$$

By Assumption 1, for the majority group of size p, there exists X_p, c_p, C_p such that, with probability

at least $1-4\delta$ and for all α , 930

$$f_1(\alpha), f_3(\alpha) \in [X_p - c_p, X_p + C_p]. \tag{29}$$

For any minority relabeling, $A_{\mu,\psi}^{m,\pm} \subseteq A_{\mu,\psi}^m$ and $-B_{\mu,\psi}^{m,\pm} \subseteq -B_{\mu,\psi}^m$, so the same assumption gives 931

$$f_{2,+}(\alpha), f_{4,+}(\alpha) \le X_m + C_m,$$
 (30)

where $X_m := X_{\text{max}}(m, \delta, 2)$ and $C_m := C(m, \delta, 2)$. Using Equation (28), we get

$$X_p - X_m \ge 2\|\psi\| + c_p + C_m. \tag{31}$$

Combining Equations (29) to (31), still uniformly in α , we obtain

$$f_1(\alpha) - f_{2,\pm}(\alpha) \ge 2\|\psi\|, \quad f_1(\alpha) - f_{4,\pm}(\alpha) \ge 2\|\psi\|, \quad f_3(\alpha) - f_{2,\pm}(\alpha) \ge 2\|\psi\|, \quad f_3(\alpha) - f_{4,\pm}(\alpha) \ge 2\|\psi\|.$$
(32)

- Case 1: $\sigma = 1$. Consider the *first* inner maximum inside $F(\alpha)$. Compare the majority entry 934
- (associated with $f_1(\alpha)$) to the minority entries $(f_2 (\alpha), f_4 (\alpha))$: 935

$$[f_1(\alpha) - \alpha \|\mu\| + \beta \|\psi\|] - [f_{2,-}(\alpha) - \alpha \|\mu\| - \beta \|\psi\|] = (f_1(\alpha) - f_{2,-}(\alpha)) + 2\beta \|\psi\| \ge 2\|\psi\| + 2\beta \|\psi\| > 0,$$

- and similarly against $f_{4,-}$. Hence the first maximum equals $f_1 \alpha \|\mu\| + \beta \|\psi\|$. For the second inner maximum, the same comparison yields the majority term $f_3 \alpha \|\mu\| + \beta \|\psi\|$. Summing then 936
- 937
- for $\sigma = 1$ we have, 938

$$F_{+}(\alpha) = f_{1}(\alpha) + f_{3}(\alpha) - 2\alpha \|\mu\| + 2\beta \|\psi\|.$$

Case 2: $\sigma = -1$. For the first inner max in $F(\alpha)$, 939

$$\left[f_{1}(\alpha) - \alpha \|\mu\| - \beta \|\psi\|\right] - \left[f_{2,-}(\alpha) - \alpha \|\mu\| + \beta \|\psi\|\right] = \left(f_{1}(\alpha) - f_{2,-}(\alpha)\right) - 2\beta \|\psi\| \ge 2\|\psi\| - 2\beta \|\psi\| \ge 0,$$

- and likewise against $f_{4,-}$. Thus the first maximum equals $f_1 \alpha \|\mu\| \beta \|\psi\|$. The second inner max
- is analogous and equals $f_3 \alpha \|\mu\| \beta \|\psi\|$. Therefore, 941

$$F_{-}(\alpha) = f_{1}(\alpha) + f_{3}(\alpha) - 2\alpha \|\mu\| - 2\beta \|\psi\|$$

For every α , $F_+(\alpha) = (f_1 + f_3) - 2\alpha \|\mu\| + 2\beta \|\psi\|$ and $F_-(\alpha) = (f_1 + f_3) - 2\alpha \|\mu\| - 2\beta \|\psi\|$, so $F_+(\alpha) \ge F_-(\alpha)$. Hence the optimal sign is $\sigma = -1$ and the objective reduces to

943

$$F(\alpha) = (f_1(\alpha) + f_3(\alpha)) - 2\alpha \|\mu\| - 2\beta \|\psi\|.$$

Maximizing $\alpha \|\mu\| + \beta \|\psi\|$ is equivalent to minimizing $F(\alpha)$ up to the bounded change (due to As-

sumption 1) of $f_1(\alpha) + f_3(\alpha)$. Next, we use the following lemma.

- **Lemma 1** (Approximate Maximization Lemma I, Lemma 14 from Chaudhuri et al. [34]). 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 947
- $\alpha_F \in \operatorname{argmax}_{\alpha} F(\alpha)$, 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
- $max_{\alpha}F(\alpha) > \sqrt{u^2 + v^2} L.$ 950
- Applying Lemma 1 with $u = \|\mu\|$ and $v = \|\psi\|$ shows that (α, β) approaches 951

$$(\alpha_g, \beta_g) = \left(\frac{\|\mu\|}{\sqrt{\|\mu\|^2 + \|\psi\|^2}}, \ \frac{\|\psi\|}{\sqrt{\|\mu\|^2 + \|\psi\|^2}}\right)$$

as $p \to \infty$. Thus 952

$$w^* \longrightarrow w^*_{\text{spu}} = \alpha_g \,\hat{\mu} + \beta_g \,\hat{\psi},$$

- . $w^* \ \longrightarrow \ w^*_{\rm spu} = \alpha_g \, \hat{\mu} + \beta_g \, \hat{\psi},$ independently of how the minority samples were relabeled in training. 953
- Under a centrally symmetric D(0) and if $\|\mu\| = \|\psi\|$, the majority group (a = +1) separates perfectly 954
- in the limit, while the minority group (a = -1) has symmetric measure about the threshold, giving 955
- $\text{TPR}_{a=+1} \to 1$, $\text{FPR}_{a=+1} \to 0$, and $\text{TPR}_{a=-1} = \text{FPR}_{a=-1} \to \frac{1}{2}$. Hence $\text{EqOd}(w_{\text{spu}}^*) \to 0.5$. 956
- 957
- This result can also be extended to \mathbb{R}^d using techniques similar to those in Chaudhuri et al. [34]. This 958
- result also encompasses the 4-Gaussian mixture model \mathbb{P}_{4GMM} used in Appendix C as a special case, 959
- leading to the following. 960
- **Proposition 2** (Informal). Consider a dataset sampled from the distribution \mathbb{P}_{4GMM} described above, 961
- where every minority point participates in the collective action by flipping all y=0 labels to y=1. 962
- Then, under sufficiently separable clusters, with high probability, the EqOd of the ERM classifier 963
- minimizing the logistic loss will asymptotically approach 0.5. 964

E.3 Success Bound With Label Error 965

- The following proof uses Lemma 11 from Hardt et al. [10]. 966
- **Lemma 2** (Lemma 11 from Hardt et al. [10]). Suppose that P, P' are two distributions such that 967
- $\mathrm{TV}(P,P') \leq \epsilon$. Take any two events E_1,E_2 measurable under P,P'. If $P(E_1) > P(E_2) + \frac{\epsilon}{1-\epsilon}$, 968
- then $P'(E_1) > P'(E_2)$. 969
- **Proposition 3.** With algorithm A(x) with label error ρ , the success of the collective is bounded by 970

$$S(\alpha) \ge 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}.$$
 (15)

- *Proof.* This proof follows closely the proof of Theorem 5 by Hardt et al. [10]. We start under the 971
- assumption of an optimal Bayes classifier, setting $\epsilon = 0$. 972
- When the new label y' is wrong with probability ρ , then we can think of the collective as being union 973
- of two sub-collectives: one with the correct label and one with the incorrect label. In the binary case 974
- this can be formulated with correct subcollective P^+ as having label $y' = \arg \max_{y} P_0(y|g(x))$ and 975
- the incorrect subcollective P^{-} as with label $y' = \arg\min_{y} P_{0}\left(y | g\left(x\right)\right)$. Then we can write the train 976
- distribution as 977

$$P_{\alpha} = \alpha \left(\rho P^{-} + (1 - \rho) P^{+} \right) + (1 - \alpha) P_{0}$$

= $\alpha \rho P^{-} + (1 - \rho) \alpha P^{+} + (1 - \alpha) P_{0}$. (33)

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

$$P_{\alpha}(y^*|x) = \alpha \rho P^{-}(y^*|x) + (1 - \rho) \alpha P^{+}(y^*|x) + (1 - \alpha) P_{0}(y^*|x)$$

= $(1 - \rho) \alpha + (1 - \alpha) P_{0}(y^*|x),$ (34)

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

$$P_{\alpha}(y|x) = \alpha \rho P^{-}(y|x) + (1 - \rho) \alpha P^{+}(y|x) + (1 - \alpha) P_{0}(y|x)$$

= \alpha \rho + (1 - \alpha) P_{0}(y|x), (35)

where $P^+(y^*|x) = 1$, $P^-(y^*|x) = 0$, $P^+(y^*|x) = 0$, $P^-(y^*|x) = 1$ by definition. 980

A Bayes classifier h returns the most probable label $h(x) = \arg \max_{y} P(y|x)$. Therefore, a Bayes 981 classifier will output y^* if the probability is greater, which can be written as the condition 982

$$P_{\alpha}(y^{*}|x) > P_{\alpha}(y|x)$$

$$(1 - \rho)\alpha + (1 - \alpha)P_{0}(y^{*}|x) > \alpha\rho + (1 - \alpha)P_{0}(y|x)$$

$$(1 - 2\rho)\alpha > (1 - \alpha)(P_{0}(y|x) - P_{0}(y^{*}|x)).$$
(36)

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

$$P_{0}(y|x) - P_{0}(y^{*}|x) \leq P_{0}(y|x) - P_{0}(y|g(x)) + P_{0}(y^{*}|g(x)) - P_{0}(y^{*}|x)$$

$$\leq 2\tau(x).$$
(37)

With that, the condition in Equation (36) can be written as

$$(1 - 2\rho)\alpha > 2(1 - \alpha)\tau(x). \tag{38}$$

With that, the success can be bounded as

$$S = P_{0} [f(x) = f(g(x))]$$

$$= P_{0} [f(x) = y^{*}(x)]$$

$$\geq P_{0} [(1 - 2\rho) \alpha > 2 (1 - \alpha) \tau(x)]$$

$$= P_{0} \left[1 - \frac{2(1 - \alpha)}{(1 - 2\rho) \alpha} \tau(x) > 0 \right]$$

$$= \mathbb{E}_{x \sim P_{0}} \left[1 \left\{ 1 - \frac{2(1 - \alpha)}{(1 - 2\rho) \alpha} \tau(x) > 0 \right\} \right]$$

$$\geq \mathbb{E}_{x \sim P_{0}} \left[1 - \frac{2(1 - \alpha)}{(1 - 2\rho) \alpha} \tau(x) \right]$$

$$= 1 - \frac{2(1 - \alpha)}{(1 - 2\rho) \alpha} \tau$$
(39)

With sub-optimality $\epsilon > 0$ A result of Lemma 2 is to write the condition in Equation (38) as

$$(1 - 2\rho)\alpha > 2(1 - \alpha)\tau(x) + \frac{\epsilon}{1 - \epsilon},\tag{40}$$

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

$$S(\alpha) \ge 1 - \frac{2(1-\alpha)}{(1-2\rho)\alpha}\tau - \frac{\epsilon}{(1-\epsilon)(1-2\rho)\alpha}.$$
(41)

988

E.4 Label Error With Better Representation

989

For the following we assume a similar setting as in Appendix E.2, visualised as a 2D distribution 990 in Figure 5. We are given the majority data, and tasked with labeling the minority data. Assume 991

all labels are distributed equally $\mathbb{P}[Y=1]=\mathbb{P}[Y=-1]=\frac{1}{2}$. The minority features X_{\min} are 992

distributed as $X_{\min} \sim \mathcal{N}(y\mu_{\min}, \Sigma_{\min})$ with $X_{\min} \in \mathbb{R}^d$. The label $\hat{y}_{1\mathrm{NN}}^{(n)}$ 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 $X_- \sim \mathcal{N}(-\mu, \Sigma)$. 993

994

995

Theorem 1. 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 $X_{\min} \sim$ 996 997 $\mathcal{N}(y\mu_{min}, \Sigma_{min}).$ 998

Then, there exists a projection $P \in \mathbb{R}^{d \times d}$ such that asymptotically for $n \to \infty$, $err_{INN}^{rep} < err_{INN}^{raw}$ 999

Proof. Consider the projection on the hyperplane perpendicular to w, where $w = \frac{\mu - \mu_{\min}}{2}$. The 1000 projection matrix associated with this transformation is $P = I - \frac{ww^\top}{w^\top w}$

Let us denote the symbols after the projection as $\bar{\mu}:=P\mu$, $\bar{\mu}_{\min}:=P\mu_{\min}$, $\bar{v}:=(P\Sigma P^T)^+\bar{\mu}$ and $\bar{\Sigma}_{\min}:=P\Sigma_{\min}P^T$. Here we denoted using A^+ the pseudoinverse of the matrix A. Note that since P is an orthogonal projection matrix, it holds that PP=P and $P^T=P$.

We apply Lemma 3 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:

$$\operatorname{err}_{1\text{NN}} = \frac{1}{2} \mathbb{P}_{X_{\min}|y=1}[\hat{y}_{1\text{NN}} = -1] + \frac{1}{2} \mathbb{P}_{X_{\min}|y=-1}[\hat{y}_{1\text{NN}} = 1]$$
 (42)

$$= \frac{1}{2} \left(1 - \Phi \left(\frac{v^{\top} \mu_{\min}}{\sqrt{v^{\top} \Sigma_{\min} v}} \right) \right) + \frac{1}{2} \Phi \left(\frac{-v^{\top} \mu_{\min}}{\sqrt{v^{\top} \Sigma_{\min} v}} \right)$$
(43)

$$=1-\Phi\left(\frac{v^{\top}\mu_{\min}}{\sqrt{v^{\top}\Sigma_{\min}v}}\right) \tag{44}$$

$$=1-\Phi(SNR),\tag{45}$$

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

Similarly, let us denote the SNR corresponding to 1NN applied on the projected representation as

1010 follows:
$$SNR_{\text{proj}} := \frac{\bar{v}^{\top} \bar{\mu}_{\min}}{\sqrt{\bar{v}^{\top} \bar{\Sigma}_{\min} \bar{v}}}$$

To show that $err_{1NN} > err_{1NN}^{rep}$ it suffices to prove that $SNR < SNR_{proj}$

We begin by rewriting the numerator of SNR_{proj} . Since $\mu \in Im(P)$ and because on Im(P) the operators Σ^{-1} and $(P\Sigma P^{\top})^+$ represent the same transformation, it follows that:

$$\bar{v} = (P\Sigma P^{\top})^{+} \bar{\mu} = \Sigma^{-1} \bar{\mu}.$$

Moving on the the denominator of SNR_{proj} , we have that:

$$\begin{split} \bar{v}^{\top} \bar{\Sigma}_{\min} \bar{v} &= \bar{v}^{\top} (P \Sigma_{\min} P^{\top})^{+} \bar{v} \\ &= \bar{v}^{\top} (P P^{\top})^{+} \bar{v} \\ &= \bar{v}^{\top} P^{+} \bar{v} \\ &= \bar{v}^{\top} P \bar{v} \\ &= \bar{v}^{\top} \bar{v} \\ &= \|\bar{v}\|^{2}. \end{split}$$

In the second line we used the fact that $\Sigma_{\min} = I$, in the third line we use the identity $P^2 = P$ due to P being a projection matrix, in the forth line we use $P^+ = P$ since P is an orthogonal projection (i.e. P is symmetric) and in the fifth line we use the fact that $\bar{v} \in Im(P)$, and hence, $P\bar{v} = \bar{v}$.

Putting everything together, and using the fact that Σ (and thus, Σ^{-1}) is positive definite (i.e. $x^{\top}\Sigma^{-1}x>0, \forall x\in\mathbb{R}^d$) we get that:

$$SNR_{\text{proj}} = \frac{\bar{\mu}^{\top} \Sigma^{-1} \bar{\mu}}{\|\bar{v}\|^2} > 0 > \frac{\mu^{\top} \Sigma^{-1} \mu_{\min}}{\|\Sigma^{-1} \mu\|^2} = SNR.$$

1020

Lemma 3. For a unimodal minority distribution $X_{min} \sim \mathcal{N}(\mu_{min}, \Sigma_{min})$ it holds that:

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

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

1022 *Proof.* Let us denote $\hat{y}_{1\mathrm{NN}} := \lim_{n \to \infty} \hat{y}_{1\mathrm{NN}}^{(n)}$ and let p_+ and p_- be the densities of two class1023 conditional distribution. Notice that the two class conditional training distributions are supported on
1024 the entire domain of \mathbb{R}^d . Therefore, in the asymptotic regime, the label $\hat{y}_{1\mathrm{NN}}$ at a test point x is given
1025 according to the class-conditional distribution that has higher density. Namely, we have:

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

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

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

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

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

1031

Note that the error from Theorem 1 is defined the same as ρ (Equation (14). This leads to the following.

Proposition 4 (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, $\rho_{FRL} \leq \rho_{plain}$.

F Technical Details

F.1 Datasets

1039

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 reoffended and the sensitive attribute is the race. We follow the same data cleaning and pre-processing
as Alghamdi et al. [38].

Adult The Adult dataset [23] 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. [38].

HSLS The High School Longitudinal Study of 2009 (HSLS) [24] 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. [38].

ACS-Income Ding et al. [25] 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 [26] 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 [27] 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.

F.2 Training

1064

1065

1066

1067

1068

1071

1072

All classification experiments were trained with *scikit-learn*'s histogram-based gradient boosting classification tree with the default parameters ³. 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 [8] we used the implementation in *Holistic AI*'s open source library ⁴ with the default parameters. For FARE [28] we used the official implementation ⁵ with hyperparameters

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

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

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

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

⁵https://github.com/eth-sri/fare

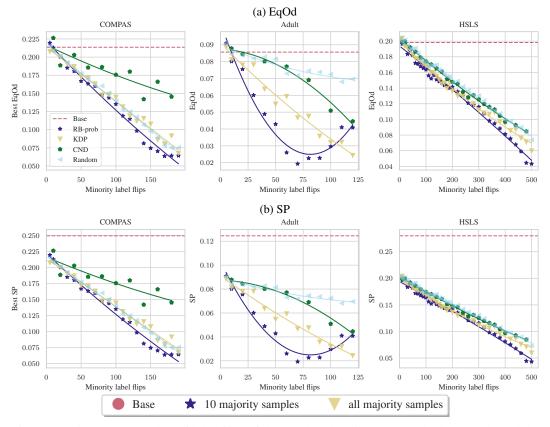


Figure 8: Fairness per number of label flips of the Random baseline, our method RB-prob, and the existing methods KDP [7] and CND [6]. Our method is more efficient than prior work, requiring less flips to achieve the same level of fairness. Note that in this experiment CND could flip any label, while all other methods were restricted to the labels of 30% of the minority.

1080 $\gamma=0.85, k=200$ and n=100. For all distance computation we used the Euclidean norm ℓ^2 -norm as $d\left(v,u\right)=\left\|v-u\right\|_2=\sqrt{\sum_i\left(v_i-u_i\right)^2}.$

G Additional Results

1082

1083

1084

1085

1086

1087

G.1 Comparison with prior work

We compare our method RB-prob with the existing methods KDP [7] and CND [6] in Figure 8. Figure 8 shows that our method, motivated by the counterfactual labeling, is more efficient in terms of required number of label flips, than the existing works.

G.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 (Equation (2)) and SP (Equation (8)) as a measure of unfairness

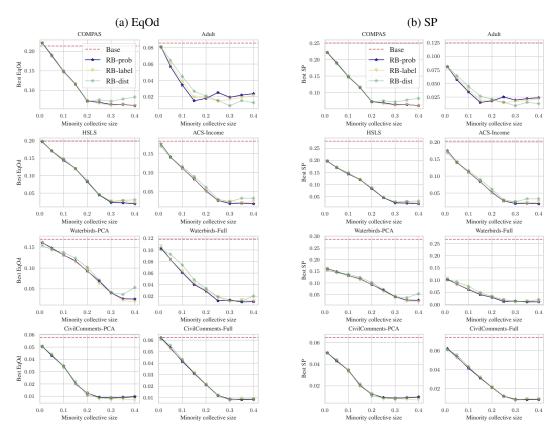


Figure 9: The lowest EqOd violation a collective can achieve greatly improves as the collective size increases, up to a certain point. Each point is a mean of 10 runs, with the standard deviation being smaller than the markers. In all the datasets we experimented on, the lowest EqOd violation converges around $\alpha=0.3$. Additional results are presented in Figure 9 in the appendix.

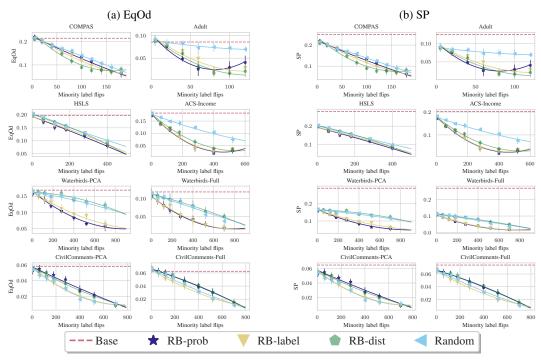


Figure 10: Our proposed methods are consistently more efficient than randomly flipping labels, requiring less label flips to attain the same level of EqOd. Each marker is the mean of 10 random runs with a specific number of label flips. The standard deviation is presented by the error bars. The dashed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

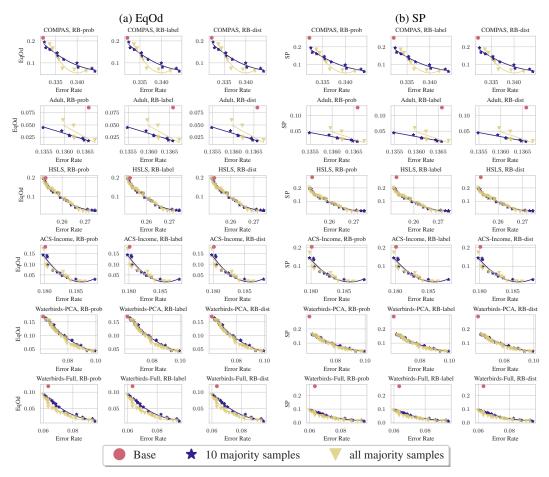


Figure 11: Limiting the knowledge of the collective about the majority does not significantly harm the Pareto front. Each point is the mean of 10 runs and the curves are fitted to guide the eye.

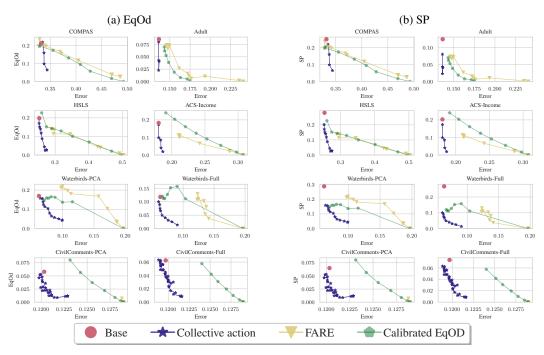


Figure 12: The firm-side pre-processing method FARE [28] and the post-processing method calibrated equalized odds [29] attain 0 EqOd with large error, while RB-prob with $\alpha=0.3$ (Section 3) has much smaller error and less unfairness than the base classifier, but unable to get 0 EqOd.

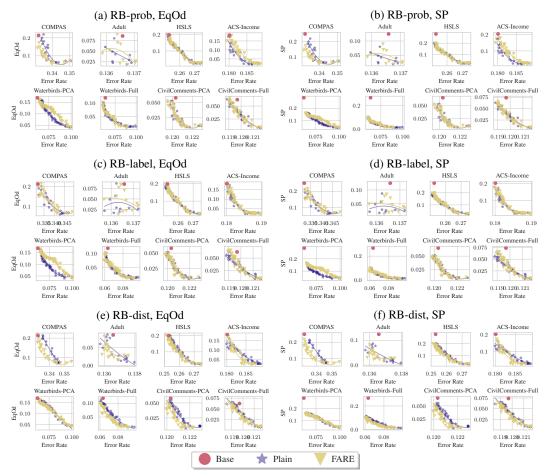


Figure 13: The Pareto fronts for using a fair representation when computing the KNN for RB-dist dominate the Pareto fronts for KNN computed on untransformed features. The blue stars represent the KNN without transforming the data, and the yellow triangles represent the KNN when the data is transformed using FARE [28]. The lines are fitted by a polynomial of degree 2 to guide the eye.