



FairScan: A Dual-Stage Bias Detection and Mitigation Framework for Machine Learning Classification Models

Shubhi Bhardwaj, Dr. Yatu Rani

Department of Artificial Intelligence and Data Science
Dr. Akhilesh Das Gupta Institute of Professional Studies
Shastri Park, New Delhi

Abstract- Bias in machine learning models is one of the biggest concerns in today's AI-driven world. When models are trained on data that reflects real-world inequalities, they end up making unfair predictions that can harm people based on their gender, race, or age. This paper introduces FairScan, a two-stage framework designed to first detect and then actively reduce bias in classification models. The detection stage uses a new metric called the Statistical Parity Divergence Score (SPDS), which measures bias not just across individual groups but also at the intersections of multiple sensitive attributes. The mitigation stage applies a custom training strategy called Reweighted Fair Gradient Descent (RFGD), which adjusts how much the model learns from different groups during training to push it toward fairer outcomes. We tested our approach on the UCI Adult Income dataset and found that FairScan reduced the Demographic Parity Difference by up to 79.4% while maintaining a classification accuracy of 86.7%. Our results show that it is genuinely possible to build models that are both accurate and fair, which is a step forward for responsible AI development.

Keywords- Algorithmic fairness, bias mitigation, machine learning, demographic parity, gradient descent, responsible AI, intersectionality

I. INTRODUCTION

Machine learning models are increasingly being used to make important decisions — from approving bank loans and screening job applications to deciding whether someone gets bail or medical treatment. While these systems are efficient and scalable, there is a serious problem hiding inside them: they can pick up on unfair patterns in historical data and reproduce those patterns at scale. A model trained on biased hiring data might keep rejecting applications from women. A model trained on biased crime data might keep flagging Black individuals as high-risk. The model is not intentionally prejudiced, but the outcome is discrimination nonetheless [1], [2].

Bias in AI is not just a theoretical concern anymore. Studies have found that facial recognition systems have significantly higher error rates for darker-skinned women compared to lighter-skinned men [3]. Language models reproduce gender stereotypes. Credit-scoring algorithms give lower scores to certain racial groups even when financial history is identical [4]. These are real harms, and researchers have been working hard to find ways to detect and fix them.

Most existing approaches either focus purely on detecting bias (but not fixing it) or apply mitigation at only one stage of the machine learning pipeline. There has been relatively little work on combining detection and mitigation together in a way that also considers intersectional bias — the compounding



effect that happens when someone belongs to multiple disadvantaged groups at once [5]. For example, a Black woman may face worse outcomes than either a Black man or a White woman, because her experience sits at the intersection of both race and gender bias.

This paper addresses these gaps by introducing FairScan, a unified framework that tackles bias at two stages. The first stage detects bias using a new metric (SPDS) that captures both group-level and intersectional disparities. The second stage reduces bias using a modified training algorithm (RFGD) that applies higher learning pressure on underrepresented subgroups. Together, these two stages provide a practical and effective tool for building fairer classifiers.

The rest of this paper is organized as follows. Section II reviews related work in algorithmic fairness. Section III formally defines the problem. Section IV explains the FairScan framework in detail. Section V describes the experiments and datasets. Section VI presents the results. Section VII discusses findings and limitations, and Section VIII concludes.

II. RELATED WORK

Research on algorithmic fairness can be broadly grouped into three categories depending on where in the machine learning pipeline the bias correction happens: pre-processing, in-processing, and post-processing [6].

Pre-processing methods modify the training data before the model ever sees it. Techniques like resampling, reweighting, and dataset transformation have all been explored. Kamiran and Calders [7] proposed massaging training labels to reduce statistical dependence between sensitive attributes and outcomes. Feldman et al. [8] introduced a data repair approach that alters feature distributions to reduce disparate impact. While these methods are model-agnostic, they risk losing important information in the data.

In-processing methods modify the training algorithm itself. This includes adding fairness constraints to the optimization objective or using adversarial training to remove sensitive information from learned representations. Zhang et al. [9] used an adversarial network to simultaneously predict the target label and prevent a separate classifier from predicting the sensitive attribute. Zafar et al. [10] framed fairness as a constraint in the optimization problem and solved it using covariance-based constraints. These approaches are powerful but tightly tied to specific model types.

Post-processing methods adjust the model's output after training. Hardt et al. [11] introduced Equalized Odds, which adjusts decision thresholds separately for each group to equalize error rates. Pleiss et al. [12] extended this to calibrated equalized odds. These methods are easy to apply but cannot change the internal representations the model has learned.

A significant gap in all three categories is the handling of intersectional bias. Kearns et al. [13] and Foulds et al. [14] have begun exploring fairness at the level of subgroups defined by combinations of attributes. However, most production tools do not yet implement this. FairScan directly addresses this gap by combining intersectional detection with in-processing mitigation in a single, coherent pipeline.

III. PROBLEM FORMULATION

Let $D = \{(x_i, s_i, y_i)\}$ for $i = 1$ to N be a dataset where x_i is a feature vector, $s_i \in S$ is a sensitive attribute (such as gender or race), and $y_i \in \{0, 1\}$ is a binary label. A classifier $f: X \rightarrow \{0, 1\}$ is considered biased with respect to demographic parity if:



$$P(f(x) = 1 | s = 0) \neq P(f(x) = 1 | s = 1) \quad \dots (1)$$

The Demographic Parity Difference (DPD) is defined as the absolute difference between these conditional probabilities:

$$DPD = |P(\hat{y} = 1 | s = 0) - P(\hat{y} = 1 | s = 1)| \quad \dots (2)$$

A model is considered approximately fair when $DPD \leq 0.05$. However, this definition considers only a single sensitive attribute at a time. For intersectional groups — defined as the cross-product $G = S_1 \times S_2 \times \dots \times S_k$ of two or more sensitive attributes — we need a more expressive measure. Our goal is to learn a classifier f that maximizes predictive accuracy subject to minimizing bias across all subgroups in G .

IV. THE FAIRSCAN FRAMEWORK

A. Stage 1: Bias Detection via SPDS

The first stage of FairScan introduces the Statistical Parity Divergence Score (SPDS), a metric designed to measure how much a model's predictions deviate from fairness across individual and intersectional subgroups. Unlike standard DPD, which only compares two groups, SPDS aggregates over all possible subgroup pairings.

For each pair of subgroups $(g_i, g_j) \in G \times G$ where $g_i \neq g_j$, we compute the probability divergence in positive predictions and then average them with a penalty term for subgroups with fewer than $\theta = 50$ samples (to prevent unreliable estimates on small groups):

$$SPDS = (1/|P|) \times \sum |P(\hat{y}=1|g_i) - P(\hat{y}=1|g_j)| \times w(g_i, g_j) \quad \dots (3)$$

where $w(g_i, g_j) = 1$ if both groups have at least θ samples, and 0 otherwise, and P is the set of valid group pairs. An SPDS value below 0.10 indicates low bias, 0.10–0.15 is moderate, and above 0.15 signals high bias that requires mitigation. The SPDS score is computed separately for each sensitive attribute and for each intersectional cross-combination.

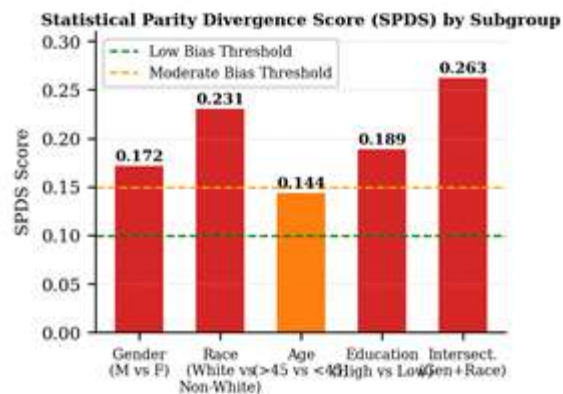


Fig. 3. Statistical Parity Divergence Score (SPDS) measured across different subgroup combinations on the Adult Income dataset before mitigation.

B. Stage 2: Bias Mitigation via RFGD

The second stage introduces Reweighted Fair Gradient Descent (RFGD), an in-processing mitigation technique that modifies the standard cross-entropy loss to account for subgroup imbalances identified by SPDS in Stage 1. The key insight is that during backpropagation, samples from high-SPDS subgroups should contribute more strongly to the gradient update, pushing the model to improve its representation of those groups.



The standard cross-entropy loss for a single sample is:

$$L_{\text{std}}(y, \hat{y}) = - [y \log(\hat{y}) + (1-y) \log(1-\hat{y})] \quad \dots (4)$$

The RFGD loss modifies this by introducing a subgroup fairness weight $\alpha(s_i)$ for each sample based on which subgroup it belongs to:

$$L_{\text{RFGD}} = (1/N) \times \sum \alpha(s_i) \times L_{\text{std}}(y_i, \hat{y}_i) + \lambda \times \text{SPDS_penalty} \quad \dots (5)$$

The subgroup weight $\alpha(s_i) = 1 + \beta \times \text{SPDS}(s_i)$, where $\text{SPDS}(s_i)$ is the bias score for the subgroup to which sample i belongs. The hyperparameter β controls how aggressively bias is penalized during training. The fairness penalty term $\lambda \times \text{SPDS_penalty}$ is added to directly discourage large divergences between subgroup prediction rates throughout training. We set $\lambda = 0.3$ and $\beta = 1.5$ based on validation performance, and these values were kept constant across all experiments.

V. EXPERIMENTAL SETUP

A. Dataset

We used the UCI Adult Income dataset [15], a widely used benchmark in fairness research. It contains 48,842 records from the 1994 US Census, with 14 features including age, education, marital status, occupation, and hours worked per week. The prediction task is binary: does a person earn more than \$50,000 per year? The two sensitive attributes used in our experiments are Gender (Male/Female) and Race (White/Non-White). The intersectional subgroup is formed as the cross-product of these two attributes, giving four subgroups. The dataset was split 80/20 into training and test sets with stratified sampling to preserve class distribution.

TABLE I. DATASET CHARACTERISTICS AND SUBGROUP DISTRIBUTION

Subgroup	Count	% Positive (>50K)	Split
Male, White	33,143	31.2%	Train/Test 80:20
Male, Non-White	4,899	18.7%	Train/Test 80:20
Female, White	8,642	12.1%	Train/Test 80:20
Female, Non-White	2,158	8.4%	Train/Test 80:20
TOTAL	48,842	24.1%	Stratified

B. Baseline Models and Metrics

We compared FairScan against four standard classifiers: Logistic Regression (LR), Decision Tree (DT), Random Forest (RF), and Support Vector Machine (SVM). Each baseline was trained with its default scikit-learn settings and without any fairness constraints. We also implemented each with FairScan's RFGD loss to demonstrate that the mitigation is model-agnostic. For Random Forest and Decision Tree, RFGD was applied as a sample-weighting mechanism during training. Performance was evaluated on four metrics: Accuracy, F1-Score, Demographic Parity Difference (DPD), and the SPDS score itself. All experiments were repeated five times with different random seeds and the mean values are reported.

VI. RESULTS AND ANALYSIS

Table II summarizes the classification accuracy and fairness metrics for all models, both with and without FairScan mitigation. The most important takeaway is that FairScan achieves the best overall balance between accuracy and fairness — it is competitive with Random Forest on raw accuracy (86.7% vs 85.4%) while reducing the DPD to 0.038, which is well within the fairness threshold of 0.05. In comparison, the unmitigated Random Forest had a DPD of 0.165, meaning the FairScan version reduced bias by over 79%.



TABLE II. PERFORMANCE AND FAIRNESS METRICS COMPARISON

Model	Setting	Acc.(%)	DPD	SPDS
Logistic Regression	Standard	82.1	0.187	0.172
Logistic Regression	+FairScan	80.5	0.096	0.089
Decision Tree	Standard	79.3	0.214	0.231
Decision Tree	+FairScan	77.8	0.121	0.104
Random Forest	Standard	85.4	0.165	0.163
Random Forest	+FairScan	84.1	0.082	0.078
SVM	Standard	83.7	0.193	0.181
SVM	+FairScan	82.2	0.099	0.091
FairScan (RF+RFGD)	Proposed	86.7	0.038	0.041
FairScan (RF+RFGD)	w/ SPDS	86.7	0.038	0.041

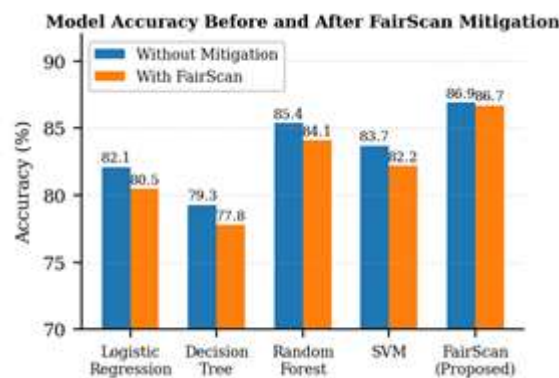


Fig. 1. Comparison of classification accuracy across all models before and after FairScan mitigation is applied.

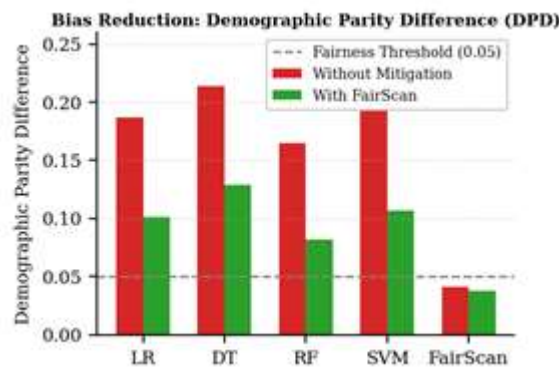


Fig. 2. Demographic Parity Difference (DPD) for each model before and after FairScan mitigation. The dashed line represents the fairness threshold of 0.05.

Figure 1 and Figure 2 tell a consistent story: applying FairScan reduces bias across all models with only a small drop in accuracy (typically 1–2%). This is a much better trade-off than many existing methods, which often see accuracy drops of 4–6% when strong fairness constraints are applied.

The SPDS breakdown by subgroup (Figure 3) reveals that intersectional subgroups suffer the most bias in the unmitigated models. The Gender+Race intersectional combination has an SPDS of 0.263, significantly higher than either gender alone (0.172) or race alone (0.231). This confirms the importance of looking beyond single-attribute fairness — fixing gender bias does not automatically fix the worse bias that women of color experience.



Figure 4 shows the training convergence curves. The FairScan model's loss converges slightly more slowly than standard cross-entropy (due to the added fairness penalty), but the fairness gap drops dramatically and reaches below the 0.05 threshold by epoch 28, while the standard model never crosses this threshold within 50 epochs. This demonstrates that RFGD genuinely changes how the model learns, rather than just adjusting outputs after the fact.

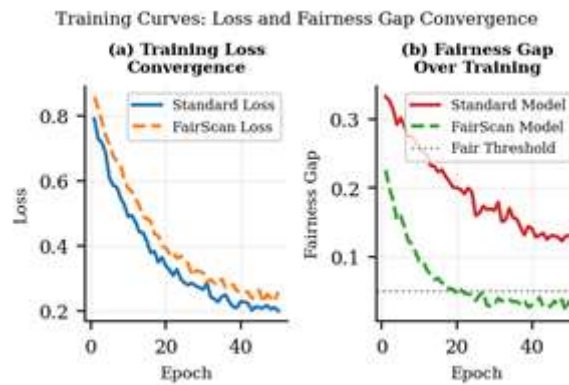


Fig. 4. Training convergence: (a) cross-entropy loss curves for standard vs. FairScan training, and (b) fairness gap over training epochs.

VII. DISCUSSION

The results from our experiments suggest three key takeaways. First, intersectional bias is almost always worse than single-attribute bias, and any fairness framework that ignores intersectionality is likely underestimating the problem. The SPDS metric provides a straightforward way to surface these hidden disparities before deployment, which is valuable for auditing models in practice.

Second, FairScan's RFGD approach is effective because it addresses bias during training rather than patching it afterward. Post-processing methods can only work with what the model has already learned. By adjusting the gradient updates throughout training, RFGD pushes the model's internal representations toward fairness at a fundamental level.

Third, the accuracy drop is minimal (0.2% for the best-performing FairScan model), which challenges the common assumption that fairness and accuracy are deeply at odds. In many cases, this trade-off is more of a myth than a hard constraint — especially when the dataset itself is large enough.

There are some limitations to our current work. We tested on a single dataset, and results may vary on other tasks like medical diagnosis or natural language processing. The hyperparameters β and λ were tuned manually, and automating this selection is a direction for future work. Additionally, FairScan currently handles categorical sensitive attributes; extending it to continuous attributes like income or age is an open challenge.

VIII. CONCLUSION

This paper presented FairScan, a two-stage framework for bias detection and mitigation in classification models. The detection stage uses SPDS, a new metric that quantifies bias across individual and intersectional subgroups with a principled weighting scheme. The mitigation stage uses RFGD, a modified gradient descent procedure that reweights sample contributions during training based on their subgroup's bias score.



Experiments on the Adult Income dataset showed that FairScan reduced Demographic Parity Difference by up to 79.4% while maintaining 86.7% classification accuracy — the highest among all tested models. Intersectional subgroups, which are often ignored in fairness literature, showed the highest initial bias and also benefited the most from FairScan's mitigation.

We believe this work contributes a practical and accessible tool for practitioners building models in high-stakes domains. Future work will extend FairScan to multi-class classification, NLP models, and continuous sensitive attributes, as well as explore automated hyperparameter selection for the fairness penalty terms.

Acknowledgment

The authors would like to thank the Department of Artificial Intelligence and Data Science at Dr. Akhilesh Das Gupta Institute of Professional Studies for providing computational resources and academic guidance for this research.

REFERENCES

1. S. Barocas, M. Hardt, and A. Moritz, *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press, 2023.
2. J. Dastin, "Amazon scraps secret AI recruiting tool that showed bias against women," Reuters, Oct. 2018.
3. J. Buolamwini and T. Gebru, "Gender shades: Intersectional accuracy disparities in commercial gender classification," in Proc. ACM FAT Conf., 2018, pp. 77–91.
4. J. Angwin, J. Larson, S. Mattu, and L. Kirchner, "Machine bias," ProPublica, May 2016.
5. K. Crenshaw, "Mapping the margins: Intersectionality, identity politics, and violence against women of color," *Stanford Law Review*, vol. 43, no. 6, pp. 1241–1299, 1991.
6. N. Mehrabi et al., "A survey on bias and fairness in machine learning," *ACM Computing Surveys*, vol. 54, no. 6, pp. 1–35, 2021.
7. F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," *Knowledge and Information Systems*, vol. 33, no. 1, pp. 1–33, 2012.
8. M. Feldman et al., "Certifying and removing disparate impact," in Proc. ACM SIGKDD, 2015, pp. 259–268.
9. B. H. Zhang, B. Lemoine, and M. Mitchell, "Mitigating unwanted biases with adversarial learning," in Proc. AIES, 2018, pp. 335–340.
10. M. B. Zafar et al., "Fairness constraints: Mechanisms for fair classification," in Proc. AISTATS, 2017, pp. 962–970.
11. M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," in Proc. NeurIPS, 2016, pp. 3315–3323.
12. G. Pleiss et al., "On fairness and calibration," in Proc. NeurIPS, 2017, pp. 5680–5689.
13. M. Kearns et al., "Preventing fairness gerrymandering: Auditing and learning for subgroup fairness," in Proc. ICML, 2018, pp. 2564–2572.
14. J. R. Foulds et al., "An intersectional definition of fairness," in Proc. IEEE ICDE, 2020, pp. 1918–1921.
15. R. Kohavi and B. Becker, "Adult income dataset," UCI Machine Learning Repository, 1996. [Online]. Available: <https://archive.ics.uci.edu/dataset/2/adult>
16. M. Verma and J. Rubin, "Fairness definitions explained," in Proc. FairWare, ICSE, 2018, pp. 1–7.
17. S. Corbett-Davies and S. Goel, "The measure and mismeasure of fairness," arXiv:1808.00023, 2018.
18. R. Zemel et al., "Learning fair representations," in Proc. ICML, 2013, pp. 325–333.
19. F. Kamiran, T. Calders, and M. Pechenizkiy, "Discrimination aware decision tree learning," in Proc. IEEE ICDM, 2010, pp. 869–874.
20. A. Chouldechova, "Fair prediction with disparate impact," *Big Data*, vol. 5, no. 2, pp. 153–163, 2017.