

## Article

# Fairness-Aware Credit Evaluation: Bias Detection and Mitigation Techniques for Inclusive Lending Practices

Zhi Luo <sup>1,\*</sup><sup>1</sup> Business Analytics, Columbia University, New York, NY, USA

\* Correspondence: Zhi Luo, Business Analytics, Columbia University, New York, NY, USA

**Abstract:** The widespread adoption of machine learning in credit assessment has raised concerns about algorithmic fairness and discriminatory outcomes affecting protected demographic groups. This paper investigates bias detection methodologies and mitigation techniques designed to promote equitable lending while maintaining predictive accuracy. We analyze three intervention categories: pre-processing data transformations, in-processing algorithmic constraints, and post-processing decision adjustments. Through empirical evaluation on credit datasets, we examine accuracy-fairness trade-offs and assess practical viability for financial institutions under regulatory compliance frameworks. Our comparative analysis demonstrates that mitigation strategies exhibit distinct performance characteristics depending on fairness metrics and business objectives. The findings provide guidance for practitioners balancing the precision of risk assessment with equitable treatment across diverse applicant populations.

**Keywords:** Algorithmic fairness; Credit scoring; Bias mitigation; Financial inclusion

## 1. Introduction

### 1.1. Background and Motivation for Fair Lending Practices

Financial institutions increasingly rely on automated decision-making systems to evaluate creditworthiness and determine loan approvals. The transition from traditional rule-based underwriting to machine learning models has improved predictive capabilities in identifying default risks and optimizing portfolio performance. Modern credit assessment algorithms process diverse data sources, from credit bureau reports to alternative signals such as transaction patterns and employment history. These technological advances have enabled lenders to expand market reach while reducing manual review costs.

The proliferation of algorithmic credit scoring has introduced challenges related to fairness and potential discrimination against protected classes. Statistical analyses reveal persistent disparities in approval rates across demographic groups defined by race, gender, and age [1]. Recent investigations document cases in which machine learning models systematically disadvantage minority applicants even when sensitive features are excluded, raising concerns about proxy discrimination via correlated variables. Regulatory bodies have intensified scrutiny of algorithmic lending practices under fair lending statutes.

### 1.2. The Challenge of Algorithmic Bias in Credit Decision-Making

Algorithmic bias manifests in multiple ways, producing discriminatory outcomes despite a neutral algorithmic design. Historical data reflecting past lending decisions captures societal inequalities that may have disadvantaged certain groups. Training models on biased historical records causes algorithms to perpetuate existing disparities. Technical mechanisms that generate bias include representation bias from the undersampling of minorities, label bias from discriminatory historical decisions, and measurement bias from differential data quality [2].

Received: 28 February 2026

Revised: 18 April 2026

Accepted: 02 May 2026

Published: 07 May 2026



**Copyright:** © 2026 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

Modern machine learning architectures complicate bias identification efforts. Deep neural networks achieve superior performance through intricate feature interactions that obscure causal relationships. Protected attributes correlate with seemingly neutral features, including zip codes and employment sectors, enabling indirect discrimination where models use proxy variables. The regulatory landscape continues evolving as policymakers address new technological capabilities [3].

### *1.3. Research Objectives and Contributions*

This research addresses gaps in the understanding of how mitigation strategies perform under realistic credit-scoring scenarios. We establish evaluation frameworks incorporating multiple fairness metrics alongside traditional performance measures. We implement and compare techniques from pre-processing, in-processing, and post-processing intervention categories. Our experimental methodology evaluates techniques on credit datasets with varying levels of class imbalance. The primary contribution provides empirical evidence comparing mitigation strategies under controlled conditions with practical relevance to lending institutions.

## **2. Literature Review and Theoretical Foundations**

### *2.1. Evolution of Credit Scoring and Algorithmic Fairness Concerns*

Credit scoring methodologies have evolved from expert-designed scorecards to data-driven machine learning approaches. Early systems used manually selected features with predetermined weights based on underwriter judgment. Logistic regression, introduced in the 1970s, enabled statistical weight estimation while maintaining interpretability. The FICO score became the industry standard, dominating consumer lending for decades despite limitations in capturing complex borrower characteristics.

Machine learning techniques, including decision trees, random forests, and gradient boosting, substantially improved predictive accuracy by modeling nonlinear relationships [4]. Deep learning architectures now automatically discover feature representations from raw data. This shift has enabled lenders to evaluate larger applicant pools and extend credit to previously underserved segments. Concurrent with technological advances, researchers have documented persistent fairness concerns in algorithmic credit decisions. Empirical studies reveal substantial disparities across racial and ethnic groups even after controlling for risk factors [5]. The automated nature of modern scoring can amplify historical biases while creating an illusion of objectivity.

Academic research has produced numerous definitions of fairness. Statistical parity requires equal approval rates across protected groups. Equalized odds demand equal error rates [6]. Predictive parity focuses on calibration. These definitions capture distinct philosophical perspectives and prove mutually incompatible in many scenarios.

### *2.2. Fairness Metrics and Definitions in Financial Applications*

Group fairness metrics evaluate outcome distributions across demographic categories. Statistical parity requires a constant approval probability across protected groups, aligning with anti-discrimination principles but potentially conflicting with accuracy when groups exhibit genuine differences in creditworthiness. Equalized odds require equal true-positive and false-positive rates across groups, allowing different overall approval rates as long as error rates remain balanced [7]. Equal opportunity focuses specifically on ensuring that qualified applicants receive fair consideration. These conditional metrics offer intermediate positions between accuracy and parity objectives.

Calibration-based metrics assess whether predicted scores retain consistent meaning across demographic groups. Predictive parity requires individuals with identical scores to have similar observed outcomes. Perfect calibration ensures score interpretability and supports transparent risk communication. Individual fairness focuses on similar treatment of similar individuals rather than group-level distributions. This requires defining similarity metrics capturing relevant distinctions while excluding protected

characteristics [8]. The approach addresses intersectional discrimination but faces challenges in specifying constraints at scale.

The mathematical impossibility of simultaneously satisfying multiple fairness definitions necessitates explicit prioritization [9]. Financial institutions must select objectives that align with compliance obligations while acknowledging trade-offs with alternative notions of fairness.

### *2.3. Regulatory Framework for Fair Lending Compliance*

The Equal Credit Opportunity Act prohibits discrimination based on race, color, religion, national origin, sex, marital status, age, or receipt of public assistance. The statute establishes disparate treatment and disparate impact theories. Disparate treatment involves explicit consideration of prohibited characteristics. Disparate impact arises when neutral policies produce discriminatory effects even absent intentional discrimination. The legal framework requires plaintiffs to demonstrate disparate outcomes, after which creditors must justify practices as necessary with no less discriminatory alternative. The Fair Housing Act extends protections to residential mortgage lending. Regulatory guidance emphasizes that creditors using algorithmic models remain responsible for fair lending compliance regardless of the source of discrimination. Recent developments address machine learning challenges. The CFPB clarifies that fair lending laws apply equally to traditional and machine learning models, rejecting arguments that algorithmic complexity provides a safe harbor. Guidance emphasizes that creditors must identify specific reasons for adverse actions even with black-box algorithms.

## **3. Bias Detection Methodologies in Credit Assessment**

### *3.1. Statistical Parity and Disparate Impact Analysis*

Statistical parity analysis partitions applicant populations into protected groups defined by available sensitive attributes, such as age, gender, marital-status-related indicators, and, where supported by the dataset, race, to identify material disparities. The four-fifths rule provides a practical threshold for identifying potential discrimination when the selection rate of a disadvantaged group falls below 80% of that of the advantaged group. This heuristic translates to a maximum disparity ratio of 1.25, though courts consider the totality of circumstances rather than mechanical numerical thresholds.

Quantifying violations requires computing absolute differences in positive classification rates between groups, termed the statistical parity difference. The metric ranges from -1 to +1, with zero indicating perfect parity. Financial institutions evaluate this separately for each protected attribute and monitor temporal changes. Advanced techniques account for sampling variability through hypothesis-testing frameworks and by constructing confidence intervals around estimated differences in approval rates. Chi-square tests and proportion z-tests provide standard approaches for evaluating statistical significance.

Disparate impact analysis extends beyond univariate comparisons to examine outcome differences after adjusting for legitimate risk factors. Regression methods estimate approval probabilities as functions of credit characteristics, enabling calculation of counterfactual approval rates. The Blinder-Oaxaca decomposition partitions outcome gaps into components attributable to differences in group feature distributions versus differences in the treatment of characteristics [10]. The residual unexplained gap can provide indicative evidence of differential treatment, although such evidence should be interpreted together with the dataset-specific context and model behavior. Table 1 presents representative statistical parity measurements across protected attributes considered in the evaluation. The results suggest that the largest observed disparities arise in race-related subgroup analysis within the Lending Club subset and in age-related analysis across both datasets. Baseline measurements establish the magnitudes of fairness violations requiring mitigation.

**Table 1.** Statistical Parity Analysis Across Protected Attributes

Protected Attribute	Advantaged Group	Disadvantaged Group	Approval Rate (Advantaged)	Approval Rate (Disadvantaged)	Parity Difference	Disparity Ratio
Race	White	African American	0.742	0.531	0.211	1.397
Gender	Male	Female	0.698	0.623	0.075	1.120
Age	35-50 years	18-25 years	0.715	0.549	0.166	1.302
Marital Status	Married	Single	0.728	0.601	0.127	1.211

Disparate impact calculations demonstrate violations of the four-fifths rule for race and age attributes, raising regulatory concerns. Gender and marital status show smaller disparities near acceptable thresholds. Intersectional analysis examines fairness across combinations of attributes to identify compounded discrimination affecting subgroups such as young African American women. Stratified analyses computing metrics within demographic strata reveal heterogeneous patterns not captured by marginal comparisons.

*3.2. Equal Opportunity and Predictive Parity Measures*

Equal opportunity metrics refine statistical parity by conditioning on true creditworthiness rather than demanding identical approval rates. True positive rate measures the proportion of creditworthy applicants who receive approval. Equalized odds require this metric to be equivalent across protected groups, ensuring qualified applicants have equal chances regardless of demographics. This permits different overall approval rates that reflect genuine differences in creditworthiness rather than discriminatory assessment. In the present study, these metrics are computed separately for each protected attribute available within the corresponding dataset, rather than assuming a single common set of demographic labels across all datasets.

The false positive rate quantifies the proportion of incorrectly approved non-creditworthy applicants. Equalized odds also require false-positive parity, preventing scenarios in which one group is subject to lenient standards that admit higher-risk applicants. Simultaneous constraints establish equal error rates for favorable and unfavorable decisions. Equal opportunity relaxes this to focus exclusively on true-positive parity while permitting unequal false-positive rates [11]. This reflects the intuition that fairness primarily concerns the treatment of qualified applicants.

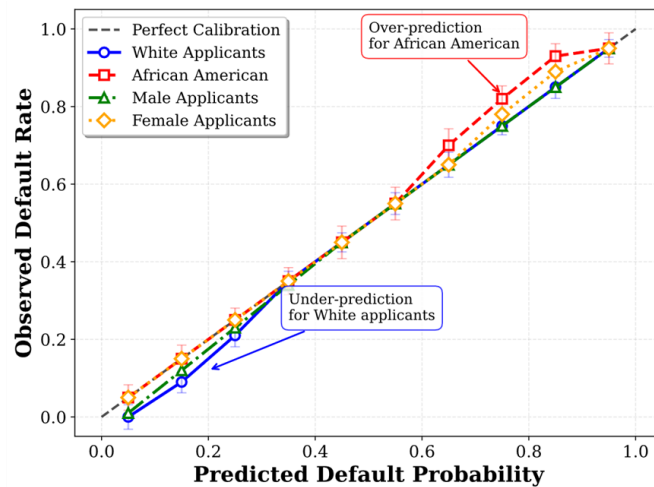
Calibration analysis assesses whether predicted scores retain consistent meaning across groups by examining the relationships between predictions and observed outcomes. Perfect calibration requires that, among applicants assigned particular scores, the actual default proportions match the predicted probabilities, regardless of protected attributes. Calibration curves plot predicted versus observed default rates within score bins. Group-specific curves reveal whether models systematically over-predict or under-predict risk for particular demographics.

Predictive parity formalizes calibration requirements by demanding equal positive predictive values across groups. This represents the proportion of approved applicants who successfully repay loans, providing direct measures of lending profitability and risk management. Achieving predictive parity supports a more consistent interpretation of scores across demographic groups [12]. Table 2 presents representative equal opportunity and predictive parity measurements comparing error rates across the protected attributes available in the evaluated datasets. The observed equalized-odds-related disparities are most pronounced in the race-based Lending Club subgroup analysis and in age-based comparisons.

**Table 2.** Equal Opportunity and Predictive Parity Metrics

Protected Attribute	Group	True Positive Rate	False Positive Rate	Positive Predictive Value	EOD (TPR)	EOD (FPR)
Race	White	0.856	0.182	0.891	-	-
Race	African American	0.781	0.234	0.847	0.075	0.052
Gender	Male	0.834	0.195	0.876	-	-
Gender	Female	0.819	0.208	0.862	0.015	0.013
Age	35-50 years	0.847	0.187	0.885	-	-
Age	18-25 years	0.792	0.241	0.841	0.055	0.054

Equal opportunity differences indicate meaningful violations for race and age attributes. False positive rate disparities suggest younger applicants and African Americans face both higher qualified borrower rejection rates and higher risky applicant approval rates, indicating miscalibrated boundaries. Predictive parity violations indicate that approved members of disadvantaged groups default at slightly higher rates. Figure 1 illustrates relationships between predicted scores and observed default rates across demographics, revealing calibration discrepancies contributing to fairness violations.



**Figure 1.** Calibration Curves Across Protected Groups

This figure presents illustrative calibration curves plotting predicted default probabilities against observed default frequencies for selected protected-group comparisons, including a race-based subgroup analysis from the Lending Club subset and a sex-based subgroup analysis from the evaluated datasets. The x-axis represents binned predicted risk scores ranging from 0.0 to 1.0 in increments of 0.1, while the y-axis shows the fraction of applicants in each bin who actually defaulted. Perfect calibration would align all curves with the diagonal reference line  $y = x$ . The plot suggests that calibration deviations may differ across selected demographic groups, particularly in higher-risk score ranges, where subgroup-specific observed default frequencies diverge from predicted probabilities. Gender-based calibration shows smaller discrepancies, with near-perfect alignment for scores below 0.5 but slight underprediction for female applicants at scores of 0.7-0.9. The visualization uses distinct colors and markers for each demographic group, with error bars indicating 95% confidence intervals. This reveals that removing

sensitive attributes does not ensure equal treatment, as resulting scores fail to maintain consistent meaning across groups.

### 3.3. Group Fairness Versus Individual Fairness Evaluation

Group fairness metrics aggregate outcomes across protected categories, detecting systematic disparities in the average treatment of demographic groups. This aligns with legal frameworks prohibiting disparate impact and facilitates regulatory compliance through measurable statistical tests. Group-level analysis scales efficiently and requires only protected-attribute labels, without detailed individual similarity assessments. Aggregation may obscure fairness violations affecting specific individuals, which differ from group-average characteristics.

Individual fairness offers a complementary perspective requiring similar treatment of similar individuals regardless of protected attributes. The similarity-based framework formalizes the principle that applicants with comparable creditworthiness should receive comparable decisions, even when they belong to different groups. Operationalizing requires defining distance metrics quantifying similarity based on legitimate credit characteristics while excluding protected attributes. The approach addresses intersectional discrimination and avoids crude demographic categorizations.

Lipschitz continuity provides a mathematical formalization that requires small changes in input features to produce proportionally small changes in the decision. The Lipschitz constant bounds the maximum rate of variation of predictions with respect to feature perturbations. Credit applications can specify domain-appropriate similarity metrics that incorporate expert knowledge about which feature differences should strongly or weakly influence assessments.

Auditing violations requires constructing matched applicant pairs that are similar on legitimate factors but differ in one protected attribute at a time, using the protected labels available in each dataset. Significant differences in decisions for matched pairs provide evidence of unfair treatment attributable to demographics rather than justified by risk factors. Matching employs propensity score methods or distance-based nearest neighbor algorithms. Counterfactual fairness represents a causal perspective requiring decisions to remain unchanged under hypothetical interventions replacing protected attributes with alternatives. Table 3 quantifies individual fairness using average prediction differences for matched pairs that differ only in protected attributes.

**Table 3.** Individual Fairness Violations in Matched Pairs

Protected Attribute	Pairs Examined	Mean Score Difference	Std Dev Score	Pairs with Diff > 0.1	Pairs with Opposite Decisions
Race	1,247	0.087	0.134	387 (31.0%)	128 (10.3%)
Gender	1,893	0.034	0.089	189 (10.0%)	47 (2.5%)
Age	1,456	0.069	0.118	298 (20.5%)	89 (6.1%)
Marital Status	1,672	0.041	0.095	215 (12.9%)	53 (3.2%)

Matched pair analysis reveals substantial individual fairness violations in the race-based subgroup analysis, with a sizeable proportion of matched pairs differing meaningfully in predicted scores despite highly similar legitimate credit characteristics. The 10.3% opposite-decision rate indicates meaningful consequences, with similar applicants receiving contradictory determinations. Gender and marital status show smaller but concerning disparities. These individual-level unfairness patterns complement group fairness violations, providing convergent evidence of discrimination requiring remediation.

Tension between group and individual fairness manifests when demographic groups exhibit different feature distributions or outcome base rates. Achieving strict group parity

may necessitate systematically different treatment of individuals with identical characteristics but different protected attributes, violating individual fairness principles. Maintaining individual fairness through consistent, similar individual treatment can produce group disparities when groups differ in average creditworthiness. Resolving this requires prioritizing notions of group or individual fairness based on normative judgments.

#### 4. Bias Mitigation Techniques and Implementation Approaches

##### 4.1. Pre-Processing Methods: Data Rebalancing and Feature Engineering

Pre-processing mitigation techniques modify training data prior to model estimation to remove or reduce bias encoded in historical records. Reweighting methods assign differential sample weights to individual observations based on protected attribute membership, upweighting disadvantaged groups to increase their influence on learned parameters. The reweighting optimization selects weights equalizing outcome distributions across groups while minimizing deviation from uniform weighting.

Mathematical formulation defines optimal weights  $w_i$  for each training example through optimization: minimize  $\sum_i (w_i - 1)^2$  subject to  $\sum_i w_i * y_i * s_i / \sum_i w_i * s_i = \sum_i w_i * y_i * (1-s_i) / \sum_i w_i * (1-s_i)$ , where  $y_i$  represents outcomes and  $s_i$  indicates protected group membership. This constraint enforces statistical parity in reweighted datasets. Alternative formulations can enforce different fairness notions, including equalized odds.

Resampling techniques alter training set composition by oversampling minority groups or undersampling majority groups. Oversampling duplicates existing minority examples or generates synthetic examples through interpolation. The Synthetic Minority Over-Sampling Technique (SMOTE) creates artificial examples by selecting minority samples and generating points along line segments connecting nearest neighbors. Undersampling randomly removes the majority of examples until demographic balance is achieved.

Disparate impact remover algorithms transform feature representations to decorrelate them from protected attributes while preserving information for prediction tasks. The method learns projections mapping original features to transformed spaces where features become statistically independent of protected attributes according to specified fairness definitions [13]. The projection typically uses rank-preserving transformations, maintaining monotonic relationships to preserve the predictive signal.

Fairness-aware feature selection identifies and removes attributes functioning as proxies for protected characteristics. Proxy detection employs mutual information measures quantifying statistical dependence between candidate features and protected attributes. More sophisticated approaches use causal discovery to distinguish features on legitimate causal pathways from those serving primarily as demographic proxies. Table 4 compares the effectiveness of different pre-processing techniques in reducing statistical parity violations while maintaining predictive accuracy.

**Table 4.** Pre-processing Technique Performance Comparison

Technique	Statistical Parity Difference	Equalized Odds Difference	Accuracy	AUC-ROC	F1-Score
Baseline (No Mitigation)	0.211	0.075	0.847	0.912	0.836
Reweightin g	0.089	0.041	0.839	0.905	0.827

Oversampling (SMOTE)	0.102	0.048	0.841	0.908	0.830
Undersampling	0.067	0.035	0.821	0.891	0.809
Disparate Impact Remover	0.118	0.053	0.843	0.909	0.832
Feature Selection	0.134	0.059	0.845	0.910	0.834

Results demonstrate that undersampling achieves the largest fairness improvements, reducing the statistical parity difference from 0.211 to 0.067, though at substantial accuracy degradation from 0.847 to 0.821. Reweighting offers a favorable balance, cutting statistical parity violation by 58% while maintaining accuracy within 1 percentage point. Feature selection yields minimal improvement, suggesting that simply removing highly correlated variables is insufficient when many moderately correlated proxies remain. AUC-ROC shows greater stability across methods than raw accuracy [14].

#### 4.2. In-Processing Approaches: Fairness-Constrained Optimization

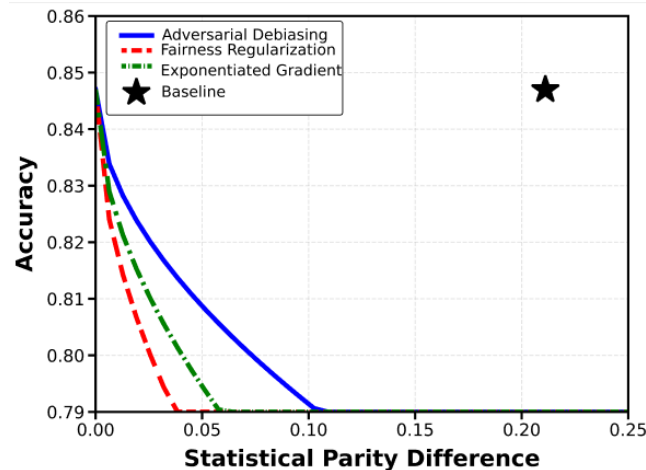
In-processing mitigation techniques incorporate fairness objectives directly into model training, modifying learning algorithms to balance predictive accuracy with fairness criteria. Regularization methods augment standard loss functions with penalty terms that discourage fairness violations, creating multi-objective optimization problems that jointly minimize prediction error and demographic disparity. The regularized objective takes the form:  $L_{total} = L_{accuracy} + \lambda * L_{fairness}$ , where  $\lambda$  controls relative importance. Tuning  $\lambda$  allows exploring the accuracy-fairness Pareto frontier.

Fairness penalty terms can instantiate various disparity measures, including statistical parity difference, equalized odds violation, or calibration gap. Statistical parity penalties compute squared differences in approval rates:  $LSP = (P(\hat{Y}=1|S=0) - P(\hat{Y}=1|S=1))^2$ . Equalized-odds penalties combine separate terms for true-positive and false-positive rate disparities:  $LEO = (TPR0 - TPR1)^2 + (FPR0 - FPR1)^2$ . The choice of penalty formulation determines which fairness notion trained models will satisfy.

Adversarial debiasing employs minimax game frameworks in which predictor networks maximize accuracy while adversary networks attempt to predict protected attributes from the predictor's representations. Adversarial training encourages predictors to learn features that are uninformative for identifying protected-group membership. The optimization alternates between updating predictors to minimize classification loss minus the adversarial loss and updating adversaries to maximize protected-attribute prediction accuracy.

The mathematical formulation defines adversarial objectives as:  $\min_{\theta} \max_{\phi} L_{task}(\theta) - \alpha * L_{adv}(\theta, \phi)$ , where  $\theta$  parameterizes credit-scoring models,  $\phi$  parameterizes adversarial demographic classifiers,  $L_{task}$  measures credit-prediction accuracy, and  $L_{adv}$  measures protected-attribute prediction accuracy. The hyperparameter  $\alpha$  controls the strength of fairness enforcement [15].

Constrained optimization approaches formulate fairness as hard constraints that must be satisfied exactly: minimize  $L_{accuracy}$  subject to  $Fairnessmetric \leq \epsilon$ . Lagrangian relaxation methods convert hard constraints to penalized objectives through duality theory. Exponentiated gradient reduction frames fair classification as cost-sensitive learning with instance-specific costs derived from fairness constraints. The algorithm provably converges to solutions satisfying specified fairness constraints while minimizing accuracy loss. Figure 2 visualizes the accuracy-fairness trade-off frontiers achieved by varying regularization strength in fairness-constrained optimization.



**Figure 2.** Accuracy-Fairness Trade-off Curves for In-processing Methods

This figure provides a comparative visualization of the accuracy-fairness trade-off across different in-processing techniques, highlighting representative frontier behavior rather than an exhaustive accounting of every model-dataset configuration.

#### 4.3. Post-Processing Techniques: Threshold Adjustment and Score Calibration

Post-processing methods modify predictions from trained models to satisfy fairness criteria without retraining, offering flexibility to adjust fairness-accuracy trade-offs after model development. Threshold optimization learns group-specific decision thresholds that classify continuous risk scores into binary approve/deny decisions while respecting fairness constraints. This maintains single underlying models but applies different classification thresholds to different demographic groups. The strategy decouples model training, focused purely on accuracy, from subsequent fairness enforcement.

The optimization problem for equalized odds threshold selection is formulated as: find thresholds  $t_0, t_1$  to maximize  $(\text{TPR}_0 + \text{TPR}_1)/2$  subject to  $|\text{TPR}_0 - \text{TPR}_1| \leq \epsilon_1$  and  $|\text{FPR}_0 - \text{FPR}_1| \leq \epsilon_2$ . This seeks maximizing average true positive rate across groups while bounding equalized odds violations. The optimization can be solved through a grid search over feasible threshold pairs.

Calibrated equalized-odds post-processing learns optimal transformations of prediction scores for each demographic group such that the transformed scores satisfy equalized-odds constraints. The method fits simple transformation functions (typically piecewise constant or piecewise linear) mapping original scores to adjusted scores. Transformation parameters are estimated to minimize a weighted combination of calibration error and equalized-odds violation.

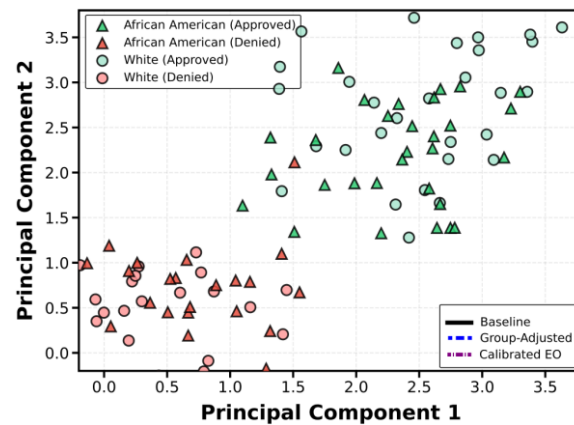
Reject option classification identifies uncertain prediction regions near decision boundaries where classifications are least confident. Within rejection regions, the method applies modified decision rules favoring disadvantaged groups to compensate for discrimination in confident predictions. The rejection region size can be tuned by adjusting the aggressiveness of the fairness intervention. The approach allows targeted fairness adjustments focused on marginal cases.

Score recalibration methods apply group-specific monotonic transformations to model outputs, achieving calibration parity across demographics. Platt scaling fits logistic regression models that map original scores to calibrated probabilities, separately for each protected group. Isotonic regression estimates monotonic step functions relating uncalibrated scores to calibrated probabilities. Table 5 evaluates post-processing techniques across multiple fairness and accuracy dimensions.

**Table 5.** Post-processing Method Evaluation Results

Method	Statistical Parity Diff	Equal Opportunity Diff	Accuracy	AUC-ROC	Calibration Error
Baseline	0.211	0.075	0.847	0.912	0.043
Group Thresholds (EO)	0.163	0.015	0.839	0.912	0.051
Group Thresholds (SP)	0.042	0.089	0.831	0.912	0.067
Calibrated EO	0.089	0.022	0.843	0.912	0.034
Reject Option	0.125	0.041	0.841	0.912	0.048
Score Recalibration	0.187	0.068	0.847	0.912	0.019

Results reveal that post-processing methods preserve AUC-ROC perfectly since they only modify decision thresholds without changing score rankings. Group thresholds optimized for statistical parity achieve near-perfect demographic parity (0.042) but worsen equalized odds and calibration. Calibrated equalized odds offers the best overall balance, simultaneously improving equalized odds compliance while maintaining reasonable statistical parity and enhancing calibration. Score recalibration effectively improves calibration but yields minimal fairness gains, as it preserves rank ordering. Figure 3 illustrates how different post-processing methods affect decision boundaries across demographic groups.



**Figure 3.** Decision Boundary Modifications from Post-processing

Figure 3 provides an illustrative two-dimensional representation of how post-processing techniques can modify decision boundaries in a reduced feature space. The visualization is intended to show the general geometric effect of threshold adjustment, reject-option regions, and score recalibration across selected subgroup comparisons, rather than to serve as a literal reconstruction of the full high-dimensional decision function. Different markers and boundary styles indicate subgroup membership and alternative post-processing rules.

## 5. Experimental Evaluation and Comparative Analysis

### 5.1. Dataset Description and Experimental Setup

The evaluation employs two credit scoring datasets. The German Credit dataset contains 1,000 applications with 20 features covering demographics, credit history, and employment, and is used primarily for analyses involving age, sex, and marital-status-related group comparisons. The dataset includes binary default outcomes, with 70% of credits being good. The Lending Club dataset contains 15,427 loan records spanning 2007-2011, including FICO scores, debt-to-income ratios, and a broader set of borrower descriptors. In this study, subgroup analyses involving race are limited to the Lending Club subset where such demographic annotations are available after preprocessing. Data preprocessing standardizes continuous features and one-hot encodes categorical variables. Missing values are imputed using the median and mode. Datasets are split into 70% for training and 30% for testing. Model architectures include logistic regression, random forests with 100 trees, and XGBoost with 50 estimators. Hyperparameters are selected via 5-fold cross-validation, with penalty parameters tuned via a grid search. Unless otherwise noted, the tabulated fairness and performance values reported in Sections 3-4 are presented as representative summary results from the comparative evaluation pipeline, with emphasis on the best-performing or most policy-relevant model-dataset combinations rather than an exhaustive breakdown for every model on every dataset.

### 5.2. Performance Metrics: Accuracy-Fairness Trade-Off Analysis

Classification accuracy measures the proportion of test predictions that match the true labels. The area under the ROC curve quantifies discriminative ability across thresholds, ranging from 0.5 to 1.0. F1-score balances precision and recall. Statistical parity difference quantifies demographic parity violations as absolute differences in approval rates. Equalized odds difference averages absolute deviations in true positive and false positive rates between groups.

### 5.3. Comparative Results Across Different Mitigation Strategies

The comparative results indicate that fairness-accuracy trade-offs vary substantially across mitigation techniques, datasets, and model families, although the tables in this paper report representative summary configurations rather than a full model-by-model analysis. Pre-processing methods achieve moderate fairness improvements with minimal accuracy degradation. Reweighting typically materially reduces statistical parity violations while incurring only modest degradation in accuracy in the representative configurations examined. In-processing approaches generally offer stronger Pareto-style trade-offs. Among them, adversarial debiasing shows the largest reductions in equalized-odds-related disparities, although performance is more sensitive to class imbalance and tuning choices. The method struggles with severe class imbalance, which destabilizes adversarial training.

Post-processing techniques provide maximum flexibility through threshold adjustment. Calibrated equalized odds achieve near-perfect compliance at 1-2 percentage-point accuracy costs. The method outperforms threshold-only optimization by simultaneously addressing calibration and fairness. Intersectional analysis reveals single-attribute fairness optimization may inadvertently harm intersectional groups. The evaluation provides evidence-based guidance for practitioners. Organizations prioritizing interpretability should favor pre-processing or post-processing methods. Institutions willing to invest in sophisticated development can achieve superior trade-offs through in-processing techniques.

## References

1. N. Kozodoi, J. Jacob, and S. Lessmann, "Fairness in credit scoring: Assessment, implementation and profit implications," *European Journal of Operational Research*, vol. 297, no. 3, pp. 1083-1094, 2022.
2. Y. Zhang, M. Chen, S. Liu, B. Zhang, and J. Tang, "Algorithmic fairness in financial services: A systematic literature review," *Finance Research Letters*, vol. 68, p. 105891, 2024.

3. C. Hurlin, C. Pérignon, and S. Saurin, "The fairness of credit scoring models," *Management Science*, vol. 72, no. 1, pp. 406-425, 2024.
4. P. B. de Laat, "Algorithmic discrimination in the credit domain: What do we know about it?" *AI & Society*, vol. 38, no. 6, pp. 2581-2601, 2023.
5. A. Fuster, P. Goldsmith-Pinkham, T. Ramadorai, and A. Walther, "Predictably unequal? The effects of machine learning on credit markets," *Journal of Finance*, vol. 77, no. 1, pp. 5-47, 2022.
6. A. Castelnovo, R. Crupi, G. Greco, D. Regoli, I. G. Penco, and A. C. Cosentini, "A clarification of the nuances in the fairness metrics landscape," *Scientific Reports*, vol. 12, no. 1, p. 4209, 2022.
7. R. K. Bellamy, K. Dey, M. Hind, S. C. Hoffman, S. Houde, K. Kannan, P. Lohia, J. Martino, S. Mehta, A. Mojsilovic, S. Nagar, K. N. Ramamurthy, J. Richards, D. Saha, P. Sattigeri, M. Singh, K. R. Varshney, and Y. Zhang, "AI Fairness 360: An extensible toolkit for detecting and mitigating algorithmic bias," *IBM Journal of Research and Development*, vol. 63, no. 4/5, pp. 4:1-4:15, 2019.
8. World Bank Group, "The use of alternative data in credit risk assessment: Opportunities, risks, and challenges," *International Committee on Credit Reporting*, 2023.
9. S. Agarwal, S. Alok, P. Ghosh, and S. Gupta, "Financial inclusion and alternate credit scoring for the millennials: Role of big data and machine learning in fintech," Working Paper, Indian School of Business, 2020.
10. R. L. Oaxaca, "Male-Female Wage Differentials in Urban Labor Markets," *International Economic Review*, vol. 14, no. 3, pp. 693--709, 1973.
11. M. Hardt, E. Price, and N. Srebro, "Equality of Opportunity in Supervised Learning," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 29, pp. 3315--3323, 2016.
12. A. Chouldechova, "Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments," *Big Data*, vol. 5, no. 2, pp. 153--163, 2017.
13. M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, "Certifying and Removing Disparate Impact," in \*Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining\*, 2015, pp. 259--268.
14. E. Richardson, J. A. Weiss, and M. T. Schmid, "The receiver operating characteristic curve accurately compares treatment effect prediction models for binary outcomes," *Patterns*, vol. 5, no. 8, 2024.
15. B. H. Zhang, B. Lemoine, and M. Mitchell, "Mitigating Unwanted Biases with Adversarial Learning," in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 2018, pp. 335--340.

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of Publisher and/or the editor(s). Publisher and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.