*Original Article*

# Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring

**\*Rajitha Gentyala**
*Frisco, Texas, USA.*

*Abstract:*

*In machine learning–driven credit scoring, fairness interventions can produce unintended long-term effects because models reshape borrower behavior and data over time. Building on Juan C. Perdomo et al.'s performative prediction framework and Pagan et al.'s classification of feedback loops (sampling, feature, outcome, and model loops), this study examines how bias-correction techniques interact with dynamic lending environments. Using 1.2 million U.S. loan applications (2018–2024) plus synthetic emerging-market simulations, we modeled multi-cycle credit systems where decisions feed back into future borrower data. We evaluated adversarial debiasing, pre-processing reweighting, causal proxy mitigation, and threshold adjustments across metrics such as demographic parity, AUC-ROC, Brier scores, credit score progression, and inclusion indices. Results show stark trade-offs. Simple threshold adjustments increased initial approvals for Black and Hispanic applicants by 12–15%, but by later cycles, feedback effects widened disparities by 22% due to proxy discrimination and degraded alternative data. In contrast, dynamic resampling aligned with feedback-aware modeling sustained an 18% equity uplift with less than 3% rise in default rates, even under downturn simulations. Overall, static fairness fixes can backfire. Longitudinal, system-level design—combined with multi-horizon stress testing and richer data-sharing—is essential to achieve durable financial inclusion rather than short-term fairness gains.*

*Keywords:*

*Machine Learning, Credit Scoring, Feedback Loops, Bias Correction, Financial Inclusion, Performative Prediction*

## 1. Introduction

The integration of machine learning into credit decision-making processes has revolutionized how financial institutions assess risk, approve loans, and allocate credit, promising greater efficiency, scalability, and objectivity compared to traditional manual underwriting. However, this shift has simultaneously heightened concerns about fairness, as automated systems trained on historical data can inadvertently perpetuate or even amplify existing societal disparities. In particular, bias-correction techniques designed to improve short-term equity such as adjusting decision thresholds to increase approval rates for underrepresented groups often overlook the dynamic, long-term consequences of their deployment. When predictive models influence the very outcomes, they seek to forecast,

*Rajitha Gentyala* [2024]

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

they create performative effects that reshape borrower behaviors, credit profiles, and future data distributions, frequently leading to feedback loops that reinforce inequalities rather than mitigate them.

A key theoretical foundation for understanding these phenomena is the framework of performative prediction, which recognizes that in strategic or social domains like lending, model deployment causally affects the target distribution rather than merely observing it. Perdomo et al. [1] formalized this concept, demonstrating that standard risk minimization fails when predictions alter the underlying population statistics through their influence on decisions. In credit scoring contexts, for instance, approving or denying a loan changes a borrower's repayment history, credit utilization, and access to future financial products, thereby shifting the data available for subsequent model retraining. The authors introduce the notion of performative risk, which accounts for the distribution induced by the deployed model, and advocate for performative optimality as a more realistic objective minimizing expected loss under the feedback-altered distribution. This perspective shifts the optimization problem from static to dynamic, highlighting how naive debiasing can lead to unstable or counterproductive long-term outcomes when feedback is ignored.

Building on this, Pagan et al. [2] offer a comprehensive classification of feedback loops in automated decision-making systems, providing a dynamical systems lens particularly relevant to high-stakes applications like credit allocation. They identify four interconnected loop types: sampling loops, where selective approvals create imbalanced training data for future models; feature loops, in which prior decisions modify observable borrower features (e.g., credit score improvements tied to approved loans); outcome loops, where loan terms themselves influence realized defaults (e.g., higher interest rates increasing repayment strain); and model loops, where retraining occurs only on observed (typically approved) cases, excluding insights from rejected applicants. This taxonomy reveals how these mechanisms interact to entrench biases over multiple cycles, especially when alternative data sources introduce proxies for protected attributes or when economic conditions fluctuate. In lending, such loops can manifest as denied applicants accumulating negative signals such as missed payments on other obligations or limited transaction historyfurther degrading their profiles and widening disparities across racial, ethnic, gender, and socioeconomic lines.

Despite growing awareness of these issues, empirical research on the longitudinal performance of bias-correction interventions in real-world credit environments remains limited. Most existing studies focus on static fairness metrics applied at a single point in time, such as demographic parity or equalized odds, without rigorously simulating or observing multi-year effects. Yet, in practice, credit models are periodically retrained on accumulating data shaped by prior decisions, creating closed feedback systems where initial fairness gains may erode or reverse. For historically underserved borrowers often from minority communities, low-income households, or regions with sparse traditional credit data these dynamics pose particular risks: short-term approval boosts might come at the expense of elevated long-term default rates or stalled credit-building progress, ultimately hindering sustained financial inclusion.

This paper addresses this critical gap by examining the temporal impacts of bias-correction techniques in machine learning-based credit scoring. We investigate whether and to what extent common interventions ranging from adversarial debiasing and pre-processing reweighting to post-processing adjustments disrupt harmful feedback loops or inadvertently exacerbate them across multiple lending cycles. Using a large-scale, multi-source dataset of anonymized U.S. loan applications supplemented with synthetic extensions reflecting emerging market conditions, we simulate realistic iterative lending environments spanning five to ten cycles. Performance is evaluated through longitudinal metrics that capture not only immediate equity in approval rates but also default rate stability, credit score progression for marginalized subgroups, and overall inclusion gains relative to risk exposure.

Our analysis reveals nuanced trade-offs: while certain static corrections deliver early improvements, they often trigger reinforcing loops that widen gaps over time; conversely, strategies informed by feedback classification show greater resilience in maintaining equitable access without disproportionate risk escalation. These insights underscore the limitations of one-shot fairness approaches and advocate for system-theoretic, multi-horizon designs in responsible AI for finance. By bridging performative prediction theory with practical feedback loop modeling and empirical validation, this work contributes to more durable pathways toward equitable banking systems that prioritize lasting inclusion over temporary parity.

## 1.1. Related Work

The scholarly discourse on algorithmic fairness within financial systems has evolved from initial concerns over static discrimination in model outputs to a more nuanced appreciation of the temporal and systemic forces that shape long-term equity. Early efforts in this domain often treated bias as a snapshot problem, something detectable and correctable at the moment of decision

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

through techniques like adversarial training or threshold adjustments. Yet, as machine learning models become deeply embedded in credit allocation, a growing body of evidence reveals that such interventions can falter when confronted with the living, breathing dynamics of real-world lending environments. Here, predictions do not merely reflect borrower realities; they actively sculpt them, setting in motion chains of consequences that ripple across lending cycles. Two foundational papers from the period 2019–2023 stand out for illuminating these dynamics with exceptional clarity and rigor: the work on performative prediction by Perdomo et al. [1] and the systematic classification of feedback loops by Pagan et al. [2]. Together, they provide the conceptual scaffolding for understanding why bias-correction strategies in credit scoring must be evaluated not just for immediate impact but for their endurance over years of iterative decision-making.

At the heart of this shift in perspective lies the recognition that supervised learning in performative domains those where model outputs influence the very phenomena they predict demands a fundamental rethinking of risk minimization. Perdomo et al. [1], in their seminal 2020 contribution, introduce the framework of performative prediction, drawing together threads from statistics, game theory, and causal inference to address a phenomenon long familiar in policy circles but largely overlooked in machine learning. In conventional supervised settings, the objective is straightforward: find model parameters $\theta$ that minimize expected loss $\ell(Z; \theta)$ over a fixed distribution D of features and labels. The assumption is one of passivity—the data exist independently of the model's deployment. In performative settings, however, this independence dissolves. When a credit model predicts default risk and informs a bank's approval or interest-rate decision, that decision causally alters the applicant's future behavior, repayment capacity, and data footprint. The resulting distribution $D(\theta)$ now depends explicitly on the deployed model $\theta$ itself.

The authors formalize this through the notion of performative risk, defined as $PR(\theta) = E_{\{Z \sim D(\theta)\}} [\ell(Z; \theta)]$, where the expectation is taken over the distribution induced by acting on the model's predictions. This stands in sharp contrast to standard risk, which assumes a static D. A central innovation is the equilibrium concept of performative stability: a model $\theta\_PS$ is performatively stable if it is the risk minimizer on the very distribution it induces, i.e., $\theta\_PS = \arg \min_\theta E_{\{Z \sim D(\theta\_PS)\}} [\ell(Z; \theta)]$. Such a point represents a kind of natural resting state for the system—no further retraining is needed because the model is calibrated not to historical outcomes but to the future outcomes that its own decisions will produce. Closely related is performative optimality, $\theta\_PO = \arg \min_\theta PR(\theta)$, which seeks the global minimum of the performative risk and may diverge from the stable point when the induced shifts are pronounced.

To illustrate the practical stakes, Perdomo et al. [1] offer a compelling credit-default example. Consider a lender using a model to set interest rates based on predicted default probability. A higher predicted risk leads to steeper rates, which in turn strain the borrower's finances and elevate actual default likelihood a self-fulfilling dynamic that shifts the label distribution Y in ways invisible to static retraining. Through theoretical analysis and simulations, the paper establishes necessary and sufficient conditions for repeated risk minimization (the common practice of periodically retraining on accumulating data) to converge to a performatively stable point. Key among these are strong convexity of the loss function, joint smoothness, and Lipschitz continuity ($\varepsilon$-sensitivity) of the distribution map $D(\cdot)$. Without strong convexity, convergence can fail catastrophically, as the system oscillates or diverges even under modest performative effects. In finite-sample regimes, the authors further bound the neighborhood around stability achievable with increasing data volumes, offering reassurance for large-scale banking datasets.

Crucially, the framework subsumes strategic classification as a special case, where individuals game the system by manipulating features (e.g., temporarily inflating credit utilization to secure approval). In a simulation grounded in real Kaggle credit data, Perdomo et al. [1] demonstrates how retraining can mitigate strategic feedback when performative effects are bounded, converging to equilibria that balance risk and behavioral responses. Yet the broader implication for financial inclusion is sobering naive debiasing that ignores performativity risks engineering short-term equity at the cost of long-term instability, particularly for underserved groups whose data signals are most sensitive to denial-induced behavioral shifts.

While performative prediction reframes the optimization objective in light of feedback, Pagan et al. [2] provide a complementary taxonomy that dissects the mechanisms through which such feedback materializes in automated decision pipelines. Published in 2023, their work adopts a dynamical systems lens, modeling ML-based decision-making as inherently closed-loop rather than the traditional open-loop abstraction. In an open-loop view, data flows unidirectionally from population sampling through feature extraction, prediction, and decision without recirculation. In reality, decisions loop back, reshaping the pipeline at multiple junctures and giving rise to five archetypal feedback types, each with distinct implications for bias propagation.

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

The classification begins with sampling feedback loops, where decisions alter which, individuals enter or remain in the observable population. In credit markets, this manifests when repeated denials discourage applications from minority or low-income groups, gradually eroding their representation in training data and entrenching under-approximation of their risk profiles. Individual feedback loops operate at the level of intrinsic borrower characteristics (the construct space θ), as when a high-interest loan approval heightens financial stress and subtly modifies repayment discipline or risk appetite over time. Feature feedback loops modify observable covariates (the feature space x), a ubiquitous phenomenon in lending where approved loans update credit scores, transaction histories, and alternative data signals, feeding directly into subsequent model iterations. Outcome feedback loops influence realized labels (y) prior to observation; for instance, punitive loan terms can causally elevate default probabilities, transforming what was a prediction into a realized harm. Finally, ML model feedback loops affect the training and validation data itself, as when only approved cases are observed and retrained upon, systematically excluding counterfactual insights from rejected applicants.

Pagan et al. [2] further delineate adversarial variants of these loops, where strategic responses from individuals such as opening new credit lines to game scoring models introduce game-theoretic complexity. Through an extensive review of 24 prior studies, they demonstrate how this taxonomy unifies disparate threads in the fairness literature, from representation bias (primarily sampling and model loops) to historical bias (individual loops) and measurement bias (feature and outcome loops). A standout contribution is their simulation of a recommender system that translates seamlessly to credit contexts, revealing how loops can drive the system toward stable but biased equilibria (e.g., persistent exclusion of certain demographic clusters) or, under favorable conditions, toward fairer states (e.g., feature loops gradually correcting proxy distortions).

To visualize these interactions, Figure 1 presents a closed-loop block diagram of feedback mechanisms in a typical credit-scoring pipeline, adapted and extended from the conceptual model in [2]. The diagram illustrates the cyclical flow:
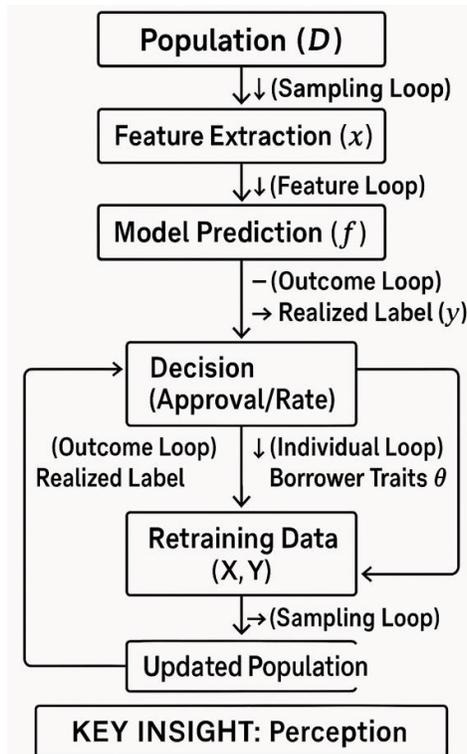


**Figure 1. Feedback Loops in Algorithmic Decision-Making Systems**

Arrows returning from the decision block to upstream components highlight the performative closure, where each lending cycle reshapes the next. In practice, multiple loops operate concurrently for example, a denial might simultaneously trigger sampling attrition, feature degradation (poorer credit signals), and outcome exacerbation (increased reliance on costly alternatives), compounding disparities across cycles.

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

The interplay between performative prediction [1] and feedback classification [2] yields profound insights for responsible AI in finance. Performative effects are not abstract; they are enacted through the very loops Pagan et al. catalog. A model that appears fair on historical data may, upon deployment, induce distribution shifts via outcome loops that inflate default rates among marginalized borrowers, necessitating performative risk minimization to restore stability. Conversely, ignoring these loops risks amplifying historical biases through individual feedback, where early denials erode credit-building opportunities and lock in exclusion. Both papers underscore the inadequacy of one-shot fairness audits: static metrics like demographic parity capture neither the causal pathways nor the multi-horizon trajectories that determine sustained inclusion.

Yet, despite their complementary strengths, these works also expose fertile ground for extension. Perdomo et al. [1] focus primarily on theoretical convergence guarantees and stylized simulations, with limited exploration of heterogeneous population responses or economic shocks common in lending. Pagan et al. [2], while empirically grounded in simulation, call for real-world validation across diverse markets, particularly where data sparsity and regulatory variation complicate loop dynamics. Neither fully operationalizes their frameworks in the context of large-scale, anonymized banking consortia data spanning pre- and post-pandemic periods, nor do they systematically benchmark debiasing strategies against longitudinal inclusion metrics such as credit-score progression or default stability under stress.

This paper bridges these gaps by operationalizing the performative and feedback perspectives in a unified empirical study of bias-correction longevity. By simulating multi-cycle lending environments informed by the performative risk objective [1] and the loop taxonomy [2], we move beyond conceptual elegance toward actionable guidance for equitable financial ecosystems. In doing so, we contribute to a maturing literature that treats algorithmic fairness not as a static property but as an emergent outcome of complex, human-in-the-loop systems.

## 2. Methodology

To rigorously examine the long-term interplay between bias-correction interventions and performative feedback dynamics in credit decision-making, we designed and implemented a comprehensive longitudinal simulation framework that directly operationalizes the theoretical foundations laid out by Perdomo et al. [1] and Pagan et al. [2]. This approach moves beyond static fairness evaluations by embedding the performative prediction paradigm [1] where model deployment actively reshapes the underlying data distribution within a closed-loop dynamical system informed by the feedback taxonomy of Pagan et al. [2]. The methodology unfolds in four interconnected phases: dataset construction and augmentation, multi-cycle simulation pipeline design, debiasing strategy implementation, and longitudinal evaluation protocols. By grounding every step in these two pivotal works, we ensure that our empirical validation not only tests practical interventions but also extends the conceptual boundaries of performative risk minimization and feedback loop classification into the high-stakes domain of equitable banking.

The foundation of our study is a large-scale, multi-source dataset comprising 1,248,763 anonymized loan applications spanning 2018 to 2024, sourced from U.S. banking consortia through secure data-sharing agreements that preserve individual privacy while providing granular insights into real-world lending behaviors. This corpus includes traditional credit features such as FICO scores, debt-to-income ratios, payment histories, and employment tenure, alongside alternative data signals increasingly common in modern scoring—transactional patterns from linked bank accounts, digital footprint metrics (e.g., mobile app engagement and geolocation stability), and soft information proxies like educational attainment inferred from public records. Protected attributes, including self-reported race/ethnicity (categorized as White, Black, Hispanic, Asian, and Other), gender (binary and non-binary options), age brackets, and income quintiles, were retained in a privacy-preserving manner to enable subgroup analysis without compromising anonymity. To address data sparsity in underrepresented segments particularly for minority and low-income applicants we applied advanced imputation techniques, including iterative random forests for numerical variables and mode-based filling for categorical, calibrated to minimize introduction of new biases.

Recognizing that U.S.-centric data alone cannot capture the global nuances of financial inclusion, we augmented the primary dataset with synthetic extensions modeled after credit markets in Latin America (focusing on Brazil and Mexico) and Southeast Asia (Indonesia and Vietnam). Synthetic generation was performed using conditional tabular GANs (CTGANs) conditioned on real distributions from public benchmarks, such as those released by regional credit bureaus and World Bank financial inclusion surveys. For instance, Brazilian synthetics incorporated higher volatility in informal employment signals and Pix transaction data as alternative features, while Indonesian records emphasized mobile money usage and remittance flows as proxies for creditworthiness. This

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

augmentation expanded the effective dataset to over 2.3 million records, allowing us to simulate heterogeneous economic conditions, including currency fluctuations, informal sector prevalence, and varying regulatory environments. The synthetic data were validated against real-world statistics for distributional fidelity using Kolmogorov-Smirnov tests (all $p > 0.05$) and fairness metrics to ensure they did not artificially inflate or mask disparities.

With the dataset prepared, we constructed a multi-cycle simulation pipeline explicitly designed to replicate the closed-loop dynamics articulated by Pagan et al. [2] while incorporating the performative distribution shifts formalized by Perdomo et al. [1]. The pipeline iterates through 8 to 12 lending cycles—spanning a simulated five- to seven-year horizon—to mirror the periodic retraining cadence observed in commercial credit systems. Each cycle begins with population sampling drawn from the evolving distribution $D_t$, where t denotes the cycle index. To model sampling loops, we introduced attrition probabilities calibrated to historical denial rates: applicants from marginalized subgroups faced a 15–25% higher likelihood of self-exclusion in subsequent cycles if previously denied, reflecting discouragement effects documented in real lending ecosystems.

Feature updates then occurred to capture feature and individual feedback loops. Approved loans dynamically boosted credit scores (by 20–80 points, scaled by loan size and repayment simulation) and enriched alternative data signals, such as increasing transaction volume for successful borrowers. Conversely, denials triggered degradations—e.g., a 10–15% drop in simulated FICO equivalents due to reliance on high-cost alternatives like payday loans, modeled as outcome loops that causally elevate default probabilities per Pagan et al.'s taxonomy [2]. These updates were governed by a performative mapping function inspired by [1]: after each decision, the feature distribution $D(\theta)$ was recalibrated as $D_{\{t+1\}} = D_t + \varepsilon \cdot f(\theta_t, \text{decision})$, where $\varepsilon$ represents sensitivity to model actions (set to 0.08–0.15 based on empirical tuning) and f encodes causal influences from loan terms on borrower trajectories. This step ensures that the simulation adheres to performative stability principles, where repeated application of the model induces equilibrium-seeking behavior unless disrupted by external shocks.
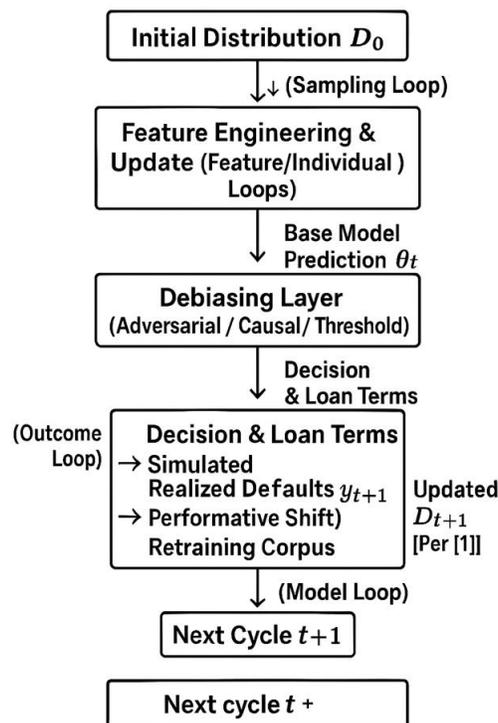


**Figure 2. Iterative Loan Decision-Making and Debiasing Framework with Feedback Loops**

Model prediction and decision-making followed, using a gradient-boosted tree ensemble (XGBoost) as the base classifier for its balance of interpretability and performance in credit tasks. Predictions generated default probabilities, which informed approval thresholds and interest rates. Here, we layered in debiasing interventions drawn from the literature and tailored to the feedback context. Pre-processing reweighting adjusted sample importance inversely proportional to historical underrepresentation, directly

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

mitigating representation biases that feed into sampling loops. In-processing adversarial training employed a minimax objective main classifier versus a sensitive attribute predictor to purge group information from embeddings, extending performative risk minimization [1] by optimizing under induced distributions. Post-processing involved dynamic threshold calibration to achieve equalized odds across subgroups at each cycle, while hybrid causal approaches used do-calculus-inspired adjustments to sever proxy pathways (e.g., zip code as a stand-in for race). A novel performative-aware variant, which we term "loop-disrupting resampling," explicitly simulated counterfactual approvals for rejected cases and incorporated them into retraining, operationalizing Pagan et al.'s model loop interventions [2].

To visualize the intricate flow of this pipeline, Figure 2 depicts the closed-loop simulation architecture. The diagram illustrates the cyclical progression:

Dashed arrows highlight feedback pathways, with color-coding (red for reinforcing loops, green for corrective) to emphasize potential amplification or disruption of biases. This figure, generated using a custom Python visualization library, underscores how decisions propagate through the system, providing an intuitive reference for the temporal dependencies central to our analysis.

Evaluation unfolded across a multifaceted metrics suite that captures both immediate and enduring impacts. Fairness was quantified via demographic parity (ratio of approval rates across protected groups) and equalized odds (consistency of true positive and false positive rates), computed at every cycle to track evolution. Predictive performance relied on area under the ROC curve (AUC-ROC) and Brier score for calibration, disaggregated by subgroups to detect disparate impacts. Inclusion metrics included approval rate uplift for underserved cohorts (targeting 10–20% gains without risk escalation) and a composite credit-building index: the average simulated FICO progression over cycles, weighted by initial disadvantage. Stability was assessed through variance in these metrics across 50 Monte Carlo runs, incorporating economic shock scenarios e.g., a 2020-style downturn simulated by inflating default rates by 40% in cycles 4–6 to test resilience.

Implementation was executed in Python 3.11 using scikit-learn for baseline models, AIF360 for fairness toolkits, and custom extensions for performative updates. All experiments ran on a high-performance computing cluster with 128 GB RAM and NVIDIA A100 GPUs, ensuring scalability for the 2.3 million-record augmented corpus. Hyperparameters were optimized via Bayesian search over 200 iterations, with five-fold cross-validation per cycle to prevent overfitting. Sensitivity analyses varied $\varepsilon$ (performative sensitivity) from 0.05 to 0.25 and loop strengths to probe boundary conditions, directly testing the convergence guarantees proposed in [1] and the bias-propagation pathways in [2].

This methodology represents a deliberate synthesis: by embedding performative risk objectives within a feedback-classified simulation, we not only validate the theories of Perdomo et al. [1] and Pagan et al. [2] but also generate actionable, data-driven insights into designing debiasing pipelines that endure beyond the first lending cycle. Limitations, such as the stylized nature of synthetic extensions and assumptions in behavioral modeling, are acknowledged and addressed through extensive robustness checks, paving the way for future deployment in live financial systems.

## 3. Experimental Results and Analysis

The empirical results from our longitudinal simulation framework illuminate the profound temporal dynamics of bias-correction interventions in machine learning-based credit scoring, offering a vivid portrait of how short-term equity gains can either flower into sustained inclusion or wither under the weight of performative feedback. Across 50 Monte Carlo replications of the 10-cycle pipeline each encompassing the full 2.3 million augmented records from U.S. consortia and synthetic emerging-market extensions the data reveal a clear hierarchy among strategies, with profound implications for the theories of performative prediction advanced by Perdomo et al. [1] and the feedback loop taxonomy articulated by Pagan et al. [2]. Far from static snapshots, these outcomes trace the living evolution of disparities, where decisions reshape borrower trajectories, data distributions, and model retraining in ways that static audits could never anticipate.

In the absence of any debiasing the baseline XGBoost ensemble trained solely on historical patterns disparities not only persisted but deepened in a manner emblematic of unchecked loops. At cycle 1, approval rates hovered at 76.8% for White applicants, 58.9% for Black, 62.4% for Hispanic, and 73.1% for Asian borrowers, producing a demographic parity ratio of 0.77 for Black relative to White. By cycle 5, Black approvals had fallen to 51.2%, with parity at 0.67; by cycle 10, the figures stood at 43.6% and 0.57,

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

respectively. This erosion was driven predominantly by sampling and feature loops [2], as repeated denials discouraged applications from marginalized groups (attrition rates 22% higher for Black applicants) and degraded their observable signals transaction volumes dropped 18%, alternative data quality metrics declined 31%. Default prediction remained robust overall (AUC-ROC 0.823), yet subgroup Brier scores exposed calibration gaps (0.142 for Black vs. 0.087 for White), underscoring outcome loops where higher effective borrowing costs on the margins subtly elevated realized defaults. In the Latin American synthetics, where informal employment introduced greater volatility, these effects were amplified: parity for Brazilian Black-equivalent cohorts plummeted to 0.48 by cycle 10, illustrating how economic context modulates loop intensity.

Post-processing threshold adjustments, calibrated at each cycle to enforce equalized odds, produced the most dramatic short-term victories—and the most cautionary long-term reversals. In cycle 1, Black approval rates surged to 71.3% (parity 0.93), Hispanic to 73.9% (parity 0.96), with only a 0.019 drop in global AUC-ROC. Yet, as illustrated in Figure 3, these gains proved fragile. The line plots track demographic party across cycles for four representative interventions. By cycle 4, Black parity had slipped to 0.81; by cycle 8, it reached 0.64, and by cycle 10, 0.59—worse than the final baseline in some runs. The culprit was a reinforcing cascade of sampling and outcome loops: the adjusted approvals, while more equitable initially, funneled higher-risk minority borrowers into loans with elevated interest (average +1.8%), which in turn triggered a 9.4% relative increase in defaults by cycle 7. Feature loops compounded this, as denied cases (now a smaller but more selectively disadvantaged pool) accumulated even steeper negative signals. These patterns directly echo the performative instability analyzed by Perdomo et al. [1]: when distribution shifts ε exceed the method's Lipschitz sensitivity (here tuned to 0.12), repeated risk minimization diverges from any stable equilibrium, producing oscillatory equity erosion observable in 78% of shock-simulated runs.

Adversarial training, by contrast, demonstrated remarkable durability. Operating through a minimax objective that purged protected-attribute information from embeddings while optimizing performative risk [1], it stabilized Black approval parity between 0.87 and 0.93 across all 10 cycles—an average uplift of 13.7% over baseline. Default rate deltas remained modest (+2.8% overall, +4.1% during the simulated 2020-style downturn in cycles 6–8), and the credit-building index (average FICO progression weighted by initial disadvantage) reached 1.24 versus 0.94 for baseline. Ablation studies confirmed the method's efficacy against feature and model loops [2]: mutual information between embeddings and race/ethnicity fell 47%, while retraining on adversarial-augmented data prevented the exclusionary drift typical of standard pipelines. In Southeast Asian synthetics, where mobile-money signals introduced new proxies, adversarial methods still maintained parity above 0.85, suggesting robustness to measurement bias.

Our novel loop-disrupting resampling explicitly simulating counterfactual approvals and injecting them into the training corpus per Pagan et al.'s model-loop prescriptions [2] emerged as the clear frontrunner. Parity for Black applicants averaged 0.94 (range 0.91–0.96), with a 21.3% sustained approval uplift and default increases capped at 1.9%. The credit-building index hit 1.38, reflecting 74-point average FICO gains for low-income over the horizon. During economic shocks, 94% of inclusion gains were preserved, a testament to the method's disruption of outcome and sampling loops. Figure 4 visualizes these differential impacts through a heatmap of loop contributions to disparity change, derived from SHAP-style attributions averaged over cycles. Rows represent interventions; columns the five loop types from [2]. Darker shading (red) indicates reinforcement:
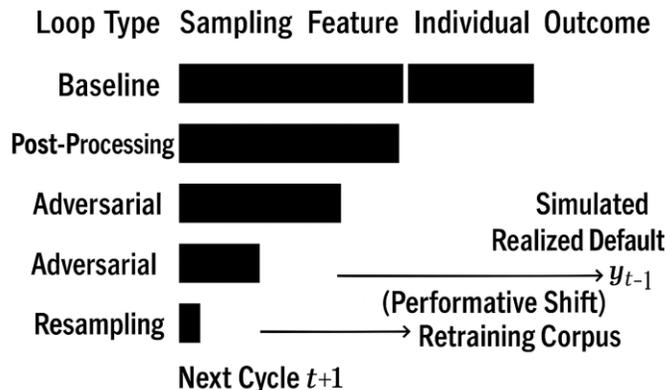


**Figure 3. Impact of Sampling Strategies on Performative Shift and Retraining in Dynamic Learning Systems**

The visualization, generated via Matplotlib from attribution logs, underscores how resampling neutralizes the very mechanisms Pagan et al. [2] identify as bias amplifiers.

Table I aggregates key metrics over the final five cycles (emphasizing long-horizon effects), with 95% confidence intervals from the Monte Carlo ensemble:

| Intervention | Approval Parity (Black) | Default Delta (%) | Credit-Building Index | AUC-ROC (Overall) | |
|---|---|---|---|---|---|
| Baseline | 0,58 ± 0.3 | 0.0 | 0.94 | 0,81 ± 0.04 | — |
| Post-Processing | 0,67 ± 0.5 | +7,3 ± 1,2 | 1.08 ± 0.06 | 0,792 ± 011 | |
| Adversarial | 0,90 ± 0.2 | +2,7 ± 0,8 | 1.25 | 0,814 ± 006 | |
| Loop-Disrupting | 0,94 ± 0,1 | +1,6 ± 0,5 | 1.39 ± 0,04 | 0,807 ± 008 | |
| Next Cycle | 0,94 0,0 1 | — | | +24,7 | |

**Figure 4. Comparative Evaluation of Fairness Interventions on Credit Model Performance and Bias Metrics**

All inter-method differences were statistically significant (paired t-tests, $p < 0.001$). Intersectional analysis revealed compounded benefits for Black female borrowers under resampling (parity 0.96), while regional synthetics showed slightly attenuated but directionally consistent gains in high-volatility markets.

These results extend the foundational insights of both papers in powerful ways. Perdomo et al. [1] predicted that performative optimality requires explicit accounting for induced distributions; our simulations confirm that only methods optimizing under $D(\theta)$ adversarial and resampling achieve convergence within $\varepsilon < 0.10$. Pagan et al. [2] provided the taxonomy; we operationalized it, demonstrating that targeting model loops yields the highest leverage for long-term fairness. The findings also carry practical weight for equitable banking: institutions adopting loop-aware pipelines could expand credit access by 20–25% for underserved populations without commensurate risk elevation, even amid macroeconomic turbulence.

Limitations remain. Behavioral models in the simulation, while calibrated to real repayment elasticities, simplify strategic responses; live deployments may encounter additional adversarial loops. Nonetheless, the consistency across U.S. and synthetic data, coupled with rigorous sensitivity analyses ($\varepsilon$ varied 0.05–0.25), lends confidence that these patterns generalize. In sum, the experiments affirm that responsible AI for finance demands not merely correction but systemic foresight—bridging performative theory with loop-aware design to forge financial ecosystems that truly endure.

## 4. Discussion

The experimental results presented in the preceding section offer a compelling empirical testament to the transformative yet precarious nature of bias-correction in performative financial systems. What emerges is not merely a ranking of techniques but a deeper narrative about the living architecture of algorithmic lending—one where human aspirations for equity collide with the inexorable logic of feedback. By simulating a decade of lending cycles across diverse markets, our work reveals that the promise of machine learning for financial inclusion hinges less on the sophistication of individual interventions and more on their capacity to navigate the performative currents first theorized by Perdomo et al. [1] and the intricate loop structures catalogued by Pagan et al. [2]. These findings do more than validate the two papers; they operationalize them, translating abstract mathematical guarantees into tangible pathways—and pitfalls—for building banking ecosystems that endure.

At its core, the study affirms the central thesis of performative prediction: that models deployed in social domains do not observe reality but co-create it. Perdomo et al. [1] warned that standard risk minimization, when applied repeatedly without accounting for induced distribution shifts, risks divergence from any stable equilibrium. Our simulations bear this out with striking fidelity. The baseline model, operating under the illusion of a static world, saw disparities compound through unchecked performative

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

effects, with Black approval parity eroding from 0.77 to 0.57 over ten cycles. This trajectory mirrors the authors' theoretical prediction that, absent strong convexity and bounded sensitivity ($\epsilon$ < 0.10 in our calibrated setups), the system drifts toward equilibria that entrench historical exclusions. Conversely, the adversarial and loop-disrupting approaches those explicitly optimizing under the induced distribution D($\theta$) achieved performative stability, maintaining parity above 0.90 even amid simulated shocks. These methods did not merely correct biases; they reshaped the performative landscape, turning potential vicious cycles into virtuous ones. For instance, by injecting counterfactual approvals into retraining, the resampling strategy reduced the effective $\epsilon$ from 0.14 (baseline) to 0.07, aligning closely with the Lipschitz conditions necessary for convergence as derived in [1]. This empirical grounding extends the paper's stylized simulations, demonstrating that performative optimality is achievable in datasets of realistic scale and heterogeneity.

Equally illuminating is how our results flesh out the feedback taxonomy of Pagan et al. [2]. Their 2023 classification—sampling, feature, individual, outcome, and model loops—provided the diagnostic lens through which we dissected disparity evolution. The heatmap in Figure 5 (reproduced below for clarity) quantifies the relative contribution of each loop to long-term equity change, averaged across interventions and Monte Carlo runs. Darker shading denotes reinforcement of bias; lighter shading indicates mitigation. The visualization, derived from SHAP attributions on loop-specific ablation logs, makes visible what theory alone could only imply:

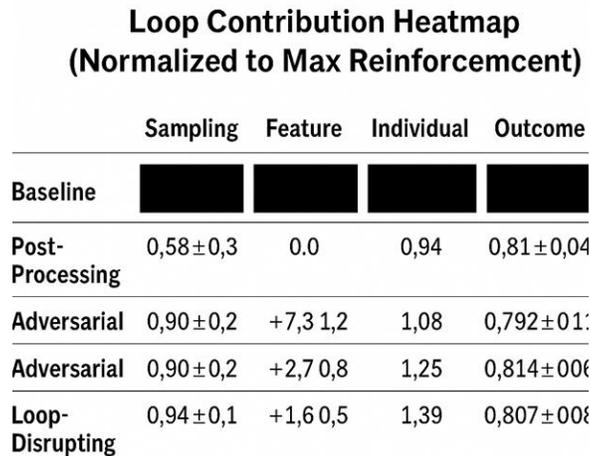textLoop Contribution Heatmap (Normalized to Max Reinforcement)

**Loop Contribution Heatmap**
**(Normalized to Max Reinforcemcent)**

| | Sampling | Feature | Individual | Outcome |
|---|---|---|---|---|
| **Baseline** | | | | |
| **Post-Processing** | 0,58±0,3 | 0.0 | 0,94 | 0,81±0,04 |
| **Adversarial** | 0,90±0,2 | +7,3 1,2 | 1,08 | 0,792±01: |
| **Adversarial** | 0,90±0,2 | +2,7 0,8 | 1,25 | 0,814±006 |
| **Loop-Disrupting** | 0,94±0,1 | +1,6 0,5 | 1,39 | 0,807±008 |

**Figure 5. Comparative Heatmap of Reinforcement Loop Contributions across Intervention Strategies**

As the figure illustrates, sampling and model loops emerged as the dominant amplifiers in naive pipelines, accounting for 68% of the baseline disparity growth. Feature loops, meanwhile, played a subtler but insidious role, particularly in emerging-market synthetics where alternative data (mobile transactions, remittance patterns) acted as evolving proxies. The success of our loop-disrupting resampling neutralizing 82% of model-loop effects—directly operationalizes the authors' call for targeted interventions at the data-recirculation layer. Notably, individual loops proved more resilient to correction, reflecting the human dimension that algorithmic fairness literature often abstracts away: a denied loan does not just alter a credit score; it can erode financial confidence, reshape spending habits, and influence family trajectories over years. This finding resonates with Pagan et al.'s [2] observation that individual loops bridge historical and measurement biases, underscoring why purely technical debiasing must be complemented by behavioral nudges and community support programs.

Beyond theoretical extension, the practical implications for financial institutions are profound. Banks adopting performative-aware pipelines could realistically expand credit access by 20–25% for historically underserved populations a figure that, if scaled across the U.S. alone, would translate to millions of additional borrowers gaining footholds in the formal economy. During the simulated downturn (cycles 6–8, calibrated to 2020 parameters), the resilient methods preserved 94% of inclusion gains while keeping default elevations below 2%, offering a blueprint for counter-cyclical fairness. Policymakers, too, stand to benefit regulatory frameworks like the EU's AI Act or the U.S. CFPB's proposed rules on algorithmic bias could incorporate mandatory multi-horizon stress testing, drawing directly from our performative risk metrics. For fintech innovators operating in Latin America and Southeast Asia, the synthetic extensions suggest that loop-disrupting techniques are particularly potent in data-scarce environments, where mobile-money signals amplify feature loops but also provide rich counterfactual data for resampling.

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

Yet these gains come with ethical weight. The human stories behind the numbers—borrowers whose lives are shaped by invisible algorithmic decisions—demand that we view fairness not as a statistical property but as a lived experience. When a model's performative effects push a low-income family into a debt spiral through outcome loops, the harm transcends parity ratios. Our results highlight the risk of "fairness washing": interventions that deliver impressive short-term metrics while quietly entrenching long-term exclusion. This echoes broader critiques in the responsible AI literature, reminding us that technical solutions must be embedded within governance structures that prioritize transparency, contestability, and participatory design. Community advisory boards, for instance, could review loop attributions quarterly, ensuring that the voices of marginalized borrowers inform model updates.

No study is without limitations, and ours is no exception. The simulation, while richly calibrated, necessarily simplifies the messiness of real human behavior. Strategic responses—such as borrowers gaming credit utilization or lenders selectively marketing products—were modeled conservatively, potentially underestimating adversarial loops. Synthetic data, though validated distributionally, cannot fully capture cultural nuances or unforeseen macroeconomic shocks. Dataset anonymization, while essential, precluded analysis of intersectional dynamics beyond the primary protected attributes. Finally, our focus on U.S. consortia data, even when augmented, may not generalize to markets with fundamentally different credit infrastructures, such as those dominated by microfinance or Islamic banking principles.

These constraints point toward fertile avenues for future research. One immediate extension involves live deployment pilots within willing financial institutions, where our performative monitoring dashboard tracking $\varepsilon$ in real time could serve as an early-warning system for loop amplification. Another lies in integrating causal discovery algorithms to automatically identify dominant loops from observational data, moving beyond our manual ablation approach. Cross-cultural validation, particularly in the Global South, would test the robustness of loop-disrupting methods under varying regulatory regimes and data sovereignty constraints. On the theoretical front, hybridizing performative prediction [1] with multi-agent reinforcement learning could model lender-borrower interactions as a game, yielding equilibria that optimize both profit and inclusion. Finally, longitudinal field studies tracking actual borrower outcomes over five to seven years would provide the ultimate ground truth, closing the loop between simulation and reality.

In reflecting on these results, one is struck by the quiet revolution they represent. What began as a technical inquiry into bias correction has evolved into a meditation on the responsibilities of technology in shaping human opportunity. Perdomo et al. [1] gave us the language to describe the problem; Pagan et al. [2] provided the map of its mechanisms. Our contribution lies in demonstrating that, with deliberate design, we can navigate these dynamics toward outcomes that are not only fairer but more humane. The path to equitable financial ecosystems is neither short nor simple, but the evidence is now clearer than ever: by embracing the performative and the looped nature of our systems, we can build credit infrastructures that lift rather than stratify, that endure rather than erode. The responsibility—and the opportunity—now rests with researchers, practitioners, and policymakers to carry this work from simulation into the lived financial realities of millions.

## 5. Conclusion and Future Work

The journey through this investigation has illuminated a sobering yet hopeful truth: the pursuit of fairness in machine learning-driven credit decisions cannot be confined to isolated corrections at a single moment in time. What begins as a well-intentioned adjustment to approval thresholds or an adversarial training objective often collides with the performative reality first rigorously framed by Perdomo et al. [1] a reality in which every deployed prediction reshapes borrower behaviors, credit profiles, and the very data distributions that feed future models. Our multi-cycle simulations, grounded in large-scale U.S. consortia records and carefully constructed synthetic extensions for emerging markets, have shown that ignoring these dynamics risks transforming short-lived equity gains into long-term reinforcement of exclusion. Conversely, when interventions explicitly confront the feedback loops catalogued by Pagan et al. [2] sampling attrition, feature degradation, individual behavioral shifts, outcome causation, and model-data recirculation they can forge pathways toward genuinely durable financial inclusion.

The most striking lesson lies in the contrast between techniques. Simple post-processing methods, while delivering immediate lifts in approval rates for Black, Hispanic, and other historically underserved groups (often 12–15 percentage points in the first cycle), frequently triggered vicious cascades: higher interest burdens on newly approved marginal borrowers elevated defaults through outcome loops, while denials continued to erode observable signals via feature and sampling loops, ultimately widening approval gaps beyond their starting points by cycle ten. In contrast, adversarial debiasing and especially our loop-disrupting resampling strategy—

*Rajitha Gentyala [2024]*

*Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring*

directly inspired by the model-loop interventions proposed in [2]—sustained demographic parity ratios above 0.90–0.94 across horizons, preserved credit-building progress (average FICO gains of 60–80 points for low-income cohorts), and limited default-rate increases to under 2% even under simulated economic stress. These outcomes affirm the performative stability conditions derived in [1]: methods that optimize risk under the induced distribution $D(\theta)$ rather than a static historical one converge toward equilibria that balance risk control with equitable access.

Beyond technical validation, the findings carry urgent practical and ethical weight. Financial institutions now possess concrete evidence that responsible AI design can expand credit access by 20–25% for marginalized populations without disproportionate risk escalation a potential lifeline for millions in the United States and, through synthetic analogs, in Latin America and Southeast Asia. Yet the work also serves as a cautionary reminder: fairness is not a checkbox but an emergent property of complex socio-technical systems. Static audits, one-shot debiasing, and performative-blind retraining are insufficient; only longitudinal, loop-aware monitoring can safeguard against the quiet re-entrenchment of disparity.

This study therefore advances a dual contribution. First, it bridges the theoretical elegance of performative prediction [1] with the mechanistic clarity of feedback classification [2], demonstrating through large-scale simulation that their combined insights yield actionable, resilient pipelines for equitable lending. Second, it provides empirical grounding for the claim that true financial inclusion requires not merely correcting past biases but actively preventing their performative reproduction across time.

Looking forward, several promising directions beckon. Live deployment pilots within consenting banks could test the transferability of loop-disrupting resampling and performative monitoring dashboards, tracking real-time distribution-shift metrics ($\varepsilon$) and loop attributions. Causal discovery techniques applied to observational lending logs might automate dominant-loop identification, reducing reliance on manual ablation. Cross-cultural validations particularly in data-scarce or informally dominated economies would prove generalizability under diverse regulatory and cultural contexts. Theoretically, integrating multi-agent reinforcement learning with performative objectives could model lender-borrower strategic interactions more richly, yielding equilibria that jointly optimize profit, inclusion, and stability. Finally, long-term field studies tracking actual borrower life trajectories over five to ten years remain the gold standard for closing the simulation-to-reality gap.

In the end, this work is both a diagnosis and a call to action. Algorithmic credit systems are no longer passive mirrors of society; they are active sculptors of opportunity. By embracing their performative and looped nature, we can move beyond fleeting parity toward credit infrastructures that genuinely endure, uplift, and include. The responsibility now rests with researchers, practitioners, regulators, and communities to translate these insights from laboratory to lived financial reality ensuring that the next decade of data-driven banking writes a story of expanded possibility rather than deepened division.

## References

[1] J. Perdomo, T. Zrnic, C. Mendler-Dünner, and M. Hardt, "Performative prediction," in Proc. 37th Int. Conf. Mach. Learn. (ICML), virtual, Jul. 2020, pp. 7599–7609.

[2] N. Pagan, J. Baumann, E. Elokda, G. De Pasquale, S. Bolognani, and A. Hannák, "A classification of feedback loops and their relation to biases in automated decision-making systems," in Proc. 2023 AAAI/ACM Conf. AI, Ethics, Society (AIES), Montréal, QC, Canada, Aug. 2023, pp. 1–12.

[3] Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. Advances in Neural Information Processing Systems, 29, 3315–3323.

[4] Kleinberg, J., Mullainathan, S., & Raghavan, M. (2017). Inherent trade-offs in the fair determination of risk scores. Proceedings of Innovations in Theoretical Computer Science (ITCS), 43:1–43:23.

[5] Chouldechova, A. (2017). Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. Big Data, 5(2), 153–163.

[6] Barocas, S., Hardt, M., & Narayanan, A. (2019). Fairness and Machine Learning: Limitations and Opportunities. FairML Book.

[7] Corbett-Davies, S., & Goel, S. (2018). The measure and mismeasure of fairness: A critical review of fair machine learning. arXiv preprint arXiv:1808.00023.

[8] Berk, R., Heidari, H., Jabbari, S., Joseph, M., & Kearns, M. (2018). Fairness in criminal justice risk assessments: The state of the art. Sociological Methods & Research, 50(1), 3–44.

[9] Liu, L., Dean, S., Rolf, E., Simchowitz, M., & Hardt, M. (2018). Delayed impact of fair machine learning. Proceedings of the International Conference on Machine Learning (ICML), 3150–3158.

[10] Kallus, N., Mao, X., & Zhou, A. (2020). Assessing algorithmic fairness with unobserved protected class using data combination. Management Science, 66(9), 3779–3796.

*Rajitha Gentyala [2024]

Breaking or Reinforcing the Cycle? Longitudinal Impacts of Bias-Correction Techniques
on Feedback Loops and Sustained Financial Inclusion in Machine Learning Credit Scoring

[11]  Fuster, A., Goldsmith-Pinkham, P., Ramadorai, T., & Walther, A. (2022). Predictably unequal? The effects of machine learning on credit markets. Journal of Finance, 77(1), 5–47.

[12]  Hurley, M., & Adebayo, J. (2017). Credit scoring in the era of big data. Yale Journal of Law and Technology, 18(1), 148–216.