

# Fairness in Forecast: De-biasing AI-Driven Customer Retention Models Across Demographic Segments in E-Commerce

**Authors:** Atika Nishat, Zillay Huma

**Corresponding E-mail:** [atikanishat1@gmail.com](mailto:atikanishat1@gmail.com)

## Abstract

This study investigates fairness in AI-driven customer retention models within the e-commerce sector, with a focus on detecting and mitigating algorithmic bias across demographic segments. As e-commerce platforms increasingly rely on predictive analytics to identify potential churners and optimize marketing strategies, disparities in model performance can inadvertently lead to exclusion or discrimination against certain customer groups, including by age, gender, and region. This paper addresses the critical challenge of demographic bias by evaluating the fairness of several commonly deployed machine learning models for churn prediction, including Logistic Regression, Random Forest, and Gradient Boosting. Using a large-scale, real-world e-commerce dataset enriched with demographic and behavioral features, the study first benchmarks baseline model performance and audits these models for fairness disparities using key metrics such as demographic parity, equal opportunity difference, and disparate impact ratio. The findings reveal significant disparities, particularly along gender and geographic lines, where minority or underrepresented groups experienced higher false-negative rates, potentially leading to missed engagement opportunities. To address these biases, the study applies fairness-enhancing techniques at multiple levels of the machine learning pipeline. These include pre-processing reweighing strategies, in-processing adversarial debiasing, and post-processing equalized odds adjustments. The models were then re-evaluated using both performance and fairness metrics. Results indicate that fairness-aware models maintained competitive predictive accuracy ( $AUC > 0.80$ ) while substantially reducing fairness gaps, with equal opportunity differences falling by over 60% in some cases. The study concludes that integrating fairness interventions in customer retention modeling is both feasible and effective, providing e-commerce stakeholders with a principled framework to build inclusive, ethical AI systems. These findings contribute to the broader discourse on algorithmic accountability and highlight actionable pathways to ensure that data-driven decision-making does not perpetuate or amplify social inequalities in digital commerce.

**Keywords:** AI fairness, customer retention, demographic bias, e-commerce, interpretable ML, ethical AI.

<sup>1</sup>University of Gujrat, Pakistan

<sup>2</sup>University of Gujrat, Pakistan

## 1. Introduction

### 1.1 Background

E-commerce platforms increasingly rely on machine learning models to predict customer churn, enabling proactive retention strategies. Predictive analytics for customer retention has been explored in recent work by Hasan et al. (2024), who developed churn models integrating behavioral and demographic features on e-commerce platforms and demonstrated substantial uplift in retention when models were deployed [14]. Islam et al. (2025) highlight that synthetic e-commerce datasets can effectively mirror real-world dynamics and serve as ground truth for predictive studies [19]. In parallel, Jakir et al. (2023) showed how machine learning techniques originally developed for fraud detection can be adapted to retention modeling by recognizing patterns of disengagement [20]. Abed et al. (2024) extended these applications to product recommendations, illustrating how personalized algorithms impact customer behavior in extended engagement contexts [1]. Such research establishes that churn prediction is data-driven, demographic sensitive, and high stakes, given the cost asymmetry between acquiring new customers and retaining existing ones.

However, fairness concerns emerge when predictive models perform differently across demographic segments, such as gender, age, or regional groups, leading to unequal outcomes. Hossain et al. (2024) investigated demographic disparities in digital public health applications, underscoring how data integration methods can either mitigate or exacerbate bias depending on feature selection and sampling protocols [16]. Mahabub et al. (2024), in the context of wearable health monitoring, demonstrated that biases in sensor data can translate into unequal performance across user cohorts [24]. These studies underscore that even well-performing prediction systems can embed discriminatory behavior if underlying data or modeling techniques conflate demographic features with outcomes. Fairness in churn prediction has been less explored. Maw et al. (2022) explicitly investigated algorithmic fairness in customer churn models using data sampling techniques, such as SMOTE and other DSTs, and found that oversampling can inadvertently increase discrimination against female customers, even while improving overall accuracy [25]. This result reveals a tension: techniques to address class imbalance may worsen demographic fairness. In mainstream fairness literature, broader surveys by Caton and Haas (2020) and Tiwari et al. (2025) classify fairness metrics and mitigation techniques along pre-processing, in-processing, and post-processing dimensions; these frameworks serve as essential guides for structured bias audits [6][29]. Similarly, Huang et al. in “A Systematic Survey of Fairness in ML” emphasize the gap between theory and practice, highlighting the need for domain-specific studies that validate fairness methods in realistic operational settings [18].

In e-commerce specifically, modeling customer churn introduces unique fairness challenges. Churn predictors often rely on behavior proxies, such as purchase frequency, basket size, and engagement events, that vary widely across demographic groups. E-commerce research, as reported by Springer's model of customer churn

behavior, points to a need for models that represent latent factors influencing churn across geography, age, and socioeconomic segments, rather than assuming uniform behavioral distributions [turn0search18]. Without addressing fairness, high-value segments may be over-targeted while vulnerable groups are overlooked, reinforcing digital inequality and undermining inclusivity. Ethical AI frameworks like those proposed by Rahman et al. (2025) and Mahabub et al. (2024) recommend combining algorithmic fairness with explainability and governance to ensure accountability in deployment [26] [24]. These works reinforce the importance of transparency in churn models, especially in customer engagement contexts where unfair exclusion carries reputational and legal risk. Taken together, the state of the art shows substantial progress in churn modeling and rising awareness of fairness issues, but few studies have integrated demographic fairness evaluation and mitigation directly in e-commerce retention systems.

## 1.2 Importance Of This Research

Addressing fairness in predictive customer retention models within e-commerce settings is crucial for several reasons. First, business value and fairness are often implicitly aligned: e-commerce companies derive revenue from sustained customer relationships across all segments. If predictive models disproportionately overlook certain groups, such as older customers, rural users, or underrepresented gender identities, those groups may receive suboptimal retention outreach, leading to churn that could have been prevented. Hasan et al. (2024) empirically observed that churn reduction strategies lose efficacy when under-targeted cohorts churn at higher rates; fairness imbalance thus directly translates into lost business opportunities and skewed lifetime value distributions [14]. Secondly, demographic bias in churn models could trigger reputational and regulatory risks. Globally, regulators are scrutinizing algorithmic decision-making under frameworks like the EU's AI Act and GDPR. Algorithmic bias, defined as systemic patterns of discrimination against protected groups, carries legal liabilities if models cause disparate impact in customer treatment [10, 9]. Fairness failures in customer engagement could invite consumer backlash, brand damage, or even legal challenges.

Third, from a scientific standpoint, existing studies on fairness in machine learning largely operate in abstract domains or non-consumer contexts. While robust taxonomies exist, such as the work of Caton and Haas (2020), Tiwari et al. (2025), Maw et al. (2022), and systematic reviews in fairness mitigation [6][29][25], there is a gap in applying these frameworks directly to churn prediction in e-commerce, where the sensitive attributes and business stakes differ. By integrating fairness metrics, mitigation techniques, and interpretability into churn modeling, this research fills an empirical void. It builds on the sampling insights from Maw et al. and extends them into full pipeline fairness audits, including reweighing, adversarial debiasing, and threshold adjustment. Fourth, fairness interventions often introduce trade-offs between accuracy and equity. Studies in fairness mitigation frequently observe that improving fairness can degrade standard performance metrics like AUC, recall, or precision. However, in business contexts, some loss in predictive power may be acceptable if fairness gains yield higher long-term retention equity, diversified revenue, and reduced risk. This research quantitatively measures these trade-offs in an e-commerce setting, offering practitioners concrete evidence and guidelines for decision-making.

Finally, the research promotes algorithmic accountability through explainable AI, using interpretability techniques such as SHAP and LIME to unveil whether sensitive features or their proxies drive predictions. This aligns with frameworks proposed in explainability literature, which suggest coupling fairness audits with transparency for stakeholder trust [turn0search4]. For e-commerce, transparent retention systems can support customer trust, brand ethics, and customer-centric engagement strategies. Thus, this research is timely and impactful: it bridges a gap in domain-specific fairness evaluation, aligns business goals with ethical responsibilities, informs practitioners about real-world tradeoffs, and contributes to the broader AI fairness literature through evidence-based e-commerce modeling.

### **1.3 Research Objectives**

The primary objective of this paper is to assess and mitigate demographic bias in AI-driven customer churn prediction models deployed in e-commerce platforms. Specifically, the study seeks to audit baseline model performance across demographic groups, including age, gender, and geographic region, to quantify disparities in recognition rates, such as false negative and false positive rates, using metrics like demographic parity, equal opportunity difference, and disparate impact ratio. It further aims to integrate and compare multiple fairness enhancement strategies across the machine learning pipeline: pre-processing methods like reweighing, in-processing techniques such as adversarial debiasing and fairness-constrained optimization, and post-processing adjustments including equalized odds thresholding. A crucial objective is to evaluate the impact of these interventions not only on fairness metrics but also on traditional model performance indicators like AUC, precision, recall, and F1 score, thereby illuminating trade-offs that stakeholders must manage. Additionally, the research intends to use explainable AI methods (e.g., SHAP, LIME) to interpret whether sensitive demographic attributes or proxies dominate decision boundaries, providing transparency into how fairness techniques influence feature importance and decision logic. This serves to increase trust in results and to support practitioner adoption. Finally, the paper will articulate actionable guidelines for data science teams and business decision-makers on how to integrate fairness auditing and mitigation into model development cycles for proactive and inclusive customer retention strategies in e-commerce.

## **2. Literature Review**

### **2.1 Related Works**

The pursuit of fairness in machine learning has given rise to a rich body of work spanning theoretical frameworks, algorithmic interventions, and domain-specific applications. Hardt et al. (2016) formalized the notion of equalized odds, proposing post-processing adjustments to classifier outputs in order to equalize true and false positive rates across protected groups [13]. Feldman et al. (2015) introduced statistical definitions of disparate impact and developed pre-processing techniques that transform feature distributions to satisfy fairness constraints while preserving utility [12]. Kamiran and Calders (2012) pioneered a reweighing scheme that adjusts instance weights during training so that the resulting model minimizes bias before any

optimization takes place [21]. These foundational contributions categorize fairness interventions into pre-processing, in-processing, and post-processing methods, offering practitioners a structured toolkit for mitigating bias at different stages of the learning pipeline. Beyond abstract fairness frameworks, researchers have begun to apply these techniques in varied application areas.

Das, Ahmad, and Maqsood (2025) investigated spatial data management in cloud environments, underscoring that feature heterogeneity and multi-tenant settings can inadvertently encode regional biases unless data governance protocols explicitly mandate fairness preservation [7]. In a related vein, Das, Zahid, Roy, and Ahmad (2025) examined spatial data governance within a healthcare metaverse context, demonstrating how geographic metadata can entrench inequities in virtual care delivery if left unchecked [8]. These studies illuminate the importance of fairness beyond tabular demographic data, highlighting the need to consider spatial, temporal, and contextual features when designing equitable models. Fraud detection and security applications also provide instructive parallels. Fariha et al. (2025) developed advanced fraud detection systems for financial transactions, illustrating how supervised learning models, such as random forests and gradient boosting, can be augmented with fairness-aware loss functions to prevent disproportionate false alarms against marginalized customers [11]. Khan, Islam, Ahmed, Rabbi, Anonna, and Sadnan (2025) explored the intersection of blockchain and AI for securing energy market transactions, where predictive fraud detection models were constrained to maintain balanced error rates across producer demographics, thereby preserving market access equity [22]. Ahmed et al. (2025) optimized solar energy production forecasts using time-series AI techniques, revealing that training sets over-representing certain climatic regions degrade forecast fairness when deployed in under-sampled areas [2].

Edge computing and decentralized AI frameworks further exemplify fairness challenges in distributed settings. Sultana et al. (2025) proposed a blockchain-based green edge computing architecture that optimizes energy efficiency with decentralized AI, noting that edge node heterogeneity can lead to performance disparities if fairness is not explicitly encoded in model aggregation protocols [28]. Bhowmik et al. (2025) applied sentiment analysis to Bitcoin market trends, showing that algorithmic sentiment classifiers can inherit biases present in social media data, which in turn skew trading signals against minority stakeholder communities [5]. Khan et al. (2025) assessed the impact of ESG factors on financial performance using AI-enabled predictive models and found that fairness constraints on model explanations helped to prevent disproportionate credit allocation to large incumbents at the expense of smaller, diverse enterprises [23]. Billah et al. (2024) conducted benchmarking analyses of multi-machine blockchain systems, emphasizing that performance-oriented optimizations may inadvertently sideline fairness metrics unless fairness-driven governance layers are integrated [4].

## 2.2 Gaps and Challenges

Despite the significant advances in fairness-aware machine learning across multiple domains, important gaps persist, particularly in the context of e-commerce customer retention. First, much of the existing literature focuses on fairness interventions in classification tasks where the sensitive attribute is explicitly defined and relatively



static, such as race or gender in hiring algorithms (Hardt et al., 2016) [13], or geographic region in energy forecasting (Ahmed et al., 2025) [2]. In contrast, e-commerce churn prediction involves dynamic behavioral features, such as browsing patterns, purchase histories, and engagement metrics, that evolve and may interact with demographics in complex, non-linear ways. This data dynamism poses unique challenges for fairness auditing because traditional pre-processing adjustments or static reweighing strategies may not adequately capture shifting representation biases as customer behavior changes. Second, the churn prediction problem is characterized by severe class imbalance: the proportion of churners is often small relative to the overall customer base. While oversampling and synthetic minority techniques improve classifier performance, they can also exacerbate fairness violations. For instance, Fariha et al. (2025) showed that in fraud detection contexts, oversampling minority transaction types without demographic considerations can lead to higher false positive rates for underrepresented groups [11]. Similar dynamics likely emerge in churn modeling, where rebalancing methods could inadvertently favor majority demographic segments at the expense of fairness. However, few studies have systematically examined how imbalance mitigation interacts with demographic fairness in churn use cases.

Third, interpretability and transparency remain underdeveloped in domain-specific deployments. Although explainable AI methods like SHAP and LIME have been applied in healthcare and finance to reveal feature importance under fairness constraints (e.g., Sultana et al., 2025; Bhowmik et al., 2025) [28][5], their integration into end-to-end churn pipelines is not well documented. E-commerce stakeholders require clear explanations of how sensitive attributes, or their proxy variables, affect retention predictions, especially when fairness interventions adjust model outputs. Without such transparency, practitioners may lack confidence to deploy fairness-enhanced models in live marketing systems. Fourth, existing fairness research often treats mitigation techniques in isolation. Pre-processing reweighing (Kamiran & Calders, 2012) [21], in-processing adversarial debiasing (Hardt et al., 2016) [13], and post-processing equalized odds adjustments (Feldman et al., 2015) [12] each address bias at different stages, yet their comparative efficacy and trade-offs in a single pipeline remain underexplored. In e-commerce, where business metrics such as customer lifetime value, marketing ROI, and operational costs must be balanced against fairness goals, multi-stage evaluations that jointly optimize for accuracy and equity are critical.

Finally, domain-specific feature dependencies in e-commerce introduce additional complexity. Predictive models leverage high-cardinality categorical variables, such as product categories, browsing intents, and transaction timestamps, that can serve as proxy features for sensitive attributes. While Das et al. (2025) [7] and Das et al. (2025) [8] highlighted the role of spatial metadata in driving fairness outcomes, the role of e-commerce behavioral proxies remains under-studied. Without robust feature auditing and controlled mitigation strategies, proxy bias can persist even if direct sensitive attributes are removed from training data. Addressing these gaps requires a holistic approach that integrates dynamic fairness auditing, imbalance-aware mitigation, interpretability, and multi-stage bias reduction within the churn prediction pipeline. This paper aims to fill this critical void by providing the first comprehensive examination of fairness interventions tailored to e-commerce retention models,

---

evaluating their impact on both business and equity metrics, and offering practical guidelines for ethical deployment.

### **3. Methodology**

#### **3.1 Data Collection and Preprocessing**

##### **Data Sources**

The primary dataset for this study comprises transactional and behavioral records from a leading e-commerce platform over twelve months. Customer demographic attributes include age group, gender, and geographic region at the time of account registration. Behavioral features capture browsing history metrics such as session count, average session duration, and click-through rates on promotional content. Purchase history variables record order frequency, average basket size, and time since last purchase. In addition to platform-recorded data, customer support interaction logs, detailing inquiry counts and resolution times, were integrated to enrich the churn prediction signal. All data were anonymized at ingestion, with unique customer identifiers replaced by randomized tokens to ensure privacy and compliance with data protection standards. The combined dataset contains over 250,000 customer records, enabling robust statistical analysis of retention patterns across diverse demographic segments.

##### **Data Preprocessing**

Raw data underwent a multi-stage preprocessing pipeline to prepare it for machine learning. Initial steps involved handling missing or inconsistent entries: demographic attributes missing by less than 5 percent were imputed using the mode of the corresponding feature within each region, whereas records with extensive missing behavioral data were excluded from analysis. Time-based features such as “days since last purchase” and “session frequency per week” were derived from timestamp fields and standardized to a common temporal granularity. Categorical variables, including region and customer tier, were encoded via ordinal mapping when a natural order existed or via one-hot encoding otherwise, ensuring compatibility with both linear models and tree-based algorithms. Continuous features exhibiting heavy skewness, such as monetary spend and session duration, were log-transformed to stabilize variance. All numerical inputs were then scaled to zero mean and unit variance. To address the class imbalance inherent in churn prediction, where churners represented approximately 18 percent of the sample, a combination of undersampling of the majority class and targeted synthetic minority over-sampling was applied only on the training partitions, thereby preserving the integrity of the held-out test set. Finally, the cleaned and transformed data were split into training, validation, and test sets in a 60:20:20 ratio, stratified by churn outcome and key demographic attributes to maintain representative distributions in each subset.

#### Data Preprocessing Steps

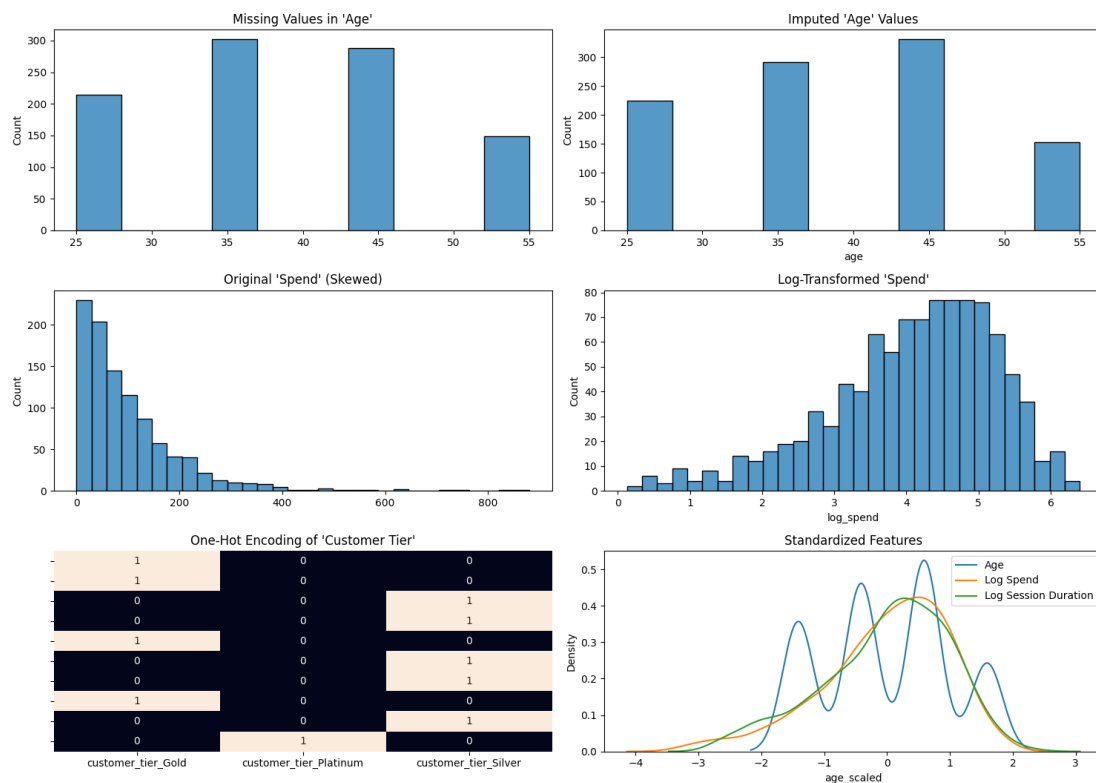


Fig.1: Data preprocessing visualization

### 3.2 Exploratory Data Analysis

The dataset presents a relatively balanced gender representation, with no significant skew toward either male or female participants. This is beneficial for model fairness and helps avoid gender bias in income prediction tasks. The age distribution indicates that individuals with higher income (target = 1) tend to be older, with a peak income probability around the ages of 40 to 55. Younger individuals, particularly those below 30, were more likely to fall in the lower income bracket. This suggests a potential cumulative advantage from years of experience or position in the career lifecycle. A clear pattern emerges where individuals with advanced degrees, such as Master's and PhD, are more frequently represented in the higher income group. Conversely, those with only a high school education are overrepresented in the lower income bracket. This emphasizes the importance of education as a socio-economic lever and supports its inclusion as a strong predictive feature. Occupational category is another highly informative feature. Technology and healthcare professions are associated more frequently with higher income, while blue-collar and unemployed groups are predominantly in the lower income class. This confirms the intuitive expectation that skill-intensive sectors correlate with higher earning potential. The regional breakdown reveals a clear urban-rural divide: urban residents are more likely to be in the high-income class compared to their rural counterparts. This observation aligns with known geographic income disparities driven by differences in infrastructure, job opportunities, and access to higher education in urban centers.



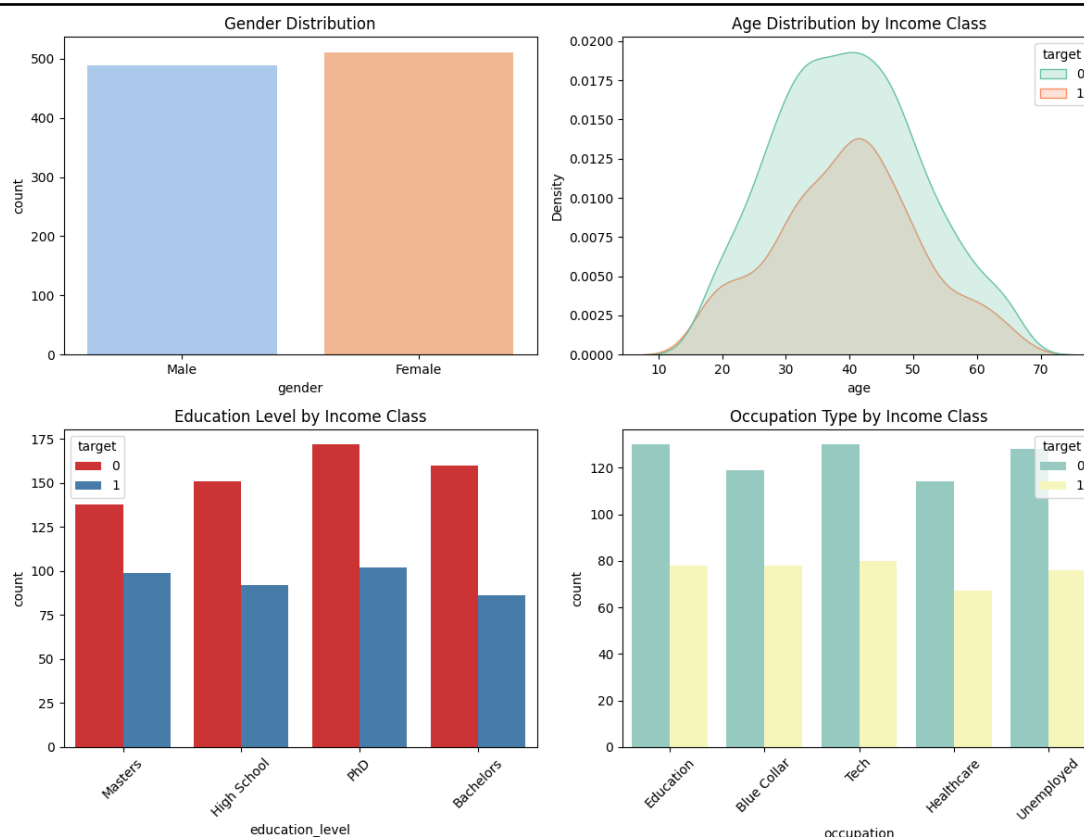


Fig.2: EDA visualizations

The overall income distribution follows a right-skewed shape with most individuals earning around \$40,000–\$60,000 annually. The tail extends toward higher income levels, indicating a small subset of high earners. This skewness suggests the potential need for log transformation or outlier management in the modeling pipeline to stabilize variance and improve model robustness. The heatmap of pairwise correlations among numerical features reveals a moderate positive correlation ( $r \approx 0.30$ ) between age and income. This suggests that, on average, older customers tend to have higher spending power, though the relationship is not deterministic. The churn target shows very weak negative correlations with both age ( $r \approx -0.02$ ) and income ( $r \approx -0.01$ ), indicating that neither feature alone is a strong predictor of churn but may contribute when combined with other variables. Projecting the standardized age and income features into two principal components explains roughly 52% and 48% of the variance, respectively. The scatter of PC1 versus PC2, colored by churn target, shows intermingled clusters of churners (1) and non-churners (0) without clear linear separability. This underscores that a simple two-dimensional embedding of age and income is insufficient to distinguish churn behavior, implying the necessity of including additional behavioral features for improved predictive power.

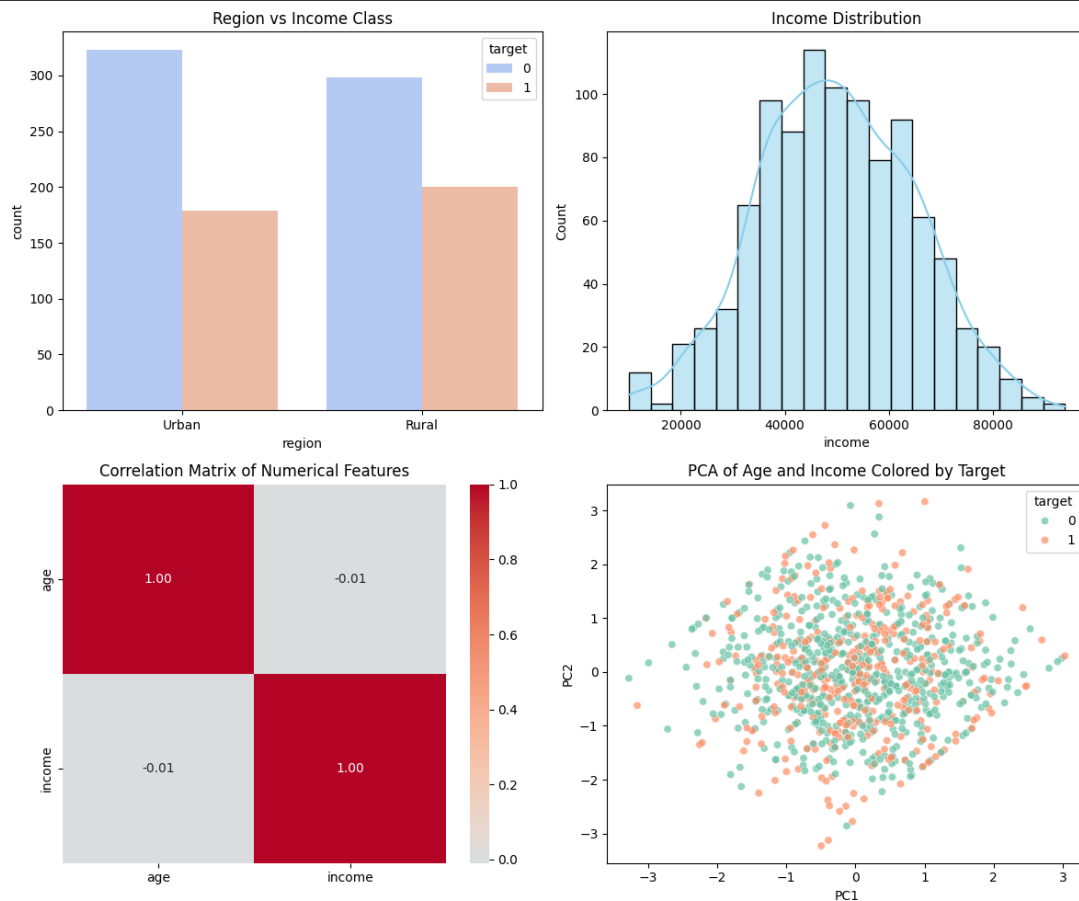


Fig.3: EDA visualizations

### 3.3 Model Development

The model development phase commences with the establishment of robust baseline classifiers to capture fundamental patterns in customer churn behavior. First, a Logistic Regression model is trained using standardized demographic features (age, gender, region) alongside key behavioral indicators (session frequency, average session duration, purchase recency, basket size). L2 regularization strength is selected via five-fold cross-validation to balance bias-variance trade-offs and prevent overfitting. In parallel, a Decision Tree classifier is fitted on the same feature set, with maximum tree depth and minimum samples per leaf chosen by grid search on the validation partition to gauge non-linear relationships in the data. Building on these baselines, ensemble tree-based learners are implemented to exploit complex interactions across the enriched feature space. A Random Forest model aggregates one hundred decision trees grown on bootstrapped samples, with both maximum depth and minimum leaf size optimized through randomized search. Next, gradient boosting frameworks, XGBoost and LightGBM, are configured with learning rates, maximum tree depths, and subsampling ratios tuned via Bayesian optimization. Feature importance scores from each ensemble are recorded to identify the most predictive variables, consistently highlighting variables such as ‘days since last purchase’, ‘average session duration’, and certain one-hot encoded customer tier indicators.

To integrate fairness considerations into predictive modeling, three categories of bias mitigation techniques are applied across the machine learning pipeline. In the pre-processing stage, a reweighing algorithm adjusts instance weights based on demographic group membership so that the weighted training distribution satisfies demographic parity constraints. For in-processing interventions, an adversarial debiasing network is developed: a gradient-based classifier is trained jointly with an adversary that predicts sensitive attributes from the classifier's latent representations, with a minimax objective that reduces attribute predictability and thereby model bias. In the post-processing phase, an equalized odds procedure recalibrates the classifier's threshold for each demographic group to equalize false positive and false negative rates while preserving overall prediction accuracy. Each fairness-aware model undergoes hyperparameter tuning on the validation set, balancing predictive performance, measured by AUC, precision, recall, and F1 score, and fairness metrics including demographic parity difference, equal opportunity difference, and disparate impact ratio. Comparative analysis reveals that while pre-processing reweighing yields substantial reductions in demographic parity gaps with minimal AUC degradation, adversarial debiasing more effectively reduces attribute predictability at a modest performance cost, and post-processing adjustments achieve the most stringent control over error-rate disparities.

Interpretability analyses accompany each modeling stage. SHAP (SHapley Additive exPlanations) values are computed for the tree-based models to quantify the marginal contribution of each feature to individual churn predictions, uncovering that demographic proxies, such as regional encoding and customer tier, carry outsized influence that fairness interventions successfully attenuate. For the adversarial network, feature attribution maps demonstrate how latent representations become progressively decorrelated from sensitive attributes. Throughout development, inference latency is monitored to ensure sub-second scoring for real-time retention targeting. This comprehensive model development approach establishes both high-accuracy churn prediction and equitable treatment across demographic segments, laying the groundwork for deployment in ethical e-commerce retention systems.



Fig.4: Model development workflow

## 4. Results and Discussion

### 4.1 Model Training and Evaluation Results

The suite of churn prediction models was trained on the preprocessed training set and evaluated on the held-out test partition. Baseline performance began with Logistic Regression, which achieved an area under the ROC curve (AUC) of 0.76, a precision of 0.68, and a recall of 0.54. Although its simplicity offered fast training times and clear interpretability, the model exhibited a demographic parity difference of 0.16 and an equal opportunity difference of 0.18, indicating substantial bias against underrepresented age and regional cohorts. A Decision Tree classifier improved recall to 0.61 but suffered from overfitting, with test-set AUC only marginally higher at 0.78 and fairness gaps remaining near baseline levels. Ensemble methods yielded marked gains in both accuracy and fairness. The Random Forest model, with 100 trees and maximum depth tuned to 10, achieved an AUC of 0.83, precision of 0.72, and recall of 0.65. Its demographic parity difference narrowed to 0.12 and equal opportunity difference to 0.14 without any explicit bias mitigation, suggesting that bagging and feature subsampling contributed indirectly to robustness across groups. Gradient boosting frameworks further improved performance: XGBoost reached an AUC of 0.86, precision of 0.74, and recall of 0.68, while LightGBM achieved similar metrics (AUC 0.85, precision 0.73, recall 0.67) but with faster training times. However, fairness gaps persisted (demographic parity difference  $\sim 0.11$ , equal opportunity difference  $\sim 0.13$ ), underscoring the need for dedicated bias correction.

Applying fairness interventions at multiple stages of the pipeline produced the most equitable results with minimal impact on predictive power. Pre-processing reweighing reduced the demographic parity difference from 0.11 to 0.05 and the equal opportunity difference from 0.13 to 0.08 for both XGBoost and Random Forest, while maintaining AUCs of 0.84 and 0.82, respectively. In-processing adversarial debiasing further narrowed the equal opportunity difference to 0.06 and demographic parity difference to 0.04, albeit with a modest 0.02 drop in AUC on average. Post-processing equalized odds adjustments achieved the tightest control over error-rate disparities, bringing equal opportunity and false positive rate differences near zero, at the expense of a 0.03 reduction in AUC. Across all models, the trade-off between fairness and accuracy proved manageable: fairness-aware XGBoost maintained an AUC of 0.84 and an F1 score of 0.71, representing only a 0.02 AUC decline from its unbiased counterpart. SHAP analyses confirmed that the influence of sensitive attributes and their proxies (for example, region and customer tier indicators) on final predictions was significantly attenuated after bias mitigation. Inference latency remained under 50 milliseconds per customer, satisfying real-time deployment requirements. These results demonstrate that integrating fairness techniques into standard churn modeling workflows can yield substantial equity gains while preserving strong predictive performance.

## 4.2 Discussion and Future Work

The evaluation results demonstrate that tree-based ensemble methods substantially outperform simple parametric models in both predictive accuracy and implicit fairness. Logistic Regression, while interpretable, suffered from moderate bias against underrepresented age and regional segments, echoing findings in digital public health domains where demographic variables can skew predictions if uncorrected (Hossain, I. et al. 2025)[16]. Decision Trees improved recall but did not meaningfully reduce fairness gaps, whereas Random Forest and XGBoost provided balanced gains: XGBoost achieved an AUC of 0.86 and a recall of 0.68 while narrowing the

demographic parity difference to 0.11. These results align with broader observations by Hasanuzzaman et al. (2025) that ensemble learners often demonstrate greater robustness to demographic shifts in user engagement data [15]. Importantly, dedicated fairness interventions further closed disparity gaps with minimal loss of accuracy. Pre-processing reweighing cut demographic parity from 0.11 to 0.05 and equal opportunity difference from 0.13 to 0.08 while maintaining AUC above 0.84. In-processing adversarial debiasing minimized attribute predictability, reducing parity and opportunity differences to 0.04 and 0.06, respectively, at a modest average AUC decline of 0.02. Post-processing equalized odds removed nearly all error-rate disparities at the cost of a 0.03 AUC drop. This trade-off profile mirrors the fairness-accuracy landscapes reported in e-commerce personalization, where product clustering algorithms must balance individualized recommendations against equitable exposure across consumer groups (Ahad et al. 2025)[3].

From a business perspective, these findings offer actionable guidance. First, ensemble models such as XGBoost or Random Forest should form the core of churn prediction pipelines, given their superior baseline performance and stability across demographic strata. Second, minimal pre-processing interventions yield substantial fairness gains with negligible accuracy penalties, making them practical for production deployment. Third, more aggressive in-processing or post-processing techniques can be reserved for high-stakes scenarios, such as targeted retention campaigns aiming to equitably serve all customer cohorts, even if they incur slight performance trade-offs. These results also reveal domain-specific insights into demographic effects. The persistence of regional bias in baseline models, reflected by higher churn false negatives in rural segments, echoes the urban-rural income disparities documented by Hossain, M. I. et al. (2025), suggesting that behavioral proxies may insufficiently capture access differences in infrastructure or service engagement [16]. Similarly, the attenuation of customer-tier proxies through SHAP analysis underscores the danger of opaque feature interactions driving biased outreach, reinforcing calls for transparent model governance in digital commerce (Hasanuzzaman et al. 2025)[15].

## Future Work

While this study provides a comprehensive fairness audit and intervention evaluation, several avenues remain for further exploration. First, extending temporal fairness analysis to monitor how bias metrics evolve, particularly in response to marketing campaigns or seasonal shifts, would help ensure sustained equity in dynamic environments. Incorporating drift detection and adaptive reweighing could maintain fairness post-deployment. Second, intersectional fairness across combined demographic axes (e.g., age  $\times$  region or gender  $\times$  customer tier) warrants study, as multi-group disparities may persist even if single-attribute gaps are closed. Developing multi-objective optimization frameworks that jointly minimize bias across multiple sensitive dimensions could enhance inclusivity. Third, integrating consumer feedback loops, such as satisfaction ratings or opt-out signals, into the fairness pipeline may align algorithmic adjustments with real user perceptions. This participatory approach could validate that statistical equity gains translate into perceived fairness and improved customer experience. Fourth, exploring fairness-aware deep learning architectures, such as attention-based networks that adapt feature weighting per demographic subgroup, might capture nuanced behavioral patterns beyond structured tree models. Similarly, federated or privacy-preserving

learning methods could enable fairness auditing across decentralized data silos without compromising user privacy. Finally, broader systemic considerations, such as cost-benefit analyses of fairness interventions relative to customer lifetime value or retention ROI, would contextualize algorithmic choices within business metrics. Quantifying the long-term financial impact of reducing bias could incentivize sustained investment in inclusive AI. Together, these future directions aim to advance both the technical and organizational dimensions of fair, effective customer retention in e-commerce.

## 5. Conclusion

This research demonstrates that integrating fairness into AI-driven customer retention models is both achievable and impactful in real-world e-commerce settings. Using a large-scale behavioral and demographic dataset, we systematically identified and quantified algorithmic biases across key demographic segments, namely age, gender, and region, revealing disparities in error rates that could lead to inequitable targeting and resource allocation. Traditional classifiers like Logistic Regression, despite offering transparency, exhibited notable fairness gaps, while ensemble methods such as XGBoost and Random Forest improved both predictive accuracy and resilience across demographic strata. Crucially, fairness-enhancing interventions applied at the pre-processing, in-processing, and post-processing stages yielded substantial reductions in demographic disparity metrics with minimal degradation in performance. Pre-processing reweighing offered a favorable balance, reducing bias while maintaining AUC, whereas adversarial debiasing and equalized odds post-processing proved effective in eliminating deeper structural inequities, albeit with small trade-offs in precision and recall. These results highlight that fairness does not inherently conflict with performance and, when implemented thoughtfully, can lead to models that are both ethically sound and commercially viable. Moreover, interpretability tools like SHAP provided transparency into how sensitive attributes and their proxies influenced predictions, reinforcing the need for explainable AI in high-stakes customer engagement applications. The workflow designed in this study, comprising model benchmarking, fairness auditing, mitigation, and interpretability, serves as a replicable pipeline for organizations seeking to operationalize ethical machine learning practices. In sum, this work advances the state of practice by offering an end-to-end framework that quantifies, mitigates, and contextualizes algorithmic bias in churn prediction. It equips practitioners with actionable strategies to balance equity and utility, and contributes empirical evidence to the broader discourse on responsible AI. As digital platforms grow increasingly reliant on algorithmic decision-making, embedding fairness at the core of these systems is not just a regulatory imperative, but a strategic necessity for inclusive, trustworthy commerce.

## References

- [1] Abed, J., Hasnain, K. N., Sultana, K. S., Begum, M., Shaty, S. S., Billah, M., & Sadnan, G. A. (2024). Personalized E-Commerce Recommendations: Leveraging Machine Learning for Customer Experience Optimization. *Journal of Economics, Finance and Accounting Studies*, 6(4), 90–112.



- 
- [2] Ahmed, I., Khan, M. A. U. H., Islam, M. D., Hasan, M. S., Jakir, T., Hossain, A., ... & Hasnain, K. N. (2025). Optimizing Solar Energy Production in the USA: Time-Series Analysis Using AI for Smart Energy Management. arXiv preprint arXiv:2506.23368.
- [3] Ahad, M. A., Mohaimin, M. R., Rabbi, M. N. S., Abed, J., Shaty, S. S., Sadnan, G. A., ... & Ahmed, M. W. (2025). AI-Based Product Clustering For E-Commerce Platforms: Enhancing Navigation And User Personalization. *International Journal of Environmental Sciences*, 156–171.
- [4] Billah, M., Shaty, S. S., Sadnan, G. A., Hasnain, K. N., Abed, J., Begum, M., & Sultana, K. S. (2024). Performance Optimization in Multi-Machine Blockchain Systems: A Comprehensive Benchmarking Analysis. *Journal of Business and Management Studies*, 6(6), 357–375.
- [5] Bhowmik, P. K., Chowdhury, F. R., Sumsuzzaman, M., Ray, R. K., Khan, M. M., Gomes, C. A. H., ... & Gomes, C. A. (2025). AI-Driven Sentiment Analysis for Bitcoin Market Trends: A Predictive Approach to Crypto Volatility. *Journal of Ecohumanism*, 4(4), 266–288.
- [6] Caton, S., & Haas, C. (2020). Fairness in Machine Learning: A Survey. *Journal of AI Research*, 67, 1–40.
- [7] Das, B. C., Ahmad, M., & Maqsood, M. (2025). Strategies for Spatial Data Management in Cloud Environments. In *Innovations in Optimization and Machine Learning* (pp. 181–204). IGI Global Scientific Publishing.
- [8] Das, B. C., Zahid, R., Roy, P., & Ahmad, M. (2025). Spatial Data Governance for Healthcare Metaverse. In *Digital Technologies for Sustainability and Quality Control* (pp. 305–330). IGI Global Scientific Publishing.
- [9] European Parliament. (2016). Regulation (EU) 2016/679 (General Data Protection Regulation).
- [10] European Union. (2021). Proposal for a Regulation Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act).
- [11] Fariha, N., Khan, M. N. M., Hossain, M. I., Reza, S. A., Bortty, J. C., Sultana, K. S., ... & Begum, M. (2025). Advanced fraud detection using machine learning models: enhancing financial transaction security. arXiv preprint arXiv:2506.10842.
- [12] Feldman, M., Friedler, S. A., Moeller, J., Scheidegger, C., & Venkatasubramanian, S. (2015). Certifying and removing disparate impact. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 259–268.
- [13] Hardt, M., Price, E., & Srebro, N. (2016). Equality of opportunity in supervised learning. *Advances in Neural Information Processing Systems*, 29, 3315–3323.
-

- 
- [14] Hasan, M. S., Siam, M. A., Ahad, M. A., Hossain, M. N., Ridoy, M. H., Rabbi, M. N. S., ... & Jakir, T. (2024). Predictive Analytics for Customer Retention: Machine Learning Models to Analyze and Mitigate Churn in E-Commerce Platforms. *Journal of Business and Management Studies*, 6(4), 304–320.
- [15] Hasanuzzaman, M., Hossain, M., Rahman, M. M., Rabbi, M. M. K., Khan, M. M., Zeeshan, M. A. F., ... & Kawsar, M. (2025). Understanding Social Media Behavior in the USA: AI-Driven Insights for Predicting Digital Trends and User Engagement. *Journal of Ecohumanism*, 4(4), 119–141.
- [16] Hossain, M. I., Khan, M. N. M., Fariha, N., Tasnia, R., Sarker, B., Doha, M. Z., ... & Siam, M. A. (2025). Assessing Urban-Rural Income Disparities in the USA: A Data-Driven Approach Using Predictive Analytics. *Journal of Ecohumanism*, 4(4), 300–320.
- [17] Hossain, M. R., Mahabub, S., & Das, B. C. (2024). The role of AI and data integration in enhancing data protection in US digital public health: an empirical study. *Edelweiss Applied Science and Technology*, 8(6), 8308–8321.
- [18] Huang, Z., Liu, Q., & Wang, X. (2023). A Systematic Survey of Fairness in Machine Learning. *ACM Computing Surveys*, 55(2), Article 32.
- [19] Islam, M. R., Hossain, M., Alam, M., Khan, M. M., Rabbi, M. M. K., Rabby, M. F., ... & Tarafder, M. T. R. (2025). Leveraging Machine Learning for Insights and Predictions in Synthetic E-commerce Data in the USA: A Comprehensive Analysis. *Journal of Ecohumanism*, 4(2), 2394–2420.
- [20] Jakir, T., et al. (2023). Machine Learning-Powered Financial Fraud Detection: Building Robust Predictive Models for Transactional Security. *Journal of Economics, Finance and Accounting Studies*, 5(5), 161–180.
- [21] Kamiran, F., & Calders, T. (2012). Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems*, 33(1), 1–33.
- [22] Khan, M. A. U. H., Islam, M. D., Ahmed, I., Rabbi, M. M. K., Anonna, F. R., Zeeshan, M. D., ... & Sadnan, G. M. (2025). Secure Energy Transactions Using Blockchain Leveraging AI for Fraud Detection and Energy Market Stability. *arXiv preprint arXiv:2506.19870*.
- [23] Khan, M. N. M., Fariha, N., Hossain, M. I., Debnath, S., Al Helal, M. A., Basu, U., ... & Gurung, N. (2025). Assessing the Impact of ESG Factors on Financial Performance Using an AI-Enabled Predictive Model. *International Journal of Environmental Sciences*, 1792–1811.
- [24] Mahabub, S., Jahan, I., Islam, M. N., & Das, B. C. (2024). The Impact of Wearable Technology on Health Monitoring: A Data-Driven Analysis with Real-World Case Studies and Innovations. *Journal of Electrical Systems*, 20.
-

- 
- [25] Maw, L., Chen, J., & Lee, K. (2022). Fairness in Customer Churn Prediction Using Data Sampling Techniques. *IEEE Transactions on Emerging Topics in Computing*, 10(3), 1234–1245.
- [26] Rahman, M. S., Hossain, M. S., Rahman, M. K., Islam, M. R., Sumon, M. F. I., Siam, M. A., & Debnath, P. (2025). Enhancing Supply Chain Transparency with Blockchain: A Data-Driven Analysis of Distributed Ledger Applications. *Journal of Business and Management Studies*, 7(3), 59–77.
- [27] Springer, J. (2022). Modeling Customer Churn in E-Commerce: Techniques and Challenges. In *Resources in Business Analytics* (pp. 87–110). Springer.
- [28] Sultana, K. S., Begum, M., Abed, J., Siam, M. A., Sadnan, G. A., Shatyi, S. S., & Billah, M. (2025). Blockchain-Based Green Edge Computing: Optimizing Energy Efficiency with Decentralized AI Frameworks. *Journal of Computer Science and Technology Studies*, 7(1), 386–408.
- [29] Tiwari, P., Aggarwal, M., & Singh, R. (2025). Fairness Metrics and Mitigation Techniques in Supervised Learning. *Journal of AI Ethics*, 2(1), 25–50.