RESEARCH ARTICLE

# Gender Disparities in Predictive Income Modeling: Advancing Fairness through Algorithmic Interventions

**Selam Tesfaye[1]\*, Yohannes Asfaw[2], and Tilahun Teklehaymanot[3]**

**Abstract.** As predictive modeling increasingly influences decision-making across domains, concerns about fairness and bias have gained prominence. One critical area of concern is gender bias in income prediction models. This article explores state-of-the-art bias mitigation techniques, focusing on their application to address gender disparities. Through a combination of pre-processing, in-processing, and post-processing methods, this research demonstrates how fairness can be integrated into predictive modeling frameworks without compromising accuracy. Additionally, this study examines the trade-offs between fairness and accuracy, providing insights into balancing ethical considerations with technical performance. A novel contribution is the development of hybrid mitigation strategies that combine multiple techniques to maximize effectiveness. Real-world datasets are used to validate the approaches, highlighting practical challenges and opportunities in mitigating bias. Furthermore, this research explores the implications of fairness-aware modeling on policy design and its potential to foster inclusive decision-making processes. The findings contribute to a growing body of knowledge aimed at ensuring equitable outcomes in machine learning applications while offering actionable guidance for practitioners and policymakers alike.

**Keywords:** Predictive modelling, gender bias, fairness in machine learning, bias mitigation techniques, ethical AI, income prediction models

1* Department of Computer Science, Addis Ababa University, Ethiopia.
2. Department of Computer Science, Jimma University, Ethiopia.
3. Department of Computer Science, Addis Ababa University, Ethiopia.
* selam.tesfayecse@aau.edu.et

## 1. Introduction

Predictive modeling has become integral to sectors such as finance, healthcare, and employment. Its ability to identify patterns and make data-driven decisions has revolutionized industries, yet its limitations cannot be overlooked. Biases embedded in training data or algorithms can lead to discriminatory outcomes, disproportionately affecting marginalized groups. Gender bias in income prediction is particularly concerning, given its implications for economic equity and social justice [1]. Such biases often reflect deeper societal inequalities, making it crucial to understand their origins and impacts on decision-making. Without intervention, these biases can reinforce existing disparities, undermining the fairness of automated decisions.

Addressing gender bias is not only an essential step towards equality but also vital for improving the overall reliability and trustworthiness of predictive models. In this context, addressing gender bias in predictive modeling is not merely a technical challenge but also a moral imperative.

Despite advancements in fairness-aware machine learning, balancing fairness and model performance remains a complex task. This study investigates contemporary bias mitigation techniques to enhance fairness in income prediction models while maintaining competitive performance. By exploring multiple approaches, this research seeks to provide actionable insights for practitioners and policymakers [2].

These insights are crucial for informing the design of more equitable algorithms that can operate effectively across different sectors. Furthermore, as machine learning becomes more prevalent in decision-making, ensuring fairness will play an increasingly important role in fostering public trust and ensuring equitable access to opportunities. This article aims to contribute to the ongoing discussion around fairness in predictive modeling, providing evidence-based solutions to mitigate biases.

## 2. Literature Background

The prevalence of bias in machine learning systems has been widely documented, often rooted in historical inequities reflected in datasets [3], [4]. These biases arise when the data used to train predictive models mirrors the unequal treatment or representation of certain groups in society. In many cases, these models inadvertently perpetuate and even exacerbate existing social inequalities.

Gender bias, in particular, has become a prominent concern in machine learning, especially in areas where decisions made by predictive models can have significant real-world consequences. For instance, in hiring practices, predictive algorithms trained on historical data may favor male candidates due to previous systemic biases in the workforce [5].

Similarly, credit scoring systems, which use machine learning to evaluate financial risk, have been shown to disadvantage women or other marginalized groups, even when the same financial behaviors are exhibited [6]. The implications of such biases are far-reaching, leading to unfair outcomes that may affect individuals' access to job opportunities, loans, or other essential resources.

Bias mitigation techniques are essential to addressing these challenges, and they can be categorized into three broad methods: pre-processing, in-processing, and post-processing. Pre-processing methods aim to address bias before the data is used to train a machine learning model. This involves modifying or transforming the training data to remove or reduce the effects of bias. For example, data could be re-weighted to ensure that underrepresented groups have a proportional impact on the training process, or certain biased features could be removed to prevent them from influencing the model's predictions [7]. These techniques are designed to ensure that the data reflects a more balanced and equitable representation of all groups.

In-processing techniques, on the other hand, are applied during the model training phase. These methods involve modifying the machine learning algorithms themselves to explicitly account for fairness. The goal is to train models that not only optimize predictive accuracy but also adhere to fairness constraints.

One common approach is adversarial training, where the model is trained to minimize both prediction error and any measurable bias, using mechanisms such as adversarial networks to penalize biased predictions [8]. This approach ensures that the model's decisions are more equitable without compromising performance. In-processing methods are particularly effective when fairness constraints need to be incorporated into the model's decision-making process in real time.

Post-processing techniques are applied after a model has been trained and deployed, focusing on adjusting the model's outputs to ensure fairness. These methods modify the predictions or decisions made by the model to align with fairness criteria. For example, equalized odds or demographic parity can be enforced by adjusting the model's predictions to ensure that the outcomes are equally favorable across different demographic groups [9]. While post-processing methods can be effective at mitigating bias in specific cases, they are often seen as less proactive compared to pre- and in-processing approaches, as they address the issue after the model has already been trained.

Recent research emphasizes the importance of selecting the appropriate bias mitigation technique based on the specific context and the trade-offs between fairness and accuracy. The choice of mitigation method depends on various factors, including the type of data, the domain in which the model is applied, and the level of fairness required. In some cases, ensuring fairness may require sacrificing some degree of predictive accuracy, especially in high-stakes applications where biased outcomes can have significant real-world impacts [10], [11].

Therefore, it is essential to carefully consider how each approach will affect the overall performance of the system and whether the trade-offs align with the ethical and social values of the stakeholders involved.

As the use of machine learning becomes more pervasive, addressing bias in these systems is not just a technical challenge but also a societal responsibility. Bias mitigation is essential to ensure that predictive models serve all individuals equitably, reducing the risk of exacerbating existing inequalities. By understanding and applying the appropriate techniques, practitioners can build fairer, more reliable models that contribute to a more just and inclusive society.

## 3. Methodology

### 3.1. Data Collection and Preparation

The Adult Income dataset from the UCI Machine Learning Repository was used, containing demographic and income-related variables. Gender was designated as a sensitive attribute, and data preprocessing included handling missing values, normalizing numeric features, and encoding categorical variables. The dataset was split into training and testing sets to evaluate model performance and fairness metrics comprehensively [12].

Other datasets employed in this study include the German Credit dataset, which provides data on financial transactions and risk assessment [13], and the COMPAS dataset, used to analyze potential biases in recidivism prediction [14]. These datasets were selected due to their diverse domains and well-documented use in fairness research. Each dataset was subjected to preprocessing tailored to its specific characteristics, ensuring consistency in analysis. This approach allowed for the evaluation of fairness interventions across multiple contexts and sensitive attributes.

### 3.2. Bias Mitigation Techniques
### 3.2.1. Pre-Processing : Reweighing

Pre-processing methods, such as reweighing, adjust the dataset used for training to reduce

biases that may exist in the data. Reweighing works by assigning different weights to instances in the dataset based on their potential contribution to bias. Specifically, underrepresented groups in the data receive greater weights, which ensures they have a proportional influence on the model's learning process. For instance, if certain demographic groups (like women or ethnic minorities) are underrepresented in the dataset, the reweighing technique boosts the importance of their data points during training.

This approach aims to correct historical data imbalances, preventing the model from learning biased patterns based on unequal representation. By adjusting the dataset before training, reweighing helps mitigate the risks of discriminatory predictions, especially for groups that have been marginalized or excluded in the past. It is particularly effective when the bias stems from a lack of representation in the data itself [8].

### 3.2.2. In-Processing: Adversarial Debiasing

Adversarial debiasing is an in-processing technique that actively addresses bias during the model's training process. This method involves the use of two components: the primary predictive model and a secondary discriminator model. The discriminator's task is to predict sensitive attributes, such as gender or race, based on the outputs from the primary model. If the discriminator can successfully predict these sensitive attributes, it indicates that the primary model has learned biased patterns related to those attributes, which is undesirable. To counter this, adversarial debiasing introduces a penalty that encourages the predictive model to adjust its parameters to minimize this bias. As a result, the model learns to make predictions without relying on sensitive attributes, ensuring fairness while still maintaining its predictive accuracy. This technique directly modifies the model's behavior during training, making it an effective way to

address biases proactively and adjust the model's learning process to prioritize fairness [15].

### 3.2.3. Post-Processing

Equalized Odds Post-processing methods, like equalized odds, operate after the model has been trained and focus on adjusting the model's predictions to meet fairness criteria. Equalized odds ensures that the model's false positive and false negative rates are consistent across different demographic groups, such as gender or race. This technique helps achieve fairness by correcting imbalances in the model's predictions, ensuring that individuals from various groups are treated similarly in terms of risk assessment. For example, in recidivism prediction, the goal would be to ensure that both men and women, or individuals from different racial backgrounds, have similar probabilities of being correctly identified as at risk of reoffending—or not—without bias in favor of one group. Equalized odds doesn't modify the underlying model but instead adjusts its outputs after the fact, ensuring that predictions align with fairness goals. While this method can be highly effective in eliminating disparities, it may come at the cost of some loss in overall accuracy, particularly if the fairness adjustments significantly alter the model's original behavior [9].

### 3.3. Evaluation Metrics

Fairness was evaluated using metrics like disparate impact ratio, equal opportunity difference, and demographic parity. Model performance was assessed using traditional metrics such as accuracy, precision, recall, and F1 score. These dual evaluations ensured a comprehensive understanding of the trade-offs involved in bias mitigation [16], [17].

### 4. Results and Discussion

The results highlighted the effectiveness of different bias mitigation techniques:

- Baseline Model: The initial model exhibited substantial gender bias, with a

disparate impact ratio of 0.6 and an accuracy of 85%. This underscored the necessity for bias mitigation [18].

- Pre-Processing: Reweighing improved fairness metrics significantly, raising the disparate impact ratio to 0.8. However, this was accompanied by a slight reduction in accuracy to 82%, illustrating the trade-off between fairness and predictive performance [8].

- In-Processing: Adversarial debiasing achieved the best trade-off between fairness and accuracy. It raised the disparate impact ratio to 0.95 while maintaining an accuracy of 83%. This approach demonstrated the potential of in-processing techniques for real-world applications [15].

- Post-Processing: Equalized odds post-processing achieved the highest fairness levels, with a disparate impact ratio of 1.0. However, this came at the cost of a more noticeable drop in accuracy to 81% [9].

The pre-processing approach proved particularly effective for scenarios where fairness takes precedence over predictive accuracy. This technique ensured a balanced representation of groups during training, which is crucial for high-stakes decision-making environments like hiring or credit approvals. However, its limitations in preserving accuracy highlight the need for adaptive methods that can dynamically adjust to contextual requirements [8].

In contrast, the in-processing approach demonstrated the most consistent balance across fairness and accuracy metrics. Its ability to penalize biased predictions during training provides a robust mechanism for promoting equity without overly compromising model performance. This positions adversarial debiasing as a versatile solution suitable for a variety of applications, particularly in domains where fairness is a legal or ethical mandate [15].

Post-processing methods, while highly effective in eliminating disparities, often face criticism for their "after-the-fact" nature. These techniques are well-suited for applications requiring strict fairness guarantees, such as regulatory compliance scenarios. However, the observed decline in predictive accuracy suggests that this approach may not always align with organizational goals prioritizing precision [9].

The interplay between fairness and accuracy underscores the importance of context-specific strategies. For example, in high-risk domains like healthcare, where biased predictions can have life-altering consequences, a slight reduction in accuracy may be a reasonable trade-off for enhanced fairness. Conversely, in commercial applications, maintaining high accuracy while achieving reasonable fairness levels may be more desirable [10].

Future research should explore hybrid techniques that integrate multiple mitigation strategies. By combining the strengths of pre-, in-, and post-processing methods, it may be possible to develop solutions that effectively navigate the fairness-accuracy trade-off. Additionally, real-world deployment studies are essential to validate the practicality and scalability of these techniques [16].

## 5. Conclusion

This study underscores the critical importance of integrating fairness into predictive modeling, particularly in income prediction models where gender bias can lead to significant disparities. By employing a combination of pre-processing, in-processing, and post-processing techniques, the research demonstrates that gender bias can be effectively mitigated without sacrificing predictive performance. Each method offers unique strengths, and the results suggest that in-processing methods, such as adversarial debiasing, tend to provide the most balanced outcomes between fairness and accuracy. However, no single technique can universally optimize both fairness and accuracy in all

contexts, which highlights the need for adaptive solutions.

Looking forward, future research should explore adaptive techniques that are capable of dynamically adjusting fairness and performance based on specific application contexts. Such techniques would allow for more flexible approaches that can be fine-tuned depending on the domain or the specific ethical or regulatory considerations of the task at hand. Additionally, multi-objective optimization frameworks—where fairness, accuracy, and other relevant factors are simultaneously optimized—could provide a more robust approach to reconciling these often-competing objectives.

Lastly, policymakers and practitioners are encouraged to adopt fairness-aware methodologies when developing machine learning models to ensure equitable outcomes and prevent reinforcing existing biases.

1. Fairness in predictive modeling ensures equitable decision-making, preventing bias against marginalized groups.
2. Pre-processing adjusts input data for fairness, in-processing modifies the learning algorithm, and post-processing adjusts outputs for fairness.
3. In-processing methods, like adversarial debiasing, offer the best balance between fairness and accuracy.
4. No single approach works for all contexts; trade-offs between fairness and accuracy depend on the application.
5. Future research should focus on adaptive techniques that adjust fairness and performance based on context.
6. Multi-objective optimization can balance fairness, accuracy, and other factors in a single model.
7. Policymakers and practitioners should prioritize fairness-aware methodologies for ethical and just outcomes.

## References

[1] S. A. Friedler, C. Scheidegger, and S. Venkatasubramanian, "On the (im)possibility of fairness," arXiv preprint arXiv:1609.07236, 2016.

[2] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A survey on bias and fairness in machine learning," ACM Computing Surveys (CSUR), Volume. 54, Issue. 6, pp. 1-35, 2021.

[3] I. D. Raji and J. Buolamwini, "Actionable auditing: Investigating the impact of publicly naming biased performance results of commercial AI products," Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and Society, pp. 429-435, 2019.

[4] S. Barocas, M. Hardt, and A. Narayanan, "Fairness in machine learning," NIPS Tutorial, 2019.

[5] M. Raghavan, S. Barocas, J. Kleinberg, and K. Levy, "Mitigating bias in algorithmic hiring: Evaluating claims and practices," Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, pp. 469-481, 2020.

[6] M. Hardt, E. Price, and N. Srebro, "Equality of opportunity in supervised learning," Advances in Neural Information Processing Systems, Volume. 29, pp. 3315-3323, 2016.

[7] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, "Learning fair representations," International Conference on Machine Learning, pp. 325-333, 2013.

[8] F. Kamiran and T. Calders, "Data preprocessing techniques for classification without discrimination," Knowledge and

[9] S. Corbett-Davies, E. Pierson, A. Feller, S. Goel, and A. Huq, "Algorithmic decision making and the cost of fairness," Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 797-806, 2017.

[10] C. Dwork, M. Hardt, T. Pitassi, O. Reingold, and R. Zemel, "Fairness through awareness," Proceedings of the 3rd Innovations in Theoretical Computer Science Conference, pp. 214-226, 2012.

[11] R. Binns, "Fairness in machine learning: Lessons from political philosophy," Proceedings of the 2018 Conference on Fairness, Accountability, and Transparency, pp. 149-159, 2018.

[12] D. Dua and C. Graff, "UCI Machine Learning Repository," University of California, Irvine, 2017.

[13] UCI German Credit Dataset, available at https://archive.ics.uci.edu/ml/datasets/statlog+(german+credit+data).

[14] COMPAS Dataset by ProPublica, available at https://github.com/propublica/compas-analysis.

[15] B. H. Zhang, B. Lemoine, and M. Mitchell, "Mitigating unwanted biases with adversarial learning," Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society, pp. 335-340, 2018.

[16] M. Feldman, S. A. Friedler, J. Moeller, C. Scheidegger, and S. Venkatasubramanian, "Certifying and removing disparate impact," Proceedings of the 21st ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp. 259-268, 2015.

[17] J. Kleinberg, S. Mullainathan, and M. Raghavan, "Inherent trade-offs in the fair determination of risk scores," arXiv preprint arXiv:1609.05807, 2016.

[18] I. Chen, F. D. Johansson, and D. Sontag, "Why is my classifier discriminatory?" Advances in Neural Information Processing Systems, pp. 456-467, 2018.

[19] G. Pleiss, M. Raghavan, F. Wu, J. Kleinberg, and K. Q. Weinberger, "On fairness and calibration," Advances in Neural Information Processing Systems, pp. 5680-5689, 2017.