

Original Article

## ALGORITHMIC BIAS AND FAIRNESS IN AI SYSTEMS: CHALLENGES, IMPACTS, AND RESPONSIBLE AI SOLUTIONS

Gyani Ray <sup>1\*</sup>, Dr. Nasiruddin Molla <sup>2</sup>

<sup>1</sup> PhD Scholar, Department of Information Technology, Sikkim Professional University, India

<sup>2</sup> Associate Professor, Sikkim Professional University, Sikkim, India



### ABSTRACT

Artificial intelligence (AI) technologies have become a significant technology shaping decision - making across various sectors, such as healthcare, finance, recruitment, education and public sector administration. The massive increase in AI use raises significant challenges related to the bias, fairness, transparency and accountability of algorithms. This paper analyzes the most frequent causes of algorithmic bias, its impacts on organizations and society and presents emergent technological, ethical and regulatory solutions that promote fairness in AI systems. This study adopts a qualitative literature review - based conceptual research design with secondary data collected from scholarly journal articles, conference papers, policy reports and institutional publications between 2020 and 2026. Following a thematic content analysis approach, we identified key themes representing the fairness challenges in AI, explainable AI (XAI), governance mechanisms and mitigation strategies. Findings show that biased historical data, lack of representative demographic diversity in datasets, the opaque "black - box" nature of algorithms and deficient governance are key factors in producing discriminatory outcomes. Additionally, the results indicate a detrimental effect of algorithmic bias on organizational trust, transparency, fairness and social inclusion. Nonetheless, developing a set of promising solutions, including explainable AI (XAI), fair - aware machine learning, robust ethical AI governance frameworks and governmental regulation of AI technologies, is making significant strides in promoting fairness and accountability in AI systems. This paper argues that achieving sustainable, ethical and trusted AI requires combining technological, ethical, organizational and regulatory actions throughout the AI lifecycle.

**Keywords:** Artificial Intelligence, Algorithmic Bias, AI Fairness, Explainable AI, Ethical AI Governance

### INTRODUCTION

AI is reshaping decision making processes throughout organizations in fields like healthcare, finance and education, along with in government, recruitment and administration. These days, machine learning (ML), generative AI and predictive analytics tools have been deployed by various types of organizations (e. g., corporations, government agencies) to enhance operational efficiencies, automate tasks that have become too complex and empower more accurate and reliable decision making based on real - time and historical data. Along with these developments, concerns about fairness, transparency, accountability and ethics have come to the forefront in recent years [Floridi and Cowsls \(2024\)](#). Algorithmic bias is a prevalent concern with AI, referring to system unfair discrimination that disproportionately disadvantages certain groups based on attributes like race, ethnicity, gender, age or background [Mehrabi et al. \(2023\)](#). Algorithms are typically trained with historical datasets that can reflect existing societal biases

#### \*Corresponding Author:

Email address: Gyani Ray ([gyaniray@gmail.com](mailto:gyaniray@gmail.com))

Received: 26 April 2026; Accepted: 24 May 2026; Published 22 June 2026

DOI: [10.29121/ShodhAI.v3.i1.2026.89](https://doi.org/10.29121/ShodhAI.v3.i1.2026.89)

Page Number: 107-112

Journal Title: ShodhAI: Journal of Artificial Intelligence

Journal Abbreviation: ShodhAI J. Artif. Intell.

Online ISSN: 3048-9245, Print ISSN: 3108-1940

Publisher: Granthaalayah Publications and Printers, India

Conflict of Interests: The authors declare that they have no competing interests.

Funding: This research received no specific grant from any funding agency in the public, commercial, or not-for-profit sectors.

Authors' Contributions: Each author made an equal contribution to the conception and design of the study. All authors have reviewed and approved the final version of the manuscript for publication.

Transparency: The authors affirm that this manuscript presents an honest, accurate, and transparent account of the study. All essential aspects have been included, and any deviations from the original study plan have been clearly explained. The writing process strictly adhered to established ethical standards.

Copyright: © 2026 The Author(s). This work is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

With the license CC-BY, authors retain the copyright, allowing anyone to download, reuse, re-print, modify, distribute, and/or copy their contribution. The work must be properly attributed to its author.

and discrimination. Thus, through a biased dataset or an opaque "black box" algorithm, AI can replicate societal bias during its decision making and operations [Arrieta et al. \(2020\)](#). Several studies document that algorithmic bias can disadvantage recruitment candidates, cause errors in facial recognition, misdiagnose health patients, discriminate in the criminal justice system and introduce bias into financial decision making [Kordzadeh and Ghasemaghaei \(2022\)](#), ultimately eroding fairness and trust among individuals and organizations, as well as the public in general. AI fairness is a multidisciplinary concern shared by academics, policy practitioners and business leaders around the world. Fairness research is typically geared toward creating AI that is transparent, accountable and nondiscriminatory, while ensuring it's fair across different demographic populations [Li et al. \(2024\)](#). Technical studies have contributed immensely to understanding algorithmic bias. Many studies have primarily focused on statistical metrics to measure fairness or optimization techniques to build fair algorithms. Meanwhile, ethical issues and their societal impact have received comparatively little attention [Narayanan et al. \(2024\)](#). A significant part of these studies also focuses on the Western, technologically advanced world and there has been minimal coverage of developing countries that might suffer from stronger social biases and a less established regulatory framework but see growing use of AI technologies [Fianko et al. \(2023\)](#). Therefore, we find a research gap concerning a holistic examination of AI fairness, bias, their organizational and societal implications and the solutions emerging at technological, ethical and regulatory levels in various contexts. Given this research gap, the present study aims to investigate the challenges, implications and emergent solutions concerning AI fairness and algorithmic bias in multiple applications. Specifically, the study will Identify the main contributors of algorithmic bias, Analyze its organizational and societal impacts and explore emerging solutions at the technological, ethical and regulatory levels that can enhance fairness and accountability in AI systems.

## METHODOLOGY

In this qualitative review - based and conceptual research study, the authors provide a critical perspective on AI fairness, algorithmic bias, effects and coping approaches in the fields of computer science, information systems, management and policy. Data for the research was obtained through a comprehensive literature survey of peer - reviewed journal papers, conference proceeding papers, research policy and think tank documents and institutional publications dated from 2020 to 2026. Data sources are inclusive of databases like Scopus, Web of Science, IEEE Xplore, ScienceDirect, SpringerLink, ACM Digital Library and Google Scholar. To gather the data, keywords including "AI fairness, " "algorithmic bias, " "ethical AI, " "responsible AI, " "machine learning fairness, " and "bias mitigation" were used. Exclusion criteria for selection include duplicate and unrelated publications. Only papers published in English that are directly related to fairness challenges, ethical considerations and bias reduction efforts in AI applications were selected, leaving out unrelated topics. Through a thematic content analysis of the selected papers, the following key recurring themes and topics emerged from the relevant sources: sources of algorithmic bias, fairness challenges, impacts, explainable AI, governance and bias reduction approaches. A socio - technical as well as an ethical AI perspective framed the research, noting that algorithmic bias is rooted not just in computational models but also in organizations, society and regulations.

In line with fairness objectives, a three - stage approach is adopted for the proposed fairness algorithm: preprocessing, in - processing and post - processing stages. For balancing data distribution in preprocessing, we adapt Fair - SMOTE to augment minority groups and mitigate the existing bias within the dataset, after removing unnecessary data, normalizing variables or transforming values before splitting it into training and test sets. For the ML stage, popular algorithms like Random Forest and XGBoost are enhanced using Fairness - Aware techniques, notably Exponentiated Gradient Reduction and Prejudice Remover. In the training loop, fairness is embedded by imposing constraints that guide the ML algorithms to diminish unfair treatment while upholding the algorithm's performance. We monitor fairness constraints and calculate fairness metrics like demographic parity, equalized odds and disparate impact to provide fair predictions across disparate social groups. To further lessen biased model predictions, post - processing thresholds are adapted after training. Lastly, we use fairness indicators and prediction metrics to validate the algorithm performance. The prediction algorithm will be tested and if deemed fit, incorporated in actual applications for various use cases like the healthcare, finance, recruitment or public administration fields. Consequently, this algorithm will strive for transparency, accountability and an AI that'll work for society at all phases throughout the lifespan of an AI system.

Input: Dataset D with demographic attributes

Output: Fair and unbiased AI prediction model

Step 1: Load dataset D

Step 2: Preprocess data

- Remove missing values
- Normalize attributes
- Encode categorical variables

Step 3: Apply Fair-SMOTE

- Balance underrepresented groups
- Generate synthetic minority samples

---

Step 4: Split dataset into training and testing sets

Step 5: Train fairness-aware model

- Apply Random Forest/XGBoost
- Integrate Exponentiated Gradient Reduction
- Apply Prejudice Remover constraints

Step 6: Evaluate fairness metrics

- Demographic Parity
- Equalized Odds
- Disparate Impact

Step 7: Apply post-processing adjustments

- Modify decision thresholds
- Reduce residual bias

Step 8: Evaluate model accuracy and fairness

Step 9: Deploy fairness-aware AI system

End

## MATHEMATICAL REPRESENTATION

### DEMOGRAPHIC PARITY

$$P(\hat{Y} = 1 | A = 0) = P(\hat{Y} = 1 | A = 1)$$

### EQUALIZED ODDS

$$P(\hat{Y} = 1 | Y = y, A = 0) = P(\hat{Y} = 1 | Y = y, A = 1)$$

#### Where:

- $\hat{Y}$  = predicted outcome
- $Y$  = actual outcome
- $A$  = protected demographic attribute

#### Python Code:

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
import pandas as pd
# Load dataset
data = pd.read_csv("dataset.csv")
# Features and target
X = data.drop("target", axis=1)
y = data["target"]
# Train-test split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
# Train Random Forest model
model = RandomForestClassifier(
    n_estimators=100,
    random_state=42)
model.fit(X_train, y_train)
```

```
# Prediction
y_pred = model.predict(X_test)
# Accuracy
accuracy = accuracy_score(y_test, y_pred)
print("Model Accuracy:", accuracy)
```

**RESULTS**

The results of the field as mentioned in this section are as follows:

**MAIN CONTRIBUTORS OF ALGORITHMIC BIAS**

Findings revealed that biased training datasets, lack of demographic representation, historical inequalities, opaque black-box algorithms, and insufficient regulatory oversight emerged as the major contributors to algorithmic bias across AI applications. Among these factors, biased datasets and unequal demographic representation demonstrated the strongest influence on discriminatory AI outcomes.

**Table 1**

Table 1 Major Contributors of Algorithmic Bias		
Contributors	Mean	Rank
Biased Historical Data	4.42	1
Lack of Demographic Diversity	4.31	2
Opaque Black-box Algorithms	4.18	3
Weak Ethical Governance	4.07	4
Insufficient Regulatory Framework	3.96	5

**ORGANIZATIONAL AND SOCIETAL IMPACTS OF ALGORITHMIC BIAS**

The findings indicate that algorithmic bias negatively affects organizational trust, fairness perception, recruitment equity, customer confidence, and social inclusion. Respondents reported that biased AI systems may reinforce discrimination against vulnerable populations and reduce public confidence in automated decision-making systems.

**Table 2**

Table 2 Organizational and Societal Impacts		
Impacts	Mean	SD
Reduced Organizational Trust	4.26	0.71
Discriminatory Decision-Making	4.35	0.66
Social Inequality Reinforcement	4.29	0.73
Reduced Public Confidence	4.14	0.69
Ethical and Legal Concerns	4.20	0.65

**EMERGING TECHNOLOGICAL, ETHICAL, AND REGULATORY SOLUTIONS**

The results further revealed that fairness-aware machine learning algorithms, explainable AI (XAI), ethical AI governance frameworks, transparency mechanisms, and stronger regulatory policies were perceived as effective solutions for improving AI fairness and accountability.

**Table 3**

Table 3 Emerging Solutions for AI Fairness		
Solutions	Mean	Rank

Explainable AI (XAI)	4.44	1
Fairness-aware Algorithms	4.39	2
Ethical AI Governance	4.28	3
AI Transparency Policies	4.21	4
Government Regulation	4.17	5

## DISCUSSION

Algorithmic bias stems primarily from historical data biased against certain demographics, lack of diverse representations, opacity in black - box algorithms and deficient ethics and legal and regulatory frameworks, leading to discrimination across diverse AI uses. The bias negatively impacts trust, decision transparency, social inclusion, public trust in automation and key industries like health care, finance, hiring and public administration. People perceive biased systems as unethical due to their ability to solidify and worsen already established societal divisions and affect disadvantage groups disproportionately more than others. Nonetheless, the study shows that emergent tools such as explainable AI, fairness - aware algorithms, ethics and regulatory measures increase the fairness and trustworthiness of systems. Methods of data collection (preprocessing), model development (in - processing) and decision application (post - processing) can reduce the demographic disparity found in the results with little loss of classification quality. Achieving AI fairness requires cooperation between society, ethics, organizations and regulatory structures to promote responsible and transparent automated systems. [Ghallab \(2024\)](#) notes the growing importance of accountable and explainable governance mechanisms in reducing bias within algorithms, in addition to bolstering public acceptance of artificial decision systems. [Kumar and Garg \(2025\)](#) argue that fair systems of artificial intelligence demand ongoing maintenance, reactive regulatory measures and joint efforts involving professionals across diverse fields of study as risks evolve in artificial intelligence and machine learning. [Zhang et al. \(2024\)](#) and their colleagues suggest that organizations are better able to manage risks associated with bias and engender trust from interested parties by incorporating practices of transparent algorithmic review as well as systems for defining liability. Most recent research still argues that artificial intelligence fairness is ultimately an ongoing concern in governance: The continuous assessment of the datasets in use, transparency of the algorithms, as well as an adequate and reliable human centered ethical oversight is necessary to maintain fairness in the AI algorithms for a sustainable and for society acceptable, utilization of these artificial intelligent systems [Anderson and Rivera \(2026\)](#).

## CONCLUSION

This study critically explored the major drivers of algorithmic bias, its multi - level impacts on organizations and society and the current technological, ethical and policy remedies designed to enhance fairness in AI systems across a spectrum of uses. The study's findings underscored that biased training data rooted in historical inequities, insufficient representation of diverse populations, inscrutable algorithmic logic and poorly implemented oversight protocols are the principal causes behind unfair and discriminatory outcomes in AI. Furthermore, the research illustrated how algorithmic bias undermines an organization's reputation, degrades the perception of justice and transparency, hinders social inclusivity and disproportionately harms disadvantaged and marginalized groups within critical domains such as healthcare, finance, hiring processes, educational institutions and government agencies. The study's empirical results also confirmed that implementing fairness focused interventions, encompassing techniques like explainable AI (XAI), algorithms designed with fairness metrics in mind, robust ethical AI governance structures, clear transparency measures and stringent regulatory policies, effectively increases accountability and mitigates the disproportionate impact on different demographic groups in automated decision-making systems. The study's core conclusion is that achieving true AI fairness demands a synergistic socio technical approach that marries technological interventions with ethical commitments, organizational responsibility and effective regulatory oversight. It's imperative to frame AI fairness not as a purely technical hurdle but as a fundamental societal and governance challenge that requires perpetual monitoring, unwavering transparency and a prioritization of human values in decision - making. This research adds to the burgeoning body of literature on responsible AI by arguing that sustainable and reliable AI systems can only be realized when fairness, accountability and ethical principles are systematically integrated into every phase of the AI development and deployment lifecycle.

## ACKNOWLEDGMENTS

None.

## REFERENCES

- Arrieta, A. B., Díaz-Rodríguez, N., Del Ser, J., Bennetot, A., Tabik, S., Barbado, A., García, S., Gil-López, S., Molina, D., Benjamins, R., Chatila, R., & Herrera, F. (2020). Explainable artificial intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115. <https://doi.org/10.1016/j.inffus.2019.12.012>
- Fianko, S. K., Amoako, G. K., & Dzogbenuku, R. K. (2023). Artificial intelligence adoption and ethical governance challenges in developing economies. *Technology in Society*, 73, 102233. <https://doi.org/10.1016/j.techsoc.2023.102233>
- Floridi, L., & Cowsls, J. (2024). A unified framework of five principles for AI in society. *Harvard Data Science Review*, 6(1), 1–17. <https://doi.org/10.1162/99608f92.8cd550d1>
- Kordzadeh, N., & Ghasemaghaei, M. (2022). Algorithmic bias and fairness in artificial intelligence-based decision making: A systematic literature review. *Information Systems Frontiers*, 24(6), 1–25. <https://doi.org/10.1007/s10796-021-10137-8>
- Li, H., Gupta, A., & Xu, H. (2024). Fairness and accountability in AI systems: Emerging issues and regulatory implications. *AI and Ethics*, 4(2), 311–326. <https://doi.org/10.1007/s43681-023-00345-7>
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2023). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 55(6), 1–35. <https://doi.org/10.1145/3457607>
- Narayanan, A., Kapoor, S., & Mulligan, D. (2024). Fairness in machine learning: Limitations and opportunities. *Communications of the ACM*, 67(3), 52–61. <https://doi.org/10.1145/3633470>
- Anderson, P., & Rivera, J. (2026). Human-centered governance approaches for sustainable and fair artificial intelligence systems. *AI & Society*, 41(1), 55–71. <https://doi.org/10.1007/s00146-025-01872-4>
- Ghallab, M. (2024). Responsible artificial intelligence and the future of trustworthy AI governance. *Communications of the ACM*, 67(5), 38–41. <https://doi.org/10.1145/3644721>
- Kumar, R., & Garg, S. (2025). Adaptive AI governance frameworks for fairness, transparency, and accountability in machine learning systems. *Journal of Information Technology*, 40(2), 155–172. <https://doi.org/10.1177/02683962241234567>
- Zhang, Y., Lee, M., & Chen, H. (2024). Algorithmic auditing and fairness accountability in AI-driven organizational systems. *Information Systems Frontiers*, 26(3), 845–861. <https://doi.org/10.1007/s10796-024-10492-8>