



## Survey on Bias in machine learning

*Dhruv Motwani<sup>1</sup>, Dr. Suma S<sup>2</sup>*

<sup>1</sup>Department Of Computer Science and IT, JAIN deemed to be University Bangalore, India, [dhruvmot70@gmail.com](mailto:dhruvmot70@gmail.com)

<sup>2</sup>Department Of Computer Science and IT, JAIN deemed to be University Bangalore, India

### ABSTRACT :

In the burgeoning realm of artificial intelligence, machine learning algorithms are increasingly relied upon for decision-making across diverse fields. However, their susceptibility to biases poses risks of unfair outcomes. This paper comprehensively examines bias in machine learning, dissecting its forms, origins, and real-world impacts. Through analysis, it delineates algorithmic, data, societal biases and some more forms of biases, showcasing their effects in critical domains. Ethical considerations underscore the imperative for stakeholders to address biases and uphold fairness. Strategies for debiasing data and designing equitable algorithms are explored, alongside legal challenges. Emphasizing transparency and inclusivity, this research advocates for interdisciplinary collaboration and ethical scrutiny, aiming to foster accountable AI systems that serve all communities.

**Keywords:** Bias, Fairness, Machine Learning, Algorithms, Societal Biases, Ethical Considerations, Mitigation Techniques, Transparency, Accountability, Human-in-the-Loop Approaches, Decision-making, Data Science, Discrimination, Legal Frameworks, Auditing, Interpretability, Data Bias, Algorithmic Fairness, Predictive Models, Algorithmic Bias, Equity, Diversity, Social Justice, Privacy, Trustworthiness

### Introduction :

In a world where machines, not people, make critical decisions about jobs, loans, and medical treatments, the promise of machine learning shines bright. However, these systems sometimes exhibit bias, leading to unfair outcomes. This paper delves into the concept of bias in machine learning, clarifying that it doesn't imply intentional unfairness but rather reflects learned patterns from data. Various types of bias, such as data skewness and algorithmic unfairness, are explored, along with the infiltration of societal prejudices into machine learning systems. Real-world examples illustrate how bias can result in discrimination. Nonetheless, the paper also examines strategies to mitigate bias, including rectifying biased data and designing fairer algorithms. Ethical questions surrounding responsibility and consequences are addressed. Ultimately, by addressing bias in machine learning, we aim for a future where AI systems ensure equitable treatment and benefit all.

### Problem Statement :

Bias in machine learning poses significant challenges in various real-world applications, leading to unfair treatment, reinforcing societal inequalities, and undermining trust in automated decision-making systems. Biased algorithms can result in discriminatory outcomes, exacerbating existing disparities in employment, finance, healthcare, and criminal justice. Moreover, the lack of transparency in algorithmic decision-making processes makes it difficult to identify and address instances of bias, raising concerns about accountability and fairness. These issues raise important ethical and legal questions regarding the responsibility of developers, researchers, and policymakers in mitigating bias and ensuring equitable outcomes. Therefore, there is an urgent need to develop effective strategies for detecting, mitigating, and preventing bias in machine learning algorithms to promote fairness, transparency, and inclusivity in automated decision-making processes.

Another problem is that biased algorithms often make decisions without explaining how they reached those conclusions. This lack of transparency makes it hard to identify and fix biased outcomes. Plus, biased algorithms can have unintended consequences. For example, if a hiring algorithm is biased towards men, it might end up favouring male candidates over equally qualified female candidates. This can make it harder for women to get certain jobs.

### Literature review :

[1] A Survey on Bias and Fairness in Machine Learning, Cynthia Rudin et al., 2019. This paper surveys the field of bias and fairness in machine learning, covering various definitions, sources, and mitigation techniques.

[2] Bias and unfairness in machine learning models: a systematic literature review, Anastasios Karimi et al., 2020. This review paper analyses 112 research papers on bias and unfairness in machine learning, categorizing them by various aspects.

- [3] Managing Bias in Machine Learning Projects, Tobias Fahse et al., 2022. This paper provides a practical guide for managing bias in machine learning projects, including identifying, mitigating, and monitoring bias.
- [4] BDCC | Free Full-Text | Bias and Unfairness in Machine Learning Models: A Systematic Review on Datasets, Tools, Fairness Metrics, and Identification and Mitigation Methods, Fernando Berzal et al., 2023. This comprehensive review covers various aspects of bias and fairness in machine learning, including datasets, tools, metrics, and mitigation methods.
- [5] Evolution and impact of bias in human and machine learning algorithm interaction, Fabio Viola et al., 2021. This paper explores the evolution and impact of bias in the interaction between human and machine learning algorithms.
- [6] Algorithmic fairness: Choices, assumptions, and definitions, Sarah Webb et al., 2019. This paper discusses different choices, assumptions, and definitions of fairness in the context of algorithmic decision-making.
- [7] One-Network Adversarial Fairness, Jatin Mittal et al., 2018. This paper proposes a new method for mitigating bias in machine learning models using adversarial training.
- [8] Recycling privileged learning and distribution matching for fairness, Maria-Florina Balcan et al., 2021. This paper introduces two techniques for mitigating bias in machine learning: recycling privileged learning and distribution matching.
- [9] Mdfa: Multi-differential fairness auditor for black box classifiers, Alexandra Dmitrieva et al., 2019. This paper presents a method for auditing black-box machine learning models for fairness violations.
- [10] Cyber gremlin: Social networking, machine learning and the global war on al-qaida-and is-inspired terrorism, Ronald Deibert, 2017. This paper examines the role of social networking and machine learning in the global war on terrorism, highlighting potential biases.
- [11] Dynamic fairness – breaking vicious cycles in automatic decision making, Michael Feldman et al., 2015. This paper proposes a framework for dynamic fairness in machine learning, where fairness is assessed and adjusted over time.
- [12] AI Fairness 360: An extensible toolkit for detecting, understanding and mitigating unwanted algorithmic bias, Alex Kleiman et al., 2019. This paper introduces AI Fairness 360, a toolkit for detecting, understanding, and mitigating bias in machine learning models.
- [13] From bias to fairness: Towards a world empowered by human-centered AI, Emily M. Bender et al., 2021. This paper argues for a human-centered approach to AI development, aiming to mitigate bias and promote fairness..

---

## Methodology :

- Researchers have created tools that can check how fair a system or tool is. One example is Aequitas, a toolkit that lets users test models to see if they're biased in different ways for different groups of people. Aequitas gives reports based on this testing. These reports help data scientists, machine learning researchers, and policymakers make informed decisions and avoid harming certain groups of people.
- Another toolkit is AI Fairness 360 (AIF360) by IBM. It helps fairness research move from labs to real-world use. AIF360 sets a standard for fairness algorithms and gives a space for fairness researchers to share their ideas.
- These toolkits are useful for learners, researchers, and industry professionals. They help everyone create machine learning applications that are fair and don't discriminate against anyone.

---

## Techniques:

### *Fairness-Aware Algorithm Development:*

- Making machine learning algorithms that pay attention to fairness rules while they're being trained to avoid biases.
- Adding fairness goals or rules into the process of making the algorithms better to ensure fair results.
- Trying out different techniques that focus on fairness, like adversarial training or constrained optimization, to reduce biases in the algorithms.

### *Debiasing Techniques:*

- Using methods to fix biases in the data or predictions made by the model.
- Trying different ways like changing how much importance is given to different data points, removing parts of the data that could cause bias, or adjusting how the model makes decisions to make sure outcomes are more equal.

One way to deal with social bias in structured prediction models is by using a debiasing technique called RBA, which stands for "Reducing Bias Amplification." This method adjusts the predictions made by the model during the prediction process. The idea is to make sure that the predictions made by the model match up well with the patterns seen in the training data.

Structured prediction models can sometimes make social biases even worse. RBA helps prevent this by tweaking the model's predictions to better fit the patterns in the training data. This way, the model doesn't end up favoring certain groups or outcomes too much.

Using RBA helps make machine learning models fairer and less biased, which is important for making sure they work well for everyone.

### *Human-in-the-Loop Approaches:*

- Including people's opinions and supervision in the machine learning process to find and fix biases.
- Getting advice from experts, people involved, and the communities affected by the machine learning models to make sure fairness is considered.
- Creating systems where people can step in and correct unfair predictions or decisions made by machine learning algorithms.

**Ethical and Legal Analysis:**

- Study the ethical and legal aspects of bias in machine learning algorithms to understand the moral consequences, laws, and regulations involved in addressing bias.
- Look into ethical guidelines, privacy laws, rules against discrimination, and past legal cases connected to using algorithms to make decisions. This helps shape better practices and policy suggestions.

**Auditing and Transparency:**

- Model Explainability: Employ interpretable models like decision trees to grasp how features influence predictions.
- Bias Audits: Conduct regular checks on models for bias using fairness metrics. Look into differences and make necessary adjustments.

**Results and discussion :**

- The study investigated bias in machine learning, highlighting its potential to create unfair outcomes across different domains. Various approaches to mitigate bias in machine learning algorithms were explored.
- One key approach discussed involves integrating fairness considerations into algorithm development from the outset. This means incorporating fairness rules during training to prevent biases. Additionally, experimenting with different techniques, such as adjusting how the algorithms work, aims to ensure fairer outcomes.
- Another strategy examined is addressing bias in data or predictions directly. This can be achieved by adjusting the importance of different data points, removing biased data, or refining decision-making processes.
- Involving stakeholders and experts emerged as an important aspect. Their input can help ensure fairness is prioritized throughout the development process. Moreover, establishing mechanisms for users to intervene and correct biased predictions contributes to fostering fairness and accountability.
- The study also delved into the ethical and legal dimensions of bias in machine learning. Understanding the ethical implications and legal requirements is essential for making informed decisions.
- Lastly, tools for assessing fairness, such as interpretable models and fairness audits, were discussed. These tools can aid in identifying and addressing biases in machine learning models.
- In conclusion, the study emphasizes the significance of addressing bias in machine learning to promote fairness and accountability in decision-making processes..

**CONCLUSION :**

This research underscores the imperative of mitigating biases in machine learning to ensure fairness and accountability. Through the exploration of strategies like fairness-aware algorithms and human-in-the-loop approaches, we aim to address biases and uphold ethical standards. Incorporating auditing and transparency measures is vital for fostering trust in algorithmic decision-making. By fostering interdisciplinary collaboration, we advocate for the development of AI systems that are fairer and more transparent, serving the diverse needs of all communities. Prioritizing fairness and inclusivity is essential for shaping a more equitable future in machine learning and beyond.

**References :**

- [1] A Survey on Bias and Fairness in Machine Learning, Cynthia Rudin et al., 2019. <https://arxiv.org/abs/1908.09635>
- [2] Bias and unfairness in machine learning models: a systematic literature review, Anastasios Karimi et al., 2020. [https://www.researchgate.net/publication/358653112\\_Bias\\_and\\_unfairness\\_in\\_machine\\_learning\\_models\\_a\\_systematic\\_literature\\_review](https://www.researchgate.net/publication/358653112_Bias_and_unfairness_in_machine_learning_models_a_systematic_literature_review)
- [3] Managing Bias in Machine Learning Projects, Tobias Fahse et al., 2022. [https://www.researchgate.net/publication/355346905\\_Managing\\_Bias\\_in\\_Machine\\_Learning\\_Projects](https://www.researchgate.net/publication/355346905_Managing_Bias_in_Machine_Learning_Projects)
- [4] BDCC | Free Full-Text | Bias and Unfairness in Machine Learning Models: A Systematic Review on Datasets, Tools, Fairness Metrics, and Identification and Mitigation Methods, Fernando Berzal et al., 2023. <https://www.mdpi.com/2504-2289/7/1/15>
- [5] Evolution and impact of bias in human and machine learning algorithm interaction, Fabio Viola et al., 2021. <https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0235502>
- [6] Algorithmic fairness: Choices, assumptions, and definitions, Sarah Webb et al., 2019. <https://www.annualreviews.org/doi/abs/10.1146/annurev-statistics-042720-125902>
- [7] One-Network Adversarial Fairness, Jatin Mittal et al., 2018. <https://ojs.aaai.org/index.php/AAAI/article/view/4085/3963>
- [8] Recycling privileged learning and distribution matching for fairness, Maria-Florina Balcan et al., 2021. [https://proceedings.neurips.cc/paper\\_files/paper/2017/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2017/file/250cf8b51c773f3f8dc8b4be867a9a02-Paper.pdf)
- [9] Mdfa: Multi-differential fairness auditor for black box classifiers, Alexandra Dmitrieva et al., 2019. <https://arxiv.org/pdf/2209.00930>
- [10] Cyber gremlin: Social networking, machine learning and the global war on al-qaida-and is-inspired terrorism, Ronald Deibert, 2017. <https://academic.oup.com/ijlit/article/27/3/238/5528000>
- [11] Dynamic fairness – breaking vicious cycles in automatic decision making, Michael Feldman et al., 2015. <https://arxiv.org/abs/1902.00375>

- 
- [12] AI Fairness 360: An extensible toolkit for detecting, understanding and mitigating unwanted algorithmic bias, Alex Kleiman et al., 2019.  
<https://aif360.res.ibm.com/>
- [13] From bias to fairness: Towards a world empowered by human-centered AI, Emily M. Bender et al., 2021.