

# Case Study Report: Amazon's Biased Hiring Tool

## Background

In 2018, Amazon discontinued an experimental AI hiring tool after discovering it showed significant gender bias against female candidates. The tool was designed to automate résumé screening, but it consistently rated male applicants more favorably for technical roles, regardless of comparable qualifications (Dastin, 2018).

## Source of Bias

The primary source of bias was training data bias. The model was trained on ten years of historical hiring data, which was predominantly male due to existing gender imbalances in the tech industry. Consequently, the AI system learned to favor résumés that resembled those submitted by male candidates. It even penalized those with words like "women's" in context such as "women's chess club" or all-female universities, reflecting encoded societal prejudices rather than objective skill assessment (Dastin, 2018; Raji et al., 2020).

## Proposed Solutions

To mitigate the bias and create a fairer system, the following interventions are recommended:

### 1. Data Rebalancing and Debiasing

- Create a new, balanced dataset by either collecting more female résumé examples or synthetically augmenting underrepresented classes.
- Remove gender proxies such as names, gendered pronouns, and gender-specific affiliations.
- Use adversarial debiasing techniques to reduce hidden gender signals in textual data.

### 2. Fairness-Constrained Model Training

- Implement fairness-aware algorithms such as the Disparate Impact Remover or Prejudice Remover Regularizer (Bellamy et al., 2018).
- During model tuning, include fairness metrics (e.g., equal opportunity difference) alongside performance metrics.

### 3. Human-in-the-Loop Governance

- Introduce a dual-layer review process, where AI filters applications but a sample-especially those rejected from underrepresented groups-is reviewed manually.
- Establish an ethics oversight team to continuously monitor and evaluate fairness metrics during deployment.

## **Fairness Evaluation Metrics**

Post-mitigation, fairness should be assessed using metrics such as:

- Disparate Impact Ratio (DIR): Measures the ratio of selection rates between groups (female vs. male), with an ideal range between 0.8 and 1.25.
- Equal Opportunity Difference: Assesses the difference in true positive rates between demographic groups.
- False Negative Rate Gap: Important to ensure qualified female candidates are not disproportionately filtered out.
- Average Odds Difference: Measures combined differences in false positives and true positives.
- Calibration by Group: Verifies that predicted probabilities align with real-world hiring outcomes across gender groups.

## **Conclusion**

The failure of Amazon's hiring tool highlights the urgent need for ethical AI systems that are transparent, fair, and accountable. Addressing AI bias requires interventions across the full development pipeline-from dataset curation to model evaluation and governance oversight. By adopting ethical frameworks and measurable fairness metrics, organizations can ensure AI augments human decision-making without reproducing historical inequities.

## **References**

Bellamy, R. K. E., Dey, K., Hind, M., Hoffman, S. C., Houde, S., Kannan, K., ... & Zhang, Y. (2018). AI Fairness 360: An extensible toolkit for detecting, understanding, and mitigating unwanted algorithmic bias. arXiv preprint arXiv:1810.01943.

Dastin, J. (2018, October 10). Amazon scraps secret AI recruiting tool that showed bias against women. Reuters. <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight-idUSKCN1MK08G>

Raji, I. D., Binns, R., Zliobaite, I., & Veale, M. (2020). The fallacy of AI function creep: How automation widens socio-economic gaps. Proceedings of the 2020 Conference on Fairness, Accountability, and

