# Case Study Analysis: Amazon's Biased AI Recruiting Tool

## Background

Amazon's AI recruiting tool, developed to screen resumes, was found to penalize female candidates, downgrading resumes with terms like "women's" or names of women's colleges. This case study analyzes the source of bias, proposes three fixes to ensure fairness, and suggests metrics to evaluate the corrected system, ensuring equitable hiring practices.

## Source of Bias

The primary source of bias was **biased training data**. The tool was trained on historical resume data from a male-dominated tech industry, where most hires were male. This skewed dataset embedded patterns favoring male candidates, as the model learned to associate male-associated terms (e.g., "executed," "competitive") with success, while penalizing female-associated terms (e.g., "women's leadership"). Additionally, **model design** contributed to bias by relying on word embeddings (e.g., word2vec) that amplified gendered correlations without mechanisms to detect or mitigate bias. Lack of **bias auditing** during development allowed these issues to persist undetected, similar to the need for `aif360` in your hospital readmission project to check disparate impact.

## Proposed Fixes

1. **Diversify Training Data**:
   - Curate a balanced dataset with equal representation of male and female hires across roles, industries, and seniority levels. Augment data with synthetic resumes or external datasets (e.g., LinkedIn) to include diverse qualifications and remove gendered terms during preprocessing. This mirrors cleaning `A1C_Result` in your Task 3 to ensure data quality.
2. **Implement Fairness-Aware Algorithms**:
   - Use fairness-aware models like Adversarial Debiasing (available in `aif360`) to minimize gender correlations in predictions. Adjust word embeddings to neutralize gendered terms (e.g., via debiasing techniques like Hard Debiasing). This ensures the model focuses on skills, not gender proxies, similar to ensuring fairness across gender in your hospital model.
3. **Regular Bias Audits and Retraining**:
   - Conduct periodic bias audits using tools like `aif360` to check for disparate impact. Retrain the model quarterly with updated, diverse data to adapt to evolving hiring trends. This aligns with your monitoring script for detecting performance drift in Task 3.

# Fairness Metrics

To evaluate fairness post-correction, use the following metrics:

- **Disparate Impact Ratio (DIR)**: Measure the ratio of positive outcomes (e.g., shortlisting) for female vs. male candidates. Target: DIR between 0.8 and 1.25 (as in your hospital project). Computed using `aif360`'s `BinaryLabelDatasetMetric`.
- **Equal Opportunity Difference (EOD)**: Assess the difference in true positive rates (hiring qualified candidates) across genders. Target: EOD close to 0, ensuring equal chances for qualified females and males.
- **Demographic Parity Difference (DPD)**: Evaluate the difference in selection rates across genders. Target: DPD near 0 to ensure balanced shortlisting. Use `aif360` for computation.
- **Accuracy and F1-Score by Group**: Measure model performance (accuracy, F1-score) separately for male and female candidates to ensure consistent predictive quality, similar to your hospital model's evaluation.

# Conclusion

The bias in Amazon's AI tool stemmed from skewed training data and unchecked model design, leading to unfair penalization of female candidates. By diversifying data, using fairness-aware algorithms, and conducting regular audits, the tool can achieve equitable outcomes. Fairness metrics like DIR, EOD, and DPD ensure ongoing accountability, aligning with ethical AI principles from Part 3 of your project. Stakeholders (HR, leadership) can trust a fair, transparent hiring process that maximizes talent diversity.

# Case Study Analysis: Facial Recognition in Policing

## Background

A facial recognition system used in policing misidentifies minorities at higher rates, raising ethical concerns. This analysis discusses the ethical risks of such systems and recommends policies for responsible deployment, ensuring public safety while protecting civil rights. The approach draws on fairness and monitoring principles from your hospital readmission project.

## Ethical Risks

1. **Wrongful Arrests and Misidentification**:
   - Higher false positive rates for minorities (e.g., Black and Hispanic individuals) can lead to wrongful arrests, as seen in cases like Robert Williams'

misidentification in 2020. This erodes trust in law enforcement and causes harm, similar to biased predictions in your hospital model risking unfair treatment.

2. **Privacy Violations**:
   o Mass surveillance using facial recognition infringes on privacy, especially when deployed without consent in minority communities. Unregulated data collection (e.g., from public cameras) risks misuse, paralleling HIPAA concerns in your Task 3.

3. **Systemic Bias Amplification**:
   o Biased training data (e.g., overrepresentation of lighter-skinned faces) perpetuates systemic inequities, leading to disproportionate targeting of minorities. This mirrors the gender bias in your hospital project's initial model outputs.

4. **Lack of Transparency**:
   o Black-box models obscure decision-making, making it hard for affected individuals to challenge misidentifications. This aligns with the need for explainability (e.g., SHAP) in your prior work.

# Recommended Policies for Responsible Deployment

1. **Mandatory Bias Testing and Reporting**:
   o Require pre-deployment testing using fairness metrics (e.g., False Positive Rate by demographic group) with tools like `aif360`. Publish results publicly to ensure transparency, similar to your hospital project's disparate impact checks. Retrain models if bias exceeds thresholds (e.g., FPR difference > 5%).

2. **Human-in-the-Loop Oversight**:
   o Mandate human review of facial recognition outputs before actions (e.g., arrests). Train officers on bias risks and require documentation of decisions, akin to your hospital dashboard's clinician-facing predictions. Use explainable AI (e.g., saliency maps) to justify matches.

3. **Strict Data Governance and Consent**:
   o Limit data collection to consented or anonymized sources, complying with privacy laws (e.g., GDPR, CCPA). Delete data after a fixed period (e.g., 30 days) unless required for active cases. This mirrors your hospital app's HIPAA-compliant non-storage policy.

4. **Community Engagement and Bans in High-Risk Cases**:
   o Engage minority communities in policy development to build trust. Ban facial recognition in high-stakes scenarios (e.g., real-time suspect identification) unless accuracy across demographics exceeds 95%, verified by independent audits.

# Conclusion

Facial recognition in policing poses ethical risks like wrongful arrests, privacy violations, and systemic bias, undermining public trust. Policies for responsible deployment include mandatory bias testing, human oversight, strict data governance, and community engagement. These measures ensure fairness and accountability, drawing on lessons from your hospital project's

bias mitigation and monitoring strategies. Stakeholders (police, policymakers, communities) benefit from a transparent, equitable system that balances safety and civil rights.