

The COMPAS dataset exhibits clear racial disparities in recidivism risk predictions. Analysis shows that Black defendants experience higher false positive rates compared to White defendants, meaning they are more likely to be incorrectly labeled as high-risk. This disparity reflects structural biases present in historical data and underscores the importance of fairness-aware modeling. Using IBM's AI Fairness 360 toolkit, we applied the reweighing algorithm to reduce bias in the training data. Post-mitigation metrics indicate a reduction in mean difference between privileged and unprivileged groups, suggesting improved fairness in model predictions. Visualizations highlight persistent differences in false positive rates, indicating that further interventions may be necessary. Recommended remediation steps include (1) rebalancing or augmenting datasets to better represent underrepresented groups, (2) applying fairness constraints during model training, (3) implementing post-processing adjustments to correct biased predictions, and (4) ensuring continuous monitoring and auditing of AI systems in deployment. Ethical deployment also requires transparency and explainability so that stakeholders understand how risk scores are computed. By integrating these practices, AI developers can reduce the risk of unjust outcomes, enhance public trust, and comply with ethical and legal standards such as GDPR. Overall, this audit illustrates that fairness interventions are both necessary and actionable in high-stakes domains like criminal justice.