

Bias Mitigation and Fairness Reflection - The Ethical Analysis

In many real-world applications, predictive models are deployed to support critical decisions such as medical diagnostics, loan approvals, or resource allocation. In such scenarios, it is essential to not only build models that are accurate but also ****fair**** — ensuring that all groups are treated equitably. This becomes especially important when certain demographic or group attributes (like age, race, gender, or department/team) may influence how predictions are interpreted or acted upon.

Lack of Demographic Data

The original Breast Cancer Wisconsin dataset used in this project does not contain any demographic information about the patients. In order to evaluate how well the model treats different groups (a fairness concern), we introduced a simulated column called `team`. This new column randomly assigns each sample to either "Team A" (70% of the data) or "Team B" (30%). This setup mimics real-world scenarios where one group might be dominant or overrepresented, and the other is underrepresented.

Disparities in Model Performance

When we trained the Random Forest model and evaluated its predictions, we found that it performed better on samples belonging to Team A than on Team B. Specifically, the ****F1-score**** and ****recall**** were lower for Team B.

- ****F1-score**** is a balanced measure of precision and recall. Lower F1-score for a group means the model is making more mistakes when handling data from that group.
- ****Recall**** measures how many actual positive instances were correctly predicted. A lower recall for Team B suggests the model is missing more high-priority cases for that group.

This imbalance is concerning because it implies the model may **unfairly disadvantage Team B** in real-world decision-making. For instance, if this model were deployed in a healthcare setting, underpredicting severity for Team B could lead to worse patient outcomes.

Bias Mitigation Techniques Applied

To address this unfairness, we implemented three key strategies:

1. SMOTE (Synthetic Minority Oversampling Technique)

- SMOTE is used to generate new, synthetic samples for underrepresented classes in the training data. In this project, it was applied to the `priority` target variable.
- This ensures that the model learns from a balanced distribution of all priority levels and doesn't favor the most common ones.

2. Group-Wise Fairness Evaluation:

- After training the model, we split the predictions based on the `team` variable and calculated separate classification metrics (precision, recall, F1-score) for each team.
- This allowed us to directly measure whether the model treats both groups equally. Identifying significant differences here is a signal that bias may be present.

3. Feature Importance Analysis:

- Feature importance tells us which variables the model relies on most for its predictions.
- If one feature (like `radius_mean`) dominates too heavily, it could act as a **proxy** for a sensitive attribute, potentially causing indirect bias.
- By analyzing and interpreting these weights, we ensure that no single variable unfairly drives the model's decisions.

Ethical Justification and Outcome

These fairness strategies represent our ****commitment to ethical and responsible AI****. Instead of accepting performance disparities as inevitable, we took proactive steps to:

- Audit model behavior across different groups
- Apply statistical techniques to balance training data
- Increase transparency via explainable AI methods (like feature importance plots)

Ultimately, the improved model was not only accurate, but also more equitable. Although we used synthetic demographic data (the `team` column), this process mirrors best practices in real-world deployments where fairness, accountability, and inclusiveness must be part of the design and evaluation process.

Fairness is not an optional add-on — it is essential. As machine learning continues to be embedded in high-stakes decision systems, evaluating and mitigating bias is critical to ensuring that models support, rather than harm, the people they are designed to serve.