# Fairness-Aware Machine Learning for Social Bias Detection in Healthcare Research Datasets

Ilerioluwakiiye Abolade, Precious Kolawole

ML Collective, Carleton University

## Motivation

AI models in healthcare can unintentionally discriminate against vulnerable groups. This work aims to detect and quantify biases in healthcare data and predictive models before deployment.

- How can researchers identify data-level and algorithm-level biases?
- How do neural networks compare to traditional models in balancing accuracy and fairness?

## Introducing the Social Bias Detection Tool

We present a lightweight, interactive tool to:

- Detect bias in healthcare data and models
- Compute fairness metrics (SPD, EOD, DD)
- Compare traditional ML vs neural nets on both accuracy and fairness
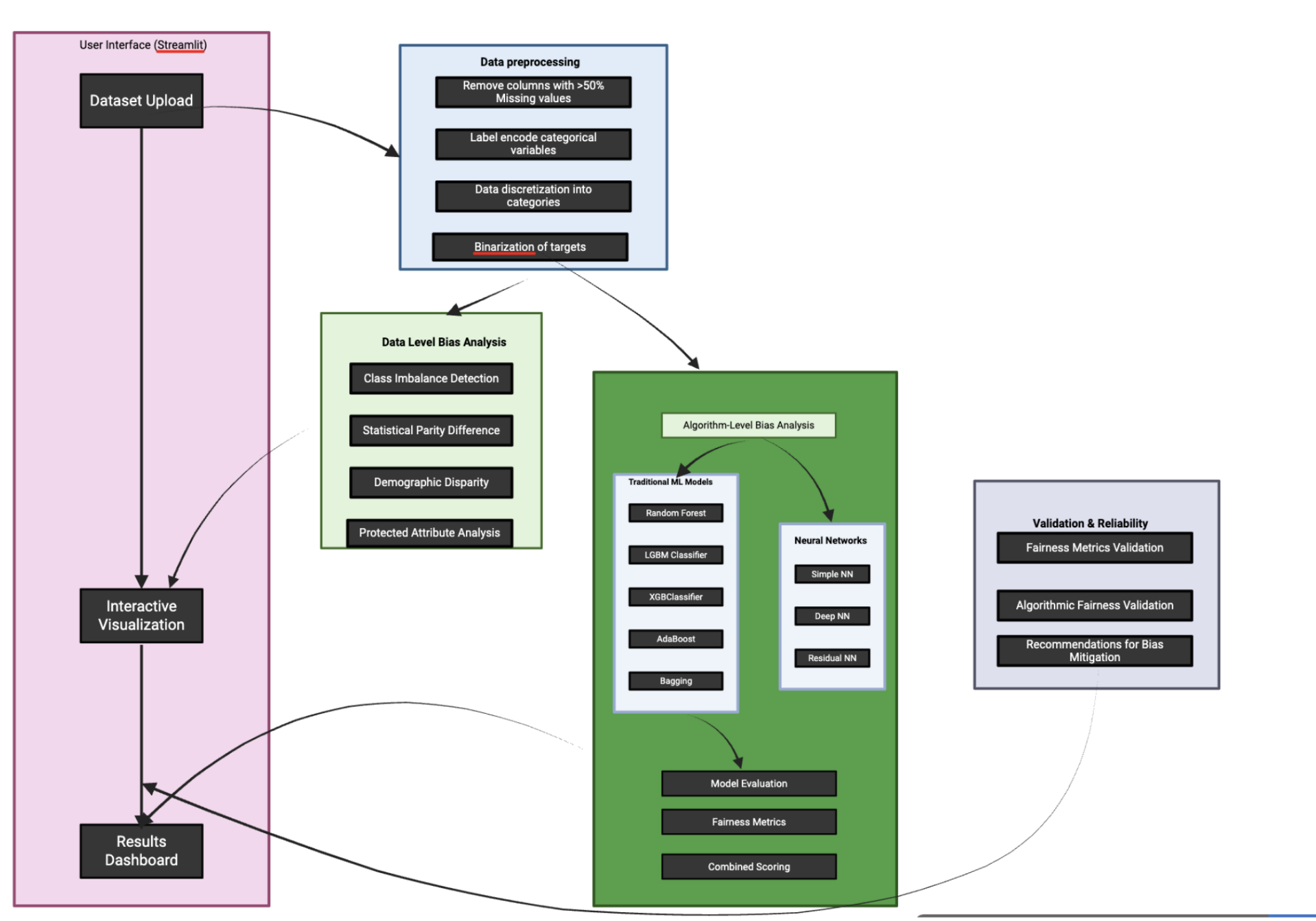- Output a "combined score" for model selection



Figure 1:Architectural overview of the Social Bias Detection framework

## Data-Level Bias Exists, Even Before Training

We evaluated two real-world healthcare datasets:

**SyntheticMass:**

- 83.6% White patients → major racial imbalance
- Age disparity: SPD = 0.82 for 0–35 vs 65+ (substantial)

**Brain Stroke Dataset:**

- Demographics more balanced overall
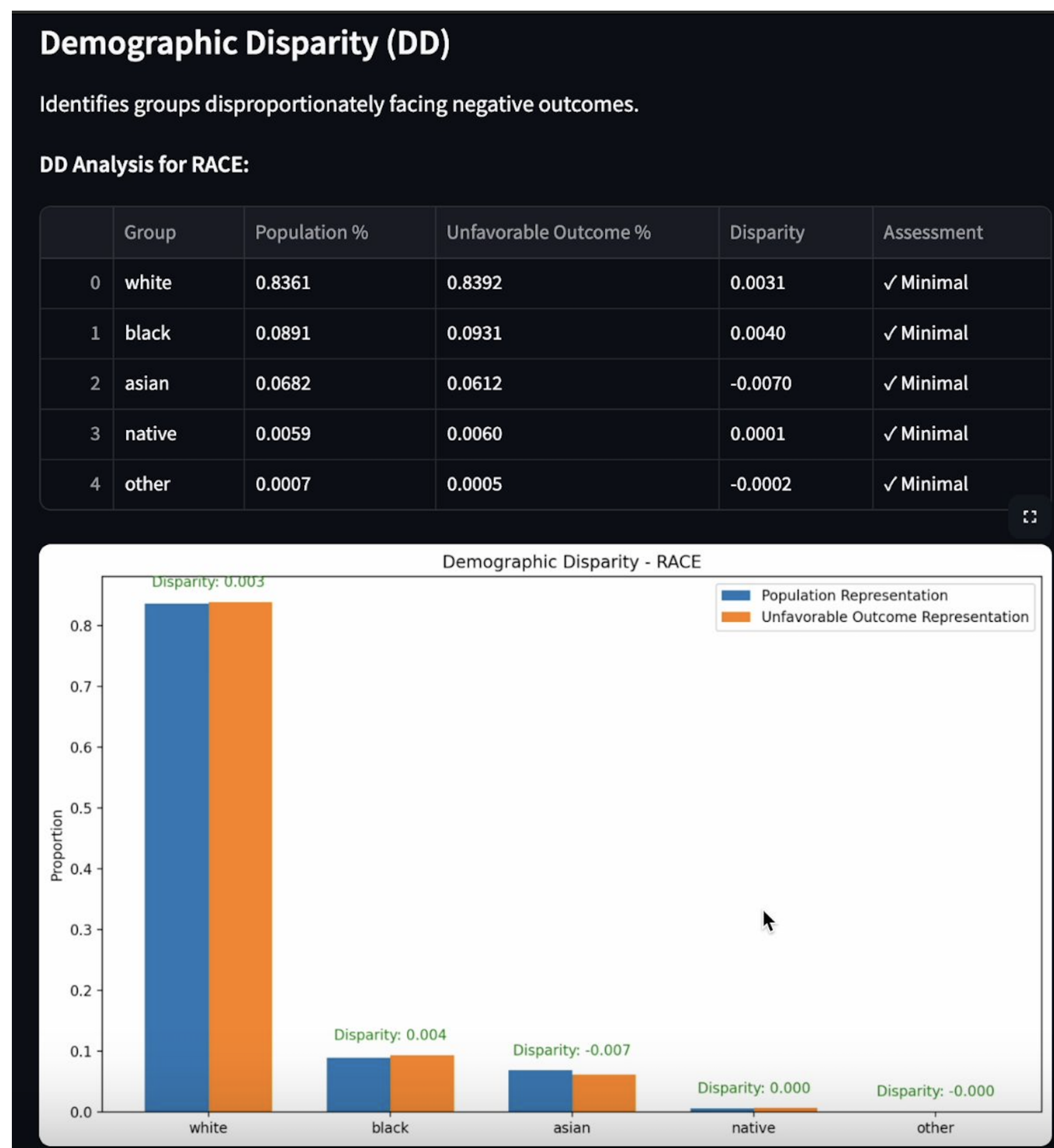- One exception: SPD = 0.10 for 65+ vs 51–65



Figure 2:Population vs. unfavorable outcomes by race, showing minimal to moderate pre-training bias.

This shows that bias can be embedded in the data itself, independent of modeling.

## Fairness Metrics Used

We assess bias using industry-standard metrics:

**1. Statistical Parity Difference (SPD)**
Difference in favorable outcomes across groups:

$$\text{SPD} = P(\hat{Y} = 1 \mid A = a) - P(\hat{Y} = 1 \mid A = b)$$

**2. Equal Opportunity Difference (EOD)**
Gap in true positive rates across groups:

$$\text{EOD} = \text{TPR}_a - \text{TPR}_b$$

**3. Average Odds Difference (AOD)**
Average gap in true *and* false positive rates:

$$\text{AOD} = \frac{1}{2}\left[(\text{FPR}_a - \text{FPR}_b) + (\text{TPR}_a - \text{TPR}_b)\right]$$

**4. Demographic Disparity (DD)**
Difference between a group's outcome share and its population share:

$$\text{DD} = P(A = a \mid Y = 1) - P(A = a)$$

*Interpretation thresholds:*
**0.00–0.05**: Minimal **0.05–0.10**: Small **> 0.10**: Substantial bias

## How the Tool Works

### 1. Data-Level Analysis:

- Class imbalance across demographic groups
- Statistical Parity Difference (SPD) between protected groups
- Demographic Disparity (DD) in outcomes vs. population share



### 2. Algorithmic Bias Analysis:

- Trains both traditional ML and neural networks
- Evaluates SPD, Equal Opportunity Difference (EOD), and Average Odds Difference (AOD)
- Computes a **Combined Score** = 0.5 $\times$ Normalized Accuracy + 0.5 $\times$ Fairness Score



## Main Findings

- Neural networks generally achieved higher fairness without sacrificing accuracy.
- ResidualNN scored highest in the SyntheticMass dataset for balancing both metrics.
- DeepNN achieved near-perfect fairness in Brain Stroke predictions.



Figure 3:Accuracy vs Fairness trade-offs across models

## How We Compare Models Fairly

High accuracy ≠ fair predictions. We introduce the **Combined Score** to balance both.

$$\text{Combined Score} = 0.5 \times \text{Normalized Accuracy} + 0.5 \times \text{Fairness Score}$$
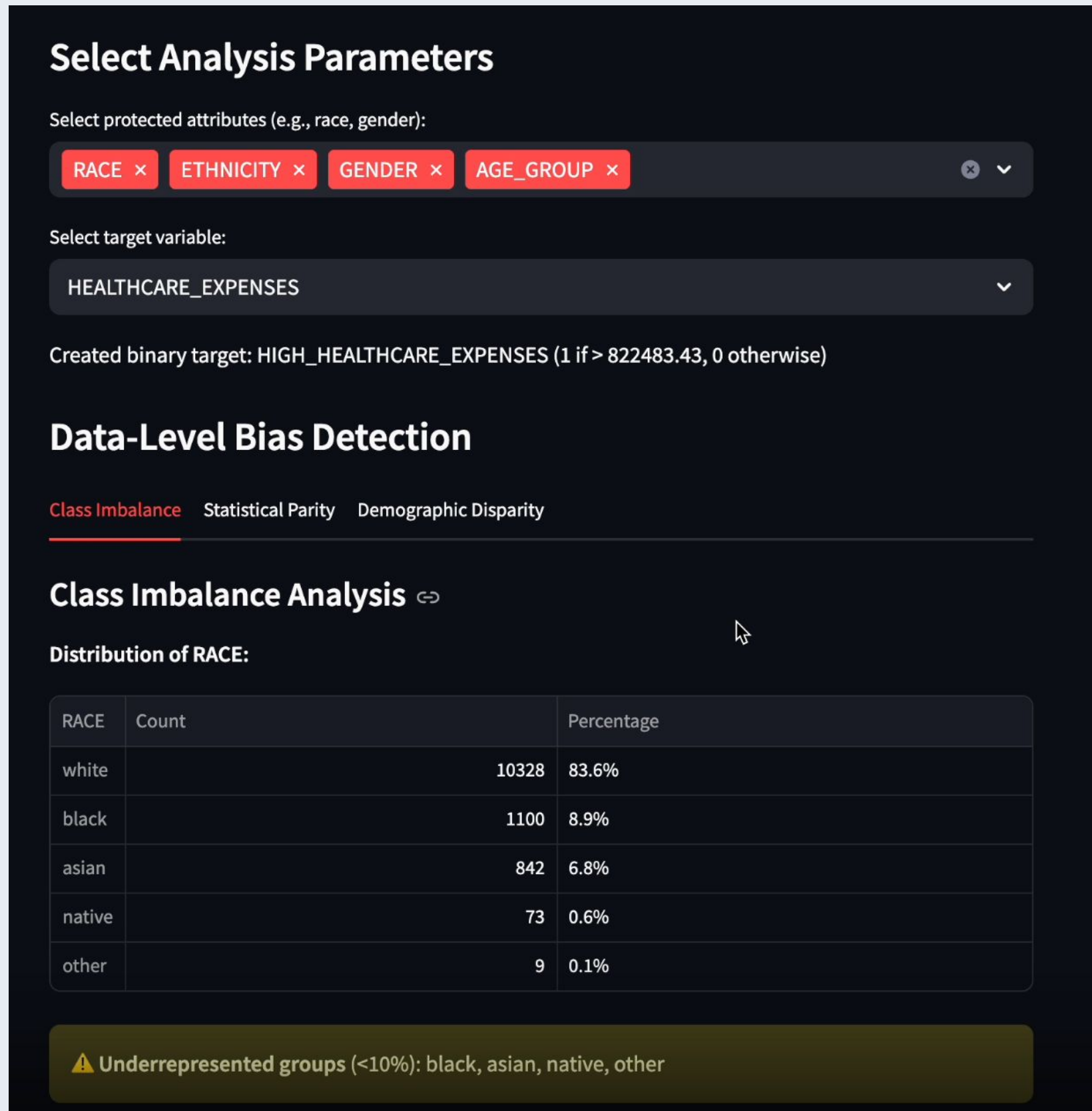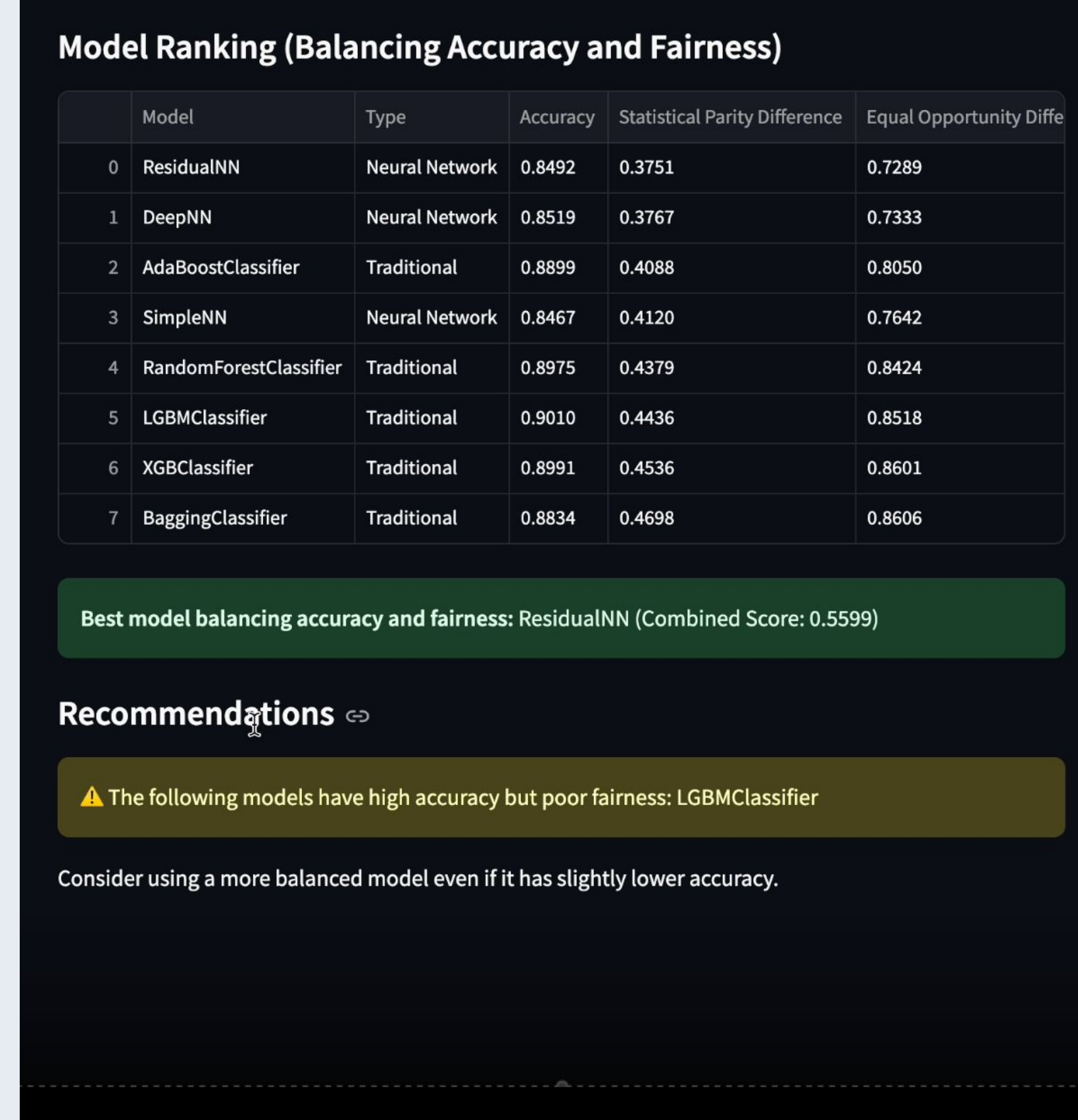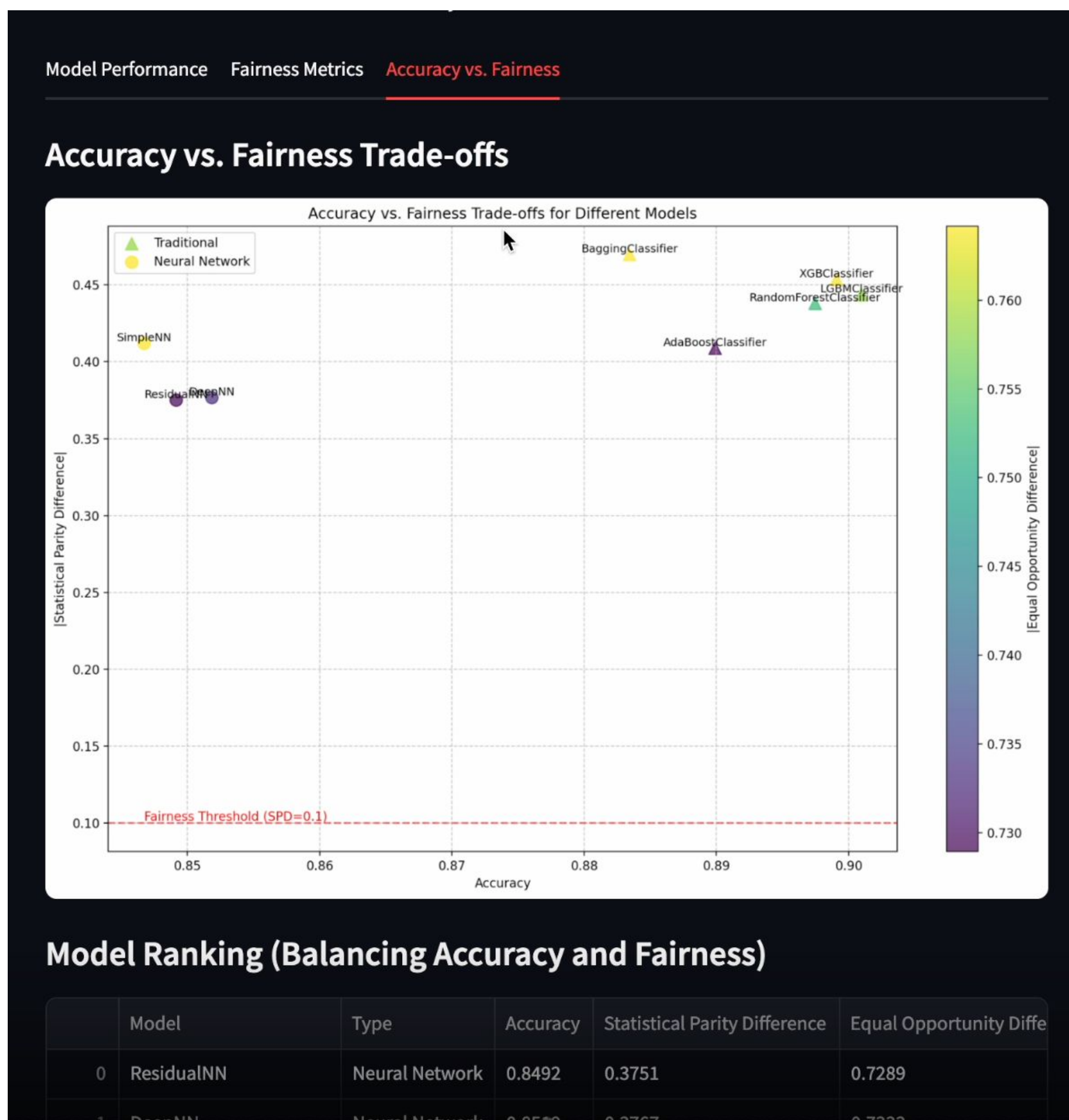
where:

$$\text{Fairness Score} = 1 - \frac{|\text{SPD}|}{|\text{SPD}_{\max}|} - \frac{|\text{EOD}|}{|\text{EOD}_{\max}|}$$

Lower SPD and EOD values increase the fairness score, rewarding models that treat groups more equally.

## Key Takeaways

- Bias is present in many healthcare datasets before training.
- Our tool enables quick, transparent bias assessment.
- Neural networks can be both accurate *and* fair.
- The Combined Score metric helps balance ethics with performance.

## Acknowledgements

## Project Repository

github.com/precillieo/social-bias-detection-tool

DEEP LEARNING INDABA