

Hw 2

XAI 2023/24 mimuw, 19.10.2023, Witold Drzewakowski

Task 1

I will start by calculating true positives, false positives, false negatives for two groups: red and blue.

$$TP_b = 0.6, FP_b = 0.05, FN_b = 0.2$$

$$TP_r = 0.25, FP_r = 0.25, FN_r = 0.25$$

Now let's calculate PPV and TPR .

$$PPV_b = \frac{TP_b}{TP_b + FP_b} = \frac{0.6}{0.65} = 12/13$$

$$PPV_r = \frac{TP_r}{TP_r + FP_r} = \frac{0.25}{0.5} = 0.5$$

$$TPR_b = \frac{TP_b}{TP_b + FN_b} = \frac{0.6}{0.8} = 0.75$$

$$TPR_r = \frac{TP_r}{TP_r + FN_r} = \frac{0.25}{0.5} = 0.5 \text{ (Alternatively)}$$

Now let's calculate probabilities of selecting a candidate from both groups:

$$R_r = P(\text{selected}|\text{red}) = 0.5$$

$$R_b = P(\text{selected}|\text{blue}) = 0.65$$

The predictive rate parity coefficient:

$$\frac{PPV_b}{PPV_r} = \frac{24}{13} \text{ — blue is more privileged}$$

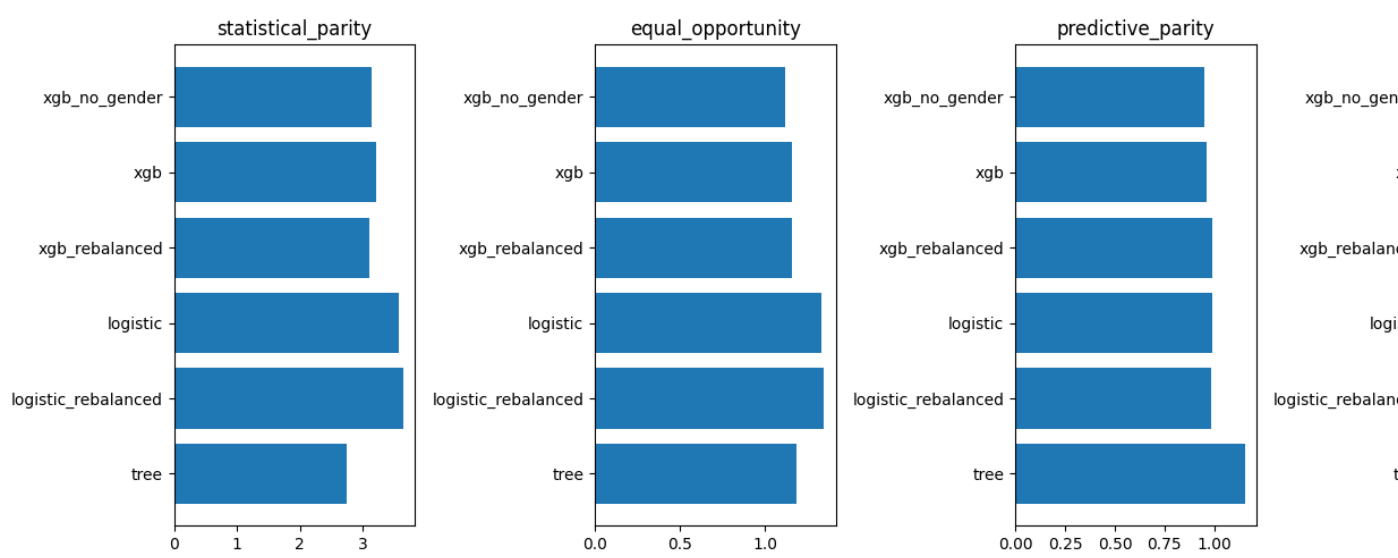
The equal opportunity coefficient:

$$\frac{TPR_b}{TPR_r} = \frac{3}{2} \text{ — blue is more privileged}$$

Demographic parity coefficient:

$$\frac{R_b}{R_r} = \frac{13}{10} \text{ — blue is more privileged}$$

I have trained logistic regression, decision tree and gradient boosted tree on income dataset. The protected attribute that I am working with is gender. I tested two strategies for mitigating bias: (i) I removed the protected column and (ii) I have tried rebalancing data. Results are presented in the following bar charts.



We can see that there seems to be negative correlation between accuracy and predictive parity. We can see that decision tree, which has worst accuracy, has best statistical parity. XGB with no gender column performs slightly better than XGB. XGB with no gender column performs slightly better than XGB in terms of fairness metrics, suggesting a slight overfitting to gender.

We can see that mitigating bias, by resampling the dataset to include equal number of men and women does not help at all.

	logistic	logistic_rebalanced	tree	xgb	xgb_no_gender
statistical_parity	3.571922	3.645472	2.741528	3.218254	3.129965
equal_opportunity	1.328129	1.346642	1.186238	1.158213	1.118033
predictive_parity	0.999806	0.999806	1.152105	0.950087	0.951020