

Table S2. Fairness Metrics

<u>Metric Name</u>	<u>General Definition</u>	<u>Example Statement</u>	<u>Equality Statement</u>	<u>Formula</u>
Independence (Demographic Parity)	A classifier (M) satisfies the independence fairness criterion, if its output (\hat{Y}) is independent of the sensitive attributes A .	The model's prediction of whether a glaucoma patient will progress to surgery is independent of that patient's sensitive characteristics; the proportion of those in the positive predicted class is the same across sensitive characteristics.	$P\{\hat{Y} = 1 A = a_1\} = P\{\hat{Y} = 1 A = a_2\}$	$ \text{Support}(M, G1) - \text{Support}(M, G2) $
Equalized Odds	A classifier (M) satisfies the equalized odds definition of fairness, if its output (\hat{Y}) is conditionally (with respect to the target output Y) independent of the sensitive attribute A .	The model's prediction of whether a glaucoma patient will progress to surgery performs with equal false positive and true positive rates for White and nonWhite patients.	$P\{\hat{Y} = 1 Y = 1, A = a_1\} = P\{\hat{Y} = 1 Y = 1, A = a_2\}$ $P\{\hat{Y} = 1 Y = 0, A = a_1\} = P\{\hat{Y} = 1 Y = 0, A = a_2\}$	$ \text{Recall}(M, G1) - \text{Recall}(M, G2) + \text{FPR}(M, G1) - \text{FPR}(M, G2) $
Overall Accuracy Equality	A classifier (M) satisfies the overall accuracy equality, if a learning system's accuracy values for both groups are equal.	The model's prediction of whether a glaucoma patient will progress to surgery is equally accurate for White and Non-White patients	$P\{\hat{Y} = Y A = a_1\} = P\{\hat{Y} = Y A = a_2\}$	$ \text{Accuracy}(M, G1) - \text{Accuracy}(M, G2) $
Sufficiency (Calibration)	A classifier (M) satisfies the sufficiency condition if the target output Y is conditionally (with respect to \hat{Y}) independent of the sensitive attribute A .	For any given predicted probability score, the actual probability of whether a glaucoma patient will progress to surgery should be equal for both White and Non-White patients.	$P\{Y = 1 \hat{Y} = \hat{y}, A = a_1\} = P\{Y = 1 \hat{Y} = \hat{y}, A = a_2\}$	$ \text{Precision}(M, G1) - \text{Precision}(M, G2) + \text{FOR}(M, G1) - \text{FOR}(M, G2) $

Let A_1, \dots, A_M be the sensitive attributes and X_1, \dots, X_m be the other (insensitive attributes). Y is the target (desired) output and \hat{Y} is the classifier output. Furthermore, assume A has two possible values a_1 and a_2 . Finally, we define $A = a_1$ as G1 and $A = a_2$ as G2.