

What is Fairness?

A model is fair if its Performance is equally good for all!

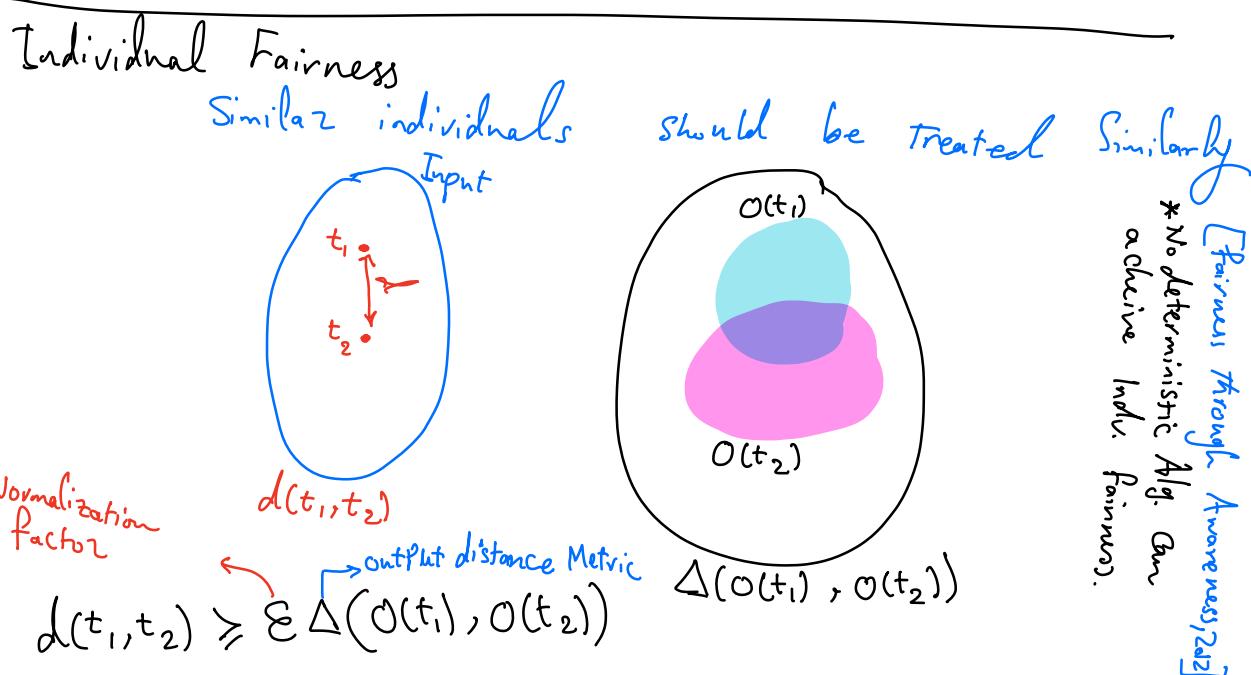
all = group memberships

$$\text{Perf}^A \underset{\sim}{\underset{\sim}{\underset{\sim}{=}}} \text{Perf}^A_{\text{group level}} \underset{\sim}{\underset{\sim}{\underset{\sim}{=}}} \text{Perf}^A_{\text{Subgroup level}} \underset{\sim}{\underset{\sim}{\underset{\sim}{=}}} \text{Perf}^A_{\text{Indiv. level}}$$

all

different groups.

Perf^A: Performance of Alg. A



Group Fairness

- We assume the existence of at least one Sensitive Attribute S (gender, race, Education, Income,...) the value of S specify the groups

GPA	SAT	gender	race
—	—	M	W
—	—	F	B
—	—	F	B
—	—	M	W
—	—

$\{M, F, W, B, \dots\}$ \leftarrow Non-Intersectional Groups

$\{WM, WF, BM, BF, \dots\}$ \leftarrow Intersectional

$\{WM, WM - \text{with College Degree}, \dots\}$

Fully Specified
Subgroup

Individual-level



$$|M| \geq |WM| \geq |WM-CD|$$

\hookrightarrow College degree

Subgroup Fairness Challenge:

→ we may not have enough samples from all subgroups to draw

Statistical Conclusions

demographic groups

Binary groups:

M vs F

B vs W

Majority vs Minority

Non-binary

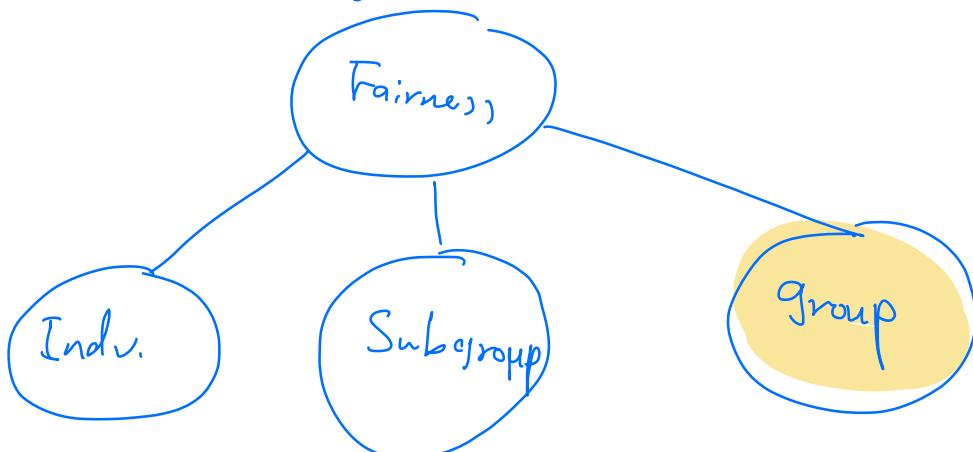
Overlapping

{F, H, B, M, ...}

Non-overlapping {B, W, H, ...}

Fairness

We Assume Non-overlapping Binary groups
for defining fairness Metrics



group Fairness (Binary Groups / Binary Classification)

$$- P(\alpha | \beta, g_1) = P(\alpha | \beta, g_2)$$

? Performance

$$- P(\alpha | \beta, g_1) = P(\alpha | \beta)$$

Decision Marker: Out of the ones labeled as Positive (+), how many are really Positive

Performance Metric

Predicted $f_\theta(x)$

		+	-
+	TP	FP	
	FN	TN	

$$\left[\frac{TP}{TP+FP} \right]_{g_2} = \left[\frac{TP}{TP+FP} \right]_{g_1}$$

y
True Label

$$P(y=1 | f=1, g_1) = P(y=1 | f=1, g_2)$$

Positive Predictive Parity (Aka Calibration)

$$- P(Y=0 | f=0, g_1) = P(Y=0 | f=0, g_2)$$

↳ Negative Predictive Parity

$$- P(Y=0 | f=1, g_1) = \dots$$

$$- P(Y=1 | f=0, g_1) = \dots$$

PP Rate Parity \rightarrow Equal
AND
NP Rate "

Decision Marker: How Likely Can an Indv. be

falsely labeled as positive.

e.g., If someone is not dangerous but mistakenly labeled as one.

Prediction $f_g(x)$	+	TP	-	FP
	-	FN	TN	
y True Label	+			
	-			

$$\left[\frac{FP}{FP+TN} \right]_{g_1} = \left[\frac{FP}{FP+TN} \right]_{g_2}$$

$$- P(f=1 | Y=0, g_1) = P(f=1 | Y=0, g_2)$$

False Positive Rate Parity

- False Neg. rate Parity
- True ~ ~ ~
- " Positive ~ ~

| Both
FPRate Parity
AND
FNRate Parity
↓
Equalized ODDS

Defendant: The likelihood of Positive Prediction
Should be equal for all groups

→ R-Scores are not fair because black
is more likely of being Predicted
as Positive (Dangerous)

	+	-
f	TP	FP
-	FN	TN

$$\left[\frac{TP + FP}{(TP + FP) + (TN + FN)} \right]_{\text{all}} = \left[\frac{TP + FP}{\text{all}} \right]_y$$

$$P(f=1 | g_1) = P(f=1 | g_2)$$

β : Null

↳ Demographic Parity
Statistical
↳ Disparate Impact

$$P(f=0 | g_1) = P(f=0 | g_2)$$

Decision Marker is of Accuracy
 { missclassification rate } is equal for the two groups

f	TP	FP
y	FN	TN

Accuracy Parity

$$\left[\frac{TP + TN}{all} \right]_{g_1} = \left[\frac{TP + TN}{all} \right]_{g_2}$$

$$P(y=f | g_1) = P(y=f | g_2)$$

Misclassification Rate (Error) Parity

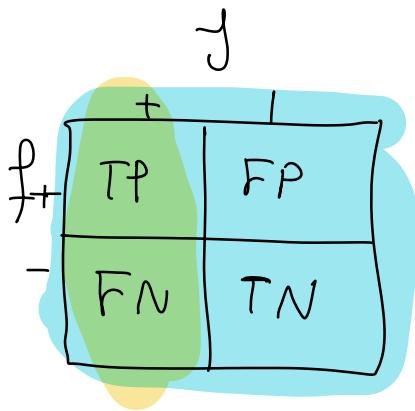
$$\left[\frac{FP + FN}{all} \right]_{g_1} = \left[\frac{FP + FN}{all} \right]_{g_2}$$

$$P(y \neq f | g_1) = P(y \neq f | g_2)$$

$$4 + 4 + 2 + 2 + 2 + 2 + 4 + 4 = 22 > 21$$

↗ rows ↗ Col ↗ Dem. Parity ↗ EEO ↗ EEO ↗ acc. (all) ↗ 1 cell ↗ D. ↗ miss

} 21 definitions of fairness



Assumption $P(y=1 | g_1) \neq P(y=1 | g_2)$
 ↳ Bias in Data.

$y \not\perp\!\!\!\perp s$

↳ The random variable y is not independent from the rand. variable s .

Simpson Paradox

admitted : $y=1$
 not admitted : $y=0$

$S=1 \quad M$

$S=0 \quad F$

UC Berkeley Admission process is sexist because

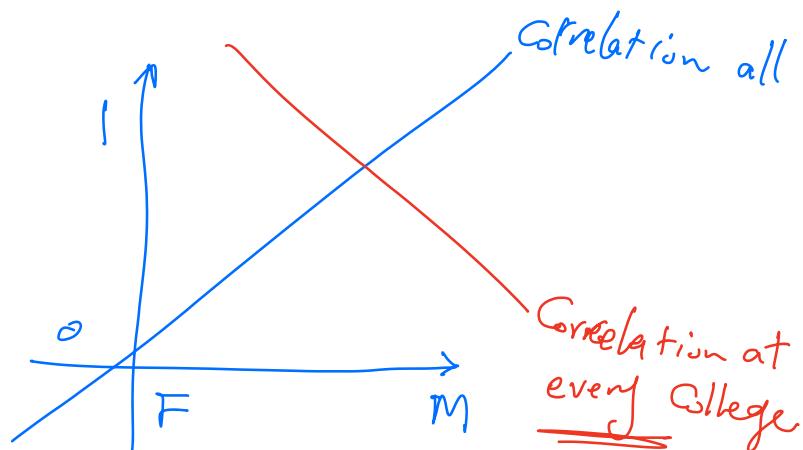
Overall : $P(f=1 | M) > P(f=1 | F)$

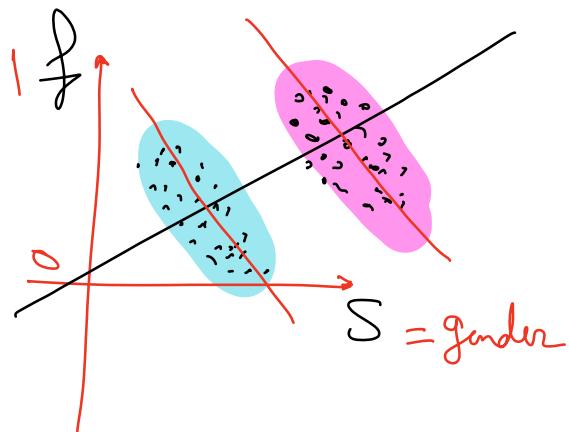
Demographic Parity

$$f \not\perp S \Rightarrow \text{Corr}(f, S) \neq 0$$

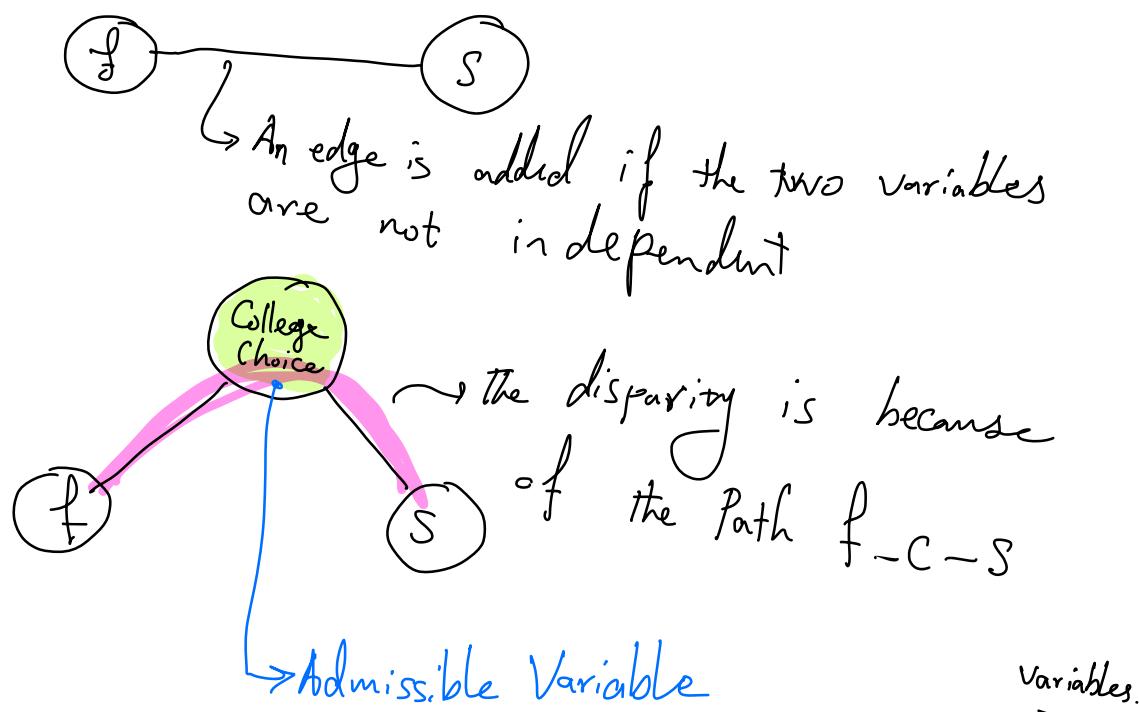
College-level : ~~almost always~~ Assume \forall College

$$P(f=1 | M) < P(f=1 | F)$$



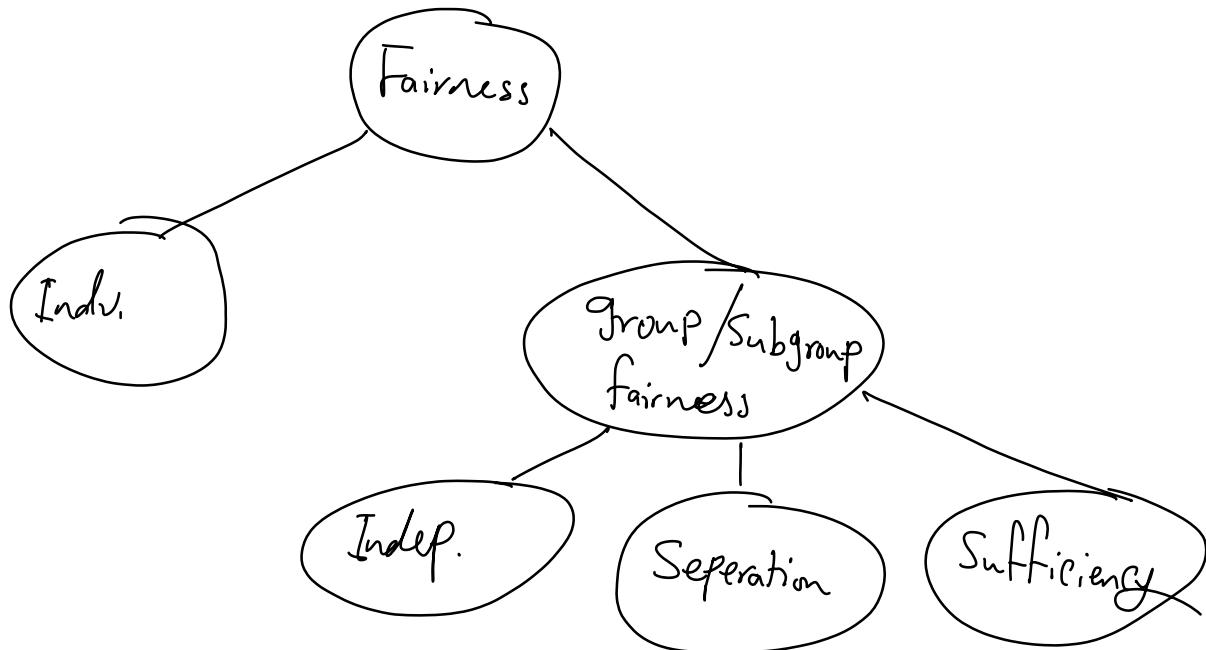


The admission process is fair because there is an admissible explanation for the overall disparity.



A disparity is unfairness iff it is not through admissible variables.

A Categorization of fairness Def.



Independence:

M satisfies independence

$S \perp\!\!\!\perp f$

e.g., demographic Parity

$I_{\text{indep.}} \Rightarrow \text{demo. Parity}$

Sufficiency:

S should be indep from f conditioned
on y

$S \perp\!\!\!\perp f \mid y$

\rightarrow $\begin{cases} TPR_f \\ FPR_P \\ TNRP \\ FNRP \\ \text{Equalized odds (EO)} \end{cases}$

Sufficiency \Rightarrow $\begin{cases} TPR_P \\ \dots \\ EO \end{cases}$

Separation:

S should be indep. from y end.

On f

$S \perp\!\!\!\perp y$

\rightarrow Pos. Predictive Parity (Calibration)

Neg. $\perp\!\!\!\perp$

Equal opportunity

Separation \Rightarrow PPPR

If data is biased: $y \not\perp s$, Assumption

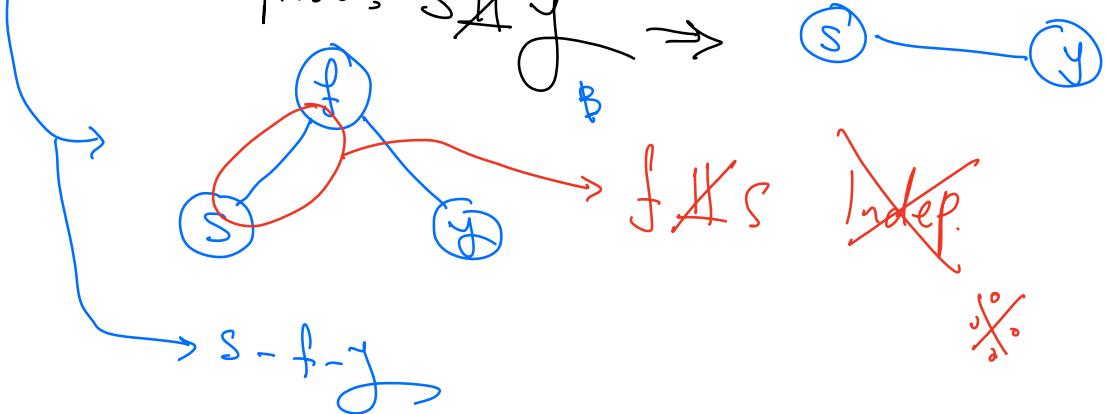
for any choice of two fairness categories, it is IMPOSSIBLE to satisfy BOTH

This: Indep. and Separation cannot be satisfied at the same time.

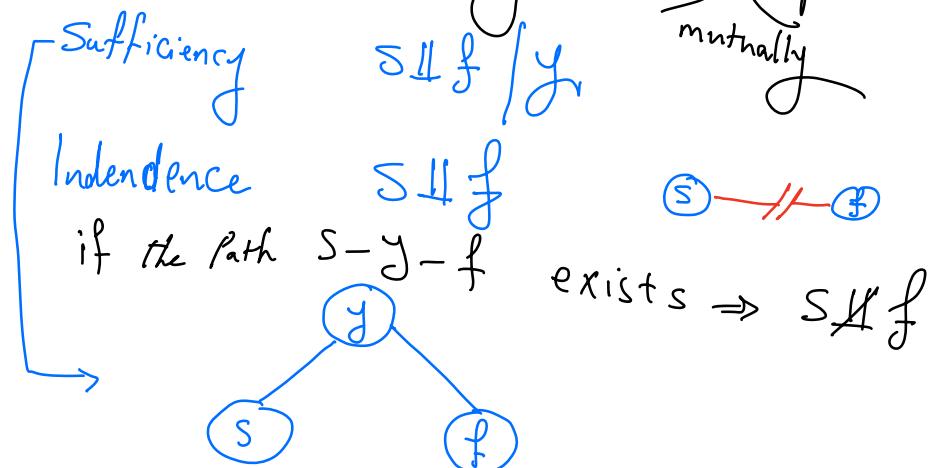
Indep. $s \perp\!\!\!\perp f$

Separation: $s \perp\!\!\!\perp y \mid f$

Assumption: $s \not\perp\!\!\!\perp y \mid f$



Independence & Sufficiency are impossible



\Rightarrow The Path $S - y - f$ does not exist.

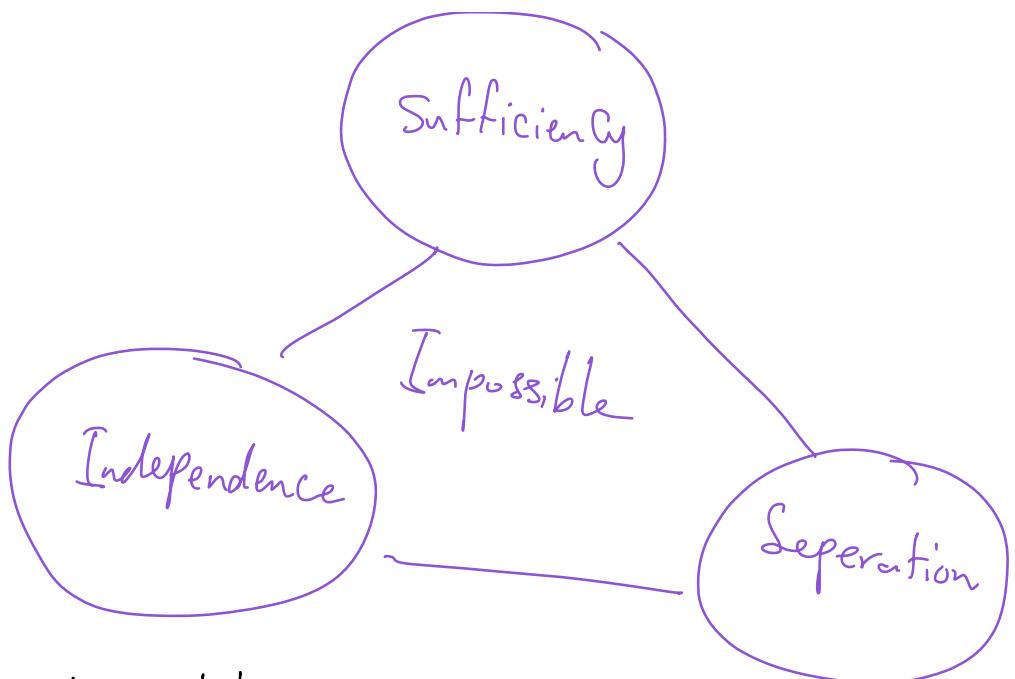
- $\Rightarrow y \perp\!\!\!\perp f \Leftarrow$ meaningless (No Prediction Power)
- $\Rightarrow S \perp\!\!\!\perp y \times \Leftarrow$ Bias Assumption

- Separation \wedge Sufficiency are mutually impossible

- $S \perp\!\!\!\perp y | f$ Separation
- $S \perp\!\!\!\perp f | y$ Sufficiency

$$S \perp\!\!\!\perp (y, f)$$

$\Rightarrow S \perp\!\!\!\perp y \times \Leftarrow$ Bias Assumption

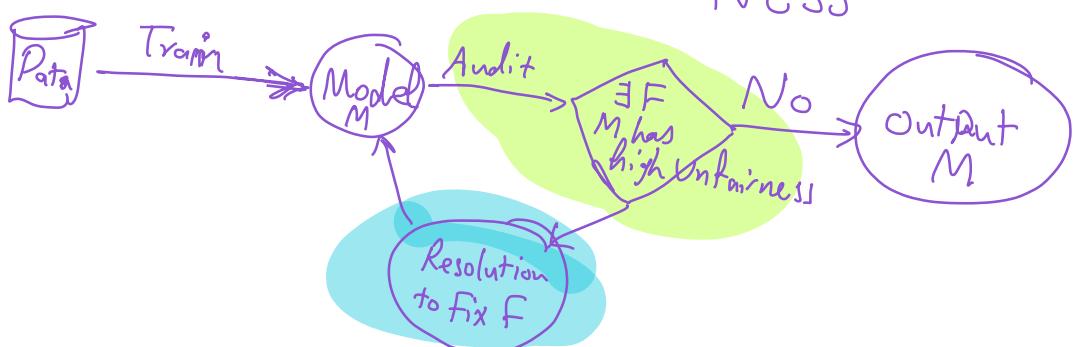


\Rightarrow Impossibility theorems only prove that exact equally is impossible for mutually exclusive definitions

we can still satisfy ALMOST EQUAL
for all definitions.

the goal should be to

MIN UNFAIRNESS



Measuring Unfairness

- If Alg. A has Perf. P_1 for g_1 ,
& P_2 for g_2 , how should we
measure the unfairness?)

$$\times - |P_2 - P_1|$$

e.g., - Unfairness, $F_1 = 0.3$

Case 2 $\frac{F_2 = 23}{}$ Different Scenarios

Case 3 $\frac{F_3 = 0.0003}{}$

$$\Rightarrow \text{Case 2: } \begin{cases} P_1 = 2368723 \\ P_2 = 2368700 \end{cases} \xrightarrow{\text{less fair}}$$

$$\Rightarrow \text{Case 3: } \begin{cases} P_1 = 0.0004 \\ P_2 = 0.0001 \end{cases} \xrightarrow{\text{less unfair}}$$

- Subtraction (Normalized)

$$F = \left| \frac{P_1 - P_2}{P_1 + P_2} \right| \rightarrow F = 0.7$$

$$- F = \frac{\min(P_1, P_2)}{\max(P_1, P_2)}$$

$$\frac{1}{1+\epsilon} < \frac{P_1}{P_2} < 1+\epsilon$$