

Definitions of Fairness,
Inherent Tradeoffs +
Impossibilities

Today

1. Definitions of fairness
 - statistical parity
 - predictive rate parity
 - equalized odds
2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)
3. Tradeoff between Fairness + Accuracy
4. Tradeoffs between simplicity and Fairness
5. Other

Running Example

COMPAS: risk assessment program

ProPublica concluded that COMPAS is biased:

The likelihood of blacks predicted to reconvict given that they did NOT is $>$ likelihood for whites

	Humans		COMPAS	
	Black %	White %	Black %	White %
Accuracy*	68.2	67.6	64.9	65.7
False Positives	37.1	27.2	40.4	25.4
False Negatives	29.2	40.3	30.9	47.9

COMPAS DATA

- Recidivism rate for blacks 51%^{do}
- Recidivism rate for whites 39%^{do}

IS COMPAS BIASED?

Today

1. Definitions of fairness

statistical parity
predictive rate parity
equalized odds

2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)
3. Tradeoff between Fairness + Accuracy
4. Tradeoffs between simplicity and Fairness
5. Other

Definitions / Notation

$x \in \mathcal{U}$ feature vector [Typically $\mathcal{U} = \mathbb{R}^d$ or discretized version]

$y \in \{0, 1\}$ actual value (we are trying to predict)

Underlying distribution is pair of r.v.'s (X, Y)

Classifier: maps x to $\hat{y} = f(x)$.

COMPAS EXAMPLE:

x : feature vector of offender

$y=1$: offender did reoffend, $y=0$ did not

\hat{y} : prediction for x

Probability Space, RV's

$U (= \mathbb{R}^d)$: set of all possible feature vectors (finite descriptions of a person - d attributes)

$p(u)$: probability (over all individuals in population) of having feature vector u

$X = p(u)$: random variable

(X, Y) : random variable

Confusion Matrix

		$\hat{y} = 0$	$\hat{y} = 1$	
$y = 0$	TNR	$\Pr[\hat{y} = 0 y = 0]$	FPR	$\Pr[\hat{y} = 1 y = 0]$
$y = 1$	FNR	$\Pr[\hat{y} = 0 y = 1]$	TPR	$\Pr[\hat{y} = 1 y = 1]$

TNR: true negative rate
FNR: false negative rate

FPR: false positive rate
TPR: true positive rate

Definitions / Notation

$x \in \mathbb{R}^d$ feature vector (may include A)

$y \in \{0,1\}$ actual value (we are trying to predict)

Underlying distribution is pair of r.v.'s (X, Y)

Classifier: maps x to $\hat{y} = f(x)$.

Sensitive variable: $A \in \{0,1\}$

Joint distribution (X, Y, A, \hat{Y})

Example: x : vector about offender
COMPAS y : whether offender will recidivate ($y=1$)
 A : black ($A=1$) or white ($A=0$)
 $\hat{y} = f(x)$: predicted value of y

Confusion Matrix $A=1$

	$\hat{y} = 0$	$\hat{y} = 1$
$y = 0$		40.4
$y = 1$	30.9	

Confusion Matrix $A=0$

	$\hat{y} = 0$	$\hat{y} = 1$
$y = 0$		25.4
$y = 1$	47.9	

	Humans		COMPAS	
	Black %	White %	Black %	White %
Accuracy*	68.2	67.6	64.9	65.7
False Positives	37.1	27.2	40.4	25.4
False Negatives	29.2	40.3	30.9	47.9

Definitions / Notation

$$\mathbb{R} = f(x)$$

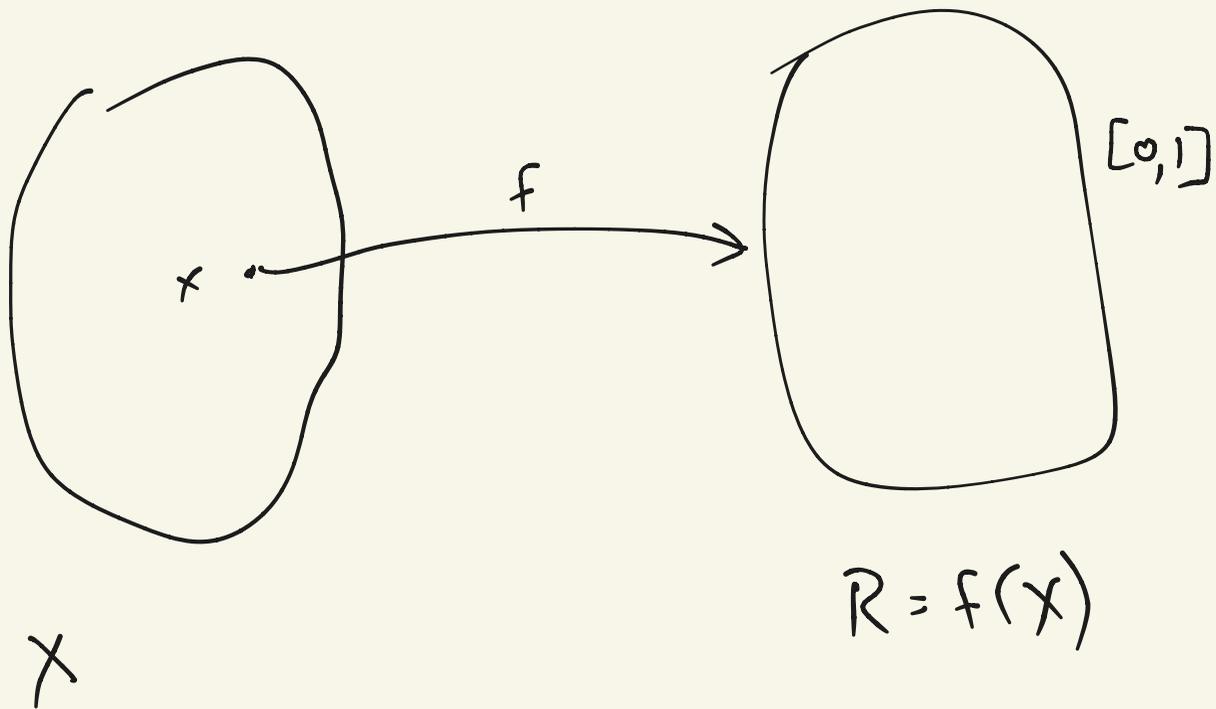
Classification often solved by first solving a regression problem to summarize data by a score, $f(x) \in \mathbb{R}$ (we assume $\in [0, 1]$)

Natural score function: $f(x) = \mathbb{E}[Y|x]$

Score to \hat{y} : Pick threshold t

$$\hat{y} = 1 \text{ iff } f(x) \geq t$$

\mathbb{R} generalizes \hat{y} - so from now on can think of \hat{y} as special case of \mathbb{R}



Definitions / Notation

1. Statistical Parity / group Parity / Independence

$$R \perp A \quad (\hat{Y} \perp A) \quad \Pr[R|A] = \Pr[R]$$

2. Predictive Rate Parity / Sufficiency

$$Y \perp A | R \quad Y \perp A | \hat{Y} \quad \Pr[Y|A, R] = \Pr[Y|R]$$

3. Equalized Odds / Separation

$$R \perp A | Y \quad \hat{Y} \perp A | Y \quad \Pr[R|A, Y] = \Pr[R|Y]$$

Definitions / Notation

1. Statistical Parity / group Parity

$$\forall b \in \{0, 1\} \quad \Pr[\hat{y} = b | A = 0] = \Pr[\hat{y} = b | A = 1]$$

2. Predictive Rate Parity

$$(*) \quad \left\{ \forall b, b' \in \{0, 1\} \quad \Pr[y = b | \hat{y} = b', A = 0] = \Pr[y = b | \hat{y} = b', A = 1] \right.$$



Equivalent to

$$(**) \quad \left\{ \begin{array}{l} \Pr[y = 1 | \hat{y} = 1, A = 0] = \Pr[y = 1 | \hat{y} = 1, A = 1] \text{ and} \\ \Pr[y = 0 | \hat{y} = 0, A = 0] = \Pr[y = 0 | \hat{y} = 0, A = 1] \end{array} \right.$$

Definitions / Notation

1. Statistical Parity / group Parity

$$\forall b \in \{0, 1\} \quad \Pr[\hat{y} = b | A = 0] = \Pr[\hat{y} = b | A = 1]$$

2. Predictive Rate Parity

$$\forall b, b' \in \{0, 1\} \quad \Pr[y = b | \hat{y} = b', A = 0] = \Pr[y = b | \hat{y} = b', A = 1]$$

3. Equalized Odds

$$\Pr[\hat{y} = 0 | y = 1, A = 0] = \Pr[\hat{y} = 0 | y = 1, A = 1]$$

$$\Pr[\hat{y} = 1 | y = 0, A = 0] = \Pr[\hat{y} = 1 | y = 0, A = 1]$$

Definitions / Notation

1. Statistical Parity / group Parity

$$\forall b \in \{0, 1\} \quad \Pr[\hat{y} = b | A = 0] = \Pr[\hat{y} = b | A = 1]$$

2. Predictive Rate Parity (PRP)

$$\forall b, b' \in \{0, 1\} \quad \Pr[y = b | \hat{y} = b', A = 0] = \Pr[y = b | \hat{y} = b', A = 1]$$

3. Equalized Odds

$$\Pr[\hat{y} = 0 | y = 1, A = 0] = \Pr[\hat{y} = 0 | y = 1, A = 1]$$

$$\Pr[\hat{y} = 1 | y = 0, A = 0] = \Pr[\hat{y} = 1 | y = 0, A = 1]$$

FNR

FPR

Today

1. Definitions of fairness

statistical parity
predictive rate parity
equalized odds

2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)

3. Tradeoff between Fairness + Accuracy

4. Tradeoffs between simplicity and Fairness

5. Other

Impossibility Theorem

any 2 of the 3 definitions of fairness
are mutually exclusive (except in degenerate cases)

Impossibility Theorem (Indep vs sufficiency)

any 2 of the 3 definitions of fairness
are mutually exclusive (except in degenerate cases)

① - ② Statistical parity & predictive rate parity
are mutually exclusive unless $A \perp Y$

$$A \perp \hat{Y} \text{ and } A \perp Y | \hat{Y} \Rightarrow A \perp Y$$

Impossibility Theorem

any 2 of the 3 definitions of fairness are mutually exclusive (except in degenerate cases)

① - ② Statistical parity & predictive rate parity are mutually exclusive unless $A \perp Y$

$$A \perp \hat{Y} \text{ and } A \perp Y | \hat{Y} \Rightarrow A \perp (Y, \hat{Y}) \Rightarrow A \perp Y$$

Pf $A \perp Y | \hat{Y} : \Pr[A | Y, \hat{Y}] = \Pr[A | \hat{Y}]$

$$A \perp \hat{Y} : \Pr[A | \hat{Y}] = \Pr[A]$$

$$\text{So } \Pr[A | Y] = \sum_b \Pr[\hat{Y}=b] \Pr[A | Y, \hat{Y}=b] = \sum_b \Pr(\hat{Y}=b) \Pr[A | \hat{Y}=b] \\ = \Pr(A)$$

Impossibility Theorem (Separation vs Sufficiency)

any 2 of the 3 definitions of fairness are mutually exclusive (except in degenerate cases)

②-③ Predictive Rate parity and Equalized odds are mutually exclusive unless $A \perp Y$

$$A \perp \hat{Y} | Y \text{ and } A \perp Y | \hat{Y} \Rightarrow A \perp (\hat{Y}, Y) \Rightarrow A \perp Y$$

$$\Pr(A | \hat{Y}, Y) = \Pr(A | \hat{Y}) \quad \text{and}$$

$$\Pr(A | \hat{Y}, Y) = \Pr(A | Y)$$

$$\Pr(A | \hat{Y}=1) = \Pr(Y=0) \Pr(A | \hat{Y}=1, Y=0) + \Pr(Y=1) \Pr(A | \hat{Y}=1, Y=1)$$

$$= \Pr(Y=0) \Pr(A | Y=0) + \Pr(Y=1) \Pr(A | Y=1)$$

$$= \Pr(A)$$

$$\Pr(A | \hat{Y}=1, Y=0) = \Pr(A | Y=0)$$

$$\Pr(A | \hat{Y}=1, Y=1) = \Pr(A | Y=1)$$

Impossibility Theorem (Indep vs Separation)

any 2 of the 3 definitions of fairness
are mutually exclusive (except in degenerate cases)

①-③ Statistical Parity + Equalized Odds (* For binary Y *)

are mutually exclusive unless $A \perp Y$
or $\hat{Y} \perp Y$

$$A \perp \hat{Y} \text{ and } A \perp \hat{Y} | Y \Rightarrow A \perp Y \text{ or } \hat{Y} \perp Y$$

Impossibility Theorem (Indep vs Separation)

any 2 of the 3 definitions of fairness are mutually exclusive (except in degenerate cases)

①-③ Statistical Parity + Equalized Odds (* For binary Y *)

are mutually exclusive unless $A \perp Y$
or $\hat{Y} \perp Y$

$$A \perp \hat{Y} \text{ and } A \perp \hat{Y} | Y \Rightarrow A \perp Y \text{ or } \hat{Y} \perp Y$$

$$\begin{aligned} \Pr[\hat{Y}=b] &= \Pr[\hat{Y}=b | A=a] = \sum_Y \Pr[\hat{Y}=b | A=a, Y=y] \Pr[Y=y | A=a] \\ &= \sum_Y \Pr[\hat{Y}=b | Y=y] \Pr[Y=y | A=a] \end{aligned}$$

$$\Pr[\hat{Y}=b] = \sum_Y \Pr[\hat{Y}=b | Y=y] \Pr[Y=y]$$

So

$$\sum_Y \Pr[\hat{Y}=b|Y=Y] \Pr[Y=Y] = \sum_Y \Pr[\hat{Y}=b|Y=Y] \Pr[Y=Y|A=a]$$

b_Y P b_Y P_a

$$p b_0 + (1-p) b_1 = P_a b_0 + (1-P_a) b_1$$

$$p(b_0 - b_1) + b_1 = P_a(b_0 - b_1) + b_1$$

$$p(b_0 - b_1) = P_a(b_0 - b_1)$$

so either $\underbrace{b_0 = b_1}_{\hat{Y} \perp Y}$ or $\underbrace{P = P_a}_{Y \perp A}$

BACK TO COMPAS

ProPublica says:

Blacks face higher false positive rates
so violates equalized odds

Northpointe's defense:

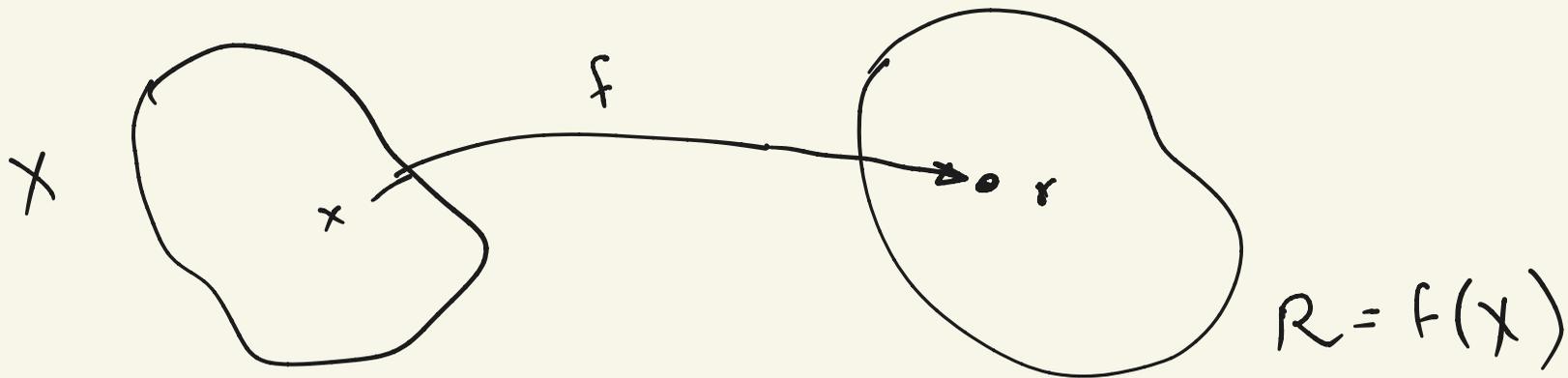
Scores satisfy predictive rate parity
and we can't have both.

CALIBRATION \approx PREDICTIVE RATE PARITY

Defn. A score R is calibrated if $\forall r$

$$\Pr[Y=1 | R=r] = r$$

*Note the natural score function $R(x) = \mathbb{E}[Y=1|x]$ is calibrated



CALIBRATION \approx PREDICTIVE RATE PARITY

Defn. A score R is calibrated by group if $\forall r$
$$\Pr[Y=1 | R=r, A] = r$$

* Note the natural score function $R(x) = \mathbb{E}[Y=1 | x, A]$
is calibrated by group

CALIBRATION \approx PREDICTIVE RATE PARITY

Defn. A score R is calibrated by group if $\forall r$
 $\Pr[Y=1 | R=r, A] = r$

Lemma

- ① R calibrated by group $\Rightarrow R$ satisfies predictive rate parity (PRP)
- ② R satisfies predictive rate parity $\Rightarrow \exists$ score function l st. $l(R)$ satisfies calibration by group

① Calibrated by group \Rightarrow PRP

$Y \perp R \mid A$

Assume $\forall r, a$

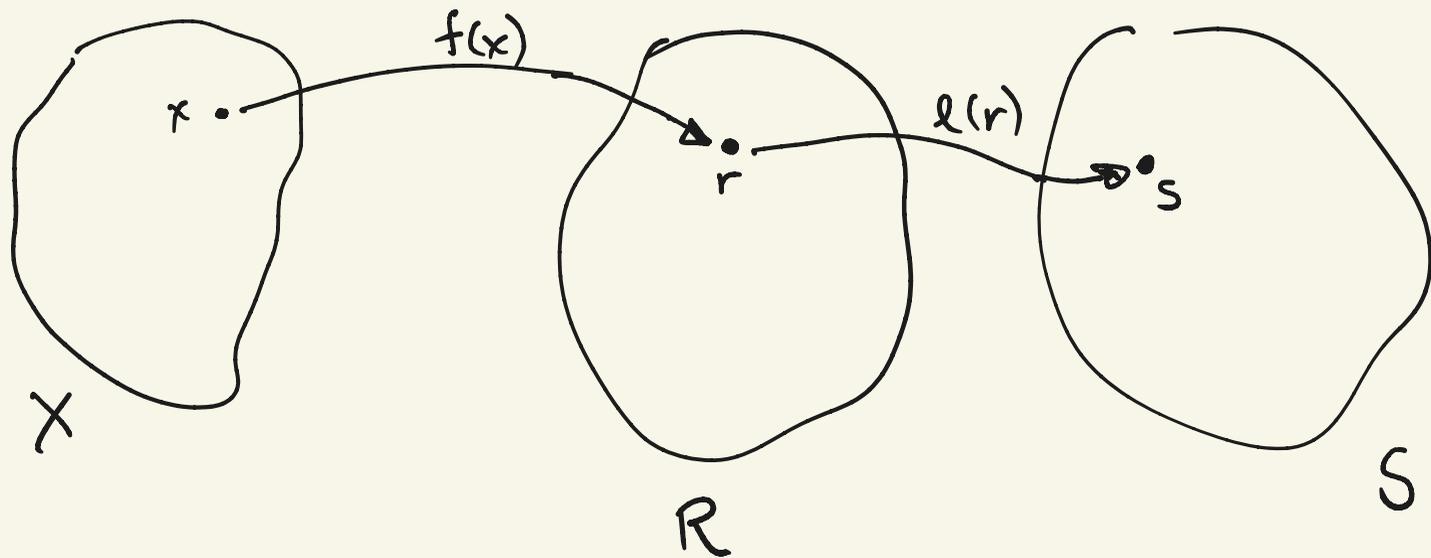
$$\Pr[Y=1 \mid R=r, A=a] = r$$

Then

$$\Pr[Y=1 \mid R, A] = \Pr[Y=1 \mid R] \quad \checkmark$$

② R satisfies PRP $\Rightarrow \exists \mathcal{L}$ s.t. $\mathcal{L}(R)$ satisfies calibration by group

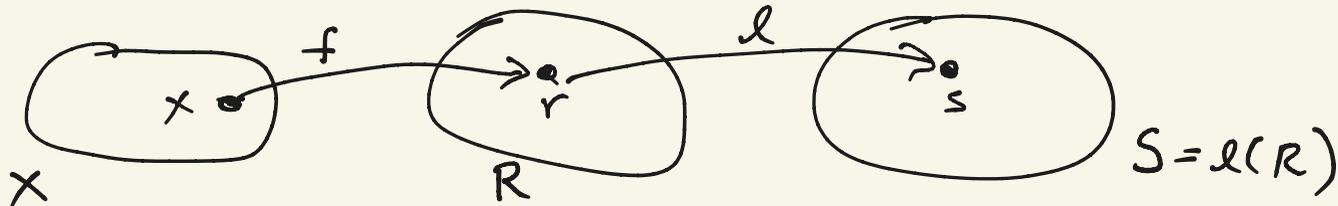
Define $\mathcal{L}(r) = \Pr[Y=1 | R=r, A=b]$



show $\Pr[Y=1 | S, A] = \Pr[Y=1 | S]$

② R satisfies PRP $\Rightarrow \exists \mathcal{L}$ s.t. $\mathcal{L}(R)$ satisfies calibration by group

Define $\mathcal{L}(r) = \Pr[Y=1 | R=r, A=b]$



$$\begin{aligned} & \Pr[Y=1 | S=s, A=a] \\ &= \Pr[Y=1 | r \in \mathcal{L}^{-1}(s), A=a] \\ &= \Pr[Y=1 | r \in \mathcal{L}^{-1}(s)] \\ &= \Pr[Y=1 | S=s] \end{aligned}$$

$\therefore \mathcal{L}(R)$ is calibrated by group.

(Semi-) Intuitive Proof (calibration vs equalized odds)

N_a = # people in group $A=a$

N_a^+ = # people in group $A=a$, with $\psi=1$

N_a^- = " " " " " " $\psi=0$

R_a = total score (sum of scores) for people in group $A=a$

R_a^+ = " " " " " " $A=a$, and $\psi=1$

R_a^- = " " " " " " $A=a$, and $\psi=0$

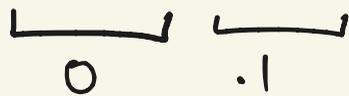
\bar{R}_a^+ = avg score of people in $A=a, \psi=1 = R_a^+ / N_a^+$

$\bar{R}_a^- = R_a^- / N_a^-$

$$R_a = N_a^+ \bar{R}_a^+ + N_a^- \bar{R}_a^-$$

Calibration : $R_a = N_a^+$

all
scores

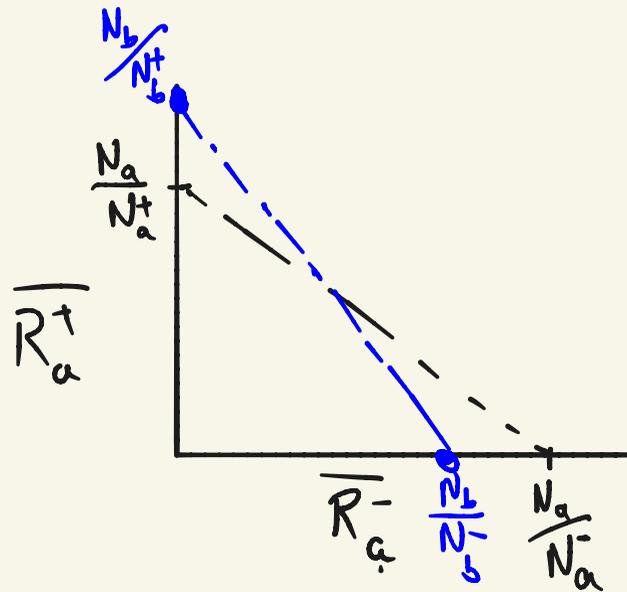


total score for bin $r = (\# \text{ people in bin } r) r$
 $= (N_{a,r}) r = N_{a,r}^+$

so

$$N_a^+ = \sum_r N_{a,r}^+ = R_a$$

$$\text{so } R_a = N_a^+ \underbrace{\bar{R}_a^+}_y + N_a^- \underbrace{\bar{R}_a^-}_x = N_a^+$$



Today

1. Definitions of fairness

statistical parity
predictive rate parity
equalized odds

2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)

3. Tradeoff between Fairness + Accuracy

4. Tradeoffs between simplicity and Fairness

5. Other

Tradeoff between Fairness + Accuracy

Example Suppose $y=1$ iff $A=1$

Then accuracy obviously at odds with fairness

"Inherent Tradeoffs in Learning Fair Representations"
[Zhao, Gordon]

- quantitative tradeoffs between statistical parity + accuracy via distance between distributions $\mathcal{D}_{A=0}(y)$ and $\mathcal{D}_{A=1}(y)$
- I.e., \hat{y} satisfies stat. parity,
error $\geq d_{TV}(\mathcal{D}_0(y), \mathcal{D}_1(y))$

Today

1. Definitions of fairness

statistical parity
predictive rate parity
equalized odds

2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)

3. Tradeoff between Fairness + Accuracy

4. Tradeoffs between simplicity and Fairness

5. Other

Simplicity/Fairness Tradeoffs

(Kleinberg, Mullainathan, 2019)

Setup:

- Set of applicants, can accept an r fraction
- $x \in \mathbb{R}^k$, $(k+1)^{\text{st}}$ dimension x_{k+1} is membership in A
- $S(x)$ gives a score to x
- s is **simple** if it doesn't depend on A
- top r percent (based on score) are admitted

Simplicity/Fairness Tradeoffs

Conditions:

① DISADVANTAGE condition on s

Let $\mu(x, A=b)$ be fraction of population with value (x, b)

For all x such that $s(x) > s(x')$

$$\frac{\mu(x, A=0)}{\mu(x, A=1)} > \frac{\mu(x', A=0)}{\mu(x', A=1)}$$

Simplicity/Fairness Tradeoffs

Conditions:

① DISADVANTAGE condition on s

Let $\mu(x, A=b)$ be fraction of population with value (x, b)

For all x such that $s(x) > s(x')$

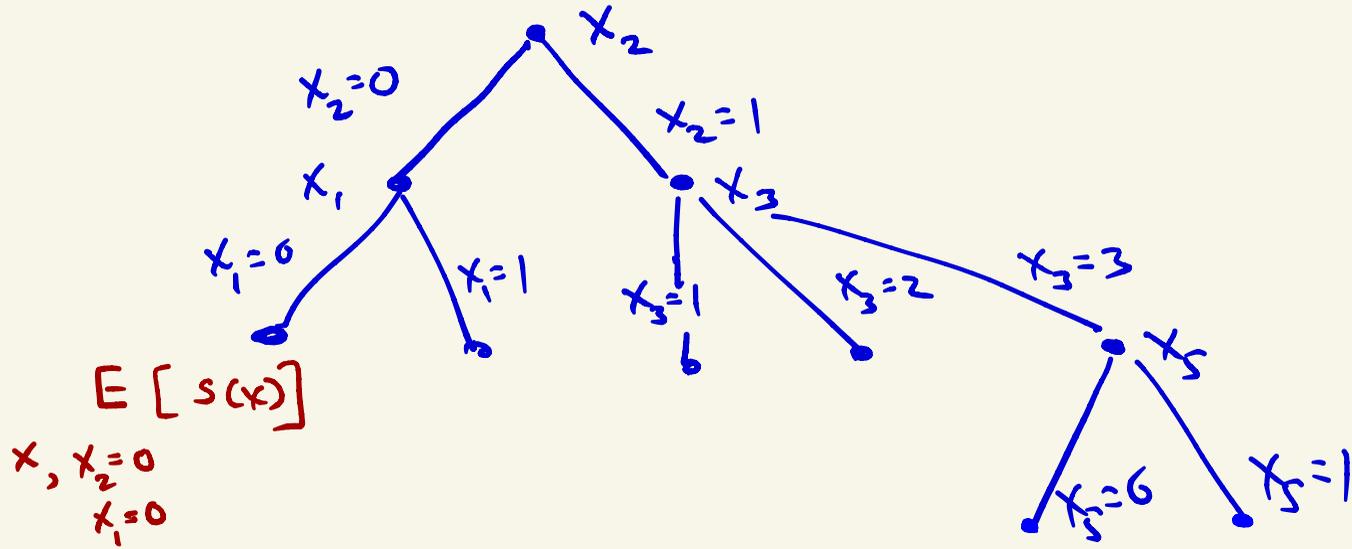
$$\frac{\mu(x, A=0)}{\mu(x, A=1)} > \frac{\mu(x', A=0)}{\mu(x', A=1)}$$

② GENERICITY condition on s for S, T subsets of applicants, $E[s(S)] \neq E[s(T)]$

↖ ensures no further simplification of score is possible

Simplicity/Fairness Tradeoffs

Simple S -approximators: Decision trees



For a path p in decision tree with partial assignment p , label leaf of p with $E[s(x)]$
 x, x consistent with p

Simplicity/Fairness Tradeoffs

Simple S -approximators: Decision trees

A decision tree f approximates S as follows:

order subcubes highest to lowest (by leaf value)

and output individuals in this order

until we reach rate r

Efficiency of f , $V_f(r)$: avg value of S
for the admitted people

Equity of f , $W_f(r)$: fraction of admitted
people who belong to
 $A=1$ (disadvantaged group)

Simplicity/Fairness Tradeoffs

Theorem 1 Let S satisfy disadvantage + genericity conditions. Then every simple S -approximator is strictly improvable:

For every nontrivial simple approximator g to S , there is a refinement h of g that is better:

$$\forall r \quad V_g(r) \leq V_h(r), \quad W_g(r) \leq W_h(r)$$

and

$$\exists r^* \text{ st } V_g(r^*) < V_h(r^*), \quad W_g(r^*) < W_h(r^*)$$

Simplicity/Fairness Tradeoffs

Theorem 2 Say that an s -approximator f (not nec. simple) is "group-agnostic" if $f(x, A=0) = f(x, A=1) \forall x$

Let f be a group agnostic approximator and let f' be the approximator to f obtained by splitting/refining each cell c_i of f according to group membership in A .

Then $V_{g'} > V_g$ but $W_{g'} < W_g$

ie. if we try to approx s by a group agnostic g , this incentivizes a rule that depends on A where value improves at expense of equity

Simplicity/Fairness Tradeoffs

Simple S -approximators: Decision trees

A decision tree f approximates S as follows:

order subcubes highest to lowest (by leaf value)

and output individuals in this order

until we reach rate r

Efficiency of f , $V_f(r)$: avg value of S
for the admitted people

Equity of f , $W_f(r)$: fraction of admitted
people who belong to
 $A = 1$ (disadvantaged group)

Today

1. Definitions of fairness
 - statistical parity
 - predictive rate parity
 - equalized odds
2. Impossibility results: any 2 of the 3 fairness conditions cannot be achieved (except in degenerate situations)
3. Tradeoff between Fairness + Accuracy
4. Tradeoffs between simplicity and Fairness
5. Other

Fairness Under Composition (Dwork, Ivento
ITCS 2019)