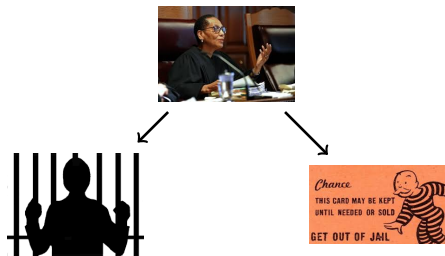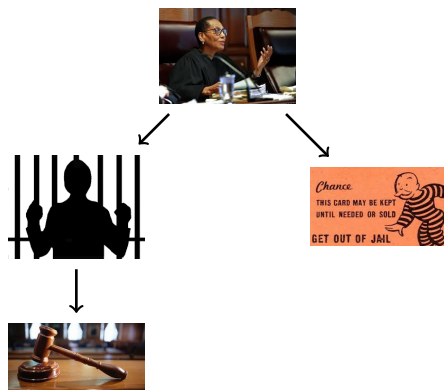# Fairness

Christos Dimitrakakis

September 27, 2018

# Bail decisions

# Bail decisions
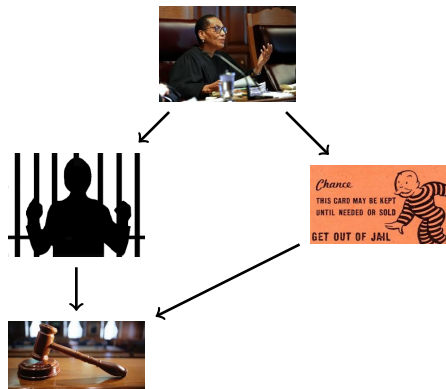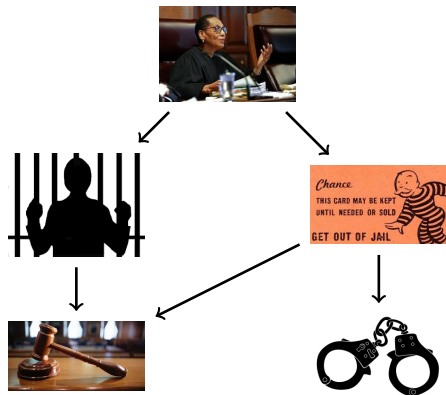
# Bail decisions

# Bail decisions

# Bail decisions

# Bail decisions

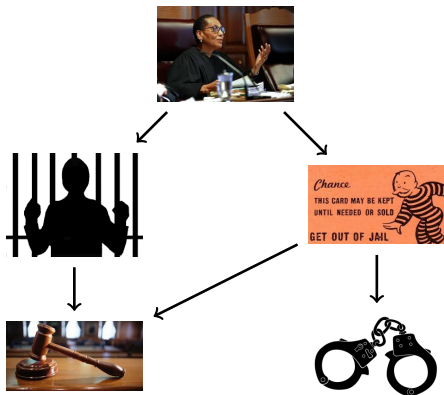# Bail decisions



His honour the machine
Prisoners released on bail*
%

Chosen by judges: 18.6

Suggested by algorithm: 14.9

of which: re-offend†

*From a representative sample of the US Department of Justice database 1990-2009
Source: Jens Ludwig, University of Chicago
†Failure to appear in court and re-arrest before trial

Economist.com

# Whites get lower scores than blacks[1]



Black                                    White

Figure: Apparent bias in risk scores towards black versus white defendants.

---

[1]Pro-publica, 2016

# But scores equally accurately predict recidivsm[2]



Figure: Recidivism rates by risk score.

# But non-offending blacks get higher scores



Figure: Score breakdown based on recidivism rates.

# Graphical models and independence

- ► Why is it not possible to be fair in all respects?
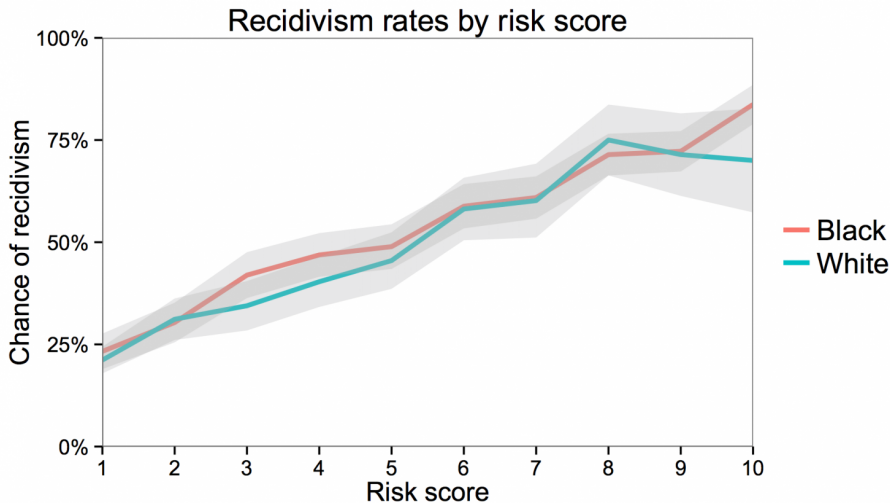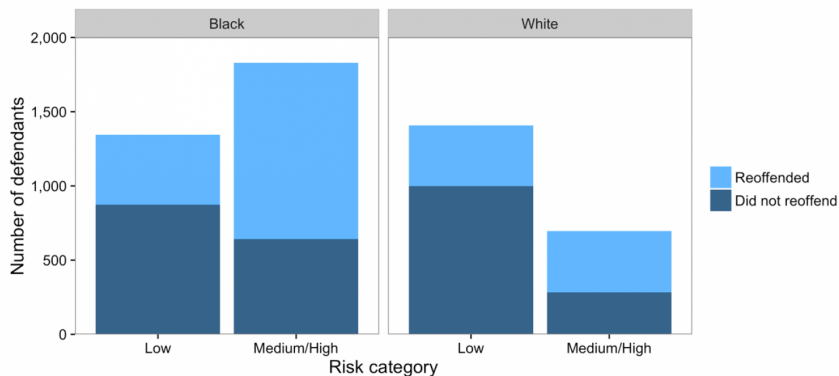- ► Different notions of conditional independence.
- ► Can only be satisfied rarely simultaneously.

# Graphical models



Figure: Graphical model (directed acyclic graph) for three variables.

## Joint probability

Let $\mathbf{x} = (x_1, \ldots, x_n)$. Then $\mathbf{x} : \Omega \to X$, $X = \prod_i X_i$ and:

$$\mathbb{P}(\mathbf{x} \in A) = P(\{\omega \in \Omega \mid \mathbf{x}(\omega) \in A\}).$$

## Factorisation

$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(\mathbf{x}_B \mid \mathbf{x}_C) \, \mathbb{P}(\mathbf{x}_C), \qquad B, C \subset [n]$$

C. Dimitrakakis                    Fairness                    September 27, 2018    7 / 30
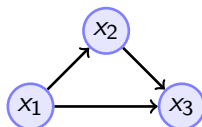
# Graphical models



Figure: Graphical model (directed acyclic graph) for three variables.

### Joint probability

Let $\mathbf{x} = (x_1, \ldots, x_n)$. Then $\mathbf{x} : \Omega \to X$, $X = \prod_i X_i$ and:

$$\mathbb{P}(\mathbf{x} \in A) = P(\{\omega \in \Omega \mid \mathbf{x}(\omega) \in A\}).$$

### Factorisation

So we can write any joint distribution as

$$\mathbb{P}(x_1) \, \mathbb{P}(x_2 \mid x_1) \, \mathbb{P}(x_3 \mid x_1, x_2) \cdots \mathbb{P}(x_n \mid x_1, \ldots, x_{n-1}).$$

# Directed graphical models



Figure: Graphical model for the factorisation $\mathbb{P}(x_3 \mid x_2)\,\mathbb{P}(x_2 \mid x_1)\,\mathbb{P}(x_1)$.

### Conditional independence

We say $x_i$ is conditionally independent of $\mathbf{x}_B$ given $\mathbf{x}_D$ and write
$x_i \mid \mathbf{x}_D \perp\!\!\!\perp \mathbf{x}_B$ iff

$$\mathbb{P}(x_i, \mathbf{x}_B \mid \mathbf{x}_D) = \mathbb{P}(x_i \mid \mathbf{x}_D)\,\mathbb{P}(\mathbf{x}_D \mid \mathbf{x}_B).$$

## Example 1 (Smoking and lung cancer)



Figure: Smoking and lung cancer graphical model, where $S$: Smoking, $C$: cancer, $A$: asbestos exposure.

## Explaining away

Even though $S, A$ are independent, they become dependent once you know $C$.

## Example 2 (Time of arrival at work)



Figure: Time of arrival at work graphical model where $T$ is a traffic jam and $x_1$ is the time John arrives at the office and $x_2$ is the time Jane arrives at the office.

### Conditional independence

Even though $x_1, x_2$ are correlated, they become independent once you know $T$.

## Example 3 (Treatment effects)



Figure: Kidney treatment model, where $x$: severity, $y$: result, $a$: treatment applied

|  | Treatment A | Treatment B |
|---|---|---|
| Small stones | 87 | 270 |
| Large stones | 263 | 80 |
| Severity | Treatment A | Treatment B |
| Small stones ) | 93% | 87% |
| Large stones | 73% | 69% |
| Average | 78% | 83% |

## Example 4 (School admission)



Figure: School admission graphical model, where $z$: gender, $s$: school applied to, $a$: whether you were admitted.

| School | Male | Female |
|--------|------|--------|
| A | 62% | 82% |
| B | 63% | 68% |
| C | 37% | 34% |
| D | 33% | 35% |
| E | 28% | 24% |
| F | 6% | 7% |
| *Average* | *45%* | *38%* |

## Exercise 1



Factorise the following graphical model.

## Exercise 1



Factorise the following graphical model.

$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1)\,\mathbb{P}(x_2 \mid x_1)\,\mathbb{P}(x_3 \mid x_1)\,\mathbb{P}(x_4)$$

Exercise 2



Factorise the following graphical model.

Exercise 2



Factorise the following graphical model.

$$Xb\, \mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1)\, \mathbb{P}(x_2 \mid x_1)\, \mathbb{P}(x_3 \mid x_1)\, \mathbb{P}(x_4 \mid x_3)$$

## Exercise 3

What dependencies does the following factorisation imply?

$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1)\,\mathbb{P}(x_2 \mid x_1)\,\mathbb{P}(x_3 \mid x_1)\,\mathbb{P}(x_4 \mid x_2, x_3)$$

## Exercise 3

What dependencies does the following factorisation imply?

$$\mathbb{P}(\mathbf{x}) = \mathbb{P}(x_1)\,\mathbb{P}(x_2 \mid x_1)\,\mathbb{P}(x_3 \mid x_1)\,\mathbb{P}(x_4 \mid x_2, x_3)$$

## Deciding conditional independence

There is an algorithm for deciding conditional independence of any two variables in a graphical model.

# Measuring independence

## Theorem 5
If $x_i \mid \mathbf{x}_D \perp\!\!\!\perp \mathbf{x}_B$ then

$$\mathbb{P}(x_i \mid \mathbf{x}_B, \mathbf{x}_D) = \mathbb{P}(x_i \mid \mathbf{x}_D)$$

## Example 6

$$\| \mathbb{P}(a \mid y, z) - \mathbb{P}(a \mid y) \|_1$$

which for discrete $a, y, z$ is:

$$\max_{i,j} \| \mathbb{P}(a \mid y = i, z = j) - \mathbb{P}(a \mid y = i) \|_1 = \max_{i,j} \| \sum_k \mathbb{P}(a = k \mid y = i, z = j) -$$

## Measuring independence

### Theorem 5
If $x_i \mid \mathbf{x}_D \perp\!\!\!\perp \mathbf{x}_B$ then

$$\mathbb{P}(x_i \mid \mathbf{x}_B, \mathbf{x}_D) = \mathbb{P}(x_i \mid \mathbf{x}_D)$$

This implies

$$\mathbb{P}(x_i \mid \mathbf{x}_B = b, \mathbf{x}_D) = \mathbb{P}(x_i \mid \mathbf{x}_B = b', \mathbf{x}_D)$$

so we can measure independence by seeing how the distribution of $x_i$ changes when we vary $\mathbf{x}_B$, keeping $\mathbf{x}_D$ fixed.

### Example 6

$$\| \mathbb{P}(a \mid y, z) - \mathbb{P}(a \mid y) \|_1$$

which for discrete $a, y, z$ is:

$$\max_{i,j} \| \mathbb{P}(a \mid y = i, z = j) - \mathbb{P}(a \mid y = i) \|_1 = \max_{i,j} \| \sum_k \mathbb{P}(a = k \mid y = i, z = j) -$$

# Coin tossing, revisited

## Example 7
The Beta-Bernoulli prior



Figure: Graphical model for a Beta-Bernoulli prior

$$\theta \sim \mathcal{Beta}(\xi_1, \xi_2), \quad \text{i.e. } \xi \text{ are Beta distribution parameters} \quad (2.1)$$
$$x \mid \theta \sim \mathcal{Bernoulli}(\theta), \quad \text{i.e. } P_\theta(x) \text{ is a Bernoulli} \quad (2.2)$$

## Example 8

An alternative model for coin-tossing This is an elaboration of Example **??** for hypothesis testing.



Figure: Graphical model for a hierarchical prior

- $\mu_1$: A Beta-Bernoulli model with $Beta(\xi_1, \xi_2)$
- $\mu_0$: The coin is fair.

$$\theta \mid \mu = \mu_0 \sim \mathcal{D}(0.5), \qquad \text{i.e. } \theta \text{ is always } 0.5 \qquad (2.3)$$

$$\theta \mid \mu = \mu_1 \sim Beta(\xi_1, \xi_2), \qquad \text{i.e. } \theta \text{ has a Beta distribution} \qquad (2.4)$$

$$x \mid \theta \sim Bernoulli(\theta), \qquad \text{i.e. } P_\theta(x) \text{ is Bernoulli} \qquad (2.5)$$

# Bayesian testing of independence



(a) $\Theta_0$ assumes independence



(b) $\Theta_1$ does not assume independence

### Example 9

Assume data $D = \{x_1^t, x_2^t, x_3^t \mid t = 1, \ldots, T\}$ with $x_i^t \in \{0, 1\}$.

$$P_\theta(D) = \prod_t P_\theta(x_3^t \mid x_2^t) P_\theta(x_2^t \mid x_1^t) P_\theta(x_1^t), \qquad \theta \in \Theta_0 \qquad (2.6)$$

$$P_\theta(D) = \prod_t P_\theta(x_3^t \mid x_2^t, x_1^t) P_\theta(x_2^t \mid x_1^t) P_\theta(x_1^t), \qquad \theta \in \Theta_1 \qquad (2.7)$$

# Bayesian testing of independence



(a) $\Theta_0$ assumes independence

(b) $\Theta_1$ does not assume independence

Example 9

$$\theta_1 \triangleq P_\theta(x_1^t = 1) \qquad\qquad (\mu_0, \mu_1)$$

$$\theta_{2|1}^i \triangleq P_\theta(x_2^t = 1 \mid x_1^t = i) \qquad\qquad (\mu_0, \mu_1)$$

$$\theta_{3|2}^j \triangleq P_\theta(x_3^t = 1 \mid x_2^t = j) \qquad\qquad (\mu_0)$$

$$\theta_{3|2,1}^{i,j} \triangleq P_\theta(x_3^t = 1 \mid x_2^t = j, x_1^t = i) \qquad\qquad (\mu_1)$$

C. Dimitrakakis                     Fairness                     September 27, 2018      20 / 30

Figure: Hierarchical model.

$$\mu_i \sim \phi \tag{2.6}$$
$$\theta \mid \mu = \mu_i \sim \xi_i \tag{2.7}$$

Marginal likelihood

$$\mathbb{P}_\phi(D) = \phi(\mu_0)\,\mathbb{P}_{\mu_0}(D) + \phi(\mu_1)\,\mathbb{P}_{\mu_1}(D) \tag{2.8}$$
$$\mathbb{P}_{\mu_i}(D) = \int_{\Theta_i} P_\theta(D)\,\mathrm{d}\xi_i(\theta). \tag{2.9}$$

Figure: Hierarchical model.

## Marginal likelihood

$$\mathbb{P}_\phi(D) = \phi(\mu_0)\,\mathbb{P}_{\mu_0}(D) + \phi(\mu_1)\,\mathbb{P}_{\mu_1}(D) \tag{2.6}$$

$$\mathbb{P}_{\mu_i}(D) = \int_{\Theta_i} P_\theta(D)\,\mathrm{d}\xi_i(\theta). \tag{2.7}$$

## Model posterior

$$\phi(\mu \mid D) = \frac{\mathbb{P}_\mu(D)\phi(\mu)}{\sum_i \mathbb{P}_{\mu_i}(D)\phi(\mu_i)} \tag{2.8}$$

# Calculating the marginal likelihood

Monte-Carlo approximation

$$\int_\Theta P_\theta(D) \, \mathrm{d}\xi(\theta) \approx \sum_{n=1}^{N} P_{\theta_n}(D) + O(1/\sqrt{N}), \qquad \theta_n \sim \xi \qquad (2.9)$$

Importance sampling

$$\int_\Theta P_\theta(D) \, \mathrm{d}\xi(\theta) \qquad\qquad\qquad\qquad (2.10)$$

# Calculating the marginal likelihood

## Monte-Carlo approximation

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\xi(\theta) \approx \sum_{n=1}^{N} P_{\theta_n}(D) + O(1/\sqrt{N}), \qquad \theta_n \sim \xi \qquad (2.9)$$

## Importance sampling

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\xi(\theta) = \int_{\Theta} P_{\theta}(D) \frac{\mathrm{d}\psi(\theta)}{\mathrm{d}\psi(\theta)} \, \mathrm{d}\xi(\theta)$$

$$(2.10)$$

# Calculating the marginal likelihood

## Monte-Carlo approximation

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\xi(\theta) \approx \sum_{n=1}^{N} P_{\theta_n}(D) + O(1/\sqrt{N}), \qquad \theta_n \sim \xi \qquad (2.9)$$

## Importance sampling

$$\int_{\Theta} P_{\theta}(D) \, \mathrm{d}\xi(\theta) = \int_{\Theta} P_{\theta}(D) \frac{\mathrm{d}\xi(\theta)}{\mathrm{d}\psi(\theta)} \, \mathrm{d}\psi(\theta)$$

$$(2.10)$$

# Calculating the marginal likelihood

Monte-Carlo approximation

$$\int_{\Theta} P_\theta(D) \, \mathrm{d}\xi(\theta) \approx \sum_{n=1}^{N} P_{\theta_n}(D) + O(1/\sqrt{N}), \qquad \theta_n \sim \xi \qquad (2.9)$$

Importance sampling

$$\int_{\Theta} P_\theta(D) \, \mathrm{d}\xi(\theta) \approx \sum_{n=1}^{N} P_\theta(D) \frac{\mathrm{d}\xi(\theta_n)}{\mathrm{d}\psi(\theta_n)}, \qquad \theta_n \sim \psi \qquad (2.10)$$

Sequential updating of the marginal likelihood

$$\mathbb{P}_\xi(D)$$

$$(2.14)$$

Example 10 (Beta-Bernoulli)

$$\mathbb{P}_\xi(x_t = 1 \mid x_1, \ldots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with $\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n, \quad \beta_t = \beta_0 + \sum_{n=1}^{t-1}(1 - x_n)$

Sequential updating of the marginal likelihood

$$\mathbb{P}_\xi(D) = \mathbb{P}_\xi(x_1, \ldots, x_T)$$

(2.14)

Example 10 (Beta-Bernoulli)

$$\mathbb{P}_\xi(x_t = 1 \mid x_1, \ldots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with $\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n, \quad \beta_t = \beta_0 + \sum_{n=1}^{t-1}(1 - x_n)$

## Sequential updating of the marginal likelihood

$$\begin{aligned}
\mathbb{P}_\xi(D) &= \mathbb{P}_\xi(x_1, \ldots, x_T) \\
&= \mathbb{P}_\xi(x_2, \ldots, x_T \mid x_1)\, \mathbb{P}_\xi(x_1)
\end{aligned} \tag{2.11}$$

$$\tag{2.14}$$

### Example 10 (Beta-Bernoulli)

$$\mathbb{P}_\xi(x_t = 1 \mid x_1, \ldots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with $\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n, \quad \beta_t = \beta_0 + \sum_{n=1}^{t-1}(1 - x_n)$

## Sequential updating of the marginal likelihood

$$\mathbb{P}_\xi(D) = \mathbb{P}_\xi(x_1, \ldots, x_T) \tag{2.11}$$

$$= \mathbb{P}_\xi(x_2, \ldots, x_T \mid x_1)\,\mathbb{P}_\xi(x_1) \tag{2.12}$$

$$= \prod_{t=1}^{T} \mathbb{P}_\xi(x_t \mid x_1, \ldots, x_{t-1})$$

$$\tag{2.14}$$

### Example 10 (Beta-Bernoulli)

$$\mathbb{P}_\xi(x_t = 1 \mid x_1, \ldots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with $\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n, \quad \beta_t = \beta_0 + \sum_{n=1}^{t-1} (1 - x_n)$

## Sequential updating of the marginal likelihood

$$\mathbb{P}_\xi(D) = \mathbb{P}_\xi(x_1, \ldots, x_T) \tag{2.11}$$

$$= \mathbb{P}_\xi(x_2, \ldots, x_T \mid x_1)\, \mathbb{P}_\xi(x_1) \tag{2.12}$$

$$= \prod_{t=1}^{T} \mathbb{P}_\xi(x_t \mid x_1, \ldots, x_{t-1}) \tag{2.13}$$

$$= \prod_{t=1}^{T} \int_\Theta P_{\theta_n}(x_t)\, \mathrm{d}\, \underbrace{\xi(\theta \mid x_1, \ldots, x_{t-1})}_{\text{posterior at time } t} \tag{2.14}$$

## Example 10 (Beta-Bernoulli)

$$\mathbb{P}_\xi(x_t = 1 \mid x_1, \ldots, x_{t-1}) = \frac{\alpha_t}{\alpha_t + \beta_t},$$

with $\alpha_t = \alpha_0 + \sum_{n=1}^{t-1} x_n, \quad \beta_t = \beta_0 + \sum_{n=1}^{t-1}(1 - x_n)$

# Further reading

### Python sources

- ▶ A simple python measure of conditional independence
  `src/fairness/ci_test.py`
- ▶ A simple test for discrete Bayesian network
  `src/fairness/DirichletTest.py`
- ▶ Using the PyMC package
  `https://docs.pymc.io/notebooks/Bayes_factor.html`

# Bail decisions, revisited

$x$



$\pi$

# Bail decisions, revisited

$x$



$\searrow \pi$



$a_1 \swarrow$



$\pi(a \mid x)$ \hspace{2cm} (policy)

# Bail decisions, revisited

$x$



$\pi$

$\pi(a \mid x)$ (policy)

$a_1$

$a_2$

# Bail decisions, revisited



$$\pi(a \mid x) \qquad \text{(policy)}$$

$$\mathbb{P}(y \mid a, x) \qquad \text{(outcome)}$$

# Bail decisions, revisited

$x$



$\pi$



$a_1$

$a_2$





$y_1$



$\pi(a \mid x)$      (policy)

$\mathbb{P}(y \mid a, x)$      (outcome)

# Bail decisions, revisited

$x$



$\pi$



$a_1$          $a_2$

        

$y_1$          $y_2$

        

$\pi(a \mid x)$      (policy)

$\mathbb{P}(y \mid a, x)$      (outcome)

# Bail decisions, revisited



$$\pi(a \mid x) \qquad \text{(policy)}$$

$$\mathbb{P}(y \mid a, x) \qquad \text{(outcome)}$$

$$U(a, y) \qquad \text{(utility)}$$
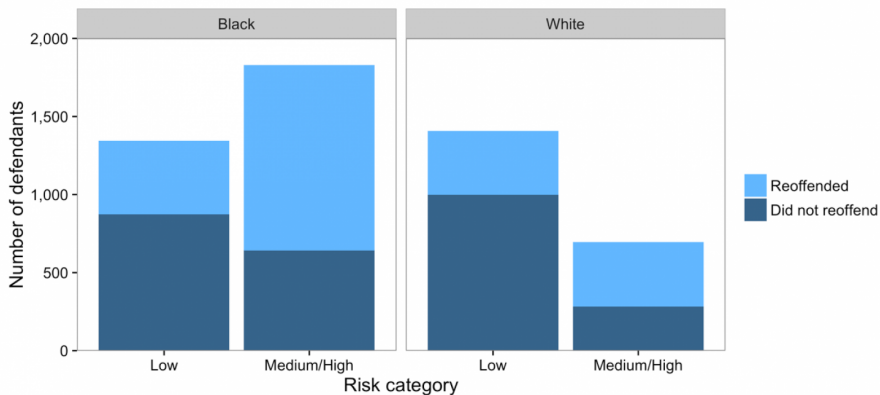
Recidivism rates by risk score

$y$  Result.

$a$  Assigned score.

$z$  Race.

$$\mathbb{P}^{\pi}(y \mid a, z) = \mathbb{P}^{\pi}(y \mid a) \qquad \text{(calibration)}$$

$$\mathbb{P}^{\pi}(a \mid y, z) = \mathbb{P}^{\pi}(a \mid y) \qquad \text{(balance)}$$
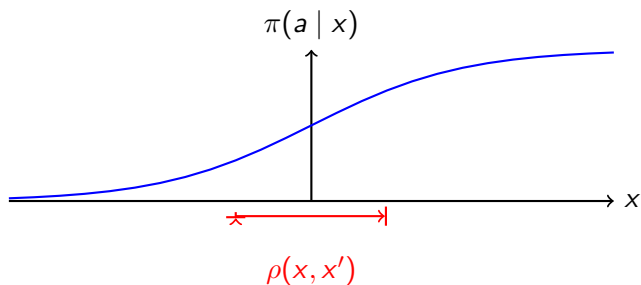
$y$ Result.

$a$ Assigned score.

$z$ Race.

$$\mathbb{P}^\pi(y \mid a, z) = \mathbb{P}^\pi(y \mid a) \qquad \text{(calibration)}$$
$$\mathbb{P}^\pi(a \mid y, z) = \mathbb{P}^\pi(a \mid y) \qquad \text{(balance)}$$

Meritocratic decision

$$a_t(\theta, x_t) \in \arg\max_a \mathbb{E}_\theta(U \mid a, x_t) = \int_{\mathcal{Y}} U(a_t, y) \, \mathbb{E}_\theta(U \mid a_t, x_t) \qquad (3.1)$$

$$D[\pi(a \mid x), \pi(a \mid x')] \leq \rho(x, x'). \tag{3.2}$$

# The Bayesian approach to fairness

### The value of a policy

Let $\lambda$ represent the trade-off between utility and fairness.

$$V(\lambda, \theta, \pi) = \lambda \overbrace{U(\theta, \pi)}^{\text{utility}} - \underbrace{(1 - \lambda)F(\theta, \pi)}_{\text{fairness violation}} \qquad (3.3)$$

$$V(\lambda, \xi, \pi) = \int_{\Theta} V(\lambda, \theta, \pi) \, d\xi(\theta). \qquad (3.4)$$

Online resources

▶ COMPAS analysis by propublica
https://github.com/propublica/compas-analysis

▶ Open policing database https://openpolicing.stanford.edu/