



# Fairness in AI

Christos Dimitrakakis

April 9, 2019

# Pervasive “intelligent” systems



Home assistants



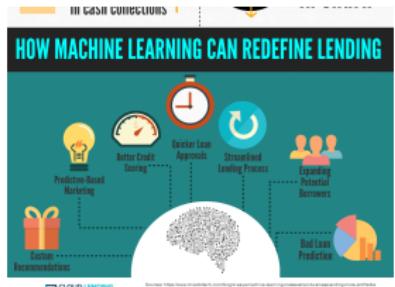
Autonomous vehicles



Web advertising



Ridesharing

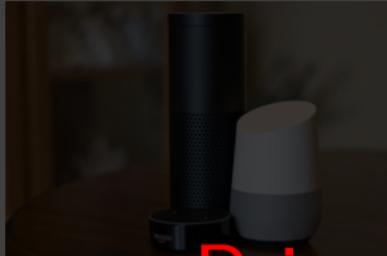


Lending



Public policy

# Pervasive “intelligent” systems



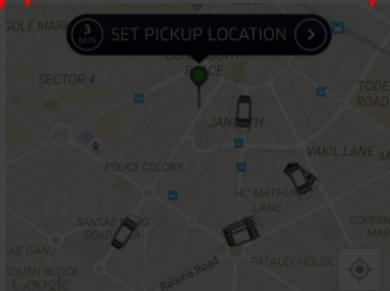
Home assistants



Autonomous vehicles



Web advertising



Ridesharing



Lending



Public policy

# Fairness

# Meritocracy

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

## Example 2 (Additional information)

- ▶ 70% of admitted Gota graduates with 4+ get their degree.
- ▶ 50% of admitted Vasa graduates with 5 get their degree.

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

## Example 2 (Additional information)

- ▶ 70% of admitted Gota graduates with 4+ get their degree.
- ▶ 50% of admitted Vasa graduates with 5 get their degree.

We still don't know how a **specific** student will do!

## Solutions

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

## Example 2 (Additional information)

- ▶ 70% of admitted Gota graduates with 4+ get their degree.
- ▶ 50% of admitted Vasa graduates with 5 get their degree.

We still don't know how a **specific** student will do!

## Solutions

- ▶ Admit **everybody**?

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

## Example 2 (Additional information)

- ▶ 70% of admitted Gota graduates with 4+ get their degree.
- ▶ 50% of admitted Vasa graduates with 5 get their degree.

We still don't know how a **specific** student will do!

## Solutions

- ▶ Admit **everybody**?
- ▶ Admit **randomly**?

# Meritocracy

## Example 1 (College admissions)

- ▶ Student *A* has a grade 4/5 from Gota Highschool.
- ▶ Student *B* has a grade 5/5 from Vasa Highschool.

## Example 2 (Additional information)

- ▶ 70% of admitted Gota graduates with 4+ get their degree.
- ▶ 50% of admitted Vasa graduates with 5 get their degree.

We still don't know how a **specific** student will do!

## Solutions

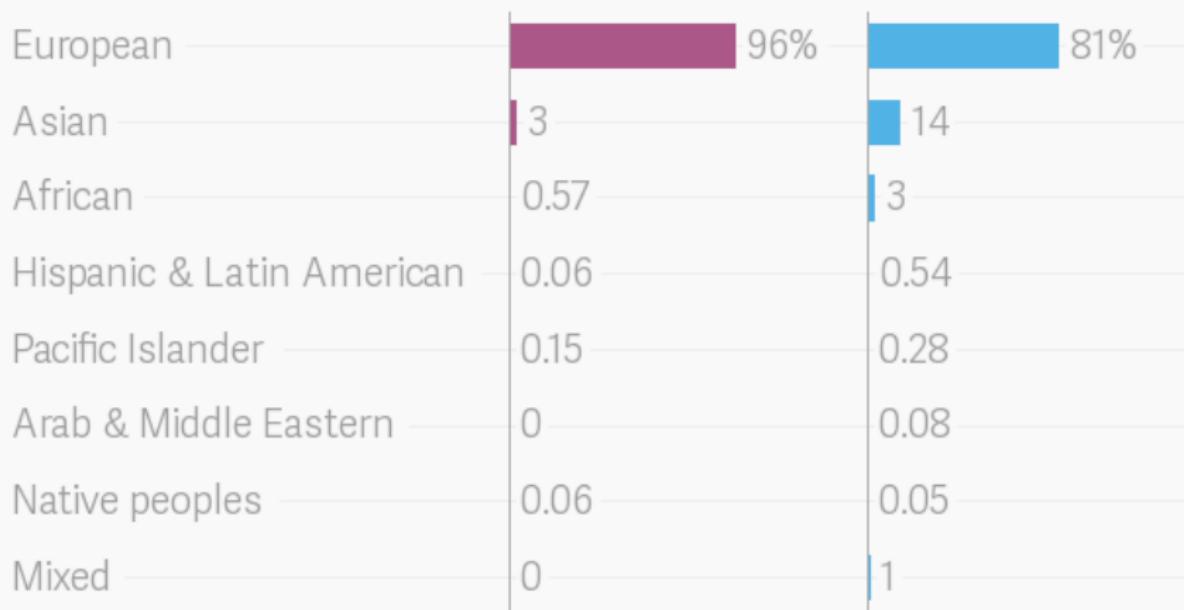
- ▶ Admit **everybody**?
- ▶ Admit **randomly**?
- ▶ Use **prediction** of individual academic performance?

# Proportional representation

Little progress is being made to improve diversity in genomics

Share of samples in genetic studies, by ancestry

■ 373 studies, up to 2009 ■ 2,511 studies, up to 2016



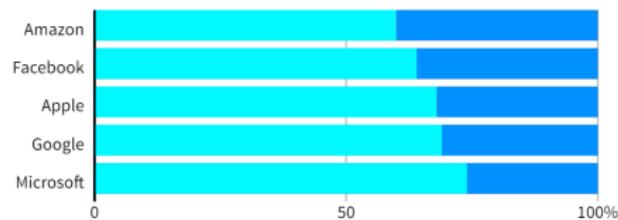
# Hiring decisions

## Dominated by men

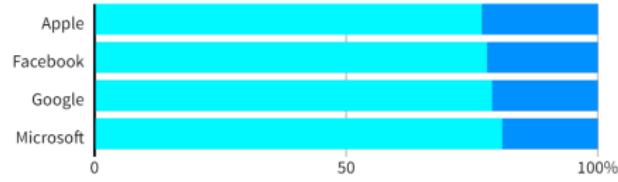
Top U.S. tech companies have yet to close the gender gap in hiring, a disparity most pronounced among technical staff such as software developers where men far outnumber women. Amazon's experimental recruiting engine followed the same pattern, learning to penalize resumes including the word "women's" until the company discovered the problem.

### GLOBAL HEADCOUNT

■ Male ■ Female



### EMPLOYEES IN TECHNICAL ROLES



C. Dimitrakakis

Fairness in AI



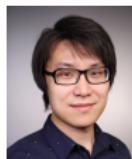
# Fairness and information

## Example 3 (College admissions data)

School	Male	Female
A	62%	82%
B	63%	68%
C	37%	34%
D	33%	35%
E	28%	24%
F	6%	7%
<i>Average</i>	45%	38%

## Bayesian Fairness<sup>1</sup>

Fairness under partial information.



<sup>1</sup>C. Dimitrakakis, Y. Liu, G. Radanovic, D. C. Parkes, AAAI 2019.

# Bail decisions

$x, z$



$\downarrow \pi$

A black arrow pointing downwards, with the Greek letter  $\pi$  written next to it, indicating a mapping or function.

- ▶  $x$  - features
- ▶  $z$  - sensitive features
- ▶  $a$  - action
- ▶  $y$  - outcome

# Bail decisions

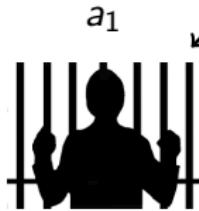
$x, z$



$\downarrow \pi$



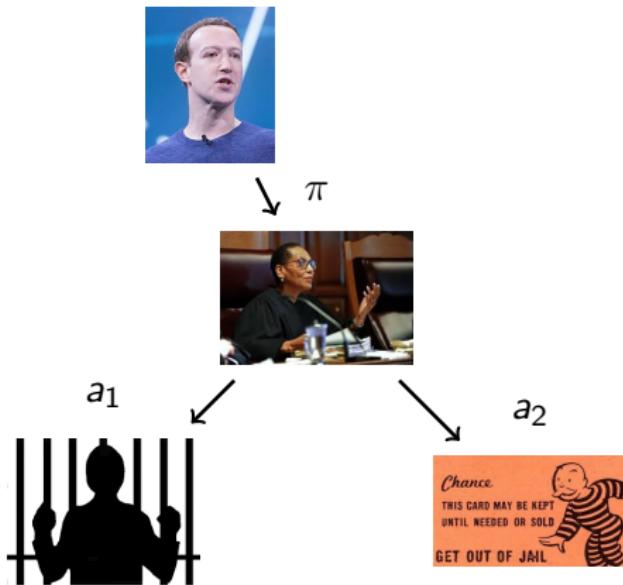
$a_1$



- ▶  $x$  - features
  - ▶  $z$  - sensitive features
  - ▶  $a$  - action
  - ▶  $y$  - outcome
- $\pi(a | x)$  (policy)

# Bail decisions

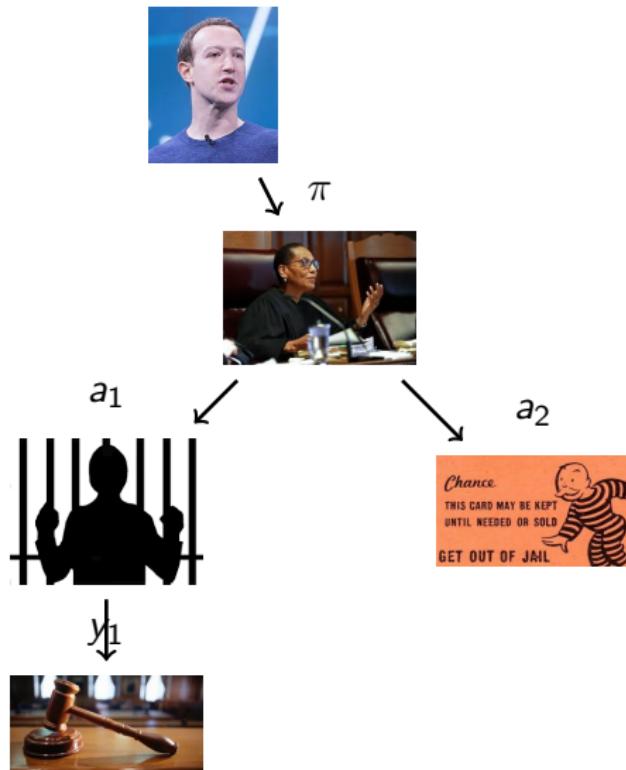
$x, z$



- ▶  $x$  - features
  - ▶  $z$  - sensitive features
  - ▶  $a$  - action
  - ▶  $y$  - outcome
- $\pi(a | x)$  (policy)

# Bail decisions

$x, z$

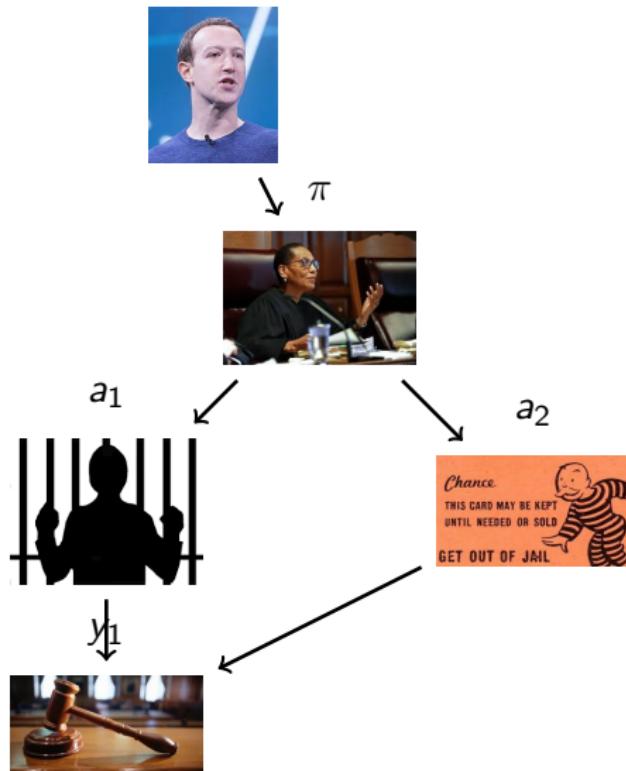


- ▶  $x$  - features
  - ▶  $z$  - sensitive features
  - ▶  $a$  - action
  - ▶  $y$  - outcome
- $\pi(a | x)$  (policy)

$$P_\theta(x, y, z) \quad (\text{model})$$

# Bail decisions

$x, z$

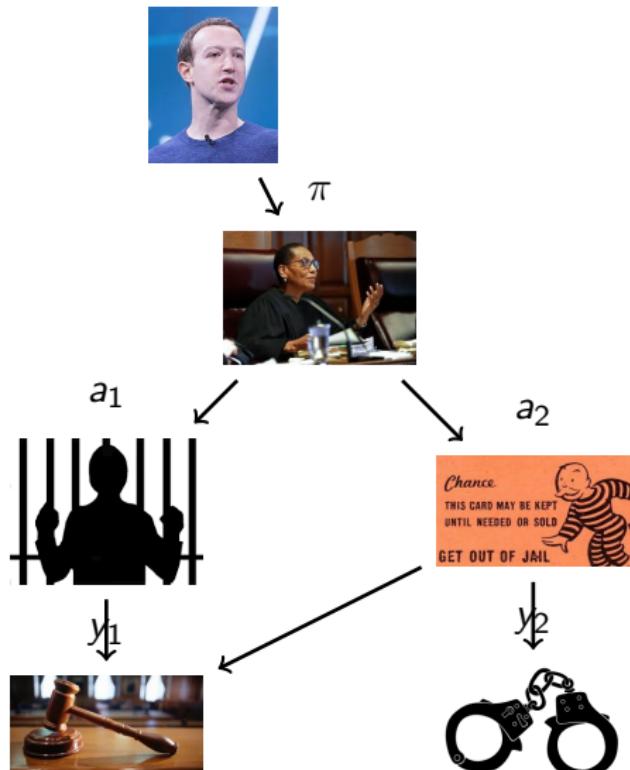


- ▶  $x$  - features
  - ▶  $z$  - sensitive features
  - ▶  $a$  - action
  - ▶  $y$  - outcome
- $$\pi(a | x) \quad (\text{policy})$$

$$P_\theta(x, y, z) \quad (\text{model})$$

# Bail decisions

$x, z$

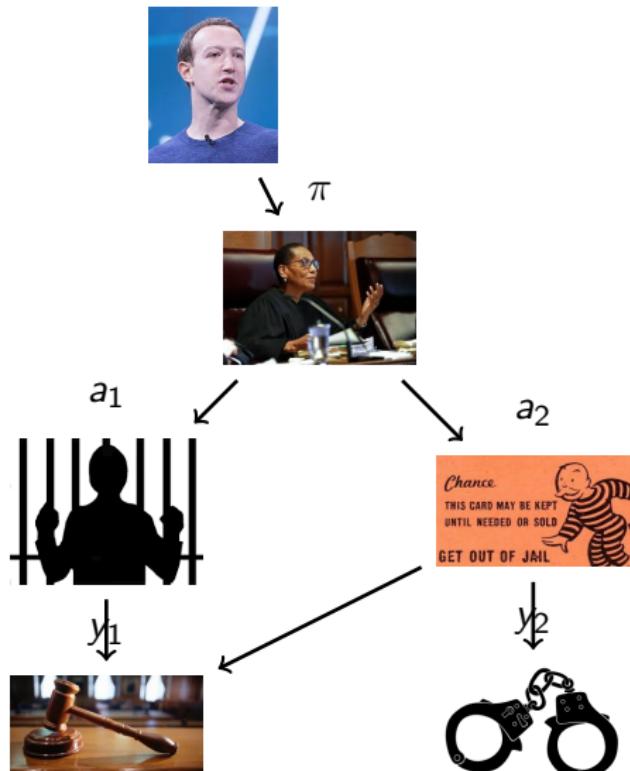


- ▶  $x$  - features
  - ▶  $z$  - sensitive features
  - ▶  $a$  - action
  - ▶  $y$  - outcome
- $\pi(a | x)$  (policy)

$$P_\theta(x, y, z) \quad (\text{model})$$

# Bail decisions

$x, z$



- $x$  - features
  - $z$  - sensitive features
  - $a$  - action
  - $y$  - outcome
- $\pi(a | x)$  (policy)

$$\pi(a | x) \quad (\text{policy})$$

$$P_\theta(x, y, z) \quad (\text{model})$$

$$U(a, y) \quad (\text{utility})$$

# Bail decisions

$x, z$



$\pi$



$a_1$



$a_2$



$y_1$



$y_2$



C. Dimitrakakis

## His honour the machine

Prisoners released on bail\*

%

Chosen by  
judges

18.6

14.9

Suggested  
by algorithm

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

$\pi(a | x)$

(policy)

$P_\theta(x, y, z)$

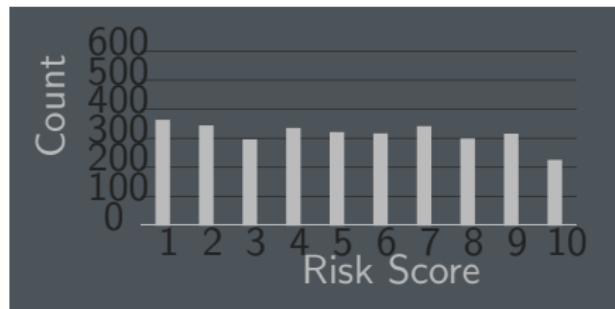
(model)

$U(a, y)$

(utility)

# The USA COMPAS System<sup>2</sup>

Assigns a **risk score**  $a_t$  to defendant  $x_t$ .



Black



White

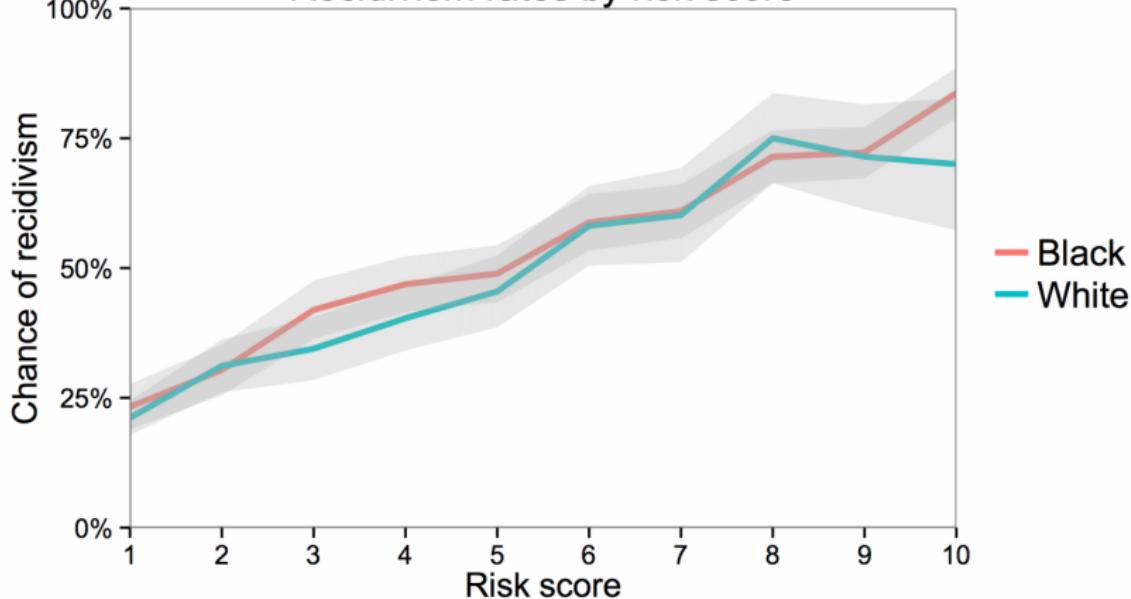
Figure : Apparent bias in risk scores towards black versus white defendants.

$$\mathbb{P}_\theta^\pi(a | z) = \mathbb{P}_\theta^\pi(a)$$

(equal treatment)

<sup>2</sup>Pro-publica, 2016

## Recidivism rates by risk score



Washington Post, 2016

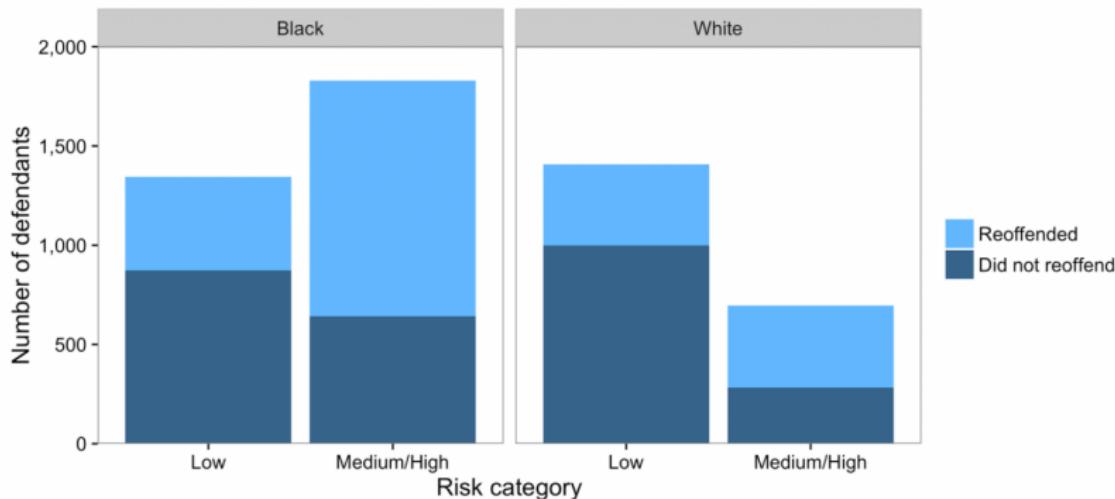
$y$  Result.

$a$  Assigned score.

$z$  Race.

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

$$\mathbb{P}_\theta^\pi(a \mid y, z) = \mathbb{P}_\theta^\pi(a \mid y) \quad (\text{balance})$$



Pro-publica, 2016

*y* Result.

*a* Assigned score.

*z* Race.

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

$$\mathbb{P}_\theta^\pi(a \mid y, z) = \mathbb{P}_\theta^\pi(a \mid y) \quad (\text{balance})$$

# The value of a policy

Fairness metrics: balance

$$F_{\text{balance}}(\theta, \pi) \triangleq \sum_{y,z,a} |\mathbb{P}_{\theta}^{\pi}(a | y, z) - \mathbb{P}_{\theta}^{\pi}(a | y)|^2 \quad (1.1)$$

# The value of a policy

Fairness metrics: balance

$$F_{\text{balance}}(\theta, \pi) \triangleq \sum_{y,z,a} |\mathbb{P}_{\theta}^{\pi}(a | y, z) - \mathbb{P}_{\theta}^{\pi}(a | y)|^2 \quad (1.1)$$

Utility: Classification accuracy

$$U(\theta, \pi) = \mathbb{P}_{\theta}^{\pi}(y_t = a_t)$$

# The value of a policy

Fairness metrics: balance

$$F_{\text{balance}}(\theta, \pi) \triangleq \sum_{y,z,a} |\mathbb{P}_{\theta}^{\pi}(a | y, z) - \mathbb{P}_{\theta}^{\pi}(a | y)|^2 \quad (1.1)$$

Utility: Classification accuracy

$$U(\theta, \pi) = \mathbb{P}_{\theta}^{\pi}(y_t = a_t)$$

Use  $\lambda$  to trade-off utility and fairness

$$V(\lambda, \theta, \pi) = (1 - \lambda) \overbrace{U(\theta, \pi)}^{\text{utility}} - \lambda \underbrace{F(\theta, \pi)}_{\text{unfairness}} \quad (1.2)$$

# Model uncertainty

$\theta$  is unknown

## Theorem 4

*A decision rule in the form of a lottery, i.e.*

$$\pi(a | x) = p_a$$

*can be the only way to satisfy balance for all possible  $\theta$ .*

## Possible solutions

- ▶ Marginalize over  $\theta$  ("expected" model)
- ▶ Use Bayesian reasoning

# The setting

## Bayesian modelling of uncertainty

- ▶ Model family  $\{P_\theta(x, y, z) \mid \theta \in \Theta\}$
- ▶ Prior belief  $\xi_0(\theta)$ .
- ▶ Data  $D \sim P_\theta(x, y, z)$ .
- ▶ Posterior belief  $\xi(\theta) \triangleq \xi_0(\theta \mid D)$ .

# The setting

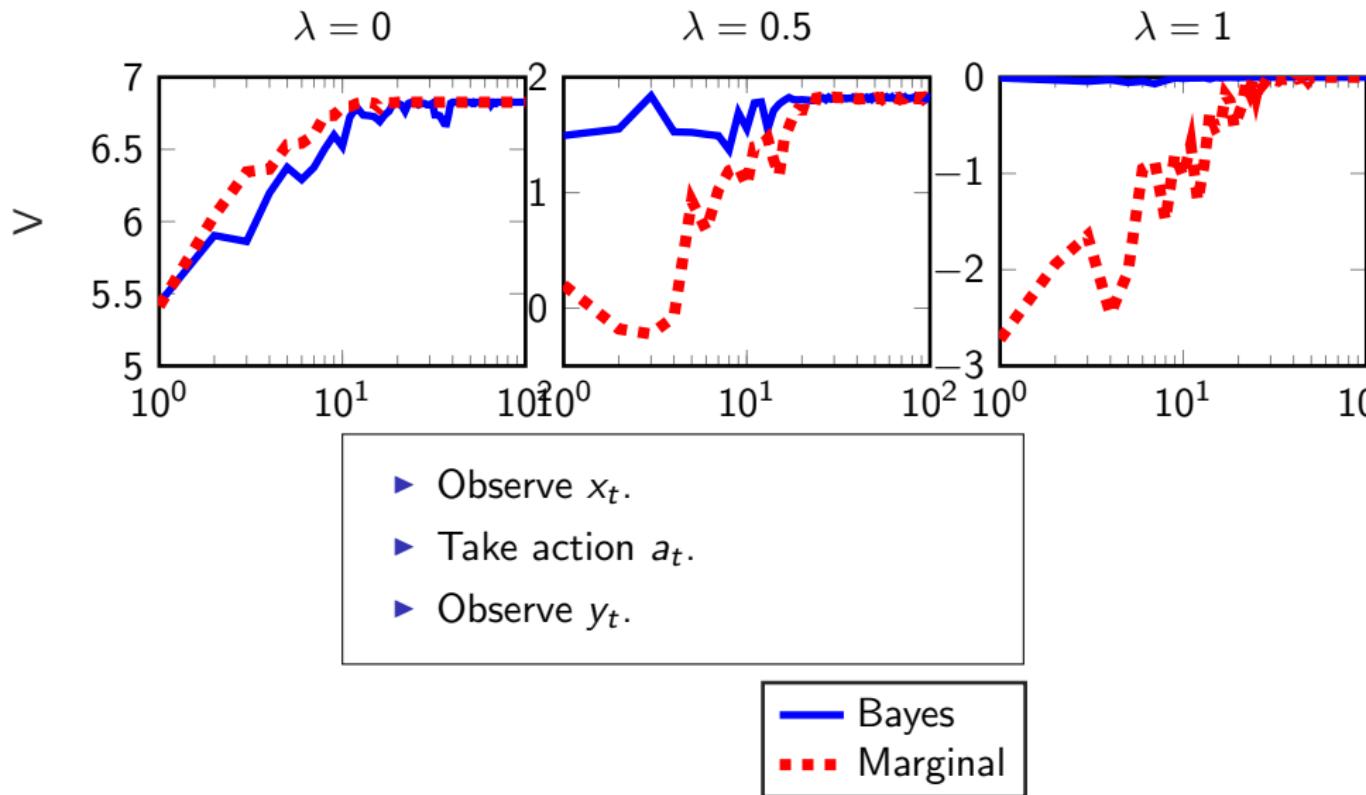
## Bayesian modelling of uncertainty

- ▶ Model family  $\{P_\theta(x, y, z) \mid \theta \in \Theta\}$
- ▶ Prior belief  $\xi_0(\theta)$ .
- ▶ Data  $D \sim P_\theta(x, y, z)$ .
- ▶ Posterior belief  $\xi(\theta) \triangleq \xi_0(\theta \mid D)$ .

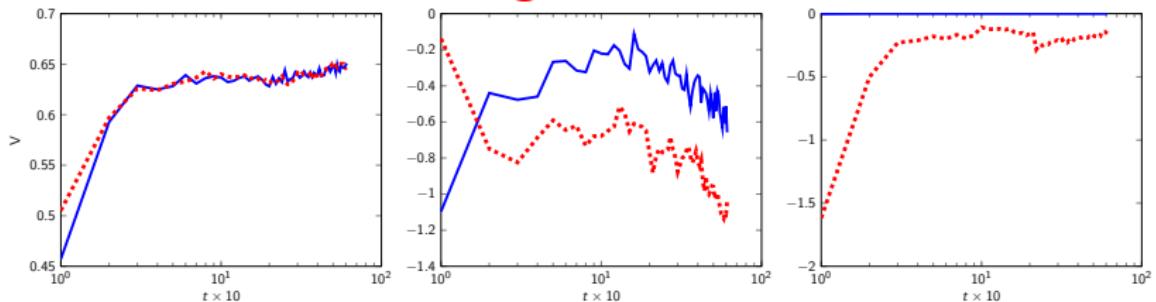
## The Bayesian value of a policy

$$V(\lambda, \xi, \pi) \triangleq \sum_{\theta \in \Theta} V(\lambda, \theta, \pi) \xi(\theta). \quad (1.3)$$

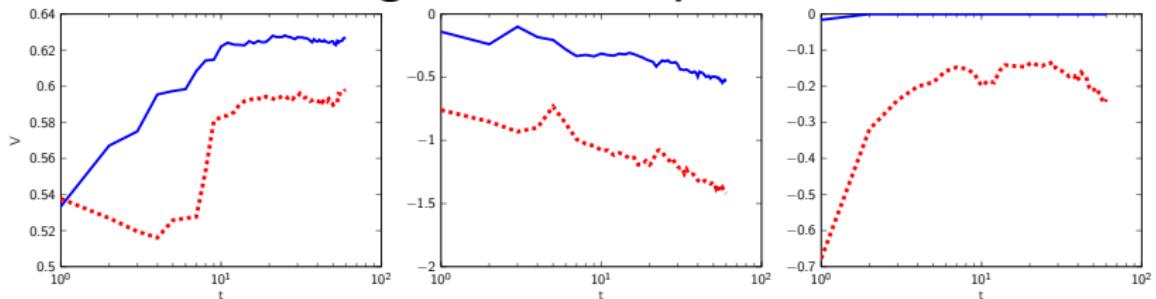
# Synthetic data



## Offline setting on COMPAS data

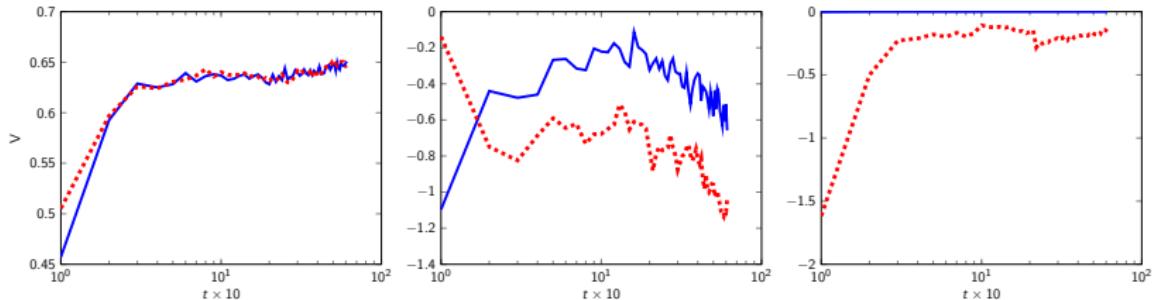


## Online setting: data seen depends on decisions

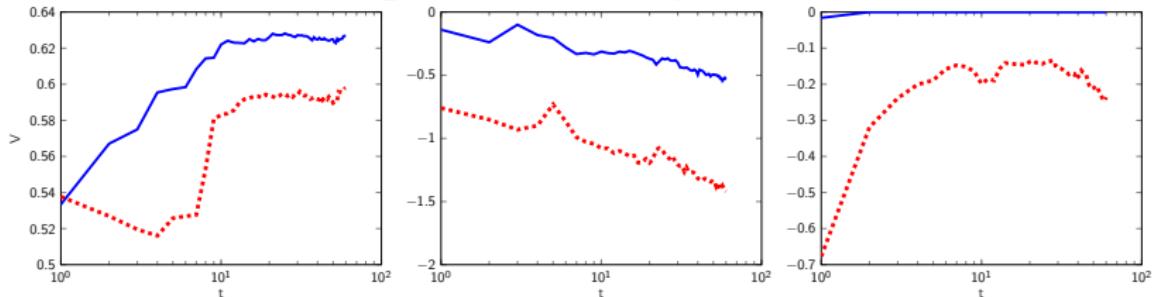


- ▶ Observe  $x_t$ .
  - ▶ Take action  $a_t$ .
  - ▶ Observe  $y_t$ .

## Offline setting on COMPAS data



**Online setting: data seen depends on decisions**



- ▶ Observe  $x_t$ .
- ▶ Take action  $a_t$ .
- ▶ Observe  $y_t$  if  $a_t = 1$ .

# Summary

- ▶ Existing criteria hard to satisfy.
- ▶ Bayesian framework:
  - ▶ Incorporates model uncertainty.
  - ▶ Considers informational aspects.
  - ▶ Lower fairness violation.

## Future work

- ▶ Non-myopic sequential decision making.
- ▶ Approximate Bayesian methods.



# Summary

- ▶ Existing criteria hard to satisfy.
- ▶ Bayesian framework:
  - ▶ Incorporates model uncertainty.
  - ▶ Considers informational aspects.
  - ▶ Lower fairness violation.

## Future work

- ▶ Non-myopic sequential decision making.
- ▶ Approximate Bayesian methods.

