

# Fairness Quantification in Machine Learning: The Power of Stochastic Satisfiability

Bishwamitra Ghosh<sup>1</sup>, Debabrota Basu<sup>2</sup>, Kuldeep S. Meel<sup>1</sup>

<sup>1</sup>National University of Singapore, Singapore, <sup>2</sup>Équipe Scool, Univ. Lille, Inria UMR 9189-CRISTAL, CNRS, Centrale Lille, France

## PROBABILISTIC FAIRNESS QUANTIFICATION

Notation: sensitive features  $\mathbf{A}$ , non-sensitive features  $\mathbf{X}$   
Given

- a binary classifier  $\mathcal{M} : (\mathbf{X}, \mathbf{A}) \rightarrow \hat{Y} \in \{0, 1\}$
- the probability distribution of features  $(\mathbf{X}, \mathbf{A}) \sim \mathcal{D}$

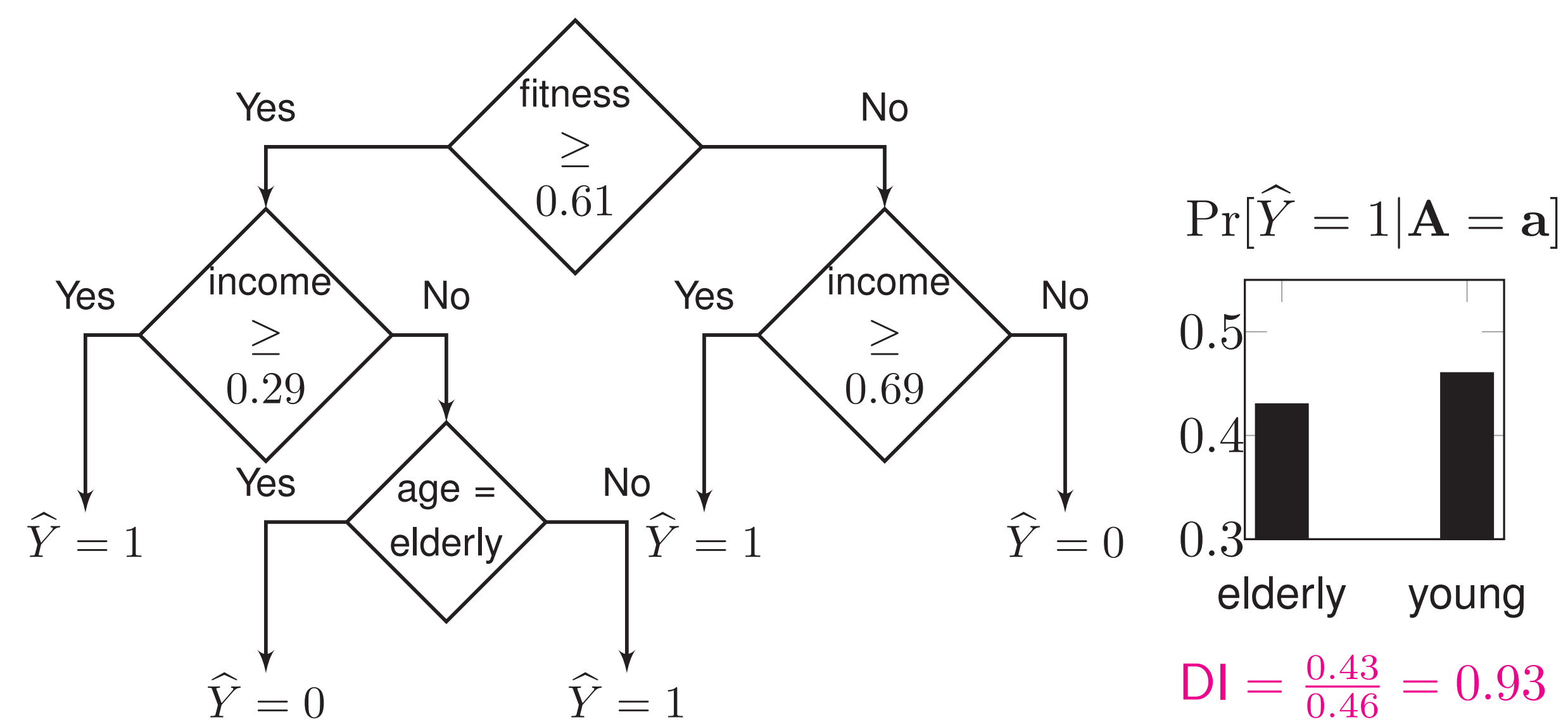
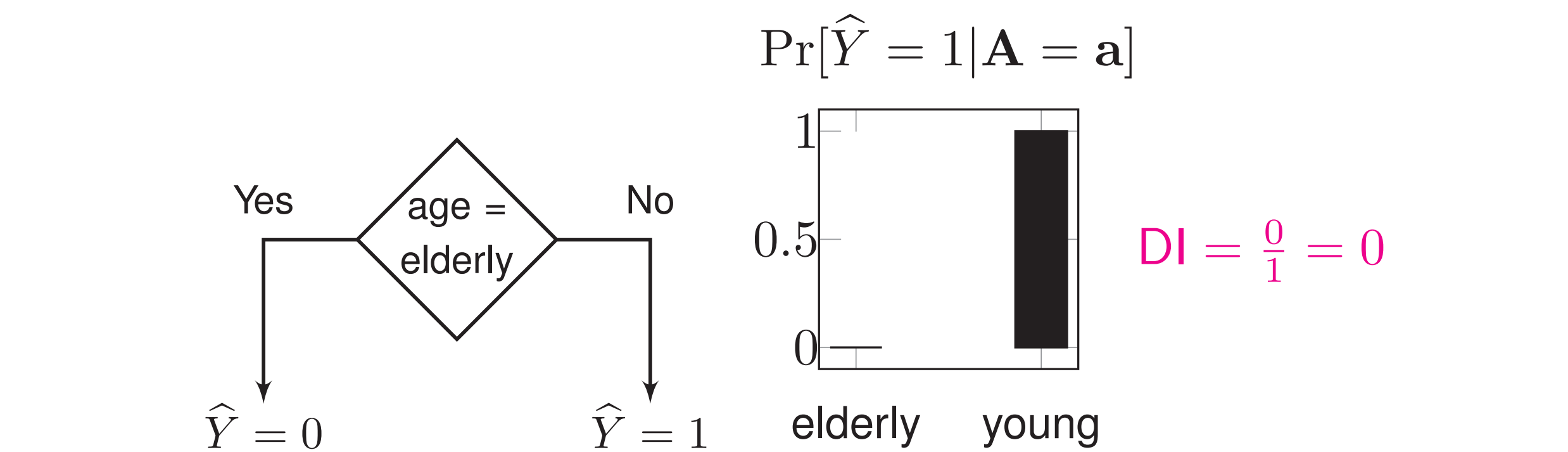
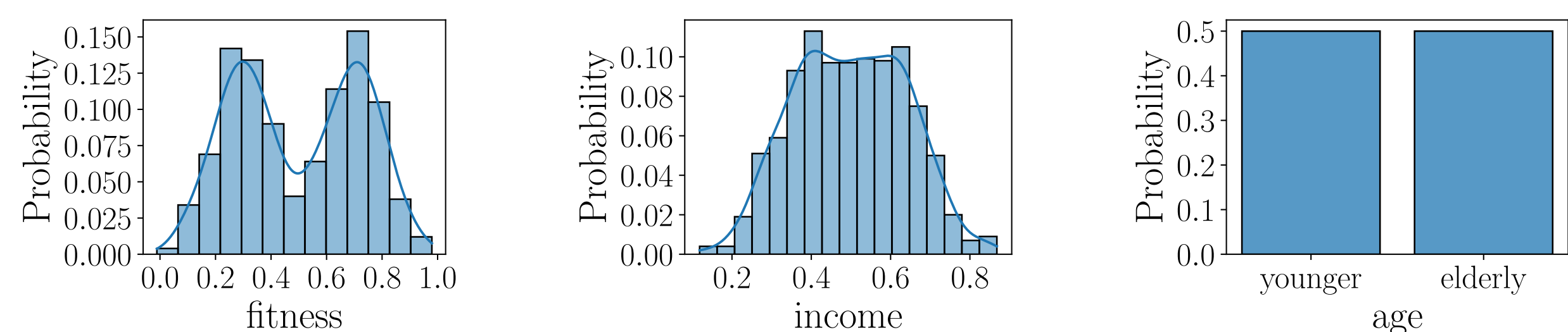
Problem: quantify the fairness of  $\mathcal{M}$  with respect to  $\mathcal{D}$

### Example of Group Fairness Metric

$$\text{Disparate impact, DI} = \frac{\min_{\mathbf{a}} \Pr[\hat{Y} = 1 | \mathbf{A} = \mathbf{a}]}{\max_{\mathbf{a}} \Pr[\hat{Y} = 1 | \mathbf{A} = \mathbf{a}]}$$

### Illustration of (Un)Fairness

- sensitive features  $\mathbf{A} = [\text{age}]$ ,  $\text{age} \in \{\text{elderly}, \text{young}\}$
- non-sensitive features  $\mathbf{X} = [\text{fitness}, \text{income}]$



## REFERENCES

[GBM-22] Algorithmic Fairness Verification with Graphical Models, AAAI-2022

[GBM-21] Justicia: A Stochastic SAT Approach to Formally Verify Fairness, AAAI-2021



Code

## CONTRIBUTION

**Justicia**, a Stochastic Satisfiability (SSAT) based framework for fairness quantification with scalable and accurate performance

**Key Idea:** Encoding conditional probability as SSAT

## STOCHASTIC SATISFIABILITY (SSAT)

Compute the probability of satisfaction of a CNF formula  $\phi$  with quantified Boolean variables

$$\Phi = \underbrace{Q_1 X_1, \dots, Q_m X_m}_{\text{prefix}}; \underbrace{\phi}_{\text{CNF}}$$

Quantifier  $Q_i \in \{\exists, \forall, \mathfrak{P}^{p_i}\}$  is an existential ( $\exists$ ), an universal ( $\forall$ ), or a random ( $\mathfrak{P}^{p_i}$ ) quantifier with  $p_i = \Pr[X_i = 1]$

### Semantics

1.  $\Pr[\mathfrak{P}^p X \Phi] = p \Pr[\Phi|_X] + (1-p) \Pr[\Phi|_{\neg X}]$  (**expectation**)
2.  $\Pr[\exists X \Phi] = \max_X \{\Pr[\Phi|_X], \Pr[\Phi|_{\neg X}]\}$  (**maximization**)
3.  $\Pr[\forall X \Phi] = \min_X \{\Pr[\Phi|_X], \Pr[\Phi|_{\neg X}]\}$  (**minimization**)
4.  $\Pr[1] = 1, \Pr[0] = 0$

### Existential-Random SSAT

$\Phi_{ER} = \exists X_2, \exists X_3, \mathfrak{P}^{0.75} X_1, (\neg X_1 \vee X_2) \wedge (X_1 \vee \neg X_2 \vee \neg X_3) \wedge X_3$   
Solution:  $\Pr[\Phi_{ER}] = 0.75, \sigma^* = \{X_2 = 1, X_3 = 1\}$

### Universal-Random SSAT

$\Phi_{UR} = \forall X_2, \forall X_3, \mathfrak{P}^{0.75} X_1, (\neg X_1 \vee X_2) \wedge (X_1 \vee \neg X_2 \vee \neg X_3) \wedge X_3$   
Solution:  $\Pr[\Phi_{UR}] = 0, \sigma^* = \{X_2 = 0, X_3 = 0\}$

## FAIRNESS QUANTIFICATION USING SSAT

For the maximum conditional probability  $\max_{\mathbf{a}} \Pr[\hat{Y} = 1 | \mathbf{A} = \mathbf{a}]$ , solve

$$\Phi_{ER} := \underbrace{\exists A_1, \dots, \exists A_n}_{\text{sensitive features}}; \underbrace{\mathfrak{P}^{p_1} X_1, \dots, \mathfrak{P}^{p_m} X_m}_{\text{non-sensitive features}}; \underbrace{\phi_{\hat{Y}}}_{\text{Classifier}}$$

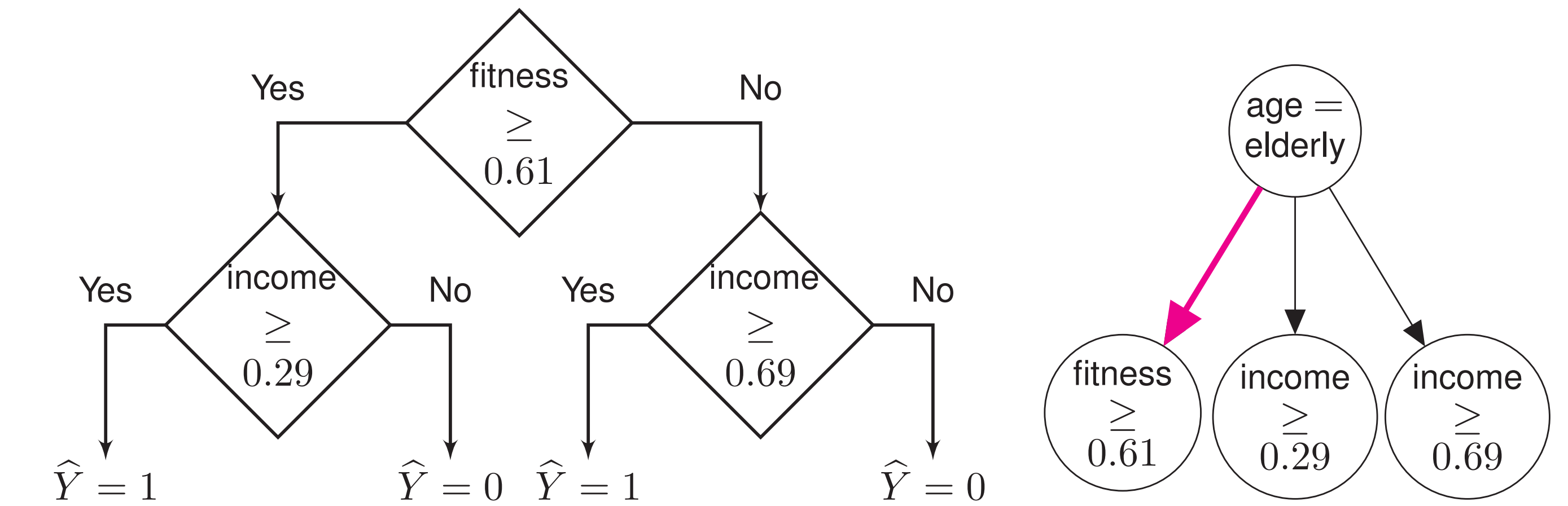
For the minimum conditional probability  $\min_{\mathbf{a}} \Pr[\hat{Y} = 1 | \mathbf{A} = \mathbf{a}]$ , solve

$$\Phi_{UR} := \underbrace{\forall A_1, \dots, \forall A_n}_{\text{sensitive features}}; \underbrace{\mathfrak{P}^{p_1} X_1, \dots, \mathfrak{P}^{p_m} X_m}_{\text{non-sensitive features}}; \underbrace{\phi_{\hat{Y}}}_{\text{Classifier}}$$

$$\text{Disparate impact} = \frac{\Pr[\Phi_{UR}]}{\Pr[\Phi_{ER}]}$$

## FEATURE CORRELATION AS BAYESIAN NETWORK

Encode conditional probability  $p = \Pr[X = 1 | \text{parent}(X) = \mathbf{x}]$  from a Bayesian Network using additional variables and clauses



### Variables

- $X_1 := \text{fitness} \geq 0.61$
- $A_1 := \text{age} = \text{elderly}$
- $Z_1 := \text{fitness} \geq 0.61 \mid \text{age} = \text{elderly}$

### Clauses

- $A_1 \wedge Z_1 \rightarrow X_1$
- $A_1 \wedge \neg Z_1 \rightarrow \neg X_1$

## EXPERIMENTAL RESULTS

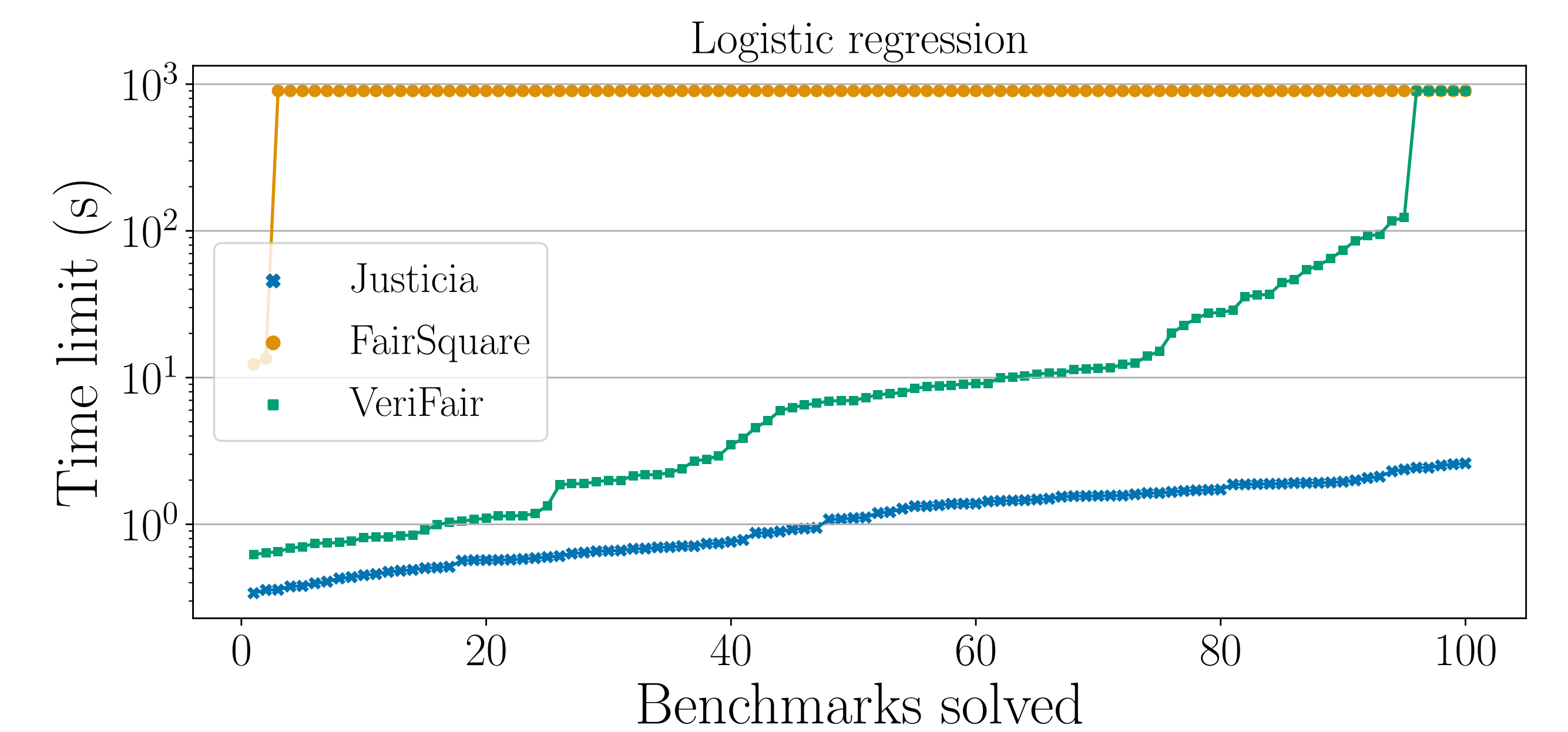


Figure 1: Higher scalability: Justicia achieves 1 to 2 orders improvement

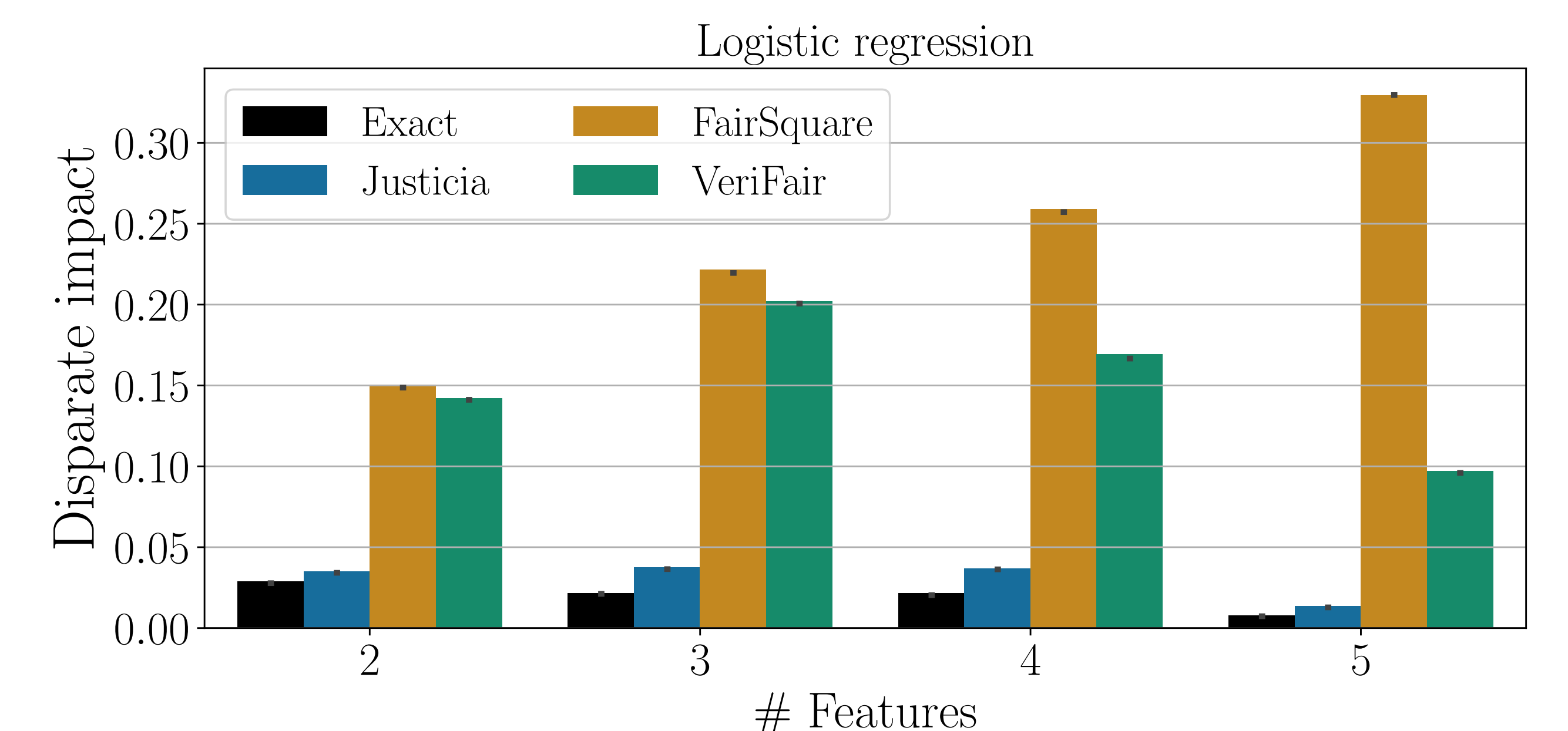


Figure 2: Improved accuracy due to encoding feature correlations