

# THE ROLE OF BAYESIAN UNCERTAINTY IN ASSESSING THE FAIRNESS OF MACHINE LEARNING MODELS

Francesca Panero<sup>1</sup>, Ernst Wit<sup>2</sup> and Marco Scutari<sup>3</sup>

<sup>1</sup> MEMOTEF, Sapienza University of Rome <sup>2</sup> Università della Svizzera Italiana <sup>3</sup> IDSIA

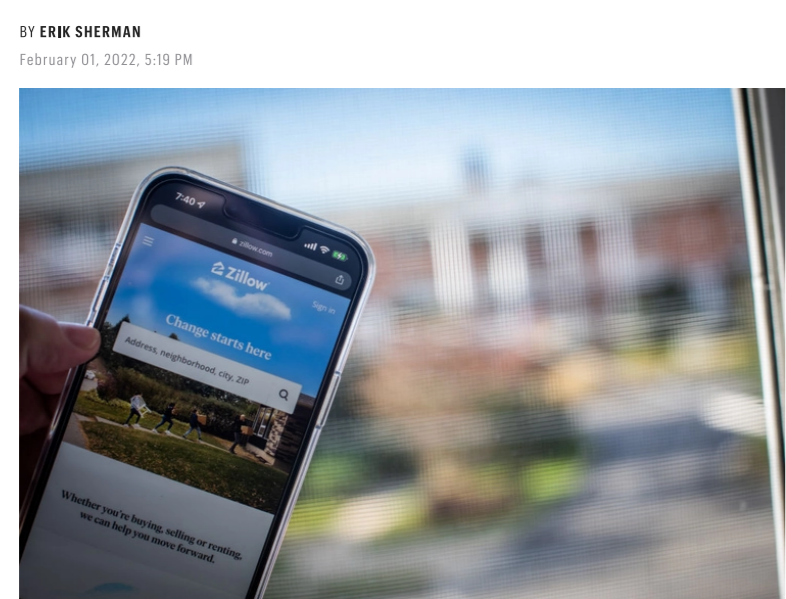


SAPIENZA  
UNIVERSITÀ DI ROMA



## WE NEED FAIRNESS!

**FORTUNE | EDUCATION**  
What Zillow's failed algorithm means for the future of data science

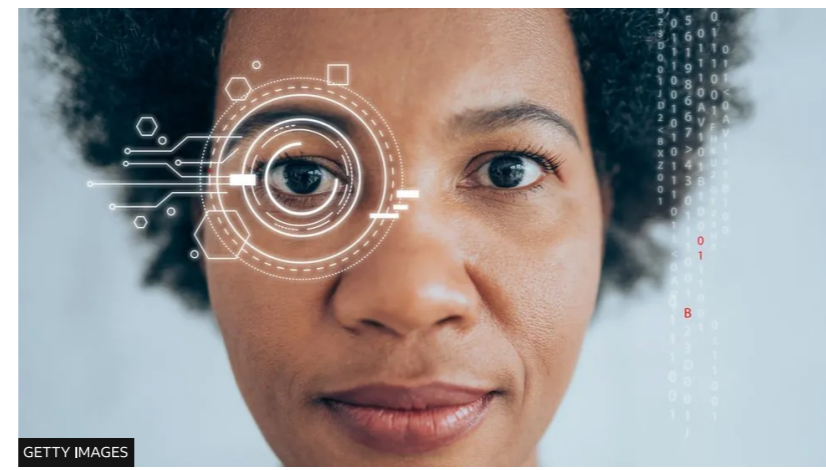


Big-data analysis told Zillow what to offer and how much to charge on the flip. Easy peasy. Except, come 2021, the wheels came off. Zillow had bought thousands of houses, and the algorithms didn't factor in repairs with the skyrocketing costs of materials and labor.



**Los Angeles**  
**LAPD ended predictive policing programs amid public outcry. A new effort shares many of their flaws**  
Documents show how data-driven policing programs reinforced harmful patterns, fueling the over-policing of Black and brown communities

**BBC | NEWS**  
Twitter finds racial bias in image-cropping AI



Twitter's automatic cropping of images had underlying issues that favoured white individuals over black people, and women over men, the company said.

**BBC | NEWS**  
Amazon scrapped 'sexist AI' tool



An algorithm that was being tested as a recruitment tool by online giant Amazon was sexist and had to be scrapped, according to a Reuters report.

## WHAT WE WANT FROM FAIR MODELS

- **Interpretability:** What do the parameters mean in the domain the data come from?
- **Explainability:** Why does the model give specific predictions?
- **Confidence:** Are the effects we measure statistically significant? What is the error margin of predictions?

## (FREQUENTIST) FAIR RIDGE REG MODEL (FRRM)

Consider a regression model with  $X$  predictors,  $S$  sensitive attributes, response  $y$  and a given level of fairness  $r \in [0, 1]$  (0 = complete fairness, 1 = no fairness constraints).

To bound the dependence between  $\hat{y}$  and  $S$  by a fairness budget  $r$ , Scutari et al. (2022) propose:

1. Compute  $\hat{U}$  (the fair predictors) as  $X = S\hat{B}_{OLS} + \hat{U}$ .
2. Estimate  $\hat{\beta}_{FRRM} = (\hat{U}^T \hat{U})^{-1} \hat{U}^T y$ .
3. To find  $\hat{\alpha}_{FRRM}$ , solve

$$\operatorname{argmin}_{\alpha, \beta} \|y - S\alpha - \hat{U}\beta\|_2^2 + \lambda(r)\|\alpha\|_2^2,$$

where  $\lambda(r)$  satisfies

$$R_S^2(\alpha, \beta) = \frac{\operatorname{VAR}(S\alpha)}{\operatorname{VAR}(S\alpha + \hat{U}\hat{\beta}_{FRRM})} = r.$$

- ✓ Single solution, computationally **inexpensive**.
- ✓ **Pluggable** fairness constraints, **interpretable** and **explainable**.
- ✓ Works for all **generalised linear models**.
- ✗ Lacking **uncertainty quantification**.

## REFERENCES

M. Scutari, F. Panero and M. Proissl (2022). Achieving Fairness with a Simple Ridge Penalty. *Statistics and Computing*, 32, 77. <https://cran.r-project.org/web/packages/fairml/>

## THE BAYESIAN ALTERNATIVE (BFRRM)

To design the Bayesian version of FRRM, we define the likelihood

$$y | \alpha, \beta, \hat{U}, S \sim \mathcal{N}(\hat{U}\beta + S\alpha, \sigma^2 I_n)$$

and the priors on  $\alpha, \beta, \sigma^2$ :

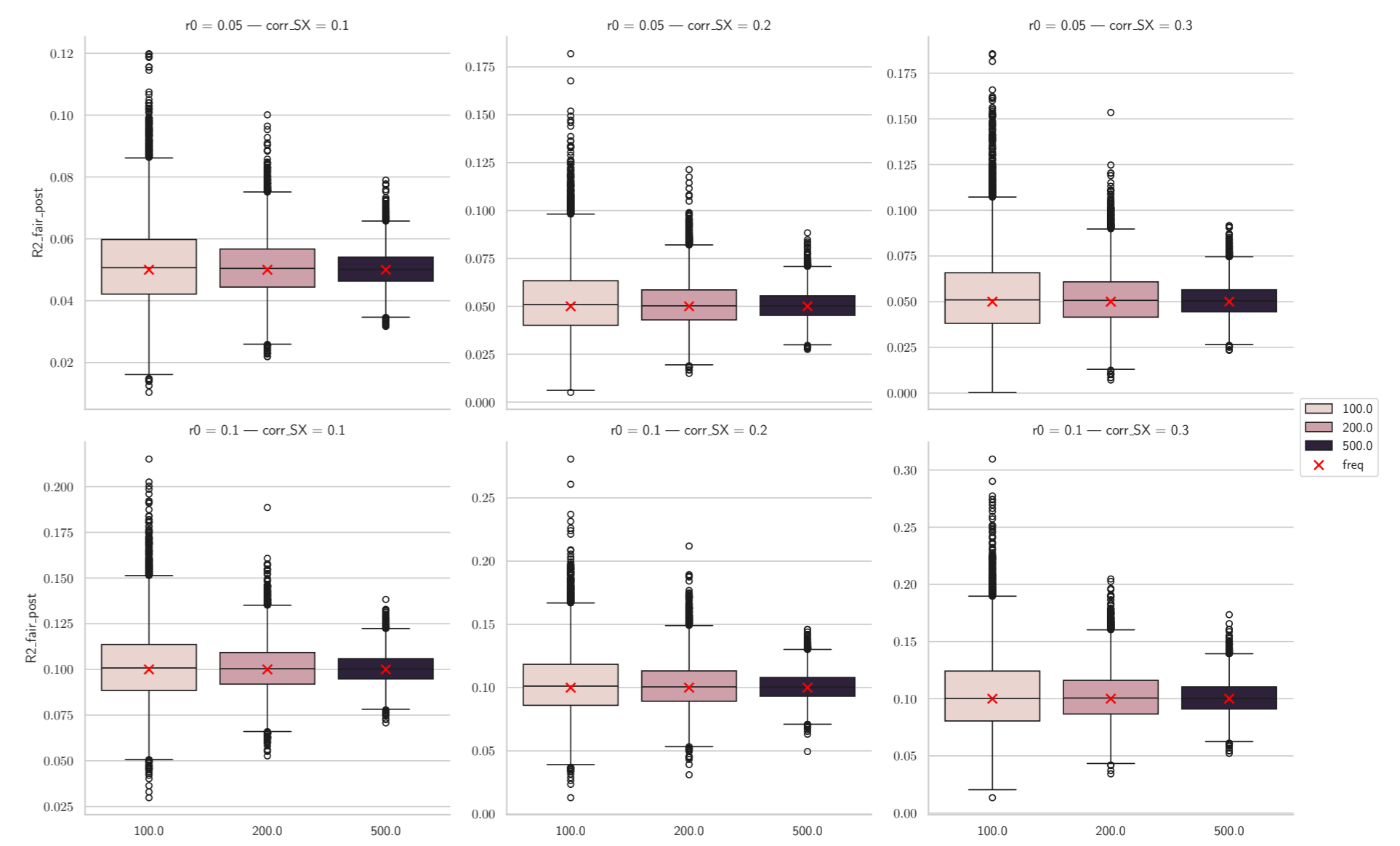
$$\beta | \sigma^2 \sim \mathcal{N}(0, g\sigma^2(\hat{U}^T \hat{U})^{-1}), \quad g > 0$$

$$\alpha | \sigma^2, \lambda(r) \sim \mathcal{N}\left(0, \frac{\sigma^2}{\lambda(r)} I_q\right)$$

$$\sigma^2 | a_0, b_0 \sim IG(a_0, b_0)$$

This formulation gives closed form posteriors for all parameters. Therefore, we can easily obtain posterior means and credible intervals for  $\alpha, \beta, \sigma^2$ . From these, we derive uncertainty measures for the unfairness measure  $R_S^2(\alpha, \beta)$ .

## POSTERIOR CALIBRATION OF BFRRM



## REAL DATA

LSAC is a survey among US law school students ( $n = 5K$ , 9 predictors): we take GPA as response, and race, age, gender as sensitive attributes. COMPAS ( $n = 6K$ , 13 predictors) comprises demographic and criminal records of offenders in Florida: we take recidivating as response and offender's age, gender, race as sensitive attributes.

