Justicia: A Stochastic SAT Approach to Formally Verify Fairness

<u>Bishwamittra Ghosh¹</u>, Debabrota Basu^{2,3}, Kuldeep S. Meel¹

¹National University of Singapore, Singapore ²Chalmers University of Technology, Goteborg, Sweden ³Scool, Inria Lille-Nord Europe, France

Unfairness in machine learning



2

Motivation

Fairness metrics

- Independence
 - Disparate impact
 - Statistical parity
- Separation
 - Equalized odds
- Sufficiency
 - Causal fairness

Fairness algorithms

- Preprocessing
- In-processing
- Postprocessing

A framework for verifying different fairness metrics and algorithms

Contribution

Fairness verification framework Justicia based on Stochastic SAT (SSAT)

- Two fairness definitions: independence and separation
- Handle compound protected groups
 - White-male, Black-female etc.
- Scalable
- Robust

Problem statement

- *X* = non-protected attributes
- *A* = protected attributes
- $Y = \text{true class label}, \hat{Y} = \text{predicted class label}$

Given

- binary classifier \mathcal{M} : $(X, A) \rightarrow \{0, 1\}$
- probability distribution $X \sim \mathcal{D}$

verify whether ${\mathcal M}$ achieves independence and separation metrics with respect to the distribution ${\mathcal D}$

Key observation

Computing positive predictive value (PPV)

$$\Pr[\widehat{Y}=1|A=a]$$

is the building block of different fairness metrics

Two approaches

- Approach 1: enumeration on each A = a
- Approach 2: learning most favored group $a_{\rm fav}$ and least favored group $a_{\rm unfav}$ based on PPV

Stochastic SAT (SSAT)

An SSAT formula has a prefix and a CNF formula $\boldsymbol{\varphi}$

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

 Q_i is either

- universal (\forall) ,
- existential (3), or
- randomized (\mathbb{R}^{p_i}) quantification with $p_i = \Pr[X_i = \text{TRUE}]$

The goal in SSAT is to compute the probability of satisfaction $\Pr[\Phi]$

Example of SSAT

 $\Phi = \mathbb{R}^{0.25} X_1, \exists X_2, \exists X_3, \ (X_1 \lor \neg X_2) \land (\neg X_1 \lor X_2 \lor X_3) \land (\neg X_1)$

Semantics of SSAT

1. $\Pr[\text{TRUE}] = 1$, $\Pr[\text{FALSE}] = 0$, 2. $\Pr[\Phi] = \max_{X} \{\Pr[\Phi|_{X}], \Pr[\Phi|_{\neg X}]\}$ if X is existentially (\exists) quantified 3. $\Pr[\Phi] = \min_{X} \{\Pr[\Phi|_{X}], \Pr[\Phi|_{\neg X}]\}$ if X is universally (\forall) quantified 4. $\Pr[\Phi] = p \Pr[\Phi|_{X}] + (1-p)\Pr[\Phi|_{\neg X}]$ if X is randomized ($\mathbb{R}^{p_{i}}$) quantified where $\Phi|_{X}$ is the substitution of left-most variable in the prefix with X = TRUE

Solution from an SSAT solver: $Pr[\Phi] = 0.75$

Approach 1: Enumeration encoding

Consider a simple case

- Attributes $X \cup A$ are Boolean
- Classifier \widehat{Y} is a CNF formula $\phi_{\widehat{Y}}$
- $p_i = \Pr[X_i]$ is known for each non-protected attribute

The computation of

$$\Pr[\hat{Y}=1|A=a]$$

is equivalent to solving

$$\Phi_{\boldsymbol{a}} \coloneqq \mathbb{R}^{p_1} X_1, \dots, \mathbb{R}^{p_m} X_m, \exists A_1, \dots, \exists A_n, \phi_{\widehat{Y}} \land (A = \boldsymbol{a})$$

non-protected protected

Example of enumeration encoding



Computation of fairness metrics

• Disparate impact:

$$\frac{\Pr[\hat{Y}=1|\text{age}\geq 40]}{\Pr[\hat{Y}=1|\text{age}< 40]} = \frac{0.43}{0.43} = 1$$

• Statistical parity:

$$|\Pr[\hat{Y} = 1| \text{age} \ge 40] - \Pr[\hat{Y} = 1| \text{age} < 40]| = |0.43 - 0.43| = 0$$

It looks like there is no discrimination

We did not consider correlation among attributes

Enumeration encoding with correlation

Use $\Pr[F|age \ge 40]$ instead of $\Pr[F]$...

$$\Phi_{\text{age} \ge 40} = \mathbb{R}^{0.01} F, \mathbb{R}^{0.99} I, \mathbb{R}^{0.18} J, \exists A, (\neg F \lor I) \land (F \lor J) \land A$$

With correlation,
$$\Pr[\hat{Y} = 1 | \text{age} \ge 40] = 0.18$$

Similarly, $\Pr[\hat{Y} = 1 | \text{age} < 40] = 0.72$

Disparate impact
$$=\frac{0.18}{0.72} \neq 1$$

Statistical parity $= 0.72 - 0.18 = 0.54 \neq 0$

Appraoch 2: Learning encoding

- Enumeration encoding has to be solved for exponential combinations of compound protected groups
- SSAT allows us to learn the assignment to existential (∃) and universal
 (∀) variables
- Learning the most favored group

$$\Phi_{fav} = \exists A, R^{0.41}F, R^{0.93}I, R^{0.09}J, (\neg F \lor I) \land (F \lor J)$$

• Learning the least favored group

$$\Phi_{\text{unfav}} = \forall A, \mathbb{R}^{0.41} F, \mathbb{R}^{0.93} I, \mathbb{R}^{0.09} J, (\neg F \lor I) \land (F \lor J)$$

Experiments

- State of the art
 - FairSquare: computes weighted volume of logical program using SMT
 - VeriFair: probabilistic verification via sampling
 - AIF360 (computes metrics on a finite dataset)
- Classifiers:
 - Linear classifier (pseudo-Boolean encoding)
 - Decision tree

Accuracy

Metric	Exact	Justicia	FairSquare	VeriFair	AIF360
Disparate impact	0.26	0.25	0.99	0.99	0.25
Stat. parity	0.53	0.54			0.54

Justicia has less than 1%-error

Scalability

Dataset	Ricci		Tita	Titanic		COMPAS		Adult	
Classifier	DT	LR	DT	LR	DT	LR	DT	LR	
Justicia	0.1	0.2	0.1	0.9	0.1	0.2	0.2	1.0	
FairSquare	4.8		16.0		36.9		—	—	
VeriFair	5.3	2.2	1.2	0.8	15.9	11.3	295.6	61.1	

DT = decision tree

LR = logistic regression classifier

Justicia reports 1 to 3 orders of magnitude speed-up

Compound protected groups



Conclusion

- A stochastic SAT-based approach to formally verify different fairness metrics and algorithms
- First method to verify compound protected groups
- More accurate, scalable and robust than state-of-the-art methods
- Python library: pip install justicia



Source code & paper