Fairness and sensitive attribute inference

Jan Aalmoes

Joint work with Antoine Boutet and Vasisht Duddu INSA - CITI laboratory - INRIA - PRIVATICS team





November 23, 2023

<□ > < □ > < □ > < Ξ > < Ξ > Ξ の Q · 1/22



Who We Support 🗸 🔹 About Us 🗸

Resources

es Blog

Contact Us

It's here! COMPAS-R Core: Transparent RNA Built Your Way

Have you ever wondered how risk and needs assessments (RNAs) are developed? Let's take a behind-the-scenes look at our newest RNA, the COMPAS-R Core



Whether you w more about ou justice solution demo, our tear out your email in touch as soo

POINT PREDICTIVE

Introducing Case Manager & Rules Engine

Case Manager provides lending teams with an intelligent layer of risk controls, automation logic, action guidance, and key metrics across their originations process. Case Manager sports an elegant user experience and is the first SasS fraud and risk management solution that has been specificably built for automotive lenders that tightly integrates Point Predictive's full complement of fraud risk accres, elets, and consortium data with existing lean origination systems (LOS) to streamline operational risk workflows for analysts and underwriters. Traditional solutions offered to the auto industry are generic workflow solutions that require vepensive customization to meet the specific needs of the auto industry. Point Predictive's Case Manager and Rules Engine are designed from the ground up as "and forist" offers' definings.



Overview of fairness notions



Demographic parity

$$P(\hat{Y} = 1 | S = 0) = P(\hat{Y} = 1 | S = 1)$$

Equality of odds¹

 $\forall y \in \{0,1\} P(\hat{Y} = 1 | S = 0, Y = y) = P(\hat{Y} = 1 | S = 1, Y = y)$

• • • •

Individual fairness

▶ Fairness Through Unawareness: Only works if X indep. of S

Fairness Through Awareness²: TODO

²Dwork,Hardt,Pitassi,Reingold,Fairness Through Awareness (IT€S) 2012 つ ۹. 4/22

^{...}

 $^{^1\}mbox{Hardt},$ Price, Srebro, Equality of Opportunity in Supervised Learning (NeurIPS) 2016



Dataset : Law school admission. Sensitive attribute : race

<□ > < □ > < □ > < Ξ > < Ξ > Ξ の Q O 5/22

- \blacktriangleright Z = 0 Non-white
- \blacktriangleright *Z* = 1 White

Attribute Inference Attack

・ロ ・ ・ 日 ・ ・ 言 ・ ・ 言 ・ う へ で 6/22

In our work



Attack target







< □ ▶ < □ ▶ < ■ ▶ < ■ ▶ < ■ ▶ = 少 Q @ 8/22

Target model

Two possible attacks

- ▶ Soft labels (s/): 0.1, 0.8, 0.2, ···
- ▶ Hard labels (*hl*): 0, 1, 0, ···

Attack(*sl*) = Sensitive attribute

or

Attack(*hI*) = Sensitive attribute

<□ > < ⓓ > < ≧ > < ≧ > ≧ < ੭ < ♡ 9/22

Fairness and attribute inference attack using **hard labels**

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三三 - のへで 10/22

Contribution

Demographic parity forbids attribute inference attack using hard labels

< □ ▶ < @ ▶ < ≧ ▶ < ≧ ▶ ≧ ⑦ Q [©] 11/22

Definition

The fairness-level of a classifier \hat{Y} , that we call $Fl(\hat{Y})$, is

$$FI(\hat{Y}) = |P(\hat{Y} = 1|S = 0) - P(\hat{Y} = 1|S = 1)|$$

 $Fl(\hat{Y}) \in [0,1]$ indicates how far \hat{Y} is from reaching demographic parity.

Theorem

For all function $a : \{0,1\} \rightarrow \{0,1\}$ the balanced accuracy of $a \circ \hat{Y}$ is lesser or equal than $\frac{1}{2} + \frac{1}{2}Fl(\hat{Y})$

Proof.

- We know that we have only four transformations of {0,1}: 0, 1, id, 1 - id.
- We compute the balanced accuracy for each of those applications. For instance for *id* the balanced accuracy is

$$\begin{split} &\frac{1}{2}(P(\hat{Y}=0|S=0)+P(\hat{Y}=1|S=1))\\ &=&\frac{1}{2}(1-P(\hat{Y}=1|S=0)+P(\hat{Y}=1|S=1))\\ &\leq&\frac{1}{2}+\frac{1}{2}FI(\hat{Y}) \end{split}$$

<□ ▶ < @ ▶ < E ▶ < E ▶ E り < C 13/22

Equality of odds doesn't forbid attribute inference attack

<ロト < 母 > < 喜 > < 喜 > 言 の へ つ 14/22

Theorem

For all function $a : \{0,1\} \rightarrow \{0,1\}$, if \hat{Y} satisfies Equality of odds, then the balanced accuracy of $a \circ \hat{Y}$ is equal to $\frac{1}{2}$ if and only if \hat{Y} is independent of Y or Y is independent of S.

So, either the prediction task is independent of the sensitive attribute or the model did not learn.

< □ ▶ < @ ▶ < ≧ ▶ < ≧ ▶ ≧ ⑦ Q ^Q 15/22

Proof.

$$\hat{S} = a \circ \hat{Y}$$

$$P(\hat{S} = 0|S = 0)$$

=P($\hat{S} = 0|S = 0Y = 0$)P(Y = 0|S = 0)
+P($\hat{S} = 0|S = 0Y = 1$)P(Y = 1|S = 0)

 $\left(P(\hat{Y} \in a^{-1}(\{1\}) | S = 1Y = 0) - P(\hat{Y} \in a^{-1}(\{1\}) | S = 1Y = 1) \right)$

▲□▶▲@▶▲분▶▲분▶ 분 외약은 16/22

$$P(\hat{S} = 0|S = 0)$$

= $P(\hat{S} = 0|S = 0Y = 0)P(Y = 0|S = 0)$
+ $P(\hat{S} = 0|S = 0Y = 1)P(Y = 1|S = 0)$

 $P(\hat{S} = 0|S = 0) + P(\hat{S} = 1|S = 1)$

=1 + (P(Y = 0|S = 0) - P(Y = 0|S = 1))

Fairness and attribute inference attack using **soft labels**

Demographic parity means "The prediction is independent of the sensitive attribute"

$$Fl(\hat{Y}) = 0 \Leftrightarrow P_{\hat{Y},S} = P_{\hat{Y}} \otimes P_S$$

No hypothesis on the support of \hat{Y}

(i.e. Demographic parity works for hard and soft labels, classification and regression problems)

Definition

 \hat{Y} satisfies demographic parity (for S) if and only if $P_{\hat{Y}}\otimes P_S$

Theorem

Let $\hat{Y} : \Omega \to E$ and $S : \Omega \to \{0,1\}$. \hat{Y} satisfies demographic parity for S if and only if, for all $a : E \to \{0,1\}$, the balanced accuracy of $a \circ \hat{Y}$ is equal to $\frac{1}{2}$

This time, we can't think in terms of number of applications. *a* results from any kind of classifier (random forest, neural network, ...)

Proof.

$$\forall a \ P(\hat{Y} \in a^{-1}(\{0\}) | S = 0) + P(\hat{Y} \in a^{-1}(\{1\}) | S = 1) = 1$$

$$\Leftrightarrow \forall a \ P(\hat{Y} \in a^{-1}(\{0\}) | S = 0) = P(\hat{Y} \in a^{-1}(\{0\}) | S = 1)$$

$$\Leftrightarrow \forall A \ P(\hat{Y} \in A) | S = 0) = P(\hat{Y} \in A | S = 1)$$

◆□ ▶ ◆ □ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ • • • ○ へ ○ 20/22

Experimental validation

Bounds for hard labels tested on tabular and image dataset

< □ ▶ < @ ▶ < ≧ ▶ < ≧ ▶ Ξ の Q @ 21/22

 Evaluation of attribute inference attacks against fairness enforcing mechanisms

Take away

 Generalized demographic parity is a tool to study attribute inference attack (for hard and soft labels)

◆□ ▶ ◆□ ▶ ◆ ■ ▶ ◆ ■ ▶ ● ■ の Q ○ 22/22

 Generalizing demographic parity introduces a notion of fairness in regression