

CS 594: 3- On Fairness and its Definitions – part 3

Abolfazl Asudeh Fall 2020 9/17/2020



Trade-offs and Impossibility Theorem

Slides prepared using: A Tutorial on Fairness in Machine Learning Ziyuan Zhong

Assumption

• Data is biased:

• Y is not independent from S





Suff.: SLY14(3)



If S and Y are not independent, then either independence holds or sufficiency, but not both
Brias: Y A S D
Tradep: J S Z

JKS But SIJI (Path S-f-J) JKY and SKJX.

Independence v.s. Separation

 If S and Y are not independent, then either independence holds or separation, but not both

() Bias: JKS Bep: JIS J (path s-y-f dues not exist) either JJLS: X. O (Bias) JLY: (vandem assignment) & useless



If S and Y are not independent, then either sufficiency holds or separation, but not both (under some assumptions...) Deriver Sty Deriver Sty Suf.: SLOP > S is indep. of ylf (Marginal on Row) 3 Sep.: SLF 14 > S is ~ ~ fly (~ ~ cl)

 \Rightarrow SL(f, y) \implies SLY X.

• Original classifier



• Fixing False positive rate disparity





• Introduces disparity on PP



Accuracy v.s. Fairness

• (Assumption) original classifier:

 $\min L(\theta)$

• Any intervention on the classifier will increase the loss

References:

- 1. Alexandra Chouldechova. "Fair prediction with disparate impact: A study of bias in recidivism prediction instruments.", arxiv.org/abs/1610.07524v1 (2016)
- 2. Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. "Inherent tradeoffs in the fair determination of risk scores." arXiv:1609.05807 (2016).
- 3. Sorelle Friedler, Carlos Scheidegger, and <u>Suresh Venkatasubramanian</u>. "On the (im) possibility of fairness." arXiv:1609.07236 (2016).

* link to related videos for [1,2] are posted in the course page

Individual Fairness and Subgroup Fairness

Individual Fairness

• Similar individual should have similar outcomes

l(x)

M

+

d(x, x')

$$\Delta\left(f_{\theta}(X_{1}), f_{\theta}(X_{2})\right) \leq \varepsilon\left(\frac{d(X_{1}, X_{2})}{\sum_{M_{on} \neq on_{i_{c}}} t_{o}}\right)$$

 Question: Can any <u>deterministic</u> (binary) classifier achieve individual fairness? (VC)

* Cynthia Dwork, et al. "Fairness through awareness." In ITCS. 2012

Subgroup fairness

- Here, fairness (parity) is defined over the intersection of demographic groups:
 - E.g.: {Hispanic, Female}
 - Hispanic Female
 - Hispanic Female under the age of 20
 - 0
- s/t between indiv. fairness and group fairness (subgroups are small)
- Challenge: Combinatorial Space ← #P problems

Fairness in Rankings

Fairness in Ranking

Ranking is a much more complex output than a binary class label.

Defining fairness is correspondingly more involved.

Toy Example





Sale -- Normalized Customer Satisfaction -- Normalized



Fairness at Top-k

The simplest measures consider the top-k, and then address it as if it were a classification task -- items ranked in the top-k have one label and the rest have another. Now we can use the entire fairness framework for classification.

This, of course, begs the question of choosing k.

Answers obtained could be very different depending on the value chosen.

Exposure-based assessment

In many ranking situations, such as in information retrieval or recommender systems, higher ranked items get more attention than lower ranked ones.

There often is not a hard cut-off. But it is possible to define a monotonically decreasing function, such as inverse rank, that quantifies how much attention an item gets.

Now we can aggregate the attention received per protected group, and make that into a criterion against which we assess fairness.

Probability-Based Assessment

If ordered lists were created separately for each protected group, and the lists were then merged at random, how likely is it that we will observe the ordering reported?

In a random merge of these sorted lists, it should be the case that for any pair of groups, if we consider all pairs of items in the list from these groups, there should be approximately as many pairs items with the group A item ranked higher as the number with group B items ranked higher. The difference in this number is a measure of bias, related to the probability measure discussed above.