

Preliminary Experiments to Examine the Stability of Bias-aware Techniques

Toshihiro Kamishima*, Shotaro Akaho*, Yukino Baba**, and Hisashi Kashima*** www.kamishima.net

*National Institute of Advanced Industrial Science and Technology (AIST), Japan **University of Tsukuba, Japan; ***Kyoto University, Japan

2nd International Workshop on Algorithmic Bias in Search and Recommendation (Bias 2021)

@ Online, 2021-04-01



- Stability : a new desirable property for bias-aware classifier
- Dataset to examine the stability
- Bias-aware models
- Discussion about the reasons why the model was unstable

Notations

S: Sensitive Feature

- Information that should not influence the prediction results
- Whether or not an item is more preferred due to cognitive bias

Y: Target Variable

- A predicted class, one item is more preferred to the other
- X : Non-sensitive features
 - A pair of items shown to the subjects
- *S* \bot *Y* : Statistical Parity
 - statistical independence between S and Y
 - No information about S influences the decision Y

Bias-aware Classifier

Bias-aware classifier : The most accurate predictor while satisfying a fairness constraint



The loss is evaluated on a biased dataset, because an unbiased dataset is unavailable However, the loss should be evaluated on an unbiased dataset

Other desirable property for bias-aware classifier

Stability of Bias-aware Classifier

Bias-aware Classifier : Removing the influence of S on Y

Stability of Bias-aware Classifier

Two datasets consist of the same information except for the information represented by *S*

A bias-aware classifier should learn the same models



Dataset to Examine the Stability



Cognitive Biases: Positional Effect

Positional Effect : the item displayed near the upper-left corner of the interface screen is more frequently chosen.



• **Baseline** : the positions of items are randomly assigned, ideal RCT

• Fixed : the positions of items are affected by the popularity of items

Cognitive Biases: Bandwagon Effect

Bandwagon Effect : the item indicated that other people prefer is more frequently chosen



 Bandwagon : the items indicated as popular are affected by the popularity of items as in the same procedure as the Fixed procedure

Measurement of Cognitive Bias



As we designed, our datasets are influenced by cognitive biases

Random	Fixed	Bandwagon
0.0229	0.0077	0.3451

* Positive value indicates the influence of cognitive biases

Bias-aware Model



Bias-aware Model

This stratification technique has connection with a generative type of a fairness-aware model

biased generative classifier

$$Pr[Y, \mathbf{X}, S] = Pr[S] Pr[Y|S] Pr[\mathbf{X}|Y, S] \qquad \text{fairness constraint} \\ = Pr[S] Pr[Y] Pr[\mathbf{X}|Y, S] \qquad (\leftarrow Y \perp S) \\ = Pr[Y] \left\{ Pr[S] Pr[\mathbf{X}|Y, S] \right\} \qquad \text{effect of } Y \text{ on } \mathbf{X} \\ \text{size of stratum} \qquad \text{total effect of } Y \text{ on } \mathbf{X} \\ \end{array}$$

Removing Biases by Stratification

Hypothesis : If a stratification technique is stable as a bias-aware classifier, similar models would be learned from three types of datasets: Random, Fixed, and Bandwagon

For each dataset, we get the probabilities $\Pr[Y = 1 | \mathbf{X} = \mathbf{x}]$ for all \mathbf{x}

If stratification is stable, these probabilities would be similar

Similarities (Frobenius norm) between a pair of probability matrices

Random vs Fixed	Random vs Bandwagon
0.0521	0.127

Contrary to our hypothesis, this stratification techniques are not stable

Next Steps

Cognitive biases may not be fully removed by the stratification

- The scheme of data processing might change the causal structure
 - \Rightarrow S may behave as a mediator as well as a confounder
- We will explore other types of causal structure

Confounders other than the controlled cognitive bias might exist

- We showed "Popular" marks irrelevant to its real popularity
 This might cause another type of cognitive bias
- We plan to collect data influenced by other cognitive bias

Conclusions

- The notion of stability of bias-aware techniques
- A dataset to examine the stability of bias-aware techniques
- The relationship between bias-aware techniques and causal inference
- Preliminary experimental results, showing the instability of our techniques
- Next steps for collecting new datasets or for developing new models