

## バイアス考慮型分類器の<br/>安定性に関する予備調査

神嶌 敏弘<sup>1</sup>,赤穂 昭太郎<sup>1</sup>,馬場 雪乃<sup>2</sup>,鹿島 久嗣<sup>3</sup> <sup>1</sup>産業技術総合研究所,<sup>2</sup>筑波大学,<sup>3</sup>京都大学 2021年度人工知能学会全国大会(第35回)@ オンライン,2021-06-10 <u>https://www.kamishima.net</u>



- 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**

# Y target variable / object variable An objective of decision making, or what to predict ex., loan approval, university admission, what to recommend Y: observed / true, Ŷ: predicted

#### sensitive feature

#### To ignore the influence to the sensitive feature from a target

ex., socially sensitive information (gender, race), items' brand

- Specified by a user or an analyst depending on his/her purpose
- It may depend on a target or other features
  - non-sensitive feature vector

All features other than a sensitive feature

#### **Independence / Statistical Parity**

[Calders+ 10, Dwork+ 12]

**Remove data bias**  $\Rightarrow$  **Independence / Statistical Parity:**  $\hat{Y} \perp S$ 



$$\hat{Y} \perp S$$
$$\parallel$$
$$\mathsf{II}$$
$$\mathsf{Pr}[\hat{Y}, S] = \mathsf{Pr}[\hat{Y}] \mathsf{Pr}[S]$$

Ratios between positives and negatives in prediction are matched among all sensitive values

### **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

## Failure of Random Assignment

#### select item pairs



#### **Cognitive biases may not be fully removed by the stratification**

- The scheme of data processing might change the causal structure
  - ➡ 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

ask subjects to choose

- 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