Fairness-aware Classifier with Prejudice Remover Regularizer

Toshihiro Kamishima*, Shotaro Akaho*, Hideki Asoh*, and Jun Sakuma†

*National Institute of Advanced Industrial Science and Technology (AIST), Japan †University of Tsukuba, Japan; and Japan Science and Technology Agency Test of Time Awards Talk @ ECMLPKDD 2022, 2022-09-21

BARCELONA, SPAIN ECMLPKDD2016

European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases









$-\Pr[\hat{Y}=1 \mid S=0] - \rightarrow$	Ŷ=1 Þ0 S=œsti	sitive $\hat{Y}=1$ mation $3=1$	$Pr[\hat{Y}=1 S=$
$- Pr[\hat{Y}=0 S=0]$	$\hat{Y}=0$ S=0	$\hat{Y}=0$ S=1	=1] \longrightarrow $\Pr[\hat{Y}=0 \mid S=1] \rightarrow$
	$\bullet Pr[S=0] \longrightarrow$	$\bullet \qquad Pr[S=1] \longrightarrow$	





Two Naïve Bayes method

Formulation as a Constrained Optimization Problem

Formulation as a Constrained Optimization Problem

Minimize Objective Function

fairness constraint -I(Y, S)

balancing fairness and utility

However, it was the beginning of the long winding road

β1 Model

This model was not converged

β2 Model

of Data Mining, 2011

"Fairness-aware Classifier with Prejudice Remover Regularizer", ECMLPKDD, 2012

Actual Independence

Actual Independence: Class labels are deterministically generated by applying a decision rule

effectiveness improvement

"Model-based and Actual Independence for Fairness-aware Classification", Data Mining and Knowledge Discovery, 2018

Fair Recommendation

Prediction Function

a prediction function is selected according to a sensitive value

$$\hat{r}(x, y, s) = \mu^{(s)} + b_x^{(s)} + c_y^{(s)} + \mathbf{p}_x^{(s)} \mathbf{q}_y^{(s)^{\top}}$$
sensitive feature

Objective Function independence parameter: control the balance between the independence and accuracy

$$\sum_{\mathcal{D}} (r_i - \hat{r}(x_i, y_i))^2 - \eta \operatorname{indep}(R, S) + \lambda \|\Theta\|^2$$

"Enhancement of the Neutrality in Recommendation", The 2nd Workshop on Human Decision Making in Recommender Systems, 2012

"Efficiency Improvement of Neutrality-enhanced Recommendation", The 3rd Workshop on Human Decision Making in Recommender Systems, 2013

"Recommendation Independence", Proc of the Conf. on Fairness, Accountability and Transparency, 2018

Interests of Data Science Communities

My survey slide about Fairness-Aware Machine Learning is available at:

https://www.kamishima.net/faml/

Acknowledgements

- We'd like to be appreciate the colleagues in our fairness-aware ML research, Hiroshi Nakagawa, Hiromi Arai, Kazuto Fukuchi, Shoko Ema, and Junpei Komiyama.
- We greatly respect the pioneering works of the fairness-aware ML by the European community, Toon Calders, Dino Pedreschi, Indrė Žliobaitė, Salvatore Ruggieri, Faisal Kamiran, and Sara Hajian.
- We'd like to express our thanks for the cooperation of the fair RecSys community: Michael Ekstrand, Robin Burke, Pierre-Nicolas Schwab, Jean Garcia-Gathright, Nasim Sonboli, Luca Belli, and Amifa Raj.
- We also thank to the FAccT community, Suresh Venkatasubramanian, Sorelle Friedler, and Solon Barocas, Krishna Gummadi, and Md. Bilal Zafar.
- Thanks for the support of the JSPS KAKENHI 科研算 Grant Number 24500194, 15K00327, 18H03300, and 21H03504, and of our affiliations; National Institute of Advanced Industrial Science and Technology (AIST) and University of Tsukuba.