

# Fair Bayes-Optimal Classifiers under Predictive Parity

---

Xianli Zeng, Edgar Dobriban, Guang Cheng  
January 30, 2023

Reviewer : Park Seok Hoon , Seoul National University IDEA LAB

# Table of Contents

- ① Contribution
- ② Preliminaries
- ③ Main Theorem
- ④ Algorithm : FairBayes-DPP
- ⑤ results

# Table of Contents

- 1 Contribution
- 2 Preliminaries
- 3 Main Theorem
- 4 Algorithm : FairBayes-DPP
- 5 results

# Contribution

- ① They identify a sufficient condition under which all fair Bayes-optimal classifiers are GWTR
- ② They proposed the FairBayes-DPP algorithm for binary fair classification.

# Table of Contents

- ① Contribution
- ② Preliminaries
- ③ Main Theorem
- ④ Algorithm : FairBayes-DPP
- ⑤ results

# Preliminaries

- $X \in \mathcal{X}$  : Feature data
- $A \in \mathcal{A}$  : protected variable
- $Y \in \{0, 1\}$  : ground truth label
- $\eta_a(x) = P(Y = 1|X = x, A = a)$

## Definition of Randomized Classifier

A randomized classifier is a measurable function  $f : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$ , indicating the probability of predicting  $\hat{Y} = 1$  when observing  $X = x$  and  $A = a$ . We denote  $\hat{Y}_f = \hat{Y}_f(x, a)$  the prediction induced by the classifier  $f$

## Definition of GWTR classifier

A classifier  $f$  is a GWTR if there are constants  $t_a, a \in \mathcal{A}$ , and functions  $r_a : \mathcal{X} \rightarrow [0, 1], a \in \mathcal{A}$  such that for all  $x \in \mathcal{X}$  and  $a \in \mathcal{A}$

$$f(x, a) = I(\eta_a(x) > t_a) + r_a(x)I(\eta_a(x) = t_a)$$



## Definition of Predictive parity and DPP

A classifier  $f$  satisfies predictive parity if for all  $a \in \mathcal{A}$ ,

$$P(Y = 1 | \hat{Y}_f = 1, A = a) = P(Y = 1 | \hat{Y}_f = 1)$$

Also,

$$DPP(f) = \sum_{a \in \mathcal{A}} |P(Y = 1 | \hat{Y}_f = 1, A = a) - P(Y = 1 | \hat{Y}_f = 1)|$$

## Definition of cost sensitive zero-one risk

When the predictive parity is commonly considered, the false positives are particularly harmful. So, we consider cost-sensitive classification.

For a cost parameter  $c \in [0, 1]$ , the cost-sensitive zero-one risk of the classifier  $f$  is as follows

$$R_c(f) = cP(\hat{Y}_f = 1, Y = 0) + (1 - c)P(\hat{Y}_f = 0, Y = 1)$$

## Main contribution in paper

In this paper, they proposed fair Bayes-optimal classifier that satisfied Predictive Parity.

$$f_{PPV}^* = \operatorname{argmin}_{f: DPP(f)=0} R_c(f)$$

# Table of Contents

- ① Contribution
- ② Preliminaries
- ③ Main Theorem**
- ④ Algorithm : FairBayes-DPP
- ⑤ results

Condition (\*)

$$\min_{a \in \mathcal{A}} P(Y = 1 | \eta_a(X) \geq c, A = a) \geq \max_{a \in \mathcal{A}} P(Y = 1 | A = a)$$

: the performances of different groups vary only moderately

## Theorem 1 - Main

Consider the cost-sensitive 0-1 risk with cost parameter  $c$ . If the condition (\*) holds, then all fair Bayes-optimal classifiers under predictive parity are GWTR. Thus, for all  $f_{PPV}^*$ , there are  $(t_a^*)_{a=1}^{|\mathcal{A}|}$  and functions  $r_a^*(x) : \mathcal{X} \rightarrow [0, 1]$  such that, for all  $x \in \mathcal{X}$  and  $a \in \mathcal{A}$ ,

$$f_{PPV}^*(x, a) = I(\eta_a(x) > t_a^*) + r_a^*(x)I(\eta_a(x) = t_a^*)$$

# Table of Contents

- ① Contribution
- ② Preliminaries
- ③ Main Theorem
- ④ Algorithm : FairBayes-DPP
- ⑤ results

- The DPP constraint is non-convex with respect to the model  $\hat{\eta}_a(x)$  parameters.
- In such cases, incorporating fairness constraints as a penalty in the training objective may be hard due to potential local minima.
- So, they consider post-processing algorithms.



## Step 1

- They apply standard ML algorithms to learn the feature and group-conditional label probabilities  $\eta_a(X)$  based on the whole datasets.
- $L$  : loss function,  $\mathcal{G} = \{g_\theta, \theta \in \Theta\}$
- $\hat{\eta}_a(X) = g_{\hat{\theta}}(x, a)$   
where  $\hat{\theta} \in \operatorname{argmin}_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n L(y_i, g_\theta(x_i, a_i))$

- If the empirical version of Condition (\*) holds,

Find threshold  $t_a, a \in \mathcal{A}$  which minimize cost-sensitive risk and satisfy sample predictive parity by using gridsearch.

# Table of Contents

- ① Contribution
- ② Preliminaries
- ③ Main Theorem
- ④ Algorithm : FairBayes-DPP
- ⑤ results**

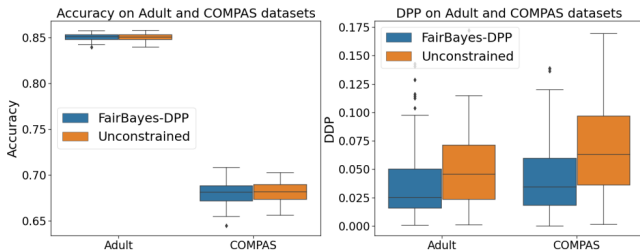
They sampled 70%, 50% and 30% as training, validation and test dataset with replacement.

- ① Adult  
: protected variable : gender
- ② Compas  
: protected variable : race
- ③ CelebA  
: protected variable : gender

- They experiments many times by sampling train, test and validation datasets.
- cost parameter  $c = 0.5$

# Adult and Compas

- The conditional probability  $\eta$  is learned via a three-layer MLP with 32 hidden neurons per layer.
- Batch size of Adult : 512, Batch size of Compas : 2048



- The conditional probability  $\eta$  is learned via a Resnet50.

ATTRIBUTES	PER-ATTRIBUTE ACCURACY		PER-ATTRIBUTE DPP	
	FAIRBAYES-DPP	UNCONSTRAINED	FAIRBAYES-DPP	UNCONSTRAINED
ARCHED EYEBROWS	0.838(0.003)	0.838(0.003)	0.027(0.015)	0.099(0.041)
ATTRACTIVE	0.825(0.002)	0.826(0.003)	0.075(0.011)	0.169(0.016)
BAGS UNDER EYES	0.853(0.002)	0.852(0.002)	0.024(0.015)	0.056(0.034)
BANGS	0.959(0.001)	0.959(0.001)	0.007(0.007)	0.069(0.029)
BIG LIPS	0.706(0.002)	0.717(0.003)	0.023(0.015)	0.115(0.027)
BIG NOSE	0.845(0.002)	0.847(0.003)	0.083(0.020)	0.145(0.023)